

# Econometric Evaluation of Social Programs

## Part II:

Using the Marginal Treatment Effect to Organize Alternative Econometric Estimators to Evaluate Social Programs, and to Forecast their Effects in New Environments

James J. Heckman  
University of Chicago  
American Bar Foundation  
University College Dublin, Ireland

Edward J. Vytlačil  
Columbia University

Econ 312, Spring 2021

## Introduction

- This part of our contribution to this Handbook reviews and extends the econometric literature on the evaluation of social policy.

## Introduction

- This part of our contribution to this Handbook reviews and extends the econometric literature on the evaluation of social policy.
- We organize our discussion around choice-theoretic models for objective and subjective outcomes of the sort discussed in Part I.

## Introduction

- This part of our contribution to this Handbook reviews and extends the econometric literature on the evaluation of social policy.
- We organize our discussion around choice-theoretic models for objective and subjective outcomes of the sort discussed in Part I.
- Specifically, we organize our discussion of the literature around the concept of the marginal treatment effect (MTE) that was introduced in Part I.



- Using the marginal treatment effect, we define a variety of treatment effects and show how they can be generated by a single economic functional, the MTE.

- Using the marginal treatment effect, we define a variety of treatment effects and show how they can be generated by a single economic functional, the MTE.
- We then show what various econometric methods assume about the MTE.

- In this part, we focus exclusively on microeconomic partial equilibrium evaluation methods, deferring analysis of general equilibrium issues to Part III.

- In this part, we focus exclusively on microeconomic partial equilibrium evaluation methods, deferring analysis of general equilibrium issues to Part III.
- Thus throughout this chapter, except when we discuss randomized evaluation of social programs, we assume that potential outcomes are not affected by interventions but choices among the potential outcomes are affected.

- In this part, we focus exclusively on microeconomic partial equilibrium evaluation methods, deferring analysis of general equilibrium issues to Part III.
- Thus throughout this chapter, except when we discuss randomized evaluation of social programs, we assume that potential outcomes are not affected by interventions but choices among the potential outcomes are affected.
- Thus, we invoke policy invariance assumptions (PI-3) and (PI-4) of Part I.

- In this part, we focus exclusively on microeconomic partial equilibrium evaluation methods, deferring analysis of general equilibrium issues to Part III.
- Thus throughout this chapter, except when we discuss randomized evaluation of social programs, we assume that potential outcomes are not affected by interventions but choices among the potential outcomes are affected.
- Thus, we invoke policy invariance assumptions (PI-3) and (PI-4) of Part I.
- We also focus primarily on mean responses, leaving analysis of distributions of responses for Part III.

- The plan of Part II is as follows.

- The plan of Part II is as follows.
- In Slide 12, we present some basic principles that underlie conventional econometric evaluation estimators.



- The plan of Part II is as follows.
- In Slide 12, we present some basic principles that underlie conventional econometric evaluation estimators.
- In Slide 90, we define the marginal treatment effect in a two potential outcome model that is a semiparametric version of the generalized Roy model.

- The plan of Part II is as follows.
- In Slide 12, we present some basic principles that underlie conventional econometric evaluation estimators.
- In Slide 90, we define the marginal treatment effect in a two potential outcome model that is a semiparametric version of the generalized Roy model.
- We then show how treatment parameters can be generated as weighted averages of the MTE.

- We carefully distinguish the definition of parameters from issues of identification.

- We carefully distinguish the definition of parameters from issues of identification.
- Slide 152 considers how instrumental variable methods that supplement the classical instrumental variable assumptions of econometrics can be used to identify treatment parameters.

- We carefully distinguish the definition of parameters from issues of identification.
- Slide 152 considers how instrumental variable methods that supplement the classical instrumental variable assumptions of econometrics can be used to identify treatment parameters.
- We discuss the crucial role of monotonicity assumptions in the recent IV literature.

- They impart an asymmetry to the admissible forms of agent heterogeneity.

- They impart an asymmetry to the admissible forms of agent heterogeneity.
- Outcomes are permitted to be heterogeneous in a general way but responses of choices to external inputs are not.

- They impart an asymmetry to the admissible forms of agent heterogeneity.
- Outcomes are permitted to be heterogeneous in a general way but responses of choices to external inputs are not.
- When heterogeneity in choices and outcomes is allowed, the IV enterprise breaks down.



- They impart an asymmetry to the admissible forms of agent heterogeneity.
- Outcomes are permitted to be heterogeneous in a general way but responses of choices to external inputs are not.
- When heterogeneity in choices and outcomes is allowed, the IV enterprise breaks down.
- Treatment parameters can still be defined but IV does not identify them.

- Slide 402 extends our analysis to consider regression discontinuity estimators introduced in ? and adapted to modern econometrics in ?.

- Slide 402 extends our analysis to consider regression discontinuity estimators introduced in ? and adapted to modern econometrics in ?.
- We interpret the regression discontinuity estimator within the MTE framework, as a special type of IV estimator.

- Slide 402 extends our analysis to consider regression discontinuity estimators introduced in ? and adapted to modern econometrics in ?.
- We interpret the regression discontinuity estimator within the MTE framework, as a special type of IV estimator.
- In Slide 412, we show how the output of the IV analysis of Slide 152 can be used to extend parameters identified in one population to other populations and to forecast the effects of new programs.

- Slide 402 extends our analysis to consider regression discontinuity estimators introduced in ? and adapted to modern econometrics in ?.
- We interpret the regression discontinuity estimator within the MTE framework, as a special type of IV estimator.
- In Slide 412, we show how the output of the IV analysis of Slide 152 can be used to extend parameters identified in one population to other populations and to forecast the effects of new programs.
- These are questions P-2 and P-3 introduced in Part I. Slides 12–402 focus solely on the problem of internal validity, which is the problem defined as P-1.

- We also develop a cost benefit analysis based on the MTE and we analyze marginal policy changes.

- We also develop a cost benefit analysis based on the MTE and we analyze marginal policy changes.
- In Slide 471, we generalize the analysis of instrumental variables to consider models with multiple outcomes.

- We also develop a cost benefit analysis based on the MTE and we analyze marginal policy changes.
- In Slide 471, we generalize the analysis of instrumental variables to consider models with multiple outcomes.
- We develop both unordered and ordered choice models linking them to an explicit choice-theoretic literature.



- In Slide 675, we consider matching as a special case of our framework.

- In Slide 675, we consider matching as a special case of our framework.
- Matching applied to estimating conditional means is a version of nonparametric least squares.

- In Slide 675, we consider matching as a special case of our framework.
- Matching applied to estimating conditional means is a version of nonparametric least squares.
- It assumes that marginal and average returns are the same whereas our general framework allows us to distinguish marginal from average returns and to identify both.

- In Slide 675, we consider matching as a special case of our framework.
- Matching applied to estimating conditional means is a version of nonparametric least squares.
- It assumes that marginal and average returns are the same whereas our general framework allows us to distinguish marginal from average returns and to identify both.
- Matching is more robust than IV to violations of conventional monotonicity assumptions but the price for this robustness is steep in terms of its economic content.

- In Slide 675, we consider matching as a special case of our framework.
- Matching applied to estimating conditional means is a version of nonparametric least squares.
- It assumes that marginal and average returns are the same whereas our general framework allows us to distinguish marginal from average returns and to identify both.
- Matching is more robust than IV to violations of conventional monotonicity assumptions but the price for this robustness is steep in terms of its economic content.
- In Slide 819, we develop randomization as an instrumental variable.

- We consider problems with compliance induced by agent self-selection decisions.

- We consider problems with compliance induced by agent self-selection decisions.
- In Slide 938, we consider how to bound the various treatment parameters when models are not identified.

- We consider problems with compliance induced by agent self-selection decisions.
- In Slide 938, we consider how to bound the various treatment parameters when models are not identified.
- Slide 1005 develops alternative methods for controlling for selection: control functions, replacement functions and proxy variables.



- We consider problems with compliance induced by agent self-selection decisions.
- In Slide 938, we consider how to bound the various treatment parameters when models are not identified.
- Slide 1005 develops alternative methods for controlling for selection: control functions, replacement functions and proxy variables.
- Slide 1028 concludes.

## The Basic Principles Underlying the Identification of the Major Econometric Evaluation Estimators

- In this section, we review the main principles underlying the major evaluation estimators used in the econometric literature.

## The Basic Principles Underlying the Identification of the Major Econometric Evaluation Estimators

- In this section, we review the main principles underlying the major evaluation estimators used in the econometric literature.
- We assume two potential outcomes  $(Y_0, Y_1)$ .

## The Basic Principles Underlying the Identification of the Major Econometric Evaluation Estimators

- In this section, we review the main principles underlying the major evaluation estimators used in the econometric literature.
- We assume two potential outcomes  $(Y_0, Y_1)$ .
- Models for multiple outcomes are developed in later sections of this chapter.

## The Basic Principles Underlying the Identification of the Major Econometric Evaluation Estimators

- In this section, we review the main principles underlying the major evaluation estimators used in the econometric literature.
- We assume two potential outcomes  $(Y_0, Y_1)$ .
- Models for multiple outcomes are developed in later sections of this chapter.
- As in Part I,  $D = 1$  if  $Y_1$  is observed, and  $D = 0$  corresponds to  $Y_0$  being observed.

- The observed objective outcome is

$$Y = DY_1 + (1 - D)Y_0. \quad (1)$$

- The observed objective outcome is

$$Y = DY_1 + (1 - D)Y_0. \quad (1)$$

- To briefly recapitulate the lessons of Part I, we distinguish two distinct econometric problems.

- The observed objective outcome is

$$Y = DY_1 + (1 - D)Y_0. \quad (1)$$

- To briefly recapitulate the lessons of Part I, we distinguish two distinct econometric problems.
- For simplicity, we focus our discussion on identification of objective outcomes.



- The observed objective outcome is

$$Y = DY_1 + (1 - D)Y_0. \quad (1)$$

- To briefly recapitulate the lessons of Part I, we distinguish two distinct econometric problems.
- For simplicity, we focus our discussion on identification of objective outcomes.
- A parallel analysis can be made for subjective outcomes.

- The *evaluation problem* arises because for each person we observe either  $Y_0$  or  $Y_1$  but not both.

- The *evaluation problem* arises because for each person we observe either  $Y_0$  or  $Y_1$  but not both.
- Thus, in general, it is not possible to identify the individual level treatment effect  $Y_1 - Y_0$  for any person.

- The *evaluation problem* arises because for each person we observe either  $Y_0$  or  $Y_1$  but not both.
- Thus, in general, it is not possible to identify the individual level treatment effect  $Y_1 - Y_0$  for any person.
- The typical solution to this problem is to reformulate the problem at the population level rather than at the individual level and to identify certain mean outcomes or quantile outcomes or various distributions of outcomes as described in Part I.

- The *evaluation problem* arises because for each person we observe either  $Y_0$  or  $Y_1$  but not both.
- Thus, in general, it is not possible to identify the individual level treatment effect  $Y_1 - Y_0$  for any person.
- The typical solution to this problem is to reformulate the problem at the population level rather than at the individual level and to identify certain mean outcomes or quantile outcomes or various distributions of outcomes as described in Part I.
- For example, a common approach is to focus attention on average treatment effects, such as  $ATE = E(Y_1 - Y_0)$ .

- If treatment is assigned or chosen on the basis of potential outcomes, so

$$(Y_0, Y_1) \not\perp D,$$

where  $\not\perp$  denotes “is not independent” and “ $\perp$ ” denotes independent, we encounter the problem of selection bias.

- If treatment is assigned or chosen on the basis of potential outcomes, so

$$(Y_0, Y_1) \not\perp D,$$

where  $\not\perp$  denotes “is not independent” and “ $\perp$ ” denotes independent, we encounter the problem of selection bias.

- Suppose that we observe people in each treatment state  $D = 0$  and  $D = 1$ .

- If treatment is assigned or chosen on the basis of potential outcomes, so

$$(Y_0, Y_1) \not\perp D,$$

where  $\not\perp$  denotes “is not independent” and “ $\perp$ ” denotes independent, we encounter the problem of selection bias.

- Suppose that we observe people in each treatment state  $D = 0$  and  $D = 1$ .
- If  $Y_j \not\perp D$ , then the observed  $Y_j$  will be selectively different from randomly assigned  $Y_j$ ,  $j = 0, 1$ .



- If treatment is assigned or chosen on the basis of potential outcomes, so

$$(Y_0, Y_1) \not\perp D,$$

where  $\not\perp$  denotes “is not independent” and “ $\perp$ ” denotes independent, we encounter the problem of selection bias.

- Suppose that we observe people in each treatment state  $D = 0$  and  $D = 1$ .
- If  $Y_j \not\perp D$ , then the observed  $Y_j$  will be selectively different from randomly assigned  $Y_j$ ,  $j = 0, 1$ .
- Thus  $E(Y_0 \mid D = 0) \neq E(Y_0)$  and  $E(Y_1 \mid D = 1) \neq E(Y_1)$ .

- Using unadjusted data to construct  $E(Y_1 - Y_0)$  will produce bias:

$$E(Y_1 | D = 1) - E(Y_0 | D = 0) \neq E(Y_1 - Y_0).$$

The selection problem is a key aspect of the problem of evaluating social programs.

- Using unadjusted data to construct  $E(Y_1 - Y_0)$  will produce bias:

$$E(Y_1 | D = 1) - E(Y_0 | D = 0) \neq E(Y_1 - Y_0).$$

The selection problem is a key aspect of the problem of evaluating social programs.

- Many methods have been proposed to solve both problems.

- Using unadjusted data to construct  $E(Y_1 - Y_0)$  will produce bias:

$$E(Y_1 | D = 1) - E(Y_0 | D = 0) \neq E(Y_1 - Y_0).$$

The selection problem is a key aspect of the problem of evaluating social programs.

- Many methods have been proposed to solve both problems.
- This chapter unifies these methods using the concept of the marginal treatment effect (MTE) introduced in Part I of this Handbook.

- The method with the greatest intuitive appeal, which is sometimes called the “gold standard” in evaluation analysis, is the method of random assignment.

- The method with the greatest intuitive appeal, which is sometimes called the “gold standard” in evaluation analysis, is the method of random assignment.
- Nonexperimental methods can be organized by how they attempt to approximate what can be obtained by an ideal random assignment.

If treatment is chosen at random with respect to  $(Y_0, Y_1)$ , or if treatments are randomly assigned and there is full compliance with the treatment assignment,

(R-1)

$(Y_0, Y_1) \perp\!\!\!\perp D.$

- It is useful to distinguish several cases where (R-1) will be satisfied.



- It is useful to distinguish several cases where (R-1) will be satisfied.
- The first is that agents (decision makers whose choices are being investigated) pick outcomes that are random with respect to  $(Y_0, Y_1)$ .

- It is useful to distinguish several cases where (R-1) will be satisfied.
- The first is that agents (decision makers whose choices are being investigated) pick outcomes that are random with respect to  $(Y_0, Y_1)$ .
- Thus agents may not know  $(Y_0, Y_1)$  at the time they make their choices to participate in treatment or at least do not act on  $(Y_0, Y_1)$ , so that  $\Pr(D = 1 | X, Y_0, Y_1) = \Pr(D = 1 | X)$  for all  $X$ .

- It is useful to distinguish several cases where (R-1) will be satisfied.
- The first is that agents (decision makers whose choices are being investigated) pick outcomes that are random with respect to  $(Y_0, Y_1)$ .
- Thus agents may not know  $(Y_0, Y_1)$  at the time they make their choices to participate in treatment or at least do not act on  $(Y_0, Y_1)$ , so that  $\Pr(D = 1 | X, Y_0, Y_1) = \Pr(D = 1 | X)$  for all  $X$ .
- Matching assumes a version of (R-1) conditional on matching variables  $X$ :  $(Y_0, Y_1) \perp\!\!\!\perp D | X$ .

- A second case arises when individuals are randomly assigned to treatment status even if they would choose to self select into no-treatment status, and they comply with the randomization protocols.

- A second case arises when individuals are randomly assigned to treatment status even if they would choose to self select into no-treatment status, and they comply with the randomization protocols.
- Let  $\xi$  be randomized assignment status.

- A second case arises when individuals are randomly assigned to treatment status even if they would choose to self select into no-treatment status, and they comply with the randomization protocols.
- Let  $\xi$  be randomized assignment status.
- With full compliance,  $\xi = 1$  implies that  $Y_1$  is observed and  $\xi = 0$  implies that  $Y_0$  is observed.

Then, under randomized assignment,

(R-2)

$$(Y_0, Y_1) \perp\!\!\!\perp \xi,$$

even if in a regime of self-selection,  $(Y_0, Y_1) \not\perp\!\!\!\perp D$ .

- If randomization is performed conditional on  $X$ , we obtain  $(Y_0, Y_1) \perp\!\!\!\perp \xi \mid X$ .



- Let  $A$  denote actual treatment status.

- Let  $A$  denote actual treatment status.
- If the randomization has full compliance among participants,  $\xi = 1 \Rightarrow A = 1$ ;  $\xi = 0 \Rightarrow A = 0$  . This is entirely consistent with a regime in which a person would choose  $D = 1$  in the absence of randomization, but would have no treatment ( $A = 0$ ) if suitably randomized, even though the agent might desire treatment.

- If treatment status is chosen by self-selection,  $D = 1 \Rightarrow A = 1$  and  $D = 0 \Rightarrow A = 0$ .

- If treatment status is chosen by self-selection,  $D = 1 \Rightarrow A = 1$  and  $D = 0 \Rightarrow A = 0$ .
- If there is imperfect compliance with randomization,  $\xi = 1 \not\Rightarrow A = 1$  because of agent choices.

- If treatment status is chosen by self-selection,  $D = 1 \Rightarrow A = 1$  and  $D = 0 \Rightarrow A = 0$ .
- If there is imperfect compliance with randomization,  $\xi = 1 \not\Rightarrow A = 1$  because of agent choices.
- In general,  $A = \xi D$  so that  $A = 1$  only if  $\xi = 1$  and  $D = 1$ .

- If treatment status is chosen by self-selection,  $D = 1 \Rightarrow A = 1$  and  $D = 0 \Rightarrow A = 0$ .
- If there is imperfect compliance with randomization,  $\xi = 1 \not\Rightarrow A = 1$  because of agent choices.
- In general,  $A = \xi D$  so that  $A = 1$  only if  $\xi = 1$  and  $D = 1$ .
- This assumes that persons randomized out of the program cannot participate in it.

- If treatment status is chosen by self-selection,  $D = 1 \Rightarrow A = 1$  and  $D = 0 \Rightarrow A = 0$ .
- If there is imperfect compliance with randomization,  $\xi = 1 \not\Rightarrow A = 1$  because of agent choices.
- In general,  $A = \xi D$  so that  $A = 1$  only if  $\xi = 1$  and  $D = 1$ .
- This assumes that persons randomized out of the program cannot participate in it.
- If treatment status is randomly assigned, either through randomization or randomized self-selection,

(R-3)

$(Y_0, Y_1) \perp\!\!\!\perp A.$



- This version of randomization can also be defined conditional on  $X$ .

- This version of randomization can also be defined conditional on  $X$ .
- Under (R-1), (R-2) or (R-3), the average treatment effect (ATE) is the same as the marginal treatment effect and the parameters treatment on the treated (TT) and treatment on the untreated (TUT) as defined in Part I:

$$TT = MTE = TUT = ATE = E(Y_1 - Y_0) = E(Y_1) - E(Y_0).$$

- Observe that even with random assignment of treatment status and full compliance, we cannot, in general, identify the distribution of the treatment effects  $(Y_1 - Y_0)$ , although we can identify the marginal distributions

$$F_1(Y_1 | A = 1, X = x) = F_1(Y_1 | X = x) \text{ and}$$

$$F_0(Y_0 | A = 0, X = x) = F_0(Y_0 | X = x).$$

- Observe that even with random assignment of treatment status and full compliance, we cannot, in general, identify the distribution of the treatment effects  $(Y_1 - Y_0)$ , although we can identify the marginal distributions

$$F_1(Y_1 | A = 1, X = x) = F_1(Y_1 | X = x) \text{ and}$$

$$F_0(Y_0 | A = 0, X = x) = F_0(Y_0 | X = x).$$

- One special assumption, common in the conventional econometrics literature, is that  $Y_1 - Y_0 = \Delta(x)$ , a constant given  $x$ .

- Observe that even with random assignment of treatment status and full compliance, we cannot, in general, identify the distribution of the treatment effects  $(Y_1 - Y_0)$ , although we can identify the marginal distributions

$$F_1(Y_1 | A = 1, X = x) = F_1(Y_1 | X = x) \text{ and}$$

$$F_0(Y_0 | A = 0, X = x) = F_0(Y_0 | X = x).$$

- One special assumption, common in the conventional econometrics literature, is that  $Y_1 - Y_0 = \Delta(x)$ , a constant given  $x$ .
- Since  $\Delta(x)$  can be identified from  $E(Y_1 | A = 1, X = x) - E(Y_0 | A = 0, X = x)$  because  $A$  is allocated by randomization, the analyst can identify the joint distribution of  $(Y_0, Y_1)$ .

- However, this approach assumes that  $(Y_0, Y_1)$  have the same distribution up to a parameter  $\Delta$  ( $Y_0$  and  $Y_1$  are perfectly dependent).

- However, this approach assumes that  $(Y_0, Y_1)$  have the same distribution up to a parameter  $\Delta$  ( $Y_0$  and  $Y_1$  are perfectly dependent).
- One can make other assumptions about the dependence across ranks from perfect positive or negative ranking to independence.

- However, this approach assumes that  $(Y_0, Y_1)$  have the same distribution up to a parameter  $\Delta$  ( $Y_0$  and  $Y_1$  are perfectly dependent).
- One can make other assumptions about the dependence across ranks from perfect positive or negative ranking to independence.
- In general, the joint distribution of  $(Y_0, Y_1)$  or of  $(Y_1 - Y_0)$  is not identified unless the analyst can pin down the dependence across  $(Y_0, Y_1)$ .



- However, this approach assumes that  $(Y_0, Y_1)$  have the same distribution up to a parameter  $\Delta$  ( $Y_0$  and  $Y_1$  are perfectly dependent).
- One can make other assumptions about the dependence across ranks from perfect positive or negative ranking to independence.
- In general, the joint distribution of  $(Y_0, Y_1)$  or of  $(Y_1 - Y_0)$  is not identified unless the analyst can pin down the dependence across  $(Y_0, Y_1)$ .
- Thus, even with data from a randomized trial one cannot, without further assumptions, identify the proportion of people who benefit from treatment in the sense of gross gain ( $\Pr(Y_1 \geq Y_0)$ ).

- However, this approach assumes that  $(Y_0, Y_1)$  have the same distribution up to a parameter  $\Delta$  ( $Y_0$  and  $Y_1$  are perfectly dependent).
- One can make other assumptions about the dependence across ranks from perfect positive or negative ranking to independence.
- In general, the joint distribution of  $(Y_0, Y_1)$  or of  $(Y_1 - Y_0)$  is not identified unless the analyst can pin down the dependence across  $(Y_0, Y_1)$ .
- Thus, even with data from a randomized trial one cannot, without further assumptions, identify the proportion of people who benefit from treatment in the sense of gross gain ( $\Pr(Y_1 \geq Y_0)$ ).
- This problem plagues all evaluation methods.

- Abbring and Heckman discuss methods for identifying joint distributions of outcomes in Part III.

- Abbring and Heckman discuss methods for identifying joint distributions of outcomes in Part III.
- Assumption (R-1) is very strong.

- Abbring and Heckman discuss methods for identifying joint distributions of outcomes in Part III.
- Assumption (R-1) is very strong.
- In many cases, it is thought that there is *selection bias* with respect to  $Y_0$ ,  $Y_1$ , so persons who select into status 1 or 0 are selectively different from randomly sampled persons in the population.

- The assumption most commonly made to circumvent problems with (R-1) is that even though  $D$  is not random with respect to potential outcomes, the analyst has access to control variables  $X$  that effectively produce a randomization of  $D$  with respect to  $(Y_0, Y_1)$  given  $X$ .

(M-1)

$(Y_0, Y_1) \perp\!\!\!\perp D \mid X.$

- The assumption most commonly made to circumvent problems with (R-1) is that even though  $D$  is not random with respect to potential outcomes, the analyst has access to control variables  $X$  that effectively produce a randomization of  $D$  with respect to  $(Y_0, Y_1)$  given  $X$ .
- This is the method of matching, which is based on the following conditional independence assumption:

(M-2)

$$(Y_0, Y_1) \perp\!\!\!\perp D \mid X.$$

- Conditioning on  $X$  randomizes  $D$  with respect to  $(Y_0, Y_1)$ .



- Conditioning on  $X$  randomizes  $D$  with respect to  $(Y_0, Y_1)$ .
- (M-1) assumes that any selective sampling of  $(Y_0, Y_1)$  can be adjusted by conditioning on observed variables.

- Conditioning on  $X$  randomizes  $D$  with respect to  $(Y_0, Y_1)$ .
- (M-1) assumes that any selective sampling of  $(Y_0, Y_1)$  can be adjusted by conditioning on observed variables.
- (R-1) and (M-1) are different assumptions and neither implies the other.

- Conditioning on  $X$  randomizes  $D$  with respect to  $(Y_0, Y_1)$ .
- (M-1) assumes that any selective sampling of  $(Y_0, Y_1)$  can be adjusted by conditioning on observed variables.
- (R-1) and (M-1) are different assumptions and neither implies the other.
- In a linear equations model, assumption (M-1) that  $D$  is independent from  $(Y_0, Y_1)$  given  $X$  justifies application of least squares on  $D$  to eliminate selection bias in mean outcome parameters.

- Conditioning on  $X$  randomizes  $D$  with respect to  $(Y_0, Y_1)$ .
- (M-1) assumes that any selective sampling of  $(Y_0, Y_1)$  can be adjusted by conditioning on observed variables.
- (R-1) and (M-1) are different assumptions and neither implies the other.
- In a linear equations model, assumption (M-1) that  $D$  is independent from  $(Y_0, Y_1)$  given  $X$  justifies application of least squares on  $D$  to eliminate selection bias in mean outcome parameters.
- For means, matching is just nonparametric regression.

In order to be able to compare  $X$ -comparable people, we must assume

(M-3)

$$0 < \Pr(D = 1 \mid X = x) < 1.$$

- Assumptions (M-1) and (M-2) justify matching.

- Assumptions (M-1) and (M-2) justify matching.
- Assumption (M-2) is required for *any* evaluation estimator that compares treated and untreated persons.

- Assumptions (M-1) and (M-2) justify matching.
- Assumption (M-2) is required for *any* evaluation estimator that compares treated and untreated persons.
- It is produced by random assignment if the randomization is conducted for all  $X = x$  and there is full compliance.



- Observe that from (M-1) and (M-2), it is possible to identify  $F_1(Y_1 | X = x)$  from the observed data  $F_1(Y_1 | D = 1, X = x)$  since we observe the left hand side of

$$\begin{aligned} F_1(Y_1 | D = 1, X = x) &= F_1(Y_1 | X = x) \\ &= F_1(Y_1 | D = 0, X = x). \end{aligned}$$

- Observe that from (M-1) and (M-2), it is possible to identify  $F_1(Y_1 | X = x)$  from the observed data  $F_1(Y_1 | D = 1, X = x)$  since we observe the left hand side of

$$\begin{aligned} F_1(Y_1 | D = 1, X = x) &= F_1(Y_1 | X = x) \\ &= F_1(Y_1 | D = 0, X = x). \end{aligned}$$

- The first equality is a consequence of conditional independence assumption (M-1).

- Observe that from (M-1) and (M-2), it is possible to identify  $F_1(Y_1 | X = x)$  from the observed data  $F_1(Y_1 | D = 1, X = x)$  since we observe the left hand side of

$$\begin{aligned} F_1(Y_1 | D = 1, X = x) &= F_1(Y_1 | X = x) \\ &= F_1(Y_1 | D = 0, X = x). \end{aligned}$$

- The first equality is a consequence of conditional independence assumption (M-1).
- The second equality comes from (M-1) and (M-2).

- By a similar argument, we observe the left hand side of

$$\begin{aligned} F_0(Y_0 \mid D = 0, X = x) &= F_0(Y_0 \mid X = x) \\ &= F_0(Y_0 \mid D = 1, X = x), \end{aligned}$$

and the equalities are a consequence of (M-1) and (M-2).

- By a similar argument, we observe the left hand side of

$$\begin{aligned} F_0(Y_0 \mid D = 0, X = x) &= F_0(Y_0 \mid X = x) \\ &= F_0(Y_0 \mid D = 1, X = x), \end{aligned}$$

and the equalities are a consequence of (M-1) and (M-2).

- Since the pair of outcomes  $(Y_0, Y_1)$  is not identified for anyone, as in the case of data from randomized trials, the joint distributions of  $(Y_0, Y_1)$  given  $X$  or of  $Y_1 - Y_0$  given  $X$  are not identified without further information.

- From the data on  $Y_1$  given  $X$  and  $D = 1$  and the data on  $Y_0$  given  $X$  and  $D = 0$ , since

$$E(Y_1 | D = 1, X = x) = E(Y_1 | X = x) = E(Y_1 | D = 0, X = x)$$

and

$$E(Y_0 | D = 0, X = x) = E(Y_0 | X = x) = E(Y_0 | D = 1, X = x),$$

we obtain

$$\begin{aligned} E(Y_1 - Y_0 | X = x) &= E(Y_1 - Y_0 | D = 1, X = x) \\ &= E(Y_1 - Y_0 | D = 0, X = x). \end{aligned}$$

- From the data on  $Y_1$  given  $X$  and  $D = 1$  and the data on  $Y_0$  given  $X$  and  $D = 0$ , since

$$E(Y_1 | D = 1, X = x) = E(Y_1 | X = x) = E(Y_1 | D = 0, X = x)$$

and

$$E(Y_0 | D = 0, X = x) = E(Y_0 | X = x) = E(Y_0 | D = 1, X = x),$$

we obtain

$$\begin{aligned} E(Y_1 - Y_0 | X = x) &= E(Y_1 - Y_0 | D = 1, X = x) \\ &= E(Y_1 - Y_0 | D = 0, X = x). \end{aligned}$$

- Effectively, we have a randomization for the subset of the support of  $X$  satisfying (M-2).

- At values of  $X$  that fail to satisfy (M-2), there is no variation in  $D$  given  $X$ .



- At values of  $X$  that fail to satisfy (M-2), there is no variation in  $D$  given  $X$ .
- We can define the residual variation in  $D$  not accounted for by  $X$  as

$$\mathcal{E}(x) = D - E(D | X = x) = D - \Pr(D = 1 | X = x).$$

- At values of  $X$  that fail to satisfy (M-2), there is no variation in  $D$  given  $X$ .
- We can define the residual variation in  $D$  not accounted for by  $X$  as

$$\mathcal{E}(x) = D - E(D | X = x) = D - \Pr(D = 1 | X = x).$$

- If the variance of  $\mathcal{E}(x)$  is zero, it is not possible to construct contrasts in outcomes by treatment status for those  $X$  values and (M-2) is violated.

- At values of  $X$  that fail to satisfy (M-2), there is no variation in  $D$  given  $X$ .
- We can define the residual variation in  $D$  not accounted for by  $X$  as

$$\mathcal{E}(x) = D - E(D | X = x) = D - \Pr(D = 1 | X = x).$$

- If the variance of  $\mathcal{E}(x)$  is zero, it is not possible to construct contrasts in outcomes by treatment status for those  $X$  values and (M-2) is violated.
- To see the consequences of this violation in a regression setting, use  $Y = Y_0 + D(Y_1 - Y_0)$  and take conditional expectations, under (M-1), to obtain

$$E(Y | X, D) = E(Y_0 | X) + D[E(Y_1 - Y_0 | X)].$$

- If  $\text{Var}(\mathcal{E}(x)) > 0$  for all  $x$  in the support of  $X$ , one can use nonparametric least squares to identify  $E(Y_1 - Y_0 | X = x) = \text{ATE}(x)$  by regressing  $Y$  on  $D$  and  $X$ .

- If  $\text{Var}(\mathcal{E}(x)) > 0$  for all  $x$  in the support of  $X$ , one can use nonparametric least squares to identify  $E(Y_1 - Y_0 \mid X = x) = \text{ATE}(x)$  by regressing  $Y$  on  $D$  and  $X$ .
- The function identified from the coefficient on  $D$  is the average treatment effect.

- If  $\text{Var}(\mathcal{E}(x)) > 0$  for all  $x$  in the support of  $X$ , one can use nonparametric least squares to identify  $E(Y_1 - Y_0 | X = x) = \text{ATE}(x)$  by regressing  $Y$  on  $D$  and  $X$ .
- The function identified from the coefficient on  $D$  is the average treatment effect.
- If  $\text{Var}(\mathcal{E}(x)) = 0$ ,  $\text{ATE}(x)$  is not identified at that  $x$  value because there is no variation in  $D$  that is not fully explained by  $X$ .

- If  $\text{Var}(\mathcal{E}(x)) > 0$  for all  $x$  in the support of  $X$ , one can use nonparametric least squares to identify  $E(Y_1 - Y_0 | X = x) = \text{ATE}(x)$  by regressing  $Y$  on  $D$  and  $X$ .
- The function identified from the coefficient on  $D$  is the average treatment effect.
- If  $\text{Var}(\mathcal{E}(x)) = 0$ ,  $\text{ATE}(x)$  is not identified at that  $x$  value because there is no variation in  $D$  that is not fully explained by  $X$ .
- A special case of matching is linear least squares where we write

$$Y_0 = X\alpha + U \qquad Y_1 = X\alpha + \beta + U,$$

$U_0 = U_1 = U$  and hence under (M-1),

$$E(Y | X, D) = X\alpha + D\beta + E(U | X).$$

- If  $D$  is perfectly predictable by  $X$ , we cannot identify  $\beta$  because of a multicollinearity problem.



- If  $D$  is perfectly predictable by  $X$ , we cannot identify  $\beta$  because of a multicollinearity problem.
- (M-2) rules out perfect collinearity.

- If  $D$  is perfectly predictable by  $X$ , we cannot identify  $\beta$  because of a multicollinearity problem.
- (M-2) rules out perfect collinearity.
- Matching is a nonparametric version of least squares that does not impose functional form assumptions on outcome equations, and that imposes support condition (M-2).

- If  $D$  is perfectly predictable by  $X$ , we cannot identify  $\beta$  because of a multicollinearity problem.
- (M-2) rules out perfect collinearity.
- Matching is a nonparametric version of least squares that does not impose functional form assumptions on outcome equations, and that imposes support condition (M-2).
- However, matching does not assume exogeneity of  $X$ .

- Conventional econometric choice models make a distinction between variables that appear in outcome equations ( $X$ ) and variables that appear in choice equations ( $Z$ ).

- Conventional econometric choice models make a distinction between variables that appear in outcome equations ( $X$ ) and variables that appear in choice equations ( $Z$ ).
- The same variables may be in ( $X$ ) and ( $Z$ ), but more typically there are some variables not in common.

- Conventional econometric choice models make a distinction between variables that appear in outcome equations ( $X$ ) and variables that appear in choice equations ( $Z$ ).
- The same variables may be in ( $X$ ) and ( $Z$ ), but more typically there are some variables not in common.
- For example, the instrumental variable estimator is based on variables that are not in  $X$  but that are in  $Z$ .

- Conventional econometric choice models make a distinction between variables that appear in outcome equations ( $X$ ) and variables that appear in choice equations ( $Z$ ).
- The same variables may be in ( $X$ ) and ( $Z$ ), but more typically there are some variables not in common.
- For example, the instrumental variable estimator is based on variables that are not in  $X$  but that are in  $Z$ .
- Matching makes no distinction between the  $X$  and the  $Z$ .

- Conventional econometric choice models make a distinction between variables that appear in outcome equations ( $X$ ) and variables that appear in choice equations ( $Z$ ).
- The same variables may be in ( $X$ ) and ( $Z$ ), but more typically there are some variables not in common.
- For example, the instrumental variable estimator is based on variables that are not in  $X$  but that are in  $Z$ .
- Matching makes no distinction between the  $X$  and the  $Z$ .
- It does not rely on exclusion restrictions.



- Conventional econometric choice models make a distinction between variables that appear in outcome equations ( $X$ ) and variables that appear in choice equations ( $Z$ ).
- The same variables may be in ( $X$ ) and ( $Z$ ), but more typically there are some variables not in common.
- For example, the instrumental variable estimator is based on variables that are not in  $X$  but that are in  $Z$ .
- Matching makes no distinction between the  $X$  and the  $Z$ .
- It does not rely on exclusion restrictions.
- The conditioning variables used to achieve conditional independence can in principle be a set of variables  $Q$  distinct from the  $X$  variables (covariates for outcomes) or the  $Z$  variables (covariates for choices).

- We use  $X$  solely to simplify the notation.

- We use  $X$  solely to simplify the notation.
- The key identifying assumption is the assumed existence of a random variable  $X$  with the properties satisfying (M-1) and (M-2).

- Conditioning on a larger vector ( $X$  augmented with additional variables) or a smaller vector ( $X$  with some components removed) may or may not produce suitably modified versions of (M-1) and (M-2).

- Conditioning on a larger vector ( $X$  augmented with additional variables) or a smaller vector ( $X$  with some components removed) may or may not produce suitably modified versions of (M-1) and (M-2).
- Without invoking further assumptions, there is no objective principle for determining what conditioning variables produce (M-1).

- Assumption (M-1) is strong.

- Assumption (M-1) is strong.
- Many economists do not have enough faith in their data to invoke it.

- Assumption (M-1) is strong.
- Many economists do not have enough faith in their data to invoke it.
- Assumption (M-2) is testable and requires no act of faith.



- Assumption (M-1) is strong.
- Many economists do not have enough faith in their data to invoke it.
- Assumption (M-2) is testable and requires no act of faith.
- To justify (M-1), it is necessary to appeal to the quality of the data.

- Using economic theory can help guide the choice of an evaluation estimator.

- Using economic theory can help guide the choice of an evaluation estimator.
- A crucial distinction is the one between the information available to the analyst and the information available to the agent whose outcomes are being studied.

- Using economic theory can help guide the choice of an evaluation estimator.
- A crucial distinction is the one between the information available to the analyst and the information available to the agent whose outcomes are being studied.
- Assumptions made about these information sets drive the properties of econometric estimators.

- Using economic theory can help guide the choice of an evaluation estimator.
- A crucial distinction is the one between the information available to the analyst and the information available to the agent whose outcomes are being studied.
- Assumptions made about these information sets drive the properties of econometric estimators.
- Analysts using matching make strong informational assumptions in terms of the data available to them.

- Using economic theory can help guide the choice of an evaluation estimator.
- A crucial distinction is the one between the information available to the analyst and the information available to the agent whose outcomes are being studied.
- Assumptions made about these information sets drive the properties of econometric estimators.
- Analysts using matching make strong informational assumptions in terms of the data available to them.
- In fact, all econometric estimators make assumptions about the presence or absence of informational asymmetries, and we exposit them in this chapter.

- To analyze the informational assumptions invoked in matching, and other econometric evaluation strategies, it is helpful to introduce five distinct information sets and establish some relationships among them.

- To analyze the informational assumptions invoked in matching, and other econometric evaluation strategies, it is helpful to introduce five distinct information sets and establish some relationships among them.
  - ① An information set  $\sigma(I_{R^*})$  with an associated random variable that satisfies conditional independence (M-1) is defined as a *relevant* information set;



- To analyze the informational assumptions invoked in matching, and other econometric evaluation strategies, it is helpful to introduce five distinct information sets and establish some relationships among them.
  - 1 An information set  $\sigma(I_{R^*})$  with an associated random variable that satisfies conditional independence (M-1) is defined as a *relevant* information set;
  - 2 The minimal information set  $\sigma(I_R)$  with associated random variable needed to satisfy conditional independence (M-1), the *minimal relevant* information set;

- To analyze the informational assumptions invoked in matching, and other econometric evaluation strategies, it is helpful to introduce five distinct information sets and establish some relationships among them.
  - 1 An information set  $\sigma(I_{R^*})$  with an associated random variable that satisfies conditional independence (M-1) is defined as a *relevant* information set;
  - 2 The minimal information set  $\sigma(I_R)$  with associated random variable needed to satisfy conditional independence (M-1), the *minimal relevant* information set;
  - 3 The information set  $\sigma(I_A)$  available to the agent at the time decisions to participate are made;

- To analyze the informational assumptions invoked in matching, and other econometric evaluation strategies, it is helpful to introduce five distinct information sets and establish some relationships among them.
  - 1 An information set  $\sigma(I_{R^*})$  with an associated random variable that satisfies conditional independence (M-1) is defined as a *relevant* information set;
  - 2 The minimal information set  $\sigma(I_R)$  with associated random variable needed to satisfy conditional independence (M-1), the *minimal relevant* information set;
  - 3 The information set  $\sigma(I_A)$  available to the agent at the time decisions to participate are made;
  - 4 The information available to the economist,  $\sigma(I_{E^*})$ ; and

- To analyze the informational assumptions invoked in matching, and other econometric evaluation strategies, it is helpful to introduce five distinct information sets and establish some relationships among them.
  - 1 An information set  $\sigma(I_{R^*})$  with an associated random variable that satisfies conditional independence (M-1) is defined as a *relevant* information set;
  - 2 The minimal information set  $\sigma(I_R)$  with associated random variable needed to satisfy conditional independence (M-1), the *minimal relevant* information set;
  - 3 The information set  $\sigma(I_A)$  available to the agent at the time decisions to participate are made;
  - 4 The information available to the economist,  $\sigma(I_{E^*})$ ; and
  - 5 The information  $\sigma(I_E)$  used by the economist in conducting an empirical analysis.

- To analyze the informational assumptions invoked in matching, and other econometric evaluation strategies, it is helpful to introduce five distinct information sets and establish some relationships among them.
  - ① An information set  $\sigma(I_{R^*})$  with an associated random variable that satisfies conditional independence (M-1) is defined as a *relevant* information set;
  - ② The minimal information set  $\sigma(I_R)$  with associated random variable needed to satisfy conditional independence (M-1), the *minimal relevant* information set;
  - ③ The information set  $\sigma(I_A)$  available to the agent at the time decisions to participate are made;
  - ④ The information available to the economist,  $\sigma(I_{E^*})$ ; and
  - ⑤ The information  $\sigma(I_E)$  used by the economist in conducting an empirical analysis.
- We will denote the random variables generated by these sets as  $I_{R^*}$ ,  $I_R$ ,  $I_A$ ,  $I_{E^*}$ , and  $I_E$ , respectively.

## Definition 1

We say that  $\sigma(I_{R^*})$  is a **relevant information set** if the information set is generated by the random variable  $I_{R^*}$ , possibly vector valued, and satisfies condition (M-1), so that

$$(Y_0, Y_1) \perp\!\!\!\perp D \mid I_{R^*}.$$

## Definition 2

We say that  $\sigma(I_R)$  is a **minimal relevant information set** if it is the intersection of all sets  $\sigma(I_{R^*})$  and satisfies  $(Y_0, Y_1) \perp\!\!\!\perp D \mid I_R$ . The associated random variable  $I_R$  is a minimum amount of information that guarantees that condition (M-1) is satisfied. There may be no such set.

- If we define a relevant information set as one that produces conditional independence, it may not be unique.



- If we define a relevant information set as one that produces conditional independence, it may not be unique.
- If the set  $\sigma(I_{R^*})$  satisfies the conditional independence condition, then the set  $\sigma(I_{R^*}, Q)$  such that  $Q \perp\!\!\!\perp (Y_0, Y_1) \mid I_{R^*}$  would also guarantee conditional independence.

- If we define a relevant information set as one that produces conditional independence, it may not be unique.
- If the set  $\sigma(I_{R^*})$  satisfies the conditional independence condition, then the set  $\sigma(I_{R^*}, Q)$  such that  $Q \perp\!\!\!\perp (Y_0, Y_1) \mid I_{R^*}$  would also guarantee conditional independence.
- For this reason, when possible, it is desirable to use the minimal relevant information set.

### Definition 3

The agent's information set,  $\sigma(I_A)$ , is defined by the information  $I_A$  used by the agent when choosing among treatments. Accordingly, we call  $I_A$  the **agent's information**.

- By the agent we mean the person making the treatment decision, not necessarily the person whose outcomes are being studied (e.g., the agent may be the parent; the person being studied may be a child).

## Definition 4

The econometrician's **full information set**,  $\sigma(I_{E^*})$ , is defined as **all** of the information available to the econometrician,  $I_{E^*}$ .

## Definition 5

The **econometrician's information set**,  $\sigma(I_E)$ , is defined by the information **used** by the econometrician when analyzing the agent's choice of treatment,  $I_E$ , in conducting an analysis.

- For the case where a unique minimal relevant information set exists, only three restrictions are implied by the structure of these sets:  $\sigma(I_R) \subseteq \sigma(I_{R^*})$ ,  $\sigma(I_R) \subseteq \sigma(I_A)$ , and  $\sigma(I_E) \subseteq \sigma(I_{E^*})$ .

- For the case where a unique minimal relevant information set exists, only three restrictions are implied by the structure of these sets:  $\sigma(I_R) \subseteq \sigma(I_{R^*})$ ,  $\sigma(I_R) \subseteq \sigma(I_A)$ , and  $\sigma(I_E) \subseteq \sigma(I_{E^*})$ .
- We have already discussed the first restriction.

- For the case where a unique minimal relevant information set exists, only three restrictions are implied by the structure of these sets:  $\sigma(I_R) \subseteq \sigma(I_{R^*})$ ,  $\sigma(I_R) \subseteq \sigma(I_A)$ , and  $\sigma(I_E) \subseteq \sigma(I_{E^*})$ .
- We have already discussed the first restriction.
- The second restriction requires that the minimal relevant information set must be part of the information the agent uses when deciding which treatment to take or assign.

- For the case where a unique minimal relevant information set exists, only three restrictions are implied by the structure of these sets:  $\sigma(I_R) \subseteq \sigma(I_{R^*})$ ,  $\sigma(I_R) \subseteq \sigma(I_A)$ , and  $\sigma(I_E) \subseteq \sigma(I_{E^*})$ .
- We have already discussed the first restriction.
- The second restriction requires that the minimal relevant information set must be part of the information the agent uses when deciding which treatment to take or assign.
- It is the information in  $\sigma(I_A)$  that gives rise to the selection problem.



- The third restriction requires that the information used by the econometrician must be part of the information that the econometrician observes.

- The third restriction requires that the information used by the econometrician must be part of the information that the econometrician observes.
- Aside from these orderings, the econometrician's information set may be different from the agent's or the relevant information set.

- The third restriction requires that the information used by the econometrician must be part of the information that the econometrician observes.
- Aside from these orderings, the econometrician's information set may be different from the agent's or the relevant information set.
- The econometrician may know something the agent doesn't know, for typically he is observing events after the decision is made.

- The third restriction requires that the information used by the econometrician must be part of the information that the econometrician observes.
- Aside from these orderings, the econometrician's information set may be different from the agent's or the relevant information set.
- The econometrician may know something the agent doesn't know, for typically he is observing events after the decision is made.
- At the same time, there may be private information known to the agent but not the econometrician.

- Assuming a minimal relevant information set exists, matching assumption (M-1) implies that  $\sigma(I_R) \subseteq \sigma(I_E)$ , so that the econometrician uses at least the minimal relevant information set, but of course he or she may use more.

- Assuming a minimal relevant information set exists, matching assumption (M-1) implies that  $\sigma(I_R) \subseteq \sigma(I_E)$ , so that the econometrician uses at least the minimal relevant information set, but of course he or she may use more.
- However, using more information is not guaranteed to produce a model with conditional independence property (M-1) satisfied for the augmented model.

- Assuming a minimal relevant information set exists, matching assumption (M-1) implies that  $\sigma(I_R) \subseteq \sigma(I_E)$ , so that the econometrician uses at least the minimal relevant information set, but of course he or she may use more.
- However, using more information is not guaranteed to produce a model with conditional independence property (M-1) satisfied for the augmented model.
- Thus an analyst can “overdo” it.

- Assuming a minimal relevant information set exists, matching assumption (M-1) implies that  $\sigma(I_R) \subseteq \sigma(I_E)$ , so that the econometrician uses at least the minimal relevant information set, but of course he or she may use more.
- However, using more information is not guaranteed to produce a model with conditional independence property (M-1) satisfied for the augmented model.
- Thus an analyst can “overdo” it.
- We present examples of the consequences of the asymmetry in agent and analyst information sets in Slide 675.



- The possibility of asymmetry in information between the agent making participation decisions and the observing economist creates the potential for a major identification problem that is ruled out by assumption (M-1).

- The possibility of asymmetry in information between the agent making participation decisions and the observing economist creates the potential for a major identification problem that is ruled out by assumption (M-1).
- The methods of control functions and instrumental variables estimators (and closely related regression discontinuity design methods) address this problem in different ways.

- The possibility of asymmetry in information between the agent making participation decisions and the observing economist creates the potential for a major identification problem that is ruled out by assumption (M-1).
- The methods of control functions and instrumental variables estimators (and closely related regression discontinuity design methods) address this problem in different ways.
- Accounting for this possibility is a more conservative approach to the selection problem than the one taken by advocates of matching.

- The possibility of asymmetry in information between the agent making participation decisions and the observing economist creates the potential for a major identification problem that is ruled out by assumption (M-1).
- The methods of control functions and instrumental variables estimators (and closely related regression discontinuity design methods) address this problem in different ways.
- Accounting for this possibility is a more conservative approach to the selection problem than the one taken by advocates of matching.
- Those advocates assume that they know the  $X$  that produces a relevant information set.

- ? show the biases that can result in matching when standard econometric model selection criteria are applied to pick the  $X$  that are used to satisfy (M-1) and we summarize their analysis in Slide 675.

- ? show the biases that can result in matching when standard econometric model selection criteria are applied to pick the  $X$  that are used to satisfy (M-1) and we summarize their analysis in Slide 675.
- Conditional independence condition (M-1) cannot be tested without maintaining other assumptions.

- ? show the biases that can result in matching when standard econometric model selection criteria are applied to pick the  $X$  that are used to satisfy (M-1) and we summarize their analysis in Slide 675.
- Conditional independence condition (M-1) cannot be tested without maintaining other assumptions.
- As noted in Part I, choosing the appropriate conditioning variables is a problem that plagues *all* econometric estimators.

- The methods of control functions, replacement functions, proxy variables and instrumental variables recognize the possibility of asymmetry in information between the agent being studied and the econometrician and further recognize that even after conditioning on  $X$  (variables in the outcome equation) and  $Z$  (variables affecting treatment choices, which may include the  $X$ ), analysts may fail to satisfy conditional independence condition (M-1).



- These methods postulate the existence of some unobservables  $\theta$ , which may be vector valued, with the property that

(U-1)

$$(Y_0, Y_1) \perp\!\!\!\perp D \mid X, Z, \theta,$$

but allow for the possibility that

(U-2)

$$(Y_0, Y_1) \not\perp\!\!\!\perp D \mid X, Z.$$

- In the event (U-2) holds, these approaches model the relationship of the unobservable  $\theta$  with  $(Y_0, Y_1)$  and  $D$  in various ways.

- In the event (U-2) holds, these approaches model the relationship of the unobservable  $\theta$  with  $(Y_0, Y_1)$  and  $D$  in various ways.
- The content in the control function principle is to specify the exact nature of the dependence on the relationship between observables and unobservables in a nontrivial fashion that is consistent with economic theory.

- In the event (U-2) holds, these approaches model the relationship of the unobservable  $\theta$  with  $(Y_0, Y_1)$  and  $D$  in various ways.
- The content in the control function principle is to specify the exact nature of the dependence on the relationship between observables and unobservables in a nontrivial fashion that is consistent with economic theory.
- We present examples of models that satisfy (U-1) but not (U-2) in Slide 675.

- The early literature focused on mean outcomes conditional on covariates (????) and assumes a weaker version of (U-1) based on conditional mean independence rather than full conditional independence.

- The early literature focused on mean outcomes conditional on covariates (????) and assumes a weaker version of (U-1) based on conditional mean independence rather than full conditional independence.
- More recent work analyzes distributions of outcomes (e.g., ??).

- The early literature focused on mean outcomes conditional on covariates (????) and assumes a weaker version of (U-1) based on conditional mean independence rather than full conditional independence.
- More recent work analyzes distributions of outcomes (e.g., ??).
- Abbring and Heckman review this work in Part III.

- The normal Roy model discussed in Part I makes distributional assumptions and identifies the joint distribution of outcomes.



- The normal Roy model discussed in Part I makes distributional assumptions and identifies the joint distribution of outcomes.
- (Recall the discussion in section 6.1 of Part I.) A large literature surveyed in ? makes alternative assumptions to satisfy (U-1) in nonparametric settings.

- The normal Roy model discussed in Part I makes distributional assumptions and identifies the joint distribution of outcomes.
- (Recall the discussion in section 6.1 of Part I.) A large literature surveyed in ? makes alternative assumptions to satisfy (U-1) in nonparametric settings.
- Replacement functions (?) are methods that proxy  $\theta$ .

- The normal Roy model discussed in Part I makes distributional assumptions and identifies the joint distribution of outcomes.
- (Recall the discussion in section 6.1 of Part I.) A large literature surveyed in ? makes alternative assumptions to satisfy (U-1) in nonparametric settings.
- Replacement functions (?) are methods that proxy  $\theta$ .
- They substitute out for  $\theta$  using observables.

- ??, ??, ?, and ?? develop methods that integrate out  $\theta$  from the model assuming  $\theta \perp\!\!\!\perp (X, Z)$ , or invoking weaker mean independence assumptions, and assuming access to proxy measurements for  $\theta$ .

- ??, ??, ?, and ?? develop methods that integrate out  $\theta$  from the model assuming  $\theta \perp\!\!\!\perp (X, Z)$ , or invoking weaker mean independence assumptions, and assuming access to proxy measurements for  $\theta$ .
- They also consider methods for estimating the distributions of treatment effects.

- ??, ??, ?, and ?? develop methods that integrate out  $\theta$  from the model assuming  $\theta \perp\!\!\!\perp (X, Z)$ , or invoking weaker mean independence assumptions, and assuming access to proxy measurements for  $\theta$ .
- They also consider methods for estimating the distributions of treatment effects.
- These methods are discussed in Part III.

- The normal selection model discussed in section 6.1 of Part I produces partial identification of a generalized Roy model and full identification of a Roy model under separability and normality.

- The normal selection model discussed in section 6.1 of Part I produces partial identification of a generalized Roy model and full identification of a Roy model under separability and normality.
- It models the conditional expectation of  $U_0$  and  $U_1$  given  $X, Z$ , and  $D$ .



- The normal selection model discussed in section 6.1 of Part I produces partial identification of a generalized Roy model and full identification of a Roy model under separability and normality.
- It models the conditional expectation of  $U_0$  and  $U_1$  given  $X, Z$ , and  $D$ .
- In terms of (U-1), it models the conditional mean dependence of  $Y_0, Y_1$  on  $D$  and  $\theta$  given  $X$  and  $Z$ .

- The normal selection model discussed in section 6.1 of Part I produces partial identification of a generalized Roy model and full identification of a Roy model under separability and normality.
- It models the conditional expectation of  $U_0$  and  $U_1$  given  $X, Z$ , and  $D$ .
- In terms of (U-1), it models the conditional mean dependence of  $Y_0, Y_1$  on  $D$  and  $\theta$  given  $X$  and  $Z$ .
- ? and ? surveys methods for identifying semiparametric versions of these models.

- The normal selection model discussed in section 6.1 of Part I produces partial identification of a generalized Roy model and full identification of a Roy model under separability and normality.
- It models the conditional expectation of  $U_0$  and  $U_1$  given  $X, Z$ , and  $D$ .
- In terms of (U-1), it models the conditional mean dependence of  $Y_0, Y_1$  on  $D$  and  $\theta$  given  $X$  and  $Z$ .
- ? and ? surveys methods for identifying semiparametric versions of these models.
- Appendix B of Part I presents a prototypical identification proof for a general selection model that implements (U-1) by estimating the distribution of  $\theta$ , assuming  $\theta \perp\!\!\!\perp (X, Z)$ , and invoking support conditions on  $(X, Z)$ .

- Central to both the selection approach and the instrumental variable approach for a model with heterogenous responses is the probability of selection.

- Central to both the selection approach and the instrumental variable approach for a model with heterogenous responses is the probability of selection.
- Let  $Z$  denote variables in the choice equation.

- Central to both the selection approach and the instrumental variable approach for a model with heterogenous responses is the probability of selection.
- Let  $Z$  denote variables in the choice equation.
- Fixing  $Z$  at different values (denoted  $z$ ), we define  $D(z)$  as an indicator function that is “1” when treatment is selected at the fixed value of  $z$  and that is “0” otherwise.

- Central to both the selection approach and the instrumental variable approach for a model with heterogenous responses is the probability of selection.
- Let  $Z$  denote variables in the choice equation.
- Fixing  $Z$  at different values (denoted  $z$ ), we define  $D(z)$  as an indicator function that is “1” when treatment is selected at the fixed value of  $z$  and that is “0” otherwise.
- In terms of the separable index model introduced in Part I, for a fixed value of  $z$ ,

$$D(z) = \mathbf{1}(\mu_D(z) \geq V)$$

where  $Z \perp\!\!\!\perp V \mid X$ .

- Thus fixing  $Z = z$ , values of  $z$  do not affect the realizations of  $V$  for any value of  $X$ .



- Thus fixing  $Z = z$ , values of  $z$  do not affect the realizations of  $V$  for any value of  $X$ .
- An alternative way of representing the independence between  $Z$  and  $V$  given  $X$ , due to ?, writes that  $D(z) \perp\!\!\!\perp Z \mid X$  for all  $z \in \mathcal{Z}$ , where  $\mathcal{Z}$  is the support of  $Z$ .

- Thus fixing  $Z = z$ , values of  $z$  do not affect the realizations of  $V$  for any value of  $X$ .
- An alternative way of representing the independence between  $Z$  and  $V$  given  $X$ , due to ?, writes that  $D(z) \perp\!\!\!\perp Z \mid X$  for all  $z \in \mathcal{Z}$ , where  $\mathcal{Z}$  is the support of  $Z$ .
- The Imbens-Angrist independence condition for IV is

$$\{D(z)\}_{z \in \mathcal{Z}} \perp\!\!\!\perp Z \mid X.$$

- Thus fixing  $Z = z$ , values of  $z$  do not affect the realizations of  $V$  for any value of  $X$ .
- An alternative way of representing the independence between  $Z$  and  $V$  given  $X$ , due to ?, writes that  $D(z) \perp\!\!\!\perp Z \mid X$  for all  $z \in \mathcal{Z}$ , where  $\mathcal{Z}$  is the support of  $Z$ .
- The Imbens-Angrist independence condition for IV is

$$\{D(z)\}_{z \in \mathcal{Z}} \perp\!\!\!\perp Z \mid X.$$

- Thus the probabilities that  $D(z) = 1$ ,  $z \in \mathcal{Z}$  are independent of  $Z$ .

The method of instrumental variables (IV) postulates that

(IV-1)

$(Y_0, Y_1, \{D(z)\}_{z \in \mathcal{Z}}) \perp\!\!\!\perp Z \mid X. (\textit{Independence})$

- One consequence of this assumption is that  $E(D | Z) = P(Z)$ , the propensity score, is random with respect to potential outcomes.

- One consequence of this assumption is that  $E(D | Z) = P(Z)$ , the propensity score, is random with respect to potential outcomes.
- Thus  $(Y_0, Y_1) \perp\!\!\!\perp P(Z) | X$ .

- One consequence of this assumption is that  $E(D | Z) = P(Z)$ , the propensity score, is random with respect to potential outcomes.
- Thus  $(Y_0, Y_1) \perp\!\!\!\perp P(Z) | X$ .
- So are all other functions of  $Z$  given  $X$ .

The method of instrumental variables also assumes that

(IV-2)

$E(D | X, Z) = P(X, Z)$  is a nondegenerate function of  $Z$  given  $X$ .  
(Rank Condition)

- Alternatively, we can write that  
 $\text{Var}(E(D | X, Z)) \neq \text{Var}(E(D | X))$ .



- Comparing (IV-1) to (M-1), in the method of instrumental variables,  $Z$  is independent of  $(Y_0, Y_1)$  given  $X$  whereas in matching,  $D$  is independent of  $(Y_0, Y_1)$  given  $X$ .

- Comparing (IV-1) to (M-1), in the method of instrumental variables,  $Z$  is independent of  $(Y_0, Y_1)$  given  $X$  whereas in matching,  $D$  is independent of  $(Y_0, Y_1)$  given  $X$ .
- So in (IV-1),  $Z$  plays the role of  $D$  in matching condition (M-1).

- Comparing (IV-1) to (M-1), in the method of instrumental variables,  $Z$  is independent of  $(Y_0, Y_1)$  given  $X$  whereas in matching,  $D$  is independent of  $(Y_0, Y_1)$  given  $X$ .
- So in (IV-1),  $Z$  plays the role of  $D$  in matching condition (M-1).
- Comparing (IV-2) with (M-2), in the method of IV, the choice probability  $\Pr(D = 1 | X, Z)$  is assumed to vary conditional on  $X$  whereas in matching,  $D$  varies conditional on  $X$ .

- Comparing (IV-1) to (M-1), in the method of instrumental variables,  $Z$  is independent of  $(Y_0, Y_1)$  given  $X$  whereas in matching,  $D$  is independent of  $(Y_0, Y_1)$  given  $X$ .
- So in (IV-1),  $Z$  plays the role of  $D$  in matching condition (M-1).
- Comparing (IV-2) with (M-2), in the method of IV, the choice probability  $\Pr(D = 1 | X, Z)$  is assumed to vary conditional on  $X$  whereas in matching,  $D$  varies conditional on  $X$ .
- Unlike the method of control functions, no explicit model of the relationship between  $D$  and  $(Y_0, Y_1)$  is required in applying IV.

- Comparing (IV-1) to (M-1), in the method of instrumental variables,  $Z$  is independent of  $(Y_0, Y_1)$  given  $X$  whereas in matching,  $D$  is independent of  $(Y_0, Y_1)$  given  $X$ .
- So in (IV-1),  $Z$  plays the role of  $D$  in matching condition (M-1).
- Comparing (IV-2) with (M-2), in the method of IV, the choice probability  $\Pr(D = 1 \mid X, Z)$  is assumed to vary conditional on  $X$  whereas in matching,  $D$  varies conditional on  $X$ .
- Unlike the method of control functions, no explicit model of the relationship between  $D$  and  $(Y_0, Y_1)$  is required in applying IV.
- We exposit the implicit model of the relationship between  $D$  and  $(Y_0, Y_1)$  used in instrumental variables in this chapter.

- (IV-2) is a rank condition and can be empirically verified.

- (IV-2) is a rank condition and can be empirically verified.
- (IV-1) is not testable as it involves assumptions about counterfactuals.

- (IV-2) is a rank condition and can be empirically verified.
- (IV-1) is not testable as it involves assumptions about counterfactuals.
- In a conventional common coefficient regression model

$$Y = \alpha + \beta D + U,$$

where  $\beta$  is a constant and where we allow for  $\text{Cov}(D, U) \neq 0$ , (IV-1) and (IV-2) identify  $\beta$ .



- (IV-2) is a rank condition and can be empirically verified.
- (IV-1) is not testable as it involves assumptions about counterfactuals.
- In a conventional common coefficient regression model

$$Y = \alpha + \beta D + U,$$

where  $\beta$  is a constant and where we allow for  $\text{Cov}(D, U) \neq 0$ , (IV-1) and (IV-2) identify  $\beta$ .

- When  $\beta$  varies in the population and is correlated with  $D$ , additional assumptions must be invoked for IV to identify interpretable parameters.

- (IV-2) is a rank condition and can be empirically verified.
- (IV-1) is not testable as it involves assumptions about counterfactuals.
- In a conventional common coefficient regression model

$$Y = \alpha + \beta D + U,$$

where  $\beta$  is a constant and where we allow for  $\text{Cov}(D, U) \neq 0$ , (IV-1) and (IV-2) identify  $\beta$ .

- When  $\beta$  varies in the population and is correlated with  $D$ , additional assumptions must be invoked for IV to identify interpretable parameters.
- We discuss these conditions in Slide 152 of this chapter, drawing on and extending the analysis of ??? and ?.

- Assumptions (IV-1) and (IV-2), with additional assumptions in the case where  $\beta$  varies in the population which we discuss in this chapter, can be used to identify mean treatment parameters.

- Assumptions (IV-1) and (IV-2), with additional assumptions in the case where  $\beta$  varies in the population which we discuss in this chapter, can be used to identify mean treatment parameters.
- Replacing  $Y_1$  with  $\mathbf{1}(Y_1 \leq t)$  and  $Y_0$  with  $\mathbf{1}(Y_0 \leq t)$ , where  $t$  is a constant, the IV approach allows us to identify marginal distributions  $F_1(y_1 | X)$  or  $F_0(y_0 | X)$ .

- In matching, the variation in  $D$  that arises after conditioning on  $X$  provides the source of randomness that switches people across treatment status.

- In matching, the variation in  $D$  that arises after conditioning on  $X$  provides the source of randomness that switches people across treatment status.
- Nature is assumed to provide an experimental manipulation conditional on  $X$  that replaces the randomization assumed in (R-1)-(R-3).

- In matching, the variation in  $D$  that arises after conditioning on  $X$  provides the source of randomness that switches people across treatment status.
- Nature is assumed to provide an experimental manipulation conditional on  $X$  that replaces the randomization assumed in (R-1)-(R-3).
- When  $D$  is perfectly predictable by  $X$ , there is no variation in it conditional on  $X$ , and the randomization by nature breaks down.

- In matching, the variation in  $D$  that arises after conditioning on  $X$  provides the source of randomness that switches people across treatment status.
- Nature is assumed to provide an experimental manipulation conditional on  $X$  that replaces the randomization assumed in (R-1)-(R-3).
- When  $D$  is perfectly predictable by  $X$ , there is no variation in it conditional on  $X$ , and the randomization by nature breaks down.
- Heuristically, matching assumes a residual  $\mathcal{E}(X) = D - E(D | X)$  that is nondegenerate and is one manifestation of the randomness that causes persons to switch status.



- In the IV method, it is the choice probability  $E(D | X, Z) = P(X, Z)$  that is random with respect to  $(Y_0, Y_1)$ , not components of  $D$  not predictable by  $(X, Z)$ .

- In the IV method, it is the choice probability  $E(D | X, Z) = P(X, Z)$  that is random with respect to  $(Y_0, Y_1)$ , not components of  $D$  not predictable by  $(X, Z)$ .
- Variation in  $Z$  for a fixed  $X$  provides the required variation in  $D$  that switches treatment status and still produces the required conditional independence:

$$(Y_0, Y_1) \perp\!\!\!\perp P(X, Z) | X.$$

- Variation in  $P(X, Z)$  produces variations in  $D$  that switch treatment status.

- Variation in  $P(X, Z)$  produces variations in  $D$  that switch treatment status.
- Components of variation in  $D$  not predictable by  $(X, Z)$  do not produce the required independence.

- Variation in  $P(X, Z)$  produces variations in  $D$  that switch treatment status.
- Components of variation in  $D$  not predictable by  $(X, Z)$  do not produce the required independence.
- Instead, the predicted component provides the required independence.

- Variation in  $P(X, Z)$  produces variations in  $D$  that switch treatment status.
- Components of variation in  $D$  not predictable by  $(X, Z)$  do not produce the required independence.
- Instead, the predicted component provides the required independence.
- It is just the opposite in matching.

- Variation in  $P(X, Z)$  produces variations in  $D$  that switch treatment status.
- Components of variation in  $D$  not predictable by  $(X, Z)$  do not produce the required independence.
- Instead, the predicted component provides the required independence.
- It is just the opposite in matching.
- Versions of the method of control functions use measurements to proxy  $\theta$  in (U-1) and (U-2) and remove spurious dependence that gives rise to selection problems.

- Variation in  $P(X, Z)$  produces variations in  $D$  that switch treatment status.
- Components of variation in  $D$  not predictable by  $(X, Z)$  do not produce the required independence.
- Instead, the predicted component provides the required independence.
- It is just the opposite in matching.
- Versions of the method of control functions use measurements to proxy  $\theta$  in (U-1) and (U-2) and remove spurious dependence that gives rise to selection problems.
- These are called replacement functions (see ?) or control variates (see ?).



- Table 1 summarizes some of the main lessons of this section.

- Table 1 summarizes some of the main lessons of this section.
- We stress that the stated conditions are necessary conditions.

- Table 1 summarizes some of the main lessons of this section.
- We stress that the stated conditions are necessary conditions.
- There are many versions of the IV and control functions principle and extensions of these ideas which refine these basic postulates more fully and we exposit them in this Handbook.

- Table 1 summarizes some of the main lessons of this section.
- We stress that the stated conditions are necessary conditions.
- There are many versions of the IV and control functions principle and extensions of these ideas which refine these basic postulates more fully and we exposit them in this Handbook.
- We start with the method of instrumental variables and analyze the general case where responses to treatment are heterogeneous and persons select into treatment status in response to the heterogeneity in treatment response.

# Table 1: Identifying Assumptions Under Commonly Used Methods

	Identifying Assumptions	Identifies marginal distributions?	Exclusion condition needed?
Random Assignment	$(Y_0, Y_1) \perp\!\!\!\perp \xi$ , $\xi = 1 \implies A = 1, \xi = 0 \implies A = 0$ (full compliance). Alternatively, if self-selection is random with respect to outcomes, $(Y_0, Y_1) \perp\!\!\!\perp D$ . Assignment can be conditional on $X$ .	Yes	No
Matching	$(Y_0, Y_1) \not\perp\!\!\!\perp D$ , but $(Y_0, Y_1) \perp\!\!\!\perp D \mid X$ , $0 < \Pr(D = 1 \mid X) < 1$ for all $X$ . So $D$ conditional on $X$ is a nondegenerate random variable.	Yes	No
Control Functions and Extensions	$(Y_0, Y_1) \not\perp\!\!\!\perp D \mid X, Z$ , but $(Y_0, Y_1) \perp\!\!\!\perp D \mid X, Z, \theta$ . The method models dependence induced by $\theta$ or else proxies $\theta$ (replacement function). Version (i). Replacement functions (substitute out $\theta$ by observables) (Blundell and Powell, 2003; Heckman and Robb, 1985; Olley and Pakes, 1996). Factor models (Carneiro, Hansen and Heckman, 2003) allow for measurement error in the proxies. Version (ii). Integrate out $\theta$ assuming $\theta \perp\!\!\!\perp (X, Z)$ (Aakvik, Heckman, and Vytlacil, 2005; Carneiro, Hansen, and Heckman, 2003). Version (iii). For separable models for mean response expect out $\theta$ conditional on $X, Z, D$ as in standard selection models (control functions in the same sense of Heckman and Robb).	Yes	Yes (for semi-parametric models)  No (under some parametric assumptions)
IV	$(Y_0, Y_1) \not\perp\!\!\!\perp D \mid X, Z$ , but $(Y_1, Y_0) \perp\!\!\!\perp Z \mid X$ , $\Pr(D = 1 \mid Z)$ is a nondegenerate function of $Z$ .	Yes	Yes

Notes:  $(Y_0, Y_1)$  are potential outcomes that depend on  $X$ ;

$$D = \begin{cases} 1 & \text{if assigned (or choose) status 1,} \\ 0 & \text{otherwise;} \end{cases}$$

$Z$  are determinants of  $D$ ,  $\theta$  is a vector of unobservables. For random assignments,  $A$  is a vector of actual treatment status.  $A = 1$  if treated;  $A = 0$  if not.  $\xi = 1$  if a person is randomized to treatment status;  $\xi = 0$  otherwise.

- Our strategy in this chapter is to anchor all of our analysis around the economic theory of choice as embodied in discrete choice theory and versions of the generalized Roy model developed in Part I.

- Our strategy in this chapter is to anchor all of our analysis around the economic theory of choice as embodied in discrete choice theory and versions of the generalized Roy model developed in Part I.
- We next show how recent developments allow analysts to define treatment parameters within a well posed economic framework but without the strong assumptions maintained in the early literature on selection models.



- Our strategy in this chapter is to anchor all of our analysis around the economic theory of choice as embodied in discrete choice theory and versions of the generalized Roy model developed in Part I.
- We next show how recent developments allow analysts to define treatment parameters within a well posed economic framework but without the strong assumptions maintained in the early literature on selection models.
- To focus our discussion, we first consider the analysis of a prototypical policy evaluation program.

## A Prototypical Policy Evaluation Problem

- To motivate our discussion in this chapter, consider the following prototypical policy problem.

## A Prototypical Policy Evaluation Problem

- To motivate our discussion in this chapter, consider the following prototypical policy problem.
- Suppose a policy is proposed for adoption in a country.

## A Prototypical Policy Evaluation Problem

- To motivate our discussion in this chapter, consider the following prototypical policy problem.
- Suppose a policy is proposed for adoption in a country.
- It has been tried in other countries and we know outcomes there.

## A Prototypical Policy Evaluation Problem

- To motivate our discussion in this chapter, consider the following prototypical policy problem.
- Suppose a policy is proposed for adoption in a country.
- It has been tried in other countries and we know outcomes there.
- We also know outcomes in countries where it was not adopted.

## A Prototypical Policy Evaluation Problem

- To motivate our discussion in this chapter, consider the following prototypical policy problem.
- Suppose a policy is proposed for adoption in a country.
- It has been tried in other countries and we know outcomes there.
- We also know outcomes in countries where it was not adopted.
- From the historical record, what can we conclude about the likely effectiveness of the policy in countries that have not implemented it?

- To answer questions of this sort, economists build models of counterfactuals.

- To answer questions of this sort, economists build models of counterfactuals.
- Consider the following model.



- To answer questions of this sort, economists build models of counterfactuals.
- Consider the following model.
- Let  $Y_0$  be the outcome of a country (e.g., GDP) under a no-policy regime.

- To answer questions of this sort, economists build models of counterfactuals.
- Consider the following model.
- Let  $Y_0$  be the outcome of a country (e.g., GDP) under a no-policy regime.
- $Y_1$  is the outcome if the policy is implemented.

- To answer questions of this sort, economists build models of counterfactuals.
- Consider the following model.
- Let  $Y_0$  be the outcome of a country (e.g., GDP) under a no-policy regime.
- $Y_1$  is the outcome if the policy is implemented.
- $(Y_1 - Y_0)$  is the “treatment effect” of the policy.

- To answer questions of this sort, economists build models of counterfactuals.
- Consider the following model.
- Let  $Y_0$  be the outcome of a country (e.g., GDP) under a no-policy regime.
- $Y_1$  is the outcome if the policy is implemented.
- $(Y_1 - Y_0)$  is the “treatment effect” of the policy.
- It may vary among countries.

- To answer questions of this sort, economists build models of counterfactuals.
- Consider the following model.
- Let  $Y_0$  be the outcome of a country (e.g., GDP) under a no-policy regime.
- $Y_1$  is the outcome if the policy is implemented.
- $(Y_1 - Y_0)$  is the “treatment effect” of the policy.
- It may vary among countries.
- We observe characteristics  $X$  of various countries (e.g., level of democracy, level of population literacy, etc.).

- To answer questions of this sort, economists build models of counterfactuals.
- Consider the following model.
- Let  $Y_0$  be the outcome of a country (e.g., GDP) under a no-policy regime.
- $Y_1$  is the outcome if the policy is implemented.
- $(Y_1 - Y_0)$  is the “treatment effect” of the policy.
- It may vary among countries.
- We observe characteristics  $X$  of various countries (e.g., level of democracy, level of population literacy, etc.).
- It is convenient to decompose  $Y_1$  into its mean given  $X$ ,  $\mu_1(X)$ , and deviation from mean  $U_1$ .

- We can make a similar decomposition for  $Y_0$ :

$$Y_1 = \mu_1(X) + U_1 \quad (2)$$

$$Y_0 = \mu_0(X) + U_0.$$

- We can make a similar decomposition for  $Y_0$ :

$$\begin{aligned} Y_1 &= \mu_1(X) + U_1 \\ Y_0 &= \mu_0(X) + U_0. \end{aligned} \tag{2}$$

- We do not need to assume additive separability but it is convenient and we initially adopt it to simplify the exposition and establish a parallel regression notation that serves to link the statistical literature on treatment effects with the economic literature.



- We can make a similar decomposition for  $Y_0$ :

$$\begin{aligned} Y_1 &= \mu_1(X) + U_1 \\ Y_0 &= \mu_0(X) + U_0. \end{aligned} \tag{2}$$

- We do not need to assume additive separability but it is convenient and we initially adopt it to simplify the exposition and establish a parallel regression notation that serves to link the statistical literature on treatment effects with the economic literature.
- We develop more general nonseparable models in later sections of this chapter.

- It may happen that controlling for the  $X$ ,  $Y_1 - Y_0$  is the same for all countries.

- It may happen that controlling for the  $X$ ,  $Y_1 - Y_0$  is the same for all countries.
- This is the case of homogeneous treatment effects given  $X$ .

- It may happen that controlling for the  $X$ ,  $Y_1 - Y_0$  is the same for all countries.
- This is the case of homogeneous treatment effects given  $X$ .
- More likely, countries vary in their responses to the policy even after controlling for  $X$ .

- Figure 1 plots the distribution of  $Y_1 - Y_0$  for a benchmark  $X$ .

- Figure 1 plots the distribution of  $Y_1 - Y_0$  for a benchmark  $X$ .
- It also displays the various treatment parameters introduced in Part I.

- Figure 1 plots the distribution of  $Y_1 - Y_0$  for a benchmark  $X$ .
- It also displays the various treatment parameters introduced in Part I.
- We use a special form of the generalized Roy model with constant cost  $C$  of adopting the policy.

- Figure 1 plots the distribution of  $Y_1 - Y_0$  for a benchmark  $X$ .
- It also displays the various treatment parameters introduced in Part I.
- We use a special form of the generalized Roy model with constant cost  $C$  of adopting the policy.
- This is called the “extended Roy model”. We use this model because it is simple and intuitive.

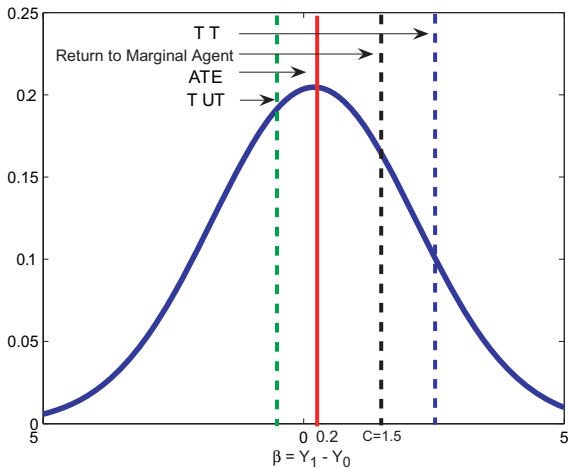


- Figure 1 plots the distribution of  $Y_1 - Y_0$  for a benchmark  $X$ .
- It also displays the various treatment parameters introduced in Part I.
- We use a special form of the generalized Roy model with constant cost  $C$  of adopting the policy.
- This is called the “extended Roy model”. We use this model because it is simple and intuitive.
- (The precise parameterization of the extended Roy model used to generate the figure and the treatment effects is given at the base of figure 1.)

- Figure 1 plots the distribution of  $Y_1 - Y_0$  for a benchmark  $X$ .
- It also displays the various treatment parameters introduced in Part I.
- We use a special form of the generalized Roy model with constant cost  $C$  of adopting the policy.
- This is called the “extended Roy model”. We use this model because it is simple and intuitive.
- (The precise parameterization of the extended Roy model used to generate the figure and the treatment effects is given at the base of figure 1.)
- The special case of homogeneity in  $Y_1 - Y_0$  arises when the distribution collapses to its mean.

Figure 1: Distribution of Gains in the Roy Economy

$$U_1 - U_0 \not\propto D$$



$$TT = 2.666, TUT = -0.632$$

$$\text{Return to Marginal Agent} = C = 1.5$$

$$ATE = \mu_1 - \mu_0 = \bar{\beta} = 0.2$$

---

---

### The Model

---

Outcomes

Choice Model

$$Y_1 = \mu_1 + U_1 = \alpha + \bar{\beta} + U_1$$
$$Y_0 = \mu_0 + U_0 = \alpha + U_0$$
$$D = \begin{cases} 1 & \text{if } D^* \geq 0 \\ 0 & \text{if } D^* < 0 \end{cases}$$

---

General Case

---

$$(U_1 - U_0) \not\propto D$$
$$\text{ATE} \neq \text{TT} \neq \text{TUT}$$

---

---

The researcher observes  $(Y, D, C)$ .

$$Y = \alpha + \beta D + U_0 \text{ where } \beta = Y_1 - Y_0.$$

### Parameterization

$$\begin{aligned} \alpha = 0.67 & \quad (U_1, U_0) \sim N(\mathbf{0}, \Sigma) & D^* = Y_1 - Y_0 - C \\ \bar{\beta} = 0.2 & \quad = \begin{bmatrix} 1 & -0.9 \\ -0.9 & 1 \end{bmatrix} & C = 1.5 \end{aligned}$$

*Source:* Heckman, Urzua and Vytlacil (2006)

- It would be ideal if we could estimate the distribution of  $Y_1 - Y_0$  given  $X$  and there is research that does this.

- It would be ideal if we could estimate the distribution of  $Y_1 - Y_0$  given  $X$  and there is research that does this.
- Abbring and Heckman survey methods for doing so in Part III.

- More often, economists focus on some mean of the distribution displayed in figure 1 and use a regression framework to interpret the data.



- More often, economists focus on some mean of the distribution displayed in figure 1 and use a regression framework to interpret the data.
- To turn (2) into a regression model, it is conventional to use the switching regression framework.

- More often, economists focus on some mean of the distribution displayed in figure 1 and use a regression framework to interpret the data.
- To turn (2) into a regression model, it is conventional to use the switching regression framework.
- Define  $D = 1$  if a country adopts a policy;  $D = 0$  if it does not.

- More often, economists focus on some mean of the distribution displayed in figure 1 and use a regression framework to interpret the data.
- To turn (2) into a regression model, it is conventional to use the switching regression framework.
- Define  $D = 1$  if a country adopts a policy;  $D = 0$  if it does not.
- The observed outcome  $Y$  is the switching regression model (1).

- Substituting (2) into this expression, and keeping all  $X$  implicit, we obtain

$$\begin{aligned} Y &= Y_0 + (Y_1 - Y_0)D & (3) \\ &= \mu_0 + (\mu_1 - \mu_0 + U_1 - U_0)D + U_0. \end{aligned}$$

- Substituting (2) into this expression, and keeping all  $X$  implicit, we obtain

$$\begin{aligned} Y &= Y_0 + (Y_1 - Y_0)D \\ &= \mu_0 + (\mu_1 - \mu_0 + U_1 - U_0)D + U_0. \end{aligned} \quad (3)$$

- Using conventional regression notation,

$$Y = \alpha + \beta D + \varepsilon \quad (4)$$

where  $\alpha = \mu_0$ ,  $\beta = (Y_1 - Y_0) = \mu_1 - \mu_0 + U_1 - U_0$  and  $\varepsilon = U_0$ .

- Substituting (2) into this expression, and keeping all  $X$  implicit, we obtain

$$\begin{aligned} Y &= Y_0 + (Y_1 - Y_0)D \\ &= \mu_0 + (\mu_1 - \mu_0 + U_1 - U_0)D + U_0. \end{aligned} \tag{3}$$

- Using conventional regression notation,

$$Y = \alpha + \beta D + \varepsilon \tag{4}$$

where  $\alpha = \mu_0$ ,  $\beta = (Y_1 - Y_0) = \mu_1 - \mu_0 + U_1 - U_0$  and  $\varepsilon = U_0$ .

- We will also use the notation that  $\eta = U_1 - U_0$ , letting  $\bar{\beta} = \mu_1 - \mu_0$  and  $\beta = \bar{\beta} + \eta$ .

- Substituting (2) into this expression, and keeping all  $X$  implicit, we obtain

$$\begin{aligned} Y &= Y_0 + (Y_1 - Y_0)D \\ &= \mu_0 + (\mu_1 - \mu_0 + U_1 - U_0)D + U_0. \end{aligned} \tag{3}$$

- Using conventional regression notation,

$$Y = \alpha + \beta D + \varepsilon \tag{4}$$

where  $\alpha = \mu_0$ ,  $\beta = (Y_1 - Y_0) = \mu_1 - \mu_0 + U_1 - U_0$  and  $\varepsilon = U_0$ .

- We will also use the notation that  $\eta = U_1 - U_0$ , letting  $\bar{\beta} = \mu_1 - \mu_0$  and  $\beta = \bar{\beta} + \eta$ .
- Throughout this section we use treatment effect and regression notation interchangeably.

- Substituting (2) into this expression, and keeping all  $X$  implicit, we obtain

$$\begin{aligned} Y &= Y_0 + (Y_1 - Y_0)D \\ &= \mu_0 + (\mu_1 - \mu_0 + U_1 - U_0)D + U_0. \end{aligned} \quad (3)$$

- Using conventional regression notation,

$$Y = \alpha + \beta D + \varepsilon \quad (4)$$

where  $\alpha = \mu_0$ ,  $\beta = (Y_1 - Y_0) = \mu_1 - \mu_0 + U_1 - U_0$  and  $\varepsilon = U_0$ .

- We will also use the notation that  $\eta = U_1 - U_0$ , letting  $\bar{\beta} = \mu_1 - \mu_0$  and  $\beta = \bar{\beta} + \eta$ .
- Throughout this section we use treatment effect and regression notation interchangeably.
- The coefficient on  $D$  is the treatment effect.



- The case where  $\beta$  is the same for every country is the case conventionally assumed.

- The case where  $\beta$  is the same for every country is the case conventionally assumed.
- More elaborate versions assume that  $\beta$  depends on  $X$  ( $\beta(X)$ ) and estimates interactions of  $D$  with  $X$ .

- The case where  $\beta$  is the same for every country is the case conventionally assumed.
- More elaborate versions assume that  $\beta$  depends on  $X$  ( $\beta(X)$ ) and estimates interactions of  $D$  with  $X$ .
- The case where  $\beta$  varies even after accounting for  $X$  is called the “random coefficient” or “heterogenous treatment effect” case.

- The case where  $\beta$  is the same for every country is the case conventionally assumed.
- More elaborate versions assume that  $\beta$  depends on  $X$  ( $\beta(X)$ ) and estimates interactions of  $D$  with  $X$ .
- The case where  $\beta$  varies even after accounting for  $X$  is called the “random coefficient” or “heterogenous treatment effect” case.
- The case where  $\eta = U_1 - U_0$  depends on  $D$  is the case of essential heterogeneity analyzed by ?.

- The case where  $\beta$  is the same for every country is the case conventionally assumed.
- More elaborate versions assume that  $\beta$  depends on  $X$  ( $\beta(X)$ ) and estimates interactions of  $D$  with  $X$ .
- The case where  $\beta$  varies even after accounting for  $X$  is called the “random coefficient” or “heterogenous treatment effect” case.
- The case where  $\eta = U_1 - U_0$  depends on  $D$  is the case of essential heterogeneity analyzed by ?.
- This case arises when treatment choices depend at least in part on the idiosyncratic return to treatment.

- The case where  $\beta$  is the same for every country is the case conventionally assumed.
- More elaborate versions assume that  $\beta$  depends on  $X$  ( $\beta(X)$ ) and estimates interactions of  $D$  with  $X$ .
- The case where  $\beta$  varies even after accounting for  $X$  is called the “random coefficient” or “heterogenous treatment effect” case.
- The case where  $\eta = U_1 - U_0$  depends on  $D$  is the case of essential heterogeneity analyzed by ?.
- This case arises when treatment choices depend at least in part on the idiosyncratic return to treatment.
- A great deal of attention has been focused on this case in recent decades and we develop the implications of this model in this chapter.

## An Index Model of Choice and Treatment Effects: Definitions and Unifying Principles

- We now present the model of treatment effects developed in [???](#) and [?](#), which relaxes the normality, separability and exogeneity assumptions invoked in the traditional economic selection models.

## An Index Model of Choice and Treatment Effects: Definitions and Unifying Principles

- We now present the model of treatment effects developed in  $???$  and  $?$ , which relaxes the normality, separability and exogeneity assumptions invoked in the traditional economic selection models.
- It is rich enough to generate all of the treatment effects displayed in figure 1 as well as many other policy parameters.



## An Index Model of Choice and Treatment Effects: Definitions and Unifying Principles

- We now present the model of treatment effects developed in ??? and ?, which relaxes the normality, separability and exogeneity assumptions invoked in the traditional economic selection models.
- It is rich enough to generate all of the treatment effects displayed in figure 1 as well as many other policy parameters.
- It does not require separability.

- It is a nonparametric generalized Roy model with testable restrictions that can be used to unify the treatment effect literature, identify different treatment effects, link the literature on treatment effects to the literature in structural econometrics and interpret the implicit economic assumptions underlying instrumental variables, regression discontinuity design methods, control functions and matching methods.

- It is a nonparametric generalized Roy model with testable restrictions that can be used to unify the treatment effect literature, identify different treatment effects, link the literature on treatment effects to the literature in structural econometrics and interpret the implicit economic assumptions underlying instrumental variables, regression discontinuity design methods, control functions and matching methods.
- We follow ?? and ? in considering binary treatments.

- It is a nonparametric generalized Roy model with testable restrictions that can be used to unify the treatment effect literature, identify different treatment effects, link the literature on treatment effects to the literature in structural econometrics and interpret the implicit economic assumptions underlying instrumental variables, regression discontinuity design methods, control functions and matching methods.
- We follow ?? and ? in considering binary treatments.
- We analyze multiple treatments in Slide 471.

- It is a nonparametric generalized Roy model with testable restrictions that can be used to unify the treatment effect literature, identify different treatment effects, link the literature on treatment effects to the literature in structural econometrics and interpret the implicit economic assumptions underlying instrumental variables, regression discontinuity design methods, control functions and matching methods.
- We follow ?? and ? in considering binary treatments.
- We analyze multiple treatments in Slide 471.
- ? develop a model with a continuum of treatments and we briefly survey that work at the end of Slide 471.

- $Y$  is the measured outcome variable.

- $Y$  is the measured outcome variable.
- It is produced from the switching regression model (1).

- $Y$  is the measured outcome variable.
- It is produced from the switching regression model (1).
- Outcomes are general nonlinear, nonseparable functions of observables and unobservables:

$$Y_1 = \mu_1(X, U_1) \quad (5)$$

$$Y_0 = \mu_0(X, U_0). \quad (6)$$



- $Y$  is the measured outcome variable.
- It is produced from the switching regression model (1).
- Outcomes are general nonlinear, nonseparable functions of observables and unobservables:

$$Y_1 = \mu_1(X, U_1) \quad (5)$$

$$Y_0 = \mu_0(X, U_0). \quad (6)$$

- Examples of models that can be written in this form include conventional latent variable models for discrete choice that are generated by a latent variable crossing a threshold:

$$Y_i = \mathbf{1}(Y_i^* \geq 0), \text{ where } Y_i^* = \mu_i(X) + U_i, i = 0, 1.$$

- $Y$  is the measured outcome variable.
- It is produced from the switching regression model (1).
- Outcomes are general nonlinear, nonseparable functions of observables and unobservables:

$$Y_1 = \mu_1(X, U_1) \quad (5)$$

$$Y_0 = \mu_0(X, U_0). \quad (6)$$

- Examples of models that can be written in this form include conventional latent variable models for discrete choice that are generated by a latent variable crossing a threshold:

$$Y_i = \mathbf{1}(Y_i^* \geq 0), \text{ where } Y_i^* = \mu_i(X) + U_i, i = 0, 1.$$

- Notice that in the general case,  $\mu_i(X, U_i) - E(Y_i | X) \neq U_i$ ,  $i = 0, 1$ .

- As defined in Part I, the individual treatment effect associated with moving an otherwise identical person from “0” to “1” is  $Y_1 - Y_0 = \Delta$  and is defined as the causal effect on  $Y$  of a *ceteris paribus* move from “0” to “1”.

- As defined in Part I, the individual treatment effect associated with moving an otherwise identical person from “0” to “1” is  $Y_1 - Y_0 = \Delta$  and is defined as the causal effect on  $Y$  of a *ceteris paribus* move from “0” to “1”.
- To link this framework to the literature on economic choice models, we characterize the decision rule for program participation by an index model:

$$\begin{aligned} D^* &= \mu_D(Z) - V; & D &= 1 \quad \text{if } D^* \geq 0; \\ & & D &= 0 \quad \text{otherwise,} \end{aligned} \quad (7)$$

where, from the point of view of the econometrician,  $(Z, X)$  is observed and  $(U_0, U_1, V)$  is unobserved.

- As defined in Part I, the individual treatment effect associated with moving an otherwise identical person from “0” to “1” is  $Y_1 - Y_0 = \Delta$  and is defined as the causal effect on  $Y$  of a *ceteris paribus* move from “0” to “1”.
- To link this framework to the literature on economic choice models, we characterize the decision rule for program participation by an index model:

$$D^* = \mu_D(Z) - V; \quad D = 1 \quad \text{if} \quad D^* \geq 0; \\ D = 0 \quad \text{otherwise,} \quad (7)$$

where, from the point of view of the econometrician,  $(Z, X)$  is observed and  $(U_0, U_1, V)$  is unobserved.

- The random variable  $V$  may be a function of  $(U_0, U_1)$ .

- For example, in the original Roy Model,  $\mu_1$  and  $\mu_0$  are additively separable in  $U_1$  and  $U_0$  respectively, and  $V = -[U_1 - U_0]$ .

- For example, in the original Roy Model,  $\mu_1$  and  $\mu_0$  are additively separable in  $U_1$  and  $U_0$  respectively, and  $V = -[U_1 - U_0]$ .
- In the original formulations of the generalized Roy model, outcome equations are separable and  $V = -[U_1 - U_0 - U_C]$ , where  $U_C$  arises from the cost function (recall the discussion in section 3.3 of Part I).

- For example, in the original Roy Model,  $\mu_1$  and  $\mu_0$  are additively separable in  $U_1$  and  $U_0$  respectively, and  $V = -[U_1 - U_0]$ .
- In the original formulations of the generalized Roy model, outcome equations are separable and  $V = -[U_1 - U_0 - U_C]$ , where  $U_C$  arises from the cost function (recall the discussion in section 3.3 of Part I).
- Without loss of generality, we define  $Z$  so that it includes all of the elements of  $X$  as well as any additional variables unique to the choice equation.



- We invoke the following assumptions that are weaker than those used in the conventional literature on structural econometrics or the recent literature on semiparametric selection models and at the same time can be used both to define and to identify different treatment parameters.

The assumptions are:

(A-1)

$(U_0, U_1, V)$  are independent of  $Z$  conditional on  $X$  (**Independence**);

(A-2)

$\mu_D(Z)$  is a nondegenerate random variable conditional on  $X$   
(**Rank Condition**);

(A-3)

The distribution of  $V$  is continuous;

(A-4)

The values of  $E(Y_1)$  and  $E(Y_0)$  are finite (**Finite Means**);

(A-5)

$0 < \Pr(D = 1 | X) < 1$ .

- (A-1) assumes that  $V$  is independent of  $Z$  given  $X$ , and is used below to generate counterfactuals.

- (A-1) assumes that  $V$  is independent of  $Z$  given  $X$ , and is used below to generate counterfactuals.
- For the definition of treatment effects, we do not need either (A-1) or (A-2).

- (A-1) assumes that  $V$  is independent of  $Z$  given  $X$ , and is used below to generate counterfactuals.
- For the definition of treatment effects, we do not need either (A-1) or (A-2).
- Our definitions of treatment effects and their unification through MTE do not require any elements of  $Z$  that are not elements of  $X$  or independence assumptions.

- (A-1) assumes that  $V$  is independent of  $Z$  given  $X$ , and is used below to generate counterfactuals.
- For the definition of treatment effects, we do not need either (A-1) or (A-2).
- Our definitions of treatment effects and their unification through MTE do not require any elements of  $Z$  that are not elements of  $X$  or independence assumptions.
- However, our analysis of instrumental variables requires that  $Z$  contain at least one element not in  $X$ .

- (A-1) assumes that  $V$  is independent of  $Z$  given  $X$ , and is used below to generate counterfactuals.
- For the definition of treatment effects, we do not need either (A-1) or (A-2).
- Our definitions of treatment effects and their unification through MTE do not require any elements of  $Z$  that are not elements of  $X$  or independence assumptions.
- However, our analysis of instrumental variables requires that  $Z$  contain at least one element not in  $X$ .
- Assumptions (A-1) or (A-2) justify application of instrumental variables methods and nonparametric selection or control function methods.



- (A-1) assumes that  $V$  is independent of  $Z$  given  $X$ , and is used below to generate counterfactuals.
- For the definition of treatment effects, we do not need either (A-1) or (A-2).
- Our definitions of treatment effects and their unification through MTE do not require any elements of  $Z$  that are not elements of  $X$  or independence assumptions.
- However, our analysis of instrumental variables requires that  $Z$  contain at least one element not in  $X$ .
- Assumptions (A-1) or (A-2) justify application of instrumental variables methods and nonparametric selection or control function methods.
- Some parameters in the recent IV literature are defined by an instrument so we make assumptions about instruments up front, noting where they are not needed.

- Assumption (A-4) is needed to satisfy standard integration conditions.

- Assumption (A-4) is needed to satisfy standard integration conditions.
- It guarantees that the mean treatment parameters are well defined.

- Assumption (A-4) is needed to satisfy standard integration conditions.
- It guarantees that the mean treatment parameters are well defined.
- Assumption (A-5) is the assumption in the population of both a treatment and a control group for each  $X$ .

- Assumption (A-4) is needed to satisfy standard integration conditions.
- It guarantees that the mean treatment parameters are well defined.
- Assumption (A-5) is the assumption in the population of both a treatment and a control group for each  $X$ .
- Observe that there are no exogeneity requirements for  $X$ .

- Assumption (A-4) is needed to satisfy standard integration conditions.
- It guarantees that the mean treatment parameters are well defined.
- Assumption (A-5) is the assumption in the population of both a treatment and a control group for each  $X$ .
- Observe that there are no exogeneity requirements for  $X$ .
- This is in contrast with the assumptions commonly made in the conventional structural literature and the semiparametric selection literature (see, e.g., ?).

- A counterfactual “no feedback” condition facilitates interpretability so that conditioning on  $X$  does not mask the effects of  $D$ .

Letting  $X_d$  denote a value of  $X$  if  $D$  is set to  $d$ , a sufficient condition that rules out feedback from  $D$  to  $X$  is:

(A-6)

*Let  $X_0$  denote the counterfactual value of  $X$  that would be observed if  $D$  is set to 0.  $X_1$  is defined analogously. Assume  $X_d = X$  for  $d = 0, 1$ . (The  $X_D$  are invariant to counterfactual manipulations.)*



- Condition (A-6) is not strictly required to formulate an evaluation model, but it enables an analyst who conditions on  $X$  to capture the “total” or “full effect” of  $D$  on  $Y$  (see ?).

- Condition (A-6) is not strictly required to formulate an evaluation model, but it enables an analyst who conditions on  $X$  to capture the “total” or “full effect” of  $D$  on  $Y$  (see ?).
- This assumption imposes the requirement that  $X$  is an external variable determined outside the model and is not affected by counterfactual manipulations of  $D$ .

- Condition (A-6) is not strictly required to formulate an evaluation model, but it enables an analyst who conditions on  $X$  to capture the “total” or “full effect” of  $D$  on  $Y$  (see ?).
- This assumption imposes the requirement that  $X$  is an external variable determined outside the model and is not affected by counterfactual manipulations of  $D$ .
- However, the assumption allows for  $X$  to be freely correlated with  $U_1$ ,  $U_0$  and  $V$  so it can be endogenous.

- Condition (A-6) is not strictly required to formulate an evaluation model, but it enables an analyst who conditions on  $X$  to capture the “total” or “full effect” of  $D$  on  $Y$  (see ?).
- This assumption imposes the requirement that  $X$  is an external variable determined outside the model and is not affected by counterfactual manipulations of  $D$ .
- However, the assumption allows for  $X$  to be freely correlated with  $U_1$ ,  $U_0$  and  $V$  so it can be endogenous.
- Until we discuss the problems of external validity and policy forecasting in Slide 412, we analyze treatment effects conditional on  $X$ , and maintain assumption (A-6).

- In this notation,  $P(Z)$  is the probability of receiving treatment given  $Z$ , or the “propensity score”

$P(Z) \equiv \Pr(D = 1 \mid Z) = F_{V|X}(\mu_D(Z))$ , where  $F_{V|X}(\cdot)$  denotes the distribution of  $V$  conditional on  $X$ .

- In this notation,  $P(Z)$  is the probability of receiving treatment given  $Z$ , or the “propensity score”  
 $P(Z) \equiv \Pr(D = 1 \mid Z) = F_{V|X}(\mu_D(Z))$ , where  $F_{V|X}(\cdot)$  denotes the distribution of  $V$  conditional on  $X$ .
- We sometimes denote  $P(Z)$  by  $P$ , suppressing the  $Z$  argument.

- In this notation,  $P(Z)$  is the probability of receiving treatment given  $Z$ , or the “propensity score”  
 $P(Z) \equiv \Pr(D = 1 \mid Z) = F_{V|X}(\mu_D(Z))$ , where  $F_{V|X}(\cdot)$  denotes the distribution of  $V$  conditional on  $X$ .
- We sometimes denote  $P(Z)$  by  $P$ , suppressing the  $Z$  argument.
- We also work with  $U_D$ , a uniform random variable ( $U_D \sim \text{Unif}[0, 1]$ ) defined by  $U_D = F_{V|X}(V)$ .

- In this notation,  $P(Z)$  is the probability of receiving treatment given  $Z$ , or the “propensity score”  
 $P(Z) \equiv \Pr(D = 1 \mid Z) = F_{V|X}(\mu_D(Z))$ , where  $F_{V|X}(\cdot)$  denotes the distribution of  $V$  conditional on  $X$ .
- We sometimes denote  $P(Z)$  by  $P$ , suppressing the  $Z$  argument.
- We also work with  $U_D$ , a uniform random variable ( $U_D \sim \text{Unif}[0, 1]$ ) defined by  $U_D = F_{V|X}(V)$ .
- The separability between  $V$  and  $\mu_D(Z)$  or  $D(Z)$  and  $U_D$  is conventional.



- In this notation,  $P(Z)$  is the probability of receiving treatment given  $Z$ , or the “propensity score”  
 $P(Z) \equiv \Pr(D = 1 \mid Z) = F_{V|X}(\mu_D(Z))$ , where  $F_{V|X}(\cdot)$  denotes the distribution of  $V$  conditional on  $X$ .
- We sometimes denote  $P(Z)$  by  $P$ , suppressing the  $Z$  argument.
- We also work with  $U_D$ , a uniform random variable ( $U_D \sim \text{Unif}[0, 1]$ ) defined by  $U_D = F_{V|X}(V)$ .
- The separability between  $V$  and  $\mu_D(Z)$  or  $D(Z)$  and  $U_D$  is conventional.
- It plays a crucial role in justifying instrumental variable estimators in the general models analyzed in this chapter.

- ? establishes that assumptions (A-1)–(A-5) for selection model (1) and (5)–(7) are equivalent to the assumptions used to generate the LATE model of ? which are developed below in Slide 152.

- ? establishes that assumptions (A-1)–(A-5) for selection model (1) and (5)–(7) are equivalent to the assumptions used to generate the LATE model of ? which are developed below in Slide 152.
- Thus the nonparametric selection model for treatment effects developed by Heckman and Vytlacil is implied by the assumptions of the Imbens-Angrist instrumental variable model for treatment effects.

- ? establishes that assumptions (A-1)–(A-5) for selection model (1) and (5)–(7) are equivalent to the assumptions used to generate the LATE model of ? which are developed below in Slide 152.
- Thus the nonparametric selection model for treatment effects developed by Heckman and Vytlacil is implied by the assumptions of the Imbens-Angrist instrumental variable model for treatment effects.
- Our approach links the IV literature to the literature on economic choice models explicated in Part I.

- ? establishes that assumptions (A-1)–(A-5) for selection model (1) and (5)–(7) are equivalent to the assumptions used to generate the LATE model of ? which are developed below in Slide 152.
- Thus the nonparametric selection model for treatment effects developed by Heckman and Vytlacil is implied by the assumptions of the Imbens-Angrist instrumental variable model for treatment effects.
- Our approach links the IV literature to the literature on economic choice models exposted in Part I.
- Our latent variable model is a version of the standard sample selection bias model.

- ? establishes that assumptions (A-1)–(A-5) for selection model (1) and (5)–(7) are equivalent to the assumptions used to generate the LATE model of ? which are developed below in Slide 152.
- Thus the nonparametric selection model for treatment effects developed by Heckman and Vytlacil is implied by the assumptions of the Imbens-Angrist instrumental variable model for treatment effects.
- Our approach links the IV literature to the literature on economic choice models explicated in Part I.
- Our latent variable model is a version of the standard sample selection bias model.
- We weave together two strands of the literature often thought to be distinct (see e.g., ?).

- The model of equations (5)-(7) and assumptions (A-1)–(A-5) impose two testable restrictions on the distribution of  $(Y, D, Z, X)$ .

- The model of equations (5)-(7) and assumptions (A-1)–(A-5) impose two testable restrictions on the distribution of  $(Y, D, Z, X)$ .
- First, it imposes an index sufficiency restriction: for any set  $\mathcal{A}$  and for  $j = 0, 1$ ,

$$\Pr(Y_j \in \mathcal{A} \mid X, Z, D = j) = \Pr(Y_j \in \mathcal{A} \mid X, P(Z), D = j).$$



- The model of equations (5)-(7) and assumptions (A-1)–(A-5) impose two testable restrictions on the distribution of  $(Y, D, Z, X)$ .
- First, it imposes an index sufficiency restriction: for any set  $\mathcal{A}$  and for  $j = 0, 1$ ,

$$\Pr(Y_j \in \mathcal{A} \mid X, Z, D = j) = \Pr(Y_j \in \mathcal{A} \mid X, P(Z), D = j).$$

- $Z$  (given  $X$ ) enters the model only through the propensity score  $P(Z)$ .

- The model of equations (5)-(7) and assumptions (A-1)–(A-5) impose two testable restrictions on the distribution of  $(Y, D, Z, X)$ .
- First, it imposes an index sufficiency restriction: for any set  $\mathcal{A}$  and for  $j = 0, 1$ ,

$$\Pr(Y_j \in \mathcal{A} \mid X, Z, D = j) = \Pr(Y_j \in \mathcal{A} \mid X, P(Z), D = j).$$

- $Z$  (given  $X$ ) enters the model only through the propensity score  $P(Z)$ .
- This restriction has empirical content when  $Z$  contains two or more variables not in  $X$ .

- The model of equations (5)-(7) and assumptions (A-1)–(A-5) impose two testable restrictions on the distribution of  $(Y, D, Z, X)$ .
- First, it imposes an index sufficiency restriction: for any set  $\mathcal{A}$  and for  $j = 0, 1$ ,

$$\Pr(Y_j \in \mathcal{A} \mid X, Z, D = j) = \Pr(Y_j \in \mathcal{A} \mid X, P(Z), D = j).$$

- $Z$  (given  $X$ ) enters the model only through the propensity score  $P(Z)$ .
- This restriction has empirical content when  $Z$  contains two or more variables not in  $X$ .
- Second, the model also imposes monotonicity in  $p$  for  $E(YD \mid X = x, P = p)$  and  $E(Y(1 - D) \mid X = x, P = p)$ .

- The model of equations (5)-(7) and assumptions (A-1)–(A-5) impose two testable restrictions on the distribution of  $(Y, D, Z, X)$ .
- First, it imposes an index sufficiency restriction: for any set  $\mathcal{A}$  and for  $j = 0, 1$ ,

$$\Pr(Y_j \in \mathcal{A} \mid X, Z, D = j) = \Pr(Y_j \in \mathcal{A} \mid X, P(Z), D = j).$$

- $Z$  (given  $X$ ) enters the model only through the propensity score  $P(Z)$ .
- This restriction has empirical content when  $Z$  contains two or more variables not in  $X$ .
- Second, the model also imposes monotonicity in  $p$  for  $E(YD \mid X = x, P = p)$  and  $E(Y(1 - D) \mid X = x, P = p)$ .
- ?, appendix A develop this condition further, and show that it is testable.

- Even though the model of treatment effects we exposit is not the most general possible model, it has testable implications and hence empirical content.

- Even though the model of treatment effects we exposit is not the most general possible model, it has testable implications and hence empirical content.
- It unites various literatures and produces a nonparametric version of the selection model, and links the treatment literature to economic choice theory.

- Even though the model of treatment effects we exposit is not the most general possible model, it has testable implications and hence empirical content.
- It unites various literatures and produces a nonparametric version of the selection model, and links the treatment literature to economic choice theory.
- We compare the assumptions used to identify IV with the assumptions used in matching in Slide 675.

## Definitions of Treatment Effects in the Two Outcome Model

- As developed in Part I, the difficulty of observing the same individual in both treated and untreated states leads to the use of various population level treatment effects widely used in the biostatistics literature and often applied in economics.



## Definitions of Treatment Effects in the Two Outcome Model

- As developed in Part I, the difficulty of observing the same individual in both treated and untreated states leads to the use of various population level treatment effects widely used in the biostatistics literature and often applied in economics.
- The most commonly invoked treatment effect is the Average Treatment Effect (ATE):  $\Delta^{\text{ATE}}(x) \equiv E(\Delta \mid X = x)$  where  $\Delta = Y_1 - Y_0$ .

## Definitions of Treatment Effects in the Two Outcome Model

- As developed in Part I, the difficulty of observing the same individual in both treated and untreated states leads to the use of various population level treatment effects widely used in the biostatistics literature and often applied in economics.
- The most commonly invoked treatment effect is the Average Treatment Effect (ATE):  $\Delta^{\text{ATE}}(x) \equiv E(\Delta \mid X = x)$  where  $\Delta = Y_1 - Y_0$ .
- This is the effect of assigning treatment randomly to everyone of type  $X$  assuming full compliance, and ignoring general equilibrium effects.

- The average impact of treatment on persons who actually take the treatment is Treatment on the Treated (TT):  
 $\Delta^{TT}(x) \equiv E(\Delta \mid X = x, D = 1).$

- The average impact of treatment on persons who actually take the treatment is Treatment on the Treated (TT):

$$\Delta^{\text{TT}}(x) \equiv E(\Delta \mid X = x, D = 1).$$

- This parameter can also be defined conditional on  $P(Z)$ :

$$\Delta^{\text{TT}}(x, p) \equiv E(\Delta \mid X = x, P(Z) = p, D = 1).$$

- The mean effect of treatment on those for whom  $X = x$  and  $U_D = u_D$ , the Marginal Treatment Effect (MTE), plays a fundamental role in the analysis of this chapter:

$$\Delta^{\text{MTE}}(x, u_D) \equiv E(\Delta \mid X = x, U_D = u_D). \quad (8)$$

- The mean effect of treatment on those for whom  $X = x$  and  $U_D = u_D$ , the Marginal Treatment Effect (MTE), plays a fundamental role in the analysis of this chapter:

$$\Delta^{\text{MTE}}(x, u_D) \equiv E(\Delta \mid X = x, U_D = u_D). \quad (8)$$

- This parameter is defined independently of any instrument.

- The mean effect of treatment on those for whom  $X = x$  and  $U_D = u_D$ , the Marginal Treatment Effect (MTE), plays a fundamental role in the analysis of this chapter:

$$\Delta^{\text{MTE}}(x, u_D) \equiv E(\Delta \mid X = x, U_D = u_D). \quad (8)$$

- This parameter is defined independently of any instrument.
- We separate the definition of parameters from their identification.

- The mean effect of treatment on those for whom  $X = x$  and  $U_D = u_D$ , the Marginal Treatment Effect (MTE), plays a fundamental role in the analysis of this chapter:

$$\Delta^{\text{MTE}}(x, u_D) \equiv E(\Delta \mid X = x, U_D = u_D). \quad (8)$$

- This parameter is defined independently of any instrument.
- We separate the definition of parameters from their identification.
- The MTE is the expected effect of treatment conditional on observed characteristics  $X$  and conditional on  $U_D$ , the unobservables from the first stage decision rule.



- For  $u_D$  evaluation points close to zero,  $\Delta^{\text{MTE}}(x, u_D)$  is the expected effect of treatment on individuals with the value of unobservables that make them most likely to participate in treatment and who would participate even if the mean scale utility  $\mu_D(Z)$  is small.

- For  $u_D$  evaluation points close to zero,  $\Delta^{\text{MTE}}(x, u_D)$  is the expected effect of treatment on individuals with the value of unobservables that make them most likely to participate in treatment and who would participate even if the mean scale utility  $\mu_D(Z)$  is small.
- If  $U_D$  is large,  $\mu_D(Z)$  would have to be large to induce people to participate.

- One can also interpret  $E(\Delta \mid X = x, U_D = u_D)$  as the mean gain in terms of  $Y_1 - Y_0$  for persons with observed characteristics  $X$  who would be indifferent between treatment or not if they were randomly assigned a value of  $Z$ , say  $z$ , such that  $\mu_D(z) = u_D$ .

- One can also interpret  $E(\Delta \mid X = x, U_D = u_D)$  as the mean gain in terms of  $Y_1 - Y_0$  for persons with observed characteristics  $X$  who would be indifferent between treatment or not if they were randomly assigned a value of  $Z$ , say  $z$ , such that  $\mu_D(z) = u_D$ .
- When  $Y_0$  and  $Y_1$  are value outcomes, MTE is a mean willingness-to-pay measure.

- One can also interpret  $E(\Delta \mid X = x, U_D = u_D)$  as the mean gain in terms of  $Y_1 - Y_0$  for persons with observed characteristics  $X$  who would be indifferent between treatment or not if they were randomly assigned a value of  $Z$ , say  $z$ , such that  $\mu_D(z) = u_D$ .
- When  $Y_0$  and  $Y_1$  are value outcomes, MTE is a mean willingness-to-pay measure.
- MTE is a choice-theoretic building block that unites the treatment effect, selection, matching and control function literatures.

- A third interpretation is that MTE conditions on  $X$  and the residual defined by subtracting the expectation of  $D^*$  from  $D^*$ :  
$$\tilde{U}_D = D^* - E(D^* | Z, X).$$

- A third interpretation is that MTE conditions on  $X$  and the residual defined by subtracting the expectation of  $D^*$  from  $D^*$ :  
$$\tilde{U}_D = D^* - E(D^* | Z, X).$$
- This is a “replacement function” interpretation in the sense of ? and ?, or “control function” interpretation in the sense of ?.

- A third interpretation is that MTE conditions on  $X$  and the residual defined by subtracting the expectation of  $D^*$  from  $D^*$ :  
$$\tilde{U}_D = D^* - E(D^* | Z, X).$$
- This is a “replacement function” interpretation in the sense of ? and ?, or “control function” interpretation in the sense of ?.
- These three interpretations are equivalent under separability in  $D^*$ , i.e., when (7) characterizes the choice equation, but lead to three different definitions of MTE when a more general nonseparable model is developed.



- This point is developed in Slide 370 where we discuss a general nonseparable model.

- This point is developed in Slide 370 where we discuss a general nonseparable model.
- The additive separability of equation (7) in terms of observables and unobservables plays a crucial role in the justification of instrumental variable methods.

- The LATE parameter of  $\tau$  is a version of MTE.

- The LATE parameter of  $\tau$  is a version of MTE.
- We present their full conditions for identification in Slide 152.

- The LATE parameter of  $\tau$  is a version of MTE.
- We present their full conditions for identification in Slide 152.
- Here we define it in the notation used in this chapter.

- The LATE parameter of ? is a version of MTE.
- We present their full conditions for identification in Slide 152.
- Here we define it in the notation used in this chapter.
- LATE is defined by an instrument in their analysis.

- The LATE parameter of ? is a version of MTE.
- We present their full conditions for identification in Slide 152.
- Here we define it in the notation used in this chapter.
- LATE is defined by an instrument in their analysis.
- As in Part I, we define LATE independently of any instrument after first presenting the Imbens–Angrist definition.

- The LATE parameter of ? is a version of MTE.
- We present their full conditions for identification in Slide 152.
- Here we define it in the notation used in this chapter.
- LATE is defined by an instrument in their analysis.
- As in Part I, we define LATE independently of any instrument after first presenting the Imbens–Angrist definition.
- Define  $D(z)$  as a counterfactual choice variable, with  $D(z) = 1$  if state 1 ( $D = 1$ ) would have been chosen if  $Z$  had been set to  $z$ , and  $D(z) = 0$  otherwise.



- The LATE parameter of ? is a version of MTE.
- We present their full conditions for identification in Slide 152.
- Here we define it in the notation used in this chapter.
- LATE is defined by an instrument in their analysis.
- As in Part I, we define LATE independently of any instrument after first presenting the Imbens–Angrist definition.
- Define  $D(z)$  as a counterfactual choice variable, with  $D(z) = 1$  if state 1 ( $D = 1$ ) would have been chosen if  $Z$  had been set to  $z$ , and  $D(z) = 0$  otherwise.
- Let  $\mathcal{Z}(x)$  denote the support of the distribution of  $Z$  conditional on  $X = x$ .

- For any  $(z, z') \in \mathcal{Z}(x) \times \mathcal{Z}(x)$  such that  $P(z) > P(z')$ , LATE is  $E(\Delta \mid X = x, D(z) = 1, D(z') = 0) = E(Y_1 - Y_0 \mid X = x, D(z) = 1, D(z') = 0)$ , the mean gain to persons who would be induced to switch from  $D = 0$  to  $D = 1$  if  $Z$  were manipulated externally from  $z'$  to  $z$ .

- For any  $(z, z') \in \mathcal{Z}(x) \times \mathcal{Z}(x)$  such that  $P(z) > P(z')$ , LATE is  $E(\Delta \mid X = x, D(z) = 1, D(z') = 0) = E(Y_1 - Y_0 \mid X = x, D(z) = 1, D(z') = 0)$ , the mean gain to persons who would be induced to switch from  $D = 0$  to  $D = 1$  if  $Z$  were manipulated externally from  $z'$  to  $z$ .
- In an example of the returns to education,  $z'$  could be the base level of tuition and  $z$  a reduced tuition level.

- Using the latent index model, developed in Part I and defined in the introduction to this section, ?? show that LATE can be written as

$$\begin{aligned} & E(Y_1 - Y_0 \mid X = x, D(z) = 1, D(z') = 0) \\ &= E(Y_1 - Y_0 \mid X = x, u'_D < U_D \leq u_D) \\ &= \Delta^{\text{LATE}}(x, u_D, u'_D) \end{aligned}$$

for  $u_D = \Pr(D(z) = 1) = P(z)$ ,  $u'_D = \Pr(D(z') = 1) = P(z')$ , where assumption (A-1) implies that  $\Pr(D(z) = 1) = \Pr(D = 1 \mid Z = z)$  and  $\Pr(D(z') = 1) = \Pr(D = 1 \mid Z = z')$ .

- Using the latent index model, developed in Part I and defined in the introduction to this section, ?? show that LATE can be written as

$$\begin{aligned} E(Y_1 - Y_0 \mid X = x, D(z) = 1, D(z') = 0) \\ &= E(Y_1 - Y_0 \mid X = x, u'_D < U_D \leq u_D) \\ &= \Delta^{\text{LATE}}(x, u_D, u'_D) \end{aligned}$$

for  $u_D = \Pr(D(z) = 1) = P(z)$ ,  $u'_D = \Pr(D(z') = 1) = P(z')$ , where assumption (A-1) implies that  $\Pr(D(z) = 1) = \Pr(D = 1 \mid Z = z)$  and  $\Pr(D(z') = 1) = \Pr(D = 1 \mid Z = z')$ .

- Imbens and Angrist define the LATE parameter as the probability limit of an estimator.

- Using the latent index model, developed in Part I and defined in the introduction to this section, ?? show that LATE can be written as

$$\begin{aligned} E(Y_1 - Y_0 \mid X = x, D(z) = 1, D(z') = 0) \\ &= E(Y_1 - Y_0 \mid X = x, u'_D < U_D \leq u_D) \\ &= \Delta^{\text{LATE}}(x, u_D, u'_D) \end{aligned}$$

for  $u_D = \Pr(D(z) = 1) = P(z)$ ,  $u'_D = \Pr(D(z') = 1) = P(z')$ , where assumption (A-1) implies that  $\Pr(D(z) = 1) = \Pr(D = 1 \mid Z = z)$  and  $\Pr(D(z') = 1) = \Pr(D = 1 \mid Z = z')$ .

- Imbens and Angrist define the LATE parameter as the probability limit of an estimator.
- Their analysis conflates issues of definition of parameters with issues of identification.

- Our representation of LATE allows us to separate these two conceptually distinct matters and to define the LATE parameter more generally.

- Our representation of LATE allows us to separate these two conceptually distinct matters and to define the LATE parameter more generally.
- One can, in principle, evaluate the right hand side of the preceding equation at any  $u_D, u'_D$  points in the unit interval and not only at points in the support of the distribution of the propensity score  $P(Z)$  conditional on  $X = x$  where it is identified.



- Our representation of LATE allows us to separate these two conceptually distinct matters and to define the LATE parameter more generally.
- One can, in principle, evaluate the right hand side of the preceding equation at any  $u_D, u'_D$  points in the unit interval and not only at points in the support of the distribution of the propensity score  $P(Z)$  conditional on  $X = x$  where it is identified.
- From assumptions (A-1), (A-3), and (A-4),  $\Delta^{\text{LATE}}(x, u_D, u'_D)$  is continuous in  $u_D$  and  $u'_D$  and
$$\lim_{u'_D \uparrow u_D} \Delta^{\text{LATE}}(x, u_D, u'_D) = \Delta^{\text{MTE}}(x, u_D).$$

- ? use assumptions (A-1)–(A-5) and the latent index structure to develop the relationship between MTE and the various treatment effect parameters shown in the first three lines of table 2A.

- ? use assumptions (A-1)–(A-5) and the latent index structure to develop the relationship between MTE and the various treatment effect parameters shown in the first three lines of table 2A.
- Appendix, Slide 1030, presents the formal derivation of the parameters and associated weights and graphically illustrates the relationship between ATE and TT.

- ? use assumptions (A-1)–(A-5) and the latent index structure to develop the relationship between MTE and the various treatment effect parameters shown in the first three lines of table 2A.
- Appendix, Slide 1030, presents the formal derivation of the parameters and associated weights and graphically illustrates the relationship between ATE and TT.
- There we establish that all treatment parameters may be expressed as weighted averages of the MTE:

$$\text{Treatment Parameter } (j) = \int_0^1 \Delta^{\text{MTE}}(x, u_D) \omega_j(x, u_D) du_D$$

where  $\omega_j(x, u_D)$  is the weighting function for the MTE and the integral is defined over the full support of  $u_D$ .

**Table 2:** A. Treatment effects and estimands as weighted averages of the marginal treatment effect

---

---

$$\text{ATE}(x) = E(Y_1 - Y_0 \mid X = x) = \int_0^1 \Delta^{\text{MTE}}(x, u_D) du_D$$

$$\text{TT}(x) = E(Y_1 - Y_0 \mid X = x, D = 1) = \int_0^1 \Delta^{\text{MTE}}(x, u_D) \omega_{\text{TT}}(x, u_D) du_D$$

$$\text{TUT}(x) = E(Y_1 - Y_0 \mid X = x, D = 0) = \int_0^1 \Delta^{\text{MTE}}(x, u_D) \omega_{\text{TUT}}(x, u_D) du_D$$

$$\text{PRTE}(x) = E(Y_{a'} \mid X = x) - E(Y_a \mid X = x) = \int_0^1 \Delta^{\text{MTE}}(x, u_D) \omega_{\text{PRTE}}(x, u_D) du_D$$

for two policies  $a$  and  $a'$  that affect the  $Z$  but not the  $X$

$$\text{IV}_J(x) = \int_0^1 \Delta^{\text{MTE}}(x, u_D) \omega_{\text{IV}}^J(x, u_D) du_D, \text{ given instrument } J$$

$$\text{OLS}(x) = \int_0^1 \Delta^{\text{MTE}}(x, u_D) \omega_{\text{OLS}}(x, u_D) du_D$$

---

## B. Weights

---

$$\omega_{\text{ATE}}(x, u_D) = 1$$

$$\omega_{\text{TT}}(x, u_D) = \left[ \int_{u_D}^1 f_{P|X}(p | X = x) dp \right] \frac{1}{E(P | X = x)}$$

$$\omega_{\text{TUT}}(x, u_D) = \left[ \int_0^{u_D} f_{P|X}(p | X = x) dp \right] \frac{1}{E((1 - P) | X = x)}$$

$$\omega_{\text{PRTE}}(x, u_D) = \left[ \frac{F_{P_{a'}} | X(u_D | x) - F_{P_a} | X(u_D | x)}{\Delta \bar{P}(x)} \right], \text{ where } \Delta \bar{P}(x) = E(P_a | X = x) - E(P_{a'} | X = x)$$

$$\omega_{\text{IV}}^J(x, u_D) = \left[ \int_{u_D}^1 \int (J(Z) - E(J(Z) | X = x)) f_{J,P|X}(j, t | X = x) dj dt \right] \frac{1}{\text{Cov}(J(Z), D | X = x)}$$

$$\omega_{\text{OLS}}(x, u_D) = 1 + \frac{E(U_1 | X = x, U_D = u_D) \omega_1(x, u_D) - E(U_0 | X = x, U_D = u_D) \omega_0(x, u_D)}{\Delta^{\text{MTE}}(x, u_D)}$$

$$\omega_1(x, u_D) = \left[ \int_{u_D}^1 f_{P|X}(p | X = x) dp \right] \frac{1}{E(P | X = x)}$$

$$\omega_0(x, u_D) = \left[ \int_0^{u_D} f_{P|X}(p | X = x) dp \right] \frac{1}{E((1 - P) | X = x)}$$

---

Source: ?

- Except for the OLS weights, the weights in the table all integrate to one, although in some cases the weights for IV may be negative.

- Except for the OLS weights, the weights in the table all integrate to one, although in some cases the weights for IV may be negative.
- We analyze how negative weights for IV might arise in Slide 152.



- In table 2A,  $\Delta^{\text{TT}}(x)$  is shown as a weighted average of  $\Delta^{\text{MTE}}$ :

$$\Delta^{\text{TT}}(x) = \int_0^1 \Delta^{\text{MTE}}(x, u_D) \omega_{\text{TT}}(x, u_D) du_D,$$

where

$$\omega_{\text{TT}}(x, u_D) = \frac{1 - F_{P|X}(u_D | x)}{\int_0^1 (1 - F_{P|X}(t | x)) dt} = \frac{S_{P|X}(u_D | x)}{E(P(Z) | X = x)}, \quad (9)$$

and  $S_{P|X}(u_D | x)$  is  $\Pr(P(Z) > u_D | X = x)$  and  $\omega_{\text{TT}}(x, u_D)$  is a weighted distribution.

- In table 2A,  $\Delta^{\text{TT}}(x)$  is shown as a weighted average of  $\Delta^{\text{MTE}}$ :

$$\Delta^{\text{TT}}(x) = \int_0^1 \Delta^{\text{MTE}}(x, u_D) \omega_{\text{TT}}(x, u_D) du_D,$$

where

$$\omega_{\text{TT}}(x, u_D) = \frac{1 - F_{P|X}(u_D | x)}{\int_0^1 (1 - F_{P|X}(t | x)) dt} = \frac{S_{P|X}(u_D | x)}{E(P(Z) | X = x)}, \quad (9)$$

and  $S_{P|X}(u_D | x)$  is  $\Pr(P(Z) > u_D | X = x)$  and  $\omega_{\text{TT}}(x, u_D)$  is a weighted distribution.

- The parameter  $\Delta^{\text{TT}}(x)$  oversamples  $\Delta^{\text{MTE}}(x, u_D)$  for those individuals with low values of  $u_D$  that make them more likely to participate in the program being evaluated.

- Treatment on the untreated (TUT) is defined symmetrically with TT and oversamples those least likely to participate.

- Treatment on the untreated (TUT) is defined symmetrically with TT and oversamples those least likely to participate.
- The various weights are displayed in table 2B.

- Treatment on the untreated (TUT) is defined symmetrically with TT and oversamples those least likely to participate.
- The various weights are displayed in table 2B.
- The other weights, treatment effects and estimands shown in this table are discussed later.

- Treatment on the untreated (TUT) is defined symmetrically with TT and oversamples those least likely to participate.
- The various weights are displayed in table 2B.
- The other weights, treatment effects and estimands shown in this table are discussed later.
- A central theme of this chapter is that under our assumptions all estimators and estimands can be written as weighted averages of MTE.

- Treatment on the untreated (TUT) is defined symmetrically with TT and oversamples those least likely to participate.
- The various weights are displayed in table 2B.
- The other weights, treatment effects and estimands shown in this table are discussed later.
- A central theme of this chapter is that under our assumptions all estimators and estimands can be written as weighted averages of MTE.
- This allows us to unify the treatment effect literature using a common functional  $\Delta^{\text{MTE}}(x, u_D)$ .

**Table 3:** Treatment parameters and estimands in the generalized Roy example

Treatment on the Treated	0.2353
Treatment on the Untreated	0.1574
Average Treatment Effect	0.2000
Sorting Gain <sup>a</sup>	0.0353
Policy Relevant Treatment Effect (PRTE)	0.1549
Selection Bias <sup>b</sup>	-0.0628
Linear Instrumental Variables <sup>c</sup>	0.2013
Ordinary Least Squares	0.1725

$${}^a TT - ATE = E(Y_1 - Y_0 | D = 1) - E(Y_1 - Y_0)$$

$${}^b OLS - TT = E(Y_0 | D = 1) - E(Y_0 | D = 0)$$

<sup>c</sup>Using Propensity Score  $P(Z)$  as the instrument.

Note: The model used to create Table 3 is the same as those used to create Figures 2A and 2B. The PRTE is computed using a policy  $t$  characterized as follows:

If  $Z > 0$  then  $D = 1$  if  $Z(1 + t) - V \geq 0$ .

If  $Z \leq 0$  then  $D = 1$  if  $Z - V \geq 0$ .

For this example  $t$  is set equal to 0.2.

Source: Heckman and Vytlačil (2005)



- Observe that if

$E(Y_1 - Y_0 \mid X = x, U_D = u_D) = E(Y_1 - Y_0 \mid X = x)$ , so

$\Delta = Y_1 - Y_0$  is mean independent of  $U_D$  given  $X = x$ , then

$\Delta^{\text{MTE}} = \Delta^{\text{ATE}} = \Delta^{\text{TT}} = \Delta^{\text{LATE}}$ .

- Observe that if

$E(Y_1 - Y_0 \mid X = x, U_D = u_D) = E(Y_1 - Y_0 \mid X = x)$ , so

$\Delta = Y_1 - Y_0$  is mean independent of  $U_D$  given  $X = x$ , then

$$\Delta^{\text{MTE}} = \Delta^{\text{ATE}} = \Delta^{\text{TT}} = \Delta^{\text{LATE}}.$$

- Therefore, in cases where there is no heterogeneity in terms of unobservables in MTE ( $\Delta$  constant conditional on  $X = x$ ) or agents do not act on it so that  $U_D$  drops out of the conditioning set, marginal treatment effects are average treatment effects, so that all of the evaluation parameters are the same.

- Observe that if

$E(Y_1 - Y_0 \mid X = x, U_D = u_D) = E(Y_1 - Y_0 \mid X = x)$ , so

$\Delta = Y_1 - Y_0$  is mean independent of  $U_D$  given  $X = x$ , then

$$\Delta^{\text{MTE}} = \Delta^{\text{ATE}} = \Delta^{\text{TT}} = \Delta^{\text{LATE}}.$$

- Therefore, in cases where there is no heterogeneity in terms of unobservables in MTE ( $\Delta$  constant conditional on  $X = x$ ) or agents do not act on it so that  $U_D$  drops out of the conditioning set, marginal treatment effects are average treatment effects, so that all of the evaluation parameters are the same.
- Otherwise, they are different.

- Observe that if

$E(Y_1 - Y_0 \mid X = x, U_D = u_D) = E(Y_1 - Y_0 \mid X = x)$ , so

$\Delta = Y_1 - Y_0$  is mean independent of  $U_D$  given  $X = x$ , then

$$\Delta^{\text{MTE}} = \Delta^{\text{ATE}} = \Delta^{\text{TT}} = \Delta^{\text{LATE}}.$$

- Therefore, in cases where there is no heterogeneity in terms of unobservables in MTE ( $\Delta$  constant conditional on  $X = x$ ) or agents do not act on it so that  $U_D$  drops out of the conditioning set, marginal treatment effects are average treatment effects, so that all of the evaluation parameters are the same.
- Otherwise, they are different.
- Only in the case where the marginal treatment effect is the average treatment effect will the “effect” of treatment be uniquely defined.

- Figure 2A plots weights for a parametric normal generalized Roy model generated from the parameters shown at the base of figure 2B.

- Figure 2A plots weights for a parametric normal generalized Roy model generated from the parameters shown at the base of figure 2B.
- This is an instance of the general model developed in Part I, section 5.

- Figure 2A plots weights for a parametric normal generalized Roy model generated from the parameters shown at the base of figure 2B.
- This is an instance of the general model developed in Part I, section 5.
- The model allows for costs to vary in the population and is more general than the extended Roy model.

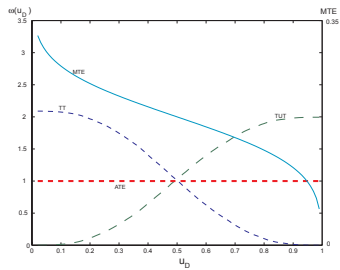
- Figure 2A plots weights for a parametric normal generalized Roy model generated from the parameters shown at the base of figure 2B.
- This is an instance of the general model developed in Part I, section 5.
- The model allows for costs to vary in the population and is more general than the extended Roy model.
- We discuss the weights for IV depicted in figure 2B in Slide 152 and the weights for OLS in Slide 675.



- Figure 2A plots weights for a parametric normal generalized Roy model generated from the parameters shown at the base of figure 2B.
- This is an instance of the general model developed in Part I, section 5.
- The model allows for costs to vary in the population and is more general than the extended Roy model.
- We discuss the weights for IV depicted in figure 2B in Slide 152 and the weights for OLS in Slide 675.
- A high  $u_D$  is associated with higher cost, relative to return, and less likelihood of choosing  $D = 1$ .

- Figure 2A plots weights for a parametric normal generalized Roy model generated from the parameters shown at the base of figure 2B.
- This is an instance of the general model developed in Part I, section 5.
- The model allows for costs to vary in the population and is more general than the extended Roy model.
- We discuss the weights for IV depicted in figure 2B in Slide 152 and the weights for OLS in Slide 675.
- A high  $u_D$  is associated with higher cost, relative to return, and less likelihood of choosing  $D = 1$ .
- The decline of MTE in terms of higher values of  $u_D$  means that people with higher  $u_D$  have lower gross returns.

Figure 2: A. Weights for the marginal treatment effect for different parameters



- TT overweights low values of  $u_D$  (i.e., it oversamples  $U_D$  that make it likely to have  $D = 1$ ).

- TT overweights low values of  $u_D$  (i.e., it oversamples  $U_D$  that make it likely to have  $D = 1$ ).
- ATE samples  $U_D$  uniformly.

- TT overweights low values of  $u_D$  (i.e., it oversamples  $U_D$  that make it likely to have  $D = 1$ ).
- ATE samples  $U_D$  uniformly.
- Treatment on the Untreated ( $E(Y_1 - Y_0 | X = x, D = 0)$ ), or TUT, oversamples the values of  $U_D$  which make it unlikely to have  $D = 1$ .

- Table 3 shows the treatment parameters produced from the different weighting schemes for the model used to generate the weights in figures 2A and 2B.



- Table 3 shows the treatment parameters produced from the different weighting schemes for the model used to generate the weights in figures 2A and 2B.
- Given the decline of the MTE in  $u_D$ , it is not surprising that  $TT > ATE > TUT$ .

- Table 3 shows the treatment parameters produced from the different weighting schemes for the model used to generate the weights in figures 2A and 2B.
- Given the decline of the MTE in  $u_D$ , it is not surprising that  $TT > ATE > TUT$ .
- This is the generalized Roy version of the principle of diminishing returns.

- Table 3 shows the treatment parameters produced from the different weighting schemes for the model used to generate the weights in figures 2A and 2B.
- Given the decline of the MTE in  $u_D$ , it is not surprising that  $TT > ATE > TUT$ .
- This is the generalized Roy version of the principle of diminishing returns.
- Those most likely to self select into the program benefit the most from it.

- Table 3 shows the treatment parameters produced from the different weighting schemes for the model used to generate the weights in figures 2A and 2B.
- Given the decline of the MTE in  $u_D$ , it is not surprising that  $TT > ATE > TUT$ .
- This is the generalized Roy version of the principle of diminishing returns.
- Those most likely to self select into the program benefit the most from it.
- The difference between TT and ATE is a sorting gain:  $E(Y_1 - Y_0 | X, D = 1) - E(Y_1 - Y_0 | X)$ , the average gain experienced by people who sort into treatment compared to what the average person would experience.

- Purposive selection on the basis of gains should lead to positive sorting gains of the kind found in the table.

- Purposive selection on the basis of gains should lead to positive sorting gains of the kind found in the table.
- If there is negative sorting on the gains, then  $TUT \geq ATE \geq TT$ .

- Purposive selection on the basis of gains should lead to positive sorting gains of the kind found in the table.
- If there is negative sorting on the gains, then  $TUT \geq ATE \geq TT$ .
- Later in this chapter, we return to this table to discuss the other numbers in it.

- Table 4 reproduced from ? presents evidence on the nonconstancy of the MTE in  $U_D$  drawn from a variety of studies of schooling, job training, migration and unionism.



- Table 4 reproduced from ? presents evidence on the nonconstancy of the MTE in  $U_D$  drawn from a variety of studies of schooling, job training, migration and unionism.
- Most of the evidence is obtained using parametric normal selection models or variants of such models.

- Table 4 reproduced from ? presents evidence on the nonconstancy of the MTE in  $U_D$  drawn from a variety of studies of schooling, job training, migration and unionism.
- Most of the evidence is obtained using parametric normal selection models or variants of such models.
- With the exception of studies of unionism, a common finding in the empirical literature is the nonconstancy of MTE given  $X$ .

- Table 4 reproduced from ? presents evidence on the nonconstancy of the MTE in  $U_D$  drawn from a variety of studies of schooling, job training, migration and unionism.
- Most of the evidence is obtained using parametric normal selection models or variants of such models.
- With the exception of studies of unionism, a common finding in the empirical literature is the nonconstancy of MTE given  $X$ .
- The evidence from the literature suggests that different treatment parameters measure different effects, and persons participate in programs based on heterogeneity in responses to the program being studied.

- Table 4 reproduced from ? presents evidence on the nonconstancy of the MTE in  $U_D$  drawn from a variety of studies of schooling, job training, migration and unionism.
- Most of the evidence is obtained using parametric normal selection models or variants of such models.
- With the exception of studies of unionism, a common finding in the empirical literature is the nonconstancy of MTE given  $X$ .
- The evidence from the literature suggests that different treatment parameters measure different effects, and persons participate in programs based on heterogeneity in responses to the program being studied.
- The phenomenon of nonconstancy of the MTE that we analyze in this chapter is of substantial empirical interest.

**Table 4:** Evidence of selection on unobservables and constancy of the MTE for separable models.

Study	Method	Finding on the Hypothesis of Constancy of the MTE
Unionism		
Lee (1978)	Normal Selection Model ( $H_0 : \sigma_{1V} = \sigma_{0V}$ )	$\sigma_{1V} = \sigma_{0V}$ Do not reject
Farber (1983)	Normal Selection Model ( $H_0 : \sigma_{1V} = \sigma_{0V}$ )	$\sigma_{1V} = \sigma_{0V}$ Do not reject
Duncan and Leigh (1985)	Normal Selection Model ( $H_0 : \sigma_{1V} = \sigma_{0V}$ )	$\sigma_{1V} = \sigma_{0V}$ Do not reject
Robinson (1989)	Normal Selection Model ( $(\mu_1 - \mu_0)_{IV} = (\mu_1 - \mu_0)_{\text{normal}}$ )	$\sigma_{1V} \neq \sigma_{0V}$ Do not reject
Schooling (College vs. High School)		
Willis and Rosen (1979)	Normal Selection Model ( $H_0 : \sigma_{1V} = \sigma_{0V}$ )	$\sigma_{1V} \neq \sigma_{0V}$ Reject
Heckman, Tobias and Vytlacil (2003)	Normal Selection Model ( $H_0 : \sigma_{1V} = \sigma_{0V}$ )	$\sigma_{1V} \neq \sigma_{0V}$ Reject

Study	Method	Finding on the Hypothesis of Constancy of the MTE
Job Training		
Björklund and Moffitt (1987)	Normal Selection Model ( $H_0 : \sigma_{1V} = \sigma_{0V}$ )	$\sigma_{1V} \neq \sigma_{0V}$ Reject
Heckman, Ichimura, Smith and Todd (1998; Supplement)	$E(U_1 - U_0   D = 1, Z, X) = E(U_1 - U_0   D = 1, X)$	Reject selection on unobservables
Sectoral Choice		
Heckman and Sedlacek (1990)	Normal Selection Model ( $H_0 : \sigma_{1V} = \sigma_{0V}$ )	$\sigma_{1V} \neq \sigma_{0V}$ Reject
Migration		
Pessino (1991)	Normal Selection Model ( $H_0 : \sigma_{1V} = \sigma_{0V}$ )	$\sigma_{1V} \neq \sigma_{0V}$ Reject
Tunali (2000)	$H_0 : E(U_1 - U_0   D = 1) = 0$ (estimated using robust selection)	Cannot reject

Notes:  $Y = DY_1 + (1 - D)Y_0$ .

$$Y_1 = \mu_1(X) + U_1$$

$$Y_0 = \mu_0(X) + U_0$$

$$Z \perp\!\!\!\perp (U_0, U_1), Z \not\perp D$$

$D = \mathbf{1}(\mu_D(Z) - V \geq 0)$ , where  $\mu_D(Z) - V$  is the index determining selection into "1" or "0"

Hypothesis: No Selection on Unobservables (Constancy of the MTE)

$H_0 : E(U_1 - U_0 | D = 1, Z, X)$  does not depend on  $D$  where

Source: Heckman (2001)

- The additively separable latent index model for  $D$  (equation (7)) and assumptions (A-1)–(A-5) are far stronger than what is required to define the parameters in terms of the MTE.



- The additively separable latent index model for  $D$  (equation (7)) and assumptions (A-1)–(A-5) are far stronger than what is required to define the parameters in terms of the MTE.
- The representations of treatment effects defined in table 2A remain valid even if  $Z$  is not independent of  $U_D$ , if there are no variables in  $Z$  that are not also contained in  $X$ , or if a more general nonseparable choice model generates  $D$  (so  $D^* = \mu_D(Z, U_D)$ ).

- The additively separable latent index model for  $D$  (equation (7)) and assumptions (A-1)–(A-5) are far stronger than what is required to define the parameters in terms of the MTE.
- The representations of treatment effects defined in table 2A remain valid even if  $Z$  is not independent of  $U_D$ , if there are no variables in  $Z$  that are not also contained in  $X$ , or if a more general nonseparable choice model generates  $D$  (so  $D^* = \mu_D(Z, U_D)$ ).
- An important advantage of our approach over other approaches to the analysis of instrumental variables in the recent literature is that no instrument  $Z$  is needed to define the parameters.

- The additively separable latent index model for  $D$  (equation (7)) and assumptions (A-1)–(A-5) are far stronger than what is required to define the parameters in terms of the MTE.
- The representations of treatment effects defined in table 2A remain valid even if  $Z$  is not independent of  $U_D$ , if there are no variables in  $Z$  that are not also contained in  $X$ , or if a more general nonseparable choice model generates  $D$  (so  $D^* = \mu_D(Z, U_D)$ ).
- An important advantage of our approach over other approaches to the analysis of instrumental variables in the recent literature is that no instrument  $Z$  is needed to define the parameters.
- We separate the tasks of definition and identification of parameters as discussed in table 1 of Part I, and present an analysis more closely rooted in economics.

- Appendices, Slides 1030 and 1049, define the treatment parameters for both separable (Appendix, Slide 1030) and nonseparable choice equations (Appendix, Slide 1049).

- Appendices, Slides 1030 and 1049, define the treatment parameters for both separable (Appendix, Slide 1030) and nonseparable choice equations (Appendix, Slide 1049).
- We show that the treatment parameters can be defined even if there is no instrument or if instrumental variables methods break down as they do in nonseparable models.

- As noted in Part I, the literature on structural econometrics is clear about the basic parameters of interest although it is not always clear about the exact combinations of parameters needed to answer specific policy problems.

- As noted in Part I, the literature on structural econometrics is clear about the basic parameters of interest although it is not always clear about the exact combinations of parameters needed to answer specific policy problems.
- The literature on treatment effects offers a variety of evaluation parameters.

- As noted in Part I, the literature on structural econometrics is clear about the basic parameters of interest although it is not always clear about the exact combinations of parameters needed to answer specific policy problems.
- The literature on treatment effects offers a variety of evaluation parameters.
- Missing from that literature is an algorithm for defining treatment effects that answer precisely formulated economic policy questions.



- As noted in Part I, the literature on structural econometrics is clear about the basic parameters of interest although it is not always clear about the exact combinations of parameters needed to answer specific policy problems.
- The literature on treatment effects offers a variety of evaluation parameters.
- Missing from that literature is an algorithm for defining treatment effects that answer precisely formulated economic policy questions.
- The MTE provides a framework for developing such an algorithm.

- As noted in Part I, the literature on structural econometrics is clear about the basic parameters of interest although it is not always clear about the exact combinations of parameters needed to answer specific policy problems.
- The literature on treatment effects offers a variety of evaluation parameters.
- Missing from that literature is an algorithm for defining treatment effects that answer precisely formulated economic policy questions.
- The MTE provides a framework for developing such an algorithm.
- In the next section, we present one well defined policy parameter that can be used to generate Benthamite policy evaluations as discussed in section 5 of Part I.

## Policy Relevant Treatment Parameters

- The conventional treatment parameters do not always answer economically interesting questions.

## Policy Relevant Treatment Parameters

- The conventional treatment parameters do not always answer economically interesting questions.
- Their link to cost-benefit analysis and interpretable economic frameworks is sometimes obscure.

## Policy Relevant Treatment Parameters

- The conventional treatment parameters do not always answer economically interesting questions.
- Their link to cost-benefit analysis and interpretable economic frameworks is sometimes obscure.
- Each answers a different question.

## Policy Relevant Treatment Parameters

- The conventional treatment parameters do not always answer economically interesting questions.
- Their link to cost-benefit analysis and interpretable economic frameworks is sometimes obscure.
- Each answers a different question.
- Many investigators estimate a treatment effect and hope that it answers an interesting question.

- A more promising approach for defining parameters is to postulate a policy question or decision problem of interest and to derive the treatment parameter that answers it.

- A more promising approach for defining parameters is to postulate a policy question or decision problem of interest and to derive the treatment parameter that answers it.
- Taking this approach does not in general produce the conventional treatment parameters or the estimands produced from instrumental variables.



- Consider a class of policies that affect  $P$ , the probability of participation in a program, but do not affect  $\Delta^{\text{MTE}}$ .

- Consider a class of policies that affect  $P$ , the probability of participation in a program, but do not affect  $\Delta^{\text{MTE}}$ .
- The policies analyzed in the treatment effect literature that change the  $Z$  not in  $X$  are more restrictive than the general policies that shift  $X$  and  $Z$  analyzed in the structural literature.

- Consider a class of policies that affect  $P$ , the probability of participation in a program, but do not affect  $\Delta^{\text{MTE}}$ .
- The policies analyzed in the treatment effect literature that change the  $Z$  not in  $X$  are more restrictive than the general policies that shift  $X$  and  $Z$  analyzed in the structural literature.
- An example from the schooling literature would be policies that change tuition or distance to school but do not directly affect the gross returns to schooling (?).

- Consider a class of policies that affect  $P$ , the probability of participation in a program, but do not affect  $\Delta^{\text{MTE}}$ .
- The policies analyzed in the treatment effect literature that change the  $Z$  not in  $X$  are more restrictive than the general policies that shift  $X$  and  $Z$  analyzed in the structural literature.
- An example from the schooling literature would be policies that change tuition or distance to school but do not directly affect the gross returns to schooling (?).
- Since we ignore general equilibrium effects in this chapter, the effects on  $(Y_0, Y_1)$  from changes in the overall level of education are assumed to be negligible.

- Let  $p$  and  $p'$  denote two potential policies and let  $D_p$  and  $D_{p'}$  denote the choices that would be made under policies  $p$  and  $p'$ .

- Let  $p$  and  $p'$  denote two potential policies and let  $D_p$  and  $D_{p'}$  denote the choices that would be made under policies  $p$  and  $p'$ .
- When we discuss the Policy Relevant Treatment Effect, we use “ $p$ ” to denote the policy and distinguish it from the realized value of  $P(Z)$ .

- Let  $p$  and  $p'$  denote two potential policies and let  $D_p$  and  $D_{p'}$  denote the choices that would be made under policies  $p$  and  $p'$ .
- When we discuss the Policy Relevant Treatment Effect, we use “ $p$ ” to denote the policy and distinguish it from the realized value of  $P(Z)$ .
- Under our assumptions, the policies affect the  $Z$  given  $X$ , but not the potential outcomes.

- Let  $p$  and  $p'$  denote two potential policies and let  $D_p$  and  $D_{p'}$  denote the choices that would be made under policies  $p$  and  $p'$ .
- When we discuss the Policy Relevant Treatment Effect, we use “ $p$ ” to denote the policy and distinguish it from the realized value of  $P(Z)$ .
- Under our assumptions, the policies affect the  $Z$  given  $X$ , but not the potential outcomes.
- Let the corresponding decision rules be  $D_p = \mathbf{1}[P_p(Z_p) \geq U_D]$ ,  $D_{p'} = \mathbf{1}[P_{p'}(Z_{p'}) \geq U_D]$ , where  $P_p(Z_p) = \Pr(D_p = 1 \mid Z_p)$  and  $P_{p'}(Z_{p'}) = \Pr(D_{p'} = 1 \mid Z_{p'})$ .



- Let  $p$  and  $p'$  denote two potential policies and let  $D_p$  and  $D_{p'}$  denote the choices that would be made under policies  $p$  and  $p'$ .
- When we discuss the Policy Relevant Treatment Effect, we use “ $p$ ” to denote the policy and distinguish it from the realized value of  $P(Z)$ .
- Under our assumptions, the policies affect the  $Z$  given  $X$ , but not the potential outcomes.
- Let the corresponding decision rules be  $D_p = \mathbf{1}[P_p(Z_p) \geq U_D]$ ,  $D_{p'} = \mathbf{1}[P_{p'}(Z_{p'}) \geq U_D]$ , where  $P_p(Z_p) = \Pr(D_p = 1 \mid Z_p)$  and  $P_{p'}(Z_{p'}) = \Pr(D_{p'} = 1 \mid Z_{p'})$ .
- To simplify the exposition, we will suppress the arguments of these functions and write  $P_p$  and  $P_{p'}$  for  $P_p(Z_p)$  and  $P_{p'}(Z_{p'})$ .

- Define  $(Y_{0,p}, Y_{1,p}, U_{D,p})$  as  $(Y_0, Y_1, U_D)$  under policy  $p$ , and define  $(Y_{0,p'}, Y_{1,p'}, U_{D,p'})$  correspondingly under policy  $p'$ .

- Define  $(Y_{0,p}, Y_{1,p}, U_{D,p})$  as  $(Y_0, Y_1, U_D)$  under policy  $p$ , and define  $(Y_{0,p'}, Y_{1,p'}, U_{D,p'})$  correspondingly under policy  $p'$ .
- We assume that  $Z_p$  and  $Z_{p'}$  are independent of  $(Y_{0,p}, Y_{1,p}, U_{D,p})$  and  $(Y_{0,p'}, Y_{1,p'}, U_{D,p'})$  respectively, conditional on  $X_p$  and  $X_{p'}$ .

- Define  $(Y_{0,p}, Y_{1,p}, U_{D,p})$  as  $(Y_0, Y_1, U_D)$  under policy  $p$ , and define  $(Y_{0,p'}, Y_{1,p'}, U_{D,p'})$  correspondingly under policy  $p'$ .
- We assume that  $Z_p$  and  $Z_{p'}$  are independent of  $(Y_{0,p}, Y_{1,p}, U_{D,p})$  and  $(Y_{0,p'}, Y_{1,p'}, U_{D,p'})$  respectively, conditional on  $X_p$  and  $X_{p'}$ .
- Let  $Y_p = D_p Y_{1,p} + (1 - D_p) Y_{0,p}$  and  $Y_{p'} = D_{p'} Y_{1,p'} + (1 - D_{p'}) Y_{0,p'}$  denote the outcomes that would be observed under policies  $p$  and  $p'$ , respectively.

- $\Delta^{\text{MTE}}$  is policy invariant in the sense of Hurwicz as defined in Part I if  $E(Y_{1,p} \mid U_{D,p} = u_D, X_p = x)$  and  $E(Y_{0,p} \mid U_{D,p} = u_D, X_p = x)$  are invariant to the choice of policy  $p$  (**Policy Invariance for the Marginal Treatment Effect**).

- $\Delta^{\text{MTE}}$  is policy invariant in the sense of Hurwicz as defined in Part I if  $E(Y_{1,p} \mid U_{D,p} = u_D, X_p = x)$  and  $E(Y_{0,p} \mid U_{D,p} = u_D, X_p = x)$  are invariant to the choice of policy  $p$  (**Policy Invariance for the Marginal Treatment Effect**).
- Policy invariance can be justified by the strong assumption that the policy being investigated does not change the counterfactual outcomes, covariates, or unobservables, i.e.,  $(Y_{0,p}, Y_{1,p}, X_p, U_{D,p}) = (Y_{0,p'}, Y_{1,p'}, X_{p'}, U_{D,p'})$ .

However,  $\Delta^{\text{MTE}}$  is policy invariant if this assumption is relaxed to the weaker assumption that the policy change does not affect the distribution of these variables conditional on  $X$ :

(A-7)

*The distribution of  $(Y_{0,p}, Y_{1,p}, U_{D,p})$  conditional on  $X_p = x$  is the same as the distribution of  $(Y_{0,p'}, Y_{1,p'}, U_{D,p'})$  conditional on  $X_{p'} = x$  (**policy invariance for distribution**).*

- Assumption (A-7) guarantees that manipulations of the distribution of  $Z$  do not affect anything in the model except the choice of outcomes.



- Assumption (A-7) guarantees that manipulations of the distribution of  $Z$  do not affect anything in the model except the choice of outcomes.
- These are specialized versions of (PI-3) and (PI-4) invoked in Part I.

- For the widely used Benthamite social welfare criterion  $\Upsilon(Y)$ , where  $\Upsilon$  is a utility function, comparing policies using mean utilities of outcomes and considering the effect for individuals with a given level of  $X = x$  we obtain the *policy relevant treatment effect*, PRTE, denoted  $\Delta^{\text{PRTE}}(x)$ :

$$\begin{aligned} & E(\Upsilon(Y_p) \mid X = x) - E(\Upsilon(Y_{p'}) \mid X = x) \\ &= \int_0^1 \Delta_{\Upsilon}^{\text{MTE}}(x, u_D) \{F_{P_{p'} \mid X}(u_D \mid x) - F_{P_p \mid X}(u_D \mid x)\} du_D, \quad (10) \end{aligned}$$

where  $F_{P_p \mid X}(\cdot \mid x)$  and  $F_{P_{p'} \mid X}(\cdot \mid x)$  are the distributions of  $P_p$  and  $P_{p'}$  conditional on  $X = x$ , respectively, defined for the different policy regimes and

$$\Delta_{\Upsilon}^{\text{MTE}}(x, u_D) = E(\Upsilon(Y_{1,p}) - \Upsilon(Y_{0,p}) \mid U_{D,p} = u_D, X_p = x).$$

- The weights in expression (10) are derived in Appendix, Slide 1082 under the assumption that the policy does not change the joint distribution of outcomes.

- The weights in expression (10) are derived in Appendix, Slide 1082 under the assumption that the policy does not change the joint distribution of outcomes.
- To simplify the notation, throughout the rest of this chapter when we discuss PRTE, we assume that  $\Upsilon(Y) = Y$ .

- The weights in expression (10) are derived in Appendix, Slide 1082 under the assumption that the policy does not change the joint distribution of outcomes.
- To simplify the notation, throughout the rest of this chapter when we discuss PRTE, we assume that  $\Upsilon(Y) = Y$ .
- Modifications of our analysis for the more general case are straightforward.

- The weights in expression (10) are derived in Appendix, Slide 1082 under the assumption that the policy does not change the joint distribution of outcomes.
- To simplify the notation, throughout the rest of this chapter when we discuss PRTE, we assume that  $\Upsilon(Y) = Y$ .
- Modifications of our analysis for the more general case are straightforward.
- We also discuss the implications of noninvariance for the definition and interpretation of the PRTE in Appendix, Slide 1082.

- Define  $\Delta \bar{P}(x) = E(P_p | X = x) - E(P_{p'} | X = x)$ , the change in the proportion of people induced into the program due to the intervention.

- Define  $\Delta \bar{P}(x) = E(P_p | X = x) - E(P_{p'} | X = x)$ , the change in the proportion of people induced into the program due to the intervention.
- Assuming  $\Delta \bar{P}(x)$  is positive, we may define per person affected weights as  $\omega_{\text{PRTE}}(x, u_D) = \frac{F_{P_{p'}|X}(u_D|x) - F_{P_p|X}(u_D|x)}{\Delta \bar{P}(x)}$ .



- Define  $\Delta \bar{P}(x) = E(P_p | X = x) - E(P_{p'} | X = x)$ , the change in the proportion of people induced into the program due to the intervention.
- Assuming  $\Delta \bar{P}(x)$  is positive, we may define per person affected weights as  $\omega_{\text{PRTE}}(x, u_D) = \frac{F_{P_{p'}|X}(u_D|x) - F_{P_p|X}(u_D|x)}{\Delta \bar{P}(x)}$ .
- These weights are displayed in table 2B.

- Define  $\Delta\bar{P}(x) = E(P_p | X = x) - E(P_{p'} | X = x)$ , the change in the proportion of people induced into the program due to the intervention.
- Assuming  $\Delta\bar{P}(x)$  is positive, we may define per person affected weights as  $\omega_{\text{PRTE}}(x, u_D) = \frac{F_{P_{p'}|X}(u_D|x) - F_{P_p|X}(u_D|x)}{\Delta\bar{P}(x)}$ .
- These weights are displayed in table 2B.
- As demonstrated in the next section, in general, conventional IV weights the MTE differently than either the conventional treatment parameters ( $\Delta^{\text{ATE}}$  or  $\Delta^{\text{TT}}$ ) or the policy relevant parameter, and so does not recover these parameters.

- Instead of hoping that conventional treatment parameters or favorite estimators answer interesting economic questions, the approach developed by **????** is to estimate the MTE and weight it by the appropriate weight determined by how the policy changes the distribution of  $P$  to construct  $\Delta^{\text{PRTE}}$ .

- Instead of hoping that conventional treatment parameters or favorite estimators answer interesting economic questions, the approach developed by **????** is to estimate the MTE and weight it by the appropriate weight determined by how the policy changes the distribution of  $P$  to construct  $\Delta^{\text{PRTE}}$ .
- In **?**, we also develop an alternative approach that produces a policy weighted instrument to identify  $\Delta^{\text{PRTE}}$  by standard instrumental variables.

- We elaborate our discussion of policy analysis based in the MTE and develop other policy parameters for local and global perturbations of policy in Slide 412 after developing the instrumental variable estimator and the related regression discontinuity estimator.

- We elaborate our discussion of policy analysis based in the MTE and develop other policy parameters for local and global perturbations of policy in Slide 412 after developing the instrumental variable estimator and the related regression discontinuity estimator.
- The analyses of Slides 152 and 402 give us tools to make specific the discussion of alternative approaches to policy evaluation.

## Instrumental Variables

- The method of instrumental variables (IV) is currently the most widely used method in economics for estimating economic models when unobservables are present that violate the matching assumption (M-1).

## Instrumental Variables

- The method of instrumental variables (IV) is currently the most widely used method in economics for estimating economic models when unobservables are present that violate the matching assumption (M-1).
- We first present an intuitive exposition of the method and then present a more formal development.



## Instrumental Variables

- The method of instrumental variables (IV) is currently the most widely used method in economics for estimating economic models when unobservables are present that violate the matching assumption (M-1).
- We first present an intuitive exposition of the method and then present a more formal development.
- We analyze a model with two outcomes.

## Instrumental Variables

- The method of instrumental variables (IV) is currently the most widely used method in economics for estimating economic models when unobservables are present that violate the matching assumption (M-1).
- We first present an intuitive exposition of the method and then present a more formal development.
- We analyze a model with two outcomes.
- We generalize the analysis to multiple outcomes in Slide 471.

- Return to the policy adoption example presented at the end of Slide 12.

- Return to the policy adoption example presented at the end of Slide 12.
- The distribution of returns to adoption is depicted in figure 1.

- Return to the policy adoption example presented at the end of Slide 12.
- The distribution of returns to adoption is depicted in figure 1.
- First, consider the method of IV, where  $\beta$  (given  $X$ ), which is the same as  $Y_1 - Y_0$  given  $X$ , is the same for every country.

- Return to the policy adoption example presented at the end of Slide 12.
- The distribution of returns to adoption is depicted in figure 1.
- First, consider the method of IV, where  $\beta$  (given  $X$ ), which is the same as  $Y_1 - Y_0$  given  $X$ , is the same for every country.
- This is the familiar case and we develop it first.

- Return to the policy adoption example presented at the end of Slide 12.
- The distribution of returns to adoption is depicted in figure 1.
- First, consider the method of IV, where  $\beta$  (given  $X$ ), which is the same as  $Y_1 - Y_0$  given  $X$ , is the same for every country.
- This is the familiar case and we develop it first.
- The model is

$$Y = \alpha + \beta D + \varepsilon, \quad (11)$$

where conditioning on  $X$  is implicit.

- A simple least squares regression of  $Y$  on  $D$  (equivalently a mean difference in outcomes between countries with  $D = 1$  and countries with  $D = 0$ ) is possibly subject to a selection bias on  $Y_0$ .



- A simple least squares regression of  $Y$  on  $D$  (equivalently a mean difference in outcomes between countries with  $D = 1$  and countries with  $D = 0$ ) is possibly subject to a selection bias on  $Y_0$ .
- Countries that adopt the policy may be atypical in terms of their  $Y_0$  ( $= \alpha + \varepsilon$ ).

- A simple least squares regression of  $Y$  on  $D$  (equivalently a mean difference in outcomes between countries with  $D = 1$  and countries with  $D = 0$ ) is possibly subject to a selection bias on  $Y_0$ .
- Countries that adopt the policy may be atypical in terms of their  $Y_0$  ( $= \alpha + \varepsilon$ ).
- Thus if countries that would have done well in terms of unobservable  $\varepsilon$  ( $= U_0$ ) even in the absence of the policy are the ones that adopt the policy,  $\beta$  estimated from OLS (or its semiparametric version – matching) is upward biased because  $\text{Cov}(D, \varepsilon) > 0$ .

- If there is an instrument  $Z$ , with the properties that

$$\text{Cov}(Z, D) \neq 0 \quad (12)$$

$$\text{Cov}(Z, \varepsilon) = 0, \quad (13)$$

then standard IV identifies  $\beta$ , at least in large samples,

$$\text{plim } \hat{\beta}_{\text{IV}} = \frac{\text{Cov}(Z, Y)}{\text{Cov}(Z, D)} = \beta.$$

- If there is an instrument  $Z$ , with the properties that

$$\text{Cov}(Z, D) \neq 0 \quad (12)$$

$$\text{Cov}(Z, \varepsilon) = 0, \quad (13)$$

then standard IV identifies  $\beta$ , at least in large samples,

$$\text{plim } \hat{\beta}_{\text{IV}} = \frac{\text{Cov}(Z, Y)}{\text{Cov}(Z, D)} = \beta.$$

- If other instruments exist, each identifies  $\beta$ .

- If there is an instrument  $Z$ , with the properties that

$$\text{Cov}(Z, D) \neq 0 \quad (12)$$

$$\text{Cov}(Z, \varepsilon) = 0, \quad (13)$$

then standard IV identifies  $\beta$ , at least in large samples,

$$\text{plim } \hat{\beta}_{\text{IV}} = \frac{\text{Cov}(Z, Y)}{\text{Cov}(Z, D)} = \beta.$$

- If other instruments exist, each identifies  $\beta$ .
- $Z$  produces a controlled variation in  $D$  relative to  $\varepsilon$ .

- If there is an instrument  $Z$ , with the properties that

$$\text{Cov}(Z, D) \neq 0 \quad (12)$$

$$\text{Cov}(Z, \varepsilon) = 0, \quad (13)$$

then standard IV identifies  $\beta$ , at least in large samples,

$$\text{plim } \hat{\beta}_{\text{IV}} = \frac{\text{Cov}(Z, Y)}{\text{Cov}(Z, D)} = \beta.$$

- If other instruments exist, each identifies  $\beta$ .
- $Z$  produces a controlled variation in  $D$  relative to  $\varepsilon$ .
- Randomization of assignment with full compliance to experimental protocols is an example of an instrument.

- If there is an instrument  $Z$ , with the properties that

$$\text{Cov}(Z, D) \neq 0 \quad (12)$$

$$\text{Cov}(Z, \varepsilon) = 0, \quad (13)$$

then standard IV identifies  $\beta$ , at least in large samples,

$$\text{plim } \hat{\beta}_{\text{IV}} = \frac{\text{Cov}(Z, Y)}{\text{Cov}(Z, D)} = \beta.$$

- If other instruments exist, each identifies  $\beta$ .
- $Z$  produces a controlled variation in  $D$  relative to  $\varepsilon$ .
- Randomization of assignment with full compliance to experimental protocols is an example of an instrument.
- From the instrumental variable estimators, we can identify the effect of adopting the policy in any country since all countries respond to the policy in the same way controlling for their  $X$ .

- If  $\beta (= Y_1 - Y_0)$  varies in the population even after controlling for  $X$ , there is a distribution of responses that cannot in general be summarized by a single number.



- If  $\beta (= Y_1 - Y_0)$  varies in the population even after controlling for  $X$ , there is a distribution of responses that cannot in general be summarized by a single number.
- Even if we are interested in the mean of the distribution, a new phenomenon distinct from selection bias might arise.

- If  $\beta (= Y_1 - Y_0)$  varies in the population even after controlling for  $X$ , there is a distribution of responses that cannot in general be summarized by a single number.
- Even if we are interested in the mean of the distribution, a new phenomenon distinct from selection bias might arise.
- This is a problem of sorting on the gain, which is distinct from sorting on levels.

- If  $\beta (= Y_1 - Y_0)$  varies in the population even after controlling for  $X$ , there is a distribution of responses that cannot in general be summarized by a single number.
- Even if we are interested in the mean of the distribution, a new phenomenon distinct from selection bias might arise.
- This is a problem of sorting on the gain, which is distinct from sorting on levels.
- If  $\beta$  varies, even after controlling for  $X$ , there may be sorting on the gain ( $\text{Cov}(\beta, D) \neq 0$ ).

- If  $\beta (= Y_1 - Y_0)$  varies in the population even after controlling for  $X$ , there is a distribution of responses that cannot in general be summarized by a single number.
- Even if we are interested in the mean of the distribution, a new phenomenon distinct from selection bias might arise.
- This is a problem of sorting on the gain, which is distinct from sorting on levels.
- If  $\beta$  varies, even after controlling for  $X$ , there may be sorting on the gain ( $\text{Cov}(\beta, D) \neq 0$ ).
- This is the model of **essential heterogeneity** as defined by ?.

- If  $\beta (= Y_1 - Y_0)$  varies in the population even after controlling for  $X$ , there is a distribution of responses that cannot in general be summarized by a single number.
- Even if we are interested in the mean of the distribution, a new phenomenon distinct from selection bias might arise.
- This is a problem of sorting on the gain, which is distinct from sorting on levels.
- If  $\beta$  varies, even after controlling for  $X$ , there may be sorting on the gain ( $\text{Cov}(\beta, D) \neq 0$ ).
- This is the model of **essential heterogeneity** as defined by ?.
- It is also called a correlated random coefficient model (?).

- The application of instrumental variables to this case is more problematic.

- The application of instrumental variables to this case is more problematic.
- Suppose that we augment the standard instrumental variable assumptions (12) and (13) by the following assumption:

$$\text{Cov}(Z, \beta) = 0. \quad (14)$$

- The application of instrumental variables to this case is more problematic.
- Suppose that we augment the standard instrumental variable assumptions (12) and (13) by the following assumption:

$$\text{Cov}(Z, \beta) = 0. \quad (14)$$

- Can we identify the mean of  $(Y_1 - Y_0)$  using IV?



- The application of instrumental variables to this case is more problematic.
- Suppose that we augment the standard instrumental variable assumptions (12) and (13) by the following assumption:

$$\text{Cov}(Z, \beta) = 0. \quad (14)$$

- Can we identify the mean of  $(Y_1 - Y_0)$  using IV?
- In general we cannot.

- To see why, let  $\bar{\beta} = (\mu_1 - \mu_0)$  be the mean treatment effect (the mean of the distribution in figure 1).

- To see why, let  $\bar{\beta} = (\mu_1 - \mu_0)$  be the mean treatment effect (the mean of the distribution in figure 1).
- $\beta = \bar{\beta} + \eta$ , where  $U_1 - U_0 = \eta$  and  $\bar{\beta} = \mu_1 - \mu_0$  and we keep the conditioning on  $X$  implicit.

- To see why, let  $\bar{\beta} = (\mu_1 - \mu_0)$  be the mean treatment effect (the mean of the distribution in figure 1).
- $\beta = \bar{\beta} + \eta$ , where  $U_1 - U_0 = \eta$  and  $\bar{\beta} = \mu_1 - \mu_0$  and we keep the conditioning on  $X$  implicit.
- Write equation (11) in terms of these parameters:

$$Y = \alpha + \bar{\beta}D + [\varepsilon + \eta D].$$

The error term of this equation ( $\varepsilon + \eta D$ ) contains two components.

- To see why, let  $\bar{\beta} = (\mu_1 - \mu_0)$  be the mean treatment effect (the mean of the distribution in figure 1).
- $\beta = \bar{\beta} + \eta$ , where  $U_1 - U_0 = \eta$  and  $\bar{\beta} = \mu_1 - \mu_0$  and we keep the conditioning on  $X$  implicit.
- Write equation (11) in terms of these parameters:

$$Y = \alpha + \bar{\beta}D + [\varepsilon + \eta D].$$

The error term of this equation ( $\varepsilon + \eta D$ ) contains two components.

- By assumption,  $Z$  is uncorrelated with  $\varepsilon$  and  $\eta$ .

- To see why, let  $\bar{\beta} = (\mu_1 - \mu_0)$  be the mean treatment effect (the mean of the distribution in figure 1).
- $\beta = \bar{\beta} + \eta$ , where  $U_1 - U_0 = \eta$  and  $\bar{\beta} = \mu_1 - \mu_0$  and we keep the conditioning on  $X$  implicit.
- Write equation (11) in terms of these parameters:

$$Y = \alpha + \bar{\beta}D + [\varepsilon + \eta D].$$

The error term of this equation ( $\varepsilon + \eta D$ ) contains two components.

- By assumption,  $Z$  is uncorrelated with  $\varepsilon$  and  $\eta$ .
- But to identify  $\bar{\beta}$ , we need IV to be uncorrelated with  $[\varepsilon + \eta D]$ .

- To see why, let  $\bar{\beta} = (\mu_1 - \mu_0)$  be the mean treatment effect (the mean of the distribution in figure 1).
- $\beta = \bar{\beta} + \eta$ , where  $U_1 - U_0 = \eta$  and  $\bar{\beta} = \mu_1 - \mu_0$  and we keep the conditioning on  $X$  implicit.
- Write equation (11) in terms of these parameters:

$$Y = \alpha + \bar{\beta}D + [\varepsilon + \eta D].$$

The error term of this equation ( $\varepsilon + \eta D$ ) contains two components.

- By assumption,  $Z$  is uncorrelated with  $\varepsilon$  and  $\eta$ .
- But to identify  $\bar{\beta}$ , we need IV to be uncorrelated with  $[\varepsilon + \eta D]$ .
- That requires  $Z$  to be uncorrelated with  $\eta D$ .

- If policy adoption is made without knowledge of  $\eta$  ( $= U_1 - U_0$ ), the idiosyncratic gain to policy adoption after controlling for the observables, then  $\eta$  and  $D$  are statistically independent and hence uncorrelated, and IV identifies  $\bar{\beta}$ .



- If policy adoption is made without knowledge of  $\eta$  ( $= U_1 - U_0$ ), the idiosyncratic gain to policy adoption after controlling for the observables, then  $\eta$  and  $D$  are statistically independent and hence uncorrelated, and IV identifies  $\bar{\beta}$ .
- If, however, policy adoption is made with partial or full knowledge of  $\eta$ , IV does not identify  $\bar{\beta}$  because  $E(\eta D | Z) = E(\eta | D = 1, Z) \Pr(D = 1 | Z)$  and if there is sorting on the unobserved gain  $\eta$ , the first term is not zero.

- If policy adoption is made without knowledge of  $\eta$  ( $= U_1 - U_0$ ), the idiosyncratic gain to policy adoption after controlling for the observables, then  $\eta$  and  $D$  are statistically independent and hence uncorrelated, and IV identifies  $\bar{\beta}$ .
- If, however, policy adoption is made with partial or full knowledge of  $\eta$ , IV does not identify  $\bar{\beta}$  because  $E(\eta D | Z) = E(\eta | D = 1, Z) \Pr(D = 1 | Z)$  and if there is sorting on the unobserved gain  $\eta$ , the first term is not zero.
- Similar calculations show that IV does not identify the mean gain to the countries that adopt the policy ( $E(\beta | D = 1)$ ) and many other summary treatment parameters.

- Whether  $\eta (= U_1 - U_0)$  is correlated with  $D$  depends on the quality of the data available to the empirical economist and cannot be settled *a priori*.

- Whether  $\eta (= U_1 - U_0)$  is correlated with  $D$  depends on the quality of the data available to the empirical economist and cannot be settled *a priori*.
- The conservative position is to allow for such correlation.

- Whether  $\eta (= U_1 - U_0)$  is correlated with  $D$  depends on the quality of the data available to the empirical economist and cannot be settled *a priori*.
- The conservative position is to allow for such correlation.
- However, this rules out IV as an interesting econometric strategy for identifying any of the familiar mean treatment parameters.

- In light of the negative conclusions about IV in the literature preceding their paper, it is remarkable that Heckman and Ichimura (1995) establish that under certain conditions, in the model with essential heterogeneity, IV can identify an interpretable parameter.

- In light of the negative conclusions about IV in the literature preceding their paper, it is remarkable that Angriston and Pischke establish that under certain conditions, in the model with essential heterogeneity, IV can identify an interpretable parameter.
- The parameter they identify is a discrete approximation to the marginal gain parameter introduced by Angriston and Pischke.

- In light of the negative conclusions about IV in the literature preceding their paper, it is remarkable that Angrast-Kotlikoff and Pischke establish that under certain conditions, in the model with essential heterogeneity, IV can identify an interpretable parameter.
- The parameter they identify is a discrete approximation to the marginal gain parameter introduced by Angrast-Kotlikoff.
- The Angrast-Kotlikoff-Moffitt parameter is a version of MTE for a parametric normal selection model.



- We derive their parameter from a selection model in Slide 338.

- We derive their parameter from a selection model in Slide 338.
- ? demonstrate how to use a selection model to identify the marginal gain to persons induced into a treatment status by a marginal change in the cost of treatment.

- We derive their parameter from a selection model in Slide 338.
- ? demonstrate how to use a selection model to identify the marginal gain to persons induced into a treatment status by a marginal change in the cost of treatment.
- ? show how to estimate a discrete approximation to the Björklund-Moffitt parameter using instrumental variables.

- ? assume the existence of an instrument  $Z$  that takes two or more distinct values.

- ? assume the existence of an instrument  $Z$  that takes two or more distinct values.
- This is implicit in (12).

- ? assume the existence of an instrument  $Z$  that takes two or more distinct values.
- This is implicit in (12).
- If  $Z$  assumes only one value, the covariance in (12) would be zero.

- ? assume the existence of an instrument  $Z$  that takes two or more distinct values.
- This is implicit in (12).
- If  $Z$  assumes only one value, the covariance in (12) would be zero.
- Strengthening the covariance conditions of equations (13) and (14), they assume (IV-1) and (IV-2) (independence and rank respectively) and that  $Z$  is independent of  $\beta = (Y_1 - Y_0)$  and  $Y_0$ .

- ? assume the existence of an instrument  $Z$  that takes two or more distinct values.
- This is implicit in (12).
- If  $Z$  assumes only one value, the covariance in (12) would be zero.
- Strengthening the covariance conditions of equations (13) and (14), they assume (IV-1) and (IV-2) (independence and rank respectively) and that  $Z$  is independent of  $\beta = (Y_1 - Y_0)$  and  $Y_0$ .
- Recall that we denote by  $D(z)$  the random variable indicating receipt of treatment when  $Z$  is set to  $z$ .



- ? assume the existence of an instrument  $Z$  that takes two or more distinct values.
- This is implicit in (12).
- If  $Z$  assumes only one value, the covariance in (12) would be zero.
- Strengthening the covariance conditions of equations (13) and (14), they assume (IV-1) and (IV-2) (independence and rank respectively) and that  $Z$  is independent of  $\beta = (Y_1 - Y_0)$  and  $Y_0$ .
- Recall that we denote by  $D(z)$  the random variable indicating receipt of treatment when  $Z$  is set to  $z$ .
- ( $D(z) = 1$  if treatment is received;  $D(z) = 0$  otherwise).

- ? assume the existence of an instrument  $Z$  that takes two or more distinct values.
- This is implicit in (12).
- If  $Z$  assumes only one value, the covariance in (12) would be zero.
- Strengthening the covariance conditions of equations (13) and (14), they assume (IV-1) and (IV-2) (independence and rank respectively) and that  $Z$  is independent of  $\beta = (Y_1 - Y_0)$  and  $Y_0$ .
- Recall that we denote by  $D(z)$  the random variable indicating receipt of treatment when  $Z$  is set to  $z$ .
- ( $D(z) = 1$  if treatment is received;  $D(z) = 0$  otherwise).
- The Imbens-Angrist independence and rank assumptions are (IV-1) and (IV-2).

- They supplement the standard IV assumptions with what they call a “monotonicity” assumption.

- They supplement the standard IV assumptions with what they call a “monotonicity” assumption.
- It is a condition across persons.

- They supplement the standard IV assumptions with what they call a “monotonicity” assumption.
- It is a condition across persons.
- The assumption maintains that if  $Z$  is fixed first at one and then at the other of two distinct values, say  $Z = z$  and  $Z = z'$ , then all persons respond in their choice of  $D$  to the change in  $Z$  in the same way.

- They supplement the standard IV assumptions with what they call a “monotonicity” assumption.
- It is a condition across persons.
- The assumption maintains that if  $Z$  is fixed first at one and then at the other of two distinct values, say  $Z = z$  and  $Z = z'$ , then all persons respond in their choice of  $D$  to the change in  $Z$  in the same way.
- In our policy adoption example, this condition states that a movement from  $z$  to  $z'$ , causes all countries to move toward (or against) adoption of the public policy being studied.

- They supplement the standard IV assumptions with what they call a “monotonicity” assumption.
- It is a condition across persons.
- The assumption maintains that if  $Z$  is fixed first at one and then at the other of two distinct values, say  $Z = z$  and  $Z = z'$ , then all persons respond in their choice of  $D$  to the change in  $Z$  in the same way.
- In our policy adoption example, this condition states that a movement from  $z$  to  $z'$ , causes all countries to move toward (or against) adoption of the public policy being studied.
- If some adopt, others do not drop the policy in response to the same change.

- More formally, letting  $D_i(z)$  be the indicator (= 1 if adopted; = 0 if not) for adoption of a policy if  $Z = z$  for country  $i$ , then for any distinct values  $z$  and  $z'$  ? assume:

(IV-3)

$D_i(z) \geq D_i(z')$  for all  $i$ , or  $D_i(z) \leq D_i(z')$  for all  $i = 1, \dots, I$ .  
**(Monotonicity or Uniformity)**



- The content in this assumption is not in the order for any person.

- The content in this assumption is not in the order for any person.
- Rather, the responses have to be uniform across people for a given choice of  $z$  and  $z'$ .

- The content in this assumption is not in the order for any person.
- Rather, the responses have to be uniform across people for a given choice of  $z$  and  $z'$ .
- One possibility allowed under (IV-3) is the existence of three values of  $z < z' < z''$  such that for all  $i$ ,  $D_i(z) \geq D_i(z')$  but  $D_i(z') \leq D_i(z'')$ . The standard usage of the term monotonicity rules out this possibility by requiring that one of the following hold for all  $i$ : (a)  $z < z'$  componentwise implies  $D_i(z) \geq D_i(z')$  or (b)  $z < z'$  componentwise implies  $D_i(z) \leq D_i(z')$ .

- The content in this assumption is not in the order for any person.
- Rather, the responses have to be uniform across people for a given choice of  $z$  and  $z'$ .
- One possibility allowed under (IV-3) is the existence of three values of  $z < z' < z''$  such that for all  $i$ ,  $D_i(z) \geq D_i(z')$  but  $D_i(z') \leq D_i(z'')$ . The standard usage of the term monotonicity rules out this possibility by requiring that one of the following hold for all  $i$ : (a)  $z < z'$  componentwise implies  $D_i(z) \geq D_i(z')$  or (b)  $z < z'$  componentwise implies  $D_i(z) \leq D_i(z')$ .
- Of course, if the  $D_i(z)$  are monotonic in  $Z$  in the same direction for all  $i$ , they are monotonic in the sense of Imbens and Angrist.

- For any value of  $z'$  in the domain of definition of  $Z$ , from (IV-1) and (IV-2) and the definition of  $D(z)$ ,  $(Y_0, Y_1, D(z'))$  is independent of  $Z$ .

- For any value of  $z'$  in the domain of definition of  $Z$ , from (IV-1) and (IV-2) and the definition of  $D(z)$ ,  $(Y_0, Y_1, D(z'))$  is independent of  $Z$ .
- For any two values of the instrument  $Z = z$  and  $Z = z'$ , we may write

$$\begin{aligned} & E(Y | Z = z) - E(Y | Z = z') \\ &= E(Y_1 D + Y_0 (1 - D) | Z = z) - E(Y_1 D + Y_0 (1 - D) | Z = z') \\ &= E(Y_0 + D(Y_1 - Y_0) | Z = z) - E(Y_0 + D(Y_1 - Y_0) | Z = z'). \end{aligned}$$

- For any value of  $z'$  in the domain of definition of  $Z$ , from (IV-1) and (IV-2) and the definition of  $D(z)$ ,  $(Y_0, Y_1, D(z'))$  is independent of  $Z$ .
- For any two values of the instrument  $Z = z$  and  $Z = z'$ , we may write

$$\begin{aligned} & E(Y | Z = z) - E(Y | Z = z') \\ &= E(Y_1 D + Y_0 (1 - D) | Z = z) - E(Y_1 D + Y_0 (1 - D) | Z = z') \\ &= E(Y_0 + D(Y_1 - Y_0) | Z = z) - E(Y_0 + D(Y_1 - Y_0) | Z = z'). \end{aligned}$$

- From the independence condition (IV-1), and the definition of  $D(z)$  and  $D(z')$ , we may write this expression as  $E[(Y_1 - Y_0)(D(z) - D(z'))]$ .

- Using the law of iterated expectations,

$$\begin{aligned} E(Y | Z = z) - E(Y | Z = z') & \quad (15) \\ &= E(Y_1 - Y_0 | D(z) - D(z') = 1) \Pr(D(z) - D(z') = 1) \\ &\quad - E(Y_1 - Y_0 | D(z) - D(z') = -1) \Pr(D(z) - D(z') = -1). \end{aligned}$$



- Using the law of iterated expectations,

$$\begin{aligned} E(Y | Z = z) - E(Y | Z = z') & \\ &= E(Y_1 - Y_0 | D(z) - D(z') = 1) \Pr(D(z) - D(z') = 1) \\ &\quad - E(Y_1 - Y_0 | D(z) - D(z') = -1) \Pr(D(z) - D(z') = -1). \end{aligned} \tag{15}$$

- By the monotonicity condition (IV-3), we eliminate one or the other term in the final expression.

- Using the law of iterated expectations,

$$\begin{aligned} E(Y | Z = z) - E(Y | Z = z') & \\ &= E(Y_1 - Y_0 | D(z) - D(z') = 1) \Pr(D(z) - D(z') = 1) \\ &\quad - E(Y_1 - Y_0 | D(z) - D(z') = -1) \Pr(D(z) - D(z') = -1). \end{aligned} \tag{15}$$

- By the monotonicity condition (IV-3), we eliminate one or the other term in the final expression.
- Suppose that  $\Pr(D(z) - D(z') = -1) = 0$ , then

$$\begin{aligned} E(Y | Z = z) - E(Y | Z = z') & \\ &= E(Y_1 - Y_0 | D(z) - D(z') = 1) \Pr(D(z) - D(z') = 1). \end{aligned}$$

- Using the law of iterated expectations,

$$\begin{aligned}
 E(Y | Z = z) - E(Y | Z = z') & \qquad \qquad \qquad (15) \\
 &= E(Y_1 - Y_0 | D(z) - D(z') = 1) \Pr(D(z) - D(z') = 1) \\
 &\quad - E(Y_1 - Y_0 | D(z) - D(z') = -1) \Pr(D(z) - D(z') = -1).
 \end{aligned}$$

- By the monotonicity condition (IV-3), we eliminate one or the other term in the final expression.
- Suppose that  $\Pr(D(z) - D(z') = -1) = 0$ , then

$$\begin{aligned}
 E(Y | Z = z) - E(Y | Z = z') \\
 &= E(Y_1 - Y_0 | D(z) - D(z') = 1) \Pr(D(z) - D(z') = 1).
 \end{aligned}$$

- Observe that, by monotonicity,  $\Pr(D(z) - D(z') = 1) = \Pr(D = 1 | Z = z) - \Pr(D = 1 | Z = z')$ .

- For values of  $z$  and  $z'$  that produce distinct propensity scores  $\Pr(D = 1 | Z = z)$ , using monotonicity once more, we obtain LATE:

$$\begin{aligned} \text{LATE} &= \frac{E(Y | Z = z) - E(Y | Z = z')}{\Pr(D = 1 | Z = z) - \Pr(D = 1 | Z = z')} \\ &= E(Y_1 - Y_0 | D(z) - D(z') = 1). \end{aligned} \quad (16)$$

- For values of  $z$  and  $z'$  that produce distinct propensity scores  $\Pr(D = 1 | Z = z)$ , using monotonicity once more, we obtain LATE:

$$\begin{aligned} \text{LATE} &= \frac{E(Y | Z = z) - E(Y | Z = z')}{\Pr(D = 1 | Z = z) - \Pr(D = 1 | Z = z')} \\ &= E(Y_1 - Y_0 | D(z) - D(z') = 1). \end{aligned} \quad (16)$$

- This is the mean gain to those induced to switch from “0” to “1” by a change in  $Z$  from  $z'$  to  $z$ .

- This is not the mean of  $Y_1 - Y_0$  (average treatment effect) unless the  $Z$  assume values  $(z, z')$  such that  $\Pr(D(z) = 1) = 1$  and  $\Pr(D(z') = 1) = 0$ .

- This is not the mean of  $Y_1 - Y_0$  (average treatment effect) unless the  $Z$  assume values  $(z, z')$  such that  $\Pr(D(z) = 1) = 1$  and  $\Pr(D(z') = 1) = 0$ .
- It is also not the effect of treatment on the treated ( $E(Y_1 - Y_0 | D = 1) = E(\beta | D = 1)$ ) unless the analyst has access to one or more values of  $Z$  such that  $\Pr(D(z) = 1) = 1$ .

- The LATE parameter is defined by a hypothetical manipulation of instruments.



- The LATE parameter is defined by a hypothetical manipulation of instruments.
- It depends on the particular instrument used.

- The LATE parameter is defined by a hypothetical manipulation of instruments.
- It depends on the particular instrument used.
- If monotonicity (uniformity) is violated, IV estimates an average response of those induced to switch into the program and those induced to switch out of the program by the change in the instrument because both terms in (15) are present.

- In an application to wage equations, ?? interprets the LATE estimator as identifying returns to marginal persons.

- In an application to wage equations, **??** interprets the LATE estimator as identifying returns to marginal persons.
- **?** notes that the actual margin of choice selected by the IV estimator is not identified by the instrument.

- In an application to wage equations, Angrist and Pischke interprets the LATE estimator as identifying returns to marginal persons.
- Angrist notes that the actual margin of choice selected by the IV estimator is not identified by the instrument.
- It is unclear as to which segment of the population the return estimated by LATE applies.

- If the analyst is interested in knowing the average response ( $\bar{\beta}$ ), the effect of the policy on the outcomes of countries that adopt it ( $E(\beta | D = 1)$ ) or the effect of the policy if a particular country adopts it, there is no guarantee that the IV estimator comes any closer to the desired target than the OLS estimator and indeed it may be more biased than OLS.

- If the analyst is interested in knowing the average response ( $\bar{\beta}$ ), the effect of the policy on the outcomes of countries that adopt it ( $E(\beta | D = 1)$ ) or the effect of the policy if a particular country adopts it, there is no guarantee that the IV estimator comes any closer to the desired target than the OLS estimator and indeed it may be more biased than OLS.
- Because different instruments define different parameters, having a wealth of different strong instruments does not improve the precision of the estimate of any particular parameter.

- If the analyst is interested in knowing the average response ( $\bar{\beta}$ ), the effect of the policy on the outcomes of countries that adopt it ( $E(\beta | D = 1)$ ) or the effect of the policy if a particular country adopts it, there is no guarantee that the IV estimator comes any closer to the desired target than the OLS estimator and indeed it may be more biased than OLS.
- Because different instruments define different parameters, having a wealth of different strong instruments does not improve the precision of the estimate of any particular parameter.
- This is in stark contrast with the traditional model with  $\beta \perp\!\!\!\perp D$ .



- If the analyst is interested in knowing the average response ( $\bar{\beta}$ ), the effect of the policy on the outcomes of countries that adopt it ( $E(\beta | D = 1)$ ) or the effect of the policy if a particular country adopts it, there is no guarantee that the IV estimator comes any closer to the desired target than the OLS estimator and indeed it may be more biased than OLS.
- Because different instruments define different parameters, having a wealth of different strong instruments does not improve the precision of the estimate of any particular parameter.
- This is in stark contrast with the traditional model with  $\beta \perp\!\!\!\perp D$ .
- In that case, all valid instruments identify  $\bar{\beta}$ .

- The  $\chi^2$  test for the validity of extra instruments applies to the traditional model.

- The  $H_0$  test for the validity of extra instruments applies to the traditional model.
- In the more general case with essential heterogeneity, because different instruments estimate different parameters, no clear inference emerges from such specification tests.

- When there are more than two distinct values of  $Z$ , Imbens and Angrist draw on the analysis of  $\tau$ , which was refined in  $\tau$  and  $\tau$ , to produce a weighted average of pairwise LATE parameters where the scalars  $Z$  are ordered to define the LATE parameter.

- When there are more than two distinct values of  $Z$ , Imbens and Angrist draw on the analysis of  $\tau$ , which was refined in  $\tau$  and  $\tau$ , to produce a weighted average of pairwise LATE parameters where the scalars  $Z$  are ordered to define the LATE parameter.
- In this case,  $IV$  is a weighted average of LATE parameters with non-negative weights.

- When there are more than two distinct values of  $Z$ , Imbens and Angrist draw on the analysis of Angrist and Pischke (2009), which was refined in Angrist and Pischke (2015) and Angrist and Pischke (2015), to produce a weighted average of pairwise LATE parameters where the scalars  $Z$  are ordered to define the LATE parameter.
- In this case,  $IV$  is a weighted average of LATE parameters with non-negative weights.
- Imbens and Angrist generalize this result to the case of vector  $Z$  assuming that instruments are monotonic functions of the probability of selection.

- $\tau_{\text{ATE}}$ ,  $\tau_{\text{LATE}}$  and  $\tau_{\text{HDE}}$  generalize the analysis of  $\tau_{\text{ATE}}$  in several ways and we report their results in this chapter.

- $\tau$  ,  $\tau$  and  $\tau$  generalize the analysis of  $\tau$  in several ways and we report their results in this chapter.
- Using a choice-theoretic parameter (the marginal treatment effect or MTE) introduced into the literature on selection models by  $\tau$ , they relate the parameters estimated by IV to well formulated choice models.



- $\tau$ ,  $\tau$  and  $\tau$  generalize the analysis of  $\tau$  in several ways and we report their results in this chapter.
- Using a choice-theoretic parameter (the marginal treatment effect or MTE) introduced into the literature on selection models by  $\tau$ , they relate the parameters estimated by IV to well formulated choice models.
- This allows treatment parameters to be defined independent of any values assumed by instruments.

- It is possible to generate all treatment effects as different weighted averages of the MTE.

- It is possible to generate all treatment effects as different weighted averages of the MTE.
- IV can also be interpreted as a weighted average of MTE.

- It is possible to generate all treatment effects as different weighted averages of the MTE.
- IV can also be interpreted as a weighted average of MTE.
- Different instruments weight different segments of the MTE differently.

- It is possible to generate all treatment effects as different weighted averages of the MTE.
- IV can also be interpreted as a weighted average of MTE.
- Different instruments weight different segments of the MTE differently.
- Using the nonparametric generalized Roy model, MTE is a limit form of LATE.

- Using MTE, we overcome a problem that plagues the LATE literature.

- Using MTE, we overcome a problem that plagues the LATE literature.
- LATE estimates marginal returns at an unidentified margin (or intervals of margins).

- Using MTE, we overcome a problem that plagues the LATE literature.
- LATE estimates marginal returns at an unidentified margin (or intervals of margins).
- We show how to use the MTE to unify diverse instrumental variables estimates and to determine what margins (or intervals of margins) they identify.



- Using MTE, we overcome a problem that plagues the LATE literature.
- LATE estimates marginal returns at an unidentified margin (or intervals of margins).
- We show how to use the MTE to unify diverse instrumental variables estimates and to determine what margins (or intervals of margins) they identify.
- Instead of reporting a marginal return for unidentified persons, we show how to report marginal returns for all persons identified by their location on the scale of a latent variable that arises from a well defined choice model and is related to the propensity of persons to make the choice being studied.

- Using MTE, we overcome a problem that plagues the LATE literature.
- LATE estimates marginal returns at an unidentified margin (or intervals of margins).
- We show how to use the MTE to unify diverse instrumental variables estimates and to determine what margins (or intervals of margins) they identify.
- Instead of reporting a marginal return for unidentified persons, we show how to report marginal returns for all persons identified by their location on the scale of a latent variable that arises from a well defined choice model and is related to the propensity of persons to make the choice being studied.
- We can interpret the margins of choice identified by various instruments and place diverse instruments on a common interpretive footing.

- ?? establish the central role of the propensity score  $(\Pr(D = 1 | Z = z) = P(z))$  in both selection and IV models.

- ?? establish the central role of the propensity score ( $\Pr(D = 1 \mid Z = z) = P(z)$ ) in both selection and IV models.
- They show that with vector  $Z$  and a scalar instrument  $J(Z)$  constructed from vector  $Z$ , the weights on LATE and MTE that are implicit in standard IV are not guaranteed to be non-negative.

- ?? establish the central role of the propensity score ( $\Pr(D = 1 | Z = z) = P(z)$ ) in both selection and IV models.
- They show that with vector  $Z$  and a scalar instrument  $J(Z)$  constructed from vector  $Z$ , the weights on LATE and MTE that are implicit in standard IV are not guaranteed to be non-negative.
- Thus IV can be negative even though all pairwise LATEs and pointwise MTEs are positive.

- ?? establish the central role of the propensity score ( $\Pr(D = 1 \mid Z = z) = P(z)$ ) in both selection and IV models.
- They show that with vector  $Z$  and a scalar instrument  $J(Z)$  constructed from vector  $Z$ , the weights on LATE and MTE that are implicit in standard IV are not guaranteed to be non-negative.
- Thus IV can be negative even though all pairwise LATEs and pointwise MTEs are positive.
- Thus the treatment effects for any pair of  $(z, z')$  can be positive but the IV can be negative.

- We present examples below.

- We present examples below.
- Certain instruments produce positive weights and avoid this particular interpretive problem.



- We present examples below.
- Certain instruments produce positive weights and avoid this particular interpretive problem.
- Our analysis generalizes the analyses of weights on treatment effects by Yitzhaki and Imbens-Angrist, who analyze a special case where all weights are positive.

- We establish the special status of  $P(z)$  as an instrument.

- We establish the special status of  $P(z)$  as an instrument.
- It always produces non-negative weights for MTE and LATE.

- We establish the special status of  $P(z)$  as an instrument.
- It always produces non-negative weights for MTE and LATE.
- It enables analysts to identify MTE or LATE.

- We establish the special status of  $P(z)$  as an instrument.
- It always produces non-negative weights for MTE and LATE.
- It enables analysts to identify MTE or LATE.
- With knowledge of  $P(z)$ , and the MTE or LATE, we can decompose any IV estimate into identifiable MTEs (at points) or LATEs (over intervals) and identifiable weights on MTE (or LATE) where the weights can be constructed from data.

- We establish the special status of  $P(z)$  as an instrument.
- It always produces non-negative weights for MTE and LATE.
- It enables analysts to identify MTE or LATE.
- With knowledge of  $P(z)$ , and the MTE or LATE, we can decompose any IV estimate into identifiable MTEs (at points) or LATEs (over intervals) and identifiable weights on MTE (or LATE) where the weights can be constructed from data.
- The ability to decompose IV into interpretable components allows analysts to determine the response to treatment of persons at different levels of unobserved factors that determine treatment status.

- We present a simple test for essential heterogeneity ( $\beta$  dependent on  $D$ ) that allows analysts to determine whether or not they can avoid the complexities of the more general model with heterogeneity in response to treatments.

- We present a simple test for essential heterogeneity ( $\beta$  dependent on  $D$ ) that allows analysts to determine whether or not they can avoid the complexities of the more general model with heterogeneity in response to treatments.
- In Slide 471, we generalize the analysis of IV in the two-outcome model to a multiple outcome model, analyzing both ordered and unordered choice cases.



- We present a simple test for essential heterogeneity ( $\beta$  dependent on  $D$ ) that allows analysts to determine whether or not they can avoid the complexities of the more general model with heterogeneity in response to treatments.
- In Slide 471, we generalize the analysis of IV in the two-outcome model to a multiple outcome model, analyzing both ordered and unordered choice cases.
- We also demonstrate the fundamental asymmetry in the recent IV literature for models with heterogeneous outcomes.

- We present a simple test for essential heterogeneity ( $\beta$  dependent on  $D$ ) that allows analysts to determine whether or not they can avoid the complexities of the more general model with heterogeneity in response to treatments.
- In Slide 471, we generalize the analysis of IV in the two-outcome model to a multiple outcome model, analyzing both ordered and unordered choice cases.
- We also demonstrate the fundamental asymmetry in the recent IV literature for models with heterogeneous outcomes.
- Responses to treatment are permitted to be heterogeneous in a general way.

- Responses of choices to instruments are not.

- Responses of choices to instruments are not.
- When heterogeneity in choice is allowed for in a general way, IV and local IV do not estimate parameters that can be interpreted as weighted averages of MTEs or LATEs.

- Responses of choices to instruments are not.
- When heterogeneity in choice is allowed for in a general way, IV and local IV do not estimate parameters that can be interpreted as weighted averages of MTEs or LATEs.
- We now turn to an analysis of the two-outcome model.

## IV in Choice Models

- A key contribution of the analysis of Heckman and Vytlacil is to adjoin choice equation (7) to the outcome equations (1), (5) and (6).

## IV in Choice Models

- A key contribution of the analysis of Heckman and Vytlacil is to adjoin choice equation (7) to the outcome equations (1), (5) and (6).
- A standard binary threshold cross model for  $D$  is  $D = \mathbf{1}(D^* \geq 0)$ , where  $\mathbf{1}(\cdot)$  is an indicator ( $\mathbf{1}(A) = 1$  if  $A$  is true, 0 otherwise).

## IV in Choice Models

- A key contribution of the analysis of Heckman and Vytlacil is to adjoin choice equation (7) to the outcome equations (1), (5) and (6).
- A standard binary threshold cross model for  $D$  is  $D = \mathbf{1}(D^* \geq 0)$ , where  $\mathbf{1}(\cdot)$  is an indicator ( $\mathbf{1}(A) = 1$  if  $A$  is true, 0 otherwise).
- A familiar version of (7) sets  $\mu_D(Z) = Z\gamma$  and writes

$$D^* = Z\gamma - V, \quad (17)$$

where  $(V \perp\!\!\!\perp Z) \mid X$ .



## IV in Choice Models

- A key contribution of the analysis of Heckman and Vytlacil is to adjoin choice equation (7) to the outcome equations (1), (5) and (6).
- A standard binary threshold cross model for  $D$  is  $D = \mathbf{1}(D^* \geq 0)$ , where  $\mathbf{1}(\cdot)$  is an indicator ( $\mathbf{1}(A) = 1$  if  $A$  is true, 0 otherwise).
- A familiar version of (7) sets  $\mu_D(Z) = Z\gamma$  and writes

$$D^* = Z\gamma - V, \quad (17)$$

where  $(V \perp\!\!\!\perp Z) \mid X$ .

- ( $V$  is independent of  $Z$  given  $X$ ).

- In this notation, the propensity score or choice probability is

$$P(z) = \Pr(D = 1 \mid Z = z) = \Pr(Z\gamma \geq V) = F_V(Z\gamma)$$

where  $F_V$  is the distribution of  $V$  which is assumed to be continuous.

- In this notation, the propensity score or choice probability is

$$P(z) = \Pr(D = 1 \mid Z = z) = \Pr(Z\gamma \geq V) = F_V(Z\gamma)$$

where  $F_V$  is the distribution of  $V$  which is assumed to be continuous.

- In terms of the generalized Roy model where  $C$  is the cost of participation in sector 1,  $D = \mathbf{1}[Y_1 - Y_0 - C > 0]$ .

- In this notation, the propensity score or choice probability is

$$P(z) = \Pr(D = 1 \mid Z = z) = \Pr(Z\gamma \geq V) = F_V(Z\gamma)$$

where  $F_V$  is the distribution of  $V$  which is assumed to be continuous.

- In terms of the generalized Roy model where  $C$  is the cost of participation in sector 1,  $D = \mathbf{1}[Y_1 - Y_0 - C > 0]$ .
- For a separable model in outcomes and in costs,

$$C = \mu_D(W) + U_C,$$

we have  $Z = (X, W)$ ,  $\mu_D(Z) = \mu_1(X) - \mu_0(X) - \mu_D(W)$ , and  $V = -(U_1 - U_0 - U_C)$ .

- In constructing many of our examples, we work with a special version where  $U_C = 0$ .

- In constructing many of our examples, we work with a special version where  $U_C = 0$ .
- We call this version the extended Roy model.

- In constructing many of our examples, we work with a special version where  $U_C = 0$ .
- We call this version the extended Roy model.
- It is the model used to produce figure 1.

- In constructing many of our examples, we work with a special version where  $U_C = 0$ .
- We call this version the extended Roy model.
- It is the model used to produce figure 1.
- Our analysis, however, applies to more general models, and we also offer examples of generalized Roy models, as we have in figure 2 and table 3.



- In the case where  $\beta$  (given  $X$ ) is a constant, under (IV-1) and (IV-2) it is not necessary to specify the choice model to identify  $\beta$ .

- In the case where  $\beta$  (given  $X$ ) is a constant, under (IV-1) and (IV-2) it is not necessary to specify the choice model to identify  $\beta$ .
- In a general model with heterogenous responses, the specification of  $P(z)$  and its relationship with the instrument play crucial roles.

- In the case where  $\beta$  (given  $X$ ) is a constant, under (IV-1) and (IV-2) it is not necessary to specify the choice model to identify  $\beta$ .
- In a general model with heterogenous responses, the specification of  $P(z)$  and its relationship with the instrument play crucial roles.
- To see this, study the covariance between  $Z$  and  $\eta D$  discussed in the introduction to this section.

- In the case where  $\beta$  (given  $X$ ) is a constant, under (IV-1) and (IV-2) it is not necessary to specify the choice model to identify  $\beta$ .
- In a general model with heterogenous responses, the specification of  $P(z)$  and its relationship with the instrument play crucial roles.
- To see this, study the covariance between  $Z$  and  $\eta D$  discussed in the introduction to this section.
- By the law of iterated expectations, letting  $\bar{Z}$  denote the mean of  $Z$ ,

$$\begin{aligned}\text{Cov}(Z, \eta D) &= E((Z - \bar{Z}) D \eta) \\ &= E((Z - \bar{Z}) \eta \mid D = 1) \Pr(D = 1) \\ &= E((Z - \bar{Z}) \eta \mid Z\gamma > V) \Pr(Z\gamma \geq V).\end{aligned}$$

- Thus, even if  $Z$  and  $\eta$  are independent, they are not independent conditional on  $D = \mathbf{1}[Z\gamma \geq V]$  if  $\eta = (U_1 - U_0)$  is dependent on  $V$  (i.e., if the decision maker has partial knowledge of  $\eta$  and acts on it).

- Thus, even if  $Z$  and  $\eta$  are independent, they are not independent conditional on  $D = \mathbf{1}[Z\gamma \geq V]$  if  $\eta = (U_1 - U_0)$  is dependent on  $V$  (i.e., if the decision maker has partial knowledge of  $\eta$  and acts on it).
- Selection models allow for this dependence (see ???; and ?).

- Keeping  $X$  implicit, assuming that

$$(U_1, U_0, V) \perp\!\!\!\perp Z \quad (18)$$

(alternatively, assuming that  $(\varepsilon, \eta) \perp\!\!\!\perp Z$ ), we obtain

$$\begin{aligned} E(Y | D = 0, Z = z) &= E(Y_0 | D = 0, Z = z) \\ &= \alpha + E(U_0 | z\gamma < V), \end{aligned}$$

where  $\alpha$  and possibly  $E(U_0 | z\gamma < V)$  depend on  $X$ , which can be written as

$$E(Y | D = 0, Z = z) = \alpha + K_0(P(z))$$

where the functional form of  $K_0$  is produced from the distribution of  $(U_0, V)$ .

- Keeping  $X$  implicit, assuming that

$$(U_1, U_0, V) \perp\!\!\!\perp Z \quad (18)$$

(alternatively, assuming that  $(\varepsilon, \eta) \perp\!\!\!\perp Z$ ), we obtain

$$\begin{aligned} E(Y \mid D = 0, Z = z) &= E(Y_0 \mid D = 0, Z = z) \\ &= \alpha + E(U_0 \mid z\gamma < V), \end{aligned}$$

where  $\alpha$  and possibly  $E(U_0 \mid z\gamma < V)$  depend on  $X$ , which can be written as

$$E(Y \mid D = 0, Z = z) = \alpha + K_0(P(z))$$

where the functional form of  $K_0$  is produced from the distribution of  $(U_0, V)$ .

- Focusing on means, the conventional selection approach models the conditional mean dependence between  $(U_0, U_1)$  and  $V$ .



- Similarly,

$$\begin{aligned} E(Y | D = 1, Z = z) &= E(Y_1 | D = 1, Z = z) \\ &= \alpha + \bar{\beta} + E(U_1 | z\gamma \geq V) \\ &= \alpha + \bar{\beta} + K_1(P(z)) \end{aligned}$$

where  $\alpha, \bar{\beta}$  and  $K_1(P(z))$  may depend on  $X$ .

- Similarly,

$$\begin{aligned} E(Y | D = 1, Z = z) &= E(Y_1 | D = 1, Z = z) \\ &= \alpha + \bar{\beta} + E(U_1 | z\gamma \geq V) \\ &= \alpha + \bar{\beta} + K_1(P(z)) \end{aligned}$$

where  $\alpha, \bar{\beta}$  and  $K_1(P(z))$  may depend on  $X$ .

- $K_0(P(z))$  and  $K_1(P(z))$  are control functions in the sense of ??.

- Similarly,

$$\begin{aligned} E(Y | D = 1, Z = z) &= E(Y_1 | D = 1, Z = z) \\ &= \alpha + \bar{\beta} + E(U_1 | z\gamma \geq V) \\ &= \alpha + \bar{\beta} + K_1(P(z)) \end{aligned}$$

where  $\alpha, \bar{\beta}$  and  $K_1(P(z))$  may depend on  $X$ .

- $K_0(P(z))$  and  $K_1(P(z))$  are control functions in the sense of ??.
- The control functions expect out the unobservables  $\theta$  that give rise to selection bias (see (U-1)).

- Similarly,

$$\begin{aligned} E(Y | D = 1, Z = z) &= E(Y_1 | D = 1, Z = z) \\ &= \alpha + \bar{\beta} + E(U_1 | z\gamma \geq V) \\ &= \alpha + \bar{\beta} + K_1(P(z)) \end{aligned}$$

where  $\alpha, \bar{\beta}$  and  $K_1(P(z))$  may depend on  $X$ .

- $K_0(P(z))$  and  $K_1(P(z))$  are control functions in the sense of ??.
- The control functions expect out the unobservables  $\theta$  that give rise to selection bias (see (U-1)).
- Under standard conditions developed in the literature, analysts can identify  $\bar{\beta}$ .

- Similarly,

$$\begin{aligned}
 E(Y | D = 1, Z = z) &= E(Y_1 | D = 1, Z = z) \\
 &= \alpha + \bar{\beta} + E(U_1 | z\gamma \geq V) \\
 &= \alpha + \bar{\beta} + K_1(P(z))
 \end{aligned}$$

where  $\alpha, \bar{\beta}$  and  $K_1(P(z))$  may depend on  $X$ .

- $K_0(P(z))$  and  $K_1(P(z))$  are control functions in the sense of ??.
- The control functions expect out the unobservables  $\theta$  that give rise to selection bias (see (U-1)).
- Under standard conditions developed in the literature, analysts can identify  $\bar{\beta}$ .
- ? discusses semiparametric identification.

- Because we condition on  $Z = z$  (or  $P(z)$ ), correct specification of the  $Z$  plays an important role in econometric selection methods.

- Because we condition on  $Z = z$  (or  $P(z)$ ), correct specification of the  $Z$  plays an important role in econometric selection methods.
- This sensitivity to the full set of instruments in  $Z$  appears to be absent from the IV method.

- If  $\beta$  is a constant (given  $X$ ), or if  $\eta (= \beta - \bar{\beta})$  is independent of  $V$ , only one instrument from vector  $Z$  needs to be used to identify the parameter.



- If  $\beta$  is a constant (given  $X$ ), or if  $\eta (= \beta - \bar{\beta})$  is independent of  $V$ , only one instrument from vector  $Z$  needs to be used to identify the parameter.
- Missing or unused instruments play no role in identifying mean responses but may affect the efficiency of the IV estimators.

- If  $\beta$  is a constant (given  $X$ ), or if  $\eta (= \beta - \bar{\beta})$  is independent of  $V$ , only one instrument from vector  $Z$  needs to be used to identify the parameter.
- Missing or unused instruments play no role in identifying mean responses but may affect the efficiency of the IV estimators.
- In a model where  $\beta$  is variable and not independent of  $V$ , misspecification of  $Z$  plays an important role in interpreting what IV estimates analogous to its role in selection models.

- If  $\beta$  is a constant (given  $X$ ), or if  $\eta (= \beta - \bar{\beta})$  is independent of  $V$ , only one instrument from vector  $Z$  needs to be used to identify the parameter.
- Missing or unused instruments play no role in identifying mean responses but may affect the efficiency of the IV estimators.
- In a model where  $\beta$  is variable and not independent of  $V$ , misspecification of  $Z$  plays an important role in interpreting what IV estimates analogous to its role in selection models.
- Misspecification of  $Z$  affects both approaches to identification.

- If  $\beta$  is a constant (given  $X$ ), or if  $\eta (= \beta - \bar{\beta})$  is independent of  $V$ , only one instrument from vector  $Z$  needs to be used to identify the parameter.
- Missing or unused instruments play no role in identifying mean responses but may affect the efficiency of the IV estimators.
- In a model where  $\beta$  is variable and not independent of  $V$ , misspecification of  $Z$  plays an important role in interpreting what IV estimates analogous to its role in selection models.
- Misspecification of  $Z$  affects both approaches to identification.
- This is a new phenomenon in models with heterogenous  $\beta$ .

- If  $\beta$  is a constant (given  $X$ ), or if  $\eta (= \beta - \bar{\beta})$  is independent of  $V$ , only one instrument from vector  $Z$  needs to be used to identify the parameter.
- Missing or unused instruments play no role in identifying mean responses but may affect the efficiency of the IV estimators.
- In a model where  $\beta$  is variable and not independent of  $V$ , misspecification of  $Z$  plays an important role in interpreting what IV estimates analogous to its role in selection models.
- Misspecification of  $Z$  affects both approaches to identification.
- This is a new phenomenon in models with heterogenous  $\beta$ .
- We now review results from the recent literature on instrumental variables in the model with essential heterogeneity.

## Instrumental Variables and Local Instrumental Variables

- In this section, we use  $\Delta^{\text{MTE}}$  defined in Slide 90 for a general nonseparable model (5)–(7) to organize the literature on econometric evaluation estimators.

## Instrumental Variables and Local Instrumental Variables

- In this section, we use  $\Delta^{\text{MTE}}$  defined in Slide 90 for a general nonseparable model (5)–(7) to organize the literature on econometric evaluation estimators.
- In terms of our simple regression model,

$$\begin{aligned}
 \Delta^{\text{MTE}}(x, u_D) &= E(\Delta \mid X = x, U_D = u_D) \\
 &= E(\beta \mid X = x, V = F_V^{-1}(u_D)) \\
 &= \bar{\beta}(x) + E(\eta \mid X = x, V = v).
 \end{aligned}$$

where  $v = F_V^{-1}(u_D)$ .

## Instrumental Variables and Local Instrumental Variables

- In this section, we use  $\Delta^{\text{MTE}}$  defined in Slide 90 for a general nonseparable model (5)–(7) to organize the literature on econometric evaluation estimators.
- In terms of our simple regression model,

$$\begin{aligned}\Delta^{\text{MTE}}(x, u_D) &= E(\Delta \mid X = x, U_D = u_D) \\ &= E(\beta \mid X = x, V = F_V^{-1}(u_D)) \\ &= \bar{\beta}(x) + E(\eta \mid X = x, V = v).\end{aligned}$$

where  $v = F_V^{-1}(u_D)$ .

- We assume policy invariance in the sense of Hurwicz for mean parameters (Assumption (A-7)).



## Instrumental Variables and Local Instrumental Variables

- In this section, we use  $\Delta^{\text{MTE}}$  defined in Slide 90 for a general nonseparable model (5)–(7) to organize the literature on econometric evaluation estimators.
- In terms of our simple regression model,

$$\begin{aligned}\Delta^{\text{MTE}}(x, u_D) &= E(\Delta \mid X = x, U_D = u_D) \\ &= E(\beta \mid X = x, V = F_V^{-1}(u_D)) \\ &= \bar{\beta}(x) + E(\eta \mid X = x, V = v).\end{aligned}$$

where  $v = F_V^{-1}(u_D)$ .

- We assume policy invariance in the sense of Hurwicz for mean parameters (Assumption (A-7)).
- For simplicity, we suppress the  $a$  and  $a'$  subscripts that indicate specific policies.

- We focus primarily on instrumental variable estimators and review the method of local instrumental variables.

- We focus primarily on instrumental variable estimators and review the method of local instrumental variables.
- Slide 184 demonstrated in a simple but familiar case that well established intuitions about instrumental variable identification strategies break down when  $\Delta^{\text{MTE}}$  is nonconstant in  $u_D$  given  $X$  ( $\beta \not\perp D \mid X$ ).

- We focus primarily on instrumental variable estimators and review the method of local instrumental variables.
- Slide 184 demonstrated in a simple but familiar case that well established intuitions about instrumental variable identification strategies break down when  $\Delta^{\text{MTE}}$  is nonconstant in  $u_D$  given  $X$  ( $\beta \not\perp D \mid X$ ).
- We acquire the probability of selection  $P(z)$  as a determinant of the IV covariance relationships.

- Two sets of instrumental variable conditions are presented in the current literature for this more general case: those associated with conventional instrumental variable assumptions, which are implied by the assumption of “no selection on heterogeneous gains,” ( $\beta \perp\!\!\!\perp D \mid X$ ) and those which permit selection on heterogeneous gains.

- Two sets of instrumental variable conditions are presented in the current literature for this more general case: those associated with conventional instrumental variable assumptions, which are implied by the assumption of “no selection on heterogeneous gains,” ( $\beta \perp\!\!\!\perp D \mid X$ ) and those which permit selection on heterogeneous gains.
- Neither set of assumptions implies the other, nor does either identify the policy relevant treatment effect or other economically interpretable parameters in the general case.

- Two sets of instrumental variable conditions are presented in the current literature for this more general case: those associated with conventional instrumental variable assumptions, which are implied by the assumption of “no selection on heterogeneous gains,” ( $\beta \perp\!\!\!\perp D \mid X$ ) and those which permit selection on heterogeneous gains.
- Neither set of assumptions implies the other, nor does either identify the policy relevant treatment effect or other economically interpretable parameters in the general case.
- Each set of conditions identifies different treatment parameters.

- In place of standard instrumental variables methods, ??? advocate a new approach to estimating policy impacts by estimating  $\Delta^{\text{MTE}}$  using local instrumental variables (LIV) to identify all of the treatment parameters from a generator  $\Delta^{\text{MTE}}$  that can be weighted in different ways to answer different policy questions.



- In place of standard instrumental variables methods, ??? advocate a new approach to estimating policy impacts by estimating  $\Delta^{\text{MTE}}$  using local instrumental variables (LIV) to identify all of the treatment parameters from a generator  $\Delta^{\text{MTE}}$  that can be weighted in different ways to answer different policy questions.
- For certain classes of policy interventions covered by Assumption (A-7) and analyzed in Slide 412,  $\Delta^{\text{MTE}}$  possesses an invariance property analogous to the invariant parameters of traditional structural econometrics.

## Conditions on the MTE that Justify the Application of Conventional Instrumental Variables

- In the general case where  $\Delta^{\text{MTE}}(x, u_D)$  is nonconstant in  $u_D$  ( $E(\beta | X = x, V = v)$  depends on  $v$ ), IV does not in general estimate any of the treatment effects defined in Slide 90.

## Conditions on the MTE that Justify the Application of Conventional Instrumental Variables

- In the general case where  $\Delta^{\text{MTE}}(x, u_D)$  is nonconstant in  $u_D$  ( $E(\beta \mid X = x, V = v)$  depends on  $v$ ), IV does not in general estimate any of the treatment effects defined in Slide 90.
- We consider a scalar instrument  $J(Z)$  constructed from  $Z$  which may be vector valued.

## Conditions on the MTE that Justify the Application of Conventional Instrumental Variables

- In the general case where  $\Delta^{\text{MTE}}(x, u_D)$  is nonconstant in  $u_D$  ( $E(\beta \mid X = x, V = v)$  depends on  $v$ ), IV does not in general estimate any of the treatment effects defined in Slide 90.
- We consider a scalar instrument  $J(Z)$  constructed from  $Z$  which may be vector valued.
- We sometimes denote  $J(Z)$  by  $J$ , leaving implicit that  $J$  is a function of  $Z$ .

## Conditions on the MTE that Justify the Application of Conventional Instrumental Variables

- In the general case where  $\Delta^{\text{MTE}}(x, u_D)$  is nonconstant in  $u_D$  ( $E(\beta \mid X = x, V = v)$  depends on  $v$ ), IV does not in general estimate any of the treatment effects defined in Slide 90.
- We consider a scalar instrument  $J(Z)$  constructed from  $Z$  which may be vector valued.
- We sometimes denote  $J(Z)$  by  $J$ , leaving implicit that  $J$  is a function of  $Z$ .
- If  $Z$  is a vector,  $J(Z)$  can be one coordinate of  $Z$ , say  $Z_1$ .

## Conditions on the MTE that Justify the Application of Conventional Instrumental Variables

- In the general case where  $\Delta^{\text{MTE}}(x, u_D)$  is nonconstant in  $u_D$  ( $E(\beta \mid X = x, V = v)$  depends on  $v$ ), IV does not in general estimate any of the treatment effects defined in Slide 90.
- We consider a scalar instrument  $J(Z)$  constructed from  $Z$  which may be vector valued.
- We sometimes denote  $J(Z)$  by  $J$ , leaving implicit that  $J$  is a function of  $Z$ .
- If  $Z$  is a vector,  $J(Z)$  can be one coordinate of  $Z$ , say  $Z_1$ .
- We develop this particular case in presenting our examples.

- The notation is sufficiently general to make  $J(Z)$  a general function of  $Z$ .

- The notation is sufficiently general to make  $J(Z)$  a general function of  $Z$ .
- The standard conditions  $J(Z) \perp\!\!\!\perp (U_0, U_1) \mid X$  and  $\text{Cov}(J(Z), D \mid X) \neq 0$  corresponding to (IV-1) and (IV-2) respectively, do not, by themselves, imply that instrumental variables using  $J(Z)$  as the instrument will identify conventional or policy relevant treatment effects.



- The notation is sufficiently general to make  $J(Z)$  a general function of  $Z$ .
- The standard conditions  $J(Z) \perp\!\!\!\perp (U_0, U_1) \mid X$  and  $\text{Cov}(J(Z), D \mid X) \neq 0$  corresponding to (IV-1) and (IV-2) respectively, do not, by themselves, imply that instrumental variables using  $J(Z)$  as the instrument will identify conventional or policy relevant treatment effects.
- When responses to treatment are heterogenous, we must supplement the standard conditions to identify interpretable parameters.

- The notation is sufficiently general to make  $J(Z)$  a general function of  $Z$ .
- The standard conditions  $J(Z) \perp\!\!\!\perp (U_0, U_1) \mid X$  and  $\text{Cov}(J(Z), D \mid X) \neq 0$  corresponding to (IV-1) and (IV-2) respectively, do not, by themselves, imply that instrumental variables using  $J(Z)$  as the instrument will identify conventional or policy relevant treatment effects.
- When responses to treatment are heterogenous, we must supplement the standard conditions to identify interpretable parameters.
- To link our analysis to conventional analyses of IV, we continue to invoke familiar-looking representations of additive separability of outcomes in terms of  $(U_0, U_1)$  so we invoke (2).

- The notation is sufficiently general to make  $J(Z)$  a general function of  $Z$ .
- The standard conditions  $J(Z) \perp\!\!\!\perp (U_0, U_1) \mid X$  and  $\text{Cov}(J(Z), D \mid X) \neq 0$  corresponding to (IV-1) and (IV-2) respectively, do not, by themselves, imply that instrumental variables using  $J(Z)$  as the instrument will identify conventional or policy relevant treatment effects.
- When responses to treatment are heterogenous, we must supplement the standard conditions to identify interpretable parameters.
- To link our analysis to conventional analyses of IV, we continue to invoke familiar-looking representations of additive separability of outcomes in terms of  $(U_0, U_1)$  so we invoke (2).
- This is not required.

- All derivations and results in this subsection hold without assuming additive separability if  $\mu_1(x)$  and  $\mu_0(x)$  are replaced by  $E(Y_1 | X = x)$  and  $E(Y_0 | X = x)$ , respectively, and  $U_1$  and  $U_0$  are replaced by  $Y_1 - E(Y_1 | X)$  and  $Y_0 - E(Y_0 | X)$ , respectively.

- All derivations and results in this subsection hold without assuming additive separability if  $\mu_1(x)$  and  $\mu_0(x)$  are replaced by  $E(Y_1 | X = x)$  and  $E(Y_0 | X = x)$ , respectively, and  $U_1$  and  $U_0$  are replaced by  $Y_1 - E(Y_1 | X)$  and  $Y_0 - E(Y_0 | X)$ , respectively.
- This highlights the point that all of our analysis of IV is conditional on  $X$  and  $X$  need not be exogenous with respect to  $(U_0, U_1)$  to identify the MTE conditional on  $X$ .

- All derivations and results in this subsection hold without assuming additive separability if  $\mu_1(x)$  and  $\mu_0(x)$  are replaced by  $E(Y_1 | X = x)$  and  $E(Y_0 | X = x)$ , respectively, and  $U_1$  and  $U_0$  are replaced by  $Y_1 - E(Y_1 | X)$  and  $Y_0 - E(Y_0 | X)$ , respectively.
- This highlights the point that all of our analysis of IV is conditional on  $X$  and  $X$  need not be exogenous with respect to  $(U_0, U_1)$  to identify the MTE conditional on  $X$ .
- To simplify the notation, we keep the conditioning on  $X$  implicit unless it is useful to break it out separately.

- Two distinct sets of instrumental variable conditions in the literature are those due to  $??$  and  $?$ , and those due to  $?$  which we previously discussed.

- Two distinct sets of instrumental variable conditions in the literature are those due to  $??$  and  $?$ , and those due to  $?$  which we previously discussed.
- We review the conditions of  $??$  and  $?$  in Appendix, Slide 1163, which is presented in the context of our discussion of matching in Slide 675, where we compare IV and matching.



- Two distinct sets of instrumental variable conditions in the literature are those due to  $??$  and  $?$ , and those due to  $?$  which we previously discussed.
- We review the conditions of  $??$  and  $?$  in Appendix, Slide 1163, which is presented in the context of our discussion of matching in Slide 675, where we compare IV and matching.
- In the case where  $\Delta^{\text{MTE}}$  is nonconstant in  $u_D$ , standard IV estimates different parameters depending on which assumptions are maintained.

- Two distinct sets of instrumental variable conditions in the literature are those due to  $??$  and  $?$ , and those due to  $?$  which we previously discussed.
- We review the conditions of  $??$  and  $?$  in Appendix, Slide 1163, which is presented in the context of our discussion of matching in Slide 675, where we compare IV and matching.
- In the case where  $\Delta^{\text{MTE}}$  is nonconstant in  $u_D$ , standard IV estimates different parameters depending on which assumptions are maintained.
- We have already shown that when responses to treatment are heterogeneous, and choices are made on the basis of this heterogeneity, standard IV does not identify  $\mu_1 - \mu_0 = \bar{\beta}$ .

- There are two important cases of the variable response model.

(C-1)

$D \perp\!\!\!\perp \Delta \implies E(\Delta \mid U_D) = E(\Delta)$ ,  $\Delta^{MTE}(u_D)$  is constant in  $u_D$  and  $\Delta^{MTE} = \Delta^{ATE} = \Delta^{TT} = \Delta^{LATE}$ , i.e.,  $E(\beta \mid D = 1) = E(\beta)$ , because  $\beta \perp\!\!\!\perp D$ .

- There are two important cases of the variable response model.
- The first case arises when responses are heterogeneous, but conditional on  $X$ , people do not base their participation on these responses.

(C-2)

$D \perp\!\!\!\perp \Delta \implies E(\Delta \mid U_D) = E(\Delta)$ ,  $\Delta^{MTE}(u_D)$  is constant in  $u_D$  and  $\Delta^{MTE} = \Delta^{ATE} = \Delta^{TT} = \Delta^{LATE}$ , i.e.,  $E(\beta \mid D = 1) = E(\beta)$ , because  $\beta \perp\!\!\!\perp D$ .

- There are two important cases of the variable response model.
- The first case arises when responses are heterogeneous, but conditional on  $X$ , people do not base their participation on these responses.
- In this case, keeping the conditioning on  $X$  implicit,

(C-3)

$D \perp\!\!\!\perp \Delta \implies E(\Delta \mid U_D) = E(\Delta)$ ,  $\Delta^{MTE}(u_D)$  is constant in  $u_D$  and  $\Delta^{MTE} = \Delta^{ATE} = \Delta^{TT} = \Delta^{LATE}$ , i.e.,  $E(\beta \mid D = 1) = E(\beta)$ , because  $\beta \perp\!\!\!\perp D$ .

- In this case, all mean treatment parameters are the same.

(C-4)

$D \perp\!\!\!\perp \Delta$  and  $E(\Delta | U_D) \neq E(\Delta)$  (i.e.,  $\beta \not\perp\!\!\!\perp D$ ).

- In this case, all mean treatment parameters are the same.
- The second case arises when selection into treatment depends on  $\beta$ :

(C-5)

$D \not\perp \Delta$  and  $E(\Delta | U_D) \neq E(\Delta)$  (i.e.,  $\beta \not\perp D$ ).

- In this case,  $\Delta^{\text{MTE}}$  is nonconstant, and in general, the treatment parameters differ among each other.



- In this case,  $\Delta^{\text{MTE}}$  is nonconstant, and in general, the treatment parameters differ among each other.
- In this case (IV-1) and (IV-2) for general instruments do not identify  $\bar{\beta}$  (as shown in Slide 184) or  $E(\beta \mid D = 1)$ .

- A sufficient condition that generates (C-1) is the information condition that decisions to participate in the program are not made on the basis of  $U_1 - U_0 (= \eta)$  (in the notation of Slide 184):

(I-1)

$$\Pr(D = 1 \mid Z, U_1 - U_0) = \Pr(D = 1 \mid Z) \text{ (i.e.,} \\ \Pr(D = 1 \mid Z, \beta) = \Pr(D = 1 \mid Z)).$$

- Before we investigate what standard instrumental variables estimators identify, we first present the local instrumental variables estimator which directly estimates the MTE.

- Before we investigate what standard instrumental variables estimators identify, we first present the local instrumental variables estimator which directly estimates the MTE.
- It is a limit form of LATE.

## Estimating the MTE Using Local Instrumental Variables

- ??? develop the Local Instrumental Variable (LIV) estimator to recover  $\Delta^{\text{MTE}}$  pointwise.

## Estimating the MTE Using Local Instrumental Variables

- ??? develop the Local Instrumental Variable (LIV) estimator to recover  $\Delta^{\text{MTE}}$  pointwise.
- LIV is the derivative of the conditional expectation of  $Y$  with respect to  $P(Z) = p$ .

## Estimating the MTE Using Local Instrumental Variables

- ??? develop the Local Instrumental Variable (LIV) estimator to recover  $\Delta^{\text{MTE}}$  pointwise.
- LIV is the derivative of the conditional expectation of  $Y$  with respect to  $P(Z) = p$ .
- This is defined as

$$\Delta^{\text{LIV}}(p) \equiv \frac{\partial E(Y | P(Z) = p)}{\partial p}. \quad (19)$$

## Estimating the MTE Using Local Instrumental Variables

- ??? develop the Local Instrumental Variable (LIV) estimator to recover  $\Delta^{\text{MTE}}$  pointwise.
- LIV is the derivative of the conditional expectation of  $Y$  with respect to  $P(Z) = p$ .
- This is defined as

$$\Delta^{\text{LIV}}(p) \equiv \frac{\partial E(Y | P(Z) = p)}{\partial p}. \quad (19)$$

- It is the population mean response to a policy change embodied in changes in  $P(Z)$  analyzed by ?.



- $E(Y | P(Z))$  is well-defined as a consequence of assumption (A-4), and  $E(Y | P(Z))$  can be recovered over the support of  $P(Z)$ .

- $E(Y | P(Z))$  is well-defined as a consequence of assumption (A-4), and  $E(Y | P(Z))$  can be recovered over the support of  $P(Z)$ .
- Under our assumptions, LIV identifies MTE at all points of continuity in  $P(Z)$  (conditional on  $X$ ).

- $E(Y | P(Z))$  is well-defined as a consequence of assumption (A-4), and  $E(Y | P(Z))$  can be recovered over the support of  $P(Z)$ .
- Under our assumptions, LIV identifies MTE at all points of continuity in  $P(Z)$  (conditional on  $X$ ).
- This expression does not require additive separability of  $\mu_1(X, U_1)$  or  $\mu_0(X, U_0)$ .

- Under standard regularity conditions, a variety of nonparametric methods can be used to estimate the derivative of  $E(Y | P(Z))$  and thus to estimate  $\Delta^{\text{MTE}}$ .

- Under standard regularity conditions, a variety of nonparametric methods can be used to estimate the derivative of  $E(Y | P(Z))$  and thus to estimate  $\Delta^{\text{MTE}}$ .
- With  $\Delta^{\text{MTE}}$  in hand, if the support of the distribution of  $P(Z)$  conditional on  $X$  is the full unit interval, one can generate all the treatment parameters defined in Slide 90 as well as the policy relevant treatment parameter presented in Slide 139 as weighted versions of  $\Delta^{\text{MTE}}$ .

- Under standard regularity conditions, a variety of nonparametric methods can be used to estimate the derivative of  $E(Y | P(Z))$  and thus to estimate  $\Delta^{\text{MTE}}$ .
- With  $\Delta^{\text{MTE}}$  in hand, if the support of the distribution of  $P(Z)$  conditional on  $X$  is the full unit interval, one can generate all the treatment parameters defined in Slide 90 as well as the policy relevant treatment parameter presented in Slide 139 as weighted versions of  $\Delta^{\text{MTE}}$ .
- When the support of the distribution of  $P(Z)$  conditional on  $X$  is not full, it is still possible to identify some parameters.

- Under standard regularity conditions, a variety of nonparametric methods can be used to estimate the derivative of  $E(Y | P(Z))$  and thus to estimate  $\Delta^{\text{MTE}}$ .
- With  $\Delta^{\text{MTE}}$  in hand, if the support of the distribution of  $P(Z)$  conditional on  $X$  is the full unit interval, one can generate all the treatment parameters defined in Slide 90 as well as the policy relevant treatment parameter presented in Slide 139 as weighted versions of  $\Delta^{\text{MTE}}$ .
- When the support of the distribution of  $P(Z)$  conditional on  $X$  is not full, it is still possible to identify some parameters.
- ? show that to identify ATE under Assumptions (A-1)–(A-5), it is necessary and sufficient that the support of the distribution of  $P(Z)$  include 0 and 1.

- Thus, identification of ATE does not require that the distribution of  $P(Z)$  be the full unit interval or that the distribution of  $P(Z)$  be continuous.



- Thus, identification of ATE does not require that the distribution of  $P(Z)$  be the full unit interval or that the distribution of  $P(Z)$  be continuous.
- But the support must include  $\{0, 1\}$ .

- Thus, identification of ATE does not require that the distribution of  $P(Z)$  be the full unit interval or that the distribution of  $P(Z)$  be continuous.
- But the support must include  $\{0, 1\}$ .
- Sharp bounds on the treatment parameters can be constructed under the same assumptions imposed in this chapter without imposing full support conditions.

- Thus, identification of ATE does not require that the distribution of  $P(Z)$  be the full unit interval or that the distribution of  $P(Z)$  be continuous.
- But the support must include  $\{0, 1\}$ .
- Sharp bounds on the treatment parameters can be constructed under the same assumptions imposed in this chapter without imposing full support conditions.
- The resulting bounds are simple and easy to apply compared with those presented in the previous literature.

- Thus, identification of ATE does not require that the distribution of  $P(Z)$  be the full unit interval or that the distribution of  $P(Z)$  be continuous.
- But the support must include  $\{0, 1\}$ .
- Sharp bounds on the treatment parameters can be constructed under the same assumptions imposed in this chapter without imposing full support conditions.
- The resulting bounds are simple and easy to apply compared with those presented in the previous literature.
- We discuss these and other bounds in Slide 938.

- To establish the relationship between LIV and ordinary IV based on  $P(Z)$  and to motivate how LIV identifies  $\Delta^{\text{MTE}}$ , notice that from the definition of  $Y$ , the conditional expectation of  $Y$  given  $P(Z)$  is, recalling that  $\Delta = Y_1 - Y_0$ ,

$$\begin{aligned} E(Y | P(Z) = p) \\ = E(Y_0 | P(Z) = p) + E(\Delta | P(Z) = p, D = 1)p, \end{aligned}$$

where we keep the conditioning on  $X$  implicit.

- To establish the relationship between LIV and ordinary IV based on  $P(Z)$  and to motivate how LIV identifies  $\Delta^{\text{MTE}}$ , notice that from the definition of  $Y$ , the conditional expectation of  $Y$  given  $P(Z)$  is, recalling that  $\Delta = Y_1 - Y_0$ ,

$$\begin{aligned} E(Y | P(Z) = p) \\ = E(Y_0 | P(Z) = p) + E(\Delta | P(Z) = p, D = 1)p, \end{aligned}$$

where we keep the conditioning on  $X$  implicit.

- Our model and conditional independence assumption (A-1) imply

$$E(Y | P(Z) = p) = E(Y_0) + E(\Delta | p \geq U_D)p.$$

- Applying the IV (Wald) estimator for two different values of  $P(Z)$ ,  $p$  and  $p'$ , for  $p \neq p'$ , we obtain:

$$\begin{aligned} & \frac{E(Y | P(Z) = p) - E(Y | P(Z) = p')}{p - p'} \quad (20) \\ = & \Delta^{\text{ATE}} + \frac{E(U_1 - U_0 | p \geq U_D)p - E(U_1 - U_0 | p' \geq U_D)p'}{p - p'}, \end{aligned}$$

where this particular expression is obtained under the assumption of additive separability in the outcomes.

- Applying the IV (Wald) estimator for two different values of  $P(Z)$ ,  $p$  and  $p'$ , for  $p \neq p'$ , we obtain:

$$\begin{aligned} & \frac{E(Y | P(Z) = p) - E(Y | P(Z) = p')}{p - p'} \\ &= \Delta^{\text{ATE}} + \frac{E(U_1 - U_0 | p \geq U_D)p - E(U_1 - U_0 | p' \geq U_D)p'}{p - p'}, \end{aligned} \quad (20)$$

where this particular expression is obtained under the assumption of additive separability in the outcomes.

- Exactly the same equation holds without additive separability if one replaces  $U_1$  and  $U_0$  with  $Y_1 - E(Y_1|X)$  and  $Y_0 - E(Y_0|X)$ .



- When  $U_1 \equiv U_0$  or  $(U_1 - U_0) \perp\!\!\!\perp U_D$ , (case (C-1)), IV based on  $P(Z)$  estimates  $\Delta^{\text{ATE}}$  because the second term on the right hand side of the expression (20) vanishes.

- When  $U_1 \equiv U_0$  or  $(U_1 - U_0) \perp\!\!\!\perp U_D$ , (case (C-1)), IV based on  $P(Z)$  estimates  $\Delta^{\text{ATE}}$  because the second term on the right hand side of the expression (20) vanishes.
- Otherwise, IV estimates a combination of MTE parameters which we analyze further below.

- Assuming additive separability of the outcome equations, another representation of  $E(Y|P(Z) = p)$  reveals the index structure.

- Assuming additive separability of the outcome equations, another representation of  $E(Y|P(Z) = p)$  reveals the index structure.
- It writes (keeping the conditioning on  $X$  implicit) that

$$E(Y|P(Z) = p) = E(Y_0) + \Delta^{\text{ATE}} p + \int_0^p E(U_1 - U_0|U_D = u_D) du_D. \quad (21)$$

- Assuming additive separability of the outcome equations, another representation of  $E(Y|P(Z) = p)$  reveals the index structure.
- It writes (keeping the conditioning on  $X$  implicit) that

$$E(Y|P(Z) = p) = E(Y_0) + \Delta^{\text{ATE}}p + \int_0^p E(U_1 - U_0|U_D = u_D)du_D. \quad (21)$$

- We can differentiate with respect to  $p$  and use LIV to identify  $\Delta^{\text{MTE}}$ :

$$\Delta^{\text{MTE}}(p) = \frac{\partial E(Y|P(Z) = p)}{\partial p} = \Delta^{\text{ATE}} + E(U_1 - U_0|U_D = p).$$

- Notice that IV estimates  $\Delta^{\text{ATE}}$  when  $E(Y | P(Z) = p)$  is a linear function of  $p$  so the third term on the right hand side of (21) vanishes.

- Notice that IV estimates  $\Delta^{\text{ATE}}$  when  $E(Y | P(Z) = p)$  is a linear function of  $p$  so the third term on the right hand side of (21) vanishes.
- Thus a test of the linearity of  $E(Y | P(Z) = p)$  in  $p$  is a test of the validity of linear IV for  $\Delta^{\text{ATE}}$ , i.e., it is a test of whether or not the data are consistent with a correlated random coefficient model ( $\beta \not\perp D$ ).

- Notice that IV estimates  $\Delta^{\text{ATE}}$  when  $E(Y | P(Z) = p)$  is a linear function of  $p$  so the third term on the right hand side of (21) vanishes.
- Thus a test of the linearity of  $E(Y | P(Z) = p)$  in  $p$  is a test of the validity of linear IV for  $\Delta^{\text{ATE}}$ , i.e., it is a test of whether or not the data are consistent with a correlated random coefficient model ( $\beta \not\perp D$ ).
- The nonlinearity of  $E(Y | P(Z) = p)$  in  $p$  provides a way to distinguish whether Case (C-1) or Case (C-2) describes the data.



- Notice that IV estimates  $\Delta^{\text{ATE}}$  when  $E(Y | P(Z) = p)$  is a linear function of  $p$  so the third term on the right hand side of (21) vanishes.
- Thus a test of the linearity of  $E(Y | P(Z) = p)$  in  $p$  is a test of the validity of linear IV for  $\Delta^{\text{ATE}}$ , i.e., it is a test of whether or not the data are consistent with a correlated random coefficient model ( $\beta \not\perp D$ ).
- The nonlinearity of  $E(Y | P(Z) = p)$  in  $p$  provides a way to distinguish whether Case (C-1) or Case (C-2) describes the data.
- It is also a test of whether or not agents can at least partially anticipate future unobserved (by the econometrician) gains (the  $Y_1 - Y_0$  given  $X$ ) at the time they make their participation decisions.

- The levels and derivatives of  $E(Y | P(Z) = p)$  and standard errors can be estimated using a variety of semiparametric methods.

- The levels and derivatives of  $E(Y | P(Z) = p)$  and standard errors can be estimated using a variety of semiparametric methods.
- ? present an algorithm for estimating  $\Delta^{\text{MTE}}$  using local linear regression.

- This analysis generalizes to the nonseparable outcomes case.

- This analysis generalizes to the nonseparable outcomes case.
- We use separability in outcomes only to simplify the exposition and link to more traditional models.

- This analysis generalizes to the nonseparable outcomes case.
- We use separability in outcomes only to simplify the exposition and link to more traditional models.
- In particular, exactly the same expression holds with exactly the same derivation for the nonseparable case if we replace  $U_1$  and  $U_0$  with  $Y_1 - E(Y_1|X)$  and  $Y_0 - E(Y_0|X)$ , respectively.

- This analysis generalizes to the nonseparable outcomes case.
- We use separability in outcomes only to simplify the exposition and link to more traditional models.
- In particular, exactly the same expression holds with exactly the same derivation for the nonseparable case if we replace  $U_1$  and  $U_0$  with  $Y_1 - E(Y_1|X)$  and  $Y_0 - E(Y_0|X)$ , respectively.
- This simple test for the absence of general heterogeneity based on linearity of  $E(Y|Z)$  in  $P(Z)$  applies to the case of LATE for any pair of instruments.

- This analysis generalizes to the nonseparable outcomes case.
- We use separability in outcomes only to simplify the exposition and link to more traditional models.
- In particular, exactly the same expression holds with exactly the same derivation for the nonseparable case if we replace  $U_1$  and  $U_0$  with  $Y_1 - E(Y_1|X)$  and  $Y_0 - E(Y_0|X)$ , respectively.
- This simple test for the absence of general heterogeneity based on linearity of  $E(Y|Z)$  in  $P(Z)$  applies to the case of LATE for any pair of instruments.
- An equivalent way is to check that all pairwise LATEs are the same over the sample support of  $Z$ .



- Figure 3A plots two cases of  $E(Y | P(Z) = p)$  based on the generalized Roy model used to generate the example in figure 2A and 2B.

- Figure 3A plots two cases of  $E(Y | P(Z) = p)$  based on the generalized Roy model used to generate the example in figure 2A and 2B.
- Recall that in this model, there are unobserved components of cost.

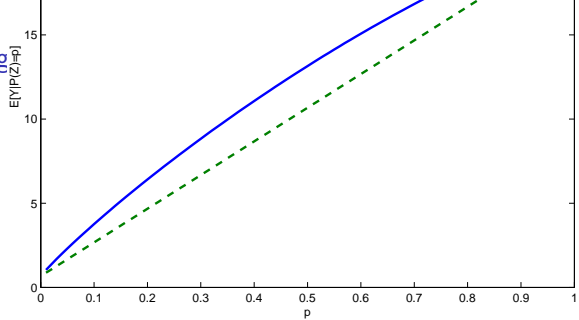
- Figure 3A plots two cases of  $E(Y | P(Z) = p)$  based on the generalized Roy model used to generate the example in figure 2A and 2B.
- Recall that in this model, there are unobserved components of cost.
- When  $\Delta^{\text{MTE}} (= E(\beta | X = x, V = v))$  does not depend on  $u_D$  (or  $v$ ) the expectation is a straight line.

- Figure 3A plots two cases of  $E(Y | P(Z) = p)$  based on the generalized Roy model used to generate the example in figure 2A and 2B.
- Recall that in this model, there are unobserved components of cost.
- When  $\Delta^{\text{MTE}} (= E(\beta | X = x, V = v))$  does not depend on  $u_D$  (or  $v$ ) the expectation is a straight line.
- This is Case (C-1).

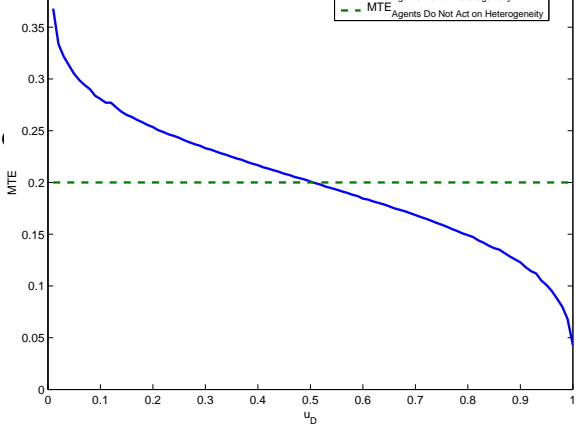
- Figure 3A plots two cases of  $E(Y | P(Z) = p)$  based on the generalized Roy model used to generate the example in figure 2A and 2B.
- Recall that in this model, there are unobserved components of cost.
- When  $\Delta^{\text{MTE}} (= E(\beta | X = x, V = v))$  does not depend on  $u_D$  (or  $v$ ) the expectation is a straight line.
- This is Case (C-1).
- Figure 3B plots the derivatives of the two curves in figure 3A.

- Figure 3A plots two cases of  $E(Y | P(Z) = p)$  based on the generalized Roy model used to generate the example in figure 2A and 2B.
- Recall that in this model, there are unobserved components of cost.
- When  $\Delta^{\text{MTE}} (= E(\beta | X = x, V = v))$  does not depend on  $u_D$  (or  $v$ ) the expectation is a straight line.
- This is Case (C-1).
- Figure 3B plots the derivatives of the two curves in figure 3A.
- When  $\Delta^{\text{MTE}}$  depends on  $u_D$  (or  $v$ ) (Case (C-2)), people sort into the program being studied positively on the basis of gains from the program, and one obtains the curved line depicted in figure 3A.

Fig



## B. Plot of the i



re 3A



*Note:* Parameters for the general heterogeneous case are the same as those used in Figures 2A and 2B. For the homogeneous case we impose

$$U_1 = U_0 (\sigma_1 = \sigma_0 = 0.012).$$

Source: Heckman and Vytlačil (2005).

## What Does Linear IV Estimate?

- It is instructive to determine what linear IV estimates when  $\Delta^{\text{MTE}}$  is nonconstant and conditions (A-1)–(A-5) hold.

## What Does Linear IV Estimate?

- It is instructive to determine what linear IV estimates when  $\Delta^{\text{MTE}}$  is nonconstant and conditions (A-1)–(A-5) hold.
- We analyze the general nonseparable case.

## What Does Linear IV Estimate?

- It is instructive to determine what linear IV estimates when  $\Delta^{\text{MTE}}$  is nonconstant and conditions (A-1)–(A-5) hold.
- We analyze the general nonseparable case.
- We consider instrumental variables conditional on  $X = x$  using a general function of  $Z$  as an instrument.

## What Does Linear IV Estimate?

- It is instructive to determine what linear IV estimates when  $\Delta^{\text{MTE}}$  is nonconstant and conditions (A-1)–(A-5) hold.
- We analyze the general nonseparable case.
- We consider instrumental variables conditional on  $X = x$  using a general function of  $Z$  as an instrument.
- We then specialize our result using  $P(Z)$  as the instrument.

## What Does Linear IV Estimate?

- It is instructive to determine what linear IV estimates when  $\Delta^{\text{MTE}}$  is nonconstant and conditions (A-1)–(A-5) hold.
- We analyze the general nonseparable case.
- We consider instrumental variables conditional on  $X = x$  using a general function of  $Z$  as an instrument.
- We then specialize our result using  $P(Z)$  as the instrument.
- As before, let  $J(Z)$  be any function of  $Z$  such that  $\text{Cov}(J(Z), D) \neq 0$ .

- Define the IV estimator:

$$\beta_{IV}(J) \equiv \frac{\text{Cov}(J(Z), Y)}{\text{Cov}(J(Z), D)},$$

where to simplify the notation we keep the conditioning on  $X$  implicit.

- Define the IV estimator:

$$\beta_{IV}(J) \equiv \frac{\text{Cov}(J(Z), Y)}{\text{Cov}(J(Z), D)},$$

where to simplify the notation we keep the conditioning on  $X$  implicit.

- Appendix, Slide 1090, derives a representation of this expression in terms of weighted averages of the MTE displayed in table 2B.



- Define the IV estimator:

$$\beta_{IV}(J) \equiv \frac{\text{Cov}(J(Z), Y)}{\text{Cov}(J(Z), D)},$$

where to simplify the notation we keep the conditioning on  $X$  implicit.

- Appendix, Slide 1090, derives a representation of this expression in terms of weighted averages of the MTE displayed in table 2B.
- We exposit this expression in this section.

- In Appendix, Slide 1090, we establish that:

$$\begin{aligned} \text{Cov}(J(Z), Y) & \\ &= \int_0^1 \Delta^{\text{MTE}}(u_D) E(J(Z) - E(J(Z)) \mid P(Z) \geq u_D) \Pr(P(Z) \geq u_D) du_D. \end{aligned} \tag{22}$$

- In Appendix, Slide 1090, we establish that:

$$\begin{aligned} \text{Cov}(J(Z), Y) & \\ &= \int_0^1 \Delta^{\text{MTE}}(u_D) E(J(Z) - E(J(Z)) \mid P(Z) \geq u_D) \Pr(P(Z) \geq u_D) du_D. \end{aligned} \tag{22}$$

- By the law of iterated expectations,  
 $\text{Cov}(J(Z), D) = \text{Cov}(J(Z), P(Z)).$

- In Appendix, Slide 1090, we establish that:

$$\begin{aligned} \text{Cov}(J(Z), Y) \\ = \int_0^1 \Delta^{\text{MTE}}(u_D) E(J(Z) - E(J(Z)) \mid P(Z) \geq u_D) \Pr(P(Z) \geq u_D) du_D. \end{aligned} \quad (22)$$

- By the law of iterated expectations,  
 $\text{Cov}(J(Z), D) = \text{Cov}(J(Z), P(Z)).$

- Thus

$$\beta_{\text{IV}}(J) = \int_0^1 \Delta^{\text{MTE}}(u_D) \omega_{\text{IV}}(u_D \mid J) du_D,$$

where

$$\omega_{\text{IV}}(u_D \mid J) = \frac{E(J(Z) - E(J(Z)) \mid P(Z) \geq u_D) \Pr(P(Z) \geq u_D)}{\text{Cov}(J(Z), P(Z))}, \quad (23)$$

assuming the standard rank condition (IV-2) holds:  
 $\text{Cov}(J(Z), P(Z)) \neq 0.$

- The weights integrate to one,

$$\int_0^1 \omega_{IV}(u_D | J) du_D = 1,$$

and can be constructed from the data on  $P(Z)$ ,  $J(Z)$  and  $D$ .

- The weights integrate to one,

$$\int_0^1 \omega_{IV}(u_D | J) du_D = 1,$$

and can be constructed from the data on  $P(Z)$ ,  $J(Z)$  and  $D$ .

- Assumptions about the properties of the weights are testable.

- We discuss additional properties of the weights for the special case where the propensity score is the instrument  $J(Z) = P(Z)$ .

- We discuss additional properties of the weights for the special case where the propensity score is the instrument  $J(Z) = P(Z)$ .
- We then analyze the properties of the weights for a general instrument  $J(Z)$ .



- We discuss additional properties of the weights for the special case where the propensity score is the instrument  $J(Z) = P(Z)$ .
- We then analyze the properties of the weights for a general instrument  $J(Z)$ .
- When  $J(Z) = P(Z)$ , equation (23) specializes to

$$\begin{aligned} \omega_{IV}(u_D | P(Z)) \\ &= \frac{[E(P(Z) | P(Z) \geq u_D) - E(P(Z))] \Pr(P(Z) \geq u_D)}{\text{Var}(P(Z))}. \end{aligned}$$

- We discuss additional properties of the weights for the special case where the propensity score is the instrument  $J(Z) = P(Z)$ .
- We then analyze the properties of the weights for a general instrument  $J(Z)$ .
- When  $J(Z) = P(Z)$ , equation (23) specializes to

$$\begin{aligned} \omega_{IV}(u_D | P(Z)) \\ = \frac{[E(P(Z) | P(Z) \geq u_D) - E(P(Z))]\Pr(P(Z) \geq u_D)}{\text{Var}(P(Z))}. \end{aligned}$$

- Figure 4A plots the IV weight for  $J(Z) = P(Z)$  and the MTE for our generalized Roy model example developed in figures 2 and 3 and table 3.

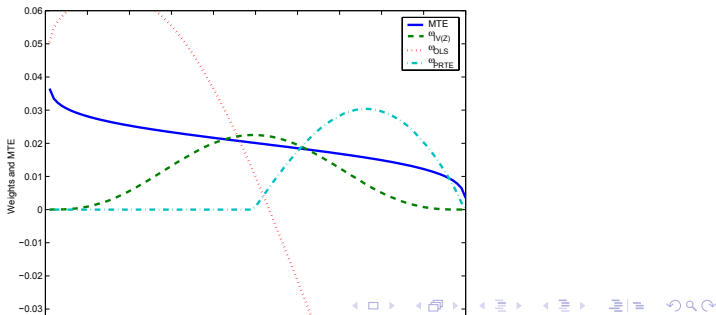
- We discuss additional properties of the weights for the special case where the propensity score is the instrument  $J(Z) = P(Z)$ .
- We then analyze the properties of the weights for a general instrument  $J(Z)$ .
- When  $J(Z) = P(Z)$ , equation (23) specializes to

$$\begin{aligned} \omega_{IV}(u_D | P(Z)) \\ = \frac{[E(P(Z) | P(Z) \geq u_D) - E(P(Z))]\Pr(P(Z) \geq u_D)}{\text{Var}(P(Z))}. \end{aligned}$$

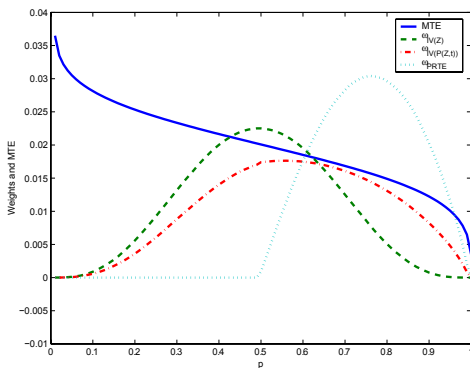
- Figure 4A plots the IV weight for  $J(Z) = P(Z)$  and the MTE for our generalized Roy model example developed in figures 2 and 3 and table 3.
- The weights are positive and peak at the mean of  $P$ .

Figure 4: A. Marginal Treatment Effect vs Linear Instrumental Variables, Ordinary Least Squares, and Policy Relevant Treatment Effect Weights: When  $P(Z)$  is the Instrument

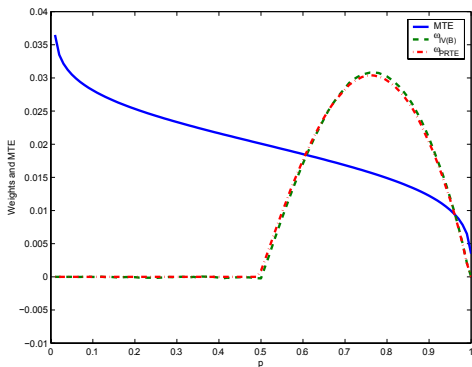
The Policy is Given at the Base of Table 3. The model parameters are given at the base of Figure 2.



B. Marginal Treatment Effect vs. Linear IV with  $Z$  as an Instrument, Linear IV with  $P(Z(1 + t(\mathbf{1}[Z > 0]))) = \tilde{P}(z, t)$  as an Instrument, and Policy Relevant Treatment Effect Weights. For The Policy Defined at the Base of Table 3. The model parameters are given at the base of Figure 2.



### C. Marginal Treatment Effect vs. IV Policy and Policy Relevant Treatment Effect Weights For The Policy Defined at the Base of Table 3.



- Figure 4A also plots the OLS weight given in table 2 and the weight for a policy exercise described below table 3 and discussed further below.

- Figure 4A also plots the OLS weight given in table 2 and the weight for a policy exercise described below table 3 and discussed further below.
- Let  $p^{\text{Min}}$  and  $p^{\text{Max}}$  denote the minimum and maximum points in the support of the distribution of  $P(Z)$  (conditional on  $X = x$ ).



- Figure 4A also plots the OLS weight given in table 2 and the weight for a policy exercise described below table 3 and discussed further below.
- Let  $p^{\text{Min}}$  and  $p^{\text{Max}}$  denote the minimum and maximum points in the support of the distribution of  $P(Z)$  (conditional on  $X = x$ ).
- The weights on MTE when  $P(Z)$  is the instrument are nonnegative for all evaluation points, are strictly positive for  $u_D \in (p^{\text{Min}}, p^{\text{Max}})$  and are zero for  $u_D < p^{\text{Min}}$  and for  $u_D > p^{\text{Max}}$ .

- The properties of the weights for general  $J(Z)$  depend on the conditional relationship between  $J(Z)$  and  $P(Z)$ .

- The properties of the weights for general  $J(Z)$  depend on the conditional relationship between  $J(Z)$  and  $P(Z)$ .
- From the general expression for (23), it is clear that the IV estimator with  $J(Z)$  as an instrument satisfies the following properties:

- The properties of the weights for general  $J(Z)$  depend on the conditional relationship between  $J(Z)$  and  $P(Z)$ .
- From the general expression for (23), it is clear that the IV estimator with  $J(Z)$  as an instrument satisfies the following properties:
  - (i) Two instruments  $J$  and  $J^*$  weight MTE equally at all values of  $u_D$  if and only if they have the same (centered) conditional expectation of  $J$  given  $P$ , i.e.,
$$E(J|P(Z) = p) - E(J) = E(J^* | P(Z) = p) - E(J^*)$$
 for all  $p$  in the support of the distribution of  $P(Z)$ .

- (ii) The support of  $\omega_{IV}(u_D | J)$  is contained in  $[p^{\text{Min}}, p^{\text{Max}}]$  the minimum and maximum value of  $p$  in the population (given  $x$ ). Therefore  $\omega_{IV}(t | J) = 0$  for  $t < p^{\text{Min}}$  and for  $t > p^{\text{Max}}$ . Using any instrument other than  $P(Z)$  leads to nonzero weights only on a subset of  $[p^{\text{Min}}, p^{\text{Max}}]$ , and using the propensity score as an instrument leads to nonnegative weights on a larger range of evaluation points than using any other instrument.

- (iii)  $\omega_{IV}(u_D | J)$  is nonnegative for all  $u_D$  if  $E(J | P(Z) \geq p)$  is weakly monotonic in  $p$ . Using  $J$  as an instrument yields nonnegative weights on  $\Delta^{\text{MTE}}$  if  $E(J | P(Z) \geq p)$  is weakly monotonic in  $p$ . This condition is satisfied when  $J(Z) = P(Z)$ . More generally, if  $J$  is a monotonic function of  $P(Z)$ , then using  $J$  as the instrument will lead to nonnegative weights on  $\Delta^{\text{MTE}}$ . There is no guarantee that the weights for a general  $J(Z)$  will be nonnegative for all  $u_D$ , although the weights integrate to unity and thus must be positive over some range of evaluation points. We produce examples below where the instrument leads to negative weights for some evaluation points. ? assume that  $J(Z)$  is monotonic in  $P(Z)$  and thus produce positive weights. Our analysis is more general.

- The propensity score plays a central role in determining the properties of the weights.

- The propensity score plays a central role in determining the properties of the weights.
- The IV weighting formula critically depends on the conditional mean dependence between instrument  $J(Z)$  and the propensity score.



- The interpretation placed on the IV estimand depends on the specification of  $P(Z)$  even if only  $Z_1$  (e.g., the first coordinate of  $Z$ ) is used as the instrument.

- The interpretation placed on the IV estimand depends on the specification of  $P(Z)$  even if only  $Z_1$  (e.g., the first coordinate of  $Z$ ) is used as the instrument.
- This drives home the point about the difference between IV in the traditional model and IV in the more general model with heterogeneous responses analyzed in this chapter.

- The interpretation placed on the IV estimand depends on the specification of  $P(Z)$  even if only  $Z_1$  (e.g., the first coordinate of  $Z$ ) is used as the instrument.
- This drives home the point about the difference between IV in the traditional model and IV in the more general model with heterogeneous responses analyzed in this chapter.
- In the traditional model, the choice of any valid instrument and the specification of instruments in  $P(Z)$  not used to construct a particular IV estimator does not affect the IV estimand.

- The interpretation placed on the IV estimand depends on the specification of  $P(Z)$  even if only  $Z_1$  (e.g., the first coordinate of  $Z$ ) is used as the instrument.
- This drives home the point about the difference between IV in the traditional model and IV in the more general model with heterogeneous responses analyzed in this chapter.
- In the traditional model, the choice of any valid instrument and the specification of instruments in  $P(Z)$  not used to construct a particular IV estimator does not affect the IV estimand.
- In the more general model, these choices matter.

- Two economists, using the same  $J(Z) = Z_1$ , will obtain the same IV point estimate, but the interpretation placed on that estimate will depend on the specification of the  $Z$  in  $P(Z)$  even if  $P(Z)$  is not used as an instrument.

- Two economists, using the same  $J(Z) = Z_1$ , will obtain the same IV point estimate, but the interpretation placed on that estimate will depend on the specification of the  $Z$  in  $P(Z)$  even if  $P(Z)$  is not used as an instrument.
- The weights can be positive for one instrument and negative for another.

- Two economists, using the same  $J(Z) = Z_1$ , will obtain the same IV point estimate, but the interpretation placed on that estimate will depend on the specification of the  $Z$  in  $P(Z)$  even if  $P(Z)$  is not used as an instrument.
- The weights can be positive for one instrument and negative for another.
- We show some examples after developing the properties of the IV weights.

## Further properties of the IV Weights

- Expression (23) for the weights does not impose any support conditions on the distribution of  $P(Z)$ , and thus does not require either that  $P(Z)$  be continuous or discrete.



## Further properties of the IV Weights

- Expression (23) for the weights does not impose any support conditions on the distribution of  $P(Z)$ , and thus does not require either that  $P(Z)$  be continuous or discrete.
- To demonstrate this, consider two extreme special cases: (i) when  $P(Z)$  is a continuous random variable, and (ii) when  $P(Z)$  is a discrete random variable.

- To simplify the exposition, initially assume that  $J(Z)$  and  $P(Z)$  are jointly continuous random variables.

- To simplify the exposition, initially assume that  $J(Z)$  and  $P(Z)$  are jointly continuous random variables.
- This assumption plays no essential role in any of the results of this chapter and we develop the discrete case after developing the continuous case.

- To simplify the exposition, initially assume that  $J(Z)$  and  $P(Z)$  are jointly continuous random variables.
- This assumption plays no essential role in any of the results of this chapter and we develop the discrete case after developing the continuous case.
- The weights defined in equation (23) can be written as

$$\omega_{IV}(u_D) = \frac{\int (j - E(J(Z))) \int_{u_D}^1 f_{J,P}(j, t) dt dj}{\text{Cov}(J(Z), D)}, \quad (24)$$

where  $f_{J,P}$  is the joint density of  $J(Z)$  and  $P(Z)$  and we implicitly condition on  $X$ .

- To simplify the exposition, initially assume that  $J(Z)$  and  $P(Z)$  are jointly continuous random variables.
- This assumption plays no essential role in any of the results of this chapter and we develop the discrete case after developing the continuous case.
- The weights defined in equation (23) can be written as

$$\omega_{IV}(u_D) = \frac{\int (j - E(J(Z))) \int_{u_D}^1 f_{J,P}(j, t) dt dj}{\text{Cov}(J(Z), D)}, \quad (24)$$

where  $f_{J,P}$  is the joint density of  $J(Z)$  and  $P(Z)$  and we implicitly condition on  $X$ .

- The weights can be negative or positive.

- Observe that  $\omega(0) = 0$  and  $\omega(1) = 0$ .

- Observe that  $\omega(0) = 0$  and  $\omega(1) = 0$ .
- The weights integrate to 1 because as shown in Appendix, Slide 1090,

$$\int \int (j - E(J(Z))) \int_{u_D}^1 f_{J,P}(j, t) dt dj du_D = \text{Cov}(J(Z), D),$$

so even if the weight is negative over some intervals, it must be positive over other intervals.

- Observe that  $\omega(0) = 0$  and  $\omega(1) = 0$ .
- The weights integrate to 1 because as shown in Appendix, Slide 1090,

$$\int \int (j - E(J(Z))) \int_{u_D}^1 f_{J,P}(j, t) dt dj du_D = \text{Cov}(J(Z), D),$$

so even if the weight is negative over some intervals, it must be positive over other intervals.

- Observe that when there is one instrument ( $Z$  is a scalar), and assumptions (A-1)–(A-5) are satisfied, the weights are always positive provided  $J(Z)$  is a monotonic function of the scalar  $Z$ .



- Observe that  $\omega(0) = 0$  and  $\omega(1) = 0$ .
- The weights integrate to 1 because as shown in Appendix, Slide 1090,

$$\int \int (j - E(J(Z))) \int_{u_D}^1 f_{J,P}(j, t) dt dj du_D = \text{Cov}(J(Z), D),$$

so even if the weight is negative over some intervals, it must be positive over other intervals.

- Observe that when there is one instrument ( $Z$  is a scalar), and assumptions (A-1)–(A-5) are satisfied, the weights are always positive provided  $J(Z)$  is a monotonic function of the scalar  $Z$ .
- In this case, which is covered by (23), but excluded in deriving (24),  $J(Z)$  and  $P(Z)$  have the same distribution and  $f_{J,P}(j, t)$  collapses to a univariate distribution.

- The possibility of negative weights arises when  $J(Z)$  is not a monotonic function of  $P(Z)$ .

- The possibility of negative weights arises when  $J(Z)$  is not a monotonic function of  $P(Z)$ .
- It also arises when there are two or more instruments, and the analyst computes estimates with only one instrument or a combination of the  $Z$  instruments that is not a monotonic function of  $P(Z)$  so that  $J(Z)$  and  $P(Z)$  are not perfectly dependent.

- The possibility of negative weights arises when  $J(Z)$  is not a monotonic function of  $P(Z)$ .
- It also arises when there are two or more instruments, and the analyst computes estimates with only one instrument or a combination of the  $Z$  instruments that is not a monotonic function of  $P(Z)$  so that  $J(Z)$  and  $P(Z)$  are not perfectly dependent.
- If the instrument is  $P(Z)$  (so  $J(Z) = P(Z)$ ) then the weights are everywhere non-negative because from (24),  
$$E(P(Z) | P(Z) > u_D) - E(P(Z)) \geq 0.$$

- The possibility of negative weights arises when  $J(Z)$  is not a monotonic function of  $P(Z)$ .
- It also arises when there are two or more instruments, and the analyst computes estimates with only one instrument or a combination of the  $Z$  instruments that is not a monotonic function of  $P(Z)$  so that  $J(Z)$  and  $P(Z)$  are not perfectly dependent.
- If the instrument is  $P(Z)$  (so  $J(Z) = P(Z)$ ) then the weights are everywhere non-negative because from (24),  
$$E(P(Z) | P(Z) > u_D) - E(P(Z)) \geq 0.$$
- In this case, the density of  $(P(Z), J(Z))$  collapses to the density of  $P(Z)$ .

- For any scalar  $Z$ , we can define  $J(Z)$  and  $P(Z)$  so that they are perfectly dependent, provided that  $J(Z)$  and  $P(Z)$  are monotonic in  $Z$ .

- For any scalar  $Z$ , we can define  $J(Z)$  and  $P(Z)$  so that they are perfectly dependent, provided that  $J(Z)$  and  $P(Z)$  are monotonic in  $Z$ .
- Generally, the weight (23) is positive if  $E(J(Z) | P(Z) > u_D)$  is weakly monotonic in  $u_D$ .

- For any scalar  $Z$ , we can define  $J(Z)$  and  $P(Z)$  so that they are perfectly dependent, provided that  $J(Z)$  and  $P(Z)$  are monotonic in  $Z$ .
- Generally, the weight (23) is positive if  $E(J(Z) | P(Z) > u_D)$  is weakly monotonic in  $u_D$ .
- Nonmonotonicity of this expression can produce negative weights.



## Constructing the Weights from Data

- Observe that the weights can be constructed from data on  $(J, P, D)$ .

## Constructing the Weights from Data

- Observe that the weights can be constructed from data on  $(J, P, D)$ .
- Data on  $(J(Z), P(Z))$  pairs and  $(J(Z), D)$  pairs (for each  $X$  value) are all that is required.

## Constructing the Weights from Data

- Observe that the weights can be constructed from data on  $(J, P, D)$ .
- Data on  $(J(Z), P(Z))$  pairs and  $(J(Z), D)$  pairs (for each  $X$  value) are all that is required.
- We can use a smoothed sample frequency to estimate the joint density  $f_{J,P}$ .

## Constructing the Weights from Data

- Observe that the weights can be constructed from data on  $(J, P, D)$ .
- Data on  $(J(Z), P(Z))$  pairs and  $(J(Z), D)$  pairs (for each  $X$  value) are all that is required.
- We can use a smoothed sample frequency to estimate the joint density  $f_{J,P}$ .
- Thus, given our maintained assumptions, any property of the weight, including its positivity at any point  $(x, u_D)$ , can be examined with data.

## Constructing the Weights from Data

- Observe that the weights can be constructed from data on  $(J, P, D)$ .
- Data on  $(J(Z), P(Z))$  pairs and  $(J(Z), D)$  pairs (for each  $X$  value) are all that is required.
- We can use a smoothed sample frequency to estimate the joint density  $f_{J,P}$ .
- Thus, given our maintained assumptions, any property of the weight, including its positivity at any point  $(x, u_D)$ , can be examined with data.
- We present examples of this approach below.

- As is evident from tables 2A and 2B and figures 2A and 2B, the weights on  $\Delta^{\text{MTE}}(u_D)$  generating  $\Delta^{\text{IV}}$  are different from the weights on  $\Delta^{\text{MTE}}(u_D)$  that generate the average treatment effect which is widely regarded as an important policy parameter (see, e.g., ?) or from the weights associated with the policy relevant treatment parameter which answers well-posed policy questions (??).

- As is evident from tables 2A and 2B and figures 2A and 2B, the weights on  $\Delta^{\text{MTE}}(u_D)$  generating  $\Delta^{\text{IV}}$  are different from the weights on  $\Delta^{\text{MTE}}(u_D)$  that generate the average treatment effect which is widely regarded as an important policy parameter (see, e.g., ?) or from the weights associated with the policy relevant treatment parameter which answers well-posed policy questions (??).
- It is not obvious why the weighted average of  $\Delta^{\text{MTE}}(u_D)$  produced by IV is of any economic interest.

- As is evident from tables 2A and 2B and figures 2A and 2B, the weights on  $\Delta^{\text{MTE}}(u_D)$  generating  $\Delta^{\text{IV}}$  are different from the weights on  $\Delta^{\text{MTE}}(u_D)$  that generate the average treatment effect which is widely regarded as an important policy parameter (see, e.g., ?) or from the weights associated with the policy relevant treatment parameter which answers well-posed policy questions (??).
- It is not obvious why the weighted average of  $\Delta^{\text{MTE}}(u_D)$  produced by IV is of any economic interest.
- Since the weights can be negative for some values of  $u_D$ ,  $\Delta^{\text{MTE}}(u_D)$  can be positive everywhere in  $u_D$  but IV can be negative.



- Thus, IV may not estimate a treatment effect for any person.

- Thus, IV may not estimate a treatment effect for any person.
- We present some examples of IV models with negative weights below.

- Thus, IV may not estimate a treatment effect for any person.
- We present some examples of IV models with negative weights below.
- A basic question is why estimate the model with IV at all given the lack of any clear economic interpretation of the IV estimator in the general case.

## Discrete Instruments

- The representation (23) can be specialized to cover discrete instruments,  $J(Z)$ .

## Discrete Instruments

- The representation (23) can be specialized to cover discrete instruments,  $J(Z)$ .
- Consider the case where the distribution of  $P(Z)$  (conditional on  $X$ ) is discrete.

## Discrete Instruments

- The representation (23) can be specialized to cover discrete instruments,  $J(Z)$ .
- Consider the case where the distribution of  $P(Z)$  (conditional on  $X$ ) is discrete.
- The support of the distribution of  $P(Z)$  contains a finite number of values  $p_1 < p_2 < \dots < p_K$  and the support of the instrument  $J(Z)$  is also discrete taking  $I$  distinct values where  $I$  and  $K$  may be distinct.

## Discrete Instruments

- The representation (23) can be specialized to cover discrete instruments,  $J(Z)$ .
- Consider the case where the distribution of  $P(Z)$  (conditional on  $X$ ) is discrete.
- The support of the distribution of  $P(Z)$  contains a finite number of values  $p_1 < p_2 < \dots < p_K$  and the support of the instrument  $J(Z)$  is also discrete taking  $I$  distinct values where  $I$  and  $K$  may be distinct.
- $E(J(Z)|P(Z) \geq u_D)$  is constant in  $u_D$ , for  $u_D$  within any  $(p_\ell, p_{\ell+1})$  interval, and  $\Pr(P(Z) \geq u_D)$  is constant in  $u_D$ , for  $u_D$  within any  $(p_\ell, p_{\ell+1})$  interval, and thus  $\omega_{IV}^J(u_D)$  is constant in  $u_D$  over any  $(p_\ell, p_{\ell+1})$  interval.

- Let  $\lambda_\ell$  denote the weight on LATE for the interval  $(\ell, \ell + 1)$ .



- Let  $\lambda_\ell$  denote the weight on LATE for the interval  $(\ell, \ell + 1)$ .
- In this notation,

$$\begin{aligned}
 \Delta_J^{IV} &= \int E(Y_1 - Y_0 | U_D = u_D) \omega_{IV}^J(u_D) du_D \\
 &= \sum_{\ell=1}^{K-1} \lambda_\ell \int_{p_\ell}^{p_{\ell+1}} E(Y_1 - Y_0 | U_D = u_D) \frac{1}{(p_{\ell+1} - p_\ell)} du_D \\
 &= \sum_{\ell=1}^{K-1} \Delta^{\text{LATE}}(p_\ell, p_{\ell+1}) \lambda_\ell.
 \end{aligned} \tag{25}$$

- Let  $j_i$  be the  $i^{\text{th}}$  smallest value of the support of  $J(Z)$ .

- Let  $j_i$  be the  $i^{\text{th}}$  smallest value of the support of  $J(Z)$ .
- The discrete version of equation (23) is

$$\lambda_\ell = \frac{\sum_{i=1}^I (j_i - E(J(Z))) \sum_{t>\ell}^K (f(j_i, p_t))}{\text{Cov}(J(Z), D)} (p_{\ell+1} - p_\ell) \quad (26)$$

where  $f$  is the probability frequency of  $(j_i, p_t)$ : the probability that  $J(Z) = j_i$  and  $P(Z) = p_t$ .

- Let  $j_i$  be the  $i^{\text{th}}$  smallest value of the support of  $J(Z)$ .
- The discrete version of equation (23) is

$$\lambda_\ell = \frac{\sum_{i=1}^I (j_i - E(J(Z))) \sum_{t>\ell}^K (f(j_i, p_t))}{\text{Cov}(J(Z), D)} (p_{\ell+1} - p_\ell) \quad (26)$$

where  $f$  is the probability frequency of  $(j_i, p_t)$ : the probability that  $J(Z) = j_i$  and  $P(Z) = p_t$ .

- There is no presumption that high values of  $J(Z)$  are associated with high values of  $P(Z)$ .

- Let  $j_i$  be the  $i^{\text{th}}$  smallest value of the support of  $J(Z)$ .
- The discrete version of equation (23) is

$$\lambda_\ell = \frac{\sum_{i=1}^I (j_i - E(J(Z))) \sum_{t>\ell}^K (f(j_i, p_t))}{\text{Cov}(J(Z), D)} (p_{\ell+1} - p_\ell) \quad (26)$$

where  $f$  is the probability frequency of  $(j_i, p_t)$ : the probability that  $J(Z) = j_i$  and  $P(Z) = p_t$ .

- There is no presumption that high values of  $J(Z)$  are associated with high values of  $P(Z)$ .
- $J(Z)$  can be one coordinate of  $Z$  that may be positively or negatively dependent on  $P(Z)$ , which depends on the full vector.

- In the case of scalar  $Z$ , as long as  $J(Z)$  and  $P(Z)$  are monotonic in  $Z$  there is perfect dependence between  $J(Z)$  and  $P(Z)$ .

- In the case of scalar  $Z$ , as long as  $J(Z)$  and  $P(Z)$  are monotonic in  $Z$  there is perfect dependence between  $J(Z)$  and  $P(Z)$ .
- In this case, the joint probability density collapses to a univariate density and the weights have to be positive, exactly as in the case for continuous instruments previously discussed.

- In the case of scalar  $Z$ , as long as  $J(Z)$  and  $P(Z)$  are monotonic in  $Z$  there is perfect dependence between  $J(Z)$  and  $P(Z)$ .
- In this case, the joint probability density collapses to a univariate density and the weights have to be positive, exactly as in the case for continuous instruments previously discussed.
- Our expression for the weight on LATE generalizes the expression presented by ? who in their analysis of the case of vector  $Z$  only consider the case where  $J(Z)$  and  $P(Z)$  are perfectly dependent because  $J(Z)$  is a monotonic function of  $P(Z)$ .



- In the case of scalar  $Z$ , as long as  $J(Z)$  and  $P(Z)$  are monotonic in  $Z$  there is perfect dependence between  $J(Z)$  and  $P(Z)$ .
- In this case, the joint probability density collapses to a univariate density and the weights have to be positive, exactly as in the case for continuous instruments previously discussed.
- Our expression for the weight on LATE generalizes the expression presented by ? who in their analysis of the case of vector  $Z$  only consider the case where  $J(Z)$  and  $P(Z)$  are perfectly dependent because  $J(Z)$  is a monotonic function of  $P(Z)$ .
- More generally, the weights can be positive or negative for any  $\ell$  but they must sum to 1 over all  $\ell$ .

- Monotonicity or uniformity is a property needed with just two values of  $Z$ ,  $Z = z_1$  and  $Z = z_2$ , to guarantee that IV estimates a treatment effect.

- Monotonicity or uniformity is a property needed with just two values of  $Z$ ,  $Z = z_1$  and  $Z = z_2$ , to guarantee that IV estimates a treatment effect.
- With more than two values of  $Z$ , we need to weight the LATEs and MTEs.

- Monotonicity or uniformity is a property needed with just two values of  $Z$ ,  $Z = z_1$  and  $Z = z_2$ , to guarantee that IV estimates a treatment effect.
- With more than two values of  $Z$ , we need to weight the LATEs and MTEs.
- If the instrument  $J(Z)$  shifts  $P(Z)$  in the same way for everyone, it shifts  $D$  in the same way for everyone since  $D = \mathbf{1}[P(Z) \geq U_D]$  and  $Z$  is independent of  $U_D$ .

- Monotonicity or uniformity is a property needed with just two values of  $Z$ ,  $Z = z_1$  and  $Z = z_2$ , to guarantee that IV estimates a treatment effect.
- With more than two values of  $Z$ , we need to weight the LATEs and MTEs.
- If the instrument  $J(Z)$  shifts  $P(Z)$  in the same way for everyone, it shifts  $D$  in the same way for everyone since  $D = \mathbf{1}[P(Z) \geq U_D]$  and  $Z$  is independent of  $U_D$ .
- If  $J(Z)$  is not monotonic in  $P(Z)$ , it may shift  $P(Z)$  in different ways for different people.

- Monotonicity or uniformity is a property needed with just two values of  $Z$ ,  $Z = z_1$  and  $Z = z_2$ , to guarantee that IV estimates a treatment effect.
- With more than two values of  $Z$ , we need to weight the LATEs and MTEs.
- If the instrument  $J(Z)$  shifts  $P(Z)$  in the same way for everyone, it shifts  $D$  in the same way for everyone since  $D = \mathbf{1}[P(Z) \geq U_D]$  and  $Z$  is independent of  $U_D$ .
- If  $J(Z)$  is not monotonic in  $P(Z)$ , it may shift  $P(Z)$  in different ways for different people.
- Negative weights are a tip-off of two-way flows.

- Monotonicity or uniformity is a property needed with just two values of  $Z$ ,  $Z = z_1$  and  $Z = z_2$ , to guarantee that IV estimates a treatment effect.
- With more than two values of  $Z$ , we need to weight the LATEs and MTEs.
- If the instrument  $J(Z)$  shifts  $P(Z)$  in the same way for everyone, it shifts  $D$  in the same way for everyone since  $D = \mathbf{1}[P(Z) \geq U_D]$  and  $Z$  is independent of  $U_D$ .
- If  $J(Z)$  is not monotonic in  $P(Z)$ , it may shift  $P(Z)$  in different ways for different people.
- Negative weights are a tip-off of two-way flows.
- We present examples below.

## Identifying Margins of Choice Associated With Each Instrument and Unifying Diverse Instruments Within a Common Framework

- We have just established that different instruments weight the MTE differently.



## Identifying Margins of Choice Associated With Each Instrument and Unifying Diverse Instruments Within a Common Framework

- We have just established that different instruments weight the MTE differently.
- Using  $P(Z)$  in the local IV estimator, we can identify the MTE.

## Identifying Margins of Choice Associated With Each Instrument and Unifying Diverse Instruments Within a Common Framework

- We have just established that different instruments weight the MTE differently.
- Using  $P(Z)$  in the local IV estimator, we can identify the MTE.
- We can construct the weights associated with each instrument from the joint distribution of  $(J(Z), P(Z))$  given  $X$ .

## Identifying Margins of Choice Associated With Each Instrument and Unifying Diverse Instruments Within a Common Framework

- We have just established that different instruments weight the MTE differently.
- Using  $P(Z)$  in the local IV estimator, we can identify the MTE.
- We can construct the weights associated with each instrument from the joint distribution of  $(J(Z), P(Z))$  given  $X$ .
- By plotting the weights for each instrument, we can determine the margins identified by the different instruments.

## Identifying Margins of Choice Associated With Each Instrument and Unifying Diverse Instruments Within a Common Framework

- We have just established that different instruments weight the MTE differently.
- Using  $P(Z)$  in the local IV estimator, we can identify the MTE.
- We can construct the weights associated with each instrument from the joint distribution of  $(J(Z), P(Z))$  given  $X$ .
- By plotting the weights for each instrument, we can determine the margins identified by the different instruments.
- Using  $P(Z)$  as the instrument enables us to extend the support associated with any single instrument, and to determine which segment of the MTE is identified by any particular instrument.

## Identifying Margins of Choice Associated With Each Instrument and Unifying Diverse Instruments Within a Common Framework

- We have just established that different instruments weight the MTE differently.
- Using  $P(Z)$  in the local IV estimator, we can identify the MTE.
- We can construct the weights associated with each instrument from the joint distribution of  $(J(Z), P(Z))$  given  $X$ .
- By plotting the weights for each instrument, we can determine the margins identified by the different instruments.
- Using  $P(Z)$  as the instrument enables us to extend the support associated with any single instrument, and to determine which segment of the MTE is identified by any particular instrument.
- As before, we keep conditioning on  $X$  implicit.

## Yitzhaki's Derivation of the Weights

- An alternative and in some ways more illuminating way to derive the weights used in IV is to follow ?? and ? who prove for a general regression function  $E(Y | P(Z) = p)$  that a linear regression of  $Y$  on  $P$  estimates

$$\beta_{Y,P} = \int_0^1 \left[ \frac{\partial E(Y | P(Z) = p)}{\partial p} \right] \omega(p) dp,$$

where

$$\omega(p) = \frac{\int_p^1 (t - E(P)) dF_P(t)}{\text{Var}(P)},$$

which is exactly the weight (23) when  $P$  is the instrument.

- Thus we can interpret (23) as the weight on  $\frac{\partial E(Y|P(Z)=p)}{\partial p}$  when two-stage least squares (2SLS) based on  $P(Z)$  is used to estimate the “causal effect” of  $D$  on  $Y$ .

- Thus we can interpret (23) as the weight on  $\frac{\partial E(Y|P(Z)=p)}{\partial p}$  when two-stage least squares (2SLS) based on  $P(Z)$  is used to estimate the “causal effect” of  $D$  on  $Y$ .
- Under uniformity,

$$\begin{aligned} \frac{\partial E(Y | P(Z) = p)}{\partial p} \Big|_{p=u_D} &= E(Y_1 - Y_0 | U_D = u_D) \\ &= \Delta^{\text{MTE}}(u_D). \end{aligned}$$



- Our analysis is more general than that of ? or ? because we allow for instruments that are not monotonic functions of  $P(Z)$ , whereas the Yitzhaki weighting formula only applies to instruments that are monotonic functions of  $P(Z)$ .

- Our analysis is more general than that of ? or ? because we allow for instruments that are not monotonic functions of  $P(Z)$ , whereas the Yitzhaki weighting formula only applies to instruments that are monotonic functions of  $P(Z)$ .
- The analysis of ? is more general than that of ?, because he does not impose uniformity (monotonicity).

- We present some further examples of these weights after discussing the role of  $P(Z)$  and the role of monotonicity and uniformity.

- We present some further examples of these weights after discussing the role of  $P(Z)$  and the role of monotonicity and uniformity.
- We present Yitzhaki's Theorem and the relationship of our analysis to Yitzhaki's analysis in Appendices, Slides 1098 and 1104.

## The Central Role of the Propensity Score

- Observe that both (23) and (24) (and their counterparts for LATE (25) and (26)) contain expressions involving the propensity score  $P(Z)$ , the probability of selection into treatment.

## The Central Role of the Propensity Score

- Observe that both (23) and (24) (and their counterparts for LATE (25) and (26)) contain expressions involving the propensity score  $P(Z)$ , the probability of selection into treatment.
- Under our assumptions, it is a monotonic function of the mean utility of treatment,  $\mu_D(Z)$ .

## The Central Role of the Propensity Score

- Observe that both (23) and (24) (and their counterparts for LATE (25) and (26)) contain expressions involving the propensity score  $P(Z)$ , the probability of selection into treatment.
- Under our assumptions, it is a monotonic function of the mean utility of treatment,  $\mu_D(Z)$ .
- The propensity score plays a central role in selection models as a determinant of control functions in selection models (see ??) as noted in Slide 184.

- In matching models, it provides a computationally convenient way to condition on  $Z$  (see, e.g., ??, and the discussion in Slide 675).



- In matching models, it provides a computationally convenient way to condition on  $Z$  (see, e.g., ??, and the discussion in Slide 675).
- For the IV weight to be correctly constructed and interpreted, we need to know the correct model for  $P(Z)$ , i.e., we need to know exactly which  $Z$  determine  $P(Z)$ .

- In matching models, it provides a computationally convenient way to condition on  $Z$  (see, e.g., ??, and the discussion in Slide 675).
- For the IV weight to be correctly constructed and interpreted, we need to know the correct model for  $P(Z)$ , i.e., we need to know exactly which  $Z$  determine  $P(Z)$ .
- As previously noted, this feature is not required in the traditional model for instrumental variables based on response heterogeneity.

- In matching models, it provides a computationally convenient way to condition on  $Z$  (see, e.g., ??, and the discussion in Slide 675).
- For the IV weight to be correctly constructed and interpreted, we need to know the correct model for  $P(Z)$ , i.e., we need to know exactly which  $Z$  determine  $P(Z)$ .
- As previously noted, this feature is not required in the traditional model for instrumental variables based on response heterogeneity.
- In that simpler framework, any instrument will identify  $\mu_1(X) - \mu_0(X)$  and the choice of a particular instrument affects efficiency but not identifiability.

- One can be casual about the choice model in the traditional setup, but not in the model of choice of treatment with essential heterogeneity.

- One can be casual about the choice model in the traditional setup, but not in the model of choice of treatment with essential heterogeneity.
- Thus, unlike the application of IV to traditional models under condition (C-1), IV applied in the model of essential heterogeneity depends on (a) the choice of the instrument  $J(Z)$ , (b) its dependence with  $P(Z)$ , the true propensity score or choice probability and (c) the specification of the propensity score (i.e., what variables go into  $Z$ ).

- One can be casual about the choice model in the traditional setup, but not in the model of choice of treatment with essential heterogeneity.
- Thus, unlike the application of IV to traditional models under condition (C-1), IV applied in the model of essential heterogeneity depends on (a) the choice of the instrument  $J(Z)$ , (b) its dependence with  $P(Z)$ , the true propensity score or choice probability and (c) the specification of the propensity score (i.e., what variables go into  $Z$ ).
- Using the propensity score one can identify LIV and LATE and the marginal returns at values of the unobserved  $U_D$ .

- One can be casual about the choice model in the traditional setup, but not in the model of choice of treatment with essential heterogeneity.
- Thus, unlike the application of IV to traditional models under condition (C-1), IV applied in the model of essential heterogeneity depends on (a) the choice of the instrument  $J(Z)$ , (b) its dependence with  $P(Z)$ , the true propensity score or choice probability and (c) the specification of the propensity score (i.e., what variables go into  $Z$ ).
- Using the propensity score one can identify LIV and LATE and the marginal returns at values of the unobserved  $U_D$ .
- From the MTE identified by  $P(Z)$  and the weights that can be constructed from the joint distribution of  $(J(Z), P(Z))$  given  $X$ , we can identify the segment of the MTE identified by any IV.

## Monotonicity, Uniformity and Conditional Instruments

- Monotonicity, or uniformity condition (IV-3), is a condition on a collection of counterfactuals for each person and hence is not testable, since we know only one element of the collection for any person.



## Monotonicity, Uniformity and Conditional Instruments

- Monotonicity, or uniformity condition (IV-3), is a condition on a collection of counterfactuals for each person and hence is not testable, since we know only one element of the collection for any person.
- It rules out general heterogeneous responses to treatment choices in response to changes in vector  $Z$ .

## Monotonicity, Uniformity and Conditional Instruments

- Monotonicity, or uniformity condition (IV-3), is a condition on a collection of counterfactuals for each person and hence is not testable, since we know only one element of the collection for any person.
- It rules out general heterogeneous responses to treatment choices in response to changes in vector  $Z$ .
- The recent literature on instrumental variables with heterogeneous responses is thus asymmetric.

## Monotonicity, Uniformity and Conditional Instruments

- Monotonicity, or uniformity condition (IV-3), is a condition on a collection of counterfactuals for each person and hence is not testable, since we know only one element of the collection for any person.
- It rules out general heterogeneous responses to treatment choices in response to changes in vector  $Z$ .
- The recent literature on instrumental variables with heterogeneous responses is thus asymmetric.
- Outcome equations can be heterogeneous in a general way while choice equations cannot be.

## Monotonicity, Uniformity and Conditional Instruments

- Monotonicity, or uniformity condition (IV-3), is a condition on a collection of counterfactuals for each person and hence is not testable, since we know only one element of the collection for any person.
- It rules out general heterogeneous responses to treatment choices in response to changes in vector  $Z$ .
- The recent literature on instrumental variables with heterogeneous responses is thus asymmetric.
- Outcome equations can be heterogeneous in a general way while choice equations cannot be.
- If  $\mu_D(Z) = Z\gamma$ , where  $\gamma$  is a common coefficient shared by everyone, the choice model satisfies the uniformity property.

- On the other hand, if  $\gamma$  is a random coefficient (i.e., has a nondegenerate distribution) that can take both negative and positive values, and there are two or more variables in  $Z$  with nondegenerate  $\gamma$  coefficients, uniformity can be violated.

- On the other hand, if  $\gamma$  is a random coefficient (i.e., has a nondegenerate distribution) that can take both negative and positive values, and there are two or more variables in  $Z$  with nondegenerate  $\gamma$  coefficients, uniformity can be violated.
- Different people can respond to changes in  $Z$  differently, so there can be non-uniformity.

- On the other hand, if  $\gamma$  is a random coefficient (i.e., has a nondegenerate distribution) that can take both negative and positive values, and there are two or more variables in  $Z$  with nondegenerate  $\gamma$  coefficients, uniformity can be violated.
- Different people can respond to changes in  $Z$  differently, so there can be non-uniformity.
- The uniformity condition can be violated even when all components of  $\gamma$  are of the same sign if  $Z$  is a vector and  $\gamma$  is a nondegenerate random variable.

- Changing one coordinate of  $Z$ , holding the other coordinates at different values across people is *not* the experiment that defines monotonicity or uniformity.



- Changing one coordinate of  $Z$ , holding the other coordinates at different values across people is *not* the experiment that defines monotonicity or uniformity.
- Changing one component of  $Z$ , allowing the other coordinates of  $Z$  to vary across people, does not necessarily produce uniform flows toward or against participation in the treatment status.

- Changing one coordinate of  $Z$ , holding the other coordinates at different values across people is *not* the experiment that defines monotonicity or uniformity.
- Changing one component of  $Z$ , allowing the other coordinates of  $Z$  to vary across people, does not necessarily produce uniform flows toward or against participation in the treatment status.
- For example, let  $\mu_D(z) = \gamma_0 + \gamma_1 z_1 + \gamma_2 z_2 + \gamma_3 z_1 z_2$ , where  $\gamma_0$ ,  $\gamma_1$ ,  $\gamma_2$  and  $\gamma_3$  are constants, and consider changing  $z_1$  from a common base state while holding  $z_2$  fixed at different values across people.

- If  $\gamma_3 < 0$ , then  $\mu_D(z)$  does not necessarily satisfy the uniformity condition.

- If  $\gamma_3 < 0$ , then  $\mu_D(z)$  does not necessarily satisfy the uniformity condition.
- If we move  $(z_1, z_2)$  as a pair from the same base values to the same destination values  $z'$ , uniformity is satisfied even if  $\gamma_3 < 0$ , although  $\mu_D(z)$  is not a monotonic function of  $z$ .

- Positive weights and uniformity are distinct issues.

- Positive weights and uniformity are distinct issues.
- Under uniformity, and assumptions (A-1)–(A-5), the weights on MTE for any particular instrument may be positive or negative.

- Positive weights and uniformity are distinct issues.
- Under uniformity, and assumptions (A-1)–(A-5), the weights on MTE for any particular instrument may be positive or negative.
- The weights for MTE using  $P(Z)$  must be positive as we have shown so the propensity score has a special status as an instrument.

- Positive weights and uniformity are distinct issues.
- Under uniformity, and assumptions (A-1)–(A-5), the weights on MTE for any particular instrument may be positive or negative.
- The weights for MTE using  $P(Z)$  must be positive as we have shown so the propensity score has a special status as an instrument.
- Negative weights associated with the use of  $J(Z)$  as an instrument do not necessarily imply failure of uniformity in  $Z$ .



- Positive weights and uniformity are distinct issues.
- Under uniformity, and assumptions (A-1)–(A-5), the weights on MTE for any particular instrument may be positive or negative.
- The weights for MTE using  $P(Z)$  must be positive as we have shown so the propensity score has a special status as an instrument.
- Negative weights associated with the use of  $J(Z)$  as an instrument do not necessarily imply failure of uniformity in  $Z$ .
- Even if uniformity is satisfied for  $Z$ , it is not necessarily satisfied for  $J(Z)$ .

- Condition (IV-3) is an assumption about a vector.

- Condition (IV-3) is an assumption about a vector.
- Fixing one combination of  $Z$  (when  $J$  is a function of  $Z$ ) or one coordinate of  $Z$  does not guarantee uniformity in  $J$  even if there is uniformity in  $Z$ .

- Condition (IV-3) is an assumption about a vector.
- Fixing one combination of  $Z$  (when  $J$  is a function of  $Z$ ) or one coordinate of  $Z$  does not guarantee uniformity in  $J$  even if there is uniformity in  $Z$ .
- The flow created by changing one coordinate of  $Z$  can be reversed by the flow created by the other components of  $Z$  if there is negative dependence among components even if *ceteris paribus* all components of  $Z$  affect  $D$  in the same direction.

- Condition (IV-3) is an assumption about a vector.
- Fixing one combination of  $Z$  (when  $J$  is a function of  $Z$ ) or one coordinate of  $Z$  does not guarantee uniformity in  $J$  even if there is uniformity in  $Z$ .
- The flow created by changing one coordinate of  $Z$  can be reversed by the flow created by the other components of  $Z$  if there is negative dependence among components even if *ceteris paribus* all components of  $Z$  affect  $D$  in the same direction.
- We present some examples below.

- The issues of positive weights and the existence of one way flows in response to an intervention are conceptually distinct.

- The issues of positive weights and the existence of one way flows in response to an intervention are conceptually distinct.
- Even with two values for a scalar  $Z$ , flows may be two way (see equation (15)).

- The issues of positive weights and the existence of one way flows in response to an intervention are conceptually distinct.
- Even with two values for a scalar  $Z$ , flows may be two way (see equation (15)).
- If we satisfy (IV-3) for a vector, so uniformity applies, weights for a particular instrument may be negative for certain intervals of  $U_D$  (i.e., for some of the LATE parameters).



- If we condition on  $Z_2 = z_2, \dots, Z_K = z_K$  using  $Z_1$  as an instrument, then a uniform flow condition is satisfied.

- If we condition on  $Z_2 = z_2, \dots, Z_K = z_K$  using  $Z_1$  as an instrument, then a uniform flow condition is satisfied.
- We call this *conditional uniformity*.

- If we condition on  $Z_2 = z_2, \dots, Z_K = z_K$  using  $Z_1$  as an instrument, then a uniform flow condition is satisfied.
- We call this *conditional uniformity*.
- By conditioning, we effectively convert the problem back to that of a scalar instrument where the weights must be positive.

- If we condition on  $Z_2 = z_2, \dots, Z_K = z_K$  using  $Z_1$  as an instrument, then a uniform flow condition is satisfied.
- We call this *conditional uniformity*.
- By conditioning, we effectively convert the problem back to that of a scalar instrument where the weights must be positive.
- If uniformity holds for  $Z_1$ , fixing the other  $Z$  at common values, then one dimensional LATE/MTE analysis applies.

- If we condition on  $Z_2 = z_2, \dots, Z_K = z_K$  using  $Z_1$  as an instrument, then a uniform flow condition is satisfied.
- We call this *conditional uniformity*.
- By conditioning, we effectively convert the problem back to that of a scalar instrument where the weights must be positive.
- If uniformity holds for  $Z_1$ , fixing the other  $Z$  at common values, then one dimensional LATE/MTE analysis applies.
- Clearly, the weights have to be defined conditionally.

- The concept of conditioning on other instruments to produce positive weights for the selected instrument is a new idea, not yet appreciated in the empirical IV literature and has no counterpart in the traditional IV model.

- The concept of conditioning on other instruments to produce positive weights for the selected instrument is a new idea, not yet appreciated in the empirical IV literature and has no counterpart in the traditional IV model.
- In the conventional model, the choice of a valid instrument affects efficiency but not the definition of the parameters as it does in the more general case.

- In summary, nothing in the economics of choice guarantees that if  $Z$  is changed from  $z$  to  $z'$ , that people respond in the same direction to the change.



- In summary, nothing in the economics of choice guarantees that if  $Z$  is changed from  $z$  to  $z'$ , that people respond in the same direction to the change.
- See the general expression (15).

- In summary, nothing in the economics of choice guarantees that if  $Z$  is changed from  $z$  to  $z'$ , that people respond in the same direction to the change.
- See the general expression (15).
- The condition that people respond to choices in the same direction for the same change in  $Z$  does not imply that  $D(z)$  is monotonic in  $z$  for any person in the usual mathematical usage of the term monotonicity.

- In summary, nothing in the economics of choice guarantees that if  $Z$  is changed from  $z$  to  $z'$ , that people respond in the same direction to the change.
- See the general expression (15).
- The condition that people respond to choices in the same direction for the same change in  $Z$  does not imply that  $D(z)$  is monotonic in  $z$  for any person in the usual mathematical usage of the term monotonicity.
- If  $D(z)$  is monotonic in the usual usage of this term and responses are in the same direction for all people, then “monotonicity” or better “uniformity” condition (IV-3) would be satisfied.

- If responses to a common change of  $Z$  are heterogenous in a general way, we obtain (15) as the general case.

- If responses to a common change of  $Z$  are heterogenous in a general way, we obtain (15) as the general case.
- Vytlacil's ? Theorem breaks down and IV cannot be expressed in terms of a weighted average of MTE terms.

- If responses to a common change of  $Z$  are heterogenous in a general way, we obtain (15) as the general case.
- Vytlačil's ? Theorem breaks down and IV cannot be expressed in terms of a weighted average of MTE terms.
- Nonetheless, Yitzhaki's characterization of IV, derived in Appendix, Slide 1090, remains valid and the weights on  $\frac{\partial E(Y|P=p)}{\partial p}$  are positive and of the same form as the weights obtained for MTE (or LATE) when the monotonicity condition holds.

- If responses to a common change of  $Z$  are heterogenous in a general way, we obtain (15) as the general case.
- Vytlacil's ? Theorem breaks down and IV cannot be expressed in terms of a weighted average of MTE terms.
- Nonetheless, Yitzhaki's characterization of IV, derived in Appendix, Slide 1090, remains valid and the weights on  $\frac{\partial E(Y|P=p)}{\partial p}$  are positive and of the same form as the weights obtained for MTE (or LATE) when the monotonicity condition holds.
- IV can still be written as a weighted average of LIV terms, even though LIV does not identify the MTE.

## Treatment Effects vs. Policy Effects

- Even if uniformity condition (IV-3) fails, IV may answer relevant policy questions.



## Treatment Effects vs. Policy Effects

- Even if uniformity condition (IV-3) fails, IV may answer relevant policy questions.
- By Yitzhaki's analysis, summarized in Slide 250, IV or 2SLS estimates a weighted average of marginal responses which may be pointwise positive or negative.

## Treatment Effects vs. Policy Effects

- Even if uniformity condition (IV-3) fails, IV may answer relevant policy questions.
- By Yitzhaki's analysis, summarized in Slide 250, IV or 2SLS estimates a weighted average of marginal responses which may be pointwise positive or negative.
- Policies may induce some people to switch into and others to switch out of choices, as is evident from equation (15).

## Treatment Effects vs. Policy Effects

- Even if uniformity condition (IV-3) fails, IV may answer relevant policy questions.
- By Yitzhaki's analysis, summarized in Slide 250, IV or 2SLS estimates a weighted average of marginal responses which may be pointwise positive or negative.
- Policies may induce some people to switch into and others to switch out of choices, as is evident from equation (15).
- These net effects are of interest in many policy analyses.

## Treatment Effects vs. Policy Effects

- Even if uniformity condition (IV-3) fails, IV may answer relevant policy questions.
- By Yitzhaki's analysis, summarized in Slide 250, IV or 2SLS estimates a weighted average of marginal responses which may be pointwise positive or negative.
- Policies may induce some people to switch into and others to switch out of choices, as is evident from equation (15).
- These net effects are of interest in many policy analyses.
- Thus, subsidized housing in a region supported by higher taxes may attract some to migrate to the region and cause others to leave.

- The net effect from the policy is all that is required to perform cost benefit calculations of the policy on outcomes.

- The net effect from the policy is all that is required to perform cost benefit calculations of the policy on outcomes.
- If the housing subsidy is the instrument, and the net effect of the subsidy is the parameter of interest, the issue of monotonicity is a red herring.

- The net effect from the policy is all that is required to perform cost benefit calculations of the policy on outcomes.
- If the housing subsidy is the instrument, and the net effect of the subsidy is the parameter of interest, the issue of monotonicity is a red herring.
- If the subsidy is exogenously imposed, IV estimates the net effect of the policy on mean outcomes.

- The net effect from the policy is all that is required to perform cost benefit calculations of the policy on outcomes.
- If the housing subsidy is the instrument, and the net effect of the subsidy is the parameter of interest, the issue of monotonicity is a red herring.
- If the subsidy is exogenously imposed, IV estimates the net effect of the policy on mean outcomes.
- Only if the effect of migration on outcomes induced by the subsidy on outcomes is the question of interest, and not the effect of the subsidy, does uniformity emerge as an interesting condition.



## Some Examples of Weights in the Generalized Roy Model and the Extended Roy Model

- It is useful to develop intuition about the properties of the IV estimator and the structure of the weights for two prototypical choice models.

## Some Examples of Weights in the Generalized Roy Model and the Extended Roy Model

- It is useful to develop intuition about the properties of the IV estimator and the structure of the weights for two prototypical choice models.
- We develop the weights for a generalized Roy model where unobserved cost components are present and an extended Roy model where cost components are observed but there are no unobserved cost components.

## Some Examples of Weights in the Generalized Roy Model and the Extended Roy Model

- It is useful to develop intuition about the properties of the IV estimator and the structure of the weights for two prototypical choice models.
- We develop the weights for a generalized Roy model where unobserved cost components are present and an extended Roy model where cost components are observed but there are no unobserved cost components.
- The extended Roy model is used to generate figure 1 and was introduced at the end of Slide 12.

- Table 3 presents the IV estimand for the generalized Roy model used to generate figures 2A and 2B using  $P(Z)$  as the instrument.

- Table 3 presents the IV estimand for the generalized Roy model used to generate figures 2A and 2B using  $P(Z)$  as the instrument.
- The model generating  $D = \mathbf{1}[Z\gamma \geq V]$  is given at the base of figure 2B ( $Z$  is a scalar,  $\gamma$  is 1,  $V$  is normal,  $U_D = \Phi\left(\frac{V}{\sigma_V}\right)$ ).

- Table 3 presents the IV estimand for the generalized Roy model used to generate figures 2A and 2B using  $P(Z)$  as the instrument.
- The model generating  $D = \mathbf{1}[Z\gamma \geq V]$  is given at the base of figure 2B ( $Z$  is a scalar,  $\gamma$  is 1,  $V$  is normal,  $U_D = \Phi\left(\frac{V}{\sigma_V}\right)$ ).
- We compare the IV estimand with the policy relevant treatment effect for a policy precisely defined at the base of table 3.

- Table 3 presents the IV estimand for the generalized Roy model used to generate figures 2A and 2B using  $P(Z)$  as the instrument.
- The model generating  $D = \mathbf{1}[Z\gamma \geq V]$  is given at the base of figure 2B ( $Z$  is a scalar,  $\gamma$  is 1,  $V$  is normal,  $U_D = \Phi\left(\frac{V}{\sigma_V}\right)$ ).
- We compare the IV estimand with the policy relevant treatment effect for a policy precisely defined at the base of table 3.
- This policy has the structure that if  $Z > 0$ , persons get a bonus  $Zt$  for participation in the program, where  $t > 0$ .

- Table 3 presents the IV estimand for the generalized Roy model used to generate figures 2A and 2B using  $P(Z)$  as the instrument.
- The model generating  $D = \mathbf{1}[Z\gamma \geq V]$  is given at the base of figure 2B ( $Z$  is a scalar,  $\gamma$  is 1,  $V$  is normal,  $U_D = \Phi\left(\frac{V}{\sigma_V}\right)$ ).
- We compare the IV estimand with the policy relevant treatment effect for a policy precisely defined at the base of table 3.
- This policy has the structure that if  $Z > 0$ , persons get a bonus  $Zt$  for participation in the program, where  $t > 0$ .
- The decision rule for program participation for  $Z > 0$  is  $D = \mathbf{1}[Z(1+t) \geq V]$ .



- Table 3 presents the IV estimand for the generalized Roy model used to generate figures 2A and 2B using  $P(Z)$  as the instrument.
- The model generating  $D = \mathbf{1}[Z\gamma \geq V]$  is given at the base of figure 2B ( $Z$  is a scalar,  $\gamma$  is 1,  $V$  is normal,  $U_D = \Phi\left(\frac{V}{\sigma_V}\right)$ ).
- We compare the IV estimand with the policy relevant treatment effect for a policy precisely defined at the base of table 3.
- This policy has the structure that if  $Z > 0$ , persons get a bonus  $Zt$  for participation in the program, where  $t > 0$ .
- The decision rule for program participation for  $Z > 0$  is  $D = \mathbf{1}[Z(1 + t) \geq V]$ .
- People are not forced into participation in the program but are rather induced into it by the bonus.

- Given the assumed distribution of  $Z$ , and the other parameters of the model, we obtain the policy relevant treatment parameter weight  $\omega_{\text{PRTE}}(u_D)$  as plotted in figures 4A–4C (the scales of the ordinates differ across the graphs, but the weight is the same).

- Given the assumed distribution of  $Z$ , and the other parameters of the model, we obtain the policy relevant treatment parameter weight  $\omega_{\text{PRTE}}(u_D)$  as plotted in figures 4A–4C (the scales of the ordinates differ across the graphs, but the weight is the same).
- We use the per capita PRTE and consider three instruments.

- Given the assumed distribution of  $Z$ , and the other parameters of the model, we obtain the policy relevant treatment parameter weight  $\omega_{\text{PRTE}}(u_D)$  as plotted in figures 4A–4C (the scales of the ordinates differ across the graphs, but the weight is the same).
- We use the per capita PRTE and consider three instruments.
- Table 5 presents estimands for the three instruments shown in the table for the generalized Roy model in three environments.

**Table 5:** Linear instrumental variable estimands and the policy relevant treatment effect

Using Propensity Score $P(Z)$ as the Instrument	0.2013
Using Propensity Score $P(Z(1 + t(\mathbf{1}[Z > 0])))$ as the Instrument	0.1859
Using a dummy $B$ as an Instrument <sup>a</sup>	0.1549
Policy Relevant Treatment Effect (PRTE)	0.1549

<sup>a</sup>The dummy  $B$  is such that  $B = 1$  if an individual belongs to a randomly assigned eligible population, 0 otherwise.

Source: Heckman and Vytlačil (2005)

- The first instrument we consider for this example is  $P(Z)$ , which assumes that there is no policy in place ( $t = 0$ ).

- The first instrument we consider for this example is  $P(Z)$ , which assumes that there is no policy in place ( $t = 0$ ).
- It is identified (estimated) on a sample with no policy in place but otherwise the model is the same as the one with the policy in place.

- The first instrument we consider for this example is  $P(Z)$ , which assumes that there is no policy in place ( $t = 0$ ).
- It is identified (estimated) on a sample with no policy in place but otherwise the model is the same as the one with the policy in place.
- The weight on this instrument is plotted in figure 4A.



- The first instrument we consider for this example is  $P(Z)$ , which assumes that there is no policy in place ( $t = 0$ ).
- It is identified (estimated) on a sample with no policy in place but otherwise the model is the same as the one with the policy in place.
- The weight on this instrument is plotted in figure 4A.
- That figure also displays the OLS weight as well as the MTE that is being weighted to generate the estimate.

- The first instrument we consider for this example is  $P(Z)$ , which assumes that there is no policy in place ( $t = 0$ ).
- It is identified (estimated) on a sample with no policy in place but otherwise the model is the same as the one with the policy in place.
- The weight on this instrument is plotted in figure 4A.
- That figure also displays the OLS weight as well as the MTE that is being weighted to generate the estimate.
- It also shows the weight used to generate PRTE.

- The IV weights for  $P(Z)$  and the weights for  $\Delta^{\text{PRTE}}$  differ.

- The IV weights for  $P(Z)$  and the weights for  $\Delta^{\text{PRTE}}$  differ.
- This is as it should be because  $\Delta^{\text{PRTE}}$  is making a comparison across regimes but the IV in this case makes comparisons within a no policy regime.

- The IV weights for  $P(Z)$  and the weights for  $\Delta^{\text{PRTE}}$  differ.
- This is as it should be because  $\Delta^{\text{PRTE}}$  is making a comparison across regimes but the IV in this case makes comparisons within a no policy regime.
- Given the shape of  $\Delta^{\text{MTE}}(u_D)$ , it is not surprising that the estimand for IV based on  $P(Z)$  is so much above the  $\Delta^{\text{PRTE}}$  which weights a lower-valued segment of  $\Delta^{\text{MTE}}(u_D)$  more heavily.

- The second instrument we consider exploits the variation induced by the policy in place and fits it on samples where the policy is in place (i.e., the  $t$  is the same as that used to generate the PRTE).

- The second instrument we consider exploits the variation induced by the policy in place and fits it on samples where the policy is in place (i.e., the  $t$  is the same as that used to generate the PRTE).
- On intuitive grounds, this instrument might be thought to work well in identifying the PRTE, but in fact it does not.

- The second instrument we consider exploits the variation induced by the policy in place and fits it on samples where the policy is in place (i.e., the  $t$  is the same as that used to generate the PRTE).
- On intuitive grounds, this instrument might be thought to work well in identifying the PRTE, but in fact it does not.
- The instrument is  $\tilde{P}(Z, t) = P(Z(1 + t\mathbf{1}[Z > 0]))$  which jumps in value when  $Z$  switches from  $Z < 0$  to  $Z > 0$ .



- The second instrument we consider exploits the variation induced by the policy in place and fits it on samples where the policy is in place (i.e., the  $t$  is the same as that used to generate the PRTE).
- On intuitive grounds, this instrument might be thought to work well in identifying the PRTE, but in fact it does not.
- The instrument is  $\tilde{P}(Z, t) = P(Z(1 + t\mathbf{1}[Z > 0]))$  which jumps in value when  $Z$  switches from  $Z < 0$  to  $Z > 0$ .
- This is the choice probability in the regime with the policy in place.

- Figure 4B plots the weight for this IV along with the weight for  $P(Z)$  as an IV and the weight for PRTE (repeated from figure 4A).

- Figure 4B plots the weight for this IV along with the weight for  $P(Z)$  as an IV and the weight for PRTE (repeated from figure 4A).
- While this weight looks a bit more like the weight for  $\Delta^{\text{PRTE}}$  than the previous instrument, it is clearly different.

- Figure 4C plots the weight for an ideal instrument for PRTE: a randomization of eligibility.

- Figure 4C plots the weight for an ideal instrument for PRTE: a randomization of eligibility.
- This compares the outcomes in one population where the policy is in place with outcomes in a regime where the policy is not in place.

- Figure 4C plots the weight for an ideal instrument for PRTE: a randomization of eligibility.
- This compares the outcomes in one population where the policy is in place with outcomes in a regime where the policy is not in place.
- Thus we use an instrument  $B$  such that

$$B = \begin{cases} 1 & \text{if a person is eligible to participate in the program} \\ 0 & \text{otherwise.} \end{cases}$$

- Persons for whom  $B = 1$ , make their participation choices under the policy with a jump in  $Z$ ,  $t\mathbf{1}(Z > 0)$ , in their choice sets.

- Persons for whom  $B = 1$ , make their participation choices under the policy with a jump in  $Z$ ,  $t\mathbf{1}(Z > 0)$ , in their choice sets.
- If  $B = 0$ , persons are embargoed from the policy and cannot receive a bonus.



- Persons for whom  $B = 1$ , make their participation choices under the policy with a jump in  $Z$ ,  $t\mathbf{1}(Z > 0)$ , in their choice sets.
- If  $B = 0$ , persons are embargoed from the policy and cannot receive a bonus.
- The  $B = 0$  case is a prepolicy regime.

- Persons for whom  $B = 1$ , make their participation choices under the policy with a jump in  $Z$ ,  $t\mathbf{1}(Z > 0)$ , in their choice sets.
- If  $B = 0$ , persons are embargoed from the policy and cannot receive a bonus.
- The  $B = 0$  case is a prepolicy regime.
- We assume  $\Pr[B = 1 \mid Y_0, Y_1, V, Z] = \Pr[B = 1] = 0.5$ , so all persons are equally likely to receive or not receive eligibility for the bonus and assignment does not depend on model unobservables in the outcome equation.

- The Wald estimator in this case is

$$\frac{E(Y | B = 1) - E(Y | B = 0)}{\Pr(D = 1 | B = 1) - \Pr(D = 1 | B = 0)}.$$

The IV weight for this estimator is a special case of equation (23):

$$\omega_{IV}(u_D | B) = \frac{E(B - E(B) | \hat{P}(Z) \geq u_D) \Pr(\hat{P}(Z) \geq u_D)}{\text{Cov}(B, \hat{P}(Z))},$$

where  $\hat{P}(Z) = P(Z(1 + t\mathbf{1}[Z > 0]))^B P(Z)^{(1-B)}$ .

- The Wald estimator in this case is

$$\frac{E(Y | B = 1) - E(Y | B = 0)}{\Pr(D = 1 | B = 1) - \Pr(D = 1 | B = 0)}.$$

The IV weight for this estimator is a special case of equation (23):

$$\omega_{IV}(u_D | B) = \frac{E\left(B - E(B) \mid \hat{P}(Z) \geq u_D\right) \Pr\left(\hat{P}(Z) \geq u_D\right)}{\text{Cov}\left(B, \hat{P}(Z)\right)},$$

where  $\hat{P}(Z) = P(Z(1 + t\mathbf{1}[Z > 0]))^B P(Z)^{(1-B)}$ .

- Here, the IV is eligibility for a policy and IV is equivalent to a social experiment that identifies the mean gain per participant who switches to participation in the program.

- The Wald estimator in this case is

$$\frac{E(Y | B = 1) - E(Y | B = 0)}{\Pr(D = 1 | B = 1) - \Pr(D = 1 | B = 0)}.$$

The IV weight for this estimator is a special case of equation (23):

$$\omega_{IV}(u_D | B) = \frac{E\left(B - E(B) \mid \hat{P}(Z) \geq u_D\right) \Pr\left(\hat{P}(Z) \geq u_D\right)}{\text{Cov}\left(B, \hat{P}(Z)\right)},$$

where  $\hat{P}(Z) = P(Z(1 + t\mathbf{1}[Z > 0]))^B P(Z)^{(1-B)}$ .

- Here, the IV is eligibility for a policy and IV is equivalent to a social experiment that identifies the mean gain per participant who switches to participation in the program.
- It is to be expected that this IV weight and  $\omega_{PRTE}$  are identical.

## Further Examples within the Extended Roy Model

- To gain a further understanding of how to construct the weights, and to understand how negative weights can arise, it is useful to return to the policy adoption model presented at the end of Slide 12.

## Further Examples within the Extended Roy Model

- To gain a further understanding of how to construct the weights, and to understand how negative weights can arise, it is useful to return to the policy adoption model presented at the end of Slide 12.
- The only unobservables in this model are in the outcome equations.

## Further Examples within the Extended Roy Model

- To gain a further understanding of how to construct the weights, and to understand how negative weights can arise, it is useful to return to the policy adoption model presented at the end of Slide 12.
- The only unobservables in this model are in the outcome equations.
- To simplify the analysis, we use an extended Roy model where the only unobservables are the unmeasured gains.



- In this framework, the cost  $C$  of adopting the policy is the same across all countries.

- In this framework, the cost  $C$  of adopting the policy is the same across all countries.
- Countries choose to adopt the policy if  $D^* > 0$  where  $D^*$  is the net benefit of adoption:  $D^* = (Y_1 - Y_0 - C)$  and  $ATE = E(\beta) = E(Y_1 - Y_0) = \mu_1 - \mu_0$ , while treatment on the treated is  $E(\beta | D = 1) = E(Y_1 - Y_0 | D = 1) = \mu_1 - \mu_0 + E(U_1 - U_0 | D = 1)$ .

- In this setting, the gross return to the country at the margin is  $C$ , i.e.,  
$$E(Y_1 - Y_0 \mid D^* = 0) = E(Y_1 - Y_0 \mid Y_1 - Y_0 = C) = C.$$

- In this setting, the gross return to the country at the margin is  $C$ , i.e.,  
$$E(Y_1 - Y_0 \mid D^* = 0) = E(Y_1 - Y_0 \mid Y_1 - Y_0 = C) = C.$$
- Recall that figure 1 presents the standard treatment parameters for the values of the choice parameter presented at the base of the figure.

- In this setting, the gross return to the country at the margin is  $C$ , i.e.,  
$$E(Y_1 - Y_0 \mid D^* = 0) = E(Y_1 - Y_0 \mid Y_1 - Y_0 = C) = C.$$
- Recall that figure 1 presents the standard treatment parameters for the values of the choice parameter presented at the base of the figure.
- Countries that adopt the policy are above average.

- In this setting, the gross return to the country at the margin is  $C$ , i.e.,  
$$E(Y_1 - Y_0 \mid D^* = 0) = E(Y_1 - Y_0 \mid Y_1 - Y_0 = C) = C.$$
- Recall that figure 1 presents the standard treatment parameters for the values of the choice parameter presented at the base of the figure.
- Countries that adopt the policy are above average.
- In a model where the cost varies (the generalized Roy model with  $U_C \neq 0$ ), and  $C$  is negatively correlated with the gain, adopting countries could be below average.

- In this setting, the gross return to the country at the margin is  $C$ , i.e.,  
$$E(Y_1 - Y_0 \mid D^* = 0) = E(Y_1 - Y_0 \mid Y_1 - Y_0 = C) = C.$$
- Recall that figure 1 presents the standard treatment parameters for the values of the choice parameter presented at the base of the figure.
- Countries that adopt the policy are above average.
- In a model where the cost varies (the generalized Roy model with  $U_C \neq 0$ ), and  $C$  is negatively correlated with the gain, adopting countries could be below average.
- We consider cases with discrete instruments and cases with continuous instruments.

- In this setting, the gross return to the country at the margin is  $C$ , i.e.,  
$$E(Y_1 - Y_0 \mid D^* = 0) = E(Y_1 - Y_0 \mid Y_1 - Y_0 = C) = C.$$
- Recall that figure 1 presents the standard treatment parameters for the values of the choice parameter presented at the base of the figure.
- Countries that adopt the policy are above average.
- In a model where the cost varies (the generalized Roy model with  $U_C \neq 0$ ), and  $C$  is negatively correlated with the gain, adopting countries could be below average.
- We consider cases with discrete instruments and cases with continuous instruments.
- We first turn to the discrete case.



## Discrete Instruments and Weights for LATE

- Consider what instrumental variables identify in the model of country policy adoption presented below figure 5.

## Discrete Instruments and Weights for LATE

- Consider what instrumental variables identify in the model of country policy adoption presented below figure 5.
- That figure presents three cases that we analyze in this section.

## Discrete Instruments and Weights for LATE

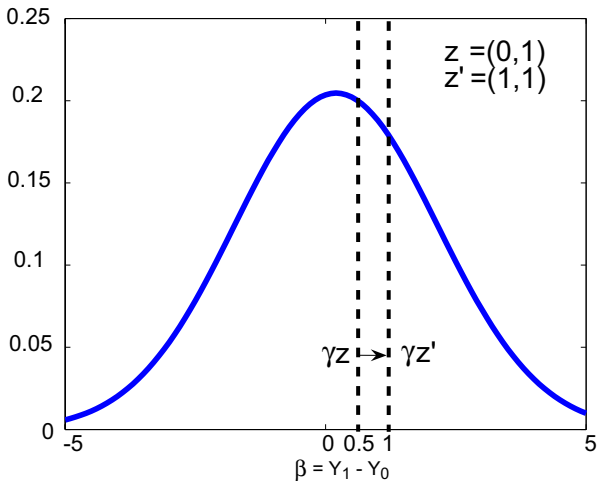
- Consider what instrumental variables identify in the model of country policy adoption presented below figure 5.
- That figure presents three cases that we analyze in this section.
- Let cost  $C = Z\gamma$  where instrument  $Z = (Z_1, Z_2)$ .

## Discrete Instruments and Weights for LATE

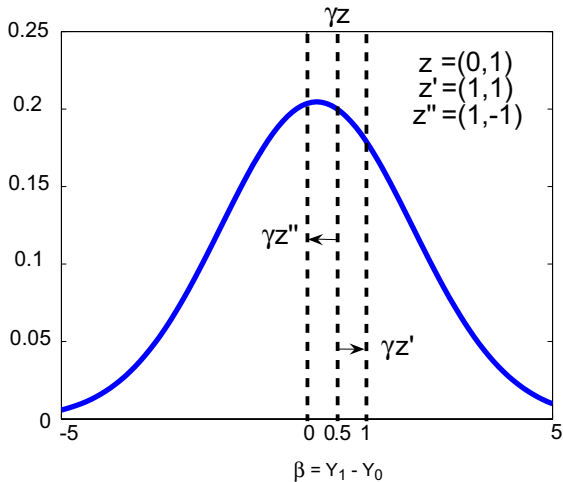
- Consider what instrumental variables identify in the model of country policy adoption presented below figure 5.
- That figure presents three cases that we analyze in this section.
- Let cost  $C = Z\gamma$  where instrument  $Z = (Z_1, Z_2)$ .
- Higher values of  $Z$  reduce the probability of adopting the policy if  $\gamma \geq 0$ , component by component.

Figure 5: Monotonicity: the extended Roy economy

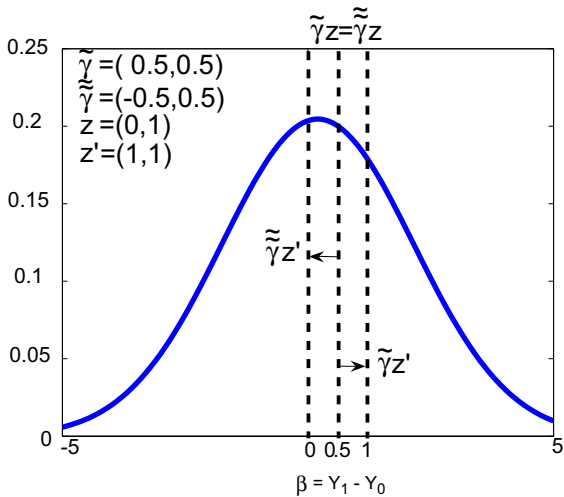
### A. Standard Case



## B. Changing $Z_1$ without Controlling for $Z_2$



### C. Random Coefficient Case



Outcomes

$$Y_1 = \alpha + \bar{\beta} + U_1$$

$$Y_0 = \alpha + U_0$$

Choice Model

$$D = \begin{cases} 1 & \text{if } Y_1 - Y_0 - \gamma Z \geq 0 \\ 0 & \text{if } Y_1 - Y_0 - \gamma Z < 0 \end{cases}$$

with  $\gamma Z = \gamma_1 Z_1 + \gamma_2 Z_2$

Parameterization

$$(U_1, U_0) \sim N(\mathbf{0}, \Sigma), \quad \Sigma = \begin{bmatrix} 1 & -0.9 \\ -0.9 & 1 \end{bmatrix}, \quad \alpha = 0.67, \quad \bar{\beta} = 0.2, \quad \gamma = (0.5, 0.5) \text{ (except in Case C)}$$

$$Z_1 = \{-1, 0, 1\} \text{ and } Z_2 = \{-1, 0, 1\}$$

A. Standard Case	B. Changing $Z_1$ without Controlling for $Z_2$	C. Random Coefficient Case
$z \rightarrow z'$ $z = (0, 1) \text{ and } z' = (1, 1)$	$z \rightarrow z' \text{ or } z \rightarrow z''$ $z = (0, 1), z' = (1, 1) \text{ and } z'' = (1, -1)$	$z \rightarrow z'$ $z = (0, 1) \text{ and } z' = (1, 1)$
		$\gamma$ is a random vector $\tilde{\gamma} = (0.5, 0.5) \text{ and } \tilde{\tilde{\gamma}} = (-0.5, 0.5)$ where $\tilde{\gamma}$ and $\tilde{\tilde{\gamma}}$ are two realizations of $\gamma$
$D(\gamma z) \geq D(\gamma z')$	$D(\gamma z) \geq D(\gamma z') \text{ or } D(\gamma z) < D(\gamma z'')$	$D(\tilde{\gamma} z) \geq D(\tilde{\tilde{\gamma}} z') \text{ and } D(\tilde{\gamma} z) < D(\tilde{\tilde{\gamma}} z')$
For all individuals	Depending on the value of $z'$ or $z''$	Depending on value of $\gamma$

Source: ?



- Consider the “standard” case depicted in figure 5A.

- Consider the “standard” case depicted in figure 5A.
- Increasing both components of discrete-valued  $Z$  raises costs and hence raises the benefit observed for the country at the margin by eliminating adoption in low return countries.

- Consider the “standard” case depicted in figure 5A.
- Increasing both components of discrete-valued  $Z$  raises costs and hence raises the benefit observed for the country at the margin by eliminating adoption in low return countries.
- It also reduces the probability that countries adopt the policy.

- Consider the “standard” case depicted in figure 5A.
- Increasing both components of discrete-valued  $Z$  raises costs and hence raises the benefit observed for the country at the margin by eliminating adoption in low return countries.
- It also reduces the probability that countries adopt the policy.
- In general a different country is at the margin when different instruments are used.

- Figure 6A plots the weights and figure 6B plots the components of the weights for the LATE values using  $P(Z)$  as an instrument for the distribution of discrete  $Z$  values shown at the base of the figure.

- Figure 6A plots the weights and figure 6B plots the components of the weights for the LATE values using  $P(Z)$  as an instrument for the distribution of discrete  $Z$  values shown at the base of the figure.
- Figure 6C presents the LATE parameter derived using  $P(Z)$  as an instrument.

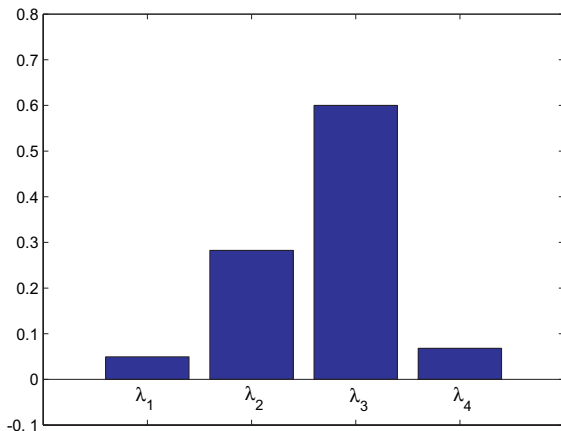
- Figure 6A plots the weights and figure 6B plots the components of the weights for the LATE values using  $P(Z)$  as an instrument for the distribution of discrete  $Z$  values shown at the base of the figure.
- Figure 6C presents the LATE parameter derived using  $P(Z)$  as an instrument.
- The weights are positive as predicted from equation (15) when  $J(Z) = P(Z)$ .

- Figure 6A plots the weights and figure 6B plots the components of the weights for the LATE values using  $P(Z)$  as an instrument for the distribution of discrete  $Z$  values shown at the base of the figure.
- Figure 6C presents the LATE parameter derived using  $P(Z)$  as an instrument.
- The weights are positive as predicted from equation (15) when  $J(Z) = P(Z)$ .
- Thus, the monotonicity condition for the weights in terms of  $u_D$  is satisfied.

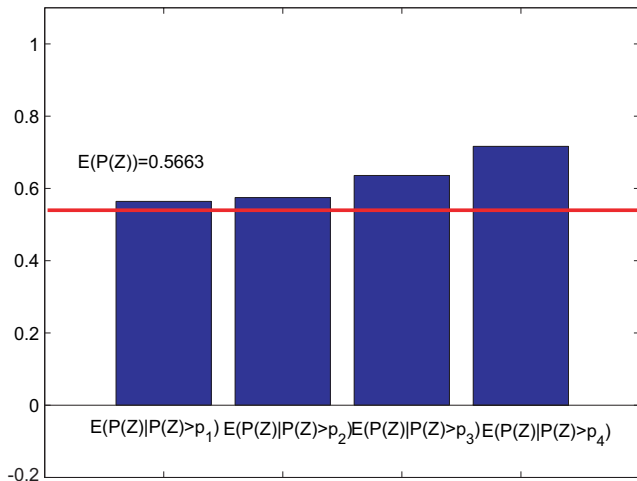


Figure 6: IV Weight and Its Components under Discrete Instruments when  $P(Z)$  is the Instrument: The Extended Roy Economy

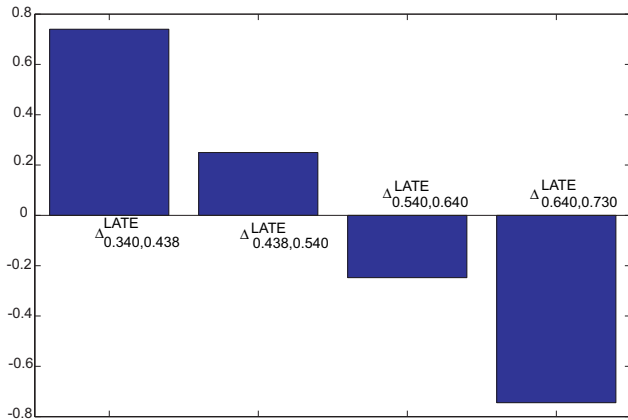
### A. IV Weights



## B. $E(P(Z)|P(Z) > p_\ell)$ and $E(P(Z))$



## C. Local Average Treatment Effects



The model is the same as the one presented below Figure 5.

$$\text{ATE} = 0.2, \text{TT} = 0.5942, \text{TUT} = -0.4823 \text{ and } \Delta_{P(Z)}^{\text{IV}} = \sum_{\ell=1}^{K-1} \Delta^{\text{LATE}}(p_{\ell}, p_{\ell+1}) \lambda_{\ell} = -0.09$$

$$\begin{aligned} \Delta^{\text{LATE}}(p_{\ell}, p_{\ell+1}) &= \frac{E(Y|P(Z) = p_{\ell+1}) - E(Y|P(Z) = p_{\ell})}{p_{\ell+1} - p_{\ell}} \\ &= \frac{\bar{\beta}(p_{\ell+1} - p_{\ell}) + \sigma_{U_1 - U_0} \left( \phi \left( \Phi^{-1}(1 - p_{\ell+1}) \right) - \phi \left( \Phi^{-1}(1 - p_{\ell}) \right) \right)}{p_{\ell+1} - p_{\ell}} \\ \lambda_{\ell} &= (p_{\ell+1} - p_{\ell}) \frac{\sum_{i=1}^K (p_i - E(P(Z))) \sum_{t>\ell}^K f(p_i, p_t)}{\text{Cov}(Z_1, D)} = (p_{\ell+1} - p_{\ell}) \frac{\sum_{t>\ell}^K (p_t - E(P(Z))) f(p_t)}{\text{Cov}(Z_1, D)} \end{aligned}$$

Joint Probability Distribution of  $(Z_1, Z_2)$  and the Propensity Score

(joint probabilities in ordinary type ( $\Pr(Z_1 = z_1, Z_2 = z_2)$ ); propensity score in italics ( $\Pr(D = 1|Z_1 = z_1, Z_2 = z_2)$ ))

$Z_1 \backslash Z_2$	-1	0	1
-1	0.02	0.02	0.36
0	<i>0.7309</i>	<i>0.6402</i>	<i>0.5409</i>
1	0.3	0.01	0.03
	<i>0.6402</i>	<i>0.5409</i>	<i>0.4388</i>
	0.2	0.05	0.01
	<i>0.5409</i>	<i>0.4388</i>	<i>0.3408</i>

$$\text{Cov}(Z_1, Z_2) = -0.5468$$

Source: ?

- The outcome and choice parameters are the same as those used to generate figures 1 and 5.

- The outcome and choice parameters are the same as those used to generate figures 1 and 5.
- The LATE parameters for each interval of  $P$  values are presented in a table just below the figures.

- The outcome and choice parameters are the same as those used to generate figures 1 and 5.
- The LATE parameters for each interval of  $P$  values are presented in a table just below the figures.
- There are four LATE parameters corresponding to the five distinct values of the propensity score for that value.

- The outcome and choice parameters are the same as those used to generate figures 1 and 5.
- The LATE parameters for each interval of  $P$  values are presented in a table just below the figures.
- There are four LATE parameters corresponding to the five distinct values of the propensity score for that value.
- The LATE parameters exhibit the declining pattern with  $u_D$  predicted by the Roy model.



- A case producing negative weights is depicted in figure 5 B.

- A case producing negative weights is depicted in figure 5 B.
- In that graph, the same  $Z$  is used to generate the choices as is used to generate figure 1B.

- A case producing negative weights is depicted in figure 5 B.
- In that graph, the same  $Z$  is used to generate the choices as is used to generate figure 1B.
- However, in this case, the analyst uses  $Z_1$  as the instrument.

- A case producing negative weights is depicted in figure 5 B.
- In that graph, the same  $Z$  is used to generate the choices as is used to generate figure 1B.
- However, in this case, the analyst uses  $Z_1$  as the instrument.
- $Z_1$  and  $Z_2$  are negatively dependent and  $E(Z_1 | P(Z) > u_D)$  is not monotonic in  $u_D$ .

- A case producing negative weights is depicted in figure 5 B.
- In that graph, the same  $Z$  is used to generate the choices as is used to generate figure 1B.
- However, in this case, the analyst uses  $Z_1$  as the instrument.
- $Z_1$  and  $Z_2$  are negatively dependent and  $E(Z_1 | P(Z) > u_D)$  is not monotonic in  $u_D$ .
- This nonmonotonicity is evident in figure 7B.

- A case producing negative weights is depicted in figure 5 B.
- In that graph, the same  $Z$  is used to generate the choices as is used to generate figure 1B.
- However, in this case, the analyst uses  $Z_1$  as the instrument.
- $Z_1$  and  $Z_2$  are negatively dependent and  $E(Z_1 | P(Z) > u_D)$  is not monotonic in  $u_D$ .
- This nonmonotonicity is evident in figure 7B.
- It produces the pattern of negative weights shown in figure 7A.

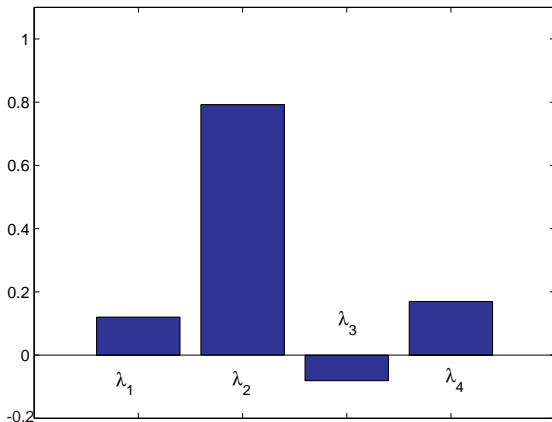
- A case producing negative weights is depicted in figure 5 B.
- In that graph, the same  $Z$  is used to generate the choices as is used to generate figure 1B.
- However, in this case, the analyst uses  $Z_1$  as the instrument.
- $Z_1$  and  $Z_2$  are negatively dependent and  $E(Z_1 | P(Z) > u_D)$  is not monotonic in  $u_D$ .
- This nonmonotonicity is evident in figure 7B.
- It produces the pattern of negative weights shown in figure 7A.
- These are associated with two way flows.

- A case producing negative weights is depicted in figure 5 B.
- In that graph, the same  $Z$  is used to generate the choices as is used to generate figure 1B.
- However, in this case, the analyst uses  $Z_1$  as the instrument.
- $Z_1$  and  $Z_2$  are negatively dependent and  $E(Z_1 | P(Z) > u_D)$  is not monotonic in  $u_D$ .
- This nonmonotonicity is evident in figure 7B.
- It produces the pattern of negative weights shown in figure 7A.
- These are associated with two way flows.
- Increasing  $Z_1$  controlling for  $Z_2$  reduces the probability of country policy adoption.

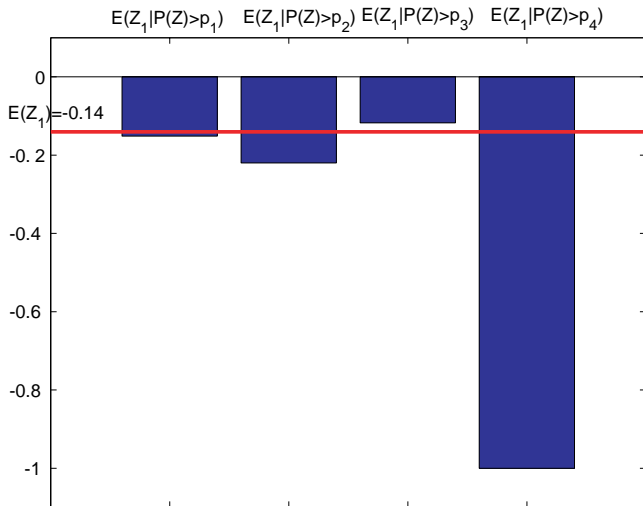


Figure 7: IV Weight and Its Components under Discrete Instruments when  $Z_1$  is the Instrument: The Extended Roy Economy

### A. IV Weights



## B. $E(Z_1|P(Z) > p_\ell)$ and $E(Z_1)$



The model is the same as the one presented below Figure 5. The values of the treatment parameters are the same as the ones presented below Figure 6.

$$\Delta_{Z_1}^{IV} = \sum_{\ell=1}^{K-1} \Delta^{\text{LATE}}(p_\ell, p_{\ell+1}) \lambda_\ell = 0.1833$$

$$\lambda_\ell = (p_{\ell+1} - p_\ell) \frac{\sum_{i=1}^I (z_{1,i} - E(Z_1)) \sum_{t>\ell}^K f(z_{1,i}, p_t)}{\text{Cov}(Z_1, D)}$$

Joint Probability Distribution of  $(Z_1, Z_2)$  and the Propensity Score  
 (joint probabilities in ordinary type ( $\Pr(Z_1 = z_1, Z_2 = z_2)$ ); propensity score in italics  
 ( $\Pr(D = 1|Z_1 = z_1, Z_2 = z_2)$ ))

$Z_1 \backslash Z_2$	-1	0	1
-1	0.02 <i>0.7309</i>	0.02 <i>0.6402</i>	0.36 <i>0.5409</i>
0	0.3 <i>0.6402</i>	0.01 <i>0.5409</i>	0.03 <i>0.4388</i>
1	0.2 <i>0.5409</i>	0.05 <i>0.4388</i>	0.01 <i>0.3408</i>

$$\text{Cov}(Z_1, Z_2) = -0.5468$$

Source: ?.

- However, we do not condition on  $Z_2$  in constructing this figure.

- However, we do not condition on  $Z_2$  in constructing this figure.
- $Z_2$  is floating.

- However, we do not condition on  $Z_2$  in constructing this figure.
- $Z_2$  is floating.
- Two way flows are induced by uncontrolled variation in  $Z_2$ .

- However, we do not condition on  $Z_2$  in constructing this figure.
- $Z_2$  is floating.
- Two way flows are induced by uncontrolled variation in  $Z_2$ .
- For some units, the strength of the associated variation in  $Z_2$  offsets the increase in  $Z_1$  and for other units it does not.

- However, we do not condition on  $Z_2$  in constructing this figure.
- $Z_2$  is floating.
- Two way flows are induced by uncontrolled variation in  $Z_2$ .
- For some units, the strength of the associated variation in  $Z_2$  offsets the increase in  $Z_1$  and for other units it does not.
- Observe that the LATE parameters defined using  $P(Z)$  are the same in both examples.



- However, we do not condition on  $Z_2$  in constructing this figure.
- $Z_2$  is floating.
- Two way flows are induced by uncontrolled variation in  $Z_2$ .
- For some units, the strength of the associated variation in  $Z_2$  offsets the increase in  $Z_1$  and for other units it does not.
- Observe that the LATE parameters defined using  $P(Z)$  are the same in both examples.
- They are just weighted differently.

- However, we do not condition on  $Z_2$  in constructing this figure.
- $Z_2$  is floating.
- Two way flows are induced by uncontrolled variation in  $Z_2$ .
- For some units, the strength of the associated variation in  $Z_2$  offsets the increase in  $Z_1$  and for other units it does not.
- Observe that the LATE parameters defined using  $P(Z)$  are the same in both examples.
- They are just weighted differently.
- We discuss the random coefficient choice model generating figure 5C in Slide 370.

- The IV estimator does not identify ATE, TT or TUT (given at the bottom of figure 6C).

- The IV estimator does not identify ATE, TT or TUT (given at the bottom of figure 6C).
- Conditioning on  $Z_2$  produces positive weights.

- The IV estimator does not identify ATE, TT or TUT (given at the bottom of figure 6C).
- Conditioning on  $Z_2$  produces positive weights.
- This is illustrated in the weights shown in table 6 that condition on  $Z_2$  using the same model that generated figure 6.

- The IV estimator does not identify ATE, TT or TUT (given at the bottom of figure 6C).
- Conditioning on  $Z_2$  produces positive weights.
- This is illustrated in the weights shown in table 6 that condition on  $Z_2$  using the same model that generated figure 6.
- Conditioning on  $Z_2$  effectively converts the problem back into one with a scalar instrument and the weights are positive for that case.

Table 6: The Conditional Instrumental Variable Estimator

	$Z_2 = -1$	$Z_2 = 0$	$Z_2 = 1$
$P(-1, Z_2) = p_3$	0.7309	0.6402	0.5409
$P(0, Z_2) = p_2$	0.6402	0.5409	0.4388
$P(1, Z_2) = p_1$	0.5409	0.4388	0.3408
$\lambda_1$	0.8418	0.5384	0.2860
$\lambda_2$	0.1582	0.4616	0.7140
$\Delta^{\text{LATE}}(p_1, p_2)$	-0.2475	0.2497	0.7470
$\Delta^{\text{LATE}}(p_2, p_3)$	-0.7448	-0.2475	0.2497
$\Delta_{Z_1 Z_2=z_2}^{\text{IV}}$	-0.3262	0.0202	0.3920

$(\Delta_{Z_1|Z_2=z_2}^{\text{IV}})$  and Conditional Local Average Treatment Effect  $(\Delta^{\text{LATE}}(p_\ell, p_{\ell+1}|Z_2 = z_2))$  when  $Z_1$  is the Instrument (given  $Z_2 = z_2$ )

The model is the same as the one presented below Figure 2

$$\Delta_{Z_1|Z_2=z_2}^{IV} = \sum_{\ell=1}^{I-1} \Delta^{\text{LATE}}(p_\ell, p_{\ell+1}|Z_2 = z_2) \lambda_{\ell|Z_2=z_2} = \sum_{\ell=1}^{I-1} \Delta^{\text{LATE}}(p_\ell, p_{\ell+1}|Z_2 = z_2) \lambda_{\ell|Z_2=z_2}$$

$$\Delta^{\text{LATE}}(p_\ell, p_{\ell+1}|Z_2 = z_2) = \frac{E(Y|P(Z) = p_{\ell+1}, Z_2 = z_2) - E(Y|P(Z) = p_\ell, Z_2 = z_2)}{p_{\ell+1} - p_\ell}$$

$$\begin{aligned} \lambda_{\ell|Z_2=z_2} &= (p_{\ell+1} - p_\ell) \frac{\sum_{i=1}^I (z_{1,i} - E(Z_1|Z_2 = z_2)) \sum_{t>\ell}^I f(z_{1,i}, p_t|Z_2 = z_2)}{\text{Cov}(Z_1, D)} \\ &= (p_{\ell+1} - p_\ell) \frac{\sum_{t>\ell}^I (z_{1,t} - E(Z_1|Z_2 = z_2)) f(z_{1,t}, p_t|Z_2 = z_2)}{\text{Cov}(Z_1, D)} \end{aligned}$$



Probability Distribution of  $Z_1$  Conditional on  $Z_2$   
( $\Pr(Z_1 = z_1 | Z_2 = z_2)$ )

$z_1$	$\Pr(Z_1 = z_1   Z_2 = -1)$	$\Pr(Z_1 = z_1   Z_2 = 0)$	$\Pr(Z_1 = z_1   Z_2 = 1)$
-1	0.0385	0.25	0.9
0	0.5769	0.125	0.075
1	0.3846	0.625	0.025

Source: ?

- From Yitzhaki's analysis, for any sample size, a regression of  $Y$  on  $P$  identifies a weighted average of slopes based on ordered regressors:

$$\frac{E(Y_\ell | p_\ell) - E(Y_{\ell-1} | p_{\ell-1})}{p_\ell - p_{\ell-1}}$$

where  $p_\ell > p_{\ell-1}$  and the weights are the positive Yitzhaki–Imbens–Angrist weights derived in ?? or in ?.

- From Yitzhaki's analysis, for any sample size, a regression of  $Y$  on  $P$  identifies a weighted average of slopes based on ordered regressors:

$$\frac{E(Y_\ell | p_\ell) - E(Y_{\ell-1} | p_{\ell-1})}{p_\ell - p_{\ell-1}}$$

where  $p_\ell > p_{\ell-1}$  and the weights are the positive Yitzhaki–Imbens–Angrist weights derived in ?? or in ?.

- The weights are positive whether or not monotonicity condition (IV-3) holds.

- From Yitzhaki's analysis, for any sample size, a regression of  $Y$  on  $P$  identifies a weighted average of slopes based on ordered regressors:

$$\frac{E(Y_\ell | p_\ell) - E(Y_{\ell-1} | p_{\ell-1})}{p_\ell - p_{\ell-1}}$$

where  $p_\ell > p_{\ell-1}$  and the weights are the positive Yitzhaki–Imbens–Angrist weights derived in ?? or in ?.

- The weights are positive whether or not monotonicity condition (IV-3) holds.
- If monotonicity holds, IV is a weighted average of LATEs.

- From Yitzhaki's analysis, for any sample size, a regression of  $Y$  on  $P$  identifies a weighted average of slopes based on ordered regressors:

$$\frac{E(Y_\ell | p_\ell) - E(Y_{\ell-1} | p_{\ell-1})}{p_\ell - p_{\ell-1}}$$

where  $p_\ell > p_{\ell-1}$  and the weights are the positive Yitzhaki–Imbens–Angrist weights derived in ?? or in ?.

- The weights are positive whether or not monotonicity condition (IV-3) holds.
- If monotonicity holds, IV is a weighted average of LATEs.
- Otherwise it is just a weighted average of ordered (by  $p_\ell$ ) estimators consistent with two way flows.

- From Yitzhaki's analysis, for any sample size, a regression of  $Y$  on  $P$  identifies a weighted average of slopes based on ordered regressors:

$$\frac{E(Y_\ell | p_\ell) - E(Y_{\ell-1} | p_{\ell-1})}{p_\ell - p_{\ell-1}}$$

where  $p_\ell > p_{\ell-1}$  and the weights are the positive Yitzhaki–Imbens–Angrist weights derived in ?? or in ?.

- The weights are positive whether or not monotonicity condition (IV-3) holds.
- If monotonicity holds, IV is a weighted average of LATEs.
- Otherwise it is just a weighted average of ordered (by  $p_\ell$ ) estimators consistent with two way flows.
- We next discuss continuous instruments.

## Continuous Instruments

- For the case of continuous  $Z$ , we present a parallel analysis for the weights associated with the MTE.

## Continuous Instruments

- For the case of continuous  $Z$ , we present a parallel analysis for the weights associated with the MTE.
- Figure 8 plots  $E(Y | P(Z))$  and MTE for the extended Roy models generated by the parameters displayed at the base of the figure.



## Continuous Instruments

- For the case of continuous  $Z$ , we present a parallel analysis for the weights associated with the MTE.
- Figure 8 plots  $E(Y | P(Z))$  and MTE for the extended Roy models generated by the parameters displayed at the base of the figure.
- In cases I and II,  $\beta \perp\!\!\!\perp D$ , so  $\Delta^{\text{MTE}}(u_D)$  is constant in  $u_D$ .

## Continuous Instruments

- For the case of continuous  $Z$ , we present a parallel analysis for the weights associated with the MTE.
- Figure 8 plots  $E(Y | P(Z))$  and MTE for the extended Roy models generated by the parameters displayed at the base of the figure.
- In cases I and II,  $\beta \perp\!\!\!\perp D$ , so  $\Delta^{\text{MTE}}(u_D)$  is constant in  $u_D$ .
- In case I, this is trivial since  $\beta$  is a constant.

## Continuous Instruments

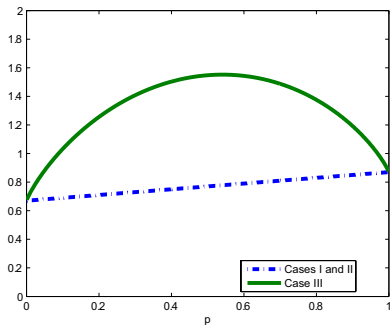
- For the case of continuous  $Z$ , we present a parallel analysis for the weights associated with the MTE.
- Figure 8 plots  $E(Y | P(Z))$  and MTE for the extended Roy models generated by the parameters displayed at the base of the figure.
- In cases I and II,  $\beta \perp\!\!\!\perp D$ , so  $\Delta^{\text{MTE}}(u_D)$  is constant in  $u_D$ .
- In case I, this is trivial since  $\beta$  is a constant.
- In case II,  $\beta$  is random but selection into  $D$  does not depend on  $\beta$ .

## Continuous Instruments

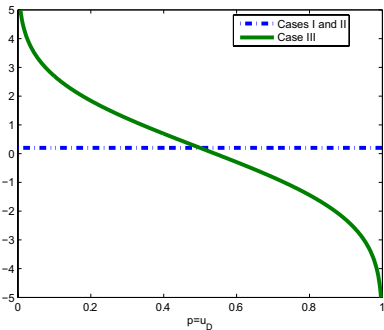
- For the case of continuous  $Z$ , we present a parallel analysis for the weights associated with the MTE.
- Figure 8 plots  $E(Y | P(Z))$  and MTE for the extended Roy models generated by the parameters displayed at the base of the figure.
- In cases I and II,  $\beta \perp\!\!\!\perp D$ , so  $\Delta^{\text{MTE}}(u_D)$  is constant in  $u_D$ .
- In case I, this is trivial since  $\beta$  is a constant.
- In case II,  $\beta$  is random but selection into  $D$  does not depend on  $\beta$ .
- Case III is the model with essential heterogeneity ( $\beta \not\perp\!\!\!\perp D$ ).

Figure 8: Conditional Expectation of  $Y$  on  $P(Z)$  and the Marginal Treatment Effect (MTE)

A.  $E(Y|P(Z) = p)$



B.  $\Delta^{\text{MTE}}(u_D)$



### Outcomes

$$Y_1 = \alpha + \bar{\beta} + U_1$$

$$Y_0 = \alpha + U_0$$

### Choice Model

$$D = \begin{cases} 1 & \text{if } D^* \geq 0 \\ 0 & \text{if } D^* < 0 \end{cases}$$

Case I	Case II	Case III
$U_1 = U_0$ $\bar{\beta} = \text{ATE} = \text{TT} = \text{TUT} = \text{IV}$	$U_1 - U_0 \perp\!\!\!\perp D$ $\bar{\beta} = \text{ATE} = \text{TT} = \text{TUT} = \text{IV}$	$U_1 - U_0 \not\perp\!\!\!\perp D$ $\bar{\beta} = \text{ATE} \neq \text{TT} \neq \text{TUT} \neq \text{IV}$

## Parameterization

Cases I, II and III	Cases II and III	Case III
$\alpha = 0.67$ $\bar{\beta} = 0.2$	$(U_1, U_0) \sim N(\mathbf{0}, )$ with $\Sigma = \begin{bmatrix} 1 & -0.9 \\ -0.9 & 1 \end{bmatrix}$	$D^* = Y_1 - Y_0 - \gamma Z$ $Z \sim N(\mu_Z, \Sigma_Z)$ $\mu_Z = (2, -2)$ and $\Sigma_Z = \begin{bmatrix} 9 & -2 \\ -2 & 9 \end{bmatrix}$ $\gamma = (0.5, 0.5)$

Source: ?

- The graph (figure 8A) depicts  $E(Y | P(Z))$  in the three cases.



- The graph (figure 8A) depicts  $E(Y | P(Z))$  in the three cases.
- Cases I and II make  $E(Y | P(Z))$  linear in  $P(Z)$ .

- The graph (figure 8A) depicts  $E(Y | P(Z))$  in the three cases.
- Cases I and II make  $E(Y | P(Z))$  linear in  $P(Z)$ .
- Case III is nonlinear in  $P(Z)$ .

- The graph (figure 8A) depicts  $E(Y | P(Z))$  in the three cases.
- Cases I and II make  $E(Y | P(Z))$  linear in  $P(Z)$ .
- Case III is nonlinear in  $P(Z)$ .
- This arises when  $\beta \not\propto D$ .

- The graph (figure 8A) depicts  $E(Y | P(Z))$  in the three cases.
- Cases I and II make  $E(Y | P(Z))$  linear in  $P(Z)$ .
- Case III is nonlinear in  $P(Z)$ .
- This arises when  $\beta \not\propto D$ .
- The derivative of  $E(Y | P(Z))$  is presented in figure 8B.

- The graph (figure 8A) depicts  $E(Y | P(Z))$  in the three cases.
- Cases I and II make  $E(Y | P(Z))$  linear in  $P(Z)$ .
- Case III is nonlinear in  $P(Z)$ .
- This arises when  $\beta \not\propto D$ .
- The derivative of  $E(Y | P(Z))$  is presented in figure 8B.
- It is a constant for cases I and II (flat MTE) but declining in  $U_D = P(Z)$  for the case with selection on the gain.

- The graph (figure 8A) depicts  $E(Y | P(Z))$  in the three cases.
- Cases I and II make  $E(Y | P(Z))$  linear in  $P(Z)$ .
- Case III is nonlinear in  $P(Z)$ .
- This arises when  $\beta \not\perp D$ .
- The derivative of  $E(Y | P(Z))$  is presented in figure 8B.
- It is a constant for cases I and II (flat MTE) but declining in  $U_D = P(Z)$  for the case with selection on the gain.
- A simple test for linearity in  $P(Z)$  in the outcome equation reveals whether or not the analyst is in cases I and II ( $\beta \perp D$ ) or case III ( $\beta \not\perp D$ ).

- The graph (figure 8A) depicts  $E(Y | P(Z))$  in the three cases.
- Cases I and II make  $E(Y | P(Z))$  linear in  $P(Z)$ .
- Case III is nonlinear in  $P(Z)$ .
- This arises when  $\beta \not\perp D$ .
- The derivative of  $E(Y | P(Z))$  is presented in figure 8B.
- It is a constant for cases I and II (flat MTE) but declining in  $U_D = P(Z)$  for the case with selection on the gain.
- A simple test for linearity in  $P(Z)$  in the outcome equation reveals whether or not the analyst is in cases I and II ( $\beta \perp D$ ) or case III ( $\beta \not\perp D$ ).
- These cases are the extended Roy counterparts to  $E(Y | P(Z) = p)$  and MTE shown for the generalized Roy model in figures 3A and 3B.

- MTE gives the mean marginal return for persons who have utility  $P(Z) = u_D$ .



- MTE gives the mean marginal return for persons who have utility  $P(Z) = u_D$ .
- Thus,  $P(Z) = u_D$  is the margin of indifference.

- MTE gives the mean marginal return for persons who have utility  $P(Z) = u_D$ .
- Thus,  $P(Z) = u_D$  is the margin of indifference.
- Those with low  $u_D$  values have high returns.

- MTE gives the mean marginal return for persons who have utility  $P(Z) = u_D$ .
- Thus,  $P(Z) = u_D$  is the margin of indifference.
- Those with low  $u_D$  values have high returns.
- Those with high  $u_D$  values have low returns.

- MTE gives the mean marginal return for persons who have utility  $P(Z) = u_D$ .
- Thus,  $P(Z) = u_D$  is the margin of indifference.
- Those with low  $u_D$  values have high returns.
- Those with high  $u_D$  values have low returns.
- Figure 8 highlights that, in the general case, MTE (and LATE) identify average returns for persons at the margin of indifference at different levels of the mean utility function ( $P(Z)$ ).

- Figure 9 plots MTE and LATE for different intervals of  $u_D$  using the model generating figure 8.

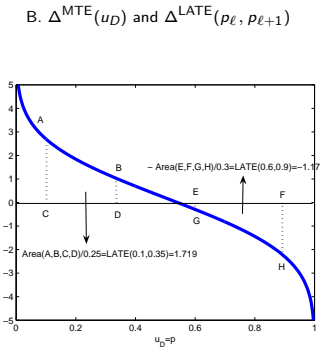
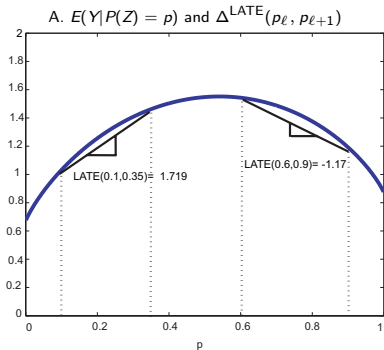
- Figure 9 plots MTE and LATE for different intervals of  $u_D$  using the model generating figure 8.
- LATE is the chord of  $E(Y | P(Z))$  evaluated at different points.

- Figure 9 plots MTE and LATE for different intervals of  $u_D$  using the model generating figure 8.
- LATE is the chord of  $E(Y | P(Z))$  evaluated at different points.
- The relationship between LATE and MTE is depicted in figure 9B.

- Figure 9 plots MTE and LATE for different intervals of  $u_D$  using the model generating figure 8.
- LATE is the chord of  $E(Y | P(Z))$  evaluated at different points.
- The relationship between LATE and MTE is depicted in figure 9B.
- LATE is the integral under the MTE curve divided by the difference between the upper and lower limits.



Figure 9: The Local Average Treatment Effect



$$\Delta^{\text{LATE}}(p_\ell, p_{\ell+1}) = \frac{E(Y|P(Z) = p_{\ell+1}) - E(Y|P(Z) = p_\ell)}{p_{\ell+1} - p_\ell} = \frac{\int_{p_\ell}^{p_{\ell+1}} \Delta^{\text{MTE}}(u_D) du_D}{p_{\ell+1} - p_\ell}$$

$$\Delta^{\text{LATE}}(0.6, 0.9) = -1.17$$

$$\Delta^{\text{LATE}}(0.1, 0.35) = 1.719$$

Outcomes

$$Y_1 = \alpha + \bar{\beta} + U_1$$

$$Y_0 = \alpha + U_0$$

Choice Model

$$D = \begin{cases} 1 & \text{if } D^* \geq 0 \\ 0 & \text{if } D^* < 0 \end{cases}$$

with  $D^* = Y_1 - Y_0 - \gamma Z$

Parameterization

$$(U_1, U_0) \sim N(\mathbf{0}, \Sigma) \text{ and } Z \sim N(\mu_Z, \sigma_Z^2)$$

$$\Sigma = \begin{bmatrix} 1 & -0.9 \\ -0.9 & 1 \end{bmatrix}, \mu_Z = (2, -2) \text{ and } \sigma_Z^2 = \begin{bmatrix} 9 & -2 \\ -2 & 9 \end{bmatrix}$$

$$\alpha = 0.67, \bar{\beta} = 0.2, \gamma = (0.5, 0.5)$$

Source: ?

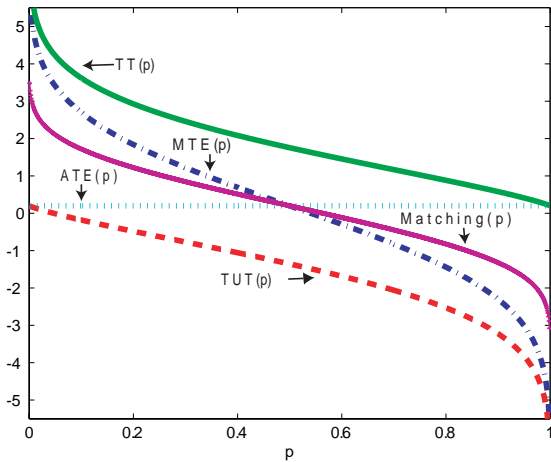
- The treatment parameters associated with case III are plotted in figure 10.

- The treatment parameters associated with case III are plotted in figure 10.
- The MTE is the same as that presented in figure 8.

- The treatment parameters associated with case III are plotted in figure 10.
- The MTE is the same as that presented in figure 8.
- ATE has the same value for all  $p$ .

- The treatment parameters associated with case III are plotted in figure 10.
- The MTE is the same as that presented in figure 8.
- ATE has the same value for all  $p$ .
- The effect of treatment on the treated for  $P(Z) = p$ ,  $\Delta^{TT}(p) = E(Y_1 - Y_0 \mid D = 1, P(Z) = p)$  declines in  $p$  (equivalently it declines in  $u_D$ ).

Figure 10: Treatment Parameters and OLS/Matching as a function of  $P(Z) = p$





Parameter	Definition	Under Assumptions (*)
Marginal Treatment Effect	$E[Y_1 - Y_0   D^* = 0, P(Z) = p]$	$\bar{\beta} + \sigma_{U_1-U_0} \Phi^{-1}(1-p)$
Average Treatment Effect	$E[Y_1 - Y_0   P(Z) = p]$	$\bar{\beta}$
Treatment on the Treated	$E[Y_1 - Y_0   D^* \geq 0, P(Z) = p]$	$\bar{\beta} + \sigma_{U_1-U_0} \frac{\phi(\Phi^{-1}(1-p))}{1-p}$
Treatment on the Untreated	$E[Y_1 - Y_0   D^* < 0, P(Z) = p]$	$\bar{\beta} - \sigma_{U_1-U_0} \frac{\phi(\Phi^{-1}(1-p))}{1-p}$
OLS/Matching on $P(Z)$	$E[Y_1   D^* \geq 0, P(Z) = p] - E[Y_0   D^* < 0, P(Z) = p]$	$\bar{\beta} + \left( \frac{\sigma_{U_1-U_0}^2}{\sqrt{\sigma_{U_1-U_0}^2}} \right) \left( \frac{1-2p}{p(1-p)} \right) \phi(\Phi^{-1}(1-p))$

Note:  $\Phi(\cdot)$  and  $\phi(\cdot)$  represent the cdf and pdf of a standard normal distribution, respectively.  $\Phi^{-1}(\cdot)$  represents the inverse of  $\Phi(\cdot)$ .

(\*): The model in this case is the same as the one presented below Figure 9.

Source: ?

- Treatment on the untreated given  $p$ ,  
 $TUT(p) = \Delta^{TUT}(p) = E(Y_1 - Y_0 \mid D = 0, P(Z) = p)$  also declines in  $p$ .

$$LATE(p, p') = \frac{\Delta^{TT}(p')p' - \Delta^{TT}(p)p}{p' - p}, \quad p' \neq p$$
$$MTE = \frac{\partial[\Delta^{TT}(p)p]}{\partial p}.$$

- Treatment on the untreated given  $p$ ,  
 $TUT(p) = \Delta^{TUT}(p) = E(Y_1 - Y_0 \mid D = 0, P(Z) = p)$  also declines in  $p$ .

$$LATE(p, p') = \frac{\Delta^{TT}(p')p' - \Delta^{TT}(p)p}{p' - p}, \quad p' \neq p$$
$$MTE = \frac{\partial[\Delta^{TT}(p)p]}{\partial p}.$$

- We can generate all of the treatment parameters from  $\Delta^{TT}(p)$ .

- Matching on  $P = p$  (which is equivalent to nonparametric regression given  $P = p$ ) produces a biased estimator of  $TT(p)$ .

- Matching on  $P = p$  (which is equivalent to nonparametric regression given  $P = p$ ) produces a biased estimator of  $TT(p)$ .
- Matching assumes a flat MTE (average return equals marginal return).

- Matching on  $P = p$  (which is equivalent to nonparametric regression given  $P = p$ ) produces a biased estimator of  $TT(p)$ .
- Matching assumes a flat MTE (average return equals marginal return).
- Therefore it is systematically biased for  $\Delta^{TT}(p)$  in a model with essential heterogeneity.

- Matching on  $P = p$  (which is equivalent to nonparametric regression given  $P = p$ ) produces a biased estimator of  $TT(p)$ .
- Matching assumes a flat MTE (average return equals marginal return).
- Therefore it is systematically biased for  $\Delta^{TT}(p)$  in a model with essential heterogeneity.
- Making observables alike makes the unobservables dissimilar.

- Matching on  $P = p$  (which is equivalent to nonparametric regression given  $P = p$ ) produces a biased estimator of  $TT(p)$ .
- Matching assumes a flat MTE (average return equals marginal return).
- Therefore it is systematically biased for  $\Delta^{TT}(p)$  in a model with essential heterogeneity.
- Making observables alike makes the unobservables dissimilar.
- Holding  $p$  constant across treatment and control groups understates  $TT(p)$  for low values of  $p$  and overstates it for high values of  $p$ .



- Matching on  $P = p$  (which is equivalent to nonparametric regression given  $P = p$ ) produces a biased estimator of  $TT(p)$ .
- Matching assumes a flat MTE (average return equals marginal return).
- Therefore it is systematically biased for  $\Delta^{TT}(p)$  in a model with essential heterogeneity.
- Making observables alike makes the unobservables dissimilar.
- Holding  $p$  constant across treatment and control groups understates  $TT(p)$  for low values of  $p$  and overstates it for high values of  $p$ .
- We develop this point further when we discuss matching in Slide 675.

- Figure 11 plots the MTE (as a function of  $u_D$  where  $u_D = F_V(v)$ ), the weights for ATE, TT and TUT and the IV weights using  $Z_1$  as the instrument for the model used to generate figure 9.

- Figure 11 plots the MTE (as a function of  $u_D$  where  $u_D = F_V(v)$ ), the weights for ATE, TT and TUT and the IV weights using  $Z_1$  as the instrument for the model used to generate figure 9.
- The distribution of the  $Z$  is assumed to be normal with generating parameters given at the base of figure 9.

- Figure 11 plots the MTE (as a function of  $u_D$  where  $u_D = F_V(v)$ ), the weights for ATE, TT and TUT and the IV weights using  $Z_1$  as the instrument for the model used to generate figure 9.
- The distribution of the  $Z$  is assumed to be normal with generating parameters given at the base of figure 9.
- The IV weight for normal  $Z$  is always nonnegative even if we use only one coordinate of vector  $Z$ .

- Figure 11 plots the MTE (as a function of  $u_D$  where  $u_D = F_V(v)$ ), the weights for ATE, TT and TUT and the IV weights using  $Z_1$  as the instrument for the model used to generate figure 9.
- The distribution of the  $Z$  is assumed to be normal with generating parameters given at the base of figure 9.
- The IV weight for normal  $Z$  is always nonnegative even if we use only one coordinate of vector  $Z$ .
- This is a consequence of the monotonicity of  $E(Z_j | P(Z) \geq u_D)$  in  $u_D$  for any component of vector  $Z$ , which is a property of normal selection models.

- Panel A of figure 11 plots the treatment weights derived by ?? and the IV weight (24), along with the MTE.

- Panel A of figure 11 plots the treatment weights derived by ?? and the IV weight (24), along with the MTE.
- The  $ATE = \Delta^{ATE}$  weight is flat (= 1).

- Panel A of figure 11 plots the treatment weights derived by ?? and the IV weight (24), along with the MTE.
- The  $ATE = \Delta^{ATE}$  weight is flat (= 1).
- TT oversamples the low  $u_D$  agents (those more likely to adopt the policies).

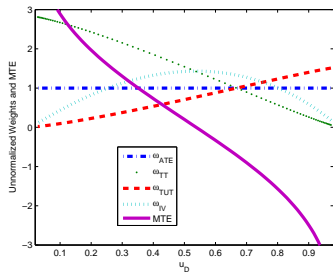


- Panel A of figure 11 plots the treatment weights derived by ?? and the IV weight (24), along with the MTE.
- The  $ATE = \Delta^{ATE}$  weight is flat (= 1).
- TT oversamples the low  $u_D$  agents (those more likely to adopt the policies).
- TUT oversamples the high  $u_D$  agents.

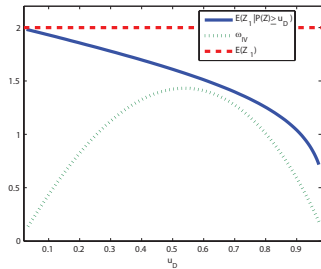
- Panel A of figure 11 plots the treatment weights derived by ?? and the IV weight (24), along with the MTE.
- The  $ATE = \Delta^{ATE}$  weight is flat (= 1).
- TT oversamples the low  $u_D$  agents (those more likely to adopt the policies).
- TUT oversamples the high  $u_D$  agents.
- The IV weight is positive as it must be when the  $Z$  are normally distributed.

Figure 11: Treatment Weights, IV Weights using  $Z_1$  as the Instrument and the Marginal Treatment Effect

A. Weights and MTE



B. IV Weights,  $E(Z_1|P(Z) \geq u_D)$  and  $E(Z_1)$



Parameter	Under Assumptions(*)
ATE	0.2
TT	1.1878
TUT	-0.9132
$IV_{Z_1}$	0.0924

- IV is far from any of the standard treatment parameters.

- IV is far from any of the standard treatment parameters.
- Panel B decomposes the weight into its numerator components  $E(Z_1 | P(Z) \geq u_D)$  and  $E(Z_1)$ , and the weight itself.

- IV is far from any of the standard treatment parameters.
- Panel B decomposes the weight into its numerator components  $E(Z_1 | P(Z) \geq u_D)$  and  $E(Z_1)$ , and the weight itself.
- The difference  $E(Z_1 | P(Z) \geq u_D) - E(Z_1)$  multiplied by  $\Pr(P(Z) \geq u_D)$  and normalized by  $\text{Cov}(Z_1, D)$  is the weight (see equation 23).

- IV is far from any of the standard treatment parameters.
- Panel B decomposes the weight into its numerator components  $E(Z_1 | P(Z) \geq u_D)$  and  $E(Z_1)$ , and the weight itself.
- The difference  $E(Z_1 | P(Z) \geq u_D) - E(Z_1)$  multiplied by  $\Pr(P(Z) \geq u_D)$  and normalized by  $\text{Cov}(Z_1, D)$  is the weight (see equation 23).
- The weight is plotted as the dotted line in figure 9B.



- Suppose that instead of assuming normality for the regressors, instrument  $Z$  is assumed to be a random vector with a distribution function given by a mixture of two normals:

$$Z \sim P_1 N(\kappa_1, \Sigma_1) + P_2 N(\kappa_2, \Sigma_2),$$

where  $P_1$  is the proportion in population 1,  $P_2$  is the proportion in population 2 and  $P_1 + P_2 = 1$ .

- Suppose that instead of assuming normality for the regressors, instrument  $Z$  is assumed to be a random vector with a distribution function given by a mixture of two normals:

$$Z \sim P_1 N(\kappa_1, \Sigma_1) + P_2 N(\kappa_2, \Sigma_2),$$

where  $P_1$  is the proportion in population 1,  $P_2$  is the proportion in population 2 and  $P_1 + P_2 = 1$ .

- This produces a model with continuous instruments, where  $E(\tilde{J}(Z) \mid P(Z) \geq u_D)$  need not be monotonic in  $u_D$  where  $\tilde{J}(Z) = J(Z) - E(J(Z))$ .

- Suppose that instead of assuming normality for the regressors, instrument  $Z$  is assumed to be a random vector with a distribution function given by a mixture of two normals:

$$Z \sim P_1 N(\kappa_1, \Sigma_1) + P_2 N(\kappa_2, \Sigma_2),$$

where  $P_1$  is the proportion in population 1,  $P_2$  is the proportion in population 2 and  $P_1 + P_2 = 1$ .

- This produces a model with continuous instruments, where  $E(\tilde{J}(Z) | P(Z) \geq u_D)$  need not be monotonic in  $u_D$  where  $\tilde{J}(Z) = J(Z) - E(J(Z))$ .
- Such a data generating process for the instrument could arise from an ecological model in which two different populations are mixed (e.g., rural and urban populations).

- Appendix, Slide 1107, derives the instrumental variable weights on  $\Delta^{\text{MTE}}$  when  $Z_1$  (the first element of  $Z$ ) is used as the instrument, i.e.,  $J(Z) = Z_1$ .

- Appendix, Slide 1107, derives the instrumental variable weights on  $\Delta^{\text{MTE}}$  when  $Z_1$  (the first element of  $Z$ ) is used as the instrument, i.e.,  $J(Z) = Z_1$ .
- For simplicity, we assume that there are no  $X$  regressors.

- Appendix, Slide 1107, derives the instrumental variable weights on  $\Delta^{\text{MTE}}$  when  $Z_1$  (the first element of  $Z$ ) is used as the instrument, i.e.,  $J(Z) = Z_1$ .
- For simplicity, we assume that there are no  $X$  regressors.
- The probability of selection is generated using  $\mu_D(Z) = Z\gamma$ .

- Appendix, Slide 1107, derives the instrumental variable weights on  $\Delta^{\text{MTE}}$  when  $Z_1$  (the first element of  $Z$ ) is used as the instrument, i.e.,  $J(Z) = Z_1$ .
- For simplicity, we assume that there are no  $X$  regressors.
- The probability of selection is generated using  $\mu_D(Z) = Z\gamma$ .
- The joint distribution of  $(Z_1, Z\gamma)$  is normal within each group.

- In our example, the dependence between  $Z_1$  and  $Z_\gamma$  ( $= F_V(Z_\gamma) = P(Z)$ ) is negative in one population and positive in another.



- In our example, the dependence between  $Z_1$  and  $Z_\gamma$  ( $= F_V(Z_\gamma) = P(Z)$ ) is negative in one population and positive in another.
- Thus in one population, as  $Z_1$  increases  $P(Z)$  increases.

- In our example, the dependence between  $Z_1$  and  $Z_\gamma$  ( $= F_V(Z_\gamma) = P(Z)$ ) is negative in one population and positive in another.
- Thus in one population, as  $Z_1$  increases  $P(Z)$  increases.
- In the other population, as  $Z_1$  increases  $P(Z)$  decreases.

- In our example, the dependence between  $Z_1$  and  $Z_\gamma$  ( $= F_V(Z_\gamma) = P(Z)$ ) is negative in one population and positive in another.
- Thus in one population, as  $Z_1$  increases  $P(Z)$  increases.
- In the other population, as  $Z_1$  increases  $P(Z)$  decreases.
- If this second population is sufficiently big ( $P_1$  is small) or the negative correlation in the second population is sufficiently big, the weights can become negative because  $E(\tilde{J}(Z) | P(Z) \geq u_D)$  is not monotonic in  $u_D$ .

- We present examples for a conventional normal outcome selection model generated by the parameters presented at the base of figure 12.

- We present examples for a conventional normal outcome selection model generated by the parameters presented at the base of figure 12.
- The discrete choice equation is a conventional probit:  
$$\Pr(D = 1 \mid Z = z) = \Phi\left(\frac{z\gamma}{\sigma_V}\right).$$

- We present examples for a conventional normal outcome selection model generated by the parameters presented at the base of figure 12.
- The discrete choice equation is a conventional probit:  
$$\Pr(D = 1 \mid Z = z) = \Phi\left(\frac{z\gamma}{\sigma_V}\right).$$
- The outcome equations are linear normal equations.

- We present examples for a conventional normal outcome selection model generated by the parameters presented at the base of figure 12.
- The discrete choice equation is a conventional probit:  
$$\Pr(D = 1 \mid Z = z) = \Phi\left(\frac{z\gamma}{\sigma_V}\right).$$
- The outcome equations are linear normal equations.
- Thus  $\Delta^{\text{MTE}}(v) = E(Y_1 - Y_0 \mid V = v)$ , is linear in  $v$ :

$$E(Y_1 - Y_0 \mid V = v) = \mu_1 - \mu_0 + \frac{\text{Cov}(U_1 - U_0, V)}{\text{Var}(V)}v.$$

- We present examples for a conventional normal outcome selection model generated by the parameters presented at the base of figure 12.
- The discrete choice equation is a conventional probit:  
$$\Pr(D = 1 \mid Z = z) = \Phi\left(\frac{z\gamma}{\sigma_V}\right).$$
- The outcome equations are linear normal equations.
- Thus  $\Delta^{\text{MTE}}(v) = E(Y_1 - Y_0 \mid V = v)$ , is linear in  $v$ :

$$E(Y_1 - Y_0 \mid V = v) = \mu_1 - \mu_0 + \frac{\text{Cov}(U_1 - U_0, V)}{\text{Var}(V)}v.$$

- At the base of the figure, we define  $\bar{\beta} = \mu_1 - \mu_0$  and  $\alpha = \mu_0$ .



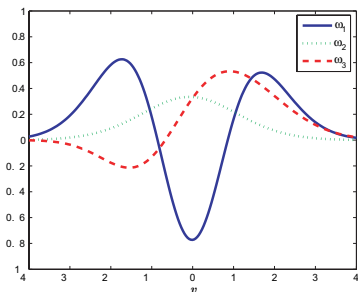
- We present examples for a conventional normal outcome selection model generated by the parameters presented at the base of figure 12.
- The discrete choice equation is a conventional probit:  
$$\Pr(D = 1 \mid Z = z) = \Phi\left(\frac{z\gamma}{\sigma_V}\right).$$
- The outcome equations are linear normal equations.
- Thus  $\Delta^{\text{MTE}}(v) = E(Y_1 - Y_0 \mid V = v)$ , is linear in  $v$ :

$$E(Y_1 - Y_0 \mid V = v) = \mu_1 - \mu_0 + \frac{\text{Cov}(U_1 - U_0, V)}{\text{Var}(V)}v.$$

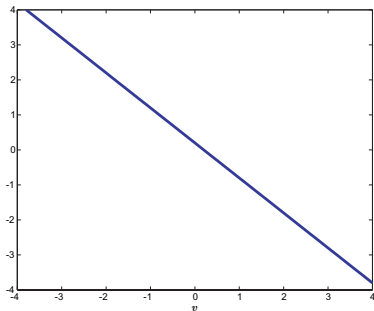
- At the base of the figure, we define  $\bar{\beta} = \mu_1 - \mu_0$  and  $\alpha = \mu_0$ .
- The average treatment effects are the same for all different distributions of the  $Z$ .

Figure 12: Marginal Treatment Effect and IV Weights using  $Z_1$  as the Instrument when  $Z = (Z_1, Z_2) \sim p_1 N(\kappa_1, \Sigma_1) + p_2 N(\kappa_2, \Sigma_2)$  for different values of  $\Sigma_2$

A. IV Weights



B.  $\Delta^{\text{MTE}}(v)$



Outcomes

$$Y_1 = \alpha + \bar{\beta} + U_1$$

$$Y_0 = \alpha + U_0$$

Choice Model

$$D = \begin{cases} 1 & \text{if } D^* \geq 0 \\ 0 & \text{if } D^* < 0 \end{cases}$$
$$D^* = Y_1 - Y_0 - \gamma Z \text{ and } V = -(U_1 - U_0)$$

Parameterization

$$(U_1, U_0) \sim N(\mathbf{0}, \Sigma), \quad \Sigma = \begin{bmatrix} 1 & -0.9 \\ -0.9 & 1 \end{bmatrix}, \quad \alpha = 0.67, \bar{\beta} = 0.2$$

$$Z = (Z_1, Z_2) \sim p_1 N(\kappa_1, \Sigma_1) + p_2 N(\kappa_2, \Sigma_2)$$

$$p_1 = 0.45, p_2 = 0.55 \quad ; \quad \Sigma_1 = \begin{bmatrix} 1.4 & 0.5 \\ 0.5 & 1.4 \end{bmatrix}$$

$$\text{Cov}(Z_1, \gamma Z) = \gamma \Sigma_1^1 = 0.98 \quad ; \quad \gamma = (0.2, 1.4)$$

Source: ?

- In each of the following examples, we show results for models with vector  $Z$  that satisfies (IV-1) and (IV-2) and with  $\gamma > 0$  componentwise where  $\gamma$  is the coefficient of  $Z$  in the cost equation.

- In each of the following examples, we show results for models with vector  $Z$  that satisfies (IV-1) and (IV-2) and with  $\gamma > 0$  componentwise where  $\gamma$  is the coefficient of  $Z$  in the cost equation.
- We vary the weights and means of the instruments.

- In each of the following examples, we show results for models with vector  $Z$  that satisfies (IV-1) and (IV-2) and with  $\gamma > 0$  componentwise where  $\gamma$  is the coefficient of  $Z$  in the cost equation.
- We vary the weights and means of the instruments.
- *Ceteris paribus*, an increase in each component of  $Z$  increases  $\Pr(D = 1 \mid Z = z)$ .

- In each of the following examples, we show results for models with vector  $Z$  that satisfies (IV-1) and (IV-2) and with  $\gamma > 0$  componentwise where  $\gamma$  is the coefficient of  $Z$  in the cost equation.
- We vary the weights and means of the instruments.
- *Ceteris paribus*, an increase in each component of  $Z$  increases  $\Pr(D = 1 \mid Z = z)$ .
- Table 7 presents the parameters treatment on the treated ( $E(Y_1 - Y_0 \mid D = 1)$ ), treatment on the untreated ( $E(Y_1 - Y_0 \mid D = 0)$ ), and the average treatment effect ( $E(Y_1 - Y_0)$ ) produced by our model for different distributions of the regressors.

Table 7: IV estimator and  $\text{Cov}(Z_2, \gamma Z)$  associated with each value of  $\Sigma_2$

Weights	$\Sigma_2$	$\kappa_1$	$\kappa_2$	IV	ATE	TT	TUT	$\text{Cov}(Z_2, \gamma Z) = \gamma \Sigma_2^1$
$\omega_1$	$\begin{bmatrix} 0.6 & -0.5 \\ -0.5 & 0.6 \end{bmatrix}$	$[ 0 \ 0 ]$	$[ 0 \ 0 ]$	0.434	0.2	1.401	-1.175	-0.58
$\omega_2$	$\begin{bmatrix} 0.6 & 0.1 \\ 0.1 & 0.6 \end{bmatrix}$	$[ 0 \ 0 ]$	$[ 0 \ 0 ]$	0.078	0.2	1.378	-1.145	0.26
$\omega_3$	$\begin{bmatrix} 0.6 & -0.3 \\ -0.3 & 0.6 \end{bmatrix}$	$[ 0 \ -1 ]$	$[ 0 \ 1 ]$	-2.261	0.2	1.310	-0.859	-0.30

Source: Heckman, Urzua and Vytlacil (2006)



- In standard IV analysis, under assumptions (IV-1) and (IV-2) the distribution of  $Z$  does not affect the probability limit of the IV estimator.

- In standard IV analysis, under assumptions (IV-1) and (IV-2) the distribution of  $Z$  does not affect the probability limit of the IV estimator.
- It only affects its sampling distribution.

- In standard IV analysis, under assumptions (IV-1) and (IV-2) the distribution of  $Z$  does not affect the probability limit of the IV estimator.
- It only affects its sampling distribution.
- Figure 12A shows three weights corresponding to the perturbations of the variances of the instruments in the second component population  $\Sigma_2$  and the means  $(\kappa_1, \kappa_2)$  shown at the table at the base of the figure.

- In standard IV analysis, under assumptions (IV-1) and (IV-2) the distribution of  $Z$  does not affect the probability limit of the IV estimator.
- It only affects its sampling distribution.
- Figure 12A shows three weights corresponding to the perturbations of the variances of the instruments in the second component population  $\Sigma_2$  and the means  $(\kappa_1, \kappa_2)$  shown at the table at the base of the figure.
- The  $\Delta_V^{\text{MTE}}$  used in all of our examples are plotted in figure 12B.

- In standard IV analysis, under assumptions (IV-1) and (IV-2) the distribution of  $Z$  does not affect the probability limit of the IV estimator.
- It only affects its sampling distribution.
- Figure 12A shows three weights corresponding to the perturbations of the variances of the instruments in the second component population  $\Sigma_2$  and the means  $(\kappa_1, \kappa_2)$  shown at the table at the base of the figure.
- The  $\Delta_V^{\text{MTE}}$  used in all of our examples are plotted in figure 12B.
- The MTE has the familiar shape, reported in ? and ? that returns are highest for those with values of  $v$  that make them more likely to get treatment (i.e., low values of  $v$ ).

- The weights  $\omega_1$  and  $\omega_3$  plotted in figure 12A correspond to the case where  $E(Z_1 - E(Z_1) | P(Z) \geq u_D)$  is not monotonic in  $u_D$ .

- The weights  $\omega_1$  and  $\omega_3$  plotted in figure 12A correspond to the case where  $E(Z_1 - E(Z_1) | P(Z) \geq u_D)$  is not monotonic in  $u_D$ .
- In these cases, the sign of the covariance between  $Z_1$  and  $Z_\gamma$  (i.e.,  $P(Z)$ ) is not the same in the two subpopulations.

- The weights  $\omega_1$  and  $\omega_3$  plotted in figure 12A correspond to the case where  $E(Z_1 - E(Z_1) | P(Z) \geq u_D)$  is not monotonic in  $u_D$ .
- In these cases, the sign of the covariance between  $Z_1$  and  $Z_\gamma$  (i.e.,  $P(Z)$ ) is not the same in the two subpopulations.
- The IV estimates reported in the table at the base of the figure range all over the place even though the parameters of the outcome and choice model are the same.



- Different distributions of  $Z$  critically affect the probability limit of the IV estimator in the model of essential heterogeneity.

- Different distributions of  $Z$  critically affect the probability limit of the IV estimator in the model of essential heterogeneity.
- The model of outcomes and choices is the same across all of these examples.

- Different distributions of  $Z$  critically affect the probability limit of the IV estimator in the model of essential heterogeneity.
- The model of outcomes and choices is the same across all of these examples.
- The MTE and ATE parameters are the same.

- Different distributions of  $Z$  critically affect the probability limit of the IV estimator in the model of essential heterogeneity.
- The model of outcomes and choices is the same across all of these examples.
- The MTE and ATE parameters are the same.
- Only the distribution of the instrument differs.

- Different distributions of  $Z$  critically affect the probability limit of the IV estimator in the model of essential heterogeneity.
- The model of outcomes and choices is the same across all of these examples.
- The MTE and ATE parameters are the same.
- Only the distribution of the instrument differs.
- The instrumental variable estimand is sometimes positive and sometimes negative, and oscillates wildly in magnitude depending on the distribution of the instruments.

- The estimated “effect” is often way off the mark for any desired treatment parameter.

- The estimated “effect” is often way off the mark for any desired treatment parameter.
- These examples show how uniformity in  $Z$  does not translate into uniformity in  $J(Z)$  ( $Z_1$  in this example).

- The estimated “effect” is often way off the mark for any desired treatment parameter.
- These examples show how uniformity in  $Z$  does not translate into uniformity in  $J(Z)$  ( $Z_1$  in this example).
- This sensitivity is a phenomenon that does not appear in the conventional homogeneous response model but is a central feature of a model with essential heterogeneity.



- The estimated “effect” is often way off the mark for any desired treatment parameter.
- These examples show how uniformity in  $Z$  does not translate into uniformity in  $J(Z)$  ( $Z_1$  in this example).
- This sensitivity is a phenomenon that does not appear in the conventional homogeneous response model but is a central feature of a model with essential heterogeneity.
- We now compare selection and IV models.

## Comparing Selection and IV Models

- We now show that local IV identifies the derivatives of a selection model.

## Comparing Selection and IV Models

- We now show that local IV identifies the derivatives of a selection model.
- Making the  $X$  explicit, in the standard selection model,  $U_1$  and  $U_0$  are scalar random variables that are additively separable in the outcome equations,  $Y_1 = \mu_1(X) + U_1$  and  $Y_0 = \mu_0(X) + U_0$ .

## Comparing Selection and IV Models

- We now show that local IV identifies the derivatives of a selection model.
- Making the  $X$  explicit, in the standard selection model,  $U_1$  and  $U_0$  are scalar random variables that are additively separable in the outcome equations,  $Y_1 = \mu_1(X) + U_1$  and  $Y_0 = \mu_0(X) + U_0$ .
- The control function approach conditions on  $Z$  and  $D$ .

- As a consequence of index sufficiency, this is equivalent to conditioning on  $P(Z)$  and  $D$ :

$$E(Y | X, D, Z) = \mu_0(X) + [\mu_1(X) - \mu_0(X)] D \\ + K_1(P(Z), X) D + K_0(P(Z), X) (1 - D),$$

where the control functions are

$$K_1(P(Z), X) = E(U_1 | D = 1, X, P(Z)) \\ K_0(P(Z), X) = E(U_0 | D = 0, X, P(Z)).$$

- The IV approach does not condition on  $D$ .

- The IV approach does not condition on  $D$ .
- It works with

$$\begin{aligned} E(Y | X, Z) &= \mu_0(X) + [\mu_1(X) - \mu_0(X)] P(Z) & (27) \\ &+ K_1(P(Z), X) P(Z) \\ &+ K_0(P(Z), X) (1 - P(Z)), \end{aligned}$$

the population mean outcome given  $X, Z$ .

- From index sufficiency,  $E(Y | X, Z) = E(Y | X, P(Z))$ .



- From index sufficiency,  $E(Y | X, Z) = E(Y | X, P(Z))$ .
- The MTE is the derivative of this expression with respect to  $P(Z)$ , which we have defined as LIV:

$$\frac{\partial E(Y | X, P(Z))}{\partial P(Z)} \Big|_{P(Z)=p} = \text{LIV}(X, p) = \text{MTE}(X, p).$$

- From index sufficiency,  $E(Y | X, Z) = E(Y | X, P(Z))$ .
- The MTE is the derivative of this expression with respect to  $P(Z)$ , which we have defined as LIV:

$$\left. \frac{\partial E(Y | X, P(Z))}{\partial P(Z)} \right|_{P(Z)=p} = \text{LIV}(X, p) = \text{MTE}(X, p).$$

- The distribution of  $P(Z)$  and the relationship between  $J(Z)$  and  $P(Z)$  determine the weight on MTE.

- From index sufficiency,  $E(Y | X, Z) = E(Y | X, P(Z))$ .
- The MTE is the derivative of this expression with respect to  $P(Z)$ , which we have defined as LIV:

$$\left. \frac{\partial E(Y | X, P(Z))}{\partial P(Z)} \right|_{P(Z)=p} = \text{LIV}(X, p) = \text{MTE}(X, p).$$

- The distribution of  $P(Z)$  and the relationship between  $J(Z)$  and  $P(Z)$  determine the weight on MTE.
- Under assumptions (A-1)–(A-5), along with rank and limit conditions (??), one can identify  $\mu_1(X)$ ,  $\mu_0(X)$ ,  $K_1(P(Z), X)$ , and  $K_0(P(Z), X)$ .

- The selection (control function) estimator identifies the conditional means

$$E(Y_1 | X, P(Z), D = 1) = \mu_1(X) + K_1(X, P(Z)) \quad (28a)$$

and

$$E(Y_0 | X, P(Z), D = 0) = \mu_0(X) + K_0(X, P(Z)). \quad (28b)$$

These can be identified from nonparametric regressions of  $Y_1$  and  $Y_0$  on  $X, Z$  in each population.

- The selection (control function) estimator identifies the conditional means

$$E(Y_1 | X, P(Z), D = 1) = \mu_1(X) + K_1(X, P(Z)) \quad (28a)$$

and

$$E(Y_0 | X, P(Z), D = 0) = \mu_0(X) + K_0(X, P(Z)). \quad (28b)$$

These can be identified from nonparametric regressions of  $Y_1$  and  $Y_0$  on  $X, Z$  in each population.

- To decompose these means and separate  $\mu_1(X)$  from  $K_1(X, P(Z))$  without invoking functional form or curvature assumptions, it is necessary to have an exclusion (a  $Z$  not in  $X$ ).

- In addition, there must exist a limit set for  $Z$  given  $X$  such that  $K_1(X, P(Z)) = 0$  for  $Z$  in that limit set.

- In addition, there must exist a limit set for  $Z$  given  $X$  such that  $K_1(X, P(Z)) = 0$  for  $Z$  in that limit set.
- Otherwise, without functional form or curvature assumptions, it is not possible to disentangle  $\mu_1(X)$  from  $K_1(X, P(Z))$  which may contain constants and functions of  $X$  that do not interact with  $P(Z)$  (see ?).

- In addition, there must exist a limit set for  $Z$  given  $X$  such that  $K_1(X, P(Z)) = 0$  for  $Z$  in that limit set.
- Otherwise, without functional form or curvature assumptions, it is not possible to disentangle  $\mu_1(X)$  from  $K_1(X, P(Z))$  which may contain constants and functions of  $X$  that do not interact with  $P(Z)$  (see ?).
- A parallel argument for  $Y_0$  shows that we require a limit set for  $Z$  given  $X$  such that  $K_0(X, P(Z)) = 0$ .



- In addition, there must exist a limit set for  $Z$  given  $X$  such that  $K_1(X, P(Z)) = 0$  for  $Z$  in that limit set.
- Otherwise, without functional form or curvature assumptions, it is not possible to disentangle  $\mu_1(X)$  from  $K_1(X, P(Z))$  which may contain constants and functions of  $X$  that do not interact with  $P(Z)$  (see ?).
- A parallel argument for  $Y_0$  shows that we require a limit set for  $Z$  given  $X$  such that  $K_0(X, P(Z)) = 0$ .
- Selection models operate by identifying the components of (28a) and (28b) and generating the treatment parameters from these components.

- In addition, there must exist a limit set for  $Z$  given  $X$  such that  $K_1(X, P(Z)) = 0$  for  $Z$  in that limit set.
- Otherwise, without functional form or curvature assumptions, it is not possible to disentangle  $\mu_1(X)$  from  $K_1(X, P(Z))$  which may contain constants and functions of  $X$  that do not interact with  $P(Z)$  (see ?).
- A parallel argument for  $Y_0$  shows that we require a limit set for  $Z$  given  $X$  such that  $K_0(X, P(Z)) = 0$ .
- Selection models operate by identifying the components of (28a) and (28b) and generating the treatment parameters from these components.
- Thus they work with levels of the  $Y$ .

- The local IV method works with derivatives of (27) and not levels and cannot directly recover the constant terms in (28a) and (28b).

- The local IV method works with derivatives of (27) and not levels and cannot directly recover the constant terms in (28a) and (28b).
- Using our analysis of LIV but applied to  $YD = Y_1D$  and  $Y(1 - D) = Y_0(1 - D)$ , it is straightforward to use LIV to estimate the components of the MTE separately.

- The local IV method works with derivatives of (27) and not levels and cannot directly recover the constant terms in (28a) and (28b).
- Using our analysis of LIV but applied to  $YD = Y_1D$  and  $Y(1 - D) = Y_0(1 - D)$ , it is straightforward to use LIV to estimate the components of the MTE separately.
- Thus we can identify

$$\mu_1(X) + E(U_1 | X, U_D = u_D)$$

and

$$\mu_0(X) + E(U_0 | X, U_D = u_D)$$

separately.

- This corresponds to what is estimated from taking the derivatives of expressions (28a) and (28b) multiplied by  $P(Z)$  and  $(1 - P(Z))$  respectively:

$$\begin{aligned} & P(Z)E(Y_1 | X, Z, D = 1) \\ = & P(Z)\mu_1(X) + P(Z)K_1(X, P(Z)) \end{aligned}$$

and

$$\begin{aligned} & (1 - P(Z))E(Y_0 | X, Z, D = 0) \\ = & (1 - P(Z))\mu_0(X) + (1 - P(Z))K_0(X, P(Z)). \end{aligned}$$

Thus the control function method works with levels, whereas the LIV approach works with slopes of combinations of the same basic functions.

- Constants that do not depend on  $P(Z)$  disappear from the estimates of the model.

- Constants that do not depend on  $P(Z)$  disappear from the estimates of the model.
- The level parameters are obtained by integration using the formulae in table 2B.



- Constants that do not depend on  $P(Z)$  disappear from the estimates of the model.
- The level parameters are obtained by integration using the formulae in table 2B.
- Misspecification of  $P(Z)$  (either its functional form or its arguments) and hence of  $K_1(P(Z), X)$  and  $K_0(P(Z), X)$ , in general, produces biased estimates of the parameters of the model under the control function approach even if semiparametric methods are used to estimate  $\mu_0, \mu_1, K_0$  and  $K_1$ .

- To implement the method, we need to know all of the arguments of  $Z$ .

- To implement the method, we need to know all of the arguments of  $Z$ .
- The terms  $K_1(P(Z), X)$  and  $K_0(P(Z), X)$  can be nonparametrically estimated so it is only necessary to know  $P(Z)$  up to a monotonic transformation.

- To implement the method, we need to know all of the arguments of  $Z$ .
- The terms  $K_1(P(Z), X)$  and  $K_0(P(Z), X)$  can be nonparametrically estimated so it is only necessary to know  $P(Z)$  up to a monotonic transformation.
- The distributions of  $U_0$ ,  $U_1$  and  $V$  do not need to be specified to estimate control function models (see ?).

- These problems with control function models have their counterparts in IV models.

- These problems with control function models have their counterparts in IV models.
- If we use a misspecified  $P(Z)$  to identify the MTE or its components, in general, we do not identify MTE or its components.

- These problems with control function models have their counterparts in IV models.
- If we use a misspecified  $P(Z)$  to identify the MTE or its components, in general, we do not identify MTE or its components.
- Misspecification of  $P(Z)$  plagues both approaches.

- One common criticism of selection models is that without invoking functional form assumptions, identification of  $\mu_1(X)$  and  $\mu_0(X)$  requires that  $P(Z) \rightarrow 1$  and  $P(Z) \rightarrow 0$  in limit sets.



- One common criticism of selection models is that without invoking functional form assumptions, identification of  $\mu_1(X)$  and  $\mu_0(X)$  requires that  $P(Z) \rightarrow 1$  and  $P(Z) \rightarrow 0$  in limit sets.
- Identification in limit sets is sometimes called “identification at infinity.” In order to identify  $ATE = E(Y_1 - Y_0|X)$ , IV methods also require that  $P(Z) \rightarrow 1$  and  $P(Z) \rightarrow 0$  in limit sets, so an identification at infinity argument is implicit when IV is used to identify this parameter.

- One common criticism of selection models is that without invoking functional form assumptions, identification of  $\mu_1(X)$  and  $\mu_0(X)$  requires that  $P(Z) \rightarrow 1$  and  $P(Z) \rightarrow 0$  in limit sets.
- Identification in limit sets is sometimes called “identification at infinity.” In order to identify  $ATE = E(Y_1 - Y_0|X)$ , IV methods also require that  $P(Z) \rightarrow 1$  and  $P(Z) \rightarrow 0$  in limit sets, so an identification at infinity argument is implicit when IV is used to identify this parameter.
- The LATE parameter avoids this problem by moving the goal posts and redefining the parameter of interest away from a level parameter like ATE or TT to a slope parameter like LATE which differences out the unidentified constants.

- Alternatively, if we define the parameter of interest to be LATE or MTE, we can use the selection model without invoking identification at infinity.

- Alternatively, if we define the parameter of interest to be LATE or MTE, we can use the selection model without invoking identification at infinity.
- The IV estimator is model dependent, just like the selection estimator, but in application, the model does not have to be fully specified to obtain  $\Delta^{\text{IV}}$  using  $Z$  (or  $J(Z)$ ).

- Alternatively, if we define the parameter of interest to be LATE or MTE, we can use the selection model without invoking identification at infinity.
- The IV estimator is model dependent, just like the selection estimator, but in application, the model does not have to be fully specified to obtain  $\Delta^{\text{IV}}$  using  $Z$  (or  $J(Z)$ ).
- However, the distribution of  $P(Z)$  and the relationship between  $P(Z)$  and  $J(Z)$  generates the weights.

- Alternatively, if we define the parameter of interest to be LATE or MTE, we can use the selection model without invoking identification at infinity.
- The IV estimator is model dependent, just like the selection estimator, but in application, the model does not have to be fully specified to obtain  $\Delta^{\text{IV}}$  using  $Z$  (or  $J(Z)$ ).
- However, the distribution of  $P(Z)$  and the relationship between  $P(Z)$  and  $J(Z)$  generates the weights.
- The interpretation placed on  $\Delta^{\text{IV}}$  in terms of weights on  $\Delta^{\text{MTE}}$  depends crucially on the specification of  $P(Z)$ .

- Alternatively, if we define the parameter of interest to be LATE or MTE, we can use the selection model without invoking identification at infinity.
- The IV estimator is model dependent, just like the selection estimator, but in application, the model does not have to be fully specified to obtain  $\Delta^{\text{IV}}$  using  $Z$  (or  $J(Z)$ ).
- However, the distribution of  $P(Z)$  and the relationship between  $P(Z)$  and  $J(Z)$  generates the weights.
- The interpretation placed on  $\Delta^{\text{IV}}$  in terms of weights on  $\Delta^{\text{MTE}}$  depends crucially on the specification of  $P(Z)$ .
- In both control function and IV approaches for the general model of heterogeneous responses,  $P(Z)$  plays a central role.

- Two economists using the same instrument will obtain the same point estimate using the same data.



- Two economists using the same instrument will obtain the same point estimate using the same data.
- Their *interpretation* of that estimate will differ depending on how they specify the arguments in  $P(Z)$ , even if neither uses  $P(Z)$  as an instrument.

- Two economists using the same instrument will obtain the same point estimate using the same data.
- Their *interpretation* of that estimate will differ depending on how they specify the arguments in  $P(Z)$ , even if neither uses  $P(Z)$  as an instrument.
- By conditioning on  $P(Z)$ , the control function approach makes the dependence of estimates on the specification of  $P(Z)$  explicit.

- Two economists using the same instrument will obtain the same point estimate using the same data.
- Their *interpretation* of that estimate will differ depending on how they specify the arguments in  $P(Z)$ , even if neither uses  $P(Z)$  as an instrument.
- By conditioning on  $P(Z)$ , the control function approach makes the dependence of estimates on the specification of  $P(Z)$  explicit.
- The IV approach is less explicit and masks the assumptions required to economically interpret the empirical output of an IV estimation.

- Two economists using the same instrument will obtain the same point estimate using the same data.
- Their *interpretation* of that estimate will differ depending on how they specify the arguments in  $P(Z)$ , even if neither uses  $P(Z)$  as an instrument.
- By conditioning on  $P(Z)$ , the control function approach makes the dependence of estimates on the specification of  $P(Z)$  explicit.
- The IV approach is less explicit and masks the assumptions required to economically interpret the empirical output of an IV estimation.
- We now turn to some empirical examples of LIV.

## Empirical Examples: “The effect” of high school graduation on wages and using IV to estimate “the effect” of the GED

- The previous examples illustrate logical possibilities.

## Empirical Examples: “The effect” of high school graduation on wages and using IV to estimate “the effect” of the GED

- The previous examples illustrate logical possibilities.
- This subsection shows that these logical possibilities arise in real data.

## Empirical Examples: “The effect” of high school graduation on wages and using IV to estimate “the effect” of the GED

- The previous examples illustrate logical possibilities.
- This subsection shows that these logical possibilities arise in real data.
- We analyze two examples: (a) the effect of graduating high school on wages, and (b) the effect of obtaining a GED on wages.

## Empirical Examples: “The effect” of high school graduation on wages and using IV to estimate “the effect” of the GED

- The previous examples illustrate logical possibilities.
- This subsection shows that these logical possibilities arise in real data.
- We analyze two examples: (a) the effect of graduating high school on wages, and (b) the effect of obtaining a GED on wages.
- We first analyze the effect of graduating high school on wages.



## Empirical Example Based on LATE: Using IV to Estimate “*The Effect*” of High School Graduation on Wages

- We first study the effects of graduating from high school on wages using data from the National Longitudinal Survey of Youth 1979 (NLSY79).

## Empirical Example Based on LATE: Using IV to Estimate “*The Effect*” of High School Graduation on Wages

- We first study the effects of graduating from high school on wages using data from the National Longitudinal Survey of Youth 1979 (NLSY79).
- This survey gathers information at multiple points in time on the labor market activities for men and women born in the years 1957–1964.

## Empirical Example Based on LATE: Using IV to Estimate “*The Effect*” of High School Graduation on Wages

- We first study the effects of graduating from high school on wages using data from the National Longitudinal Survey of Youth 1979 (NLSY79).
- This survey gathers information at multiple points in time on the labor market activities for men and women born in the years 1957–1964.
- We estimate LATE using log hourly wages at age 30 as the outcome measure.

## Empirical Example Based on LATE: Using IV to Estimate “*The Effect*” of High School Graduation on Wages

- We first study the effects of graduating from high school on wages using data from the National Longitudinal Survey of Youth 1979 (NLSY79).
- This survey gathers information at multiple points in time on the labor market activities for men and women born in the years 1957–1964.
- We estimate LATE using log hourly wages at age 30 as the outcome measure.
- Following a large body of research (see ?), we use the number of siblings and residence in the south at age 14 as instruments.

- Figure 13 plots the weights on LATE using the estimated  $P(Z)$ .

- Figure 13 plots the weights on LATE using the estimated  $P(Z)$ .
- The procedure used to derive the estimates is explained in ?.

- Figure 13 plots the weights on LATE using the estimated  $P(Z)$ .
- The procedure used to derive the estimates is explained in ?.
- The weights are derived from equation (26).

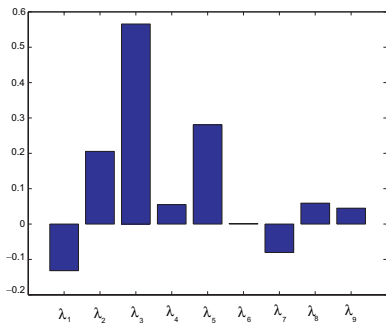
- Figure 13 plots the weights on LATE using the estimated  $P(Z)$ .
- The procedure used to derive the estimates is explained in ?.
- The weights are derived from equation (26).
- The LATE parameters are both positive and negative.



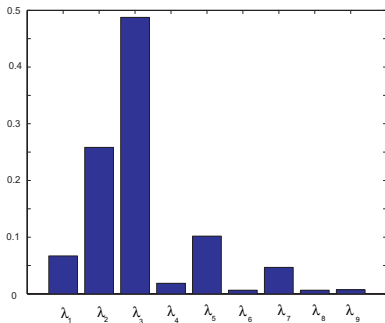
- Figure 13 plots the weights on LATE using the estimated  $P(Z)$ .
- The procedure used to derive the estimates is explained in ?.
- The weights are derived from equation (26).
- The LATE parameters are both positive and negative.
- The weights using siblings as an instrument are both positive and negative.

Figure 13: IV Weights - The Effect of Graduating from High School, Sample of High School Dropouts and High School Graduates, White Males - NLSY79

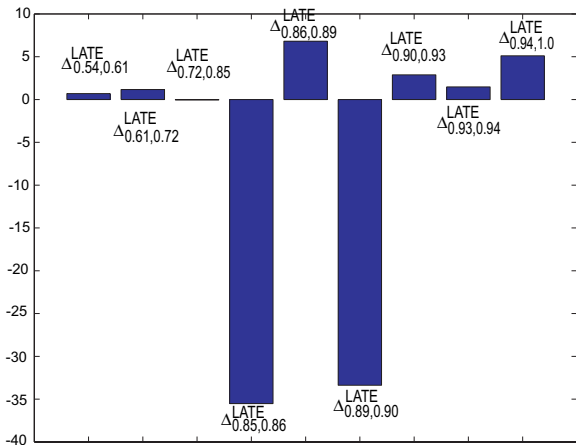
A. Weights: Number of Siblings as Instrument



B. Weights: Propensity Score as Instrument



### C. The Local Average Treatment Effects



$Y = \text{Log per-hour wage at age 30}$ ,  $Z_1 = \text{Number of Siblings in 1979}$ ,  $Z_2 = \text{Mother is a High School Graduate}$

$$D = \begin{cases} 1 & \text{if High School Graduate} \\ 0 & \text{if High School Dropout} \end{cases}$$

IV Estimates  
(bootstrap std. errors in parentheses - 100 replications)

Instrument	Value
Number of Siblings in 1979	0.115 (0.695)
Propensity Score	0.316 (0.110)

Joint Probability Distribution of  $(Z_1, Z_2)$  and the Propensity Score  
(joint probabilities  $\Pr(Z_1 = z_1, Z_2 = z_2)$  in ordinary type; propensity score  $\Pr(D = 1|Z_1 = z_1, Z_2 = z_2)$  in italics)

$Z_2 \backslash Z_1$	0	1	2	3	4
0	0.07 <i>1.0</i>	0.03 <i>0.54</i>	0.47 <i>0.86</i>	0.121 <i>0.72</i>	0.06 <i>0.61</i>
1	0.039 <i>0.94</i>	0.139 <i>0.89</i>	0.165 <i>0.90</i>	0.266 <i>0.85</i>	0.121 <i>0.93</i>

$\text{Cov}(Z_1, Z_2) = -0.066$  - Number of Observations = 1, 702

Source: Heckman, Urzua and Vytlacil (2006)

- The weights using  $P(Z)$  as an instrument are positive, as they must be following the analysis of ?.

- The weights using  $P(Z)$  as an instrument are positive, as they must be following the analysis of ?.
- The two IV estimates differ from each other because the weights are different.

- The weights using  $P(Z)$  as an instrument are positive, as they must be following the analysis of ?.
- The two IV estimates differ from each other because the weights are different.
- The overall IV estimate is a crude summary of the underlying component LATEs that are both large and positive and large and negative.

- The weights using  $P(Z)$  as an instrument are positive, as they must be following the analysis of ?.
- The two IV estimates differ from each other because the weights are different.
- The overall IV estimate is a crude summary of the underlying component LATEs that are both large and positive and large and negative.
- We next turn to analysis of the GED.



## Effect of the GED on Wages

- The GED test is used to certify high school dropouts as high school equivalents.

## Effect of the GED on Wages

- The GED test is used to certify high school dropouts as high school equivalents.
- Numerous studies document that the economic return to the GED is low (see ??).

## Effect of the GED on Wages

- The GED test is used to certify high school dropouts as high school equivalents.
- Numerous studies document that the economic return to the GED is low (see ??).
- It is estimated by the method described in ?.

## Effect of the GED on Wages

- The GED test is used to certify high school dropouts as high school equivalents.
- Numerous studies document that the economic return to the GED is low (see ??).
- It is estimated by the method described in ?.
- In this example, we study the effect of the GED on the wages of recipients compared to wages of dropouts.

- We use data from the National Longitudinal Survey of Youth 1979 (NLSY79) which gathers information at multiple points in time on the labor market activities for men and women born in the years 1957–1964.

- We use data from the National Longitudinal Survey of Youth 1979 (NLSY79) which gathers information at multiple points in time on the labor market activities for men and women born in the years 1957–1964.
- We estimate the MTE for the GED and also consider the IV weights for various instruments for a sample of males at age 25.

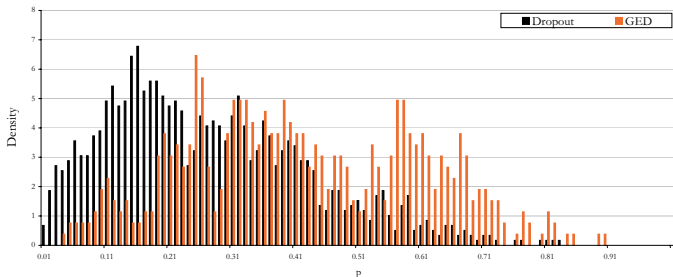
- We use data from the National Longitudinal Survey of Youth 1979 (NLSY79) which gathers information at multiple points in time on the labor market activities for men and women born in the years 1957–1964.
- We estimate the MTE for the GED and also consider the IV weights for various instruments for a sample of males at age 25.
- Figure 14 shows the sample support of  $P(Z)$  for both GEDs and high school dropouts.

- We use data from the National Longitudinal Survey of Youth 1979 (NLSY79) which gathers information at multiple points in time on the labor market activities for men and women born in the years 1957–1964.
- We estimate the MTE for the GED and also consider the IV weights for various instruments for a sample of males at age 25.
- Figure 14 shows the sample support of  $P(Z)$  for both GEDs and high school dropouts.
- It is not possible to estimate the MTE over its full support.



- We use data from the National Longitudinal Survey of Youth 1979 (NLSY79) which gathers information at multiple points in time on the labor market activities for men and women born in the years 1957–1964.
- We estimate the MTE for the GED and also consider the IV weights for various instruments for a sample of males at age 25.
- Figure 14 shows the sample support of  $P(Z)$  for both GEDs and high school dropouts.
- It is not possible to estimate the MTE over its full support.
- Thus the Average Treatment Effect (ATE) and Treatment on the Treated (TT) cannot be estimated from these data.

**Figure 14:** Frequency of the Propensity Score by Final Schooling Decision, Dropouts and GEDs—Males of the NLSY at Age 25



Note: The propensity score ( $P(D = 1|Z)$ ) is computed using as controls ( $Z$ ): Father's Highest Grade Completed, Mother's Highest Grade Completed, Number of Siblings, GED testing fee by state between 1993 and 2000, Family Income in 1979, Dropout's local wage at age 17, and High School Graduate's local unemployment at age 17. We also include two dummy variables controlling for the place of residence at age 14 (south and urban), and a set of dummies controlling for the year of birth (1957-1963).

- The list of  $Z$  variables is presented in Table 363 along with IV estimates.

- The list of  $Z$  variables is presented in Table 363 along with IV estimates.
- The IV estimates fluctuate from positive to negative.

- The list of  $Z$  variables is presented in Table 363 along with IV estimates.
- The IV estimates fluctuate from positive to negative.
- Using  $P(Z)$  as an instrument, the GED effect on log wages is in general negative.

- The list of  $Z$  variables is presented in Table 363 along with IV estimates.
- The IV estimates fluctuate from positive to negative.
- Using  $P(Z)$  as an instrument, the GED effect on log wages is in general negative.
- For other instruments, the signs and magnitudes vary.

- The list of  $Z$  variables is presented in Table 363 along with IV estimates.
- The IV estimates fluctuate from positive to negative.
- Using  $P(Z)$  as an instrument, the GED effect on log wages is in general negative.
- For other instruments, the signs and magnitudes vary.
- Figure 15 plots the estimated MTE.

- The list of  $Z$  variables is presented in Table 363 along with IV estimates.
- The IV estimates fluctuate from positive to negative.
- Using  $P(Z)$  as an instrument, the GED effect on log wages is in general negative.
- For other instruments, the signs and magnitudes vary.
- Figure 15 plots the estimated MTE.
- Details of the nonparametric estimation procedure used to produce these estimates are shown in an appendix in ?.



## Table 8: Instrumental Variables Estimates<sup>a</sup>

Sample of GED and Dropouts - Males at Age 25<sup>b</sup>

Instruments	IV-MTE
Father's Highest Grade Completed	0.146 (0.251)
Mother's Highest Grade Completed	-0.052 (0.179)
Number of Siblings	-0.052 (0.160)
GED Cost	-0.053 (0.156)
Family Income in 1979	-0.047 (0.177)
Dropout's Local Wage at Age 17	-0.013 (0.218)
High School Graduate's Local Wage at Age 17	-0.049 (0.182)
Dropout's Local Unemployment Rate at Age 17	0.443 (1.051)
High School Graduate's Local Unemployment Rate at Age 17	-0.563 (0.577)
Propensity Score <sup>c</sup>	-0.058 (0.164)

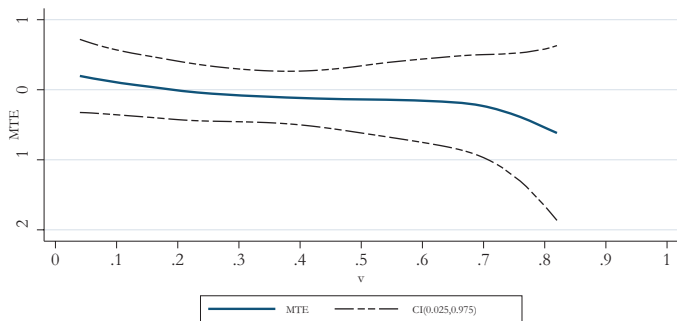
Notes: <sup>a</sup>The IV estimates are computed by taking the weighted sum of the MTE. The standard deviations (in parentheses) are computed using bootstrapping (50 draws).

<sup>b</sup>We excluded the oversample of poor whites and the military sample. The cost of the GED corresponds to the average testing fee per GED battery by state between 1993 and 2000 (Source: GED Statistical Report). Average local wage for dropouts and high school graduates correspond to the average in the place of residence for each group, respectively, and local unemployment rate corresponds to the unemployment rate in the place of residence. Average local wages, local unemployment rates, mother's and father's education refer to the level at age 17.

<sup>c</sup>The propensity score ( $P(D = 1|Z = z)$ ) is computed using as controls the instruments presented in the table, as well as two dummy variables controlling for the place of residence at age 14 (south and urban), and a set of dummy variables controlling for the year of birth (1957-1963).

Source: Heckman, Urzua and Vytlacil (2004).

Figure 15: MTE of the GED with Confidence Interval, Dropouts and GEDs—Males of the NLSY at Age 25



Note: The dependent variable in the outcome equation is the log of the average hourly wage reported between ages 24 and 26. The controls in the outcome equations are tenure, tenure squared, experience, corrected AFQT, black (dummy), Hispanic (dummy), marital status, and years of schooling. Let  $D = 0$  denote dropout status and  $D = 1$  denote GED status. The model for  $D$  (choice model) includes as controls the corrected AFQT, number of siblings, father's education, mother's education, family income at age 17, local GED costs, broken home at age 14, average local wage at age 17 for dropouts and high school graduates, local unemployment rate at age 17 for dropouts and high school graduates, the dummy variables for black and Hispanic, and a set of dummy variables controlling for year of birth. We also include two dummy variables controlling for the place of residence at age 14 (south and urban). The choice model is estimated using a probit model. In computing the MTE, the bandwidths are selected using the "leave one out" cross-validation method. We use biweight kernel functions. The confidence interval is computed from bootstrapping using 50 draws.

- Local linear regression is used to estimate the MTE implementing equation (19).

- Local linear regression is used to estimate the MTE implementing equation (19).
- While the standard error band is large, the estimated  $\Delta^{\text{MTE}}$  is in general negative, suggesting a negative marginal treatment effect for most participants.

- Local linear regression is used to estimate the MTE implementing equation (19).
- While the standard error band is large, the estimated  $\Delta^{\text{MTE}}$  is in general negative, suggesting a negative marginal treatment effect for most participants.
- However, we observe that for small values of  $u_D$  the point estimates of the marginal effect are positive.

- Local linear regression is used to estimate the MTE implementing equation (19).
- While the standard error band is large, the estimated  $\Delta^{\text{MTE}}$  is in general negative, suggesting a negative marginal treatment effect for most participants.
- However, we observe that for small values of  $u_D$  the point estimates of the marginal effect are positive.
- This analysis indicates that for people who are more likely to take the GED exam in terms of their unobservables (i.e., for people at the margin of indifference associated with a small  $u_D$ ), the marginal effect is in fact positive.

- It is instructive to examine the various IV estimates using the one instrument at a time strategy favored by many applied economists who like to do sensitivity analysis.



- It is instructive to examine the various IV estimates using the one instrument at a time strategy favored by many applied economists who like to do sensitivity analysis.
- Many of the variables used in the analysis are determined by age 17.

- It is instructive to examine the various IV estimates using the one instrument at a time strategy favored by many applied economists who like to do sensitivity analysis.
- Many of the variables used in the analysis are determined by age 17.
- Both father's highest grade completed and local unemployment rate among high school dropouts produce positive (if not precisely determined) IV estimates.

- It is instructive to examine the various IV estimates using the one instrument at a time strategy favored by many applied economists who like to do sensitivity analysis.
- Many of the variables used in the analysis are determined by age 17.
- Both father's highest grade completed and local unemployment rate among high school dropouts produce positive (if not precisely determined) IV estimates.
- A negative MTE weighted by negative IV weights produces a positive IV.

- A naive application of IV could produce the wrong causal inference, i.e., that GED certification raises wages.

- A naive application of IV could produce the wrong causal inference, i.e., that GED certification raises wages.
- Our estimates show that our theoretical examples have real world counterparts.

- A naive application of IV could produce the wrong causal inference, i.e., that GED certification raises wages.
- Our estimates show that our theoretical examples have real world counterparts.
- ? present an extensive empirical analysis of the wage returns to college attendance.

- A naive application of IV could produce the wrong causal inference, i.e., that GED certification raises wages.
- Our estimates show that our theoretical examples have real world counterparts.
- ? present an extensive empirical analysis of the wage returns to college attendance.
- They show how to unify and interpret diverse instruments within a common framework using the MTE and the weights derived in ???.

- A naive application of IV could produce the wrong causal inference, i.e., that GED certification raises wages.
- Our estimates show that our theoretical examples have real world counterparts.
- ? present an extensive empirical analysis of the wage returns to college attendance.
- They show how to unify and interpret diverse instruments within a common framework using the MTE and the weights derived in ???.
- They show negative weights on the MTE for commonly used instruments.



- ? use the MTE and the derived weights to identify the ranges of the MTE identified by different instruments in their analysis of the costs of breast cancer.

- ? use the MTE and the derived weights to identify the ranges of the MTE identified by different instruments in their analysis of the costs of breast cancer.
- We next discuss the implications of relaxing separability in the choice equations.

## Monotonicity, Uniformity, Nonseparability, Independence and Policy Invariance: The Limits of Instrumental Variables

- The analysis of this section and the entire recent literature on instrumental variables estimators for models with heterogeneous responses (i.e., models with outcomes of the forms (5) and (6)) relies critically on the assumption that the treatment choice equation has a representation in the additively separable form (7).

## Monotonicity, Uniformity, Nonseparability, Independence and Policy Invariance: The Limits of Instrumental Variables

- The analysis of this section and the entire recent literature on instrumental variables estimators for models with heterogeneous responses (i.e., models with outcomes of the forms (5) and (6)) relies critically on the assumption that the treatment choice equation has a representation in the additively separable form (7).
- From ?, we know that under assumptions (A-1)– (A-5), separability is equivalent to the assumption of monotonicity or uniformity, (IV-3).

- This uniformity condition imparts an asymmetry to the entire instrumental variable enterprise.

- This uniformity condition imparts an asymmetry to the entire instrumental variable enterprise.
- Responses are permitted to be heterogeneous in a general way, but choices of treatment are not.

- This uniformity condition imparts an asymmetry to the entire instrumental variable enterprise.
- Responses are permitted to be heterogeneous in a general way, but choices of treatment are not.
- In this section, we relax the assumption of additive separability in (7).

- This uniformity condition imparts an asymmetry to the entire instrumental variable enterprise.
- Responses are permitted to be heterogeneous in a general way, but choices of treatment are not.
- In this section, we relax the assumption of additive separability in (7).
- We establish that in the absence of additive separability or uniformity, the entire instrumental variable identification strategy in this section and the entire recent literature collapses.



- This uniformity condition imparts an asymmetry to the entire instrumental variable enterprise.
- Responses are permitted to be heterogeneous in a general way, but choices of treatment are not.
- In this section, we relax the assumption of additive separability in (7).
- We establish that in the absence of additive separability or uniformity, the entire instrumental variable identification strategy in this section and the entire recent literature collapses.
- Parameters can be defined as weighted averages of an MTE.

- MTE and the derived parameters cannot be identified using any instrumental variable strategy.

- MTE and the derived parameters cannot be identified using any instrumental variable strategy.
- Appendix, Slide 1049, presents a comprehensive discussion, which we summarize in this subsection.

- MTE and the derived parameters cannot be identified using any instrumental variable strategy.
- Appendix, Slide 1049, presents a comprehensive discussion, which we summarize in this subsection.
- One natural benchmark nonseparable model is a random coefficient model of choice  $D = \mathbf{1} [Z\gamma \geq 0]$ , where  $\gamma$  is a random coefficient vector and  $\gamma \perp\!\!\!\perp (Z, U_0, U_1)$ .

- MTE and the derived parameters cannot be identified using any instrumental variable strategy.
- Appendix, Slide 1049, presents a comprehensive discussion, which we summarize in this subsection.
- One natural benchmark nonseparable model is a random coefficient model of choice  $D = \mathbf{1} [Z\gamma \geq 0]$ , where  $\gamma$  is a random coefficient vector and  $\gamma \perp\!\!\!\perp (Z, U_0, U_1)$ .
- If  $\gamma$  is a random coefficient with a nondegenerate distribution and with components that take both positive and negative values, uniformity is clearly violated.

- MTE and the derived parameters cannot be identified using any instrumental variable strategy.
- Appendix, Slide 1049, presents a comprehensive discussion, which we summarize in this subsection.
- One natural benchmark nonseparable model is a random coefficient model of choice  $D = \mathbf{1} [Z\gamma \geq 0]$ , where  $\gamma$  is a random coefficient vector and  $\gamma \perp\!\!\!\perp (Z, U_0, U_1)$ .
- If  $\gamma$  is a random coefficient with a nondegenerate distribution and with components that take both positive and negative values, uniformity is clearly violated.
- However, it can be violated even when all components of  $\gamma$  are of the same sign if  $Z$  is a vector.

- Relax the additive separability assumption of equation (7) to consider a more general case

$$D^* = \mu_D(Z, V), \quad (29a)$$

where  $\mu_D(Z, V)$  is not necessarily additively separable in  $Z$  and  $V$ , and  $V$  is not necessarily a scalar.

- Relax the additive separability assumption of equation (7) to consider a more general case

$$D^* = \mu_D(Z, V), \quad (29a)$$

where  $\mu_D(Z, V)$  is not necessarily additively separable in  $Z$  and  $V$ , and  $V$  is not necessarily a scalar.

- In the random coefficient example,  $V = \gamma$  and  $\mu_D = z\gamma$ .



- Relax the additive separability assumption of equation (7) to consider a more general case

$$D^* = \mu_D(Z, V), \quad (29a)$$

where  $\mu_D(Z, V)$  is not necessarily additively separable in  $Z$  and  $V$ , and  $V$  is not necessarily a scalar.

- In the random coefficient example,  $V = \gamma$  and  $\mu_D = z\gamma$ .



$$D = \mathbf{1} [D^* \geq 0]. \quad (29b)$$

We maintain assumptions (A-1)–(A-5) and (A-7).

- In special cases, (29a) can be expressed in an additively separable form.

- In special cases, (29a) can be expressed in an additively separable form.
- For example, if  $D^*$  is weakly separable in  $Z$  and  $V$ ,  $D^* = \mu_D(\theta(Z), V)$  for any  $V$  where  $\theta(Z)$  is a scalar function,  $\mu_D$  is increasing in  $\theta(Z)$ , and  $V$  is a scalar, then we can write (29b) in the same form as (7):

$$D = \mathbf{1} \left[ \theta(Z) \geq \tilde{V} \right],$$

where  $\tilde{V} = \mu_D^{-1}(0; V)$  and  $\tilde{V} \perp\!\!\!\perp Z \mid X$ , and the inverse function is expressed with respect to the first argument (see ?).

- In special cases, (29a) can be expressed in an additively separable form.
- For example, if  $D^*$  is weakly separable in  $Z$  and  $V$ ,  $D^* = \mu_D(\theta(Z), V)$  for any  $V$  where  $\theta(Z)$  is a scalar function,  $\mu_D$  is increasing in  $\theta(Z)$ , and  $V$  is a scalar, then we can write (29b) in the same form as (7):

$$D = \mathbf{1} \left[ \theta(Z) \geq \tilde{V} \right],$$

where  $\tilde{V} = \mu_D^{-1}(0; V)$  and  $\tilde{V} \perp\!\!\!\perp Z \mid X$ , and the inverse function is expressed with respect to the first argument (see ?).

- ? shows that any model that does not satisfy uniformity (or “monotonicity”) will not have a representation in this form.

- In the additively separable case, the MTE (8) has three equivalent interpretations.

- In the additively separable case, the MTE (8) has three equivalent interpretations.
- (i)  $U_D = F_V(V)$  is the only unobservable in the first stage decision rule, and MTE is the average effect of treatment given the unobserved characteristics in the decision rule ( $V = v$ ).

- In the additively separable case, the MTE (8) has three equivalent interpretations.
- (i)  $U_D = F_V(V)$  is the only unobservable in the first stage decision rule, and MTE is the average effect of treatment given the unobserved characteristics in the decision rule ( $V = v$ ).
- (ii) A person with  $V = v$  would be indifferent between treatment or not if  $P(Z) = u_D$ , where  $P(Z)$  is a mean scale utility function.

- In the additively separable case, the MTE (8) has three equivalent interpretations.
- (i)  $U_D = F_V(V)$  is the only unobservable in the first stage decision rule, and MTE is the average effect of treatment given the unobserved characteristics in the decision rule ( $V = v$ ).
- (ii) A person with  $V = v$  would be indifferent between treatment or not if  $P(Z) = u_D$ , where  $P(Z)$  is a mean scale utility function.
- Thus, the MTE is the average effect of treatment given that the individual would be indifferent between treatment or not if  $P(Z) = u_D$ .



- (iii) One can also view the additively separable form (7) as intrinsic in the way we are defining the parameter and interpret the MTE (equation (8)) as an average effect conditional on the additive error term from the first stage choice model.

- (iii) One can also view the additively separable form (7) as intrinsic in the way we are defining the parameter and interpret the MTE (equation (8)) as an average effect conditional on the additive error term from the first stage choice model.
- Under all interpretations of the MTE and under the assumptions used in the preceding sections of this chapter, MTE can be identified by LIV; the MTE does not depend on  $Z$  and hence it is policy invariant and the MTE integrates up to generate all treatment effects, policy effects and all IV estimands.

- The three definitions are not the same in the general nonseparable case (29a).

- The three definitions are not the same in the general nonseparable case (29a).
- ? extend MTE in the nonseparable case using interpretation (i).

- The three definitions are not the same in the general nonseparable case (29a).
- ? extend MTE in the nonseparable case using interpretation (i).
- MTE defined this way is policy invariant to changes in  $Z$ . Appendix, Slide 1049, which summarizes their work, shows that LIV is a weighted average of the MTE with possibly negative weights and does not identify MTE.

- The three definitions are not the same in the general nonseparable case (29a).
- ? extend MTE in the nonseparable case using interpretation (i).
- MTE defined this way is policy invariant to changes in  $Z$ . Appendix, Slide 1049, which summarizes their work, shows that LIV is a weighted average of the MTE with possibly negative weights and does not identify MTE.
- If uniformity does not hold, the definition of MTE allows one to integrate MTE to obtain all of the treatment effects, but the instrumental variables estimator breaks down.

- Alternatively, one could define MTE based on (ii):

$$\Delta_{ii}^{\text{MTE}}(z) = E(Y_1 - Y_0 \mid V \in \{v : \mu_D(z, v) = 0\}).$$

This is the average treatment effect for individuals who would be indifferent between treatment or not at a given value of  $z$  (recall that we keep the conditioning on  $X$  implicit).

- Alternatively, one could define MTE based on (ii):

$$\Delta_{ii}^{\text{MTE}}(z) = E(Y_1 - Y_0 \mid V \in \{v : \mu_D(z, v) = 0\}).$$

This is the average treatment effect for individuals who would be indifferent between treatment or not at a given value of  $z$  (recall that we keep the conditioning on  $X$  implicit).

- ? show that in the nonseparable case LIV does not identify this MTE and that MTE does not change when the distribution of  $Z$  changes, provided that the support of MTE does not change.



- Alternatively, one could define MTE based on (ii):

$$\Delta_{ii}^{\text{MTE}}(z) = E(Y_1 - Y_0 \mid V \in \{v : \mu_D(z, v) = 0\}).$$

This is the average treatment effect for individuals who would be indifferent between treatment or not at a given value of  $z$  (recall that we keep the conditioning on  $X$  implicit).

- ? show that in the nonseparable case LIV does not identify this MTE and that MTE does not change when the distribution of  $Z$  changes, provided that the support of MTE does not change.
- In general, this definition of MTE does not allow one to integrate up MTE to obtain the treatment parameters.

- A third possibility is to force the index rule into an additive form by taking  $\mu_D^*(Z) = E(\mu_D(Z, V) | Z)$ , defining  $V^* = \mu_D(Z, V) - E(\mu_D(Z, V) | Z)$  and define MTE as  $E(Y_1 - Y_0 | V^* = v^*)$ .

- A third possibility is to force the index rule into an additive form by taking  $\mu_D^*(Z) = E(\mu_D(Z, V) | Z)$ , defining  $V^* = \mu_D(Z, V) - E(\mu_D(Z, V) | Z)$  and define MTE as  $E(Y_1 - Y_0 | V^* = v^*)$ .
- Note that  $V^*$  is not independent of  $Z$ , is not policy invariant and is not structural.

- A third possibility is to force the index rule into an additive form by taking  $\mu_D^*(Z) = E(\mu_D(Z, V) | Z)$ , defining  $V^* = \mu_D(Z, V) - E(\mu_D(Z, V) | Z)$  and define MTE as  $E(Y_1 - Y_0 | V^* = v^*)$ .
- Note that  $V^*$  is not independent of  $Z$ , is not policy invariant and is not structural.
- LIV does not estimate this MTE.

- A third possibility is to force the index rule into an additive form by taking  $\mu_D^*(Z) = E(\mu_D(Z, V) | Z)$ , defining  $V^* = \mu_D(Z, V) - E(\mu_D(Z, V) | Z)$  and define MTE as  $E(Y_1 - Y_0 | V^* = v^*)$ .
- Note that  $V^*$  is not independent of  $Z$ , is not policy invariant and is not structural.
- LIV does not estimate this MTE.
- With this definition of the MTE it is not possible, in general, to integrate up MTE to obtain the various treatment effects.

- For any version of the nonseparable model, except those that can be transformed to separability, index sufficiency fails.

- For any version of the nonseparable model, except those that can be transformed to separability, index sufficiency fails.
- To see this, assume that  $\mu_D(Z, V)$  is continuous.

- For any version of the nonseparable model, except those that can be transformed to separability, index sufficiency fails.
- To see this, assume that  $\mu_D(Z, V)$  is continuous.
- Define  $\Omega(z) = \{v : \mu_D(z, v) \geq 0\}$ .



- For any version of the nonseparable model, except those that can be transformed to separability, index sufficiency fails.
- To see this, assume that  $\mu_D(Z, V)$  is continuous.
- Define  $\Omega(z) = \{v : \mu_D(z, v) \geq 0\}$ .
- In the additively separable case,  
 $P(z) \equiv \Pr(D = 1 \mid Z = z) = \Pr(U_D \in \Omega(z))$ ,  
 $P(z) = P(z') \Leftrightarrow \Omega(z) = \Omega(z')$ .

- For any version of the nonseparable model, except those that can be transformed to separability, index sufficiency fails.
- To see this, assume that  $\mu_D(Z, V)$  is continuous.
- Define  $\Omega(z) = \{v : \mu_D(z, v) \geq 0\}$ .
- In the additively separable case,  
 $P(z) \equiv \Pr(D = 1 \mid Z = z) = \Pr(U_D \in \Omega(z))$ ,  
 $P(z) = P(z') \Leftrightarrow \Omega(z) = \Omega(z')$ .
- This produces index sufficiency.

- For any version of the nonseparable model, except those that can be transformed to separability, index sufficiency fails.
- To see this, assume that  $\mu_D(Z, V)$  is continuous.
- Define  $\Omega(z) = \{v : \mu_D(z, v) \geq 0\}$ .
- In the additively separable case,  
 $P(z) \equiv \Pr(D = 1 \mid Z = z) = \Pr(U_D \in \Omega(z))$ ,  
 $P(z) = P(z') \Leftrightarrow \Omega(z) = \Omega(z')$ .
- This produces index sufficiency.
- In the more general case of (29a), it is possible to have  $(z, z')$  such that  $P(z) = P(z')$  and  $\Omega(z) \neq \Omega(z')$  so index sufficiency does not hold.

## Implications of Nonseparability

- This section develops generalization (i), leaving development of the other interpretations for later research.

## Implications of Nonseparability

- This section develops generalization (i), leaving development of the other interpretations for later research.
- We focus on an analysis of PRTE, comparing two policies  $p, p' \in \mathcal{P}$ .

## Implications of Nonseparability

- This section develops generalization (i), leaving development of the other interpretations for later research.
- We focus on an analysis of PRTE, comparing two policies  $p, p' \in \mathcal{P}$ .
- Here “ $p$ ” denotes a policy and not a realization of  $P(Z)$  as in the previous sections.

## Implications of Nonseparability

- This section develops generalization (i), leaving development of the other interpretations for later research.
- We focus on an analysis of PRTE, comparing two policies  $p, p' \in \mathcal{P}$ .
- Here “ $p$ ” denotes a policy and not a realization of  $P(Z)$  as in the previous sections.
- This is our convention when we discuss PRTE.

## Implications of Nonseparability

- This section develops generalization (i), leaving development of the other interpretations for later research.
- We focus on an analysis of PRTE, comparing two policies  $p, p' \in \mathcal{P}$ .
- Here “ $p$ ” denotes a policy and not a realization of  $P(Z)$  as in the previous sections.
- This is our convention when we discuss PRTE.
- The analysis of the other treatment parameters follows by parallel arguments.



- For any  $v$  in the support of the distribution of  $V$ , define  $\Omega = \{z : \mu_D(z, v) \geq 0\}$ .

- For any  $v$  in the support of the distribution of  $V$ , define  $\Omega = \{z : \mu_D(z, v) \geq 0\}$ .
- For example, in the random coefficient case, with  $V \equiv \gamma$  and  $D = \mathbf{1}[Z\gamma \geq 0]$ , we have  $\Omega_g = \{z : zg \geq 0\}$ , where  $g$  is a realization of  $\gamma$ .

- For any  $v$  in the support of the distribution of  $V$ , define  $\Omega = \{z : \mu_D(z, v) \geq 0\}$ .
- For example, in the random coefficient case, with  $V \equiv \gamma$  and  $D = \mathbf{1}[Z\gamma \geq 0]$ , we have  $\Omega_g = \{z : zg \geq 0\}$ , where  $g$  is a realization of  $\gamma$ .
- Define  $\mathbf{1}_{\mathcal{A}}(t)$  to be the indicator function for the event  $t \in \mathcal{A}$ .

- Then, making the  $X$  explicit, Appendix, Slide 1049, derives the result that

$$E(Y_p) - E(Y_{p'}) = E[E(Y_p|X) - E(Y_{p'} | X)] \quad (30)$$

$$= \int \left[ \int E(\Delta^{\text{MTE}} | X = x, V = v) \right. \quad (31)$$

$$\times \left. \left( \begin{array}{c} \Pr[Z_p \in \Omega | X = x] \\ - \Pr[Z_{p'} \in \Omega | X = x] \end{array} \right) dF_{V|X}(v|x) \right] dF_X(x).$$

- Then, making the  $X$  explicit, Appendix, Slide 1049, derives the result that

$$E(Y_p) - E(Y_{p'}) = E[E(Y_p|X) - E(Y_{p'} | X)] \quad (30)$$

$$= \int \left[ \int E(\Delta^{\text{MTE}} | X = x, V = v) \right. \quad (31)$$

$$\times \left. \left( \begin{array}{c} \Pr[Z_p \in \Omega | X = x] \\ - \Pr[Z_{p'} \in \Omega | X = x] \end{array} \right) dF_{V|X}(v|x) \right] dF_X(x).$$

- Thus, without additive separability, we can still derive an expression for PRTE and by similar reasoning the other treatment parameters.

- Then, making the  $X$  explicit, Appendix, Slide 1049, derives the result that

$$E(Y_p) - E(Y_{p'}) = E[E(Y_p|X) - E(Y_{p'} | X)] \quad (30)$$

$$= \int \left[ \int E(\Delta^{\text{MTE}} | X = x, V = v) \right. \quad (31)$$

$$\left. \times \left( \begin{array}{c} \Pr[Z_p \in \Omega | X = x] \\ - \Pr[Z_{p'} \in \Omega | X = x] \end{array} \right) dF_{V|X}(v|x) \right] dF_X(x).$$

- Thus, without additive separability, we can still derive an expression for PRTE and by similar reasoning the other treatment parameters.
- However, to evaluate the expression requires knowledge of MTE, of  $\Pr[Z_p \in \Omega | X = x]$  and  $\Pr[Z_{p'} \in \Omega | X = x]$  for every  $(v, x)$  in the support of the distribution of  $(V, X)$ , and of the distribution of  $V$ .

- In general, if no structure is placed on the  $\mu_D$  function, one can normalize  $V$  to be unit uniform (or a vector of unit uniform random variables) so that  $F_{V|X}$  will be known.

- In general, if no structure is placed on the  $\mu_D$  function, one can normalize  $V$  to be unit uniform (or a vector of unit uniform random variables) so that  $F_{V|X}$  will be known.
- However, in this case, the  $\Omega = \{z : \mu_D(z, v) \geq 0\}$  sets will not in general be identified.



- In general, if no structure is placed on the  $\mu_D$  function, one can normalize  $V$  to be unit uniform (or a vector of unit uniform random variables) so that  $F_{V|X}$  will be known.
- However, in this case, the  $\Omega = \{z : \mu_D(z, v) \geq 0\}$  sets will not in general be identified.
- If structure is placed on the  $\mu_D$  function, one might be able to identify the  $\Omega = \{z : \mu_D(z, v) \geq 0\}$  sets but then one needs to identify the distribution of  $V$  (conditional on  $X$ ).

- In general, if no structure is placed on the  $\mu_D$  function, one can normalize  $V$  to be unit uniform (or a vector of unit uniform random variables) so that  $F_{V|X}$  will be known.
- However, in this case, the  $\Omega = \{z : \mu_D(z, v) \geq 0\}$  sets will not in general be identified.
- If structure is placed on the  $\mu_D$  function, one might be able to identify the  $\Omega = \{z : \mu_D(z, v) \geq 0\}$  sets but then one needs to identify the distribution of  $V$  (conditional on  $X$ ).
- If structure is placed on  $\mu_D$ , one cannot in general normalize the distribution of  $V$  to be unit uniform without undoing the structure being imposed on  $\mu_D$ .

- In particular, consider the random coefficient model  $D = \mathbf{1}[Z\gamma \geq 0]$  where  $V = \gamma$  is a random vector, so that  $\Omega_\gamma = \{z : z\gamma \geq 0\}$ .

- In particular, consider the random coefficient model  $D = \mathbf{1}[Z\gamma \geq 0]$  where  $V = \gamma$  is a random vector, so that  $\Omega_\gamma = \{z : z\gamma \geq 0\}$ .
- In this case, if all of the other assumptions hold, including  $Z \perp\!\!\!\perp \gamma \mid X$ , and the policy change does not affect  $(Y_1, Y_0, X, \gamma)$ , the PRTE is given by

$$\begin{aligned}
 E(Y_p) - E(Y_{p'}) &= E[E(Y_p|X) - E(Y_{p'}|X)] \\
 &= \int \left[ \int E(\Delta^{\text{MTE}} \mid X = x, \gamma = g) \right. \\
 &\quad \times \left. \begin{pmatrix} \Pr[Z_p \in \Omega_g \mid X = x] \\ -\Pr[Z_{p'} \in \Omega_g \mid X = x] \end{pmatrix} dF_{\gamma|X}(g|x) \right] dF_X(x).
 \end{aligned}$$

- Because structure has been placed on the  $\mu_D(Z, \gamma)$  function, the sets  $\Omega_\gamma$  are known.

- Because structure has been placed on the  $\mu_D(Z, \gamma)$  function, the sets  $\Omega_\gamma$  are known.
- However, evaluating the function requires knowledge of the distribution of  $\gamma$  which will not in general be identified without further assumptions.

- Because structure has been placed on the  $\mu_D(Z, \gamma)$  function, the sets  $\Omega_\gamma$  are known.
- However, evaluating the function requires knowledge of the distribution of  $\gamma$  which will not in general be identified without further assumptions.
- Normalizing the distribution of  $\gamma$  to be a vector of unit uniform random variables produces the distribution of  $\gamma$  but eliminates the assumed linear index structure on  $\mu_D$  and results in  $\Omega_\gamma$  sets that are not identified.

- Even if the weights are identified, ? show that it is not possible to use LIV to identify MTE without additive separability between  $Z$  and  $V$  in the selection rule index.



- Even if the weights are identified, ? show that it is not possible to use LIV to identify MTE without additive separability between  $Z$  and  $V$  in the selection rule index.
- Appendix, Slide 1119, develops this point for the random coefficient model.

- Even if the weights are identified, ? show that it is not possible to use LIV to identify MTE without additive separability between  $Z$  and  $V$  in the selection rule index.
- Appendix, Slide 1119, develops this point for the random coefficient model.
- Without additive separability in the latent index for the selection rule, we can still create an expression for PRTE (and the other treatment parameters) but both the weights and the MTE function are no longer identified using instrumental variables.

- One superficially plausible way to avoid these problems would be to define  $\tilde{\mu}_D(Z) = E(\mu_D(Z, V) | Z)$  and  $\tilde{V} = \mu_D(Z, V) - E(\mu_D(Z, V) | Z)$ , producing the model  $D = \mathbf{1}[\tilde{\mu}_D(Z) + \tilde{V} \geq 0]$ .

- One superficially plausible way to avoid these problems would be to define  $\tilde{\mu}_D(Z) = E(\mu_D(Z, V) | Z)$  and  $\tilde{V} = \mu_D(Z, V) - E(\mu_D(Z, V) | Z)$ , producing the model  $D = \mathbf{1}[\tilde{\mu}_D(Z) + \tilde{V} \geq 0]$ .
- We keep the conditioning on  $X$  implicit.

- One superficially plausible way to avoid these problems would be to define  $\tilde{\mu}_D(Z) = E(\mu_D(Z, V) | Z)$  and  $\tilde{V} = \mu_D(Z, V) - E(\mu_D(Z, V) | Z)$ , producing the model  $D = \mathbf{1}[\tilde{\mu}_D(Z) + \tilde{V} \geq 0]$ .
- We keep the conditioning on  $X$  implicit.
- One could redefine MTE using  $\tilde{V}$  and proceed as if the true model possessed additive separability between observables and unobservables in the latent index.

- One superficially plausible way to avoid these problems would be to define  $\tilde{\mu}_D(Z) = E(\mu_D(Z, V) | Z)$  and  $\tilde{V} = \mu_D(Z, V) - E(\mu_D(Z, V) | Z)$ , producing the model  $D = \mathbf{1}[\tilde{\mu}_D(Z) + \tilde{V} \geq 0]$ .
- We keep the conditioning on  $X$  implicit.
- One could redefine MTE using  $\tilde{V}$  and proceed as if the true model possessed additive separability between observables and unobservables in the latent index.
- This is the method pursued in approach (iii).

- For two reasons, this approach does not solve the problem of providing an adequate generalization of MTE.

- For two reasons, this approach does not solve the problem of providing an adequate generalization of MTE.
- First, with this definition,  $\tilde{V}$  is a function of  $(Z, V)$ , and a policy that changes  $Z$  will then also change  $\tilde{V}$ .



- For two reasons, this approach does not solve the problem of providing an adequate generalization of MTE.
- First, with this definition,  $\tilde{V}$  is a function of  $(Z, V)$ , and a policy that changes  $Z$  will then also change  $\tilde{V}$ .
- Thus, policy invariance of the MTE no longer holds.

- For two reasons, this approach does not solve the problem of providing an adequate generalization of MTE.
- First, with this definition,  $\tilde{V}$  is a function of  $(Z, V)$ , and a policy that changes  $Z$  will then also change  $\tilde{V}$ .
- Thus, policy invariance of the MTE no longer holds.
- Second, this approach generates a  $\tilde{V}$  that is no longer statistically independent of  $Z$  so that assumption (A-1) no longer holds when  $\tilde{V}$  is substituted for  $V$  even when (A-1) is true for  $V$ .

- For two reasons, this approach does not solve the problem of providing an adequate generalization of MTE.
- First, with this definition,  $\tilde{V}$  is a function of  $(Z, V)$ , and a policy that changes  $Z$  will then also change  $\tilde{V}$ .
- Thus, policy invariance of the MTE no longer holds.
- Second, this approach generates a  $\tilde{V}$  that is no longer statistically independent of  $Z$  so that assumption (A-1) no longer holds when  $\tilde{V}$  is substituted for  $V$  even when (A-1) is true for  $V$ .
- Lack of independence between observables and unobservables in the latent index both invalidates our expression for PRTE (and the expressions for the other treatment effects) and causes LIV to no longer identify MTE.

- The nonseparable model can also restrict the support of  $P(Z)$ .

- The nonseparable model can also restrict the support of  $P(Z)$ .
- For example, consider a standard normal random coefficient model with a scalar regressor ( $Z = (1, Z_1)$ ).

- The nonseparable model can also restrict the support of  $P(Z)$ .
- For example, consider a standard normal random coefficient model with a scalar regressor ( $Z = (1, Z_1)$ ).
- Assume  $\gamma_0 \sim N(0, \sigma_0^2)$ ,  $\gamma_1 \sim N(\bar{\gamma}_1, \sigma_1^2)$ , and  $\gamma_0 \perp\!\!\!\perp \gamma_1$ .

- The nonseparable model can also restrict the support of  $P(Z)$ .
- For example, consider a standard normal random coefficient model with a scalar regressor ( $Z = (1, Z_1)$ ).
- Assume  $\gamma_0 \sim N(0, \sigma_0^2)$ ,  $\gamma_1 \sim N(\bar{\gamma}_1, \sigma_1^2)$ , and  $\gamma_0 \perp\!\!\!\perp \gamma_1$ .
- Then

$$P(z_1) = \Phi\left(\frac{\bar{\gamma}_1 z_1}{\sqrt{\sigma_0^2 + \sigma_1^2 z_1^2}}\right),$$

where  $\Phi$  is the standard cumulative normal distribution.

- If the support of  $z_1$  is  $\mathbb{R}$ , then in the standard additive model,  $\sigma_1^2 = 0$  and  $P(z_1)$  has support  $[0, 1]$ .



- If the support of  $z_1$  is  $\mathbb{R}$ , then in the standard additive model,  $\sigma_1^2 = 0$  and  $P(z_1)$  has support  $[0, 1]$ .
- When  $\sigma_1^2 > 0$ , the support is strictly within the unit interval.

- If the support of  $z_1$  is  $\mathbb{R}$ , then in the standard additive model,  $\sigma_1^2 = 0$  and  $P(z_1)$  has support  $[0, 1]$ .
- When  $\sigma_1^2 > 0$ , the support is strictly within the unit interval.
- In the special case when  $\sigma_0^2 = 0$ , the support is one point  $\left( P(z) = \Phi \left( \frac{\bar{\gamma}_1}{\sigma_1} \right) \right)$ .

- If the support of  $z_1$  is  $\mathbb{R}$ , then in the standard additive model,  $\sigma_1^2 = 0$  and  $P(z_1)$  has support  $[0, 1]$ .
- When  $\sigma_1^2 > 0$ , the support is strictly within the unit interval.
- In the special case when  $\sigma_0^2 = 0$ , the support is one point  $\left( P(z) = \Phi\left(\frac{\bar{z}_1}{\sigma_1}\right) \right)$ .
- We cannot, in general, identify ATE, TT or any treatment effect requiring the endpoints 0 or 1.

- Thus the general models of nonuniformity presented in this section do not satisfy the index sufficiency property, and the support of the treatment effects and estimators is, in general, less than full.

- Thus the general models of nonuniformity presented in this section do not satisfy the index sufficiency property, and the support of the treatment effects and estimators is, in general, less than full.
- The random coefficient model for choice may explain the empirical support problems for  $P(Z)$  found in ? and many other evaluation studies.

## Implications of Dependence

- We next consider relaxing the independence assumption (A-1) to allow  $Z \not\perp\!\!\!\perp V \mid X$  while maintaining the assumption that  $Z \perp\!\!\!\perp (Y_0, Y_1) \mid (X, V)$ .

## Implications of Dependence

- We next consider relaxing the independence assumption (A-1) to allow  $Z \not\perp\!\!\!\perp V \mid X$  while maintaining the assumption that  $Z \perp\!\!\!\perp (Y_0, Y_1) \mid (X, V)$ .
- We maintain the other assumptions, including additive separability between  $Z$  and  $V$  in the latent index for the selection rule (equation (7)) and the assumption that the policy changes  $Z$  but does not change  $(V, Y_0, Y_1, X)$ .

## Implications of Dependence

- We next consider relaxing the independence assumption (A-1) to allow  $Z \not\perp V \mid X$  while maintaining the assumption that  $Z \perp\!\!\!\perp (Y_0, Y_1) \mid (X, V)$ .
- We maintain the other assumptions, including additive separability between  $Z$  and  $V$  in the latent index for the selection rule (equation (7)) and the assumption that the policy changes  $Z$  but does not change  $(V, Y_0, Y_1, X)$ .
- Thus we assume that the policy shift does not change the MTE function (policy invariance).



- Given these assumptions, we derive in Appendix, Slide 1082, the following expression for PRTE in the nonindependent case for policies  $p, p' \in \mathcal{P}$ :

$$E(Y_p) - E(Y_{p'}) = E[E(Y_p | X) - E(Y_{p'} | X)] \quad (32)$$

$$= \int \left[ \int E\left(\Delta^{\text{MTE}} \mid \begin{array}{l} X = x, \\ V = v \end{array} \right) \right. \quad (33)$$

$$\left. \times \left( \begin{array}{l} \Pr[\mu_D(Z_{p'}) < v \mid X = x, V = v] \\ - \Pr[\mu_D(Z_p) < v \mid X = x, V = v] \end{array} \right) dF_{V|X}(v | x) \right] dF_X(x).$$

- Given these assumptions, we derive in Appendix, Slide 1082, the following expression for PRTE in the nonindependent case for policies  $p, p' \in \mathcal{P}$ :

$$E(Y_p) - E(Y_{p'}) = E[E(Y_p | X) - E(Y_{p'} | X)] \quad (32)$$

$$= \int \left[ \int E\left(\Delta^{\text{MTE}} \mid \begin{array}{l} X = x, \\ V = v \end{array} \right) \right. \quad (33)$$

$$\times \left( \begin{array}{l} \Pr[\mu_D(Z_{p'}) < v \mid X = x, V = v] \\ - \Pr[\mu_D(Z_p) < v \mid X = x, V = v] \end{array} \right) dF_{V|X}(v | x) \Big] dF_X(x).$$

- Notice that “ $p$ ” denotes a policy and not a realized value of  $P(Z)$ .

- Although we can derive an expression for PRTE without requiring independence between  $Z$  and  $V$ , to evaluate this expression requires knowledge of MTE and of  $Pr[\mu_D(Z_{p'}) < v \mid X = x, V = v]$  and of  $Pr[\mu_D(Z_p) < v \mid X = x, V = v]$  for every  $(x, v)$  in the support of the distribution of  $(X, V)$ .

- Although we can derive an expression for PRTE without requiring independence between  $Z$  and  $V$ , to evaluate this expression requires knowledge of MTE and of  $Pr[\mu_D(Z_{p'}) < v \mid X = x, V = v]$  and of  $Pr[\mu_D(Z_p) < v \mid X = x, V = v]$  for every  $(x, v)$  in the support of the distribution of  $(X, V)$ .
- This requirement is stronger than what is needed in the case of independence since the weights no longer depend only on the distribution of  $P_p(Z_p)$  and  $P_{p'}(Z_{p'})$  conditional on  $X$ .

- Although we can derive an expression for PRTE without requiring independence between  $Z$  and  $V$ , to evaluate this expression requires knowledge of MTE and of  $Pr[\mu_D(Z_{p'}) < v \mid X = x, V = v]$  and of  $Pr[\mu_D(Z_p) < v \mid X = x, V = v]$  for every  $(x, v)$  in the support of the distribution of  $(X, V)$ .
- This requirement is stronger than what is needed in the case of independence since the weights no longer depend only on the distribution of  $P_p(Z_p)$  and  $P_{p'}(Z_{p'})$  conditional on  $X$ .
- To evaluate these weights requires knowledge of the function  $\mu_D$  and of the joint distribution of  $(V, Z_p)$  and  $(V, Z_{p'})$  conditional on  $X$ , and these will in general not be identified without further assumptions.

- Even if the weights are identified, ? show that it is not possible to use LIV to identify MTE without independence between  $Z$  and  $V$  conditional on  $X$ .

- Even if the weights are identified, ? show that it is not possible to use LIV to identify MTE without independence between  $Z$  and  $V$  conditional on  $X$ .
- Thus, without conditional independence between  $Z$  and  $V$  in the latent index for the decision rule, we can still create an expression for PRTE but both the weights and the MTE function are no longer identified without invoking further assumptions.

- One superficially appealing way to avoid these problems is to define  $\tilde{V} = F_{V|X,Z}(V)$  and  $\tilde{\mu}_D(Z) = F_{V|X,Z}(\mu_D(Z))$ , so  $D = \mathbf{1}[\mu_D(Z) - V \geq 0] = \mathbf{1}[\tilde{\mu}_D(Z) - \tilde{V} \geq 0]$  with  $\tilde{V} \sim \text{Unif}[0, 1]$  conditional on  $X$  and  $Z$  and so  $\tilde{V}$  is independent of  $X$  and  $Z$ .



- One superficially appealing way to avoid these problems is to define  $\tilde{V} = F_{V|X,Z}(V)$  and  $\tilde{\mu}_D(Z) = F_{V|X,Z}(\mu_D(Z))$ , so  $D = \mathbf{1}[\mu_D(Z) - V \geq 0] = \mathbf{1}[\tilde{\mu}_D(Z) - \tilde{V} \geq 0]$  with  $\tilde{V} \sim \text{Unif}[0, 1]$  conditional on  $X$  and  $Z$  and so  $\tilde{V}$  is independent of  $X$  and  $Z$ .
- It might seem that the previous analysis would carry over.

- One superficially appealing way to avoid these problems is to define  $\tilde{V} = F_{V|X,Z}(V)$  and  $\tilde{\mu}_D(Z) = F_{V|X,Z}(\mu_D(Z))$ , so  $D = \mathbf{1}[\mu_D(Z) - V \geq 0] = \mathbf{1}[\tilde{\mu}_D(Z) - \tilde{V} \geq 0]$  with  $\tilde{V} \sim \text{Unif}[0, 1]$  conditional on  $X$  and  $Z$  and so  $\tilde{V}$  is independent of  $X$  and  $Z$ .
- It might seem that the previous analysis would carry over.
- However, by defining  $\tilde{V} = F_{V|X,Z}(V)$ , we have defined  $\tilde{V}$  in a way that depends functionally on  $Z$  and  $X$ , and hence we violate invariance of the MTE with respect to the shifts in the distribution of  $Z$  given  $X$ .

## The Limits of Instrumental Variable Estimators

- The treatment effect literature focuses on a class of policies that move treatment choices in the same direction for everyone.

## The Limits of Instrumental Variable Estimators

- The treatment effect literature focuses on a class of policies that move treatment choices in the same direction for everyone.
- General instruments do not have universally positive weights on  $\Delta^{\text{MTE}}$ .

## The Limits of Instrumental Variable Estimators

- The treatment effect literature focuses on a class of policies that move treatment choices in the same direction for everyone.
- General instruments do not have universally positive weights on  $\Delta^{\text{MTE}}$ .
- They are not guaranteed to shift everyone in the same direction.

## The Limits of Instrumental Variable Estimators

- The treatment effect literature focuses on a class of policies that move treatment choices in the same direction for everyone.
- General instruments do not have universally positive weights on  $\Delta^{\text{MTE}}$ .
- They are not guaranteed to shift everyone in the same direction.
- They do not necessarily estimate gross treatment effects.

- However, the effect of treatment is not always the parameter of policy interest.

- However, the effect of treatment is not always the parameter of policy interest.
- Thus, in the housing subsidy example developed in Slide 268, migration is the vehicle through which the policy operates.



- However, the effect of treatment is not always the parameter of policy interest.
- Thus, in the housing subsidy example developed in Slide 268, migration is the vehicle through which the policy operates.
- One might be interested in the effect of migration (the treatment effect) or the effect of the policy (the housing subsidy).

- However, the effect of treatment is not always the parameter of policy interest.
- Thus, in the housing subsidy example developed in Slide 268, migration is the vehicle through which the policy operates.
- One might be interested in the effect of migration (the treatment effect) or the effect of the policy (the housing subsidy).
- These are separate issues unless the policy is the treatment.

- Generalizing the MTE to the case of a nonseparable choice equation that violates the monotonicity condition, we can define but cannot identify the policy parameters of interest using ordinary instrumental variables or our extension LIV.

- Generalizing the MTE to the case of a nonseparable choice equation that violates the monotonicity condition, we can define but cannot identify the policy parameters of interest using ordinary instrumental variables or our extension LIV.
- If we make the model symmetrically heterogeneous in outcome and choice equations, the method of instrumental variables and our extensions of it break down in terms of estimating economically interpretable parameters.

- Generalizing the MTE to the case of a nonseparable choice equation that violates the monotonicity condition, we can define but cannot identify the policy parameters of interest using ordinary instrumental variables or our extension LIV.
- If we make the model symmetrically heterogeneous in outcome and choice equations, the method of instrumental variables and our extensions of it break down in terms of estimating economically interpretable parameters.
- ? and ? restore symmetry in the IV analysis of treatment choice and outcome equations by imposing uniformity on both outcome and choice equations.

- The general case of heterogeneity in both treatment and choice equations is beyond the outer limits of the entire IV literature, although it captures intuitively plausible phenomena.

- The general case of heterogeneity in both treatment and choice equations is beyond the outer limits of the entire IV literature, although it captures intuitively plausible phenomena.
- More general structural methods are required.

## Regression Discontinuity Estimators and LATE

- ? developed the regression discontinuity design which is now widely used.



## Regression Discontinuity Estimators and LATE

- ? developed the regression discontinuity design which is now widely used.
- (See an early discussion of this estimator in econometrics by ?).

## Regression Discontinuity Estimators and LATE

- ? developed the regression discontinuity design which is now widely used.
- (See an early discussion of this estimator in econometrics by ?).
- ? present an exposition of the regression discontinuity estimator within a LATE framework.

## Regression Discontinuity Estimators and LATE

- ? developed the regression discontinuity design which is now widely used.
- (See an early discussion of this estimator in econometrics by ?).
- ? present an exposition of the regression discontinuity estimator within a LATE framework.
- This section expositis the regression discontinuity method within our MTE framework.

- Suppose assumptions (A-1)–(A-5) hold except that we relax independence assumption (A-1) to assume that  $(Y_1 - Y_0, U_D)$  is independent of  $Z$  conditional on  $X$ .

- Suppose assumptions (A-1)–(A-5) hold except that we relax independence assumption (A-1) to assume that  $(Y_1 - Y_0, U_D)$  is independent of  $Z$  conditional on  $X$ .
- We *do not* impose the condition that  $Y_0$  is independent of  $Z$  conditional on  $X$ .

- Suppose assumptions (A-1)–(A-5) hold except that we relax independence assumption (A-1) to assume that  $(Y_1 - Y_0, U_D)$  is independent of  $Z$  conditional on  $X$ .
- We *do not* impose the condition that  $Y_0$  is independent of  $Z$  conditional on  $X$ .
- Relaxing the assumption that  $Y_0$  is independent of  $Z$  conditional on  $X$  causes the standard LIV estimand to differ from the MTE.

- Suppose assumptions (A-1)–(A-5) hold except that we relax independence assumption (A-1) to assume that  $(Y_1 - Y_0, U_D)$  is independent of  $Z$  conditional on  $X$ .
- We *do not* impose the condition that  $Y_0$  is independent of  $Z$  conditional on  $X$ .
- Relaxing the assumption that  $Y_0$  is independent of  $Z$  conditional on  $X$  causes the standard LIV estimand to differ from the MTE.
- We show that the LIV estimand in this case equals MTE plus a bias term that depends on  $\frac{\partial}{\partial p} E(Y_0 | X = x, P(Z) = p)$ .

- Suppose assumptions (A-1)–(A-5) hold except that we relax independence assumption (A-1) to assume that  $(Y_1 - Y_0, U_D)$  is independent of  $Z$  conditional on  $X$ .
- We *do not* impose the condition that  $Y_0$  is independent of  $Z$  conditional on  $X$ .
- Relaxing the assumption that  $Y_0$  is independent of  $Z$  conditional on  $X$  causes the standard LIV estimand to differ from the MTE.
- We show that the LIV estimand in this case equals MTE plus a bias term that depends on  $\frac{\partial}{\partial p} E(Y_0 | X = x, P(Z) = p)$ .
- Likewise, we show that the discrete-difference IV formula will no longer correspond to LATE, but will now correspond to LATE plus a bias term.



- A regression discontinuity design allows analysts to recover a LATE parameter at a particular value of  $Z$ .

- A regression discontinuity design allows analysts to recover a LATE parameter at a particular value of  $Z$ .
- If  $E(Y_0|X = x, Z = z)$  is continuous in  $z$ , while  $P(z)$  is discontinuous in  $z$  at a particular point, then it will be possible to use a regression discontinuity design to recover a LATE parameter.

- A regression discontinuity design allows analysts to recover a LATE parameter at a particular value of  $Z$ .
- If  $E(Y_0|X = x, Z = z)$  is continuous in  $z$ , while  $P(z)$  is discontinuous in  $z$  at a particular point, then it will be possible to use a regression discontinuity design to recover a LATE parameter.
- While the regression discontinuity design does have the advantage of allowing  $Y_0$  to depend on  $Z$  conditional on  $X$ , it only recovers a LATE parameter at a particular value of  $Z$  and cannot in general be used to recover either other treatment parameters such as the average treatment effect or the answers to policy questions such as the PRTE.

- A regression discontinuity design allows analysts to recover a LATE parameter at a particular value of  $Z$ .
- If  $E(Y_0|X = x, Z = z)$  is continuous in  $z$ , while  $P(z)$  is discontinuous in  $z$  at a particular point, then it will be possible to use a regression discontinuity design to recover a LATE parameter.
- While the regression discontinuity design does have the advantage of allowing  $Y_0$  to depend on  $Z$  conditional on  $X$ , it only recovers a LATE parameter at a particular value of  $Z$  and cannot in general be used to recover either other treatment parameters such as the average treatment effect or the answers to policy questions such as the PRTE.
- The following discussion is motivated by the analysis of ?.

- For simplicity, assume that  $Z$  is a scalar random variable.

- For simplicity, assume that  $Z$  is a scalar random variable.
- First, consider LIV while relaxing independence assumption (A-1) to assume that  $(Y_1 - Y_0, U_D)$  is independent of  $Z$  conditional on  $X$  but without imposing that  $Y_0$  is independent of  $Z$  conditional on  $X$ .

- For simplicity, assume that  $Z$  is a scalar random variable.
- First, consider LIV while relaxing independence assumption (A-1) to assume that  $(Y_1 - Y_0, U_D)$  is independent of  $Z$  conditional on  $X$  but without imposing that  $Y_0$  is independent of  $Z$  conditional on  $X$ .
- In order to make the comparison with the regression discontinuity design easier, we will condition on  $Z$  instead of  $P(Z)$ .

- Using  $Y = Y_0 + D(Y_1 - Y_0)$ , we obtain:

$$\begin{aligned} E(Y|X = x, Z = z) &= E(Y_0|X = x, Z = z) + E(D(Y_1 - Y_0)|X = x, Z = z) \\ &= E(Y_0|X = x, Z = z) + \int_0^{P(z)} E(Y_1 - Y_0|X = x, U_D = u_D) du_D. \end{aligned}$$



- Using  $Y = Y_0 + D(Y_1 - Y_0)$ , we obtain:

$$\begin{aligned} E(Y|X = x, Z = z) &= E(Y_0|X = x, Z = z) + E(D(Y_1 - Y_0)|X = x, Z = z) \\ &= E(Y_0|X = x, Z = z) + \int_0^{P(z)} E(Y_1 - Y_0|X = x, U_D = u_D) du_D. \end{aligned}$$

- So

$$\frac{\frac{\partial}{\partial z} E(Y|X = x, Z = z)}{\frac{\partial}{\partial z} P(z)} = \frac{\frac{\partial}{\partial z} E(Y_0|X = x, Z = z)}{\frac{\partial}{\partial z} P(z)} + E(Y_1 - Y_0|X = x, U_D = P(z))$$

where we have assumed that  $\frac{\partial}{\partial z} P(z) \neq 0$  and that  $E(Y_0|X = x, Z = z)$  is differentiable in  $z$ .

- Using  $Y = Y_0 + D(Y_1 - Y_0)$ , we obtain:

$$\begin{aligned} E(Y|X = x, Z = z) &= E(Y_0|X = x, Z = z) + E(D(Y_1 - Y_0)|X = x, Z = z) \\ &= E(Y_0|X = x, Z = z) + \int_0^{P(z)} E(Y_1 - Y_0|X = x, U_D = u_D) du_D. \end{aligned}$$

- So

$$\frac{\frac{\partial}{\partial z} E(Y|X = x, Z = z)}{\frac{\partial}{\partial z} P(z)} = \frac{\frac{\partial}{\partial z} E(Y_0|X = x, Z = z)}{\frac{\partial}{\partial z} P(z)} + E(Y_1 - Y_0|X = x, U_D = P(z))$$

where we have assumed that  $\frac{\partial}{\partial z} P(z) \neq 0$  and that  $E(Y_0|X = x, Z = z)$  is differentiable in  $z$ .

- Notice that under our stronger independence condition (A-1),  $\frac{\partial}{\partial z} E(Y_0|X = x, Z = z) = 0$  so that we identify MTE as before.

- With  $Y_0$  possibly dependent on  $Z$  conditional on  $X$ , we now get MTE plus the bias term that depends on  $\frac{\partial}{\partial z} E(Y_0 | X = x, Z = z)$ .

- With  $Y_0$  possibly dependent on  $Z$  conditional on  $X$ , we now get MTE plus the bias term that depends on  $\frac{\partial}{\partial z} E(Y_0|X = x, Z = z)$ .
- Likewise, if we consider the discrete change form of IV:

$$\begin{aligned}
 & \frac{E(Y|X = x, Z = z) - E(Y|X = x, Z = z')}{P(z) - P(z')} \\
 = & \underbrace{\frac{E(Y_0|X = x, Z = z) - E(Y_0|X = x, Z = z')}{P(z) - P(z')}}_{\text{Bias for LATE}} \\
 & + \underbrace{E(Y_1 - Y_0|X = x, P(z) > U_D > P(z'))}_{\text{LATE}}
 \end{aligned}$$

so that we now recover LATE plus a bias term.

- Now consider a regression discontinuity design.

- Now consider a regression discontinuity design.
- Suppose that there exists an evaluation point  $z_0$  for  $Z$  such that  $P(\cdot)$  is discontinuous at  $z_0$ , and suppose that  $E(Y_0|X = x, Z = z)$  is continuous at  $z_0$ .

- Now consider a regression discontinuity design.
- Suppose that there exists an evaluation point  $z_0$  for  $Z$  such that  $P(\cdot)$  is discontinuous at  $z_0$ , and suppose that  $E(Y_0|X = x, Z = z)$  is continuous at  $z_0$ .
- Suppose that  $P(\cdot)$  is increasing in a neighborhood of  $z_0$ .

- Now consider a regression discontinuity design.
- Suppose that there exists an evaluation point  $z_0$  for  $Z$  such that  $P(\cdot)$  is discontinuous at  $z_0$ , and suppose that  $E(Y_0|X = x, Z = z)$  is continuous at  $z_0$ .
- Suppose that  $P(\cdot)$  is increasing in a neighborhood of  $z_0$ .
- Let

$$P(z_0-) = \lim_{\epsilon \downarrow 0} P(z_0 - \epsilon),$$

$$P(z_0+) = \lim_{\epsilon \downarrow 0} P(z_0 + \epsilon),$$

and note that the conditions that  $P(\cdot)$  is increasing in a neighborhood of  $z_0$  and discontinuous at  $z_0$  imply that  $P(z_0+) > P(z_0-)$ .



- Let

$$\mu(x, z_0-) = \lim_{\epsilon \downarrow 0} E(Y|X = x, Z = z_0 - \epsilon)$$

$$\mu(x, z_0+) = \lim_{\epsilon \downarrow 0} E(Y|X = x, Z = z_0 + \epsilon),$$

and note that

$$\mu(x, z_0-) = E(Y_0|X = x, Z = z_0) + \int_0^{P(z_0-)} E(Y_1 - Y_0|X = x, U_D = u_D) du_D$$

and

$$\mu(x, z_0+) = E(Y_0|X = x, Z = z_0) + \int_0^{P(z_0+)} E(Y_1 - Y_0|X = x, U_D = u_D) du_D,$$

where we use the fact that  $E(Y_0|X = x, Z = z)$  is continuous at  $z_0$ .

- Thus,

$$\mu(x, z_0+) - \mu(x, z_0-) = \int_{P(z_0-)}^{P(z_0+)} E(Y_1 - Y_0 | X = x, U_D = u_D) du_D$$

$$\Rightarrow \frac{\mu(x, z_0+) - \mu(x, z_0-)}{P(z_0+) - P(z_0-)} = E(Y_1 - Y_0 | X = x, P(z_0+) \geq U_D > P(z_0-))$$

so that we now recover a LATE parameter for a particular point of evaluation.

- Note that if  $P(z)$  is only discontinuous at  $z_0$ , then we only identify  $E(Y_1 - Y_0 | X = x, P(z_0+) \geq U_D > P(z_0-))$  and not any LATE or MTE at any other evaluation points.

- Note that if  $P(z)$  is only discontinuous at  $z_0$ , then we only identify  $E(Y_1 - Y_0 | X = x, P(z_0+) \geq U_D > P(z_0-))$  and not any LATE or MTE at any other evaluation points.
- While this discussion assumes that  $Z$  is a scalar, it is straightforward to generalize the discussion to allow for  $Z$  to be a vector.

- Note that if  $P(z)$  is only discontinuous at  $z_0$ , then we only identify  $E(Y_1 - Y_0 | X = x, P(z_0+) \geq U_D > P(z_0-))$  and not any LATE or MTE at any other evaluation points.
- While this discussion assumes that  $Z$  is a scalar, it is straightforward to generalize the discussion to allow for  $Z$  to be a vector.
- For more discussion of the regression discontinuity design estimator and an example, see ?.

## Policy Evaluation, Out-of-Sample Policy Forecasting, Forecasting the Effects of New Policies and Structural Models Based on the MTE

- We have thus far focused on policy problem P-1, the problem of “internal validity”. We have shown how to identify a variety of parameters but have not put them to use in evaluating policies.

## Policy Evaluation, Out-of-Sample Policy Forecasting, Forecasting the Effects of New Policies and Structural Models Based on the MTE

- We have thus far focused on policy problem P-1, the problem of “internal validity”. We have shown how to identify a variety of parameters but have not put them to use in evaluating policies.
- This section discusses policy evaluation and out-of-sample forecasting.

## Policy Evaluation, Out-of-Sample Policy Forecasting, Forecasting the Effects of New Policies and Structural Models Based on the MTE

- We have thus far focused on policy problem P-1, the problem of “internal validity”. We have shown how to identify a variety of parameters but have not put them to use in evaluating policies.
- This section discusses policy evaluation and out-of-sample forecasting.
- We discuss two distinct evaluation and forecasting problems.



## Policy Evaluation, Out-of-Sample Policy Forecasting, Forecasting the Effects of New Policies and Structural Models Based on the MTE

- We have thus far focused on policy problem P-1, the problem of “internal validity”. We have shown how to identify a variety of parameters but have not put them to use in evaluating policies.
- This section discusses policy evaluation and out-of-sample forecasting.
- We discuss two distinct evaluation and forecasting problems.
- The first problem uses the MTE to develop a cost benefit analysis.

- Corresponding to the gross benefit parameters analyzed in Slides 90–152, there is a parallel set of cost parameters that emerge from the economics of the generalized Roy model.

- Corresponding to the gross benefit parameters analyzed in Slides 90–152, there is a parallel set of cost parameters that emerge from the economics of the generalized Roy model.
- This part of our analysis works in the domain of problem P-1 to construct a cost-benefit analysis for programs in place.

- Corresponding to the gross benefit parameters analyzed in Slides 90–152, there is a parallel set of cost parameters that emerge from the economics of the generalized Roy model.
- This part of our analysis works in the domain of problem P-1 to construct a cost-benefit analysis for programs in place.
- However, these tools can be extended to new environments using the other results established in this section.

- Corresponding to the gross benefit parameters analyzed in Slides 90–152, there is a parallel set of cost parameters that emerge from the economics of the generalized Roy model.
- This part of our analysis works in the domain of problem P-1 to construct a cost-benefit analysis for programs in place.
- However, these tools can be extended to new environments using the other results established in this section.
- The second topic is the problem of constructing the PRTE in new environments in a more general way.

- Corresponding to the gross benefit parameters analyzed in Slides 90–152, there is a parallel set of cost parameters that emerge from the economics of the generalized Roy model.
- This part of our analysis works in the domain of problem P-1 to construct a cost-benefit analysis for programs in place.
- However, these tools can be extended to new environments using the other results established in this section.
- The second topic is the problem of constructing the PRTE in new environments in a more general way.
- This addresses policy problems P-2 and P-3 and considers large scale changes in policies and forecasts of new policies.

## Econometric Cost Benefit Analysis Based on the MTE

- This section complements the analysis of Slide 90.

## Econometric Cost Benefit Analysis Based on the MTE

- This section complements the analysis of Slide 90.
- There we developed gross outcome measures for a generalized Roy model.



## Econometric Cost Benefit Analysis Based on the MTE

- This section complements the analysis of Slide 90.
- There we developed gross outcome measures for a generalized Roy model.
- Here we define a parallel set of treatment parameters for the generalized Roy model corresponding to the average cost of participating in a program.

## Econometric Cost Benefit Analysis Based on the MTE

- This section complements the analysis of Slide 90.
- There we developed gross outcome measures for a generalized Roy model.
- Here we define a parallel set of treatment parameters for the generalized Roy model corresponding to the average cost of participating in a program.
- The central feature of the generalized Roy model is that the agent chooses treatment if the benefit exceeds the subjective cost perceived by the agent.

- This creates a simple relationship between the cost and benefit parameters that can be exploited for identifying or bounding the cost parameters by adapting the results of the previous sections.

- This creates a simple relationship between the cost and benefit parameters that can be exploited for identifying or bounding the cost parameters by adapting the results of the previous sections.
- The main result of this section is that cost parameters in the generalized Roy model can be identified or bounded without direct information on the costs of treatment.

- This creates a simple relationship between the cost and benefit parameters that can be exploited for identifying or bounding the cost parameters by adapting the results of the previous sections.
- The main result of this section is that cost parameters in the generalized Roy model can be identified or bounded without direct information on the costs of treatment.
- Our analysis complements and extends the analysis of ? who first noted this duality.

- Assume the outcomes  $(Y_0, Y_1)$  are generated by the additively separable system (2).

- Assume the outcomes  $(Y_0, Y_1)$  are generated by the additively separable system (2).
- Let  $C$  denote the individual-specific subjective cost of selecting into treatment.

- Assume the outcomes  $(Y_0, Y_1)$  are generated by the additively separable system (2).
- Let  $C$  denote the individual-specific subjective cost of selecting into treatment.
- We assume that  $C$  is generated by  $C = \mu_C(W) + U_C$ , where  $W$  is a (possibly vector-valued) observed random variable and  $U_C$  is an unobserved random variable.



- Assume the outcomes  $(Y_0, Y_1)$  are generated by the additively separable system (2).
- Let  $C$  denote the individual-specific subjective cost of selecting into treatment.
- We assume that  $C$  is generated by  $C = \mu_C(W) + U_C$ , where  $W$  is a (possibly vector-valued) observed random variable and  $U_C$  is an unobserved random variable.
- We assume that the agent selects into treatment if the benefit exceeds the cost, using the structure of the generalized Roy model where  $D = \mathbf{1}[Y_1 - Y_0 \geq C]$  and  $C = \mu_C(W) + U_C$ , where  $\mu_C(W)$  is nondegenerate and integrable;  $U_C$  is continuous and  $Z = (W, X)$  is independent of  $(U_C, U_0, U_1)$ .

- We do not assume any particular functional form for the functions  $\mu_0, \mu_1$  and  $\mu_C$ , and we do not assume that the distribution of  $U_0, U_1$ , or  $U_C$  is known.

- We do not assume any particular functional form for the functions  $\mu_0, \mu_1$  and  $\mu_C$ , and we do not assume that the distribution of  $U_0, U_1$ , or  $U_C$  is known.
- Let  $V \equiv U_C - (U_1 - U_0)$  and let  $F_V$  denote the distribution function of  $V$ . As before, we use the convention that  $U_D$  is the probability integral transformation of the latent variable generating choices so that  $U_D = F_V(V)$ .

- We do not assume any particular functional form for the functions  $\mu_0, \mu_1$  and  $\mu_C$ , and we do not assume that the distribution of  $U_0, U_1$ , or  $U_C$  is known.
- Let  $V \equiv U_C - (U_1 - U_0)$  and let  $F_V$  denote the distribution function of  $V$ . As before, we use the convention that  $U_D$  is the probability integral transformation of the latent variable generating choices so that  $U_D = F_V(V)$ .
- Let  $P(z) \equiv \Pr(D = 1|Z = z)$  so that  $P(z) = F_V(\mu_1(x) - \mu_0(x) - \mu_C(w))$ .

- We do not assume any particular functional form for the functions  $\mu_0, \mu_1$  and  $\mu_C$ , and we do not assume that the distribution of  $U_0, U_1$ , or  $U_C$  is known.
- Let  $V \equiv U_C - (U_1 - U_0)$  and let  $F_V$  denote the distribution function of  $V$ . As before, we use the convention that  $U_D$  is the probability integral transformation of the latent variable generating choices so that  $U_D = F_V(V)$ .
- Let  $P(z) \equiv \Pr(D = 1|Z = z)$  so that  $P(z) = F_V(\mu_1(x) - \mu_0(x) - \mu_C(w))$ .
- For convenience, we will assume that  $F_V$  is strictly increasing so that  $F_V$  will be invertible, though this assumption is not required.

- We work with  $U_D = F_V(V)$  instead of working directly with  $V$  to link our analysis to that in Slide 90.

- We work with  $U_D = F_V(V)$  instead of working directly with  $V$  to link our analysis to that in Slide 90.
- In this section we make explicit the conditioning on  $X$ ,  $Z$ , and  $W$  because it plays an important role in the analysis.

- We work with  $U_D = F_V(V)$  instead of working directly with  $V$  to link our analysis to that in Slide 90.
- In this section we make explicit the conditioning on  $X$ ,  $Z$ , and  $W$  because it plays an important role in the analysis.
- Corresponding to the treatment parameters defined in Slide 12 and tables 2A and 2B, we can define analogous cost parameters.



- We work with  $U_D = F_V(V)$  instead of working directly with  $V$  to link our analysis to that in Slide 90.
- In this section we make explicit the conditioning on  $X$ ,  $Z$ , and  $W$  because it plays an important role in the analysis.
- Corresponding to the treatment parameters defined in Slide 12 and tables 2A and 2B, we can define analogous cost parameters.
- We define the marginal cost of treatment for a person with characteristics  $W = w$  and  $U_D = u_D$  as

$$C^{\text{MTE}}(w, u_D) \equiv E(C|W = w, U_D = u_D).$$

This is a cost version of the marginal treatment effect.

- Likewise, we have an analogue average cost:

$$\begin{aligned} C^{\text{ATE}}(w) &\equiv E(C|W = w) \\ &= \int_0^1 E(C|W = w, U_D = u_D) du_D, \end{aligned} \quad (34)$$

recalling that  $dF_{U_D}(u_D) = du_D$  because  $U_D$  is uniform.

- Likewise, we have an analogue average cost:

$$\begin{aligned} C^{\text{ATE}}(w) &\equiv E(C|W = w) \\ &= \int_0^1 E(C|W = w, U_D = u_D) du_D, \end{aligned} \quad (34)$$

recalling that  $dF_{U_D}(u_D) = du_D$  because  $U_D$  is uniform.

- This is the mean subjective cost of treatment as perceived by the average agent.

- Likewise, we have an analogue average cost:

$$\begin{aligned} C^{\text{ATE}}(w) &\equiv E(C|W = w) \\ &= \int_0^1 E(C|W = w, U_D = u_D) du_D, \end{aligned} \quad (34)$$

recalling that  $dF_{U_D}(u_D) = du_D$  because  $U_D$  is uniform.

- This is the mean subjective cost of treatment as perceived by the average agent.
- We next consider

$$\begin{aligned} C^{\text{TT}}(w, P(z)) &\equiv E(C | W = w, P(Z) = P(z), D = 1) \\ &= \frac{1}{P(z)} \int_0^{P(z)} E(C | W = w, U_D = u_D) du_D. \end{aligned}$$

- This is the mean subjective cost of treatment as perceived by the treated with a given value of  $P(z)$ .

- This is the mean subjective cost of treatment as perceived by the treated with a given value of  $P(z)$ .
- Removing the conditioning on  $P(z)$ ,

$$\begin{aligned} C^{\text{TT}}(w) &\equiv E(C|W = w, D = 1) \\ &= \int_0^1 E(C|W = w, U_D = u_D) g_w(u_D) du_D, \end{aligned}$$

where  $g_w(u_D) = \frac{1 - F_{P(Z)|W=w}(u_D)}{\int (1 - F_{P(Z)|W=w}(t)) dt}$  and  $F_{P(Z)|W=w}$  denotes the distribution of  $P(Z)$  conditional on  $W = w$ .

- This is the mean subjective cost of treatment as perceived by the treated with a given value of  $P(z)$ .
- Removing the conditioning on  $P(z)$ ,

$$\begin{aligned} C^{TT}(w) &\equiv E(C|W=w, D=1) \\ &= \int_0^1 E(C|W=w, U_D=u_D) g_w(u_D) du_D, \end{aligned}$$

where  $g_w(u_D) = \frac{1-F_{P(Z)|W=w}(u_D)}{\int (1-F_{P(Z)|W=w}(t)) dt}$  and  $F_{P(Z)|W=w}$  denotes the distribution of  $P(Z)$  conditional on  $W=w$ .

- This is the mean subjective cost of treatment for the treated.

- Finally, we can derive a LATE version of the cost:

$$C^{\text{LATE}}(w, P(z), P(z')) \equiv \frac{1}{P(z) - P(z')} \int_{P(z')}^{P(z)} E(C|W = w, U_D = u_D) du_D.$$

This is the mean subjective cost of switching states for those induced to switch status by a change in the instrument.



- The generalized Roy model makes a tight link between the cost of treatment and the benefit of treatment.

- The generalized Roy model makes a tight link between the cost of treatment and the benefit of treatment.
- Thus one might expect a relationship between the gross benefit and cost parameters.

- The generalized Roy model makes a tight link between the cost of treatment and the benefit of treatment.
- Thus one might expect a relationship between the gross benefit and cost parameters.
- We show that the benefit and cost parameters coincide for MTE.

- The generalized Roy model makes a tight link between the cost of treatment and the benefit of treatment.
- Thus one might expect a relationship between the gross benefit and cost parameters.
- We show that the benefit and cost parameters coincide for MTE.
- This relationship can be used to infer information on the subjective cost of treatment by the use of local instrumental variables.

- Define  $\Delta^{\text{LIV}}(x, P(z))$  as in equation (19):

$$\Delta^{\text{LIV}}(x, P(z)) \equiv \frac{\partial E(Y|X=x, P(Z)=P(z))}{\partial P(z)}.$$

Under assumptions (A-1)–(A-5), LIV identifies MTE:

$$\Delta^{\text{LIV}}(x, P(z)) = \Delta^{\text{MTE}}(x, P(z)).$$

Note that

$$\begin{aligned} \Delta^{\text{MTE}}(x, P(z)) &= E(\Delta \mid X=x, U_D = P(z)) \\ &= E(\Delta \mid X=x, \Delta(x) = C(w)) \\ &= E(\Delta(x) \mid \Delta(x) = C(w)), \end{aligned} \quad (35)$$

where  $\Delta(x) = \mu_1(x) - \mu_0(x) + U_1 - U_0$ , and  $C(w) = \mu_C(w) + U_C$ .

- ( $\Delta(x)$  and  $C(w)$  are, respectively, the benefit and cost for the agent if the  $X$  and  $W$  are externally set to  $x$  and  $w$  without changing ( $U_1, U_0, U_D$ ) values.) We thus obtain:

$$\begin{aligned}
 E(\Delta(x) \mid \Delta(x) = C(w)) &= E(C(w) \mid \Delta(x) = C(w)) \\
 &= E(C(w) \mid W = w, U_D = P(z)) \\
 &= C^{\text{MTE}}(w, P(z)).
 \end{aligned} \tag{36}$$

Thus,

$$\Delta^{\text{LIV}}(x, P(z)) = \Delta^{\text{MTE}}(x, P(z)) = C^{\text{MTE}}(w, P(z)). \tag{37}$$

where  $\Delta^{\text{LIV}}(w, P(z))$  is  $\Delta^{\text{LIV}}(x, P(z))$  defined for the support where  $\Delta(x) = C(w)$ .

- The benefit and cost parameters coincide for the MTE parameter because at the margin, the marginal cost should equal the marginal benefit.

- The benefit and cost parameters coincide for the MTE parameter because at the margin, the marginal cost should equal the marginal benefit.
- The benefit to treatment for an agent indifferent between treatment and no treatment is equal to the cost of treatment, and thus the two parameters coincide.



- The benefit and cost parameters coincide for the MTE parameter because at the margin, the marginal cost should equal the marginal benefit.
- The benefit to treatment for an agent indifferent between treatment and no treatment is equal to the cost of treatment, and thus the two parameters coincide.
- Suppose that one has access to a large sample of  $(Y, D, X, W)$  observations.

- The benefit and cost parameters coincide for the MTE parameter because at the margin, the marginal cost should equal the marginal benefit.
- The benefit to treatment for an agent indifferent between treatment and no treatment is equal to the cost of treatment, and thus the two parameters coincide.
- Suppose that one has access to a large sample of  $(Y, D, X, W)$  observations.
- Since  $\Delta^{\text{LIV}}(x, P(z)) = \frac{\partial E(Y|X=x, P(Z)=P(z))}{\partial P(z)}$ ,  $\Delta^{\text{LIV}}(x, P(z))$  can be identified for any  $(x, P(z))$  in the support of  $(X, P(Z))$ , and thus the corresponding  $\Delta^{\text{MTE}}(x, P(z))$  and  $C^{\text{MTE}}(w, P(z))$  parameters can also be identified.

- One can thus identify the marginal cost parameter without direct information on the cost of treatment by using the structure of the Roy model and by identifying the marginal benefit parameter.

- One can thus identify the marginal cost parameter without direct information on the cost of treatment by using the structure of the Roy model and by identifying the marginal benefit parameter.
- ? establish conditions under which  $\Delta^{\text{LIV}}$  can be used to identify  $\Delta^{\text{ATE}}$  and  $\Delta^{\text{TT}}$  given large support conditions, and to bound those parameters without large support conditions if the outcome variables are bounded.

- One can thus identify the marginal cost parameter without direct information on the cost of treatment by using the structure of the Roy model and by identifying the marginal benefit parameter.
- ? establish conditions under which  $\Delta^{\text{LIV}}$  can be used to identify  $\Delta^{\text{ATE}}$  and  $\Delta^{\text{TT}}$  given large support conditions, and to bound those parameters without large support conditions if the outcome variables are bounded.
- We review their results on bounds in Slide 938.

- One can thus identify the marginal cost parameter without direct information on the cost of treatment by using the structure of the Roy model and by identifying the marginal benefit parameter.
- ? establish conditions under which  $\Delta^{\text{LIV}}$  can be used to identify  $\Delta^{\text{ATE}}$  and  $\Delta^{\text{TT}}$  given large support conditions, and to bound those parameters without large support conditions if the outcome variables are bounded.
- We review their results on bounds in Slide 938.
- We surveyed their results on identification of  $\Delta^{\text{ATE}}$  and  $\Delta^{\text{TT}}$  in Slides 90 and 152.

- From (34) and (37), we can use the same arguments to use  $C^{MTE}$  to identify or bound  $C^{ATE}$  and  $C^{TT}$ .

- From (34) and (37), we can use the same arguments to use  $C^{MTE}$  to identify or bound  $C^{ATE}$  and  $C^{TT}$ .
- Thus,  $C^{MTE}$  can be used to identify  $C^{ATE}(w)$  if the support of  $P(Z)$  conditional on  $W = w$  is the full unit interval.



- From (34) and (37), we can use the same arguments to use  $C^{\text{MTE}}$  to identify or bound  $C^{\text{ATE}}$  and  $C^{\text{TT}}$ .
- Thus,  $C^{\text{MTE}}$  can be used to identify  $C^{\text{ATE}}(w)$  if the support of  $P(Z)$  conditional on  $W = w$  is the full unit interval.
- If the support of  $P(Z)$  conditional on  $W = w$  is a proper subset of the full unit interval, then  $C^{\text{MTE}}$  can be used to bound  $C^{\text{ATE}}(x)$  if  $C$  is bounded.

- From (34) and (37), we can use the same arguments to use  $C^{\text{MTE}}$  to identify or bound  $C^{\text{ATE}}$  and  $C^{\text{TT}}$ .
- Thus,  $C^{\text{MTE}}$  can be used to identify  $C^{\text{ATE}}(w)$  if the support of  $P(Z)$  conditional on  $W = w$  is the full unit interval.
- If the support of  $P(Z)$  conditional on  $W = w$  is a proper subset of the full unit interval, then  $C^{\text{MTE}}$  can be used to bound  $C^{\text{ATE}}(x)$  if  $C$  is bounded.
- One can thus identify or bound the average cost of treatment or the cost of treatment on the treated without direct information on the cost of treatment.

- We next consider what information is available on the underlying benefit functions  $\mu_0$  and  $\mu_1$  and the underlying cost function  $\mu_C(w)$ .

- We next consider what information is available on the underlying benefit functions  $\mu_0$  and  $\mu_1$  and the underlying cost function  $\mu_C(w)$ .
- From the definitions,

$$\begin{aligned}\Delta^{\text{MTE}}(x, P(z)) &= E(\Delta \mid X = x, U_D = P(z)) \\ &= \mu_1(x) - \mu_0(x) + \Upsilon(P(z))\end{aligned}\tag{38}$$

with  $\Upsilon(P(z)) = E(U_1 - U_0 \mid U_D = P(z))$ .

- We next consider what information is available on the underlying benefit functions  $\mu_0$  and  $\mu_1$  and the underlying cost function  $\mu_C(w)$ .
- From the definitions,

$$\begin{aligned}\Delta^{\text{MTE}}(x, P(z)) &= E(\Delta \mid X = x, U_D = P(z)) \\ &= \mu_1(x) - \mu_0(x) + \Upsilon(P(z))\end{aligned}\quad (38)$$

with  $\Upsilon(P(z)) = E(U_1 - U_0 \mid U_D = P(z))$ .

- Likewise,

$$\begin{aligned}C^{\text{MTE}}(w, P(z)) &= E(C \mid W = w, U_D = P(z)) \\ &= \mu_C(w) + \Gamma(P(z)),\end{aligned}\quad (39)$$

with  $\Gamma(P(z)) = E(U_C \mid U_D = P(z))$ .

- Let  $\Delta^{\text{LIV}}(z) = \Delta^{\text{LIV}}(x, P(z))$ , and recall from the preceding analysis that  $\Delta^{\text{LIV}}(z) = \Delta^{\text{MTE}}(x, P(z)) = C^{\text{MTE}}(w, P(z))$ .

- Let  $\Delta^{\text{LIV}}(z) = \Delta^{\text{LIV}}(x, P(z))$ , and recall from the preceding analysis that  $\Delta^{\text{LIV}}(z) = \Delta^{\text{MTE}}(x, P(z)) = C^{\text{MTE}}(w, P(z))$ .
- Consider two points of evaluation  $(z, z')$  such that  $P(z) = P(z')$ .

- Let  $\Delta^{\text{LIV}}(z) = \Delta^{\text{LIV}}(x, P(z))$ , and recall from the preceding analysis that  $\Delta^{\text{LIV}}(z) = \Delta^{\text{MTE}}(x, P(z)) = C^{\text{MTE}}(w, P(z))$ .
- Consider two points of evaluation  $(z, z')$  such that  $P(z) = P(z')$ .
- Using equation (37), we obtain:

$$\Delta^{\text{LIV}}(z) - \Delta^{\text{LIV}}(z') = (\mu_1(x) - \mu_0(x)) - (\mu_1(x') - \mu_0(x')) = \mu_C(w) - \mu_C(w').$$



- Let  $\Delta^{\text{LIV}}(z) = \Delta^{\text{LIV}}(x, P(z))$ , and recall from the preceding analysis that  $\Delta^{\text{LIV}}(z) = \Delta^{\text{MTE}}(x, P(z)) = C^{\text{MTE}}(w, P(z))$ .
- Consider two points of evaluation  $(z, z')$  such that  $P(z) = P(z')$ .
- Using equation (37), we obtain:

$$\Delta^{\text{LIV}}(z) - \Delta^{\text{LIV}}(z') = (\mu_1(x) - \mu_0(x)) - (\mu_1(x') - \mu_0(x')) = \mu_C(w) - \mu_C(w').$$

- Assuming that  $X$  and  $W$  each have at least one component not in the other, we can identify  $\mu_C(w)$  up to constants within the support of  $W$  conditional on  $P(Z) = P(z)$  using  $\Delta^{\text{LIV}}(z)$ .

- Shifting  $z$  while conditioning on  $P(z)$  shifts  $(\mu_1(x) - \mu_0(x))$  and  $\mu_C(w)$  along the line  $(\mu_1(x) - \mu_0(x)) - \mu_C(w) = F_V^{-1}(p)$ . Thus, conditional on  $P(z)$ , a shift in the benefit,  $\mu_1(X) - \mu_0(X)$ , is associated with the same shift in the cost,  $\mu_C(w)$ .

- Shifting  $z$  while conditioning on  $P(z)$  shifts  $(\mu_1(x) - \mu_0(x))$  and  $\mu_C(w)$  along the line  $(\mu_1(x) - \mu_0(x)) - \mu_C(w) = F_V^{-1}(p)$ . Thus, conditional on  $P(z)$ , a shift in the benefit,  $\mu_1(X) - \mu_0(X)$ , is associated with the same shift in the cost,  $\mu_C(w)$ .
- For any  $p \in (0, 1)$ , let  $\Omega_p = \{z : P(z) = p\} = \{(w, x) : (\mu_1(x) - \mu_0(x)) - \mu_C(w) = F_V^{-1}(p)\}$ .

- Shifting  $z$  while conditioning on  $P(z)$  shifts  $(\mu_1(x) - \mu_0(x))$  and  $\mu_C(w)$  along the line  $(\mu_1(x) - \mu_0(x)) - \mu_C(w) = F_V^{-1}(p)$ . Thus, conditional on  $P(z)$ , a shift in the benefit,  $\mu_1(X) - \mu_0(X)$ , is associated with the same shift in the cost,  $\mu_C(w)$ .
- For any  $p \in (0, 1)$ , let  $\Omega_p = \{z : P(z) = p\} = \{(w, x) : (\mu_1(x) - \mu_0(x)) - \mu_C(w) = F_V^{-1}(p)\}$ .
- As we vary  $z$  within the set  $\Omega_p$ , we trace out changes in  $\mu_C(w)$  and  $\mu_1(x) - \mu_0(x)$ , where the changes in  $\mu_C(w)$  equal the changes in  $\mu_1(x) - \mu_0(x)$ .

- For the special case of the generalized Roy model where  $U_C$  is degenerate,  $\Delta^{\text{LIV}}(z) = \mu_C(w)$ .

- For the special case of the generalized Roy model where  $U_C$  is degenerate,  $\Delta^{\text{LIV}}(z) = \mu_C(w)$ .
- Thus, in the case of a deterministic cost function, LIV identifies  $\mu_C(w)$ .

- For the special case of the generalized Roy model where  $U_C$  is degenerate,  $\Delta^{\text{LIV}}(z) = \mu_C(w)$ .
- Thus, in the case of a deterministic cost function, LIV identifies  $\mu_C(w)$ .
- We plot this case in figures 5A–5C for the country policy adoption example where the cost  $C$  is a constant across all countries.

- For the special case of the generalized Roy model where  $U_C$  is degenerate,  $\Delta^{\text{LIV}}(z) = \mu_C(w)$ .
- Thus, in the case of a deterministic cost function, LIV identifies  $\mu_C(w)$ .
- We plot this case in figures 5A–5C for the country policy adoption example where the cost  $C$  is a constant across all countries.
- In the case where  $U_C$  is nondegenerate but  $U_1 - U_0$  is degenerate,  $Y_1 - Y_0 = \mu_1(X) - \mu_0(X)$ , ( $\beta = \bar{\beta}$  in the context of the model of Slide 12) and there is no variation in the gross benefit from participating in the program conditional on  $X$ .



- In that case,  $\Delta^{\text{LIV}}(z) = \mu_1(x) - \mu_0(x) = \bar{\beta}$ , where we keep the conditioning on  $X$  implicit in defining  $\Delta^{\text{LIV}}(z)$ .

- In that case,  $\Delta^{\text{LIV}}(z) = \mu_1(x) - \mu_0(x) = \bar{\beta}$ , where we keep the conditioning on  $X$  implicit in defining  $\Delta^{\text{LIV}}(z)$ .
- Thus, in the case of a deterministic benefit from participation, LIV identifies the benefit function.

- In that case,  $\Delta^{\text{LIV}}(z) = \mu_1(x) - \mu_0(x) = \bar{\beta}$ , where we keep the conditioning on  $X$  implicit in defining  $\Delta^{\text{LIV}}(z)$ .
- Thus, in the case of a deterministic benefit from participation, LIV identifies the benefit function.
- If  $U_D$  and  $U_1 - U_0$  are both degenerate, then  $\Delta^{\text{LIV}}(z)$  is not well defined.

- In that case,  $\Delta^{\text{LIV}}(z) = \mu_1(x) - \mu_0(x) = \bar{\beta}$ , where we keep the conditioning on  $X$  implicit in defining  $\Delta^{\text{LIV}}(z)$ .
- Thus, in the case of a deterministic benefit from participation, LIV identifies the benefit function.
- If  $U_D$  and  $U_1 - U_0$  are both degenerate, then  $\Delta^{\text{LIV}}(z)$  is not well defined.
- In summary, the generalized Roy model structure can be exploited to identify cost parameters without direct information on the cost of treatment.

- The MTE parameter for cost is immediately identified within the proper support, and can be used to identify or bound the average cost of treatment and the cost of treatment on the treated.

- The MTE parameter for cost is immediately identified within the proper support, and can be used to identify or bound the average cost of treatment and the cost of treatment on the treated.
- In addition, the MTE parameter allows one to infer how the cost function shifts in response to a change in observed covariates, and to completely identify the cost function if the cost of treatment is deterministic conditional on observable covariates.

- The MTE parameter for cost is immediately identified within the proper support, and can be used to identify or bound the average cost of treatment and the cost of treatment on the treated.
- In addition, the MTE parameter allows one to infer how the cost function shifts in response to a change in observed covariates, and to completely identify the cost function if the cost of treatment is deterministic conditional on observable covariates.
- Thus we can compute the costs and benefits of alternative programs for various population averages.

- The MTE parameter for cost is immediately identified within the proper support, and can be used to identify or bound the average cost of treatment and the cost of treatment on the treated.
- In addition, the MTE parameter allows one to infer how the cost function shifts in response to a change in observed covariates, and to completely identify the cost function if the cost of treatment is deterministic conditional on observable covariates.
- Thus we can compute the costs and benefits of alternative programs for various population averages.
- ? develop this analysis to consider marginal extensions of the policy relevant treatment effect (PRTE).



## Constructing the PRTE in New Environments

- In this section, we present conditions for constructing PRTE for new environments and for new programs using historical data for general changes in policies and environments.

## Constructing the PRTE in New Environments

- In this section, we present conditions for constructing PRTE for new environments and for new programs using historical data for general changes in policies and environments.
- We consider general changes in the environment and policies and not just the marginal perturbations of the  $P(Z)$  considered in the previous section.

## Constructing the PRTE in New Environments

- In this section, we present conditions for constructing PRTE for new environments and for new programs using historical data for general changes in policies and environments.
- We consider general changes in the environment and policies and not just the marginal perturbations of the  $P(Z)$  considered in the previous section.
- We address policy problems (P-2), forecasting the effects of existing policies to new environments and (P-3), forecasting the effects of new policies, never previously implemented.

- Let  $p \in \mathcal{P}$  denote a policy characterized by random vector  $Z_p$ . The usage of “ $p$ ” in this section is to be distinguished from a realized value of  $P(Z)$  as in most other sections in this chapter.

- Let  $p \in \mathcal{P}$  denote a policy characterized by random vector  $Z_p$ . The usage of “ $p$ ” in this section is to be distinguished from a realized value of  $P(Z)$  as in most other sections in this chapter.
- Let  $e \in \mathcal{E}$  denote an environment characterized by random vector  $X_e$ .

- Let  $p \in \mathcal{P}$  denote a policy characterized by random vector  $Z_p$ . The usage of “ $p$ ” in this section is to be distinguished from a realized value of  $P(Z)$  as in most other sections in this chapter.
- Let  $e \in \mathcal{E}$  denote an environment characterized by random vector  $X_e$ .
- A history,  $\mathcal{H}$ , is a collection of policy-environment  $(p, e)$  pairs that have been experienced and documented.

- Let  $p \in \mathcal{P}$  denote a policy characterized by random vector  $Z_p$ . The usage of “ $p$ ” in this section is to be distinguished from a realized value of  $P(Z)$  as in most other sections in this chapter.
- Let  $e \in \mathcal{E}$  denote an environment characterized by random vector  $X_e$ .
- A history,  $\mathcal{H}$ , is a collection of policy-environment  $(p, e)$  pairs that have been experienced and documented.
- We assume that the environment is autonomous so the choice of  $p$  does not affect  $X_e$ .

- Let  $p \in \mathcal{P}$  denote a policy characterized by random vector  $Z_p$ . The usage of “ $p$ ” in this section is to be distinguished from a realized value of  $P(Z)$  as in most other sections in this chapter.
- Let  $e \in \mathcal{E}$  denote an environment characterized by random vector  $X_e$ .
- A history,  $\mathcal{H}$ , is a collection of policy-environment  $(p, e)$  pairs that have been experienced and documented.
- We assume that the environment is autonomous so the choice of  $p$  does not affect  $X_e$ .
- Letting  $X_{e,p}$  denote the value of  $X_e$  under policy  $p$ , autonomy requires that



(A-8)

$$X_{e,p} = X_e, \quad \forall p, e \quad (\mathbf{Autonomy}).$$

- Autonomy is a more general notion than the no-feedback assumption introduced in (A-6).

- Autonomy is a more general notion than the no-feedback assumption introduced in (A-6).
- They are the same when the policy is a treatment.

- Autonomy is a more general notion than the no-feedback assumption introduced in (A-6).
- They are the same when the policy is a treatment.
- General equilibrium feedback effects can cause a failure of autonomy.

- Autonomy is a more general notion than the no-feedback assumption introduced in (A-6).
- They are the same when the policy is a treatment.
- General equilibrium feedback effects can cause a failure of autonomy.
- In this section, we will assume autonomy, in accordance with the partial equilibrium tradition in the treatment effect literature.

- Autonomy is a more general notion than the no-feedback assumption introduced in (A-6).
- They are the same when the policy is a treatment.
- General equilibrium feedback effects can cause a failure of autonomy.
- In this section, we will assume autonomy, in accordance with the partial equilibrium tradition in the treatment effect literature.
- Autonomy is a version of Hurwicz's policy invariance postulate but for a random variable and not a function.

- Evaluating a particular policy  $p'$  in environment  $e'$  is straightforward if  $(p', e') \in \mathcal{H}$ .

- Evaluating a particular policy  $p'$  in environment  $e'$  is straightforward if  $(p', e') \in \mathcal{H}$ .
- One simply looks at the associated outcomes and treatment effects formed in that policy environment and applies the methods previously discussed to obtain internally valid estimates.



- Evaluating a particular policy  $p'$  in environment  $e'$  is straightforward if  $(p', e') \in \mathcal{H}$ .
- One simply looks at the associated outcomes and treatment effects formed in that policy environment and applies the methods previously discussed to obtain internally valid estimates.
- The challenge comes in forecasting the impacts of policies ( $p'$ ) in environments ( $e'$ ) for  $(p', e')$  not in  $\mathcal{H}$ .

- We show how  $\Delta^{\text{MTE}}$  plays the role of a policy-invariant functional that aids in creating counterfactual states never previously experienced.

- We show how  $\Delta^{\text{MTE}}$  plays the role of a policy-invariant functional that aids in creating counterfactual states never previously experienced.
- We focus on the problem of constructing the policy relevant treatment effect  $\Delta^{\text{PRTE}}$  but our discussion applies more generally to the other treatment parameters.

- Given the assumptions invoked in Slide 90,  $\Delta^{\text{MTE}}$  can be used to evaluate a whole menu of policies characterized by different conditional distributions of  $P_{p'}$ .

- Given the assumptions invoked in Slide 90,  $\Delta^{\text{MTE}}$  can be used to evaluate a whole menu of policies characterized by different conditional distributions of  $P_{p'}$ .
- In addition, given our assumptions, we can focus on how policy  $p'$ , which is characterized by  $Z_{p'}$ , produces the distribution  $F_{P_{p'}|X}$  which weights an invariant  $\Delta^{\text{MTE}}$  without having to conduct a new investigation of  $(Y, X, Z)$  relationships for each proposed policy.

## Constructing Weights for New Policies in a Common Environment

- The problem of constructing  $\Delta^{\text{PRTE}}$  for policy  $p'$  (compared to baseline policy  $\bar{p}$ ) in environment  $e$  when  $(p', e) \notin \mathcal{H}$  entails constructing  $E(\Upsilon(Y_{p'}))$ .

## Constructing Weights for New Policies in a Common Environment

- The problem of constructing  $\Delta^{\text{PRTE}}$  for policy  $p'$  (compared to baseline policy  $\bar{p}$ ) in environment  $e$  when  $(p', e) \notin \mathcal{H}$  entails constructing  $E(\Upsilon(Y_{p'}))$ .
- We maintain the assumption that the baseline policy is observed, so  $(\bar{p}, e) \in \mathcal{H}$ .

## Constructing Weights for New Policies in a Common Environment

- The problem of constructing  $\Delta^{\text{PRTE}}$  for policy  $p'$  (compared to baseline policy  $\bar{p}$ ) in environment  $e$  when  $(p', e) \notin \mathcal{H}$  entails constructing  $E(\Upsilon(Y_{p'}))$ .
- We maintain the assumption that the baseline policy is observed, so  $(\bar{p}, e) \in \mathcal{H}$ .
- We also postulate instrumental variable assumptions (A-1)–(A-5), presented in Slide 90, and the policy invariance assumption (A-7), presented in Slide 139 and embedded in assumption (A-8).



- We use separable choice equation (7) to characterize choices.

- We use separable choice equation (7) to characterize choices.
- The policy is assumed not to change the distribution of  $(Y_0, Y_1, U_D)$  conditional on  $X$ .

- We use separable choice equation (7) to characterize choices.
- The policy is assumed not to change the distribution of  $(Y_0, Y_1, U_D)$  conditional on  $X$ .
- Under these conditions, Equation (10) is a valid expression for PRTE and constructing PRTE only requires identification of  $\Delta^{\text{MTE}}$  and constructing  $F_{P_{p'}|X_e}$  from the policy histories  $\mathcal{H}_e$ , defined as the elements of  $\mathcal{H}$  for a particular environment  $e$ ,  $\mathcal{H}_e = \{p : (p, e) \in \mathcal{H}\}$ .

- Associated with the policy histories  $p \in \mathcal{H}_e$  is a collection of policy variables  $\{Z_p : p \in \mathcal{H}_e\}$ .

- Associated with the policy histories  $p \in \mathcal{H}_e$  is a collection of policy variables  $\{Z_p : p \in \mathcal{H}_e\}$ .
- Suppose that a new policy  $p'$  can be written as  $Z_{p'} = T_{p',j}(Z_j)$  for some  $j \in \mathcal{H}_e$ , where  $T_{p',j}$  is a known deterministic transformation and  $Z_{p'}$  has the same list of variables as  $Z_j$ .

- Associated with the policy histories  $p \in \mathcal{H}_e$  is a collection of policy variables  $\{Z_p : p \in \mathcal{H}_e\}$ .
- Suppose that a new policy  $p'$  can be written as  $Z_{p'} = T_{p',j}(Z_j)$  for some  $j \in \mathcal{H}_e$ , where  $T_{p',j}$  is a known deterministic transformation and  $Z_{p'}$  has the same list of variables as  $Z_j$ .
- Examples of policies that can be characterized in this way are tax and subsidy policies on wages, prices and incomes that affect unit costs (wages or prices) and transfers.

- Associated with the policy histories  $p \in \mathcal{H}_e$  is a collection of policy variables  $\{Z_p : p \in \mathcal{H}_e\}$ .
- Suppose that a new policy  $p'$  can be written as  $Z_{p'} = T_{p',j}(Z_j)$  for some  $j \in \mathcal{H}_e$ , where  $T_{p',j}$  is a known deterministic transformation and  $Z_{p'}$  has the same list of variables as  $Z_j$ .
- Examples of policies that can be characterized in this way are tax and subsidy policies on wages, prices and incomes that affect unit costs (wages or prices) and transfers.
- Tuition might be shifted upward for everyone by the same amount, or tuition might be shifted according to a nonlinear function of current tuition, parents' income, and other observable characteristics in  $Z_j$ .

- Constructing  $F_{P_{p'}|X_e}$  from data in the policy history entails two distinct steps.



- Constructing  $F_{P_{p'}|X_e}$  from data in the policy history entails two distinct steps.
- From the definitions,

$$\Pr(P_{p'} \leq t \mid X_e) = \Pr(\{Z_{p'} : \Pr(D_{p'} = 1 \mid Z_{p'}, X_e) \leq t\} \mid X_e).$$

- Constructing  $F_{P_{p'}|X_e}$  from data in the policy history entails two distinct steps.
- From the definitions,

$$\Pr(P_{p'} \leq t \mid X_e) = \Pr(\{Z_{p'} : \Pr(D_{p'} = 1 \mid Z_{p'}, X_e) \leq t\} \mid X_e).$$

- If (i) we know the distribution of  $Z_{p'}$ , and (ii) we know the function  $\Pr(D_{p'} = 1 \mid Z_{p'} = z, X_e = x)$  over the appropriate support, we can then recover the distribution of  $P_{p'}$  conditional on  $X_e$ .

- Given that  $Z_{p'} = T_{p',j}(Z_j)$  for a known function  $T_{p',j}(\cdot)$ , step (i) is straightforward since we recover the distribution of  $Z_{p'}$  from the distribution of  $Z_j$  by using the fact that 
$$\Pr(Z_{p'} \leq t \mid X_e) = \Pr(\{Z_j : T_{p',j}(Z_j) \leq t\} \mid X_e).$$

- Given that  $Z_{p'} = T_{p',j}(Z_j)$  for a known function  $T_{p',j}(\cdot)$ , step (i) is straightforward since we recover the distribution of  $Z_{p'}$  from the distribution of  $Z_j$  by using the fact that 
$$\Pr(Z_{p'} \leq t \mid X_e) = \Pr(\{Z_j : T_{p',j}(Z_j) \leq t\} \mid X_e).$$
- Alternatively, part of the specification of the policy  $p'$  might be the distribution  $\Pr(Z_{p'} \leq t \mid X_e)$ .

- Given that  $Z_{p'} = T_{p',j}(Z_j)$  for a known function  $T_{p',j}(\cdot)$ , step (i) is straightforward since we recover the distribution of  $Z_{p'}$  from the distribution of  $Z_j$  by using the fact that 
$$\Pr(Z_{p'} \leq t \mid X_e) = \Pr(\{Z_j : T_{p',j}(Z_j) \leq t\} \mid X_e).$$
- Alternatively, part of the specification of the policy  $p'$  might be the distribution  $\Pr(Z_{p'} \leq t \mid X_e)$ .
- We now turn to the second step, recovering the function  $\Pr(D_{p'} = 1 \mid Z_{p'} = z, X_e = x)$  over the appropriate support.

- If  $Z_{p'}$  and  $Z_j$  contain the same elements though possibly with different distributions, then a natural approach to forecasting the new policy is to postulate that

$$P_j(z) = \Pr(D_j = 1 \mid Z_j = z, X_e) \quad (40)$$

$$= \Pr(D_{p'} = 1 \mid Z_{p'} = z, X_e) = P_{p'}(z), \quad (41)$$

i.e., that over a common support for  $Z_j$  and  $Z_{p'}$  the known conditional probability function and the desired conditional probability function agree.

- If  $Z_{p'}$  and  $Z_j$  contain the same elements though possibly with different distributions, then a natural approach to forecasting the new policy is to postulate that

$$P_j(z) = \Pr(D_j = 1 \mid Z_j = z, X_e) \quad (40)$$

$$= \Pr(D_{p'} = 1 \mid Z_{p'} = z, X_e) = P_{p'}(z), \quad (41)$$

i.e., that over a common support for  $Z_j$  and  $Z_{p'}$  the known conditional probability function and the desired conditional probability function agree.

- Condition (40) will hold, for example, if
 
$$D_j = \mathbf{1}[\mu_D(Z_j) - V \geq 0], \quad D_{p'} = \mathbf{1}[\mu_D(Z_{p'}) - V \geq 0],$$

$$Z_j \perp\!\!\!\perp V \mid X_e, \text{ and } Z_{p'} \perp\!\!\!\perp U_D \mid X_e, \text{ recalling that } U_D = F_{V|X}(V).$$

- Even if condition (40) is satisfied on a common support, the support of  $Z_j$  and  $Z_{j'}$  may not be the same.



- Even if condition (40) is satisfied on a common support, the support of  $Z_j$  and  $Z_{j'}$  may not be the same.
- If the support of the distribution of  $Z_{j'}$  is not contained in the support of the distribution of  $Z_j$ , then some form of extrapolation is needed.

- Even if condition (40) is satisfied on a common support, the support of  $Z_j$  and  $Z_{j'}$  may not be the same.
- If the support of the distribution of  $Z_{j'}$  is not contained in the support of the distribution of  $Z_j$ , then some form of extrapolation is needed.
- Alternatively, if we strengthen our assumptions so that (40) holds for all  $j \in \mathcal{H}_e$ , we can identify  $P_{j'}(z)$  for all  $z$  in  $\bigcup_{j \in \mathcal{H}_e} \text{Supp}(Z_j)$ .

- Even if condition (40) is satisfied on a common support, the support of  $Z_j$  and  $Z_{p'}$  may not be the same.
- If the support of the distribution of  $Z_{p'}$  is not contained in the support of the distribution of  $Z_j$ , then some form of extrapolation is needed.
- Alternatively, if we strengthen our assumptions so that (40) holds for all  $j \in \mathcal{H}_e$ , we can identify  $P_{p'}(z)$  for all  $z$  in  $\bigcup_{j \in \mathcal{H}_e} \text{Supp}(Z_j)$ .
- However, there is no guarantee that the support of the distribution of  $Z_{p'}$  will be contained in  $\bigcup_{j \in \mathcal{H}_e} \text{Supp}(Z_j)$ , in which case some form of extrapolation is needed.

- If extrapolation is required, one approach is to assume a parametric functional form for  $P_j(\cdot)$ .

- If extrapolation is required, one approach is to assume a parametric functional form for  $P_j(\cdot)$ .
- Given a parametric functional form, one can use the joint distribution of  $(D_j, Z_j)$  to identify the unknown parameters of  $P_j(\cdot)$  and then extrapolate the parametric functional form to evaluate  $P_j(\cdot)$  for all evaluation points in the support of  $Z_{j'}$ .

- If extrapolation is required, one approach is to assume a parametric functional form for  $P_j(\cdot)$ .
- Given a parametric functional form, one can use the joint distribution of  $(D_j, Z_j)$  to identify the unknown parameters of  $P_j(\cdot)$  and then extrapolate the parametric functional form to evaluate  $P_j(\cdot)$  for all evaluation points in the support of  $Z_{p'}$ .
- Alternatively, if there is overlap between the support of  $Z_{p'}$  and  $Z_j$ , so there is some overlap in the historical and policy  $p'$  supports of  $Z$ , we may use nonparametric methods presented in ? and extended by her in ?, based on functional restrictions (e.g., homogeneity) to construct the desired probabilities on new supports or to bound them.

- If extrapolation is required, one approach is to assume a parametric functional form for  $P_j(\cdot)$ .
- Given a parametric functional form, one can use the joint distribution of  $(D_j, Z_j)$  to identify the unknown parameters of  $P_j(\cdot)$  and then extrapolate the parametric functional form to evaluate  $P_j(\cdot)$  for all evaluation points in the support of  $Z_{p'}$ .
- Alternatively, if there is overlap between the support of  $Z_{p'}$  and  $Z_j$ , so there is some overlap in the historical and policy  $p'$  supports of  $Z$ , we may use nonparametric methods presented in ? and extended by her in ?, based on functional restrictions (e.g., homogeneity) to construct the desired probabilities on new supports or to bound them.
- Under the appropriate conditions, we may use analytic continuation to extend  $\Pr(D_j = 1|Z_j = z, X_e = x)$  to a new support for each  $X_e = x$  (?).

- The approach just presented is based on the assumption stated in equation (40).



- The approach just presented is based on the assumption stated in equation (40).
- That assumption is quite natural when  $Z_{p'}$  and  $Z_j$  both contain the same elements, say they both contain tuition and parent's income.

- The approach just presented is based on the assumption stated in equation (40).
- That assumption is quite natural when  $Z_{p'}$  and  $Z_j$  both contain the same elements, say they both contain tuition and parent's income.
- However, in some cases  $Z_{p'}$  might contain additional elements not contained in  $Z_j$ .

- The approach just presented is based on the assumption stated in equation (40).
- That assumption is quite natural when  $Z_{p'}$  and  $Z_j$  both contain the same elements, say they both contain tuition and parent's income.
- However, in some cases  $Z_{p'}$  might contain additional elements not contained in  $Z_j$ .
- As an example,  $Z_{p'}$  might include new user fees while  $Z_j$  consists of taxes and subsidies but does not include user fees.

- The approach just presented is based on the assumption stated in equation (40).
- That assumption is quite natural when  $Z_{p'}$  and  $Z_j$  both contain the same elements, say they both contain tuition and parent's income.
- However, in some cases  $Z_{p'}$  might contain additional elements not contained in  $Z_j$ .
- As an example,  $Z_{p'}$  might include new user fees while  $Z_j$  consists of taxes and subsidies but does not include user fees.
- In this case, the assumption stated in equation (40) is not expected to hold and is not even well defined if  $Z_{p'}$  and  $Z_j$  contain a different number of elements.

- A more basic approach analyzes a class of policies that operate on constraints, prices and endowments arrayed in vector  $Q$ .

- A more basic approach analyzes a class of policies that operate on constraints, prices and endowments arrayed in vector  $Q$ .
- Given the preferences and technology of the agent, a given  $Q = q$ , however arrived at, generates the same choices for the agent.

- A more basic approach analyzes a class of policies that operate on constraints, prices and endowments arrayed in vector  $Q$ .
- Given the preferences and technology of the agent, a given  $Q = q$ , however arrived at, generates the same choices for the agent.
- Thus a wage tax offset by a wage subsidy of the same amount produces a wage that has the same effect on choices as a no-policy wage.

- A more basic approach analyzes a class of policies that operate on constraints, prices and endowments arrayed in vector  $Q$ .
- Given the preferences and technology of the agent, a given  $Q = q$ , however arrived at, generates the same choices for the agent.
- Thus a wage tax offset by a wage subsidy of the same amount produces a wage that has the same effect on choices as a no-policy wage.
- Policy  $j$  affects  $Q$  (e.g., it affects prices paid, endowments and constraints).



- Define a map  $\Phi_j : Z_j \rightarrow Q_j$  which maps a policy  $j$ , described by  $Z_j$ , into its consequences  $(Q_j)$  for the baseline, fixed-dimensional vector  $Q$ .

- Define a map  $\Phi_j : Z_j \rightarrow Q_j$  which maps a policy  $j$ , described by  $Z_j$ , into its consequences ( $Q_j$ ) for the baseline, fixed-dimensional vector  $Q$ .
- A new policy  $p'$ , characterized by  $Z_{p'}$ , produces  $Q_{p'}$  that is possibly different from  $Q_j$  for all previous policies  $j \in \mathcal{H}_e$ .

- Define a map  $\Phi_j : Z_j \rightarrow Q_j$  which maps a policy  $j$ , described by  $Z_j$ , into its consequences ( $Q_j$ ) for the baseline, fixed-dimensional vector  $Q$ .
- A new policy  $p'$ , characterized by  $Z_{p'}$ , produces  $Q_{p'}$  that is possibly different from  $Q_j$  for all previous policies  $j \in \mathcal{H}_e$ .
- To construct the random variable  $P_{p'} = \Pr(D_{p'} = 1 \mid Z_{p'}, X_e)$ , we postulate that

$$\begin{aligned} \Pr(D_j = 1 \mid Z_j \in \Phi_j^{-1}(q), X_e = x) &= \Pr(D_j = 1 \mid Q_j = q, X_e = x) \\ &= \Pr(D_{p'} = 1 \mid Q_{p'} = q, X_e = x) \\ &= \Pr(D_{p'} = 1 \mid Z_{p'} \in \Phi_{p'}^{-1}(q), X_e = x), \end{aligned}$$

where  $\Phi_j^{-1}(q) = \{z : \Phi_j(z) = q\}$  and  $\Phi_{p'}^{-1}(q) = \{z : \Phi_{p'}(z) = q\}$ .

- Given these assumptions, our ability to recover  $\Pr(D_{p'} = 1 \mid Z_{p'} = z, X_e = x)$  for all  $(z, x)$  in the support of  $(Z_{p'}, X_e)$  depends on what  $\Phi_j$  functions have been historically observed and the richness of the histories of  $Q_j, j \in \mathcal{H}_e$ . For each  $z_{p'}$  evaluation point in the support of the distribution of  $Z_{p'}$ , there is a corresponding  $q = \Phi_{p'}(z_{p'})$  evaluation point in the support of the distribution of  $Q_j = \Phi_j(Z_j)$ .

- Given these assumptions, our ability to recover  $\Pr(D_{p'} = 1 \mid Z_{p'} = z, X_e = x)$  for all  $(z, x)$  in the support of  $(Z_{p'}, X_e)$  depends on what  $\Phi_j$  functions have been historically observed and the richness of the histories of  $Q_j, j \in \mathcal{H}_e$ . For each  $z_{p'}$  evaluation point in the support of the distribution of  $Z_{p'}$ , there is a corresponding  $q = \Phi_{p'}(z_{p'})$  evaluation point in the support of the distribution of  $Q_j = \Phi_j(Z_j)$ .
- If, in the policy histories, there is at least one  $j \in \mathcal{H}_e$  such that  $\Phi_j(z_j) = q$  for a  $z_j$  with  $(z_j, x)$  in the support of the distribution of  $(Z_j, X_e)$ , then we can construct the probability of the new policy from data in the policy histories.

- Given these assumptions, our ability to recover  $\Pr(D_{p'} = 1 \mid Z_{p'} = z, X_e = x)$  for all  $(z, x)$  in the support of  $(Z_{p'}, X_e)$  depends on what  $\Phi_j$  functions have been historically observed and the richness of the histories of  $Q_j, j \in \mathcal{H}_e$ . For each  $z_{p'}$  evaluation point in the support of the distribution of  $Z_{p'}$ , there is a corresponding  $q = \Phi_{p'}(z_{p'})$  evaluation point in the support of the distribution of  $Q_j = \Phi_j(Z_j)$ .
- If, in the policy histories, there is at least one  $j \in \mathcal{H}_e$  such that  $\Phi_j(z_j) = q$  for a  $z_j$  with  $(z_j, x)$  in the support of the distribution of  $(Z_j, X_e)$ , then we can construct the probability of the new policy from data in the policy histories.
- The methods used to extrapolate  $P_{p'}(\cdot)$  over new regions, discussed previously, apply here.

- If the distribution of  $Q_{p'}$  (or  $\Phi_{p'}$  and the distribution of  $Z_{p'}$ ) is known as part of the specification of the proposed policy, the distribution of  $F_{P_{p'}|X_e}$  can be constructed using the constructed  $P_{p'}$ .

- If the distribution of  $Q_{p'}$  (or  $\Phi_{p'}$  and the distribution of  $Z_{p'}$ ) is known as part of the specification of the proposed policy, the distribution of  $F_{P_{p'}|X_e}$  can be constructed using the constructed  $P_{p'}$ .
- Alternatively, if we can relate  $Q_{p'}$  to  $Q_j$  by  $Q_{p'} = \Psi_{p',j}(Q_j)$  for a known function  $\Psi_{p',j}$  or if we can relate  $Z_{p'}$  to  $Z_j$  by  $Z_{p'} = T_{p',j}(Z_j)$  for a known function  $T_{p',j}$ , and the distributions of  $Q_j$  and/or  $Z_j$  are known for some  $j \in \mathcal{H}_e$ , we can apply the method previously discussed to derive  $F_{P_{p'}|X_e}$  and hence the policy weights for the new policy.



- This approach assumes that a new policy acts on components of  $Q$  like a policy in  $\mathcal{H}_e$ , so it is possible to forecast the effect of a policy with nominally new aspects.

- This approach assumes that a new policy acts on components of  $Q$  like a policy in  $\mathcal{H}_e$ , so it is possible to forecast the effect of a policy with nominally new aspects.
- The essential idea is to recast the new aspects of policy in terms of old aspects previously measured.

- This approach assumes that a new policy acts on components of  $Q$  like a policy in  $\mathcal{H}_e$ , so it is possible to forecast the effect of a policy with nominally new aspects.
- The essential idea is to recast the new aspects of policy in terms of old aspects previously measured.
- Thus in a model of schooling, let  $D = \mathbf{1}[Y_1 - Y_0 - B \geq 0]$  where  $Y_1 - Y_0$  is the discounted gain in earnings from going to school and  $B$  is the tuition cost.

- This approach assumes that a new policy acts on components of  $Q$  like a policy in  $\mathcal{H}_e$ , so it is possible to forecast the effect of a policy with nominally new aspects.
- The essential idea is to recast the new aspects of policy in terms of old aspects previously measured.
- Thus in a model of schooling, let  $D = \mathbf{1}[Y_1 - Y_0 - B \geq 0]$  where  $Y_1 - Y_0$  is the discounted gain in earnings from going to school and  $B$  is the tuition cost.
- In this example, a decrease in a unit of cost ( $B$ ) has the same effect on choice as an increase in return ( $Y_1 - Y_0$ ).

- Historically, we might only observe variation in  $Y_1 - Y_0$  (say tuition has never previously been charged).

- Historically, we might only observe variation in  $Y_1 - Y_0$  (say tuition has never previously been charged).
- But  $B$  is on the same footing (has the same effect on choice, except for sign) as  $Y_1 - Y_0$ .

- Historically, we might only observe variation in  $Y_1 - Y_0$  (say tuition has never previously been charged).
- But  $B$  is on the same footing (has the same effect on choice, except for sign) as  $Y_1 - Y_0$ .
- The identified historical variation in  $Y_1 - Y_0$  can be used to nonparametrically forecast the effect of introducing  $B$ , provided that the support of  $P_{p'}$  is in the historical support generated by the policy histories in  $\mathcal{H}_e$ .

- Historically, we might only observe variation in  $Y_1 - Y_0$  (say tuition has never previously been charged).
- But  $B$  is on the same footing (has the same effect on choice, except for sign) as  $Y_1 - Y_0$ .
- The identified historical variation in  $Y_1 - Y_0$  can be used to nonparametrically forecast the effect of introducing  $B$ , provided that the support of  $P_{p'}$  is in the historical support generated by the policy histories in  $\mathcal{H}_e$ .
- Otherwise, some functional structure (parametric or semiparametric) must be imposed to solve the support problem for  $P_{p'}$ .



- Historically, we might only observe variation in  $Y_1 - Y_0$  (say tuition has never previously been charged).
- But  $B$  is on the same footing (has the same effect on choice, except for sign) as  $Y_1 - Y_0$ .
- The identified historical variation in  $Y_1 - Y_0$  can be used to nonparametrically forecast the effect of introducing  $B$ , provided that the support of  $P_{p'}$  is in the historical support generated by the policy histories in  $\mathcal{H}_e$ .
- Otherwise, some functional structure (parametric or semiparametric) must be imposed to solve the support problem for  $P_{p'}$ .
- We used this basic principle in constructing our econometric cost benefit analysis in Slide 414.

- As another example, following ?, consider the introduction of wage taxes in a world where there has never before been a tax.

- As another example, following ?, consider the introduction of wage taxes in a world where there has never before been a tax.
- This example is analyzed in ?.

- As another example, following ?, consider the introduction of wage taxes in a world where there has never before been a tax.
- This example is analyzed in ?.
- Let  $Z_j$  be the wage without taxes.

- As another example, following ?, consider the introduction of wage taxes in a world where there has never before been a tax.
- This example is analyzed in ?.
- Let  $Z_j$  be the wage without taxes.
- We seek to forecast a post-tax net wage  $Z_{p'} = (1 - \tau) Z_j + b$  where  $\tau$  is the tax rate and  $b$  is a constant shifter.

- As another example, following ?, consider the introduction of wage taxes in a world where there has never before been a tax.
- This example is analyzed in ?.
- Let  $Z_j$  be the wage without taxes.
- We seek to forecast a post-tax net wage  $Z_{p'} = (1 - \tau) Z_j + b$  where  $\tau$  is the tax rate and  $b$  is a constant shifter.
- Thus  $Z_{p'}$  is a known linear transformation of policy  $Z_j$ .

- As another example, following ?, consider the introduction of wage taxes in a world where there has never before been a tax.
- This example is analyzed in ?.
- Let  $Z_j$  be the wage without taxes.
- We seek to forecast a post-tax net wage  $Z_{p'} = (1 - \tau) Z_j + b$  where  $\tau$  is the tax rate and  $b$  is a constant shifter.
- Thus  $Z_{p'}$  is a known linear transformation of policy  $Z_j$ .
- We can construct  $Z_{p'}$  from  $Z_j$ .

- We can forecast under (A-1) using
$$\Pr(D_j = 1 \mid Z_j = z) = \Pr(D_{p'} = 1 \mid Z_{p'} = z).$$



- We can forecast under (A-1) using
$$\Pr(D_j = 1 \mid Z_j = z) = \Pr(D_{p'} = 1 \mid Z_{p'} = z).$$
- This assumes that the response to after tax wages is the same as the response to wages at the after tax level.

- We can forecast under (A-1) using
$$\Pr(D_j = 1 \mid Z_j = z) = \Pr(D_{p'} = 1 \mid Z_{p'} = z).$$
- This assumes that the response to after tax wages is the same as the response to wages at the after tax level.
- The issue is whether  $P_{p'|X_e}$  lies in the historical support, or whether extrapolation is needed.

- We can forecast under (A-1) using
$$\Pr(D_j = 1 \mid Z_j = z) = \Pr(D_{p'} = 1 \mid Z_{p'} = z).$$
- This assumes that the response to after tax wages is the same as the response to wages at the after tax level.
- The issue is whether  $P_{p'|X_e}$  lies in the historical support, or whether extrapolation is needed.
- Nonlinear versions of this example can be constructed.

- As a final example, environmental economists use variation in one component of cost (e.g., travel cost) to estimate the effect of a new cost (e.g., a park registration fee).

- As a final example, environmental economists use variation in one component of cost (e.g., travel cost) to estimate the effect of a new cost (e.g., a park registration fee).
- See ?.

- As a final example, environmental economists use variation in one component of cost (e.g., travel cost) to estimate the effect of a new cost (e.g., a park registration fee).
- See ?.
- Relating the costs and characteristics of new policies to the costs and characteristics of old policies is a standard, but sometimes controversial, method for forecasting the effects of new policies.

- In the context of our model, extrapolation and forecasting are confined to constructing  $P_{p'}$  and its distribution.

- In the context of our model, extrapolation and forecasting are confined to constructing  $P_{p'}$  and its distribution.
- If policy  $p'$ , characterized by vector  $Z_{p'}$ , consists of new components that cannot be related to  $Z_j, j \in \mathcal{H}_e$ , or a base set of characteristics whose variation cannot be identified, the problem is intractable.



- In the context of our model, extrapolation and forecasting are confined to constructing  $P_{p'}$  and its distribution.
- If policy  $p'$ , characterized by vector  $Z_{p'}$ , consists of new components that cannot be related to  $Z_j, j \in \mathcal{H}_e$ , or a base set of characteristics whose variation cannot be identified, the problem is intractable.
- Then  $P_{p'}$  and its distribution cannot be formed using econometric methods applied to historical data.

- When it can be applied, our approach allows us to simplify the policy forecasting problem and concentrate our attention on forecasting choice probabilities and their distribution in solving the policy forecasting problem.

- When it can be applied, our approach allows us to simplify the policy forecasting problem and concentrate our attention on forecasting choice probabilities and their distribution in solving the policy forecasting problem.
- We can use choice theory and choice data to construct these objects to forecast the impacts of new policies, by relating new policies to previously experienced policies.

## Forecasting the Effects of Policies in New Environments

- When the effects of policy  $p$  are forecast for a new environment  $e'$  from baseline environment  $e$ , and  $X_e \neq X_{e'}$ , in general both  $\Delta^{\text{MTE}}(x, u_D)$  and  $F_{P_p|X_e}$  will change.

## Forecasting the Effects of Policies in New Environments

- When the effects of policy  $p$  are forecast for a new environment  $e'$  from baseline environment  $e$ , and  $X_e \neq X_{e'}$ , in general both  $\Delta^{\text{MTE}}(x, u_D)$  and  $F_{P_p|X_e}$  will change.
- In general, neither object is environment invariant.

## Forecasting the Effects of Policies in New Environments

- When the effects of policy  $p$  are forecast for a new environment  $e'$  from baseline environment  $e$ , and  $X_e \neq X_{e'}$ , in general both  $\Delta^{\text{MTE}}(x, u_D)$  and  $F_{P_p|X_e}$  will change.
- In general, neither object is environment invariant.
- The new  $X_{e'}$  may have a different support than  $X_e$  or any other environment in  $\mathcal{H}$ .

## Forecasting the Effects of Policies in New Environments

- When the effects of policy  $p$  are forecast for a new environment  $e'$  from baseline environment  $e$ , and  $X_e \neq X_{e'}$ , in general both  $\Delta^{\text{MTE}}(x, u_D)$  and  $F_{P_p|X_e}$  will change.
- In general, neither object is environment invariant.
- The new  $X_{e'}$  may have a different support than  $X_e$  or any other environment in  $\mathcal{H}$ .
- In addition, the new  $(X_{e'}, U_D)$  stochastic relationship may be different from the historical  $(X_e, U_D)$  stochastic relationship.

- Constructing  $F_{P_p|X_{e'}}$  from  $F_{P_p|X_e}$  and  $F_{Z_p|X_{e'}}$  from  $F_{Z_p|X_e}$  can be done using (i) functional form (including semiparametric functional restrictions) or (ii) analytic continuation methods.



- Constructing  $F_{P_p|X_{e'}}$  from  $F_{P_p|X_e}$  and  $F_{Z_p|X_{e'}}$  from  $F_{Z_p|X_e}$  can be done using (i) functional form (including semiparametric functional restrictions) or (ii) analytic continuation methods.
- Notice that the maps  $T_{p,j}$  and  $\Phi_p$  may depend on  $X_e$  and so the induced changes in these transformations must also be modelled.

- Constructing  $F_{P_p|X_{e'}}$  from  $F_{P_p|X_e}$  and  $F_{Z_p|X_{e'}}$  from  $F_{Z_p|X_e}$  can be done using (i) functional form (including semiparametric functional restrictions) or (ii) analytic continuation methods.
- Notice that the maps  $T_{p,j}$  and  $\Phi_p$  may depend on  $X_e$  and so the induced changes in these transformations must also be modelled.
- There is a parallel discussion for  $\Delta^{\text{MTE}}(x, u_D)$ .

- Constructing  $F_{P_p|X_{e'}}$  from  $F_{P_p|X_e}$  and  $F_{Z_p|X_{e'}}$  from  $F_{Z_p|X_e}$  can be done using (i) functional form (including semiparametric functional restrictions) or (ii) analytic continuation methods.
- Notice that the maps  $T_{p,j}$  and  $\Phi_p$  may depend on  $X_e$  and so the induced changes in these transformations must also be modelled.
- There is a parallel discussion for  $\Delta^{\text{MTE}}(x, u_D)$ .
- The stochastic dependence between  $X_{e'}$  and  $(U_0, U_1, U_D)$  may be different from the stochastic dependence between  $X_e$  and  $(U_0, U_1, U_D)$ .

- Constructing  $F_{P_p|X_{e'}}$  from  $F_{P_p|X_e}$  and  $F_{Z_p|X_{e'}}$  from  $F_{Z_p|X_e}$  can be done using (i) functional form (including semiparametric functional restrictions) or (ii) analytic continuation methods.
- Notice that the maps  $T_{p,j}$  and  $\Phi_p$  may depend on  $X_e$  and so the induced changes in these transformations must also be modelled.
- There is a parallel discussion for  $\Delta^{\text{MTE}}(x, u_D)$ .
- The stochastic dependence between  $X_{e'}$  and  $(U_0, U_1, U_D)$  may be different from the stochastic dependence between  $X_e$  and  $(U_0, U_1, U_D)$ .
- We suppress the dependence of  $U_0$  and  $U_1$  on  $e$  and  $p$  only for convenience of exposition and make it explicit in the next paragraph.

- Forecasting new stochastic relationships between  $X_{e'}$  and  $(U_0, U_1, U_D)$  is a difficult task.

(A-9)

$$(U_{0,e,p}, U_{1,e,p}, U_{D,e,p}) \perp\!\!\!\perp (X_e, Z_p) \quad \forall e, p.$$

- Forecasting new stochastic relationships between  $X_{e'}$  and  $(U_0, U_1, U_D)$  is a difficult task.
- Some of the difficulty can be avoided if we invoke the traditional exogeneity assumptions of classical econometrics:

(A-10)

$$(U_{0,e,p}, U_{1,e,p}, U_{D,e,p}) \perp\!\!\!\perp (X_e, Z_p) \quad \forall e, p.$$

- Under (A-9), we only encounter the support problems for  $\Delta^{\text{MTE}}$  and the distribution of  $\Pr(D_p = 1 \mid Z_p, X_e)$  in constructing policy counterfactuals.

- Conditions (A-7), (A-8) and (A-9) are unnecessary if the only goal of the analysis is to establish internal validity, the standard objective of the treatment effect literature.



- Conditions (A-7), (A-8) and (A-9) are unnecessary if the only goal of the analysis is to establish internal validity, the standard objective of the treatment effect literature.
- This is problem P-1.

- Conditions (A-7), (A-8) and (A-9) are unnecessary if the only goal of the analysis is to establish internal validity, the standard objective of the treatment effect literature.
- This is problem P-1.
- Autonomy and exogeneity conditions become important issues if we seek external validity.

- Conditions (A-7), (A-8) and (A-9) are unnecessary if the only goal of the analysis is to establish internal validity, the standard objective of the treatment effect literature.
- This is problem P-1.
- Autonomy and exogeneity conditions become important issues if we seek external validity.
- An important lesson from this analysis is that as we try to make the treatment effect literature do the tasks of structural econometrics (i.e., make out-of-sample forecasts), common assumptions are invoked in the two literatures.

## A Comparison of Three Approaches to Policy Evaluation

- Table 9 compares the strengths and limitations of the three approaches to policy evaluation that we have discussed in this Handbook chapter and our contribution in Part I: the structural approach, the conventional treatment effect approach, and the approach to treatment effects based on the MTE function developed by ???.

# Table 9: Comparison of Alternative Approaches to Program Evaluation

	Structural Econometric Approach	Treatment Effect Approach	Approach Based on MTE
Interpretability	Well defined economic parameters and welfare comparisons	Link to economics and welfare comparisons obscure	Interpretable in terms of willingness to pay; weighted averages of the MTE answer well-posed economic questions
Range of Questions Addressed	Answers many counterfactual questions	Focuses on one treatment effect or narrow range of effects	With support conditions, generates all treatment parameters
Extrapolation to New Environments	Provides ingredients for extrapolation	Evaluates one program in one environment	Can be partially extrapolated; extrapolates to new policy environments with different distributions of the probability of participation due solely to differences in distributions of $Z$
Comparability Across Studies	Policy invariant parameters comparable across studies	Not generally comparable	Partially comparable; comparable across environments with different distributions of the probability of participation due solely to differences in distributions of $Z$
Key Econometric Problems	Exogeneity, policy invariance and selection bias	Selection bias	Selection bias: exogeneity and policy invariance if used for forecasting
Range of Policies that can be Evaluated	Programs with either partial or universal coverage, depending on variation in data (prices/endowments)	Programs with partial coverage (treatment and control groups)	Programs with partial coverage (treatment and control groups)
Extension to General Equilibrium Evaluation	Need to link to time series data; parameters compatible with general equilibrium theory	Difficult because link to economics is not precisely specified	Can be linked to nonparametric general equilibrium models under exogeneity and policy invariance

Source: Heckman and Vytlačil (2005)

- The approach based on the MTE function and the structural approach share interpretability of parameters.

- The approach based on the MTE function and the structural approach share interpretability of parameters.
- Like the structural approach, it addresses a range of policy evaluation questions.

- The approach based on the MTE function and the structural approach share interpretability of parameters.
- Like the structural approach, it addresses a range of policy evaluation questions.
- The MTE parameter is less comparable and less easily extrapolated across environments than are structural parameters, unless nonparametric versions of invariance and exogeneity assumptions are made.



- The approach based on the MTE function and the structural approach share interpretability of parameters.
- Like the structural approach, it addresses a range of policy evaluation questions.
- The MTE parameter is less comparable and less easily extrapolated across environments than are structural parameters, unless nonparametric versions of invariance and exogeneity assumptions are made.
- However,  $\Delta^{\text{MTE}}$  is comparable across populations with different distributions of  $P$  (conditional on  $X_e$ ) and results from one population can be applied to another population under the conditions presented in this section.

- Analysts can use  $\Delta^{\text{MTE}}$  to forecast a variety of policies.

- Analysts can use  $\Delta^{\text{MTE}}$  to forecast a variety of policies.
- This invariance property is shared with conventional structural parameters.

- Analysts can use  $\Delta^{\text{MTE}}$  to forecast a variety of policies.
- This invariance property is shared with conventional structural parameters.
- Our framework solves the problem of external validity, which is ignored in the standard treatment effect approach.

- Analysts can use  $\Delta^{\text{MTE}}$  to forecast a variety of policies.
- This invariance property is shared with conventional structural parameters.
- Our framework solves the problem of external validity, which is ignored in the standard treatment effect approach.
- The price of these advantages of the structural approach is the greater range of econometric problems that must be solved.

- Analysts can use  $\Delta^{\text{MTE}}$  to forecast a variety of policies.
- This invariance property is shared with conventional structural parameters.
- Our framework solves the problem of external validity, which is ignored in the standard treatment effect approach.
- The price of these advantages of the structural approach is the greater range of econometric problems that must be solved.
- They are avoided in the conventional treatment approach at the cost of producing parameters that cannot be linked to well-posed economic models and hence do not provide building blocks for an empirically motivated general equilibrium analysis or for investigation of the impacts of new public policies.

- $\Delta^{\text{MTE}}$  estimates the preferences of the agents being studied and provides a basis for integration with well posed economic models.

- $\Delta^{\text{MTE}}$  estimates the preferences of the agents being studied and provides a basis for integration with well posed economic models.
- If the goal of a study is to examine one policy in place (the problem of internal validity), the stronger assumptions invoked in this section of the chapter, and in structural econometrics, are unnecessary.



- $\Delta^{\text{MTE}}$  estimates the preferences of the agents being studied and provides a basis for integration with well posed economic models.
- If the goal of a study is to examine one policy in place (the problem of internal validity), the stronger assumptions invoked in this section of the chapter, and in structural econometrics, are unnecessary.
- Even if this is the only goal of the analysis, however, our approach allows the analyst to generate all treatment effects and IV estimands from a common parameter and provides a basis for unification of the treatment effect literature.

## Extension of MTE to the Analysis of More than Two Treatments and Associated Outcomes

- We have thus far analyzed models with two potential outcomes associated with receipt of binary treatments ( $D = 0$  or  $D = 1$ ).

## Extension of MTE to the Analysis of More than Two Treatments and Associated Outcomes

- We have thus far analyzed models with two potential outcomes associated with receipt of binary treatments ( $D = 0$  or  $D = 1$ ).
- Focusing on this simple case allows us to develop main ideas.

## Extension of MTE to the Analysis of More than Two Treatments and Associated Outcomes

- We have thus far analyzed models with two potential outcomes associated with receipt of binary treatments ( $D = 0$  or  $D = 1$ ).
- Focusing on this simple case allows us to develop main ideas.
- However, models with more than two outcomes are common in empirical work.

## Extension of MTE to the Analysis of More than Two Treatments and Associated Outcomes

- We have thus far analyzed models with two potential outcomes associated with receipt of binary treatments ( $D = 0$  or  $D = 1$ ).
- Focusing on this simple case allows us to develop main ideas.
- However, models with more than two outcomes are common in empirical work.
- ? analyze an ordered choice model with a single instrument that shifts people across all margins.

- We generalize their analysis in several ways.

- We generalize their analysis in several ways.
- We consider vectors of instruments, some of which may affect choices at all margins and some of which affect choices only at certain margins.

- We generalize their analysis in several ways.
- We consider vectors of instruments, some of which may affect choices at all margins and some of which affect choices only at certain margins.
- We then analyze a general unordered choice model.



## Background for our Analysis of the Ordered Choice Model

- ? extend their analysis of LATE to an ordered choice model with outcomes generated by a scalar instrument that can assume multiple values.

## Background for our Analysis of the Ordered Choice Model

- ? extend their analysis of LATE to an ordered choice model with outcomes generated by a scalar instrument that can assume multiple values.
- From their analysis of the effect of schooling on earnings, it is unclear even under a strengthened “monotonicity” condition whether IV estimates the effect of a change of schooling on earnings for a well defined margin of choice.

- To summarize their analysis, let  $\bar{S}$  be the number of possible outcome states with associated outcomes  $Y_s$  and choice indicators  $D_s$ ,  $s = 1, \dots, \bar{S}$ .

- To summarize their analysis, let  $\bar{S}$  be the number of possible outcome states with associated outcomes  $Y_s$  and choice indicators  $D_s$ ,  $s = 1, \dots, \bar{S}$ .
- The  $s$ , in their analysis, correspond to different levels of schooling.

- To summarize their analysis, let  $\bar{S}$  be the number of possible outcome states with associated outcomes  $Y_s$  and choice indicators  $D_s$ ,  $s = 1, \dots, \bar{S}$ .
- The  $s$ , in their analysis, correspond to different levels of schooling.
- For any two instrument values  $Z = z_i$  and  $Z = z_j$  with  $z_i > z_j$ , we can define associated indicators  $\{D_s(z_i)\}_{s=1}^{\bar{S}}$  and  $\{D_s(z_j)\}_{s=1}^{\bar{S}}$ , where  $D_s(z_i) = 1$  if a person assigned instrument value  $z_i$  chooses state  $s$ .

- To summarize their analysis, let  $\bar{S}$  be the number of possible outcome states with associated outcomes  $Y_s$  and choice indicators  $D_s$ ,  $s = 1, \dots, \bar{S}$ .
- The  $s$ , in their analysis, correspond to different levels of schooling.
- For any two instrument values  $Z = z_i$  and  $Z = z_j$  with  $z_i > z_j$ , we can define associated indicators  $\{D_s(z_i)\}_{s=1}^{\bar{S}}$  and  $\{D_s(z_j)\}_{s=1}^{\bar{S}}$ , where  $D_s(z_i) = 1$  if a person assigned instrument value  $z_i$  chooses state  $s$ .
- As in the two-outcome model, the instrument  $Z$  is assumed to be independent of the potential outcomes  $\{Y_s\}_{s=1}^{\bar{S}}$  as well as the associated indicator functions defined by fixing  $Z$  at  $z_i$  and  $z_j$ .

- Observed schooling for instrument  $z_j$  is  $S(z_j) = \sum_{s=1}^{\bar{S}} sD_s(z_j)$ .

- Observed schooling for instrument  $z_j$  is  $S(z_j) = \sum_{s=1}^{\bar{S}} sD_s(z_j)$ .
- Observed outcomes with this instrument are  $Y(z_j) = \sum_{s=1}^{\bar{S}} Y_s D_s(z_j)$ .



- Observed schooling for instrument  $z_j$  is  $S(z_j) = \sum_{s=1}^{\bar{S}} sD_s(z_j)$ .
- Observed outcomes with this instrument are  $Y(z_j) = \sum_{s=1}^{\bar{S}} Y_s D_s(z_j)$ .
- Angrist and Imbens show that IV (with  $Z = z_i$  and  $Z = z_j$ ) applied to  $S$  in a two stage least squares regression of  $Y$  on  $S$  identifies a “causal parameter”

$$\Delta^{\text{IV}} = \sum_{s=2}^{\bar{S}} \{E(Y_s - Y_{s-1} \mid S(z_i) \geq s > S(z_j))\} \cdot \frac{\Pr(S(z_i) \geq s > S(z_j))}{\sum_{s=2}^{\bar{S}} \Pr(S(z_i) \geq s > S(z_j))}. \quad (42)$$

- Observed schooling for instrument  $z_j$  is  $S(z_j) = \sum_{s=1}^{\bar{S}} sD_s(z_j)$ .
- Observed outcomes with this instrument are  $Y(z_j) = \sum_{s=1}^{\bar{S}} Y_s D_s(z_j)$ .
- Angrist and Imbens show that IV (with  $Z = z_i$  and  $Z = z_j$ ) applied to  $S$  in a two stage least squares regression of  $Y$  on  $S$  identifies a “causal parameter”

$$\Delta^{IV} = \sum_{s=2}^{\bar{S}} \{E(Y_s - Y_{s-1} \mid S(z_i) \geq s > S(z_j))\} \cdot \frac{\Pr(S(z_i) \geq s > S(z_j))}{\sum_{s=2}^{\bar{S}} \Pr(S(z_i) \geq s > S(z_j))}. \quad (42)$$

- This “causal parameter” is a weighted average of the gross returns from going from  $s - 1$  to  $s$  for persons induced by the change in the instrument to move from *any* schooling level below  $s$  to *any* schooling level  $s$  or above.

- Thus the conditioning set defining the  $s^{\text{th}}$  component of IV includes people who have schooling below  $s - 1$  at instrument value  $Z = z_j$  and people who have schooling above level  $s$  at instrument value  $Z = z_j$ .

- Thus the conditioning set defining the  $s^{\text{th}}$  component of IV includes people who have schooling below  $s - 1$  at instrument value  $Z = z_j$  and people who have schooling above level  $s$  at instrument value  $Z = z_j$ .
- In expression (42), the average return experienced by some of the people in the conditioning set for each component conditional expectation does not correspond to the average outcome corresponding to the gain in the argument of the expectation.

- Thus the conditioning set defining the  $s^{\text{th}}$  component of IV includes people who have schooling below  $s - 1$  at instrument value  $Z = z_j$  and people who have schooling above level  $s$  at instrument value  $Z = z_j$ .
- In expression (42), the average return experienced by some of the people in the conditioning set for each component conditional expectation does not correspond to the average outcome corresponding to the gain in the argument of the expectation.
- In the case where  $\bar{S} = 2$ , agents face only two choices and the margin of choice is well defined.

- Thus the conditioning set defining the  $s^{\text{th}}$  component of IV includes people who have schooling below  $s - 1$  at instrument value  $Z = z_j$  and people who have schooling above level  $s$  at instrument value  $Z = z_j$ .
- In expression (42), the average return experienced by some of the people in the conditioning set for each component conditional expectation does not correspond to the average outcome corresponding to the gain in the argument of the expectation.
- In the case where  $\bar{S} = 2$ , agents face only two choices and the margin of choice is well defined.
- Agents in each conditioning set are at different margins of choice.

- The weights are positive but, as noted by ?, persons can be counted multiple times in forming the weights.

- The weights are positive but, as noted by ?, persons can be counted multiple times in forming the weights.
- When they generalize their analysis to multiple-valued instruments, they use the ? weights.



- The weights are positive but, as noted by ?, persons can be counted multiple times in forming the weights.
- When they generalize their analysis to multiple-valued instruments, they use the ? weights.
- Whereas the weights in equation (42) can be constructed empirically using nonparametric discrete choice theory (see, e.g., our analysis in appendix B of Part I or ?), the terms in braces cannot be identified by any standard IV procedure.

- The weights are positive but, as noted by ?, persons can be counted multiple times in forming the weights.
- When they generalize their analysis to multiple-valued instruments, they use the ? weights.
- Whereas the weights in equation (42) can be constructed empirically using nonparametric discrete choice theory (see, e.g., our analysis in appendix B of Part I or ?), the terms in braces cannot be identified by any standard IV procedure.
- We present decompositions with components that are recoverable, whose weights can be estimated from the data and that are economically interpretable.

- In this section, we generalize LATE to a multiple outcome case where we can identify agents at different well defined margins of choice.

- In this section, we generalize LATE to a multiple outcome case where we can identify agents at different well defined margins of choice.
- Specifically, we (1) analyze both ordered and unordered choice models; (2) analyze outcomes associated with choices at various well defined margins; and (3) develop models with multiple instruments that can affect different margins of choice differently.

- In this section, we generalize LATE to a multiple outcome case where we can identify agents at different well defined margins of choice.
- Specifically, we (1) analyze both ordered and unordered choice models; (2) analyze outcomes associated with choices at various well defined margins; and (3) develop models with multiple instruments that can affect different margins of choice differently.
- With our methods, we can define and estimate a variety of economically interpretable parameters.

- In contrast, the Angrist-Imbens analysis produces a single “causal parameter” (42) that does not answer any well defined policy question such as that posed by the PRTE.

- In contrast, the Angrist-Imbens analysis produces a single “causal parameter” (42) that does not answer any well defined policy question such as that posed by the PRTE.
- We first consider an explicit ordered choice model and decompose the IV into policy-useful (identifiable) components.

## Analysis of an Ordered Choice Model

- Ordered choice models arise in many settings.



## Analysis of an Ordered Choice Model

- Ordered choice models arise in many settings.
- In schooling models, there are multiple grades.

## Analysis of an Ordered Choice Model

- Ordered choice models arise in many settings.
- In schooling models, there are multiple grades.
- One has to complete grade  $s - 1$  to proceed to grade  $s$ .

## Analysis of an Ordered Choice Model

- Ordered choice models arise in many settings.
- In schooling models, there are multiple grades.
- One has to complete grade  $s - 1$  to proceed to grade  $s$ .
- The ordered choice model has been widely used to fit data on schooling transitions (??).

## Analysis of an Ordered Choice Model

- Ordered choice models arise in many settings.
- In schooling models, there are multiple grades.
- One has to complete grade  $s - 1$  to proceed to grade  $s$ .
- The ordered choice model has been widely used to fit data on schooling transitions (??).
- Its nonparametric identifiability has been studied (??).

- It can also be used as a duration model for dynamic treatment effects with associated outcomes as in ?.

- It can also be used as a duration model for dynamic treatment effects with associated outcomes as in ?.
- It also represents the “vertical” model of the choice of product quality (???)

- It can also be used as a duration model for dynamic treatment effects with associated outcomes as in ?.
- It also represents the “vertical” model of the choice of product quality (???)
- Our analysis generalizes the analysis for the binary model in a parallel way.

- It can also be used as a duration model for dynamic treatment effects with associated outcomes as in ?.
- It also represents the “vertical” model of the choice of product quality (???)
- Our analysis generalizes the analysis for the binary model in a parallel way.
- Write potential outcomes as

$$Y_s = \mu_s(X, U_s) \quad s = 1, \dots, \bar{S}.$$



- It can also be used as a duration model for dynamic treatment effects with associated outcomes as in ?.
- It also represents the “vertical” model of the choice of product quality (???) .
- Our analysis generalizes the analysis for the binary model in a parallel way.
- Write potential outcomes as

$$Y_s = \mu_s(X, U_s) \quad s = 1, \dots, \bar{S}.$$

- The  $\bar{S}$  could be different schooling levels or product qualities.

- We define latent variables  $D_S^* = \mu_D(Z) - V$  where

$$D_s = \mathbf{1}[C_{s-1}(W_{s-1}) < \mu_D(Z) - V \leq C_s(W_s)], \quad s = 1, \dots, \bar{S},$$

and the cutoff values satisfy

$$C_{s-1}(W_{s-1}) \leq C_s(W_s), \quad C_0(W_0) = -\infty \quad \text{and} \quad C_{\bar{S}}(W_{\bar{S}}) = \infty.$$

The cutoffs used to define the intervals are allowed to depend on observed (by the economist) regressors  $W_s$ .

- We define latent variables  $D_s^* = \mu_D(Z) - V$  where

$$D_s = \mathbf{1}[C_{s-1}(W_{s-1}) < \mu_D(Z) - V \leq C_s(W_s)], \quad s = 1, \dots, \bar{S},$$

and the cutoff values satisfy

$$C_{s-1}(W_{s-1}) \leq C_s(W_s), \quad C_0(W_0) = -\infty \quad \text{and} \quad C_{\bar{S}}(W_{\bar{S}}) = \infty.$$

The cutoffs used to define the intervals are allowed to depend on observed (by the economist) regressors  $W_s$ .

- In Appendix, Slide 1127, we extend the analysis presented in the text to allow the cutoffs to depend on unobserved regressors as well, following structural analysis along these lines by ? and ?.

- We define latent variables  $D_s^* = \mu_D(Z) - V$  where

$$D_s = \mathbf{1}[C_{s-1}(W_{s-1}) < \mu_D(Z) - V \leq C_s(W_s)], \quad s = 1, \dots, \bar{S},$$

and the cutoff values satisfy

$$C_{s-1}(W_{s-1}) \leq C_s(W_s), \quad C_0(W_0) = -\infty \quad \text{and} \quad C_{\bar{S}}(W_{\bar{S}}) = \infty.$$

The cutoffs used to define the intervals are allowed to depend on observed (by the economist) regressors  $W_s$ .

- In Appendix, Slide 1127, we extend the analysis presented in the text to allow the cutoffs to depend on unobserved regressors as well, following structural analysis along these lines by ? and ?.
- Observed outcomes are:  $Y = \sum_{s=1}^{\bar{S}} Y_s D_s$ .

- The  $Z$  shift the index generally; the  $W_s$  affect  $s$ -specific transitions.

- The  $Z$  shift the index generally; the  $W_s$  affect  $s$ -specific transitions.
- Thus, in a schooling example,  $Z$  could include family background variables while  $W_s$  could include college tuition or opportunity wages for unskilled labor.

- The  $Z$  shift the index generally; the  $W_s$  affect  $s$ -specific transitions.
- Thus, in a schooling example,  $Z$  could include family background variables while  $W_s$  could include college tuition or opportunity wages for unskilled labor.
- Collect the  $W_s$  into  $W = (W_1, \dots, W_{\bar{S}})$ , and the  $U_s$  into  $U = (U_1, \dots, U_{\bar{S}})$ .

- The  $Z$  shift the index generally; the  $W_s$  affect  $s$ -specific transitions.
- Thus, in a schooling example,  $Z$  could include family background variables while  $W_s$  could include college tuition or opportunity wages for unskilled labor.
- Collect the  $W_s$  into  $W = (W_1, \dots, W_{\bar{S}})$ , and the  $U_s$  into  $U = (U_1, \dots, U_{\bar{S}})$ .
- Larger values of  $C_s(W_s)$  make it more likely that  $D_s = 1$ .



- The  $Z$  shift the index generally; the  $W_s$  affect  $s$ -specific transitions.
- Thus, in a schooling example,  $Z$  could include family background variables while  $W_s$  could include college tuition or opportunity wages for unskilled labor.
- Collect the  $W_s$  into  $W = (W_1, \dots, W_{\bar{S}})$ , and the  $U_s$  into  $U = (U_1, \dots, U_{\bar{S}})$ .
- Larger values of  $C_s(W_s)$  make it more likely that  $D_s = 1$ .
- The inequality restrictions on the  $C_s(W_s)$  functions play a critical role in defining the model and producing its statistical implications.

- Analogous to the assumptions made for the binary outcome model, we assume

(OC-1)

$(U_s, V) \perp\!\!\!\perp (Z, W) | X, s = 1, \dots, \bar{S}$ . (**Conditional Independence of the Instruments**).

(OC-2)

$\mu_D(Z)$  is a nondegenerate random variable conditional on  $X$  and  $W$ . (**Rank Condition**).

(OC-3)

*The distribution of  $V$  is continuous.*

(OC-4)

 *$E(|Y_s|) < \infty, s = 1, \dots, \bar{S}$ . (**Finite Means**).*

(OC-5)

 *$0 < \Pr(D_s = 1|X) < 1$  for  $s = 1, \dots, \bar{S}$  for all  $X$ . (**In large samples, there are some persons in each treatment state**).*

(OC-6)

*For  $s = 1, \dots, \bar{S} - 1$ , the distribution of  $C_s(W_s)$  conditional on  $X, Z$  and the other  $C_j(W_j), j = 1, \dots, \bar{S}, j \neq s$ , is nondegenerate and continuous.*

- Assumptions (OC-1)–(OC-5) play roles analogous to their counterparts in the two-outcome model, (A-1)–(A-5).

- Assumptions (OC-1)–(OC-5) play roles analogous to their counterparts in the two-outcome model, (A-1)–(A-5).
- (OC-6) is a new condition that is key to identification of the  $\Delta^{\text{MTE}}$  defined below for each transition.

- Assumptions (OC-1)–(OC-5) play roles analogous to their counterparts in the two-outcome model, (A-1)–(A-5).
- (OC-6) is a new condition that is key to identification of the  $\Delta^{\text{MTE}}$  defined below for each transition.
- It assumes that we can vary the choice sets of agents at different margins of schooling choice without affecting other margins of choice.

- Assumptions (OC-1)–(OC-5) play roles analogous to their counterparts in the two-outcome model, (A-1)–(A-5).
- (OC-6) is a new condition that is key to identification of the  $\Delta^{\text{MTE}}$  defined below for each transition.
- It assumes that we can vary the choice sets of agents at different margins of schooling choice without affecting other margins of choice.
- A necessary condition for (OC-6) to hold is that at least one element of  $W_s$  is nondegenerate and continuous conditional on  $X, Z$  and  $C_j(W_j)$  for  $j \neq s$ .

- Intuitively, one needs an instrument (or source of variability) for each transition.



- Intuitively, one needs an instrument (or source of variability) for each transition.
- The continuity of the regressor allows us to differentiate with respect to  $C_s(W_s)$ , like we differentiated with respect to  $P(Z)$  to estimate the MTE in the analysis of the two-outcome model.

- Intuitively, one needs an instrument (or source of variability) for each transition.
- The continuity of the regressor allows us to differentiate with respect to  $C_s(W_s)$ , like we differentiated with respect to  $P(Z)$  to estimate the MTE in the analysis of the two-outcome model.
- The analysis of ? discussed in the introduction to this section makes independence and monotonicity assumptions that generalize their earlier work.

- Intuitively, one needs an instrument (or source of variability) for each transition.
- The continuity of the regressor allows us to differentiate with respect to  $C_s(W_s)$ , like we differentiated with respect to  $P(Z)$  to estimate the MTE in the analysis of the two-outcome model.
- The analysis of ? discussed in the introduction to this section makes independence and monotonicity assumptions that generalize their earlier work.
- They do not consider estimation of transition-specific parameters as we do, or even transition-specific LATE.

- We present a different decomposition of the IV estimator where each component can be recovered from the data, and where the transition-specific MTEs answer well defined and economically interpretable policy evaluation questions.

- We present a different decomposition of the IV estimator where each component can be recovered from the data, and where the transition-specific MTEs answer well defined and economically interpretable policy evaluation questions.
- The probability of  $D_s = 1$  given  $X, Z$  and  $W$  is generated by an ordered choice model:

$$\begin{aligned} \Pr(D_s = 1 \mid Z, W, X) &\equiv P_s(Z, W, X) \\ &= \Pr(C_{s-1}(W_{s-1}) < \mu_D(Z) - V \leq C_s(W_s) \mid X). \end{aligned}$$

Analogous to the binary case, we can define  $U_D = F_{V|X}(V)$  so  $U_D \sim \text{Uniform}[0, 1]$  under our assumption that the distribution of  $V$  is absolutely continuous with respect to Lebesgue measure.

- The probability integral transformation used extensively in the binary choice model is somewhat less useful for analyzing ordered choices, so we work with both  $U_D$  and  $V$  in this section of the chapter.

- The probability integral transformation used extensively in the binary choice model is somewhat less useful for analyzing ordered choices, so we work with both  $U_D$  and  $V$  in this section of the chapter.
- Monotonic transformations of  $V$  induce monotonic transformations of  $\mu_D(Z) - C_S(W_S)$ , but one is not free to form arbitrary monotonic transformations of  $\mu_D(Z)$  and  $C_S(W_S)$  separately.

- The probability integral transformation used extensively in the binary choice model is somewhat less useful for analyzing ordered choices, so we work with both  $U_D$  and  $V$  in this section of the chapter.
- Monotonic transformations of  $V$  induce monotonic transformations of  $\mu_D(Z) - C_s(W_s)$ , but one is not free to form arbitrary monotonic transformations of  $\mu_D(Z)$  and  $C_s(W_s)$  separately.
- Using the probability integral transformation, the expression for choice  $s$  is  $D_s = \mathbf{1}[F_{V|X}(\mu_D(Z) - C_{s-1}(W_{s-1})) > U_D \geq F_{V|X}(\mu_D(Z) - C_s(W_s))]$ .



- Keeping the conditioning on  $X$  implicit, we define

$$P_s(Z, W) = F_V(\mu_D(Z) - C_{s-1}(W_{s-1})) - F_V(\mu_D(Z) - C_s(W_s)).$$

- Keeping the conditioning on  $X$  implicit, we define
 
$$P_s(Z, W) = F_V(\mu_D(Z) - C_{s-1}(W_{s-1})) - F_V(\mu_D(Z) - C_s(W_s)).$$
- It is convenient to work with the probability that  $S > s$ ,
 
$$\pi_s(Z, W_s) = F_V(\mu_D(Z) - C_s(W_s)) =$$

$$\Pr\left(\sum_{j=s+1}^{\bar{S}} D_j = 1 \mid Z, W_s\right), \pi_{\bar{S}}(Z, W_{\bar{S}}) = 0, \pi_0(Z, W_0) = 1$$
 and  $P_s(Z, W) = \pi_{s-1}(Z, W_{s-1}) - \pi_s(Z, W_s).$

- Keeping the conditioning on  $X$  implicit, we define  $P_s(Z, W) = F_V(\mu_D(Z) - C_{s-1}(W_{s-1})) - F_V(\mu_D(Z) - C_s(W_s))$ .
- It is convenient to work with the probability that  $S > s$ ,  $\pi_s(Z, W_s) = F_V(\mu_D(Z) - C_s(W_s)) = \Pr\left(\sum_{j=s+1}^{\bar{S}} D_j = 1 \mid Z, W_s\right)$ ,  $\pi_{\bar{S}}(Z, W_{\bar{S}}) = 0$ ,  $\pi_0(Z, W_0) = 1$  and  $P_s(Z, W) = \pi_{s-1}(Z, W_{s-1}) - \pi_s(Z, W_s)$ .
- The transition-specific  $\Delta^{\text{MTE}}$  for the transition from  $s$  to  $s+1$  is defined in terms of  $U_D$ .

$$\Delta_{s,s+1}^{\text{MTE}}(x, u_D) = E(Y_{s+1} - Y_s \mid X = x, U_D = u_D), \quad s = 1, \dots, \bar{S} - 1.$$

- Alternatively, one can condition on  $V$ .

- Alternatively, one can condition on  $V$ .
- Analogous to the analysis of the earlier sections of this chapter, when we set  $u_D = \pi_s(Z, W_s)$ , we obtain the mean return to persons indifferent between  $s$  and  $s + 1$  at mean level of utility  $\pi_s(Z, W_s)$ .

- In this notation, keeping  $X$  implicit, the mean outcome  $Y$ , conditional on  $(Z, W)$ , is the sum of the mean outcomes conditional on each state weighted by the probability of being in each state summed over all states:

$$\begin{aligned}
 E(Y|Z, W) &= \sum_{s=1}^{\bar{S}} E(Y_s | D_s = 1, Z, W) \Pr(D_s = 1 | Z, W) \quad (43) \\
 &= \sum_{s=1}^{\bar{S}} \int_{\pi_s(Z, W_s)}^{\pi_{s-1}(Z, W_{s-1})} E(Y_s | U_D = u_D) du_D,
 \end{aligned}$$

where we use conditional independence assumption (OC-1) to obtain the final expression.

- Analogous to the result for the binary outcome model, we obtain the index sufficiency restriction

$$E(Y|Z, W) = E(Y | \pi(Z, W)), \text{ where } \pi(Z, W) = [\pi_1(Z, W_1), \dots, \pi_{\bar{S}-1}(Z, W_{\bar{S}-1})].$$

- Analogous to the result for the binary outcome model, we obtain the index sufficiency restriction

$$E(Y|Z, W) = E(Y | \pi(Z, W)), \text{ where } \pi(Z, W) = [\pi_1(Z, W_1), \dots, \pi_{\bar{S}-1}(Z, W_{\bar{S}-1})].$$

- The choice probabilities encode all of the influence of  $(Z, W)$  on outcomes.



- Analogous to the result for the binary outcome model, we obtain the index sufficiency restriction

$$E(Y|Z, W) = E(Y | \pi(Z, W)), \text{ where } \pi(Z, W) = [\pi_1(Z, W_1), \dots, \pi_{\bar{S}-1}(Z, W_{\bar{S}-1})].$$

- The choice probabilities encode all of the influence of  $(Z, W)$  on outcomes.
- We can identify  $\pi_s(z, w_s)$  for  $(z, w_s)$  in the support of the distribution of  $(Z, W_s)$  from the relationship
$$\pi_s(z, w_s) = \Pr(\sum_{j=s+1}^{\bar{S}} D_j = 1 \mid Z = z, W_s = w_s).$$

- Analogous to the result for the binary outcome model, we obtain the index sufficiency restriction

$$E(Y|Z, W) = E(Y | \pi(Z, W)), \text{ where } \pi(Z, W) = [\pi_1(Z, W_1), \dots, \pi_{\bar{S}-1}(Z, W_{\bar{S}-1})].$$

- The choice probabilities encode all of the influence of  $(Z, W)$  on outcomes.
- We can identify  $\pi_s(z, w_s)$  for  $(z, w_s)$  in the support of the distribution of  $(Z, W_s)$  from the relationship
 
$$\pi_s(z, w_s) = \Pr(\sum_{j=s+1}^{\bar{S}} D_j = 1 | Z = z, W_s = w_s).$$
- Thus  $E(Y | \pi(Z, W) = \pi)$  is identified for all  $\pi$  in the support of  $\pi(Z, W)$ .

- Assumptions (OC-1), (OC-3), and (OC-4) imply that  $E(Y | \pi(Z, W) = \pi)$  is differentiable in  $\pi$ .

- Assumptions (OC-1), (OC-3), and (OC-4) imply that  $E(Y | \pi(Z, W) = \pi)$  is differentiable in  $\pi$ .
- So  $\frac{\partial}{\partial \pi} E(Y | \pi(Z, W) = \pi)$  is well-defined.

- Assumptions (OC-1), (OC-3), and (OC-4) imply that  $E(Y | \pi(Z, W) = \pi)$  is differentiable in  $\pi$ .
- So  $\frac{\partial}{\partial \pi} E(Y | \pi(Z, W) = \pi)$  is well-defined.
- Thus analogous to the result obtained in the binary case

$$\begin{aligned} \frac{\partial E(Y | \pi(Z, W) = \pi)}{\partial \pi_s} &= \Delta_{s,s+1}^{\text{MTE}}(U_D = \pi_s) & (44) \\ &= E(Y_{s+1} - Y_s | U_D = \pi_s). \end{aligned}$$

- Assumptions (OC-1), (OC-3), and (OC-4) imply that  $E(Y | \pi(Z, W) = \pi)$  is differentiable in  $\pi$ .
- So  $\frac{\partial}{\partial \pi} E(Y | \pi(Z, W) = \pi)$  is well-defined.
- Thus analogous to the result obtained in the binary case

$$\begin{aligned} \frac{\partial E(Y | \pi(Z, W) = \pi)}{\partial \pi_s} &= \Delta_{s,s+1}^{\text{MTE}}(U_D = \pi_s) & (44) \\ &= E(Y_{s+1} - Y_s | U_D = \pi_s). \end{aligned}$$

- Equation (44) is the basis for identification of the transition-specific MTE from data on  $(Y, Z, X)$ .

- From index sufficiency, we can express (43) as

$$\begin{aligned}
 E(Y \mid \pi(Z, W) = \pi) &= \sum_{s=1}^{\bar{S}} E(Y_s \mid \pi_s \leq U_D < \pi_{s-1})(\pi_{s-1} - \pi_s) & (45) \\
 &= \sum_{s=1}^{\bar{S}-1} \left[ \begin{array}{l} E(Y_{s+1} \mid \pi_{s+1} \leq U_D < \pi_s) \\ -E(Y_s \mid \pi_s \leq U_D < \pi_{s-1}) \end{array} \right] \pi_s \\
 &\quad + E(Y_1 \mid \pi_1 \leq U_D < 1) \\
 &= \sum_{s=1}^{\bar{S}-1} \{m_{s+1}(\pi_{s+1}, \pi_s) - m_s(\pi_s, \pi_{s-1})\} \pi_s \\
 &\quad + E(Y_1 \mid \pi_1 \leq U_D < 1)
 \end{aligned}$$

where  $m_s(\pi_s, \pi_{s-1}) = E[Y_s \mid \pi_s \leq U_D < \pi_{s-1}]$ .

- In general, this expression is a nonlinear function of  $(\pi_s, \pi_{s-1})$ .



- In general, this expression is a nonlinear function of  $(\pi_s, \pi_{s-1})$ .
- This model has a testable restriction of index sufficiency in the general case:  $E(Y|\pi(Z, W) = \pi)$  is a nonlinear function that is additive in functions of  $(\pi_s, \pi_{s-1})$  so there are no interactions between  $\pi_s$  and  $\pi_{s'}$  if  $|s - s'| > 1$ , i.e.,

$$\frac{\partial^2 E(Y | \pi(Z, W) = \pi)}{\partial \pi_s \partial \pi_{s'}} = 0 \quad \text{if } |s - s'| > 1.$$

- In general, this expression is a nonlinear function of  $(\pi_s, \pi_{s-1})$ .
- This model has a testable restriction of index sufficiency in the general case:  $E(Y | \pi(Z, W) = \pi)$  is a nonlinear function that is additive in functions of  $(\pi_s, \pi_{s-1})$  so there are no interactions between  $\pi_s$  and  $\pi_{s'}$  if  $|s - s'| > 1$ , i.e.,

$$\frac{\partial^2 E(Y | \pi(Z, W) = \pi)}{\partial \pi_s \partial \pi_{s'}} = 0 \quad \text{if } |s - s'| > 1.$$

- Observe that if  $U_D \perp\!\!\!\perp U_s$  for  $s = 1, \dots, \bar{S}$ ,

$$\begin{aligned} E(Y | \pi(Z, W) = \pi) &= \sum_{s=1}^{\bar{S}} E(Y_s)(\pi_{s-1} - \pi_s) \\ &= \sum_{s=1}^{\bar{S}-1} [E(Y_{s+1}) - E(Y_s)] \pi_s + E(Y_1). \end{aligned}$$

- Defining  $E(Y_{s+1}) - E(Y_s) = \Delta_{s,s+1}^{\text{ATE}}$ ,  
 $E(Y | \pi(Z, W) = \pi) = \sum_{s=1}^{\bar{S}-1} \Delta_{s,s+1}^{\text{ATE}} \pi_s + E(Y_1)$ .

- Defining  $E(Y_{s+1}) - E(Y_s) = \Delta_{s,s+1}^{\text{ATE}}$ ,  
 $E(Y | \pi(Z, W) = \pi) = \sum_{s=1}^{\bar{S}-1} \Delta_{s,s+1}^{\text{ATE}} \pi_s + E(Y_1)$ .
- Thus, under full independence, we obtain linearity of the conditional mean of  $Y$  in the  $\pi_s$ ,  $s = 1, \dots, \bar{S}$ .

- Defining  $E(Y_{s+1}) - E(Y_s) = \Delta_{s,s+1}^{\text{ATE}}$ ,  
 $E(Y | \pi(Z, W) = \pi) = \sum_{s=1}^{\bar{S}-1} \Delta_{s,s+1}^{\text{ATE}} \pi_s + E(Y_1)$ .
- Thus, under full independence, we obtain linearity of the conditional mean of  $Y$  in the  $\pi_s$ ,  $s = 1, \dots, \bar{S}$ .
- This result generalizes the test for the presence of essential heterogeneity presented in Slide 152 to the ordered case.

- Defining  $E(Y_{s+1}) - E(Y_s) = \Delta_{s,s+1}^{\text{ATE}}$ ,  
 $E(Y | \pi(Z, W) = \pi) = \sum_{s=1}^{\bar{S}-1} \Delta_{s,s+1}^{\text{ATE}} \pi_s + E(Y_1)$ .
- Thus, under full independence, we obtain linearity of the conditional mean of  $Y$  in the  $\pi_s$ ,  $s = 1, \dots, \bar{S}$ .
- This result generalizes the test for the presence of essential heterogeneity presented in Slide 152 to the ordered case.
- We can ignore the complexity induced by the model of essential heterogeneity if  $E(Y | \pi(Z, W) = \pi)$  is linear in the  $\pi_s$  and can use conventional IV estimators to identify well-defined treatment effects.

## The Policy Relevant Treatment Effect for the Ordered Choice Model

- The policy relevant treatment effect compares the mean outcome under one policy regime  $p$  with the mean outcome under policy regime  $p'$ .

## The Policy Relevant Treatment Effect for the Ordered Choice Model

- The policy relevant treatment effect compares the mean outcome under one policy regime  $p$  with the mean outcome under policy regime  $p'$ .
- It is defined analogously to the way it is defined in the binary case in Slide 139 and in ??.



## The Policy Relevant Treatment Effect for the Ordered Choice Model

- The policy relevant treatment effect compares the mean outcome under one policy regime  $p$  with the mean outcome under policy regime  $p'$ .
- It is defined analogously to the way it is defined in the binary case in Slide 139 and in ??.
- Policies  $(p, p')$  are assumed to induce different distributions of  $(Z, W)$ ,  $F^p(Z, W)$ .

## The Policy Relevant Treatment Effect for the Ordered Choice Model

- The policy relevant treatment effect compares the mean outcome under one policy regime  $p$  with the mean outcome under policy regime  $p'$ .
- It is defined analogously to the way it is defined in the binary case in Slide 139 and in ??.
- Policies  $(p, p')$  are assumed to induce different distributions of  $(Z, W)$ ,  $F^p(Z, W)$ .
- Forming  $E_p(Y) = \int E(Y | Z = z, W = w) dF^p_{Z,W}(z, w)$  for each policy  $p$ , the policy relevant treatment effect is  $E_{p'}(Y) - E_p(Y)$ .

- We can represent the PRTE as a weighted average of pairwise MTE:

$$\Delta_{p,p'}^{\text{PRTE}} = E_{p'}(Y) - E_p(Y) = \sum_{s=1}^{\bar{S}-1} \int E(Y_{s+1} - Y_s \mid V = v) \omega_{p,p'}(v) dF(v). \quad (46)$$

The weights are known functions of the data.

- We can represent the PRTE as a weighted average of pairwise MTE:

$$\Delta_{p,p'}^{\text{PRTE}} = E_{p'}(Y) - E_p(Y) = \sum_{s=1}^{\bar{S}-1} \int E(Y_{s+1} - Y_s \mid V = v) \omega_{p,p'}(v) dF(v). \quad (46)$$

The weights are known functions of the data.

- See Appendix, Slide 1139, for a derivation of the weights and expression (46).

- We can represent the PRTE as a weighted average of pairwise MTE:

$$\Delta_{p,p'}^{\text{PRTE}} = E_{p'}(Y) - E_p(Y) = \sum_{s=1}^{\bar{S}-1} \int E(Y_{s+1} - Y_s \mid V = v) \omega_{p,p'}(v) dF(v). \quad (46)$$

The weights are known functions of the data.

- See Appendix, Slide 1139, for a derivation of the weights and expression (46).
- Using the probability integral transform, we can alternatively express this in terms of  $U_D = F_{V|X}(V)$ .

## What do Instruments Identify in the Ordered Choice Model?

- We now characterize what scalar instrument  $J(Z, W)$  identifies.

## What do Instruments Identify in the Ordered Choice Model?

- We now characterize what scalar instrument  $J(Z, W)$  identifies.
- When  $Y$  is log earnings, it is common practice to regress  $Y$  on  $S$  where  $S$  is completed years of schooling and call the coefficient on  $S$  a rate of return.

## What do Instruments Identify in the Ordered Choice Model?

- We now characterize what scalar instrument  $J(Z, W)$  identifies.
- When  $Y$  is log earnings, it is common practice to regress  $Y$  on  $S$  where  $S$  is completed years of schooling and call the coefficient on  $S$  a rate of return.
- We seek an expression for the instrumental variables estimator of the effect of  $S$  on  $Y$  in the ordered choice model:

$$\frac{\text{Cov}(J(Z, W), Y)}{\text{Cov}(J(Z, W), D)}, \quad (47)$$

where  $S = \sum_{s=1}^{\bar{S}} sD_s$  is the number of years of schooling attainment.



- We keep the conditioning on  $X$  implicit.

- We keep the conditioning on  $X$  implicit.
- We now analyze the weights for IV.

- We keep the conditioning on  $X$  implicit.
- We now analyze the weights for IV.
- Their full derivation is presented in Appendix, Slide 1142.

- We keep the conditioning on  $X$  implicit.
- We now analyze the weights for IV.
- Their full derivation is presented in Appendix, Slide 1142.
- Define  $K_s(v) = E\left(\tilde{J}(Z, W) \mid \mu_D(Z) - C_s(W_s) > v\right) \Pr(\mu_D(Z) - C_s(W_s) > v)$ , where  $\tilde{J}(Z, W) = J(Z, W) - E(J(Z, W))$ .

- Thus,

$$\begin{aligned}\Delta_J^{IV} &= \frac{\text{Cov}(J, Y)}{\text{Cov}(J, S)} \\ &= \sum_{s=1}^{\bar{S}-1} \int E(Y_{s+1} - Y_s \mid V = v) \omega(s, v) f_V(v) dv,\end{aligned}\tag{48}$$

where

$$\begin{aligned}\omega(s, v) &= \frac{K_s(v)}{\sum_{s=1}^{\bar{S}} s \int [K_{s-1}(v) - K_s(v)] f_V(v) dv} \\ &= \frac{K_s(v)}{\sum_{s=1}^{\bar{S}-1} \int K_s(v) f_V(v) dv},\end{aligned}$$

and clearly  $\sum_{s=1}^{\bar{S}-1} \int \omega(s, v) f_V(v) dv = 1$ ,  $\omega(0, v) = 0$ , and  $\omega(\bar{S}, v) = 0$ .

- We can rewrite this result in terms of the MTE, expressed in terms of  $u_D$

$$\Delta_{s,s+1}^{\text{MTE}}(u_D) = E(Y_{s+1} - Y_s \mid U_D = u_D)$$

so that

$$\frac{\text{Cov}(J, Y)}{\text{Cov}(J, S)} = \sum_{s=1}^{\bar{S}-1} \int_0^1 \Delta_{s,s+1}^{\text{MTE}}(u_D) \tilde{\omega}(s, u_D) du_D,$$

where

$$\begin{aligned} \tilde{\omega}(s, u_D) &= \frac{\tilde{K}_s(u_D)}{\sum_{s=1}^{\bar{S}} s \int_0^1 [\tilde{K}_{s-1}(u_D) - \tilde{K}_s(u_D)] du_D} \\ &= \frac{\tilde{K}_s(u_D)}{\sum_{s=1}^{\bar{S}-1} \int_0^1 \tilde{K}_s(u_D) du_D} \end{aligned} \quad (49)$$

and

$$\tilde{K}_s(u_D) = E(\tilde{J}(Z, W) \mid \pi_s(Z, W_s) \geq u_D) \Pr(\pi_s(Z, W_s) \geq u_D). \quad (50)$$

- Compare equations (49) and (50) for the ordered choice model to equations (23) and (24) for the binary choice model.

- Compare equations (49) and (50) for the ordered choice model to equations (23) and (24) for the binary choice model.
- The numerator of the weights for the  $\Delta^{\text{MTE}}$  in the ordered choice model for a particular transition is exactly the numerator of the weights for the binary choice model, substituting  $\pi_s(Z, W_s) = \Pr(S > s \mid Z, W_s)$  for  $P(Z) = \Pr(D = 1 \mid Z)$ .



- Compare equations (49) and (50) for the ordered choice model to equations (23) and (24) for the binary choice model.
- The numerator of the weights for the  $\Delta^{\text{MTE}}$  in the ordered choice model for a particular transition is exactly the numerator of the weights for the binary choice model, substituting  $\pi_s(Z, W_s) = \Pr(S > s \mid Z, W_s)$  for  $P(Z) = \Pr(D = 1 \mid Z)$ .
- The numerator for the weights for IV in the binary choice model is driven by the connection between the instrument and  $P(Z)$ .

- Compare equations (49) and (50) for the ordered choice model to equations (23) and (24) for the binary choice model.
- The numerator of the weights for the  $\Delta^{\text{MTE}}$  in the ordered choice model for a particular transition is exactly the numerator of the weights for the binary choice model, substituting  $\pi_s(Z, W_s) = \Pr(S > s \mid Z, W_s)$  for  $P(Z) = \Pr(D = 1 \mid Z)$ .
- The numerator for the weights for IV in the binary choice model is driven by the connection between the instrument and  $P(Z)$ .
- The numerator for the weights for IV in the ordered choice model for a particular transition is driven by the connection between the instrument and  $\pi_s(Z, W_s)$ .

- Compare equations (49) and (50) for the ordered choice model to equations (23) and (24) for the binary choice model.
- The numerator of the weights for the  $\Delta^{\text{MTE}}$  in the ordered choice model for a particular transition is exactly the numerator of the weights for the binary choice model, substituting  $\pi_s(Z, W_s) = \Pr(S > s \mid Z, W_s)$  for  $P(Z) = \Pr(D = 1 \mid Z)$ .
- The numerator for the weights for IV in the binary choice model is driven by the connection between the instrument and  $P(Z)$ .
- The numerator for the weights for IV in the ordered choice model for a particular transition is driven by the connection between the instrument and  $\pi_s(Z, W_s)$ .
- The denominator of the weights is the covariance between the instrument and  $D$  (or  $S$ ) for the binary (or ordered) case, respectively.

- However, in the binary case the covariance between the instrument and  $D$  is completely determined by the covariance with  $S$  between the instrument and  $P(Z)$ , while in the ordered choice case the covariance depends on the relationship between the instrument and the full vector  $[\pi_1(Z, W_1), \dots, \pi_{\bar{S}-1}(Z, W_{\bar{S}-1})]$ .

- However, in the binary case the covariance between the instrument and  $D$  is completely determined by the covariance with  $S$  between the instrument and  $P(Z)$ , while in the ordered choice case the covariance depends on the relationship between the instrument and the full vector  $[\pi_1(Z, W_1), \dots, \pi_{\bar{S}-1}(Z, W_{\bar{S}-1})]$ .
- Comparing our decomposition of  $\Delta^{IV}$  to decomposition (42), ours corresponds to weighting up marginal outcomes across well defined and adjacent boundary values experienced by agents having their instruments manipulated whereas the Angrist-Imbens decomposition corresponds to outcomes not experienced by some of the persons whose instruments are being manipulated.

- From equation (50), the IV estimator using  $J(Z, W)$  as an instrument satisfies the following properties.

- From equation (50), the IV estimator using  $J(Z, W)$  as an instrument satisfies the following properties.
- (a) The numerator of the weights on  $\Delta_{s,s+1}^{\text{MTE}}(u_D)$  is non-negative for all  $u_D$  if  $E(J(Z, W_s) \mid \pi_s(Z, W_s) \geq \pi_s)$  is weakly monotonic in  $\pi_s$ .

- From equation (50), the IV estimator using  $J(Z, W)$  as an instrument satisfies the following properties.
- (a) The numerator of the weights on  $\Delta_{s,s+1}^{\text{MTE}}(u_D)$  is non-negative for all  $u_D$  if  $E(J(Z, W_s) \mid \pi_s(Z, W_s) \geq \pi_s)$  is weakly monotonic in  $\pi_s$ .
- For example, if  $\text{Cov}(\pi_s(Z, W_s), S) > 0$ , setting  $J(Z, W) = \pi_s(Z, W_s)$  will lead to nonnegative weights on  $\Delta_{s,s+1}^{\text{MTE}}(u_D)$ , though it may lead to negative weights on other transitions.



- A second property (b) is that the support of the weights on  $\Delta_{s,s+1}^{\text{MTE}}$  using  $\pi_s(Z, W_s)$  as the instrument is  $(\pi_s^{\text{Min}}, \pi_s^{\text{Max}})$  where  $\pi_s^{\text{Min}}$  and  $\pi_s^{\text{Max}}$  are the minimum and maximum values in the support of  $\pi_s(Z, W_s)$ , respectively, and the support of the weights on  $\Delta_{s,s+1}^{\text{MTE}}$  using any other instrument is a subset of  $(\pi_s^{\text{Min}}, \pi_s^{\text{Max}})$ .

- A second property (b) is that the support of the weights on  $\Delta_{s,s+1}^{\text{MTE}}$  using  $\pi_s(Z, W_s)$  as the instrument is  $(\pi_s^{\text{Min}}, \pi_s^{\text{Max}})$  where  $\pi_s^{\text{Min}}$  and  $\pi_s^{\text{Max}}$  are the minimum and maximum values in the support of  $\pi_s(Z, W_s)$ , respectively, and the support of the weights on  $\Delta_{s,s+1}^{\text{MTE}}$  using any other instrument is a subset of  $(\pi_s^{\text{Min}}, \pi_s^{\text{Max}})$ .
- A third property (c) is that the weights on  $\Delta_{s,s+1}^{\text{MTE}}$  implied by using  $J(Z, W)$  as an instrument are the same as the weights on  $\Delta_{s,s+1}^{\text{MTE}}$  implied by using  $E(J(Z, W) \mid \pi_s(Z, W_s))$  as the instrument.

- A second property (b) is that the support of the weights on  $\Delta_{s,s+1}^{\text{MTE}}$  using  $\pi_s(Z, W_s)$  as the instrument is  $(\pi_s^{\text{Min}}, \pi_s^{\text{Max}})$  where  $\pi_s^{\text{Min}}$  and  $\pi_s^{\text{Max}}$  are the minimum and maximum values in the support of  $\pi_s(Z, W_s)$ , respectively, and the support of the weights on  $\Delta_{s,s+1}^{\text{MTE}}$  using any other instrument is a subset of  $(\pi_s^{\text{Min}}, \pi_s^{\text{Max}})$ .
- A third property (c) is that the weights on  $\Delta_{s,s+1}^{\text{MTE}}$  implied by using  $J(Z, W)$  as an instrument are the same as the weights on  $\Delta_{s,s+1}^{\text{MTE}}$  implied by using  $E(J(Z, W) \mid \pi_s(Z, W_s))$  as the instrument.
- Our analysis generalizes that of ? and ? by considering multiple instruments and by introducing both transition-specific instruments (the  $W$ ) and general instruments ( $Z$ ) across all transitions.

- In general, the method of linear instrumental variables applied to  $S$  does not estimate anything that is economically interpretable.

- In general, the method of linear instrumental variables applied to  $S$  does not estimate anything that is economically interpretable.
- It is not guaranteed to estimate a positive number even if the MTE is everywhere positive since the weights can be negative.

- In general, the method of linear instrumental variables applied to  $S$  does not estimate anything that is economically interpretable.
- It is not guaranteed to estimate a positive number even if the MTE is everywhere positive since the weights can be negative.
- In contrast, we can use our generalization of LIV presented in equation (44) under conditions (OC-1)–(OC-6) to apply LIV to identify  $\Delta^{\text{MTE}}$  for each transition, which can be used to build up  $\Delta^{\text{PRTE}}$  using weights that can be estimated.

## Some Theoretical Examples of the Weights in the Ordered Choice Model

- Suppose that the distributions of  $W_s$ ,  $s = 1, \dots, \bar{S}$ , are degenerate so that the  $C_s$  are constants satisfying  $C_1 < \dots < C_{\bar{S}-1}$ .

## Some Theoretical Examples of the Weights in the Ordered Choice Model

- Suppose that the distributions of  $W_s$ ,  $s = 1, \dots, \bar{S}$ , are degenerate so that the  $C_s$  are constants satisfying  $C_1 < \dots < C_{\bar{S}-1}$ .
- This is the classical ordered choice model.



## Some Theoretical Examples of the Weights in the Ordered Choice Model

- Suppose that the distributions of  $W_s$ ,  $s = 1, \dots, \bar{S}$ , are degenerate so that the  $C_s$  are constants satisfying  $C_1 < \dots < C_{\bar{S}-1}$ .
- This is the classical ordered choice model.
- In this case,  $\pi_s(Z, W_s) = F_V(\mu_D(Z) - C_s)$  for any  $s = 1, \dots, \bar{S}$ .

## Some Theoretical Examples of the Weights in the Ordered Choice Model

- Suppose that the distributions of  $W_s$ ,  $s = 1, \dots, \bar{S}$ , are degenerate so that the  $C_s$  are constants satisfying  $C_1 < \dots < C_{\bar{S}-1}$ .
- This is the classical ordered choice model.
- In this case,  $\pi_s(Z, W_s) = F_V(\mu_D(Z) - C_s)$  for any  $s = 1, \dots, \bar{S}$ .
- For this special case, using  $J$  as an instrument will lead to nonnegative weights on all transitions if  $J(Z, W)$  is a monotonic function of  $\mu_D(Z)$ .

## Some Theoretical Examples of the Weights in the Ordered Choice Model

- Suppose that the distributions of  $W_s$ ,  $s = 1, \dots, \bar{S}$ , are degenerate so that the  $C_s$  are constants satisfying  $C_1 < \dots < C_{\bar{S}-1}$ .
- This is the classical ordered choice model.
- In this case,  $\pi_s(Z, W_s) = F_V(\mu_D(Z) - C_s)$  for any  $s = 1, \dots, \bar{S}$ .
- For this special case, using  $J$  as an instrument will lead to nonnegative weights on all transitions if  $J(Z, W)$  is a monotonic function of  $\mu_D(Z)$ .
- For example, note that  $\mu_D(Z) - C_s > v$  can be written as  $\mu_D(Z) > C_s + F_V^{-1}(u_D)$ .

- Using  $\mu_D(Z)$  as the instrument leads to weights on  $\Delta_{s,s+1}^{MTE}(u_D)$  of the form specified above with  $\tilde{K}_s(u_D) = \left[ E(\mu_D(Z) \mid \mu_D(Z) > F_V^{-1}(u_D) + C_s) - E(\mu_D(Z)) \right] \Pr(\mu_D(Z) > F_V^{-1}(u_D) + C_s)$ .

- Using  $\mu_D(Z)$  as the instrument leads to weights on  $\Delta_{s,s+1}^{MTE}(u_D)$  of the form specified above with  $\tilde{K}_s(u_D) = \left[ E(\mu_D(Z) \mid \mu_D(Z) > F_V^{-1}(u_D) + C_s) - E(\mu_D(Z)) \right] \Pr(\mu_D(Z) > F_V^{-1}(u_D) + C_s)$ .
- Clearly, these weights will be nonnegative for all points of evaluation and will be strictly positive for any evaluation point  $u_D$  such that  $1 > \Pr(\mu_D(Z) > F_V^{-1}(u_D) + C_s) > 0$ .

- Using  $\mu_D(Z)$  as the instrument leads to weights on  $\Delta_{s,s+1}^{MTE}(u_D)$  of the form specified above with  $\tilde{K}_s(u_D) = \left[ E(\mu_D(Z) \mid \mu_D(Z) > F_V^{-1}(u_D) + C_s) - E(\mu_D(Z)) \right] \Pr(\mu_D(Z) > F_V^{-1}(u_D) + C_s)$ .
- Clearly, these weights will be nonnegative for all points of evaluation and will be strictly positive for any evaluation point  $u_D$  such that  $1 > \Pr(\mu_D(Z) > F_V^{-1}(u_D) + C_s) > 0$ .
- Next consider the case where  $C_s(W_s) = W_s$ , a scalar, for  $s = 1, \dots, \bar{S} - 1$ , and where  $\mu_D(Z) = 0$ .

- Using  $\mu_D(Z)$  as the instrument leads to weights on  $\Delta_{s,s+1}^{\text{MTE}}(u_D)$  of the form specified above with  $\tilde{K}_s(u_D) = \left[ E(\mu_D(Z) \mid \mu_D(Z) > F_V^{-1}(u_D) + C_s) - E(\mu_D(Z)) \right] \Pr(\mu_D(Z) > F_V^{-1}(u_D) + C_s)$ .
- Clearly, these weights will be nonnegative for all points of evaluation and will be strictly positive for any evaluation point  $u_D$  such that  $1 > \Pr(\mu_D(Z) > F_V^{-1}(u_D) + C_s) > 0$ .
- Next consider the case where  $C_s(W_s) = W_s$ , a scalar, for  $s = 1, \dots, \bar{S} - 1$ , and where  $\mu_D(Z) = 0$ .
- Consider  $J(Z, W) = W_s$ , a purely transition-specific instrument.

- In this case, the weight on  $\Delta_{s,s+1}^{\text{MTE}}(u_D)$  is of the form given above, with

$$\tilde{K}_s(u_D) = \left[ E(W_s \mid W_s > F_V^{-1}(u_D)) - E(W_s) \right] \Pr(W_s > F_V^{-1}(u_D)),$$

which will be nonnegative for all evaluation points and strictly positive for any evaluation point such that  $1 > \Pr(W_s > F_V^{-1}(u_D)) > 0$ .



- In this case, the weight on  $\Delta_{s,s+1}^{\text{MTE}}(u_D)$  is of the form given above, with

$$\tilde{K}_s(u_D) = \left[ E(W_s \mid W_s > F_V^{-1}(u_D)) - E(W_s) \right] \Pr(W_s > F_V^{-1}(u_D)),$$

which will be nonnegative for all evaluation points and strictly positive for any evaluation point such that  $1 > \Pr(W_s > F_V^{-1}(u_D)) > 0$ .

- What are the implied weights on  $\Delta_{s',s'+1}^{\text{MTE}}(u_D)$  for  $s' \neq s$ ?

- In this case, the weight on  $\Delta_{s,s+1}^{\text{MTE}}(u_D)$  is of the form given above, with

$$\tilde{K}_s(u_D) = \left[ E(W_s \mid W_s > F_V^{-1}(u_D)) - E(W_s) \right] \Pr(W_s > F_V^{-1}(u_D)),$$

which will be nonnegative for all evaluation points and strictly positive for any evaluation point such that  $1 > \Pr(W_s > F_V^{-1}(u_D)) > 0$ .

- What are the implied weights on  $\Delta_{s',s'+1}^{\text{MTE}}(u_D)$  for  $s' \neq s$ ?
- First, consider the case where  $W_s$  is independent of  $W_{s'}$  for  $s \neq s'$ .

- In this case, the weight on  $\Delta_{s,s+1}^{\text{MTE}}(u_D)$  is of the form given above, with

$$\tilde{K}_s(u_D) = \left[ E(W_s \mid W_s > F_V^{-1}(u_D)) - E(W_s) \right] \Pr(W_s > F_V^{-1}(u_D)),$$

which will be nonnegative for all evaluation points and strictly positive for any evaluation point such that  $1 > \Pr(W_s > F_V^{-1}(u_D)) > 0$ .

- What are the implied weights on  $\Delta_{s',s'+1}^{\text{MTE}}(u_D)$  for  $s' \neq s$ ?
- First, consider the case where  $W_s$  is independent of  $W_{s'}$  for  $s \neq s'$ .
- This independence of  $W_s$  and  $W_{s'}$  is not in conflict with the requirement  $W_s > W_{s'}$  for  $s > s'$  if the supports do not overlap for any  $s' \neq s$ .

- In this case, the weight on  $\Delta_{s',s'+1}^{\text{MTE}}(u_D)$  for  $s' \neq s$  is of the form given above with

$$\tilde{K}_{s'}(u_D) = \left[ E(W_s \mid W_{s'} > F_V^{-1}(u_D)) - E(W_s) \right] \Pr(W_{s'} > F_V^{-1}(u_D)) = 0.$$

Thus, in this case, the instrument only weights the  $\Delta^{\text{MTE}}$  for the  $s$  to  $s + 1$  transition.

- In this case, the weight on  $\Delta_{s',s'+1}^{\text{MTE}}(u_D)$  for  $s' \neq s$  is of the form given above with

$$\tilde{K}_{s'}(u_D) = \left[ E(W_s \mid W_{s'} > F_V^{-1}(u_D)) - E(W_s) \right] \Pr(W_{s'} > F_V^{-1}(u_D)) = 0.$$

Thus, in this case, the instrument only weights the  $\Delta^{\text{MTE}}$  for the  $s$  to  $s + 1$  transition.

- Note that this result relies critically on the assumption that  $W_s$  is independent of  $W_{s'}$  for  $s' \neq s$ .

- In this case, the weight on  $\Delta_{s',s'+1}^{\text{MTE}}(u_D)$  for  $s' \neq s$  is of the form given above with

$$\tilde{K}_{s'}(u_D) = \left[ E(W_s \mid W_{s'} > F_V^{-1}(u_D)) - E(W_s) \right] \Pr(W_{s'} > F_V^{-1}(u_D)) = 0.$$

Thus, in this case, the instrument only weights the  $\Delta^{\text{MTE}}$  for the  $s$  to  $s + 1$  transition.

- Note that this result relies critically on the assumption that  $W_s$  is independent of  $W_{s'}$  for  $s' \neq s$ .
- Consider another version of this example where  $C_s(W_s) = W_s$ ,  $s = 1, \dots, \bar{S} - 1$ , with  $W_s$  a scalar, but now allow  $\mu_D(Z)$  to have a nondegenerate distribution and allow there to be dependence across the  $W_s$ .

- In particular, consider the case where  $W = (W_1, \dots, W_{\bar{S}-1})$  is a continuous random vector with a density given by

$$\frac{\prod_{i=1}^{\bar{S}-1} f_i(w_i) \mathbf{1}[w_1 < w_2 < \dots < w_{\bar{S}-1}]}{\int \dots \int \left[ \mathbf{1}[w_1 < w_2 < \dots < w_{\bar{S}-1}] \prod_{i=1}^{\bar{S}-1} f_i(w_i) \right] dw_1 \dots dw_{\bar{S}-1}}$$

for some marginal density functions  $f_1(w_1), f_2(w_2), \dots, f_{\bar{S}-1}(w_{\bar{S}-1})$ .

- In particular, consider the case where  $W = (W_1, \dots, W_{\bar{S}-1})$  is a continuous random vector with a density given by

$$\frac{\prod_{i=1}^{\bar{S}-1} f_i(w_i) \mathbf{1}[w_1 < w_2 < \dots < w_{\bar{S}-1}]}{\int \dots \int \left[ \mathbf{1}[w_1 < w_2 < \dots < w_{\bar{S}-1}] \prod_{i=1}^{\bar{S}-1} f_i(w_i) \right] dw_1 \dots dw_{\bar{S}-1}}$$

for some marginal density functions  $f_1(w_1), f_2(w_2), \dots, f_{\bar{S}-1}(w_{\bar{S}-1})$ .

- In this case, using  $W_j$  as the instrument, we have

$$\begin{aligned} \omega(s, v) = & \left( \int \dots \int_{-\infty < w_1 < \dots < w_{\bar{S}-1} < \infty} (w_j - E(w_j))(1 - F_{\mu_D(Z)}(w_s + v)) \right. \\ & \times f_1(w_1) \dots f_{\bar{S}-1}(w_{\bar{S}-1}) dw_1 \dots dw_{\bar{S}-1} f_V(v) dv \Big) \\ & \times \left( \sum_{s=1}^{\bar{S}-1} \int \int \dots \int_{-\infty < w_1 < \dots < w_{\bar{S}-1} < \infty} (w_j - E(w_j))(1 - F_{\mu_D(Z)}(w_s + v)) \right. \\ & \times f_1(w_1) \dots f_{\bar{S}-1}(w_{\bar{S}-1}) dw_1 \dots dw_{\bar{S}-1} f_V(v) dv \Big)^{-1}. \end{aligned}$$



- In the special case where  $\mu_D(Z) \sim \text{Uniform}(-K, K)$ , with  $Z \perp\!\!\!\perp W_s$  for  $s = 1, \dots, \bar{S} - 1$ , assuming  $-K < w_s + v < K$  for all  $w_s, v$  in the support of  $W_s$  and  $V$  respectively, the numerator

$$\int \cdots \int_{-\infty < w_1 < \cdots < w_{\bar{S}-1} < \infty} (w_j - E(w_j)) \\ \times \frac{(w_s + v + K)}{2K} f_1(w_1) \cdots f_{\bar{S}-1}(w_{\bar{S}-1}) dw_1 \cdots dw_{\bar{S}-1} f_V(v) dv$$

is

Observe that when the latent  $W_j, W_s$  are independently distributed for all  $j, s$ , by Bickel's Theorem (?), we know that this expression is positive.

- In the special case where  $\mu_D(Z) \sim \text{Uniform}(-K, K)$ , with  $Z \perp\!\!\!\perp W_s$  for  $s = 1, \dots, \bar{S} - 1$ , assuming  $-K < w_s + v < K$  for all  $w_s, v$  in the support of  $W_s$  and  $V$  respectively, the numerator

$$\int \cdots \int_{-\infty < w_1 < \cdots < w_{\bar{S}-1} < \infty} (w_j - E(w_j)) \\ \times \frac{(w_s + v + K)}{2K} f_1(w_1) \cdots f_{\bar{S}-1}(w_{\bar{S}-1}) dw_1 \cdots dw_{\bar{S}-1} f_V(v) dv$$

is

Observe that when the latent  $W_j, W_s$  are independently distributed for all  $j, s$ , by Bickel's Theorem (?), we know that this expression is positive.

- (This is trivial when  $j = s$ .) The ordering  $W_1 < \cdots < W_{\bar{S}-1}$  implies that  $W_l$  is stochastically increasing in  $W_j$  for  $l < j$  (the lower boundary is shifted to the right).

- Hence, because of the order on the  $W$  implied by the ordered discrete choice model, a positive weighting is produced.

- Hence, because of the order on the  $W$  implied by the ordered discrete choice model, a positive weighting is produced.
- This result can be overturned when  $F(w)$  has a general structure.

- Hence, because of the order on the  $W$  implied by the ordered discrete choice model, a positive weighting is produced.
- This result can be overturned when  $F(w)$  has a general structure.
- The positive dependence induced by the order on the components of  $W$  can be reversed by negative dependence in the structure of  $F(w)$ .

- Hence, because of the order on the  $W$  implied by the ordered discrete choice model, a positive weighting is produced.
- This result can be overturned when  $F(w)$  has a general structure.
- The positive dependence induced by the order on the components of  $W$  can be reversed by negative dependence in the structure of  $F(w)$ .
- We present examples of these phenomena in our discussions in figures 19 and 20 below.

## Some Numerical Examples of the IV Weights

- Figures 16–18 plot the transition-specific MTEs and the IV weights for the models and distributions of the weights at the base of each of the figures.

## Some Numerical Examples of the IV Weights

- Figures 16–18 plot the transition-specific MTEs and the IV weights for the models and distributions of the weights at the base of each of the figures.
- We consider a three outcome ( $\bar{S} = 3$ ) model with common instruments ( $Z$ ) and transition-specific ( $W_s$ ) instruments.



## Some Numerical Examples of the IV Weights

- Figures 16–18 plot the transition-specific MTEs and the IV weights for the models and distributions of the weights at the base of each of the figures.
- We consider a three outcome ( $\bar{S} = 3$ ) model with common instruments ( $Z$ ) and transition-specific ( $W_s$ ) instruments.
- The  $Z$  and  $W_s$ ,  $s = 1, \dots, \bar{S}$ , are assumed to be independent.

## Some Numerical Examples of the IV Weights

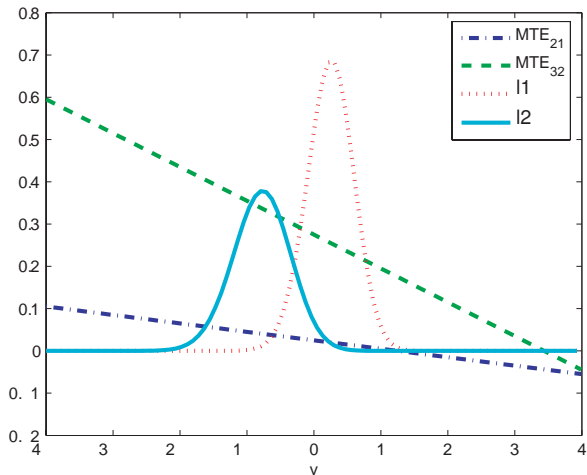
- Figures 16–18 plot the transition-specific MTEs and the IV weights for the models and distributions of the weights at the base of each of the figures.
- We consider a three outcome ( $\bar{S} = 3$ ) model with common instruments ( $Z$ ) and transition-specific ( $W_s$ ) instruments.
- The  $Z$  and  $W_s$ ,  $s = 1, \dots, \bar{S}$ , are assumed to be independent.
- The exact specification is given in the notes below figure 16.

## Some Numerical Examples of the IV Weights

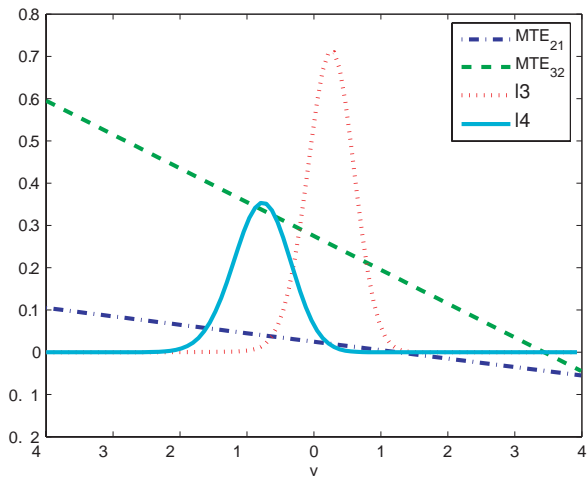
- Figures 16–18 plot the transition-specific MTEs and the IV weights for the models and distributions of the weights at the base of each of the figures.
- We consider a three outcome ( $\bar{S} = 3$ ) model with common instruments ( $Z$ ) and transition-specific ( $W_s$ ) instruments.
- The  $Z$  and  $W_s$ ,  $s = 1, \dots, \bar{S}$ , are assumed to be independent.
- The exact specification is given in the notes below figure 16.
- In this example,  $D_s$  can be interpreted as an indicator of schooling.

Figure 16: Treatment Parameters and IV – The Generalized Ordered Choice Roy Model under Normality

A. Z as Instrument



## B. $W_1$ as Instrument



Outcomes

$$\begin{aligned}
 Y_1 &= \alpha + \beta_1 + U_1 \\
 Y_2 &= \alpha + \beta_2 + U_2 \\
 Y_3 &= \alpha + \beta_3 + U_3
 \end{aligned}$$

Choice Model

$$D_s = \mathbf{1}[W_{s-1} < \gamma Z - V \leq W_s]$$

$s = 1, 2, 3$

Parameterization

$$(U_1, U_2, U_3, V) \sim N(\mathbf{0}, \Sigma_{UV}), \quad (Z, W_1, W_2) \sim N(\mu_{ZW}, \Sigma_{ZW}) \text{ and } W_0 = -\infty; W_3 = \infty.$$

$$\Sigma_{UV} = \begin{bmatrix} 1 & 0.16 & 0.2 & -0.3 \\ 0.16 & 0.64 & 0.16 & -0.32 \\ 0.2 & 0.16 & 1 & -0.4 \\ -0.3 & -0.32 & -0.4 & 1 \end{bmatrix}, \quad \mu_{ZW} = (-0.6, -1.08, 0.08) \text{ and } \Sigma_{ZW} = \begin{bmatrix} 0.1 & 0 & 0 \\ 0 & 0.1 & 0.09 \\ 0 & 0.09 & 0.25 \end{bmatrix}$$

$$\begin{aligned}
 \text{Cov}(U_2 - U_1, V) &= -0.02 & \text{Cov}(U_3 - U_2, V) &= -0.08 \\
 \beta_1 &= 0; \beta_2 &= 0.025; \beta_3 &= 0.3, \gamma = 1
 \end{aligned}$$

## IV Estimates and their Components\*

Parameter	Value
$\Delta^{IVZ}$	0.1487
$\Delta_{12}^{IVZ}$	0.0120
$\Delta_{23}^{IVZ}$	0.1367
$\Delta^{IVW_1}$	0.1406
$\Delta_{12}^{IVW_1}$	0.0126
$\Delta_{23}^{IVW_1}$	0.1280

## Treatment Parameters and their Values

Parameter	Value
$ATE_{12} = E(Y_2 - Y_1)$	0.025
$ATE_{23} = E(Y_3 - Y_2)$	0.275
$TT_{12} = E(Y_2 - Y_1   D_2 = 1)$	0.0282
$TT_{23} = E(Y_3 - Y_2   D_3 = 1)$	0.1908
$TUT_{12} = E(Y_2 - Y_1   D_1 = 1)$	0.0060
$TUT_{23} = E(Y_3 - Y_2   D_2 = 1)$	0.2956

\*  $IV^Z$  is decomposed as:

$$\begin{aligned}
 IV^Z &= \int E(Y_2 - Y_1 | V = v) \omega^Z(1, v) f_V(v) dv + \int E(Y_3 - Y_2 | V = v) \omega^Z(2, v) f_V(v) dv \\
 &= IV_{21}^Z + IV_{32}^Z.
 \end{aligned}$$

An analogous decomposition applies to  $IV^{W_1}$ .

- $Y_1$  is the potential earnings of the person as a dropout,  $Y_2$  is the potential earnings of the person as a high school graduate, and  $Y_3$  is the potential earnings of the person as a college graduate.



- $Y_1$  is the potential earnings of the person as a dropout,  $Y_2$  is the potential earnings of the person as a high school graduate, and  $Y_3$  is the potential earnings of the person as a college graduate.
- There are two transitions:  $1 \rightarrow 2$  and  $2 \rightarrow 3$ .

- $Y_1$  is the potential earnings of the person as a dropout,  $Y_2$  is the potential earnings of the person as a high school graduate, and  $Y_3$  is the potential earnings of the person as a college graduate.
- There are two transitions:  $1 \rightarrow 2$  and  $2 \rightarrow 3$ .
- The IV estimates using  $Z_1$  and  $W_1$  as instruments are reported transition by transition and overall decomposing IV representation (48) into its transition-specific components.

- $Y_1$  is the potential earnings of the person as a dropout,  $Y_2$  is the potential earnings of the person as a high school graduate, and  $Y_3$  is the potential earnings of the person as a college graduate.
- There are two transitions:  $1 \rightarrow 2$  and  $2 \rightarrow 3$ .
- The IV estimates using  $Z_1$  and  $W_1$  as instruments are reported transition by transition and overall decomposing IV representation (48) into its transition-specific components.
- The IV weights are defined by equations (49) and (50).

- In particular, when the first element of  $Z$ ,  $Z_1$ , is used as the instrument, we can decompose  $IV^{Z_1}$  as

$$\begin{aligned} IV^{Z_1} &= \sum_{s=1}^2 \int E(Y_{s+1} - Y_s \mid V = v) \omega^{Z_1(s,v)} f_V(v) dv \\ &= \int \Delta_{12}^{\text{MTE}}(v) \omega^{Z_1(1,v)} f_V(v) dv + \int \Delta_{23}^{\text{MTE}}(v) \omega^{Z_1(2,v)} f_V(v) dv \\ &= IV_{21}^{Z_1} + IV_{32}^{Z_1}. \end{aligned}$$

- In particular, when the first element of  $Z$ ,  $Z_1$ , is used as the instrument, we can decompose  $IV^{Z_1}$  as

$$\begin{aligned} IV^{Z_1} &= \sum_{s=1}^2 \int E(Y_{s+1} - Y_s \mid V = v) \omega^{Z_1(s,v)} f_V(v) dv \\ &= \int \Delta_{12}^{\text{MTE}}(v) \omega^{Z_1(1,v)} f_V(v) dv + \int \Delta_{23}^{\text{MTE}}(v) \omega^{Z_1(2,v)} f_V(v) dv \\ &= IV_{21}^{Z_1} + IV_{32}^{Z_1}. \end{aligned}$$

- The same logic applies for the decomposition of  $IV^P$  which uses  $P(Z)$  as an instrument.

- In particular, when the first element of  $Z$ ,  $Z_1$ , is used as the instrument, we can decompose  $IV^{Z_1}$  as

$$\begin{aligned}
 IV^{Z_1} &= \sum_{s=1}^2 \int E(Y_{s+1} - Y_s \mid V = v) \omega^{Z_1(s,v)} f_V(v) dv \\
 &= \int \Delta_{12}^{\text{MTE}}(v) \omega^{Z_1}(1, v) f_V(v) dv + \int \Delta_{23}^{\text{MTE}}(v) \omega^{Z_1}(2, v) f_V(v) dv \\
 &= IV_{21}^{Z_1} + IV_{32}^{Z_1}.
 \end{aligned}$$

- The same logic applies for the decomposition of  $IV^P$  which uses  $P(Z)$  as an instrument.
- These decompositions show in this case that an important component of the total values of  $IV^Z$  and  $IV^{W_1}$  comes from the  $2 \rightarrow 3$  transition.

- The bottom table presents the transition-specific treatment parameters.

- The bottom table presents the transition-specific treatment parameters.
- In figure 16, the shape of the IV weights for  $Z_1$  and  $W_1$  are nearly identical.



- The bottom table presents the transition-specific treatment parameters.
- In figure 16, the shape of the IV weights for  $Z_1$  and  $W_1$  are nearly identical.
- The IV estimates reflect this.

- The bottom table presents the transition-specific treatment parameters.
- In figure 16, the shape of the IV weights for  $Z_1$  and  $W_1$  are nearly identical.
- The IV estimates reflect this.
- The bottom table reveals that the IV estimates are far from standard treatment parameters.

- In figure 17, the IV weights for the  $Z_1$  and  $W_1$  are very different.

- In figure 17, the IV weights for the  $Z_1$  and  $W_1$  are very different.
- So, correspondingly, are the IV estimates produced from each instrument, which are far off the mark of the standard treatment parameters shown in the bottom of the table.

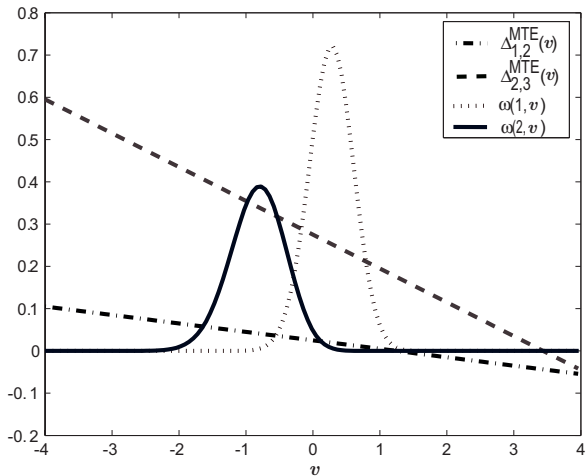
- In figure 17, the IV weights for the  $Z_1$  and  $W_1$  are very different.
- So, correspondingly, are the IV estimates produced from each instrument, which are far off the mark of the standard treatment parameters shown in the bottom of the table.
- Observe that the IV weight for  $W_1$  in the second transition is negative for an interval of values.

- In figure 17, the IV weights for the  $Z_1$  and  $W_1$  are very different.
- So, correspondingly, are the IV estimates produced from each instrument, which are far off the mark of the standard treatment parameters shown in the bottom of the table.
- Observe that the IV weight for  $W_1$  in the second transition is negative for an interval of values.
- This accounts for the dramatically lower IV estimate based on  $W_1$  as the instrument.

- In figure 17, the IV weights for the  $Z_1$  and  $W_1$  are very different.
- So, correspondingly, are the IV estimates produced from each instrument, which are far off the mark of the standard treatment parameters shown in the bottom of the table.
- Observe that the IV weight for  $W_1$  in the second transition is negative for an interval of values.
- This accounts for the dramatically lower IV estimate based on  $W_1$  as the instrument.
- Figure 18 shows a different configuration of  $(Z_1, W_1, W_2)$ .

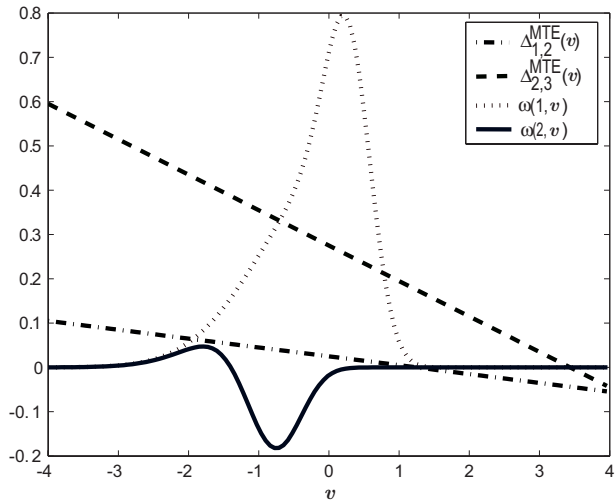
Figure 17: Treatment Parameters and IV – The Generalized Ordered Choice Roy Model under Normality: Case I

A. Z as Instrument





## B. $W_1$ as Instrument



Outcomes

$$Y_1 = \alpha + \beta_1 + U_1$$

$$Y_2 = \alpha + \beta_2 + U_2$$

$$Y_3 = \alpha + \beta_3 + U_3$$

Choice Model

$$D_s = \mathbf{1}[W_{s-1} < \gamma Z - V \leq W_s]$$

$s = 1, 2, 3$

Parameterization

$(U_1, U_2, U_3, V) \sim N(\mathbf{0}, UV)$ ,  $(Z, W_1, W_2) \sim N(\mu_{ZW}, ZW)$  and  $W_0 = -\infty$ ;  $W_3 = \infty$ .

$$UV = \begin{bmatrix} 1 & 0.16 & 0.2 & -0.3 \\ 0.16 & 0.64 & 0.16 & -0.32 \\ 0.2 & 0.16 & 1 & -0.4 \\ -0.3 & -0.32 & -0.4 & 1 \end{bmatrix}, \mu_{ZW} = (-0.6, -1.08, 0.08) \text{ and } ZW = \begin{bmatrix} 0.1 & 0 & 0 \\ 0 & 0.1 & -0.09 \\ 0 & -0.09 & 0.25 \end{bmatrix}$$

$\text{Cov}(U_2 - U_1, V) = -0.02$      $\text{Cov}(U_3 - U_2, V) = -0.08$   
 $\beta_1 = 0$ ;  $\beta_2 = 0.025$ ;  $\beta_3 = 0.3$ ;  $\gamma = 1$

## IV Estimates and Their Components\*

Parameter	Value
$\Delta^{IVZ}$	0.1489
$\Delta_{12}^{IVZ}$	0.0117
$\Delta_{23}^{IVZ}$	0.1372
$\Delta^{IVW_1}$	0.0017
$\Delta_{12}^{IVW_1}$	0.0325
$\Delta_{23}^{IVW_1}$	-0.0308

## Treatment Parameters and Their Values

Parameter	Value
$ATE_{12} = E(Y_2 - Y_1)$	0.025
$ATE_{23} = E(Y_3 - Y_2)$	0.275
$TT_{12} = E(Y_2 - Y_1   D_2 = 1)$	0.0271
$TT_{23} = E(Y_3 - Y_2   D_3 = 1)$	0.1871
$TUT_{12} = E(Y_2 - Y_1   D_1 = 1)$	0.0047
$TUT_{23} = E(Y_3 - Y_2   D_2 = 1)$	0.2854

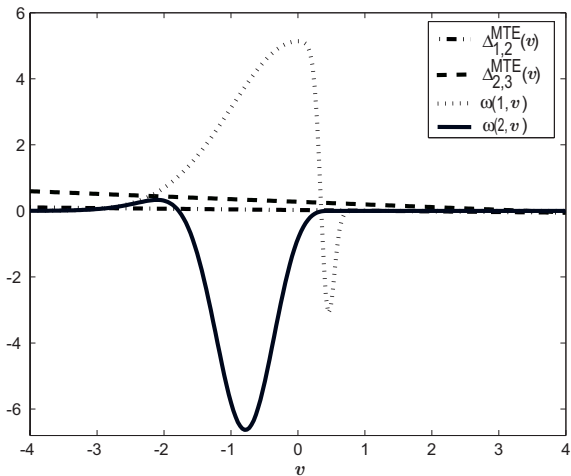
\*  $\Delta^{IVZ}$  is decomposed as:

$$\Delta^{IVZ} = \int E(Y_2 - Y_1 | V = v) \omega^Z(1, v) f_V(v) dv + \int E(Y_3 - Y_2 | V = v) \omega^Z(2, v) f_V(v) dv = \Delta_{12}^{IVZ} + \Delta_{23}^{IVZ}$$

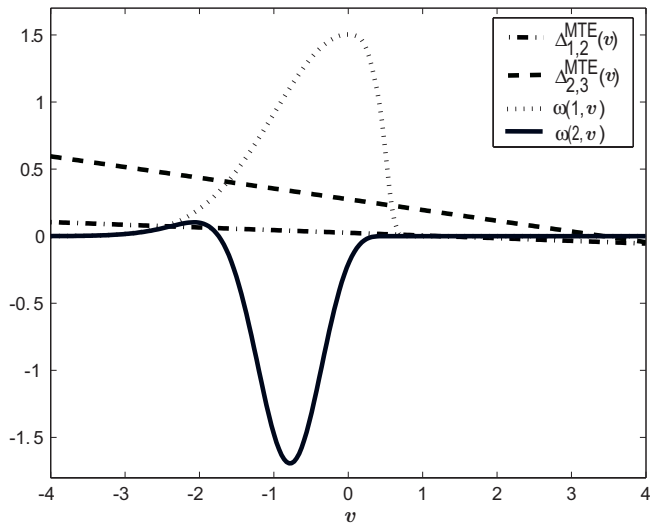
An analogous decomposition applies to  $\Delta^{IVW_1}$ .

Source: Heckman, Urzua and Vytlačil (2006)

Figure 18: A. Treatment Parameters and IV – The Generalized Ordered Choice Roy Model under Normality: Case II:  $Z$  as Instrument



## B. $W_1$ as Instrument



Outcomes

$$\begin{aligned} Y_1 &= \alpha + \beta_1 + U_1 \\ Y_2 &= \alpha + \beta_2 + U_2 \\ Y_3 &= \alpha + \beta_3 + U_3 \end{aligned}$$

Choice Model

$$\begin{aligned} D_s &= \mathbf{1}[W_{s-1} < \gamma Z - V \leq W_s] \\ s &= 1, 2, 3 \end{aligned}$$

Parameterization

$$(U_1, U_2, U_3, V) \sim N(\mathbf{0}, \Sigma_{UV}), \quad (Z, W_1, W_2) \sim N(\mu_{ZW}, \Sigma_{ZW}) \quad \text{and} \quad W_0 = -\infty; W_3 = \infty.$$

$$\Sigma_{UV} = \begin{bmatrix} 1 & 0.16 & 0.2 & -0.3 \\ 0.16 & 0.64 & 0.16 & -0.32 \\ 0.2 & 0.16 & 1 & -0.4 \\ -0.3 & -0.32 & -0.4 & 1 \end{bmatrix}, \quad \mu_{ZW} = (-0.6, -1.08, 0.08) \quad \text{and} \quad \Sigma_{ZW} = \begin{bmatrix} 0.1 & 0.092 & -0.036 \\ 0.092 & 0.1 & -0.09 \\ -0.036 & -0.09 & 0.25 \end{bmatrix}$$

$$\begin{aligned} \text{Cov}(U_2 - U_1, V) &= -0.02 & \text{Cov}(U_3 - U_2, V) &= -0.08 \\ \beta_1 &= 0; \beta_2 = 0.025; \beta_3 = 0.3; \gamma = 1 \end{aligned}$$

### IV Estimates and Their Components\*

Parameter	Value
$\Delta^{IVZ}$	-1.8091
$\Delta_{12}^{IVZ}$	0.2866
$\Delta_{23}^{IVZ}$	-2.0957
$\Delta^{IVW_1}$	-0.4284
$\Delta_{12}^{IVW_1}$	0.0909
$\Delta_{23}^{IVW_1}$	-0.5193

### Treatment Parameters and Their Values

Parameter	Value
$ATE_{12} = E(Y_2 - Y_1)$	0.025
$ATE_{23} = E(Y_3 - Y_2)$	0.275
$TT_{12} = E(Y_2 - Y_1   D_2 = 1)$	0.0283
$TT_{23} = E(Y_3 - Y_2   D_3 = 1)$	0.1754
$TUT_{12} = E(Y_2 - Y_1   D_1 = 1)$	0.0025
$TUT_{23} = E(Y_3 - Y_2   D_2 = 1)$	0.2898

\* See the footnote below Figure 16 for details of the decomposition of  $\Delta^{IVZ}$  and  $\Delta^{IVW_1}$ .  
 Source: Heckman, Urzua and Vytlačil (2006)

- This produces negative weights for  $Z_1$  for both transitions and a negative weight for  $W_1$  in the second transition.



- This produces negative weights for  $Z_1$  for both transitions and a negative weight for  $W_1$  in the second transition.
- For both instruments, IV is negative even though both MTEs are positive throughout most of their range.

- This produces negative weights for  $Z_1$  for both transitions and a negative weight for  $W_1$  in the second transition.
- For both instruments, IV is negative even though both MTEs are positive throughout most of their range.
- IV provides a misleading summary of the underlying marginal treatment effects.

- This produces negative weights for  $Z_1$  for both transitions and a negative weight for  $W_1$  in the second transition.
- For both instruments, IV is negative even though both MTEs are positive throughout most of their range.
- IV provides a misleading summary of the underlying marginal treatment effects.
- Comparing figures 16–18, it is important to recall that all are based on the same structural model.

- This produces negative weights for  $Z_1$  for both transitions and a negative weight for  $W_1$  in the second transition.
- For both instruments, IV is negative even though both MTEs are positive throughout most of their range.
- IV provides a misleading summary of the underlying marginal treatment effects.
- Comparing figures 16–18, it is important to recall that all are based on the same structural model.
- All have the same MTE and average treatment effects.

- This produces negative weights for  $Z_1$  for both transitions and a negative weight for  $W_1$  in the second transition.
- For both instruments, IV is negative even though both MTEs are positive throughout most of their range.
- IV provides a misleading summary of the underlying marginal treatment effects.
- Comparing figures 16–18, it is important to recall that all are based on the same structural model.
- All have the same MTE and average treatment effects.
- But the IV estimates are very different solely as a consequence of the differences in the distributions of instruments across the examples.

- An alternative way to benchmark what IV estimates in the ordered choice model is to compare IV estimates to the PRTE for well defined policy experiments.

- An alternative way to benchmark what IV estimates in the ordered choice model is to compare IV estimates to the PRTE for well defined policy experiments.
- We consider two such experiments, corresponding to proportional and fixed subsidies for attending different levels of schooling.

- An alternative way to benchmark what IV estimates in the ordered choice model is to compare IV estimates to the PRTE for well defined policy experiments.
- We consider two such experiments, corresponding to proportional and fixed subsidies for attending different levels of schooling.
- We use the definition of the PRTE given in equation (46).



- An alternative way to benchmark what IV estimates in the ordered choice model is to compare IV estimates to the PRTE for well defined policy experiments.
- We consider two such experiments, corresponding to proportional and fixed subsidies for attending different levels of schooling.
- We use the definition of the PRTE given in equation (46).
- The baseline model is the one used to generate figure 17.

- An alternative way to benchmark what IV estimates in the ordered choice model is to compare IV estimates to the PRTE for well defined policy experiments.
- We consider two such experiments, corresponding to proportional and fixed subsidies for attending different levels of schooling.
- We use the definition of the PRTE given in equation (46).
- The baseline model is the one used to generate figure 17.
- The weights can be constructed from data and are derived in Appendix, Slide 1139.

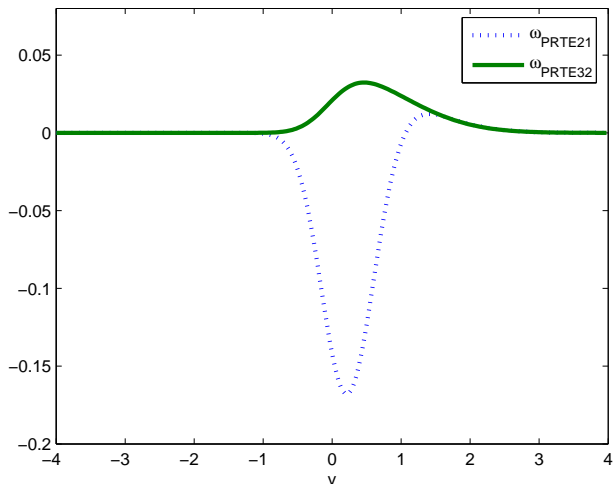
- Figure 19 plots the weights for the PRTE for each transition for a policy experiment.

- Figure 19 plots the weights for the PRTE for each transition for a policy experiment.
- We change the economy from the benchmark economy that generates figure 17 to an economy where  $W_2$  is subsidized by a proportional amount  $\tau$ .

- Figure 19 plots the weights for the PRTE for each transition for a policy experiment.
- We change the economy from the benchmark economy that generates figure 17 to an economy where  $W_2$  is subsidized by a proportional amount  $\tau$ .
- The PRTE weights for each transition are negative over certain intervals.

- Figure 19 plots the weights for the PRTE for each transition for a policy experiment.
- We change the economy from the benchmark economy that generates figure 17 to an economy where  $W_2$  is subsidized by a proportional amount  $\tau$ .
- The PRTE weights for each transition are negative over certain intervals.
- The overall PRTE is close to zero and can be decomposed into two components corresponding to a negative component on the second transition.

Figure 19: The Policy Relevant Treatment Effect Weights – The Generalized Ordered Choice Roy Model under Normality



Outcomes

$$\begin{aligned} Y_1 &= \alpha + \beta_1 + U_1 \\ Y_2 &= \alpha + \beta_2 + U_2 \\ Y_3 &= \alpha + \beta_3 + U_3 \end{aligned}$$

Choice Model

$$D_s = \mathbf{1}[W_{s-1} < \gamma Z - V \leq W_s]$$

$s = 1, 2, 3$

Parameterization

The benchmark model (regime  $p$ ) is the same as the one presented below Figure 17.

Under the new regime (regime  $p'$ ) we define  $W_1^{p'} = W_1^p(1 - \tau)$  with  $\tau = 0.5$ . Thus, under regime  $p$  we have

$$\mu_{ZW}^{p'} = (-0.6, -0.54, 0.08) \text{ and } \mu_{ZW}^p = \begin{bmatrix} 0.1 & 0 & 0 \\ 0 & 0.025 & -0.045 \\ 0 & -0.045 & 0.25 \end{bmatrix}$$

The other parameters remain at the values set under the regime  $p$ .



PRTE Estimates and their Components<sup>1</sup>

Parameter	Value
$PRTE^{p',p}$	0.0076
$PRTE_{21}^{p',p}$	-0.0032
$PRTE_{32}^{p',p}$	0.0109

<sup>1</sup>  $PRTE^{p',p}$  is decomposed as:

$$PRTE^{p',p} = \int E(Y_2 - Y_1 | V = v) \omega^{p',p}(1, v) f_V(v) dv$$

$$+ \int E(Y_3 - Y_2 | V = v) \omega^{p',p}(2, v) f_V(v) dv = PRTE_{21}^{p',p} + PRTE_{32}^{p',p}.$$

IV Estimates and Treatment Parameters under Different Regimes<sup>2</sup>

Parameter	Regime $p$	Regime $p'$
$IV^Z$	0.1489	0.1521
$IV_{12}^Z$	0.0117	0.0174
$IV_{23}^Z$	0.1372	0.1347
$IV^{W_1}$	0.0017	0.0804
$IV_{12}^{W_1}$	0.0325	0.0358
$IV_{23}^{W_1}$	-0.0308	0.0446
$ATE_{12}$	0.025	0.025
$ATE_{23}$	0.275	0.275
$TT_{12}$	0.0271	0.0327
$TT_{23}$	0.1871	0.1789
$TUT_{12}$	0.0047	0.0103
$TUT_{23}$	0.2854	0.3067

<sup>2</sup> See footnote below Figure 16 for details of the decompositions of  $IV^Z$  and  $IV^{W_1}$ .

- The IV for the benchmark regime ( $p$ ) and new regime ( $p'$ ) are given in the bottom table.

- The IV for the benchmark regime ( $p$ ) and new regime ( $p'$ ) are given in the bottom table.
- The IV based on  $Z$  are far from the PRTE parameter.

- The IV for the benchmark regime ( $p$ ) and new regime ( $p'$ ) are given in the bottom table.
- The IV based on  $Z$  are far from the PRTE parameter.
- In general, the IV estimands are far off the mark from the PRTEs.

- We next present a comparison between what IV estimates and the PRTE for a policy that consists of changing  $W_2$  to  $W_2 - t$  ( $t = 1.2$  in the simulations).

- We next present a comparison between what IV estimates and the PRTE for a policy that consists of changing  $W_2$  to  $W_2 - t$  ( $t = 1.2$  in the simulations).
- This can be thought of as a college tuition reduction policy.

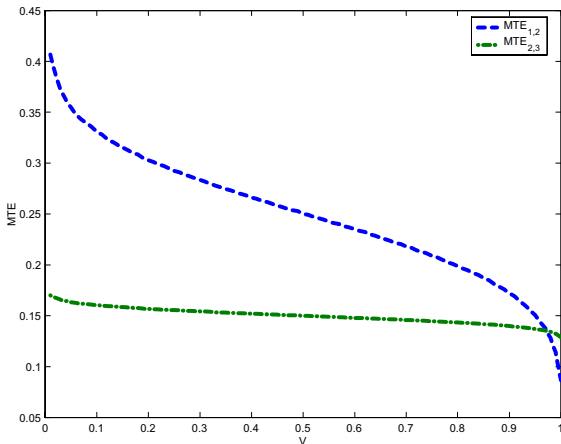
- We next present a comparison between what IV estimates and the PRTE for a policy that consists of changing  $W_2$  to  $W_2 - t$  ( $t = 1.2$  in the simulations).
- This can be thought of as a college tuition reduction policy.
- We compare the weights on PRTE with the weights on IV using  $W_1$  (figure 20) and  $Z$  (figure 21) as instruments.



- We next present a comparison between what IV estimates and the PRTE for a policy that consists of changing  $W_2$  to  $W_2 - t$  ( $t = 1.2$  in the simulations).
- This can be thought of as a college tuition reduction policy.
- We compare the weights on PRTE with the weights on IV using  $W_1$  (figure 20) and  $Z$  (figure 21) as instruments.
- The case using  $W_2$  as an instrument is similar and for the sake of brevity is not discussed.

- We next present a comparison between what IV estimates and the PRTE for a policy that consists of changing  $W_2$  to  $W_2 - t$  ( $t = 1.2$  in the simulations).
- This can be thought of as a college tuition reduction policy.
- We compare the weights on PRTE with the weights on IV using  $W_1$  (figure 20) and  $Z$  (figure 21) as instruments.
- The case using  $W_2$  as an instrument is similar and for the sake of brevity is not discussed.
- In figure 20A, we plot the transition-specific MTE for the values of the model presented at the base of the table.

Figure 20: A.  $W_2 - t$  where  $t = 1.2$  and  $W_1$  is the instrument: Marginal treatment effects by transition



$$Y_3 = \alpha + \beta_3 + U_3;$$

$$Y_2 = \alpha + \beta_2 + U_2;$$

$$Y_1 = \alpha + U_1;$$

$$I = Z - V$$

Sample size = 1500

$$D_3 = 1 \text{ if } W_2 < I < \infty;$$

$$D_2 = 1 \text{ if } W_1 < I \leq W_2;$$

$$D_1 = 1 \text{ if } -\infty < I \leq W_1;$$

$$U_3 = \sigma_3 \tau;$$

$$U_2 = \sigma_2 \tau;$$

$$U_1 = \sigma_1 \tau;$$

$$V = \sigma_V \tau;$$

$$u_D = F_V(V)$$

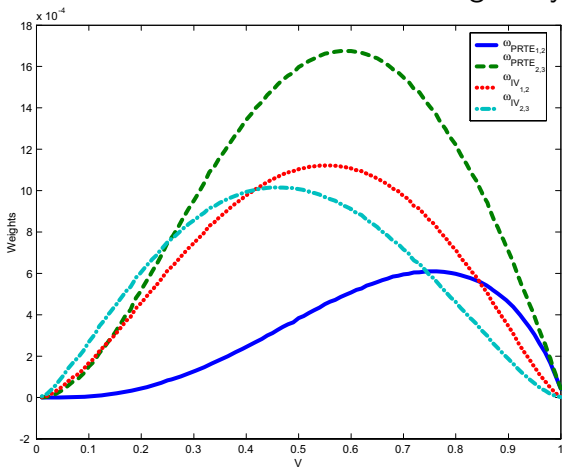
$$\sigma_3 = 0.02, \sigma_2 = 0.012, \sigma_1 = -0.05, \sigma_V = -1$$

$$\alpha = 0.67, \beta_2 = 0.25, \beta_3 = 0.4$$

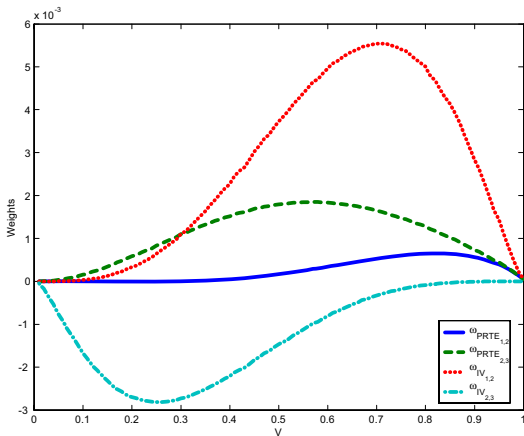
$$Z \sim N(-0.0026, 0.27) \text{ and } Z \perp\!\!\!\perp V$$

$$\tau \sim N(0, 1)$$

B.  $W_2 - t$  where  $t = 1.2$  and  $W_1$  is the instrument: Policy relevant treatment effect vs. instrumental variables weights by transition.



C.  $W_2 - t$  where  $t = 1.2$  and  $W_1$  is the instrument: Policy relevant treatment effect vs. instrumental variables weights by transition.



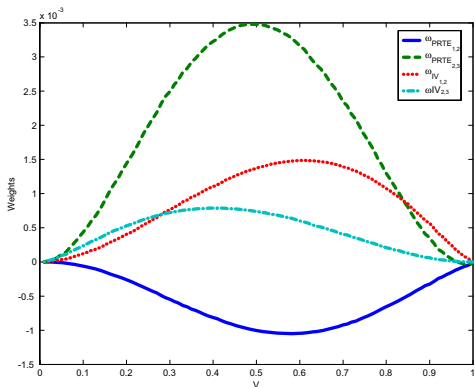
$$(W_1, W_2) \sim N \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix} \right)$$

$$\Delta^{PRTE} = 0.159 \quad IV = 0.296$$

Proportion Induced to Change from  $D_1 = 1$  to  $D_3 = 1$  = 32.1%

Proportion Induced to Change from  $D_2 = 1$  to  $D_3 = 1$  = 64.7%

D.  $W_2 - t$  where  $t = 1.2$  and  $W_1$  is the instrument: Policy relevant treatment effect vs. instrumental variables weights by transition.



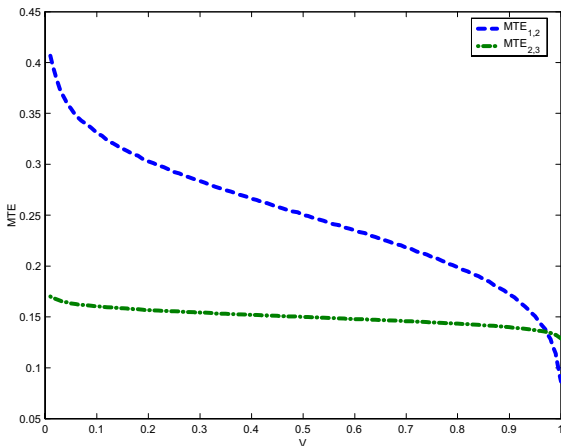
$$(W_1, W_2) \sim N \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right)$$

$$\Delta^{P RTE} = 0.110 \quad IV = 0.210$$

Proportion Induced to Change from  $D_1 = 1$  to  $D_3 = 1$  = 27.5%

Proportion Induced to Change from  $D_2 = 1$  to  $D_3 = 1$  = 76.8%

Figure 21: A.  $W_2 - t$  where  $t = 1.2$  and  $Z$  is the instrument: Marginal treatment effects by transition.



$$Y_3 = \alpha + \beta_3 + U_3;$$

$$Y_2 = \alpha + \beta_2 + U_2;$$

$$Y_1 = \alpha + U_1;$$

$$I = Z - V$$

Sample size = 1500

$$D_3 = 1 \text{ if } W_2 < I < \infty;$$

$$D_2 = 1 \text{ if } W_1 < I \leq W_2;$$

$$D_1 = 1 \text{ if } -\infty < I \leq W_1;$$

$$U_3 = \sigma_3 \tau;$$

$$U_2 = \sigma_2 \tau;$$

$$U_1 = \sigma_1 \tau;$$

$$V = \sigma_V \tau;$$

$$\sigma_3 = 0.02, \sigma_2 = 0.012, \sigma_1 = -0.05, \sigma_V = -1$$

$$\alpha = 0.67, \beta_2 = 0.25, \beta_3 = 0.4$$

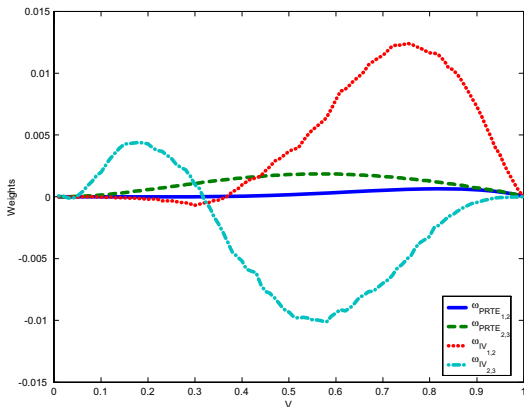
$$Z \sim N(-0.0026, 0.27) \text{ and } Z \perp\!\!\!\perp V$$

$$\tau \sim N(0, 1)$$

B.  $W_2 - t$  where  $t = 1.2$  and  $Z$  is the instrument: Policy relevant treatment effect vs. instrumental variables weights by transition.



C.  $W_2 - t$  where  $t = 1.2$  and  $Z$  is the instrument: Policy relevant treatment effect vs. instrumental variables weights by transition.



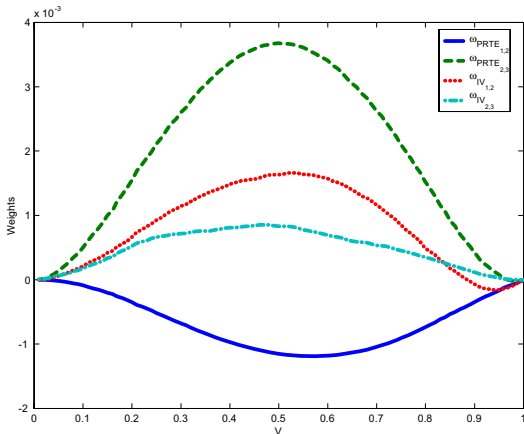
$$(W_1, W_2) \sim N \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix} \right)$$

$$\Delta^{PRTE} = 0.159 \quad IV = 0.346$$

Proportion Induced to Change from  $D_1 = 1$  to  $D_3 = 1$  = 32.1%

Proportion Induced to Change from  $D_2 = 1$  to  $D_3 = 1$  = 64.7%

D.  $W_2 - t$  where  $t = 1.2$  and  $Z$  is the instrument: Policy relevant treatment effect vs. instrumental variables weights by transition.



$$(W_1, W_2) \sim N \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right)$$

$$\Delta^{PRTE} = 0.104 \quad IV = 0.215$$

Proportion Induced to Change from  $D_1 = 1$  to  $D_3 = 1$  = 27.3%

Proportion Induced to Change from  $D_2 = 1$  to  $D_3 = 1$  = 69.3%

- These are identical to the transition-specific MTE plotted in figure 21A.

- These are identical to the transition-specific MTE plotted in figure 21A.
- Both of the  $\Delta^{\text{MTE}}$  parameters have the typical shape of declining returns for people less likely to make the transition, i.e., those who have a higher  $V = v$ .

- These are identical to the transition-specific MTE plotted in figure 21A.
- Both of the  $\Delta^{\text{MTE}}$  parameters have the typical shape of declining returns for people less likely to make the transition, i.e., those who have a higher  $V = v$ .
- Even though the levels are higher for outcomes 2 and 3, the marginal returns are higher for the transition  $1 \rightarrow 2$ .

- These are identical to the transition-specific MTE plotted in figure 21A.
- Both of the  $\Delta^{\text{MTE}}$  parameters have the typical shape of declining returns for people less likely to make the transition, i.e., those who have a higher  $V = v$ .
- Even though the levels are higher for outcomes 2 and 3, the marginal returns are higher for the transition  $1 \rightarrow 2$ .
- Figure 20B plots the policy weights for the two transitions for a policy that lowers  $W_2$  (“reduces tuition”).

- It also plots the IV weights for the two  $\Delta^{\text{MTE}}$  functions for the case where  $W_1$  is the instrument.

- It also plots the IV weights for the two  $\Delta^{\text{MTE}}$  functions for the case where  $W_1$  is the instrument.
- The correlation pattern for  $(W_1, W_2)$  is positive with specific values given below the figure.



- It also plots the IV weights for the two  $\Delta^{\text{MTE}}$  functions for the case where  $W_1$  is the instrument.
- The correlation pattern for  $(W_1, W_2)$  is positive with specific values given below the figure.
- The policy studied in figure 20B shifts 42.8% of the  $D_1 = 1$  people into the category  $D_3 = 1$  and 92.4% of  $D_2$  people into  $D_3$ .

- It also plots the IV weights for the two  $\Delta^{\text{MTE}}$  functions for the case where  $W_1$  is the instrument.
- The correlation pattern for  $(W_1, W_2)$  is positive with specific values given below the figure.
- The policy studied in figure 20B shifts 42.8% of the  $D_1 = 1$  people into the category  $D_3 = 1$  and 92.4% of  $D_2$  people into  $D_3$ .
- In this simulation, the IV weights are positive.

- It also plots the IV weights for the two  $\Delta^{\text{MTE}}$  functions for the case where  $W_1$  is the instrument.
- The correlation pattern for  $(W_1, W_2)$  is positive with specific values given below the figure.
- The policy studied in figure 20B shifts 42.8% of the  $D_1 = 1$  people into the category  $D_3 = 1$  and 92.4% of  $D_2$  people into  $D_3$ .
- In this simulation, the IV weights are positive.
- The IV weights and  $\Delta^{\text{PRTE}}$  weights are distinctly different and the IV estimate is 0.201 vs.  $\Delta^{\text{PRTE}}$  of 0.166.

- When we change the correlation structure between  $W_1$  and  $W_2$  so that they are negatively correlated (figure 20C), the IV weight for  $\Delta_{2,3}^{\text{MTE}}$  becomes *negative* while that for  $\Delta_{1,2}^{\text{MTE}}$  remains positive.

- When we change the correlation structure between  $W_1$  and  $W_2$  so that they are negatively correlated (figure 20C), the IV weight for  $\Delta_{2,3}^{\text{MTE}}$  becomes *negative* while that for  $\Delta_{1,2}^{\text{MTE}}$  remains positive.
- The contrast in these figures between negative and positive IV weights depends on the correlation structure between  $W_1$  and  $W_2$ .

- When we change the correlation structure between  $W_1$  and  $W_2$  so that they are negatively correlated (figure 20C), the IV weight for  $\Delta_{2,3}^{\text{MTE}}$  becomes *negative* while that for  $\Delta_{1,2}^{\text{MTE}}$  remains positive.
- The contrast in these figures between negative and positive IV weights depends on the correlation structure between  $W_1$  and  $W_2$ .
- The stochastic order ( $W_2 > W_1$ ) is a force toward positive weights, which can be undone when the dependence induced by the density ( $f(w_1, w_2)$ ) is sufficiently negative.

- When we change the correlation structure between  $W_1$  and  $W_2$  so that they are negatively correlated (figure 20C), the IV weight for  $\Delta_{2,3}^{\text{MTE}}$  becomes *negative* while that for  $\Delta_{1,2}^{\text{MTE}}$  remains positive.
- The contrast in these figures between negative and positive IV weights depends on the correlation structure between  $W_1$  and  $W_2$ .
- The stochastic order ( $W_2 > W_1$ ) is a force toward positive weights, which can be undone when the dependence induced by the density ( $f(w_1, w_2)$ ) is sufficiently negative.
- The discord between the IV and  $\Delta^{\text{PRTE}}$  weights is substantial and is reflected in the estimates ( $\Delta^{\text{PRTE}} = 0.159$  vs.  $\Delta^{\text{IV}} = 0.296$ ).

- As figure 20D illustrates, the weights on  $\Delta^{\text{PRTE}}$  are not guaranteed to be positive either.



- As figure 20D illustrates, the weights on  $\Delta^{\text{PRTE}}$  are not guaranteed to be positive either.
- Thus neither the IV weights nor the weights on  $\Delta^{\text{PRTE}}$  are guaranteed to be positive or negative and the relationship between the two sets of weights can be quite weak.

- As figure 20D illustrates, the weights on  $\Delta^{\text{PRTE}}$  are not guaranteed to be positive either.
- Thus neither the IV weights nor the weights on  $\Delta^{\text{PRTE}}$  are guaranteed to be positive or negative and the relationship between the two sets of weights can be quite weak.
- Figures 21A–21D present a parallel set of simulations when  $Z$  is used as an instrument.

- As figure 20D illustrates, the weights on  $\Delta^{\text{PRTE}}$  are not guaranteed to be positive either.
- Thus neither the IV weights nor the weights on  $\Delta^{\text{PRTE}}$  are guaranteed to be positive or negative and the relationship between the two sets of weights can be quite weak.
- Figures 21A–21D present a parallel set of simulations when  $Z$  is used as an instrument.
- Changes in  $Z$  shift persons across all transitions whereas  $W_1$  is a transition-specific shifter.

- Figure 21 reproduces the policy invariant  $\Delta^{\text{MTE}}$  parameters from figure 20A.

- Figure 21 reproduces the policy invariant  $\Delta^{\text{MTE}}$  parameters from figure 20A.
- Figure 21B shows that the IV weights for  $\Delta_{1,2}^{\text{MTE}}$  assume both positive and negative values.

- Figure 21 reproduces the policy invariant  $\Delta^{\text{MTE}}$  parameters from figure 20A.
- Figure 21B shows that the IV weights for  $\Delta_{1,2}^{\text{MTE}}$  assume both positive and negative values.
- The IV weights for  $\Delta_{2,3}^{\text{MTE}}$  are positive but not monotonic.

- Figure 21 reproduces the policy invariant  $\Delta^{\text{MTE}}$  parameters from figure 20A.
- Figure 21B shows that the IV weights for  $\Delta_{1,2}^{\text{MTE}}$  assume both positive and negative values.
- The IV weights for  $\Delta_{2,3}^{\text{MTE}}$  are positive but not monotonic.
- In figure 21C, where there is negative dependence between  $W_1$  and  $W_2$ , both sets of IV weights assume both positive and negative values.

- Figure 21 reproduces the policy invariant  $\Delta^{\text{MTE}}$  parameters from figure 20A.
- Figure 21B shows that the IV weights for  $\Delta_{1,2}^{\text{MTE}}$  assume both positive and negative values.
- The IV weights for  $\Delta_{2,3}^{\text{MTE}}$  are positive but not monotonic.
- In figure 21C, where there is negative dependence between  $W_1$  and  $W_2$ , both sets of IV weights assume both positive and negative values.
- In the case where  $f(w_1, w_2) = f_1(w_1)f_2(w_2)$ , the weights on  $\Delta_{1,2}^{\text{MTE}}$  for  $\Delta^{\text{PRTE}}$  are negative.



- These simulations show a rich variety of shapes and signs for the weights.

- These simulations show a rich variety of shapes and signs for the weights.
- They illustrate a main point of this chapter—that standard IV methods are not guaranteed to weight marginal treatment effects positively or to produce estimates close to policy relevant treatment effects or even to produce any gross treatment effect.

- These simulations show a rich variety of shapes and signs for the weights.
- They illustrate a main point of this chapter—that standard IV methods are not guaranteed to weight marginal treatment effects positively or to produce estimates close to policy relevant treatment effects or even to produce any gross treatment effect.
- Estimators based on LIV and its extension to the ordered model (44) identify  $\Delta^{\text{MTE}}$  for each transition and answer policy relevant questions.

- These simulations show a rich variety of shapes and signs for the weights.
- They illustrate a main point of this chapter—that standard IV methods are not guaranteed to weight marginal treatment effects positively or to produce estimates close to policy relevant treatment effects or even to produce any gross treatment effect.
- Estimators based on LIV and its extension to the ordered model (44) identify  $\Delta^{\text{MTE}}$  for each transition and answer policy relevant questions.
- We now turn to an analysis of a general unordered model.

## Extension to Multiple Treatments that are Unordered

- The previous section analyzes a multiple treatment model where the treatment choice equation is an ordered choice model.

## Extension to Multiple Treatments that are Unordered

- The previous section analyzes a multiple treatment model where the treatment choice equation is an ordered choice model.
- In this section, we develop a framework for the analysis of multiple treatments when the choice equation is a nonparametric version of the classical multinomial choice model with no order imposed.

## Extension to Multiple Treatments that are Unordered

- The previous section analyzes a multiple treatment model where the treatment choice equation is an ordered choice model.
- In this section, we develop a framework for the analysis of multiple treatments when the choice equation is a nonparametric version of the classical multinomial choice model with no order imposed.
- Appendix B of Part I, and ? analyze nonparametric and semiparametric identification of discrete choice models.

- With this framework, treatment effects can be defined as the difference in the counterfactual outcomes that would have been observed if the agent faced different general choice sets, i.e., the effect of the individual being forced to choose from one choice set instead of another.



- With this framework, treatment effects can be defined as the difference in the counterfactual outcomes that would have been observed if the agent faced different general choice sets, i.e., the effect of the individual being forced to choose from one choice set instead of another.
- We define treatment parameters for a general multiple treatment problem and present conditions for the application of instrumental variables for identifying a variety of new treatment parameters.

- With this framework, treatment effects can be defined as the difference in the counterfactual outcomes that would have been observed if the agent faced different general choice sets, i.e., the effect of the individual being forced to choose from one choice set instead of another.
- We define treatment parameters for a general multiple treatment problem and present conditions for the application of instrumental variables for identifying a variety of new treatment parameters.
- Our identification conditions are weaker than the ones used in appendix B of Part I, which establishes conditions under which it is possible to nonparametrically identify a full multinomial selection model.

- Our use of choice theory is a unique aspect of our approach to the analysis of treatment effects.

- Our use of choice theory is a unique aspect of our approach to the analysis of treatment effects.
- One particularly helpful result we draw on is the representation of the multinomial choices in terms of the choice between a particular choice and the best option among all other choices.

- Our use of choice theory is a unique aspect of our approach to the analysis of treatment effects.
- One particularly helpful result we draw on is the representation of the multinomial choices in terms of the choice between a particular choice and the best option among all other choices.
- This representation is crucial for understanding why LIV allows one to identify the MTE for the effect of one choice versus the best alternative option.

- Our use of choice theory is a unique aspect of our approach to the analysis of treatment effects.
- One particularly helpful result we draw on is the representation of the multinomial choices in terms of the choice between a particular choice and the best option among all other choices.
- This representation is crucial for understanding why LIV allows one to identify the MTE for the effect of one choice versus the best alternative option.
- The representation was introduced in ?, and has been used in the analysis of parametric multinomial selection models by ? and ?.

- Unlike those authors, we systematically explore treatment effect heterogeneity, consider nonparametric identification, and examine the application of the LIV methodology to such models.

- Unlike those authors, we systematically explore treatment effect heterogeneity, consider nonparametric identification, and examine the application of the LIV methodology to such models.
- Our analysis proceeds as follows.



- Unlike those authors, we systematically explore treatment effect heterogeneity, consider nonparametric identification, and examine the application of the LIV methodology to such models.
- Our analysis proceeds as follows.
- We first introduce our nonparametric, multinomial selection model and state our assumptions in Slide 563.

- Unlike those authors, we systematically explore treatment effect heterogeneity, consider nonparametric identification, and examine the application of the LIV methodology to such models.
- Our analysis proceeds as follows.
- We first introduce our nonparametric, multinomial selection model and state our assumptions in Slide 563.
- In Slide 576, we define treatment effects in a general unordered model as the differences in the counterfactual outcomes that would have been observed if the agent faced different choice sets, i.e., the effects observed if individuals are forced to choose from one choice set instead of another.

- Unlike those authors, we systematically explore treatment effect heterogeneity, consider nonparametric identification, and examine the application of the LIV methodology to such models.
- Our analysis proceeds as follows.
- We first introduce our nonparametric, multinomial selection model and state our assumptions in Slide 563.
- In Slide 576, we define treatment effects in a general unordered model as the differences in the counterfactual outcomes that would have been observed if the agent faced different choice sets, i.e., the effects observed if individuals are forced to choose from one choice set instead of another.
- We also define the corresponding treatment parameters.

- Treatment effects in this context exhibit a form of treatment effect heterogeneity not present in the binary treatment case.

- Treatment effects in this context exhibit a form of treatment effect heterogeneity not present in the binary treatment case.
- The new form of heterogeneity arises from agents facing different choice sets, which we discuss in Slide 589.

- Treatment effects in this context exhibit a form of treatment effect heterogeneity not present in the binary treatment case.
- The new form of heterogeneity arises from agents facing different choice sets, which we discuss in Slide 589.
- Slide 596 establishes that LIV and the nonparametric Wald-IV estimand produce identification of the MTE/LATE versions of the effect of one choice versus the best alternative option without requiring knowledge of the latent index functions generating choices or large support assumptions.

- Treatment effects in this context exhibit a form of treatment effect heterogeneity not present in the binary treatment case.
- The new form of heterogeneity arises from agents facing different choice sets, which we discuss in Slide 589.
- Slide 596 establishes that LIV and the nonparametric Wald-IV estimand produce identification of the MTE/LATE versions of the effect of one choice versus the best alternative option without requiring knowledge of the latent index functions generating choices or large support assumptions.
- Mean treatment effects comparing one option versus the best alternative are the easiest treatment effects to study using instrumental variable methods because we effectively collapse a multiple outcome model to a series of two-outcome models, picking one outcome relative to the rest.

- In Slide 621, we consider a more general case and state conditions for identifying the mean effect of the outcome associated with the best option in one choice set to the mean effect of the best option not in that choice set.



- In Slide 621, we consider a more general case and state conditions for identifying the mean effect of the outcome associated with the best option in one choice set to the mean effect of the best option not in that choice set.
- We show that identification of the corresponding MTE/LATE parameters requires knowledge of the latent index functions of the multinomial choice model.

- In Slide 621, we consider a more general case and state conditions for identifying the mean effect of the outcome associated with the best option in one choice set to the mean effect of the best option not in that choice set.
- We show that identification of the corresponding MTE/LATE parameters requires knowledge of the latent index functions of the multinomial choice model.
- Thus, to identify the parameters by using IV or LIV requires the formulation and estimation of an explicit choice model.

- In Slide 631, we analyze the identification of treatment parameters corresponding to the mean effect of one specified choice versus another specified choice.

- In Slide 631, we analyze the identification of treatment parameters corresponding to the mean effect of one specified choice versus another specified choice.
- Identification of marginal treatment parameters in this case requires the use of identification at infinity arguments relying on large support assumptions, but does not require knowledge of the latent index functions of the multinomial choice problem.

- In Slide 631, we analyze the identification of treatment parameters corresponding to the mean effect of one specified choice versus another specified choice.
- Identification of marginal treatment parameters in this case requires the use of identification at infinity arguments relying on large support assumptions, but does not require knowledge of the latent index functions of the multinomial choice problem.
- This use of large support assumptions is closely related to the need for large support assumptions to identify the full model developed in Appendix B of Part I.

- In Slide 631, we analyze the identification of treatment parameters corresponding to the mean effect of one specified choice versus another specified choice.
- Identification of marginal treatment parameters in this case requires the use of identification at infinity arguments relying on large support assumptions, but does not require knowledge of the latent index functions of the multinomial choice problem.
- This use of large support assumptions is closely related to the need for large support assumptions to identify the full model developed in Appendix B of Part I.
- We summarize our analysis in Slide 645.

## Model and Assumptions

- Consider the following model with multiple choices and multiple outcome states for a general unordered model.

## Model and Assumptions

- Consider the following model with multiple choices and multiple outcome states for a general unordered model.
- Let  $\mathcal{J}$  denote the agent's choice set, where  $\mathcal{J}$  contains a finite number of elements.



## Model and Assumptions

- Consider the following model with multiple choices and multiple outcome states for a general unordered model.
- Let  $\mathcal{J}$  denote the agent's choice set, where  $\mathcal{J}$  contains a finite number of elements.
- The value to the agent of choosing option  $j \in \mathcal{J}$  is

$$R_j(Z_j) = \vartheta_j(Z_j) - V_j, \quad (51)$$

where  $Z_j$  are the agent's observed characteristics that affect the utility from choosing choice  $j$ , and  $V_j$  is the unobserved shock to the agent's utility from choice  $j$ .

- We will sometimes suppress the argument and write  $R_j$  for  $R_j(Z_j)$ .

- We will sometimes suppress the argument and write  $R_j$  for  $R_j(Z_j)$ .
- Let  $Z$  denote the random vector containing all unique elements of  $\{Z_j\}_{j \in \mathcal{J}}$ , i.e.,  $Z = \cup_{j \in \mathcal{J}} \{Z_j\}_{j \in \mathcal{J}}$ .

- We will sometimes suppress the argument and write  $R_j$  for  $R_j(Z_j)$ .
- Let  $Z$  denote the random vector containing all unique elements of  $\{Z_j\}_{j \in \mathcal{J}}$ , i.e.,  $Z = \cup_{j \in \mathcal{J}} \{Z_j\}_{j \in \mathcal{J}}$ .
- We will also sometimes write  $R_j(Z)$  for  $R_j(Z_j)$ , leaving implicit that  $R_j(\cdot)$  only depends on those elements of  $Z$  that are contained in  $Z_j$ .

- We will sometimes suppress the argument and write  $R_j$  for  $R_j(Z_j)$ .
- Let  $Z$  denote the random vector containing all unique elements of  $\{Z_j\}_{j \in \mathcal{J}}$ , i.e.,  $Z = \cup_{j \in \mathcal{J}} \{Z_j\}_{j \in \mathcal{J}}$ .
- We will also sometimes write  $R_j(Z)$  for  $R_j(Z_j)$ , leaving implicit that  $R_j(\cdot)$  only depends on those elements of  $Z$  that are contained in  $Z_j$ .
- Let  $D_{\mathcal{J}j}$  be an indicator variable for whether the agent would choose option  $j$  if confronted with choice set  $\mathcal{J}$ :

$$D_{\mathcal{J}j} = \begin{cases} 1 & \text{if } R_j \geq R_k \quad \forall k \in \mathcal{J} \\ 0 & \text{otherwise.} \end{cases}$$

- Let  $I_{\mathcal{J}}$  denote the choice that would be made by the agent if confronted with choice set  $\mathcal{J}$ :

$$I_{\mathcal{J}} = j \iff D_{\mathcal{J},j} = 1.$$

- Let  $I_{\mathcal{J}}$  denote the choice that would be made by the agent if confronted with choice set  $\mathcal{J}$ :

$$I_{\mathcal{J}} = j \iff D_{\mathcal{J},j} = 1.$$

- Let  $Y_{\mathcal{J}}$  be the outcome variable that would be observed if the agent faced choice set  $\mathcal{J}$ , determined by

$$Y_{\mathcal{J}} = \sum_{j \in \mathcal{J}} D_{\mathcal{J},j} Y_j,$$

where  $Y_j$  is the potential outcome, observed only if option  $j$  is chosen.

- Let  $I_{\mathcal{J}}$  denote the choice that would be made by the agent if confronted with choice set  $\mathcal{J}$ :

$$I_{\mathcal{J}} = j \iff D_{\mathcal{J},j} = 1.$$

- Let  $Y_{\mathcal{J}}$  be the outcome variable that would be observed if the agent faced choice set  $\mathcal{J}$ , determined by

$$Y_{\mathcal{J}} = \sum_{j \in \mathcal{J}} D_{\mathcal{J},j} Y_j,$$

where  $Y_j$  is the potential outcome, observed only if option  $j$  is chosen.

- $Y_j$  is determined by

$$Y_j = \mu_j(X_j, U_j),$$

where  $X_j$  is a vector of the agent's observed characteristics and  $U_j$  is an unobserved random vector.



- Let  $X$  denote the random vector containing all unique elements of  $\{X_j\}_{j \in \mathcal{J}}$ , i.e.,  $X = \cup_{j \in \mathcal{J}} \{X_j\}_{j \in \mathcal{J}}$ .

- Let  $X$  denote the random vector containing all unique elements of  $\{X_j\}_{j \in \mathcal{J}}$ , i.e.,  $X = \cup_{j \in \mathcal{J}} \{X_j\}_{j \in \mathcal{J}}$ .
- $(Z, X, I_{\mathcal{J}}, Y_{\mathcal{J}})$  is assumed to be observed.

- Let  $X$  denote the random vector containing all unique elements of  $\{X_j\}_{j \in \mathcal{J}}$ , i.e.,  $X = \cup_{j \in \mathcal{J}} \{X_j\}_{j \in \mathcal{J}}$ .
- $(Z, X, I_{\mathcal{J}}, Y_{\mathcal{J}})$  is assumed to be observed.
- Define  $R_{\mathcal{J}}$  as the maximum obtainable value given choice set  $\mathcal{J}$ :

$$\begin{aligned} R_{\mathcal{J}} &= \max_{j \in \mathcal{J}} \{R_j\} \\ &= \sum_{j \in \mathcal{J}} D_{\mathcal{J},j} R_j. \end{aligned}$$

- Let  $X$  denote the random vector containing all unique elements of  $\{X_j\}_{j \in \mathcal{J}}$ , i.e.,  $X = \cup_{j \in \mathcal{J}} \{X_j\}_{j \in \mathcal{J}}$ .
- $(Z, X, I_{\mathcal{J}}, Y_{\mathcal{J}})$  is assumed to be observed.
- Define  $R_{\mathcal{J}}$  as the maximum obtainable value given choice set  $\mathcal{J}$ :

$$\begin{aligned} R_{\mathcal{J}} &= \max_{j \in \mathcal{J}} \{R_j\} \\ &= \sum_{j \in \mathcal{J}} D_{\mathcal{J},j} R_j. \end{aligned}$$

- We thus obtain the traditional representation of the decision process that choice  $j$  being optimal implies that choice  $j$  is better than the “next best” option:

$$I_{\mathcal{J}} = j \iff R_j \geq R_{\mathcal{J} \setminus j}.$$

- More generally, a choice from  $\mathcal{K}$  being optimal is equivalent to the highest value obtainable from choices in  $\mathcal{K}$  being higher than the highest value that can be obtained from choices outside that set,

$$I_{\mathcal{J}} \in \mathcal{K} \iff R_{\mathcal{K}} \geq R_{\mathcal{J} \setminus \mathcal{K}}.$$

- More generally, a choice from  $\mathcal{K}$  being optimal is equivalent to the highest value obtainable from choices in  $\mathcal{K}$  being higher than the highest value that can be obtained from choices outside that set,

$$I_{\mathcal{J}} \in \mathcal{K} \iff R_{\mathcal{K}} \geq R_{\mathcal{J} \setminus \mathcal{K}}.$$

- As we will show, this simple representation is the key intuition for understanding how nonparametric instrumental variables estimate the effect of a given choice versus the “next best” alternative.

- Analogous to our definition of  $R_{\mathcal{J}}$ , we define  $R_{\mathcal{J}}(z)$  to be the maximum obtainable value given choice set  $\mathcal{J}$  when instruments are fixed at  $Z = z$ ,

$$R_{\mathcal{J}}(z) = \max_{j \in \mathcal{J}} \{R_j(z)\}.$$

- Analogous to our definition of  $R_{\mathcal{J}}$ , we define  $R_{\mathcal{J}}(z)$  to be the maximum obtainable value given choice set  $\mathcal{J}$  when instruments are fixed at  $Z = z$ ,

$$R_{\mathcal{J}}(z) = \max_{j \in \mathcal{J}} \{R_j(z)\}.$$

- Thus, for example, a choice from  $\mathcal{K}$  is optimal when instruments are fixed at  $Z = z$  if  $R_{\mathcal{K}}(z) \geq R_{\mathcal{J} \setminus \mathcal{K}}(z)$ .



- We make the following assumptions, which generalize assumptions (A-1)–(A-5) invoked in ? and later used in ?, as developed in Slide 12.

- We make the following assumptions, which generalize assumptions (A-1)–(A-5) invoked in ? and later used in ?, as developed in Slide 12.
- We present the assumptions in a fashion parallel to (A-1)–(A-5) and (OC-1)–(OC-6).

- We make the following assumptions, which generalize assumptions (A-1)–(A-5) invoked in ? and later used in ?, as developed in Slide 12.
- We present the assumptions in a fashion parallel to (A-1)–(A-5) and (OC-1)–(OC-6).
- For that reason, we present the second assumption, which requires special attention, out of order.

- (B-1)  $\{(V_j, U_j)\}_{j \in \mathcal{J}}$  is independent of  $Z$  conditional on  $X$ .
- (B-3) The distribution of  $(\{V_j\}_{j \in \mathcal{J}})$  is continuous.
- (B-4)  $E(|Y_j|) < \infty$  for all  $j \in \mathcal{J}$ .
- (B-5)  $\Pr(I_{\mathcal{J}} = j | X) > 0$  for all  $j \in \mathcal{J}$ .

- Assumption (B-1) and (B-3) imply that  $R_j \neq R_k$  w.p.1 for  $j \neq k$ , so that  $\operatorname{argmax}\{R_j\}$  is unique w.p.1.

- Assumption (B-1) and (B-3) imply that  $R_j \neq R_k$  w.p.1 for  $j \neq k$ , so that  $\operatorname{argmax}\{R_j\}$  is unique w.p.1.
- Assumption (B-4) is required for the mean treatment parameters to be well defined.

- Assumption (B-1) and (B-3) imply that  $R_j \neq R_k$  w.p.1 for  $j \neq k$ , so that  $\operatorname{argmax}\{R_j\}$  is unique w.p.1.
- Assumption (B-4) is required for the mean treatment parameters to be well defined.
- It allows us to integrate to the limit, which will be a crucial step for all identification analysis.

- Assumption (B-1) and (B-3) imply that  $R_j \neq R_k$  w.p.1 for  $j \neq k$ , so that  $\operatorname{argmax}\{R_j\}$  is unique w.p.1.
- Assumption (B-4) is required for the mean treatment parameters to be well defined.
- It allows us to integrate to the limit, which will be a crucial step for all identification analysis.
- Assumption (B-5) requires that at least some individuals participate in each program for all  $X$ .



- Our definition and analysis of the treatment parameters only require assumptions (B-1) and (B-3)–(B-5).

- Our definition and analysis of the treatment parameters only require assumptions (B-1) and (B-3)–(B-5).
- However, we will also impose an exclusion restriction for our identification analysis.

- Our definition and analysis of the treatment parameters only require assumptions (B-1) and (B-3)–(B-5).
- However, we will also impose an exclusion restriction for our identification analysis.
- Let  $Z^{[j]}$  denote the  $j$ th components of  $Z$  that are in  $Z_j$  but not in  $Z_k$ ,  $k \neq j$ .

- Our definition and analysis of the treatment parameters only require assumptions (B-1) and (B-3)–(B-5).
- However, we will also impose an exclusion restriction for our identification analysis.
- Let  $Z^{[j]}$  denote the  $j$ th components of  $Z$  that are in  $Z_j$  but not in  $Z_k$ ,  $k \neq j$ .
- Let  $Z^{[-j]}$  denote all elements of  $Z$  except for the components in  $Z^{[j]}$ .

- Our definition and analysis of the treatment parameters only require assumptions (B-1) and (B-3)–(B-5).
- However, we will also impose an exclusion restriction for our identification analysis.
- Let  $Z^{[j]}$  denote the  $j$ th components of  $Z$  that are in  $Z_j$  but not in  $Z_k$ ,  $k \neq j$ .
- Let  $Z^{[-j]}$  denote all elements of  $Z$  except for the components in  $Z^{[j]}$ .
- We work with two alternative assumptions for the exclusion restriction.

- Consider

- Consider

(B-2a) For each  $j \in \mathcal{J}$ , there exists at least one element of  $Z$ , say  $Z^{[j]}$ , such that  $Z^{[j]}$  is not an element of  $Z_k$ ,  $k \neq j$ , and such that the distribution of  $\vartheta_j(Z_j)$  conditional on  $(X, Z^{[-j]})$  is nondegenerate,

- Consider

(B-2a) For each  $j \in \mathcal{J}$ , there exists at least one element of  $Z$ , say  $Z^{[j]}$ , such that  $Z^{[j]}$  is not an element of  $Z_k$ ,  $k \neq j$ , and such that the distribution of  $\vartheta_j(Z_j)$  conditional on  $(X, Z^{[-j]})$  is nondegenerate,

- Or consider



- Consider

(B-2a) For each  $j \in \mathcal{J}$ , there exists at least one element of  $Z$ , say  $Z^{[j]}$ , such that  $Z^{[j]}$  is not an element of  $Z_k$ ,  $k \neq j$ , and such that the distribution of  $\vartheta_j(Z_j)$  conditional on  $(X, Z^{[-j]})$  is nondegenerate,

- Or consider

(B-2b) For each  $j \in \mathcal{J}$ , there exists at least one element of  $Z$ , say  $Z^{[j]}$ , such that  $Z^{[j]}$  is not an element of  $Z_k$ ,  $k \neq j$ , and such that the distribution of  $\vartheta_j(Z_j)$  conditional on  $(X, Z^{[-j]})$  is continuous.

- Assumption (B-2a) imposes the requirement that the analyst be able to independently vary the index for the given value function.

- Assumption (B-2a) imposes the requirement that the analyst be able to independently vary the index for the given value function.
- This produces variation that affects only the value of the  $j$ th value function and causes people to enter or exit sector  $j$ .

- Assumption (B-2a) imposes the requirement that the analyst be able to independently vary the index for the given value function.
- This produces variation that affects only the value of the  $j$ th value function and causes people to enter or exit sector  $j$ .
- It imposes an exclusion restriction, that for any  $j \in \mathcal{J}$ ,  $Z$  contains an element such that (i) it is contained in  $Z_j$ ; (ii) it is not contained in any  $Z_k$  for  $k \neq j$ , and (iii)  $\vartheta_j(\cdot)$  is a nontrivial function of that element conditional on all other regressors.

- Assumption (B-2a) imposes the requirement that the analyst be able to independently vary the index for the given value function.
- This produces variation that affects only the value of the  $j$ th value function and causes people to enter or exit sector  $j$ .
- It imposes an exclusion restriction, that for any  $j \in \mathcal{J}$ ,  $Z$  contains an element such that (i) it is contained in  $Z_j$ ; (ii) it is not contained in any  $Z_k$  for  $k \neq j$ , and (iii)  $\vartheta_j(\cdot)$  is a nontrivial function of that element conditional on all other regressors.
- Assumption (B-2b) strengthens (B-2a) by adding a smoothness assumption.

- A necessary condition for (B-2b) is for the excluded variable to have a density with respect to Lebesgue measure conditional on all other regressors and for  $\vartheta_j(\cdot)$  to be a continuous and nontrivial function of the excluded variable.

- A necessary condition for (B-2b) is for the excluded variable to have a density with respect to Lebesgue measure conditional on all other regressors and for  $\vartheta_j(\cdot)$  to be a continuous and nontrivial function of the excluded variable.
- Assumption (B-2a) will be used to identify a generalization of the LATE parameter.

- A necessary condition for (B-2b) is for the excluded variable to have a density with respect to Lebesgue measure conditional on all other regressors and for  $\vartheta_j(\cdot)$  to be a continuous and nontrivial function of the excluded variable.
- Assumption (B-2a) will be used to identify a generalization of the LATE parameter.
- Assumption (B-2b) will be used to identify a generalization of the MTE parameter.



- A necessary condition for (B-2b) is for the excluded variable to have a density with respect to Lebesgue measure conditional on all other regressors and for  $\vartheta_j(\cdot)$  to be a continuous and nontrivial function of the excluded variable.
- Assumption (B-2a) will be used to identify a generalization of the LATE parameter.
- Assumption (B-2b) will be used to identify a generalization of the MTE parameter.
- For certain portions of the analysis, we strengthen (B-2b) to a large support condition, though the large support assumption will not be required for most of our analysis.

- A necessary condition for (B-2b) is for the excluded variable to have a density with respect to Lebesgue measure conditional on all other regressors and for  $\vartheta_j(\cdot)$  to be a continuous and nontrivial function of the excluded variable.
- Assumption (B-2a) will be used to identify a generalization of the LATE parameter.
- Assumption (B-2b) will be used to identify a generalization of the MTE parameter.
- For certain portions of the analysis, we strengthen (B-2b) to a large support condition, though the large support assumption will not be required for most of our analysis.
- Assumptions (B-2a) and (B-2b) mirror (A-2) for the binary choice model and are analogous to (OC-2) and (OC-6) in an ordered choice model.

## Definition of Treatment Effects and Treatment Parameters

- Treatment effects are defined as the difference in the counterfactual outcomes that would have been observed if the agent faced different choice sets.

## Definition of Treatment Effects and Treatment Parameters

- Treatment effects are defined as the difference in the counterfactual outcomes that would have been observed if the agent faced different choice sets.
- For any two choice sets,  $\mathcal{K}, \mathcal{L} \subset \mathcal{J}$ , define

$$\Delta_{\mathcal{K}, \mathcal{L}} = Y_{\mathcal{K}} - Y_{\mathcal{L}}.$$

## Definition of Treatment Effects and Treatment Parameters

- Treatment effects are defined as the difference in the counterfactual outcomes that would have been observed if the agent faced different choice sets.
- For any two choice sets,  $\mathcal{K}, \mathcal{L} \subset \mathcal{J}$ , define

$$\Delta_{\mathcal{K}, \mathcal{L}} = Y_{\mathcal{K}} - Y_{\mathcal{L}}.$$

- This is the effect of the individual being forced to choose from choice set  $\mathcal{K}$  versus choice set  $\mathcal{L}$ .

- The conventional treatment effect is defined as the difference in potential outcomes between two specified states,

$$\Delta_{k,l} = Y_k - Y_l,$$

which is nested within this framework by taking  $\mathcal{K} = \{k\}$ ,  $\mathcal{L} = \{l\}$ .

- The conventional treatment effect is defined as the difference in potential outcomes between two specified states,

$$\Delta_{k,l} = Y_k - Y_l,$$

which is nested within this framework by taking  $\mathcal{K} = \{k\}$ ,  $\mathcal{L} = \{l\}$ .

- It is the effect for the individual of having no choice except to choose state  $l$ .

- The conventional treatment effect is defined as the difference in potential outcomes between two specified states,

$$\Delta_{k,l} = Y_k - Y_l,$$

which is nested within this framework by taking  $\mathcal{K} = \{k\}$ ,  $\mathcal{L} = \{l\}$ .

- It is the effect for the individual of having no choice except to choose state  $l$ .
- $\Delta_{\mathcal{K},\mathcal{L}}$  will be zero for agents who make the same choice when confronted with choice set  $\mathcal{K}$  and choice set  $\mathcal{L}$ .



- Thus,  $I_{\mathcal{K}} = I_{\mathcal{L}}$  implies  $\Delta_{\mathcal{K},\mathcal{L}} = 0$ , and we have

$$\Delta_{\mathcal{K},\mathcal{L}} = \mathbf{1}(I_{\mathcal{L}} \neq I_{\mathcal{K}})\Delta_{\mathcal{K}\setminus\mathcal{L},\mathcal{L}} \quad (52)$$

$$= \mathbf{1}(I_{\mathcal{L}} \neq I_{\mathcal{K}}) \left( \sum_{j \in \mathcal{K} \setminus \mathcal{L}} D_{\mathcal{K},j} \Delta_{j,\mathcal{L}} \right). \quad (53)$$

- Thus,  $I_{\mathcal{K}} = I_{\mathcal{L}}$  implies  $\Delta_{\mathcal{K},\mathcal{L}} = 0$ , and we have

$$\Delta_{\mathcal{K},\mathcal{L}} = \mathbf{1}(I_{\mathcal{L}} \neq I_{\mathcal{K}})\Delta_{\mathcal{K}\setminus\mathcal{L},\mathcal{L}} \quad (52)$$

$$= \mathbf{1}(I_{\mathcal{L}} \neq I_{\mathcal{K}}) \left( \sum_{j \in \mathcal{K} \setminus \mathcal{L}} D_{\mathcal{K},j} \Delta_{j,\mathcal{L}} \right). \quad (53)$$

- Two examples will be of particular importance for our analysis.

- Thus,  $I_{\mathcal{K}} = I_{\mathcal{L}}$  implies  $\Delta_{\mathcal{K},\mathcal{L}} = 0$ , and we have

$$\Delta_{\mathcal{K},\mathcal{L}} = \mathbf{1}(I_{\mathcal{L}} \neq I_{\mathcal{K}}) \Delta_{\mathcal{K} \setminus \mathcal{L}, \mathcal{L}} \quad (52)$$

$$= \mathbf{1}(I_{\mathcal{L}} \neq I_{\mathcal{K}}) \left( \sum_{j \in \mathcal{K} \setminus \mathcal{L}} D_{\mathcal{K},j} \Delta_{j,\mathcal{L}} \right). \quad (53)$$

- Two examples will be of particular importance for our analysis.
- First, consider choice set  $\mathcal{K} = \{k\}$  versus choice set  $\mathcal{L} = \mathcal{J} \setminus \{k\}$ .

- Thus,  $I_{\mathcal{K}} = I_{\mathcal{L}}$  implies  $\Delta_{\mathcal{K},\mathcal{L}} = 0$ , and we have

$$\Delta_{\mathcal{K},\mathcal{L}} = \mathbf{1}(I_{\mathcal{L}} \neq I_{\mathcal{K}}) \Delta_{\mathcal{K} \setminus \mathcal{L}, \mathcal{L}} \quad (52)$$

$$= \mathbf{1}(I_{\mathcal{L}} \neq I_{\mathcal{K}}) \left( \sum_{j \in \mathcal{K} \setminus \mathcal{L}} D_{\mathcal{K},j} \Delta_{j,\mathcal{L}} \right). \quad (53)$$

- Two examples will be of particular importance for our analysis.
- First, consider choice set  $\mathcal{K} = \{k\}$  versus choice set  $\mathcal{L} = \mathcal{J} \setminus \{k\}$ .
- In this case,  $\Delta_{k, \mathcal{J} \setminus k}$  is the difference between the agent's potential outcome in state  $k$  versus the outcome that would have been observed if he or she had not been allowed to choose state  $k$ .

- If  $I_{\mathcal{J}} = k$ , then  $\Delta_{k, \mathcal{J} \setminus k}$  is the difference between the outcome in the agent's preferred state and the outcome in the agent's "next-best" state.

- If  $I_{\mathcal{J}} = k$ , then  $\Delta_{k, \mathcal{J} \setminus k}$  is the difference between the outcome in the agent's preferred state and the outcome in the agent's "next-best" state.
- Second, consider the set  $\mathcal{K} = \mathcal{J}$  versus choice set  $\mathcal{L} = \mathcal{J} \setminus \{k\}$ . In this case,  $\Delta_{\mathcal{J}, \mathcal{J} \setminus k}$  is the difference between the agent's best outcome and what his or her outcome would have been if state  $k$  had not been available.

- If  $I_{\mathcal{J}} = k$ , then  $\Delta_{k, \mathcal{J} \setminus k}$  is the difference between the outcome in the agent's preferred state and the outcome in the agent's "next-best" state.
- Second, consider the set  $\mathcal{K} = \mathcal{J}$  versus choice set  $\mathcal{L} = \mathcal{J} \setminus \{k\}$ . In this case,  $\Delta_{\mathcal{J}, \mathcal{J} \setminus k}$  is the difference between the agent's best outcome and what his or her outcome would have been if state  $k$  had not been available.
- Note that

$$\Delta_{\mathcal{J}, \mathcal{J} \setminus k} = D_{\mathcal{J}, k} \Delta_{k, \mathcal{J} \setminus k}.$$

- If  $I_{\mathcal{J}} = k$ , then  $\Delta_{k, \mathcal{J} \setminus k}$  is the difference between the outcome in the agent's preferred state and the outcome in the agent's "next-best" state.
- Second, consider the set  $\mathcal{K} = \mathcal{J}$  versus choice set  $\mathcal{L} = \mathcal{J} \setminus \{k\}$ . In this case,  $\Delta_{\mathcal{J}, \mathcal{J} \setminus k}$  is the difference between the agent's best outcome and what his or her outcome would have been if state  $k$  had not been available.
- Note that

$$\Delta_{\mathcal{J}, \mathcal{J} \setminus k} = D_{\mathcal{J}, k} \Delta_{k, \mathcal{J} \setminus k}.$$

- Thus, there is a trivial connection between the two parameters,  $\Delta_{\mathcal{J}, \mathcal{J} \setminus k}$  and  $\Delta_{k, \mathcal{J} \setminus k}$ .



- We will focus on  $\Delta_{k, \mathcal{J} \setminus k}$ , the effect of being forced to choose option  $k$  versus being denied option  $k$ .

- We will focus on  $\Delta_{k, \mathcal{J} \setminus k}$ , the effect of being forced to choose option  $k$  versus being denied option  $k$ .
- However, one can use equation 52 to use the results for  $\Delta_{k, \mathcal{J} \setminus k}$  to obtain results for  $\Delta_{\mathcal{J}, \mathcal{J} \setminus k}$ .

- We will focus on  $\Delta_{k, \mathcal{J} \setminus k}$ , the effect of being forced to choose option  $k$  versus being denied option  $k$ .
- However, one can use equation 52 to use the results for  $\Delta_{k, \mathcal{J} \setminus k}$  to obtain results for  $\Delta_{\mathcal{J}, \mathcal{J} \setminus k}$ .
- To fix ideas regarding these alternative definitions of treatment effects, consider the following example concerning GED certification.

- We will focus on  $\Delta_{k, \mathcal{J} \setminus k}$ , the effect of being forced to choose option  $k$  versus being denied option  $k$ .
- However, one can use equation 52 to use the results for  $\Delta_{k, \mathcal{J} \setminus k}$  to obtain results for  $\Delta_{\mathcal{J}, \mathcal{J} \setminus k}$ .
- To fix ideas regarding these alternative definitions of treatment effects, consider the following example concerning GED certification.
- The GED is an exam that certifies that high school dropouts who pass the test are the equivalents of high school graduates.

**Example: GED Certification** Consider studying the effect of GED certification on later wages. Consider the case where  $\mathcal{J} = \{ \text{GED}, \text{HS Degree}, \text{Permanent Dropout} \}$ . Let  $j = \{ \text{GED} \}$ ,  $k = \{ \text{HS Degree} \}$ , and  $l = \{ \text{Permanent Dropout} \}$ . Suppose one wishes to study the effect of the GED. Then possible definitions of the effect of the GED include:

- $\Delta_{j,k}$  is the individual's outcome if he or she received the GED versus if he or she had graduated from High School;

**Example: GED Certification** Consider studying the effect of GED certification on later wages. Consider the case where  $\mathcal{J} = \{ \text{GED}, \text{HS Degree}, \text{Permanent Dropout} \}$ . Let  $j = \{ \text{GED} \}$ ,  $k = \{ \text{HS Degree} \}$ , and  $l = \{ \text{Permanent Dropout} \}$ . Suppose one wishes to study the effect of the GED. Then possible definitions of the effect of the GED include:

- $\Delta_{j,k}$  is the individual's outcome if he or she received the GED versus if he or she had graduated from High School;
- $\Delta_{j,l}$  is the individual's outcome if he or she received the GED versus if he or she had been a permanent dropout;

**Example: GED Certification** Consider studying the effect of GED certification on later wages. Consider the case where  $\mathcal{J} = \{ \text{GED}, \text{HS Degree}, \text{Permanent Dropout} \}$ . Let  $j = \{ \text{GED} \}$ ,  $k = \{ \text{HS Degree} \}$ , and  $l = \{ \text{Permanent Dropout} \}$ . Suppose one wishes to study the effect of the GED. Then possible definitions of the effect of the GED include:

- $\Delta_{j,k}$  is the individual's outcome if he or she received the GED versus if he or she had graduated from High School;
- $\Delta_{j,l}$  is the individual's outcome if he or she received the GED versus if he or she had been a permanent dropout;
- $\Delta_{j, \mathcal{J} \setminus j}$  is the individual's outcome if he or she had received the GED versus what the outcome would have been if he or she had not had the option of receiving the GED;

- $\Delta_{\mathcal{J}, \mathcal{J} \setminus j}$  is the individual's outcome if he or she had the option of receiving the GED versus the outcome if he or she did not have the option of receiving the GED.



- $\Delta_{\mathcal{J}, \mathcal{J} \setminus j}$  is the individual's outcome if he or she had the option of receiving the GED versus the outcome if he or she did not have the option of receiving the GED.
- Notice that  $\Delta_{\mathcal{J}, \mathcal{J} \setminus j}$  is a version of an option value treatment effect.

- We now define treatment parameters for a general unordered model.

## Treatment Parameters

- The conventional definition of the average treatment effect (ATE) is

$$\Delta_{k,l}^{\text{ATE}}(x, z) = E(\Delta_{k,l} | X = x, Z = z),$$

which immediately generalizes to the class of parameters discussed in this section as

$$\Delta_{\mathcal{K},\mathcal{L}}^{\text{ATE}}(x, z) = E(\Delta_{\mathcal{K},\mathcal{L}} | X = x, Z = z).$$

## Treatment Parameters

- The conventional definition of the average treatment effect (ATE) is

$$\Delta_{k,l}^{\text{ATE}}(x, z) = E(\Delta_{k,l} | X = x, Z = z),$$

which immediately generalizes to the class of parameters discussed in this section as

$$\Delta_{\mathcal{K},\mathcal{L}}^{\text{ATE}}(x, z) = E(\Delta_{\mathcal{K},\mathcal{L}} | X = x, Z = z).$$

- Notice that the treatment parameters now depend on the value of  $Z$ .

## Treatment Parameters

- The conventional definition of the average treatment effect (ATE) is

$$\Delta_{k,l}^{\text{ATE}}(x, z) = E(\Delta_{k,l} | X = x, Z = z),$$

which immediately generalizes to the class of parameters discussed in this section as

$$\Delta_{\mathcal{K},\mathcal{L}}^{\text{ATE}}(x, z) = E(\Delta_{\mathcal{K},\mathcal{L}} | X = x, Z = z).$$

- Notice that the treatment parameters now depend on the value of  $Z$ .
- We explain the source of this dependence below.

- The conventional definition of the treatment on the treated (TT) parameter is

$$\Delta_{k,l}^{\text{TT}}(x, z) = E(\Delta_{k,l} | X = x, Z = z, I_{\mathcal{J}} = k),$$

which we generalize to

$$\Delta_{\mathcal{K},\mathcal{L}}^{\text{TT}}(x, z) = E(\Delta_{\mathcal{K},\mathcal{L}} | X = x, Z = z, I_{\mathcal{J}} \in \mathcal{K}).$$

- The conventional definition of the treatment on the treated (TT) parameter is

$$\Delta_{k,l}^{\text{TT}}(x, z) = E(\Delta_{k,l} | X = x, Z = z, I_{\mathcal{J}} = k),$$

which we generalize to

$$\Delta_{\mathcal{K},\mathcal{L}}^{\text{TT}}(x, z) = E(\Delta_{\mathcal{K},\mathcal{L}} | X = x, Z = z, I_{\mathcal{J}} \in \mathcal{K}).$$

- We also generalize the Marginal Treatment Effect (MTE) and Local Average Treatment Effect (LATE) parameters considered in ? . We generalize the MTE parameter to be the average effect conditional on being indifferent between the best option among choice set  $\mathcal{K}$  versus the best option among choice set  $\mathcal{L}$  at some fixed value of the instruments,  $Z = z$ :

$$\Delta_{\mathcal{K},\mathcal{L}}^{\text{MTE}}(x, z) = E(\Delta_{\mathcal{K},\mathcal{L}} | X = x, Z = z, R_{\mathcal{K}}(z) = R_{\mathcal{L}}(z)). \quad (54)$$

- We generalize the LATE parameter to be the average effect for someone for whom the optimal choice in choice set  $\mathcal{K}$  is preferred to the optimal choice in choice set  $\mathcal{L}$  at  $Z = \tilde{z}$ , but who prefers the optimal choice in choice set  $\mathcal{L}$  to the optimal choice in choice set  $\mathcal{K}$  at  $Z = z$ :

$$\Delta_{\mathcal{K}, \mathcal{L}}^{\text{LATE}}(x, z, \tilde{z}) = E(\Delta_{\mathcal{K}, \mathcal{L}} | X = x, Z = z, R_{\mathcal{K}}(\tilde{z}) \geq R_{\mathcal{L}}(\tilde{z}), R_{\mathcal{L}}(z) \geq R_{\mathcal{K}}(z)). \quad (55)$$



- We generalize the LATE parameter to be the average effect for someone for whom the optimal choice in choice set  $\mathcal{K}$  is preferred to the optimal choice in choice set  $\mathcal{L}$  at  $Z = \tilde{z}$ , but who prefers the optimal choice in choice set  $\mathcal{L}$  to the optimal choice in choice set  $\mathcal{K}$  at  $Z = z$ :

$$\Delta_{\mathcal{K}, \mathcal{L}}^{\text{LATE}}(x, z, \tilde{z}) = E(\Delta_{\mathcal{K}, \mathcal{L}} | X = x, Z = z, R_{\mathcal{K}}(\tilde{z}) \geq R_{\mathcal{L}}(\tilde{z}), R_{\mathcal{L}}(z) \geq R_{\mathcal{K}}(z)). \quad (55)$$

- An important special case of this parameter arises when  $z = \tilde{z}$  except for elements that enter the index functions only for choices in  $\mathcal{K}$  and not for any choice in  $\mathcal{L}$ .

- We generalize the LATE parameter to be the average effect for someone for whom the optimal choice in choice set  $\mathcal{K}$  is preferred to the optimal choice in choice set  $\mathcal{L}$  at  $Z = \tilde{z}$ , but who prefers the optimal choice in choice set  $\mathcal{L}$  to the optimal choice in choice set  $\mathcal{K}$  at  $Z = z$ :

$$\Delta_{\mathcal{K},\mathcal{L}}^{\text{LATE}}(x, z, \tilde{z}) = E(\Delta_{\mathcal{K},\mathcal{L}} | X = x, Z = z, R_{\mathcal{K}}(\tilde{z}) \geq R_{\mathcal{L}}(\tilde{z}), R_{\mathcal{L}}(z) \geq R_{\mathcal{K}}(z)). \quad (55)$$

- An important special case of this parameter arises when  $z = \tilde{z}$  except for elements that enter the index functions only for choices in  $\mathcal{K}$  and not for any choice in  $\mathcal{L}$ .
- In that special case, equation (55) simplifies to

$$\Delta_{\mathcal{K},\mathcal{L}}^{\text{LATE}}(x, z, \tilde{z}) = E(\Delta_{\mathcal{K},\mathcal{L}} | X = x, Z = z, R_{\mathcal{K}}(\tilde{z}) \geq R_{\mathcal{L}}(z) \geq R_{\mathcal{K}}(z)),$$

since  $R_{\mathcal{L}}(z) = R_{\mathcal{L}}(\tilde{z})$  in this special case.

- We have defined each of these parameters as conditional not only on  $X$  but also on the “instruments”  $Z$ .

- We have defined each of these parameters as conditional not only on  $X$  but also on the “instruments”  $Z$ .
- In general, the parameters depend on the  $Z$  evaluation point.

- We have defined each of these parameters as conditional not only on  $X$  but also on the “instruments”  $Z$ .
- In general, the parameters depend on the  $Z$  evaluation point.
- For example,  $\Delta_{\mathcal{K}, \mathcal{L}}^{\text{ATE}}(x, z)$  generally depends on the  $z$  evaluation point.

- We have defined each of these parameters as conditional not only on  $X$  but also on the “instruments”  $Z$ .
- In general, the parameters depend on the  $Z$  evaluation point.
- For example,  $\Delta_{\mathcal{K},\mathcal{L}}^{\text{ATE}}(x, z)$  generally depends on the  $z$  evaluation point.
- To see this, note that  $Y_{\mathcal{K}} = \sum_{k \in \mathcal{K}} D_{\mathcal{K},k} Y_k$ , and  $Y_{\mathcal{L}} = \sum_{l \in \mathcal{L}} D_{\mathcal{L},l} Y_l$ .

- By conditional independence assumption (B-1),  $Z \perp\!\!\!\perp \{Y_j\}_{j \in \mathcal{J}} \mid X$ , but  $D_{\mathcal{K},k}$  and  $D_{\mathcal{L},l}$  depend on  $Z$  conditional on  $X$  and thus  $Y_{\mathcal{K}} - Y_{\mathcal{L}}$ , in general, is dependent on  $Z$  conditional on  $X$ .

- By conditional independence assumption (B-1),  $Z \perp\!\!\!\perp \{Y_j\}_{j \in \mathcal{J}} \mid X$ , but  $D_{\mathcal{K},k}$  and  $D_{\mathcal{L},l}$  depend on  $Z$  conditional on  $X$  and thus  $Y_{\mathcal{K}} - Y_{\mathcal{L}}$ , in general, is dependent on  $Z$  conditional on  $X$ .
- In other words, even though  $Z$  is conditionally independent of each individual potential outcome, it is correlated with the indicator for the choice that is optimal within the sets  $\mathcal{K}$  and  $\mathcal{L}$  and thus is related to  $Y_{\mathcal{K}} - Y_{\mathcal{L}}$ .



## Heterogeneity in Treatment Effects

- Consider heterogeneity in the pairwise treatment effect  $\Delta_{j,k}$  (with  $(j, k) \in \mathcal{J}$ ) defined as

$$\Delta_{j,k} = Y_j - Y_k = \mu_j(X_j, U_j) - \mu_k(X_k, U_k),$$

which in general will vary with both observables ( $X$ ) and unobservables ( $U_j, U_k$ ).

## Heterogeneity in Treatment Effects

- Consider heterogeneity in the pairwise treatment effect  $\Delta_{j,k}$  (with  $(j, k) \in \mathcal{J}$ ) defined as

$$\Delta_{j,k} = Y_j - Y_k = \mu_j(X_j, U_j) - \mu_k(X_k, U_k),$$

which in general will vary with both observables ( $X$ ) and unobservables ( $U_j, U_k$ ).

- Since we have not assumed that the error terms are additively separable, the treatment effect will in general vary with unobservables even if  $U_j = U_k$ .

- The mean treatment parameters for  $\Delta_{j,k}$  will differ if the effect of treatment is heterogeneous and agents base participation decisions, in part, on their idiosyncratic treatment effect.

- The mean treatment parameters for  $\Delta_{j,k}$  will differ if the effect of treatment is heterogeneous and agents base participation decisions, in part, on their idiosyncratic treatment effect.
- In general, the ATE, TT, and the marginal treatment parameters for  $\Delta_{j,k}$  will differ as long as there is dependence between  $(U_j, U_k)$  and the decision rule, i.e., if there is dependence between  $(U_j, U_k)$  and  $(\{V_l\}_{l \in \mathcal{J}})$ .

- The mean treatment parameters for  $\Delta_{j,k}$  will differ if the effect of treatment is heterogeneous and agents base participation decisions, in part, on their idiosyncratic treatment effect.
- In general, the ATE, TT, and the marginal treatment parameters for  $\Delta_{j,k}$  will differ as long as there is dependence between  $(U_j, U_k)$  and the decision rule, i.e., if there is dependence between  $(U_j, U_k)$  and  $(\{V_l\}_{l \in \mathcal{J}})$ .
- If we impose that  $(\{V_l\}_{l \in \mathcal{J}})$  is independent of  $(U_j, U_k)$ , then the treatment effect will still be heterogeneous, but the average treatment effect, average effect of treatment on the treated, and the marginal average treatment effects will all coincide.

- The literature on treatment effects often imposes additive separability in outcomes between observables and unobservables.

- The literature on treatment effects often imposes additive separability in outcomes between observables and unobservables.
- In particular, it is commonly assumed that  $U_j$  and  $U_k$  are scalar random variables and that  $Y_j = \mu_j(X_j) + U_j$ ,  $Y_k = \mu_k(X_k) + U_k$ .

- The literature on treatment effects often imposes additive separability in outcomes between observables and unobservables.
- In particular, it is commonly assumed that  $U_j$  and  $U_k$  are scalar random variables and that  $Y_j = \mu_j(X_j) + U_j$ ,  $Y_k = \mu_k(X_k) + U_k$ .
- In that case, a common treatment effect model is produced if the additive error term does not vary with the treatment state:  
 $U_j = U_k$ .



- The literature on treatment effects often imposes additive separability in outcomes between observables and unobservables.
- In particular, it is commonly assumed that  $U_j$  and  $U_k$  are scalar random variables and that  $Y_j = \mu_j(X_j) + U_j$ ,  $Y_k = \mu_k(X_k) + U_k$ .
- In that case, a common treatment effect model is produced if the additive error term does not vary with the treatment state:  
 $U_j = U_k$ .
- Thus, in the special case of additive separability, the treatment parameters for  $\Delta_{j,k}$  will be the same even if there is dependence between  $\{V_j\}_{j \in \mathcal{J}}$  and  $(U_j, U_k)$  as long as  $U_j = U_k$ .

- There is an additional source of treatment heterogeneity in the more general case of  $\Delta_{\mathcal{K},\mathcal{L}}$  arising from heterogeneity in which states are being compared.

- There is an additional source of treatment heterogeneity in the more general case of  $\Delta_{\mathcal{K},\mathcal{L}}$  arising from heterogeneity in which states are being compared.
- Consider, for example,  $\Delta_{j,\mathcal{J}\setminus j}$ .

- There is an additional source of treatment heterogeneity in the more general case of  $\Delta_{\mathcal{K},\mathcal{L}}$  arising from heterogeneity in which states are being compared.
- Consider, for example,  $\Delta_{j,\mathcal{J}\setminus j}$ .
- We have that

$$\Delta_{j,\mathcal{J}\setminus j} = \sum_{k \in \mathcal{J}\setminus j} D_{\mathcal{J}\setminus j,k} \Delta_{j,k},$$

which will vary over individuals even if each individual has the same  $\Delta_{j,k}$  treatment effect.

- Consider the corresponding ATE and TT parameters:

$$\begin{aligned}
 \Delta_{j, \mathcal{J} \setminus j}^{\text{ATE}}(x, z) &= E(\Delta_{j, \mathcal{J} \setminus j} \mid X = x, Z = z) \\
 &= \sum_{k \in \mathcal{J} \setminus j} \Pr(I_{\mathcal{J} \setminus j} = k \mid X = x, Z = z) E(\Delta_{j, k} \mid X = x, Z = z, I_{\mathcal{J} \setminus j} = k)
 \end{aligned}$$

and

$$\begin{aligned}
 \Delta_{j, \mathcal{J} \setminus j}^{\text{TT}}(x, z) &= E(\Delta_{j, \mathcal{J} \setminus j} \mid X = x, Z = z, I_{\mathcal{J}} = j) \\
 &= \sum_{k \in \mathcal{J} \setminus j} \Pr(I_{\mathcal{J} \setminus j} = k \mid X = x, Z = z, I_{\mathcal{J}} = j) \\
 &\quad \times E(\Delta_{j, k} \mid X = x, Z = z, I_{\mathcal{J}} = j, I_{\mathcal{J} \setminus j} = k).
 \end{aligned}$$

- Even in the case where  $\{U_j\}_{j \in \mathcal{J}}$  is independent of  $\{V_j\}_{j \in \mathcal{J}}$ , so that  $E(\Delta_{j,k} | X = x, Z = z, I_{\mathcal{J} \setminus j} = k) = E(\Delta_{j,k} | X = x, Z = z, I_{\mathcal{J}} = j, I_{\mathcal{J} \setminus j} = k)$ , it will still in general be the case that  $\Delta_{j, \mathcal{J} \setminus j}^{\text{ATE}}(x, z) \neq \Delta_{j, \mathcal{J} \setminus j}^{\text{TT}}(x, z)$  since in general  $\Pr(I_{\mathcal{J} \setminus j} = k | X = x, Z = z) \neq \Pr(I_{\mathcal{J} \setminus j} = k | X = x, Z = z, I_{\mathcal{J}} = j)$ .

- Even in the case where  $\{U_j\}_{j \in \mathcal{J}}$  is independent of  $\{V_j\}_{j \in \mathcal{J}}$ , so that  $E(\Delta_{j,k} | X = x, Z = z, I_{\mathcal{J} \setminus j} = k) = E(\Delta_{j,k} | X = x, Z = z, I_{\mathcal{J}} = j, I_{\mathcal{J} \setminus j} = k)$ , it will still in general be the case that  $\Delta_{j, \mathcal{J} \setminus j}^{\text{ATE}}(x, z) \neq \Delta_{j, \mathcal{J} \setminus j}^{\text{TT}}(x, z)$  since in general  $\Pr(I_{\mathcal{J} \setminus j} = k | X = x, Z = z) \neq \Pr(I_{\mathcal{J} \setminus j} = k | X = x, Z = z, I_{\mathcal{J}} = j)$ .
- Thus, the ATE and TT parameters will differ in part because they place different weights on the alternative pairwise treatment effects, and thus will differ even in the case where the pairwise ( $j$  versus  $k$ ) treatment effects are common across all individuals.

- In summary,  $\Delta_{j,k}$  will be heterogeneous depending on the functional form of the  $\mu_j(\cdot)$  and  $\mu_k(\cdot)$  equations and on the pairwise dependence between the  $U_j$  and  $U_k$  terms.



- In summary,  $\Delta_{j,k}$  will be heterogeneous depending on the functional form of the  $\mu_j(\cdot)$  and  $\mu_k(\cdot)$  equations and on the pairwise dependence between the  $U_j$  and  $U_k$  terms.
- The  $\Delta_{j,k}$  mean treatment parameters will also vary depending on the dependence between  $\{V_I\}_{I \in \mathcal{J}}$  and  $(U_j, U_k)$ .

- In summary,  $\Delta_{j,k}$  will be heterogeneous depending on the functional form of the  $\mu_j(\cdot)$  and  $\mu_k(\cdot)$  equations and on the pairwise dependence between the  $U_j$  and  $U_k$  terms.
- The  $\Delta_{j,k}$  mean treatment parameters will also vary depending on the dependence between  $\{V_I\}_{I \in \mathcal{J}}$  and  $(U_j, U_k)$ .
- For  $\Delta_{j, \mathcal{J} \setminus j}$ , there is an additional source of heterogeneity arising from the variability in the optimal option in the set  $\mathcal{J} \setminus j$ .

- In summary,  $\Delta_{j,k}$  will be heterogeneous depending on the functional form of the  $\mu_j(\cdot)$  and  $\mu_k(\cdot)$  equations and on the pairwise dependence between the  $U_j$  and  $U_k$  terms.
- The  $\Delta_{j,k}$  mean treatment parameters will also vary depending on the dependence between  $\{V_I\}_{I \in \mathcal{J}}$  and  $(U_j, U_k)$ .
- For  $\Delta_{j, \mathcal{J} \setminus j}$ , there is an additional source of heterogeneity arising from the variability in the optimal option in the set  $\mathcal{J} \setminus j$ .
- Even if there is no heterogeneity in the pairwise  $\Delta_{j,k}$  terms, there will still be heterogeneity in  $\Delta_{j, \mathcal{J} \setminus j}$ , and heterogeneity in the corresponding mean treatment parameters.

## LIV and Nonparametric Wald Estimands for One Choice vs. the Best Alternative

- We first consider identification of treatment parameters corresponding to averages of  $\Delta_{j, \mathcal{J} \setminus j}$ , the effect of choosing option  $j$  versus the preferred option  $\mathcal{J}$  if  $j$  is not available.

## LIV and Nonparametric Wald Estimands for One Choice vs. the Best Alternative

- We first consider identification of treatment parameters corresponding to averages of  $\Delta_{j, \mathcal{J} \setminus j}$ , the effect of choosing option  $j$  versus the preferred option  $\mathcal{J}$  if  $j$  is not available.
- We analyze both a discrete change (Wald form for the instrumental variables estimand) and the local instrumental variables (LIV) estimand.

## LIV and Nonparametric Wald Estimands for One Choice vs. the Best Alternative

- We first consider identification of treatment parameters corresponding to averages of  $\Delta_{j, \mathcal{J} \setminus j}$ , the effect of choosing option  $j$  versus the preferred option  $\mathcal{J}$  if  $j$  is not available.
- We analyze both a discrete change (Wald form for the instrumental variables estimand) and the local instrumental variables (LIV) estimand.
- Using a concise notation, define  $Z^{[j]}$  as the vector of elements in  $Z_j$  that do not enter any other choice index, and that  $Z^{[-j]}$  is a vector of elements of  $Z$  not in  $Z^{[j]}$ .

## LIV and Nonparametric Wald Estimands for One Choice vs. the Best Alternative

- We first consider identification of treatment parameters corresponding to averages of  $\Delta_{j, \mathcal{J} \setminus j}$ , the effect of choosing option  $j$  versus the preferred option  $\mathcal{J}$  if  $j$  is not available.
- We analyze both a discrete change (Wald form for the instrumental variables estimand) and the local instrumental variables (LIV) estimand.
- Using a concise notation, define  $Z^{[j]}$  as the vector of elements in  $Z_j$  that do not enter any other choice index, and that  $Z^{[-j]}$  is a vector of elements of  $Z$  not in  $Z^{[j]}$ .
- The  $Z^{[j]}$  thus act as shifters attracting people into or out of state  $j$  but not affecting the valuations in the arguments of the other choice functions.

- For this case, we can develop an analysis of IV parallel to that given for the binary case or the ordered choice case if we condition on  $Z^{[-j]}$ .



- For this case, we can develop an analysis of IV parallel to that given for the binary case or the ordered choice case if we condition on  $Z^{[-j]}$ .
- We obtain monotonicity or uniformity in this model if the movements among states induced by  $Z^{[j]}$  are the same for all persons conditional on  $Z^{[-j]} = z^{[-j]}$  and  $X = x$ .

- For this case, we can develop an analysis of IV parallel to that given for the binary case or the ordered choice case if we condition on  $Z^{[-j]}$ .
- We obtain monotonicity or uniformity in this model if the movements among states induced by  $Z^{[j]}$  are the same for all persons conditional on  $Z^{[-j]} = z^{[-j]}$  and  $X = x$ .
- For example, *ceteris paribus* if  $Z^{[j]} = z^{[j]}$  increases,  $R_j(Z_j)$  increases but the  $R_k(Z_k)$  are not affected, so the flow is toward state  $j$ .

- For this case, we can develop an analysis of IV parallel to that given for the binary case or the ordered choice case if we condition on  $Z^{[-j]}$ .
- We obtain monotonicity or uniformity in this model if the movements among states induced by  $Z^{[j]}$  are the same for all persons conditional on  $Z^{[-j]} = z^{[-j]}$  and  $X = x$ .
- For example, *ceteris paribus* if  $Z^{[j]} = z^{[j]}$  increases,  $R_j(Z_j)$  increases but the  $R_k(Z_k)$  are not affected, so the flow is toward state  $j$ .
- Let  $D_{\mathcal{J},j}$  be an indicator variable denoting whether option  $j$  is selected.

- $$\begin{aligned} D_{\mathcal{J},j} &= \mathbf{1} \left( R_j(Z_j) \geq \max_{\ell \neq j} \{R_\ell(Z_\ell)\} \right) \\ &= \mathbf{1} \left( \vartheta_j(Z_j) \geq V_j + \max_{\ell \neq j} \{R_\ell(Z_\ell)\} \right) \\ &= \mathbf{1} \left( \vartheta_j(Z_j) \geq \tilde{V}_j \right), \end{aligned} \tag{56}$$

where  $\tilde{V}_j = V_j + \max_{\ell \neq j} \{R_\ell(Z_\ell)\}$ .



$$\begin{aligned}
 D_{\mathcal{J},j} &= \mathbf{1} \left( R_j(Z_j) \geq \max_{\ell \neq j} \{R_\ell(Z_\ell)\} \right) & (56) \\
 &= \mathbf{1} \left( \vartheta_j(Z_j) \geq V_j + \max_{\ell \neq j} \{R_\ell(Z_\ell)\} \right) \\
 &= \mathbf{1} \left( \vartheta_j(Z_j) \geq \tilde{V}_j \right),
 \end{aligned}$$

where  $\tilde{V}_j = V_j + \max_{\ell \neq j} \{R_\ell(Z_\ell)\}$ .

- Thus we obtain  $D_{\mathcal{J},j} = \mathbf{1} (P_j(Z_j) \geq U_{D_j})$ , where  $U_{D_j} = F_{\tilde{V}_j|Z^{[-j]}}(V_j + \max_{\ell \neq j} \{R_\ell(Z_\ell)\} | Z^{[-j]} = z^{[-j]})$ , where  $F_{\tilde{V}_j|Z^{[-j]}}$  is the cdf of  $\tilde{V}_j$  given  $Z^{[-j]} = z^{[-j]}$ .

- In a format parallel to the binary model, we write

$$Y = D_{\mathcal{J},j} Y_j + (1 - D_{\mathcal{J},j}) Y_{\mathcal{J}\setminus j}, \quad (57)$$

where  $Y_{\mathcal{J}\setminus j}$  is the outcome that would be observed if option  $j$  were not available.

- In a format parallel to the binary model, we write

$$Y = D_{\mathcal{J},j} Y_j + (1 - D_{\mathcal{J},j}) Y_{\mathcal{J}\setminus j}, \quad (57)$$

where  $Y_{\mathcal{J}\setminus j}$  is the outcome that would be observed if option  $j$  were not available.

- This case is just a version of the binary case developed in previous sections of the paper.

- In a format parallel to the binary model, we write

$$Y = D_{\mathcal{J},j} Y_j + (1 - D_{\mathcal{J},j}) Y_{\mathcal{J}\setminus j}, \quad (57)$$

where  $Y_{\mathcal{J}\setminus j}$  is the outcome that would be observed if option  $j$  were not available.

- This case is just a version of the binary case developed in previous sections of the paper.
- There is one crucial difference, however, and that is that the distributions of the  $\tilde{V}_j$  now depend on the excluded  $Z = z$ .



- In a format parallel to the binary model, we write

$$Y = D_{\mathcal{J},j} Y_j + (1 - D_{\mathcal{J},j}) Y_{\mathcal{J}\setminus j}, \quad (57)$$

where  $Y_{\mathcal{J}\setminus j}$  is the outcome that would be observed if option  $j$  were not available.

- This case is just a version of the binary case developed in previous sections of the paper.
- There is one crucial difference, however, and that is that the distributions of the  $\tilde{V}_j$  now depend on the excluded  $Z = z$ .
- Thus instruments and parameters have to be defined conditionally on  $Z = z$ .

- We can define MTE as

$$E\left(Y_j - Y_{\mathcal{J}\setminus j} \mid X = x, Z = z, \vartheta_j(z_j) - V_j = R_{\mathcal{J}\setminus j}(z)\right).$$

We have to condition on  $Z = z$  because the choice sets are defined over the max of elements in  $\mathcal{J} \setminus j$  (see equation (56)).

- We can define MTE as

$$E\left(Y_j - Y_{\mathcal{J}\setminus j} \mid X = x, Z = z, \vartheta_j(z_j) - V_j = R_{\mathcal{J}\setminus j}(z)\right).$$

We have to condition on  $Z = z$  because the choice sets are defined over the max of elements in  $\mathcal{J} \setminus j$  (see equation (56)).

- We now show that our identification strategies presented in the preceding part of this paper extend naturally to the identification of treatment parameters for  $\Delta_{j, \mathcal{J}\setminus j}$ .

- We can define MTE as

$$E\left(Y_j - Y_{\mathcal{J}\setminus j} \mid X = x, Z = z, \vartheta_j(z_j) - V_j = R_{\mathcal{J}\setminus j}(z)\right).$$

We have to condition on  $Z = z$  because the choice sets are defined over the max of elements in  $\mathcal{J} \setminus j$  (see equation (56)).

- We now show that our identification strategies presented in the preceding part of this paper extend naturally to the identification of treatment parameters for  $\Delta_{j,\mathcal{J}\setminus j}$ .
- In particular, it is possible to recover LATE and MTE parameters for  $\Delta_{j,\mathcal{J}\setminus j}$  by use of discrete change IV methods and local instrumental variable methods, respectively.

- We can define MTE as

$$E\left(Y_j - Y_{\mathcal{J} \setminus j} \mid X = x, Z = z, \vartheta_j(z_j) - V_j = R_{\mathcal{J} \setminus j}(z)\right).$$

We have to condition on  $Z = z$  because the choice sets are defined over the max of elements in  $\mathcal{J} \setminus j$  (see equation (56)).

- We now show that our identification strategies presented in the preceding part of this paper extend naturally to the identification of treatment parameters for  $\Delta_{j, \mathcal{J} \setminus j}$ .
- In particular, it is possible to recover LATE and MTE parameters for  $\Delta_{j, \mathcal{J} \setminus j}$  by use of discrete change IV methods and local instrumental variable methods, respectively.
- Averages of the effect of option  $j$  versus the next best alternative are the easiest effects to study using instrumental variable methods and are natural generalizations of our two-outcome analysis.

- The discrete change instrumental variables estimand will allow us to recover a version of the local average treatment effect (LATE) parameter.

- The discrete change instrumental variables estimand will allow us to recover a version of the local average treatment effect (LATE) parameter.
- Invoke assumption (B-2a).

- The discrete change instrumental variables estimand will allow us to recover a version of the local average treatment effect (LATE) parameter.
- Invoke assumption (B-2a).
- Assuming only one excluded variable  $Z^{[j]}$  in  $Z_j$ .



- The discrete change instrumental variables estimand will allow us to recover a version of the local average treatment effect (LATE) parameter.
- Invoke assumption (B-2a).
- Assuming only one excluded variable  $Z^{[j]}$  in  $Z_j$ .
- If there are more, pick any one that satisfies (B-2a).

- The discrete change instrumental variables estimand will allow us to recover a version of the local average treatment effect (LATE) parameter.
- Invoke assumption (B-2a).
- Assuming only one excluded variable  $Z^{[j]}$  in  $Z_j$ .
- If there are more, pick any one that satisfies (B-2a).
- Let  $Z^{[-j]}$  denote the excluded variable for option  $j$  with properties assumed in (B-2a).

- The discrete change instrumental variables estimand will allow us to recover a version of the local average treatment effect (LATE) parameter.
- Invoke assumption (B-2a).
- Assuming only one excluded variable  $Z^{[j]}$  in  $Z_j$ .
- If there are more, pick any one that satisfies (B-2a).
- Let  $Z^{[-j]}$  denote the excluded variable for option  $j$  with properties assumed in (B-2a).
- We let  $Z = [Z^{[-j]}, Z^{[j]}]$  and  $\tilde{Z} = [\tilde{Z}^{[-j]}, \tilde{Z}^{[j]}]$  be two values where we only manipulate scalar  $Z^{[j]}$ .



$$\Delta_j^{\text{Wald}}(x, z^{[-j]}, z^{[j]}, \tilde{z}^{[j]}) = \frac{E(Y|X = x, Z = \tilde{z}) - E(Y|X = x, Z = z)}{\Pr(D_{\mathcal{J},j} = 1|X = x, Z = \tilde{z}) - \Pr(D_{\mathcal{J},j} = 1|X = x, Z = z)},$$

where for notational convenience we are assuming that  $Z^{[j]}$  is the last element of  $Z$ .



$$\Delta_j^{\text{Wald}}(x, z^{[-j]}, z^{[j]}, \tilde{z}^{[j]}) = \frac{E(Y|X = x, Z = \tilde{z}) - E(Y|X = x, Z = z)}{\Pr(D_{\mathcal{J},j} = 1|X = x, Z = \tilde{z}) - \Pr(D_{\mathcal{J},j} = 1|X = x, Z = z)},$$

where for notational convenience we are assuming that  $Z^{[j]}$  is the last element of  $Z$ .

- Note that all components of  $z$  and  $\tilde{z}$  are the same except for the  $j^{\text{th}}$  component.



$$\Delta_j^{\text{Wald}}(x, z^{[-j]}, z^{[j]}, \tilde{z}^{[j]}) = \frac{E(Y|X = x, Z = \tilde{z}) - E(Y|X = x, Z = z)}{\Pr(D_{\mathcal{J},j} = 1|X = x, Z = \tilde{z}) - \Pr(D_{\mathcal{J},j} = 1|X = x, Z = z)},$$

where for notational convenience we are assuming that  $Z^{[j]}$  is the last element of  $Z$ .

- Note that all components of  $z$  and  $\tilde{z}$  are the same except for the  $j^{\text{th}}$  component.
- Without loss of generality, we assume that  $\vartheta_j(\tilde{z}) > \vartheta_j(z)$ .

- If there were no  $X$  regressors, and if  $Z$  were a scalar, binary random variable, then  $\Delta_j^{\text{Wald}}(x, z^{[-j]}, z^{[j]}, \tilde{z}^{[j]})$  would be the probability limit of the Wald form of two-stage least squares regression (2SLS).

- If there were no  $X$  regressors, and if  $Z$  were a scalar, binary random variable, then  $\Delta_j^{\text{Wald}}(x, z^{[-j]}, z^{[j]}, \tilde{z}^{[j]})$  would be the probability limit of the Wald form of two-stage least squares regression (2SLS).
- With  $X$  regressors, and with  $Z$  a vector possibly including continuous components, it no longer corresponds to a Wald/2SLS, but rather to a nonparametric version of the Wald estimator where the analyst nonparametrically conditions on  $X$  and on  $Z$  taking one of two specified values.



- The local instrumental variables estimator (LIV) estimand introduced in ?, and developed further in ??? and ?, will allow us to recover a version of the Marginal Treatment Effect (MTE) parameter.

- The local instrumental variables estimator (LIV) estimand introduced in ?, and developed further in ??? and ?, will allow us to recover a version of the Marginal Treatment Effect (MTE) parameter.
- Impose (B-2b), and let  $Z^{[j]}$  denote the excluded variable for option  $j$  with properties assumed in (B-2b).

- The local instrumental variables estimator (LIV) estimand introduced in ?, and developed further in ??? and ?, will allow us to recover a version of the Marginal Treatment Effect (MTE) parameter.
- Impose (B-2b), and let  $Z^{[j]}$  denote the excluded variable for option  $j$  with properties assumed in (B-2b).
- Because of the index structure, the LIV estimand will be invariant to which particular variable in  $Z^{[j]}$  satisfying (B-2b) is used if there is more than one variable with the property assumed in (B-2b).

- The local instrumental variables estimator (LIV) estimand introduced in ?, and developed further in ??? and ?, will allow us to recover a version of the Marginal Treatment Effect (MTE) parameter.
- Impose (B-2b), and let  $Z^{[j]}$  denote the excluded variable for option  $j$  with properties assumed in (B-2b).
- Because of the index structure, the LIV estimand will be invariant to which particular variable in  $Z^{[j]}$  satisfying (B-2b) is used if there is more than one variable with the property assumed in (B-2b).
- The effects are *not* invariant to variables in  $Z^{[-j]}$ .

- Define

$$\Delta_j^{\text{LIV}}(x, z) \equiv \frac{\partial}{\partial z_j} E(Y|X = x, Z = z) / \frac{\partial}{\partial z_j} \text{Pr}(D_{\mathcal{J}, j} = 1|X = x, Z = z).$$

- Define

$$\Delta_j^{\text{LIV}}(x, z) \equiv \frac{\partial}{\partial z^j} E(Y|X = x, Z = z) / \frac{\partial}{\partial z^j} \Pr(D_{\mathcal{J}, j} = 1|X = x, Z = z).$$

- $\Delta_j^{\text{LIV}}(x, z)$  is thus the limit form of  $\Delta_j^{\text{Wald}}(x, z^{[-j]}, z^{[j]}, \tilde{z}^{[j]})$  as  $\tilde{z}^{[j]}$  approaches  $z^{[j]}$ .

- Define

$$\Delta_j^{\text{LIV}}(x, z) \equiv \frac{\partial}{\partial z^{[j]}} E(Y|X = x, Z = z) \Big/ \frac{\partial}{\partial z^{[j]}} \Pr(D_{\mathcal{J}, j} = 1|X = x, Z = z).$$

- $\Delta_j^{\text{LIV}}(x, z)$  is thus the limit form of  $\Delta_j^{\text{Wald}}(x, z^{[-j]}, z^{[j]}, \tilde{z}^{[j]})$  as  $\tilde{z}^{[j]}$  approaches  $z^{[j]}$ .
- Given our previous assumptions, one can easily show that this limit exists w.p.1.

- Define

$$\Delta_j^{\text{LIV}}(x, z) \equiv \frac{\partial}{\partial z^{[j]}} E(Y|X = x, Z = z) \Big/ \frac{\partial}{\partial z^{[j]}} \Pr(D_{\mathcal{J}, j} = 1|X = x, Z = z).$$

- $\Delta_j^{\text{LIV}}(x, z)$  is thus the limit form of  $\Delta_j^{\text{Wald}}(x, z^{[-j]}, z^{[j]}, \tilde{z}^{[j]})$  as  $\tilde{z}^{[j]}$  approaches  $z^{[j]}$ .
- Given our previous assumptions, one can easily show that this limit exists w.p.1.
- LIV corresponds to a nonparametric, local version of indirect least squares.



- It is a function of the distribution of the observable data, and it can be consistently estimated using any nonparametric estimator of the derivative of a conditional expectation.

- It is a function of the distribution of the observable data, and it can be consistently estimated using any nonparametric estimator of the derivative of a conditional expectation.
- Given these definitions, we have the following identification theorem.

## Theorem 6

- 1 Assume (B-1), (B-3)–(B-5), and (B-2a). Then

$$\Delta_j^{Wald}(x, z^{[-j]}, z^{[j]}, \tilde{z}^{[j]}) = \Delta_{j, \mathcal{J} \setminus j}^{LATE}(x, z, \tilde{z})$$

where  $\tilde{z} = (z^{[-j]}, \tilde{z}^{[j]})$ .
- 2 Assume (B-1), (B-3)–(B-5), and (B-2b). Then

$$\Delta_j^{LIV}(x, z) = \Delta_{j, \mathcal{J} \setminus j}^{MTE}(x, z).$$

Proof.

See Appendix, Slide 1148.



- The intuition underlying the proof is simple.

- The intuition underlying the proof is simple.
- Under (B-1), (B-3)–(B-5), and (B-2a), we can convert the problem of comparing the outcome under  $j$  with the outcome under the next best option.

- The intuition underlying the proof is simple.
- Under (B-1), (B-3)–(B-5), and (B-2a), we can convert the problem of comparing the outcome under  $j$  with the outcome under the next best option.
- This is an IV version of the selection modeling of ?.

- The intuition underlying the proof is simple.
- Under (B-1), (B-3)–(B-5), and (B-2a), we can convert the problem of comparing the outcome under  $j$  with the outcome under the next best option.
- This is an IV version of the selection modeling of ?.
- $\Delta_{j, \mathcal{J} \setminus j}^{\text{LATE}}(x, z, \tilde{z})$  is the average effect of switching to state  $j$  from state  $l_{\mathcal{J} \setminus j}$  for individuals who would choose  $l_{\mathcal{J} \setminus j}$  at  $Z = z$  but would choose  $j$  at  $Z = \tilde{z}$ .



- The intuition underlying the proof is simple.
- Under (B-1), (B-3)–(B-5), and (B-2a), we can convert the problem of comparing the outcome under  $j$  with the outcome under the next best option.
- This is an IV version of the selection modeling of ?.
- $\Delta_{j, \mathcal{J} \setminus j}^{\text{LATE}}(x, z, \tilde{z})$  is the average effect of switching to state  $j$  from state  $l_{\mathcal{J} \setminus j}$  for individuals who would choose  $l_{\mathcal{J} \setminus j}$  at  $Z = z$  but would choose  $j$  at  $Z = \tilde{z}$ .
- $\Delta_{j, \mathcal{J} \setminus j}^{\text{MTE}}(x, z)$  is the average effect of switching to state  $j$  from state  $l_{\mathcal{J} \setminus j}$  (the best option besides state  $j$ ) for individuals who are indifferent between state  $j$  and  $l_{\mathcal{J} \setminus j}$  at the given values of the selection indices (at  $Z = z$ , i.e., at  $\{\vartheta_k(Z_k) = \vartheta_k(z_k)\}_{k \in \mathcal{J}}$ ).

- The mean effect of state  $j$  versus state  $I_{\mathcal{J} \setminus j}$  (the next best option) is a weighted average over  $k \in \mathcal{J} \setminus j$  of the effect of state  $j$  versus state  $k$ , conditional on  $k$  being the next best option, weighted by the probability that  $k$  is the next best option.

- The mean effect of state  $j$  versus state  $I_{\mathcal{J}\setminus j}$  (the next best option) is a weighted average over  $k \in \mathcal{J} \setminus j$  of the effect of state  $j$  versus state  $k$ , conditional on  $k$  being the next best option, weighted by the probability that  $k$  is the next best option.
- For example, for the LATE parameter,

$$\begin{aligned} \Delta_{j, \mathcal{J}\setminus j}^{\text{LATE}}(x, z, \tilde{z}) &= E(\Delta_{j, \mathcal{J}\setminus j} | X = x, Z = z, R_j(\tilde{z}) \geq R_{\mathcal{J}\setminus j}(z) \geq R_j(z)) \\ &= \sum_{k \in \mathcal{J}\setminus j} \left[ Pr(I_{\mathcal{J}\setminus j} = k | Z \in \{z, \tilde{z}\}, X = x, R_j(\tilde{z}) \geq R_{\mathcal{J}\setminus j}(z) \geq R_j(z)) \right. \\ &\quad \left. \times E(\Delta_{j, k} | X = x, Z \in \{z, \tilde{z}\}, R_j(\tilde{z}) \geq R_{\mathcal{J}\setminus j}(z) \geq R_j(z), I_{\mathcal{J}\setminus j} = k) \right]. \end{aligned}$$

where we use the result that  $R_{\mathcal{J}\setminus j}(z) = R_{\mathcal{J}\setminus j}(\tilde{z})$  since  $z = \tilde{z}$  except for one component that only enters the index for the  $j$ th option.

- The higher  $\vartheta_k(z_k)$ , holding the other indices constant, the larger the weight given to  $k$  as the base state.

- The higher  $\vartheta_k(z_k)$ , holding the other indices constant, the larger the weight given to  $k$  as the base state.
- Thus, how heavily each option is weighted in this average depends on the switching probability  $\Pr(I_{\mathcal{J}\setminus j} = k | Z = z, X = x, R_j(\tilde{z}_j) \geq R_k(z_k) \geq R_j(z_j))$ , which in turn depends on  $\{\vartheta_k(z_k)\}_{k \in \mathcal{J}\setminus j}$ .

- The LIV and Wald estimands depend on the  $z$  evaluation point.

- The LIV and Wald estimands depend on the  $z$  evaluation point.
- Alternatively, one can define averaged versions of the LIV and Wald estimands that will recover averaged versions of the MTE and LATE parameters,

$$\begin{aligned} \int \Delta_j^{\text{Wald}}(x, z^{[-j]}, z^{[j]}, \tilde{z}^{[j]}) dF_{Z^{[-j]}}(z^{[-j]}) \\ &= \int \Delta_{j, \mathcal{J} \setminus j}^{\text{LATE}}(x, z, \tilde{z}) dF_{Z^{[-j]}}(z^{[-j]}) \\ &= E(\Delta_{j, \mathcal{J} \setminus j} | X = x, R_j(Z^{[-j]}, \tilde{z}^{[j]}) \geq R_{\mathcal{J} \setminus j}(Z^{[-j]}) \geq R_j(Z^{[-j]}, z^{[j]})), \end{aligned}$$

and

$$\begin{aligned} \int \Delta_j^{\text{LIV}}(x, z) dF_Z(z) &= \int \Delta_{j, \mathcal{J} \setminus j}^{\text{MTE}}(x, z) dF_Z(z) \\ &= E(\Delta_{j, \mathcal{J} \setminus j} | X = x, R_j(Z) = R_{\mathcal{J} \setminus j}(Z)). \end{aligned}$$

- Thus far we have only considered identification of marginal treatment effect parameters, LATE and MTE, and not of the more standard treatment parameters like ATE and TT.



- Thus far we have only considered identification of marginal treatment effect parameters, LATE and MTE, and not of the more standard treatment parameters like ATE and TT.
- However, following ??, LATE can approximate ATE or TT arbitrarily well given the appropriate support conditions.

- Thus far we have only considered identification of marginal treatment effect parameters, LATE and MTE, and not of the more standard treatment parameters like ATE and TT.
- However, following ??, LATE can approximate ATE or TT arbitrarily well given the appropriate support conditions.
- Theorem 6 shows that we can use Wald estimands to identify LATE for  $\Delta_{j, \mathcal{J} \setminus j}$ , and we can thus adapt the analysis of ??, as reviewed in Slide 152, to identify ATE or TT for  $\Delta_{j, \mathcal{J} \setminus j}$ .

- Suppose that  $Z^{[j]}$  denotes the excluded variable for option  $j$  with properties assumed in (B-2a), and suppose that: (i) the support of the distribution of  $Z^{[j]}$  conditional on all other elements of  $Z$  is the full real line; (ii)  $\vartheta_j(z_j) \rightarrow \infty$  as  $z^{[j]} \rightarrow \infty$ , and  $\vartheta_j(z_j) \rightarrow -\infty$  as  $z^{[j]} \rightarrow -\infty$ .

- Suppose that  $Z^{[j]}$  denotes the excluded variable for option  $j$  with properties assumed in (B-2a), and suppose that: (i) the support of the distribution of  $Z^{[j]}$  conditional on all other elements of  $Z$  is the full real line; (ii)  $\vartheta_j(z_j) \rightarrow \infty$  as  $z^{[j]} \rightarrow \infty$ , and  $\vartheta_j(z_j) \rightarrow -\infty$  as  $z^{[j]} \rightarrow -\infty$ .
- Then  $\Delta_{j, \mathcal{J} \setminus j}^{\text{ATE}}(x, z)$  and  $\Delta_j^{\text{LATE}}(x, z^{[-j]}, z^{[j]}, \tilde{z}^{[j]})$  are arbitrarily close when evaluated at a sufficiently large value of  $\tilde{z}^{[j]}$  and a sufficiently small value of  $z^{[j]}$ .

- Suppose that  $Z^{[j]}$  denotes the excluded variable for option  $j$  with properties assumed in (B-2a), and suppose that: (i) the support of the distribution of  $Z^{[j]}$  conditional on all other elements of  $Z$  is the full real line; (ii)  $\vartheta_j(z_j) \rightarrow \infty$  as  $z^{[j]} \rightarrow \infty$ , and  $\vartheta_j(z_j) \rightarrow -\infty$  as  $z^{[j]} \rightarrow -\infty$ .
- Then  $\Delta_{j, \mathcal{J} \setminus j}^{\text{ATE}}(x, z)$  and  $\Delta_j^{\text{LATE}}(x, z^{[-j]}, z^{[j]}, \tilde{z}^{[j]})$  are arbitrarily close when evaluated at a sufficiently large value of  $\tilde{z}^{[j]}$  and a sufficiently small value of  $z^{[j]}$ .
- Following ?,  $\Delta_{j, \mathcal{J} \setminus j}^{\text{TT}}(x, z)$  and  $\Delta_j^{\text{LATE}}(x, z^{[-j]}, z^{[j]}, \tilde{z}^{[j]})$  are arbitrarily close for sufficiently small  $z^{[j]}$ .

- Using Theorem 6, we can use Wald estimands to identify the LATE parameters, and thus can use the Wald estimand to identify the ATE and TT parameters provided that there is sufficient support for the  $Z$ .

- Using Theorem 6, we can use Wald estimands to identify the LATE parameters, and thus can use the Wald estimand to identify the ATE and TT parameters provided that there is sufficient support for the  $Z$ .
- While this discussion has used the Wald estimands, alternatively we could also follow [?](#), as summarized in Slide 90, in expressing ATE and TT as integrated versions of MTE.

- Using Theorem 6, we can use Wald estimands to identify the LATE parameters, and thus can use the Wald estimand to identify the ATE and TT parameters provided that there is sufficient support for the  $Z$ .
- While this discussion has used the Wald estimands, alternatively we could also follow [?](#), as summarized in Slide 90, in expressing ATE and TT as integrated versions of MTE.
- By Theorem 6, we can use LIV to identify MTE and can thus express ATE and TT as integrated versions of the LIV estimand.



- For a general instrument  $J(Z^{[j]}, Z^{[-j]})$  constructed from  $(Z^{[j]}, Z^{[-j]})$ , which we denote as  $J^{[j]}$ , we can obtain a parallel construction to the characterization of standard IV given in Slide 221:

$$\Delta_{J^{[j]}}^{IV} = \int_0^1 \Delta^{\text{MTE}}(x, z, u_{D_j}) \omega_{IV}^{J^{[j]}}(u_{D_j}) du_{D_j}, \quad (58)$$

where

$$\omega_{IV}^{J^{[j]}}(u_{D_j}) = \frac{E[J^{[j]} - E(J^{[j]}) \mid P_j(Z) \geq u_{D_j}] \Pr(P_j(Z) \geq u_{D_j} \mid Z^{[-j]} = z^{[-j]})}{\text{Cov}(Z^{[j]}, D_{\mathcal{J},j})}, \quad (59)$$

where  $u_{D_j}$  is defined at the beginning of this subsection and where we keep the conditioning on  $X = x$  implicit.

- Note that from Theorem 6, we obtain that

$$\frac{\frac{\partial}{\partial z^j} E[Y | X = x, Z = z]}{\frac{\partial P_j(z)}{\partial z^j}} = \frac{\partial E[Y | X = x, Z = z]}{\partial P_j(z)}$$

$$= E[Y_j - Y_{\mathcal{J} \setminus j} | X = x, Z = z, \vartheta_j(Z_j) - V_j = R]$$

so LIV identifies MTE and linear IV is a weighted average of LIV with the weights summing to one.

- Note that from Theorem 6, we obtain that

$$\frac{\frac{\partial}{\partial z^j} E[Y | X = x, Z = z]}{\frac{\partial P_j(z)}{\partial z^j}} = \frac{\partial E[Y | X = x, Z = z]}{\partial P_j(z)}$$

$$= E[Y_j - Y_{\mathcal{J} \setminus j} | X = x, Z = z, \vartheta_j(Z_j) - V_j = R]$$

so LIV identifies MTE and linear IV is a weighted average of LIV with the weights summing to one.

- These results mirror the results established in the binary case.

- In the literature on the effects of schooling ( $S = \sum_{j \in \mathcal{J}} j D_{\mathcal{J},j}$ ) on earnings ( $Y_{\mathcal{J}}$ ), it is conventional to instrument  $S$ .

- In the literature on the effects of schooling ( $S = \sum_{j \in \mathcal{J}} j D_{\mathcal{J},j}$ ) on earnings ( $Y_{\mathcal{J}}$ ), it is conventional to instrument  $S$ .
- The website of ? presents an analysis of this case.

- In the literature on the effects of schooling ( $S = \sum_{j \in \mathcal{J}} j D_{\mathcal{J},j}$ ) on earnings ( $Y_{\mathcal{J}}$ ), it is conventional to instrument  $S$ .
- The website of ? presents an analysis of this case.
- For the general unordered case,

$$\Delta_{J^j}^{IV} = \frac{\text{Cov}(J^j, Y_{\mathcal{J}})}{\text{Cov}(J^j, S)}$$

can be decomposed into economically interpretable components where the weights can be identified but the objects being weighted cannot be identified using local instrumental variables or LATE without making large support assumptions.

- In the literature on the effects of schooling ( $S = \sum_{j \in \mathcal{J}} j D_{\mathcal{J},j}$ ) on earnings ( $Y_{\mathcal{J}}$ ), it is conventional to instrument  $S$ .
- The website of ? presents an analysis of this case.
- For the general unordered case,

$$\Delta_{j^l}^{IV} = \frac{\text{Cov}(j^l, Y_{\mathcal{J}})}{\text{Cov}(j^l, S)}$$

can be decomposed into economically interpretable components where the weights can be identified but the objects being weighted cannot be identified using local instrumental variables or LATE without making large support assumptions.

- However, the components can be identified using a structural model.

- The trick we have used in this subsection comparing outcomes in  $j$  to the next best option converts a general unordered multiple outcome model into a two-outcome setup.



- The trick we have used in this subsection comparing outcomes in  $j$  to the next best option converts a general unordered multiple outcome model into a two-outcome setup.
- This effectively partitions  $Y_{\mathcal{J}}$  into two components, as in (57).

- The trick we have used in this subsection comparing outcomes in  $j$  to the next best option converts a general unordered multiple outcome model into a two-outcome setup.
- This effectively partitions  $Y_{\mathcal{J}}$  into two components, as in (57).
- Thus we write

$$Y_{\mathcal{J}} = D_{\mathcal{J},j} Y_j + (1 - D_{\mathcal{J},j}) Y_{\mathcal{J} \setminus j},$$

where

$$Y_{\mathcal{J} \setminus j} = \sum_{\substack{\ell \neq j \\ \ell \in \mathcal{J}}} \frac{D_{\mathcal{J},\ell}}{1 - D_{\mathcal{J},j}} Y_{\ell} \cdot \mathbf{1}(D_{\mathcal{J},j} \neq 1).$$

In the more general unordered case with three or more choices, to analyze IV estimates of the effect of  $S$  on  $Y_{\mathcal{J}}$ , we must work with  $Y_{\mathcal{J}} = \sum_{k \in \mathcal{J}} D_{\mathcal{J},k} Y_k$  and make multiple comparisons across potential outcomes.

- This requires us to move outside of the LATE/LIV framework, which is inherently based on binary comparisons.

- This requires us to move outside of the LATE/LIV framework, which is inherently based on binary comparisons.
- We turn to that analysis next.

## Identification: Effect of Best Option in $\mathcal{K}$ Versus Best Option not in $\mathcal{K}$

- We just presented an analysis of identification for treatment parameters defined as averages of  $\Delta_{j, \mathcal{J} \setminus j}$ , the effect of choosing option  $j$  versus the preferred option in  $\mathcal{J}$  if  $j$  were not available.

## Identification: Effect of Best Option in $\mathcal{K}$ Versus Best Option not in $\mathcal{K}$

- We just presented an analysis of identification for treatment parameters defined as averages of  $\Delta_{j, \mathcal{J} \setminus j}$ , the effect of choosing option  $j$  versus the preferred option in  $\mathcal{J}$  if  $j$  were not available.
- We now consider identification of  $\Delta_{\mathcal{K}, \mathcal{J} \setminus \mathcal{K}}$ , the effect of choosing the preferred choice among set  $\mathcal{K}$  versus the preferred choice among  $\mathcal{J}$  if no option in  $\mathcal{K}$  were available.

## Identification: Effect of Best Option in $\mathcal{K}$ Versus Best Option not in $\mathcal{K}$

- We just presented an analysis of identification for treatment parameters defined as averages of  $\Delta_{j, \mathcal{J} \setminus j}$ , the effect of choosing option  $j$  versus the preferred option in  $\mathcal{J}$  if  $j$  were not available.
- We now consider identification of  $\Delta_{\mathcal{K}, \mathcal{J} \setminus \mathcal{K}}$ , the effect of choosing the preferred choice among set  $\mathcal{K}$  versus the preferred choice among  $\mathcal{J}$  if no option in  $\mathcal{K}$  were available.
- This is an effect where we compare sets of options, and not just a single option compared to the rest.

- We first start with an analysis that varies the  $\{\vartheta_k(\cdot)\}_{k \in \mathcal{J}}$  indices directly.



- We first start with an analysis that varies the  $\{\vartheta_k(\cdot)\}_{k \in \mathcal{J}}$  indices directly.
- This analysis would be useful if one first identifies the index function, e.g., through an identification at infinity argument using the analysis in ?, as in Appendix B of Part I or ?.

- We first start with an analysis that varies the  $\{\vartheta_k(\cdot)\}_{k \in \mathcal{J}}$  indices directly.
- This analysis would be useful if one first identifies the index function, e.g., through an identification at infinity argument using the analysis in ?, as in Appendix B of Part I or ?.
- We then perform an analysis shifting  $Z$  directly.

- We first start with an analysis that varies the  $\{\vartheta_k(\cdot)\}_{k \in \mathcal{J}}$  indices directly.
- This analysis would be useful if one first identifies the index function, e.g., through an identification at infinity argument using the analysis in ?, as in Appendix B of Part I or ?.
- We then perform an analysis shifting  $Z$  directly.
- We show that it is possible to identify MTE and LATE averages of the  $\Delta_{\mathcal{K}, \mathcal{J} \setminus \mathcal{K}}$  effect if one has knowledge of the  $\{\vartheta_k(\cdot)\}_{k \in \mathcal{J}}$  index functions but is not possible using shifts in  $Z$  without knowledge of the index functions.

- The one exception to this result is the special case already considered, when  $\mathcal{K} = k$ , i.e., the set only contains one element, in which case it is possible to identify the marginal parameters using shifts in  $Z$  directly without knowledge of the index functions.

- The one exception to this result is the special case already considered, when  $\mathcal{K} = k$ , i.e., the set only contains one element, in which case it is possible to identify the marginal parameters using shifts in  $Z$  directly without knowledge of the index functions.
- Let  $\vartheta_{\mathcal{J}}(Z)$  denote a random vector stacking the indices,

$$\vartheta_{\mathcal{J}}(Z) = \cup_{k \in \mathcal{J}} \{\vartheta_k(Z) : k \in \mathcal{J}\}.$$

- The one exception to this result is the special case already considered, when  $\mathcal{K} = k$ , i.e., the set only contains one element, in which case it is possible to identify the marginal parameters using shifts in  $Z$  directly without knowledge of the index functions.
- Let  $\vartheta_{\mathcal{J}}(Z)$  denote a random vector stacking the indices,

$$\vartheta_{\mathcal{J}}(Z) = \cup_{k \in \mathcal{J}} \{\vartheta_k(Z) : k \in \mathcal{J}\}.$$

- Let  $\vartheta_{\mathcal{J}}$  be a vector denoting a potential evaluation point of  $\vartheta_{\mathcal{J}}(Z)$ ,  $\vartheta_{\mathcal{J}} = \{\vartheta_k : k \in \mathcal{J}\}$ , so that  $\vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}$  denotes the event  $\{\vartheta_k(Z) = \vartheta_k : k \in \mathcal{J}\}$ .

- The one exception to this result is the special case already considered, when  $\mathcal{K} = k$ , i.e., the set only contains one element, in which case it is possible to identify the marginal parameters using shifts in  $Z$  directly without knowledge of the index functions.
- Let  $\vartheta_{\mathcal{J}}(Z)$  denote a random vector stacking the indices,

$$\vartheta_{\mathcal{J}}(Z) = \cup_{k \in \mathcal{J}} \{\vartheta_k(Z) : k \in \mathcal{J}\}.$$

- Let  $\vartheta_{\mathcal{J}}$  be a vector denoting a potential evaluation point of  $\vartheta_{\mathcal{J}}(Z)$ ,  $\vartheta_{\mathcal{J}} = \{\vartheta_k : k \in \mathcal{J}\}$ , so that  $\vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}$  denotes the event  $\{\vartheta_k(Z) = \vartheta_k : k \in \mathcal{J}\}$ .
- Let  $\vartheta_{\mathcal{J}} + h$  denote  $\{\vartheta_k + h : k \in \mathcal{J}\}$ , where  $h \in \mathbb{R}$ .

- We now define a version of the Wald estimand that uses the indices directly as instruments instead of using  $Z$  as instruments,

$$\begin{aligned} \tilde{\Delta}_{\mathcal{K}}^{\text{Wald}}(x, \vartheta_{\mathcal{J}}, h) &\equiv \left[ E(Y \mid X = x, \vartheta_{\mathcal{K}}(Z) = \vartheta_{\mathcal{K}} + h, \vartheta_{\mathcal{J} \setminus \mathcal{K}}(Z) = \vartheta_{\mathcal{J} \setminus \mathcal{K}}) \right. \\ &\quad \left. - E(Y \mid X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}) \right] \\ &\quad \times \left[ \Pr(I_{\mathcal{J}} \in \mathcal{K} \mid X = x, \vartheta_{\mathcal{K}}(Z) = \vartheta_{\mathcal{K}} + h, \vartheta_{\mathcal{J} \setminus \mathcal{K}}(Z) = \vartheta_{\mathcal{J} \setminus \mathcal{K}}) \right. \\ &\quad \left. - \Pr(I_{\mathcal{J}} \in \mathcal{K} \mid X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}) \right]^{-1}. \end{aligned}$$



- We now define a version of the Wald estimand that uses the indices directly as instruments instead of using  $Z$  as instruments,

$$\begin{aligned} \tilde{\Delta}_{\mathcal{K}}^{\text{Wald}}(x, \vartheta_{\mathcal{J}}, h) & \equiv \left[ E(Y \mid X = x, \vartheta_{\mathcal{K}}(Z) = \vartheta_{\mathcal{K}} + h, \vartheta_{\mathcal{J} \setminus \mathcal{K}}(Z) = \vartheta_{\mathcal{J} \setminus \mathcal{K}}) \right. \\ & \quad \left. - E(Y \mid X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}) \right] \\ & \quad \times \left[ \Pr(I_{\mathcal{J}} \in \mathcal{K} \mid X = x, \vartheta_{\mathcal{K}}(Z) = \vartheta_{\mathcal{K}} + h, \vartheta_{\mathcal{J} \setminus \mathcal{K}}(Z) = \vartheta_{\mathcal{J} \setminus \mathcal{K}}) \right. \\ & \quad \left. - \Pr(I_{\mathcal{J}} \in \mathcal{K} \mid X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}) \right]^{-1}. \end{aligned}$$

- $\tilde{\Delta}_{\mathcal{K}}^{\text{Wald}}(x, \vartheta_{\mathcal{J}}, h)$  corresponds to the effect of a shift in each index in  $\mathcal{K}$  upward by  $h$  while holding each index in  $\mathcal{J} \setminus \mathcal{K}$  constant.

- Using indices, we define a version of the LIV estimand using indices  $\tilde{\Delta}_{\mathcal{K}}^{\text{LIV}}(x, \vartheta_{\mathcal{J}})$  through a limit expression:

$$\tilde{\Delta}_{\mathcal{K}}^{\text{LIV}}(x, \vartheta_{\mathcal{J}}) = \lim_{h \rightarrow 0} \tilde{\Delta}_{\mathcal{K}}^{\text{Wald}}(x, \vartheta_{\mathcal{J}}, h).$$

- Using indices, we define a version of the LIV estimand using indices  $\tilde{\Delta}_{\mathcal{K}}^{\text{LIV}}(x, \vartheta_{\mathcal{J}})$  through a limit expression:

$$\tilde{\Delta}_{\mathcal{K}}^{\text{LIV}}(x, \vartheta_{\mathcal{J}}) = \lim_{h \rightarrow 0} \tilde{\Delta}_{\mathcal{K}}^{\text{Wald}}(x, \vartheta_{\mathcal{J}}, h).$$

- Likewise, we define versions of the LATE and MTE parameters that are functions of the  $\vartheta$  indices instead of functions of  $z$  evaluation points,

$$\tilde{\Delta}_{\mathcal{K}, \mathcal{L}}^{\text{LATE}}(x, \vartheta_{\mathcal{J}}, h) = E(\Delta_{\mathcal{K}, \mathcal{L}} | X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}, R_{\mathcal{K}}(Z) + h \geq R_{\mathcal{L}}(Z) \geq R_{\mathcal{K}}(Z))$$

$$\tilde{\Delta}_{\mathcal{K}, \mathcal{L}}^{\text{MTE}}(x, \vartheta_{\mathcal{J}}) = E(\Delta_{\mathcal{K}, \mathcal{L}} | X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}, R_{\mathcal{K}}(Z) = R_{\mathcal{L}}(Z))$$

- Using indices, we define a version of the LIV estimand using indices  $\tilde{\Delta}_{\mathcal{K}}^{\text{LIV}}(x, \vartheta_{\mathcal{J}})$  through a limit expression:

$$\tilde{\Delta}_{\mathcal{K}}^{\text{LIV}}(x, \vartheta_{\mathcal{J}}) = \lim_{h \rightarrow 0} \tilde{\Delta}_{\mathcal{K}}^{\text{Wald}}(x, \vartheta_{\mathcal{J}}, h).$$

- Likewise, we define versions of the LATE and MTE parameters that are functions of the  $\vartheta$  indices instead of functions of  $z$  evaluation points,

$$\tilde{\Delta}_{\mathcal{K}, \mathcal{L}}^{\text{LATE}}(x, \vartheta_{\mathcal{J}}, h) = E(\Delta_{\mathcal{K}, \mathcal{L}} | X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}, R_{\mathcal{K}}(Z) + h \geq R_{\mathcal{L}}(Z) \geq R_{\mathcal{K}}(Z))$$

$$\tilde{\Delta}_{\mathcal{K}, \mathcal{L}}^{\text{MTE}}(x, \vartheta_{\mathcal{J}}) = E(\Delta_{\mathcal{K}, \mathcal{L}} | X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}, R_{\mathcal{K}}(Z) = R_{\mathcal{L}}(Z))$$

- We state the following identification theorem:

## Theorem 7

1 Assume (B-1), (B-3)–(B-5), and (B-2a).

2 Then:

$$\tilde{\Delta}_{\mathcal{K}}^{Wald}(x, \vartheta_{\mathcal{J}}, h) = \tilde{\Delta}_{\mathcal{K}, \mathcal{J} \setminus \mathcal{K}}^{LATE}(x, \vartheta_{\mathcal{J}}, h),$$

3 Assume (B-1), (B-3)–(B-5), and (B-2b).

4 Then:

$$\tilde{\Delta}_{\mathcal{K}}^{LIV}(x, \vartheta_{\mathcal{J}}) = \tilde{\Delta}_{\mathcal{K}, \mathcal{J} \setminus \mathcal{K}}^{MTE}(x, \vartheta_{\mathcal{J}})$$

Proof.

Follows with trivial modifications from the proof of Theorem 6.

- Now consider the same analysis shifting  $Z$  directly instead of shifting the indices.

- Now consider the same analysis shifting  $Z$  directly instead of shifting the indices.
- First consider LATE.



- Now consider the same analysis shifting  $Z$  directly instead of shifting the indices.
- First consider LATE.
- If one knew what shifts in  $Z$  corresponded to shifting each index in  $\mathcal{K}$  upward by the same amount while holding each index in  $\mathcal{J} \setminus \mathcal{K}$  constant, then one could immediately follow the preceding analysis to recover  $E(\Delta_{\mathcal{K}, \mathcal{J} \setminus \mathcal{K}} | X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}, R_{\mathcal{K}}(Z) + h \geq R_{\mathcal{J} \setminus \mathcal{K}}(Z) \geq R_{\mathcal{K}}(Z))$ .

- Now consider the same analysis shifting  $Z$  directly instead of shifting the indices.
- First consider LATE.
- If one knew what shifts in  $Z$  corresponded to shifting each index in  $\mathcal{K}$  upward by the same amount while holding each index in  $\mathcal{J} \setminus \mathcal{K}$  constant, then one could immediately follow the preceding analysis to recover  $E(\Delta_{\mathcal{K}, \mathcal{J} \setminus \mathcal{K}} | X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}, R_{\mathcal{K}}(Z) + h \geq R_{\mathcal{J} \setminus \mathcal{K}}(Z) \geq R_{\mathcal{K}}(Z))$ .
- However, unless  $\mathcal{K}$  is a singleton, without knowledge of the index functions one does not know what shifts in  $Z$  will have this property.

- Now consider the same analysis shifting  $Z$  directly instead of shifting the indices.
- First consider LATE.
- If one knew what shifts in  $Z$  corresponded to shifting each index in  $\mathcal{K}$  upward by the same amount while holding each index in  $\mathcal{J} \setminus \mathcal{K}$  constant, then one could immediately follow the preceding analysis to recover  $E(\Delta_{\mathcal{K}, \mathcal{J} \setminus \mathcal{K}} | X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}, R_{\mathcal{K}}(Z) + h \geq R_{\mathcal{J} \setminus \mathcal{K}}(Z) \geq R_{\mathcal{K}}(Z))$ .
- However, unless  $\mathcal{K}$  is a singleton, without knowledge of the index functions one does not know what shifts in  $Z$  will have this property.
- One possible approach would be to only shift elements of  $Z$  that are elements of  $Z_j$  for  $j \in \mathcal{K}$  but are excluded from  $Z_j$  for  $j \in \mathcal{J} \setminus \mathcal{K}$ .

- However, unless the shifts move the indices for choices in  $\mathcal{K}$  all by the same amount, the shift in  $Z$  will result in movement not only from the set  $\mathcal{J} \setminus \mathcal{K}$  to the set  $\mathcal{K}$  but also cause movement between choices within  $\mathcal{K}$ .

- However, unless the shifts move the indices for choices in  $\mathcal{K}$  all by the same amount, the shift in  $Z$  will result in movement not only from the set  $\mathcal{J} \setminus \mathcal{K}$  to the set  $\mathcal{K}$  but also cause movement between choices within  $\mathcal{K}$ .
- Thus, one can use shifts in  $Z$  to recover a LATE-type parameter for  $\Delta_{\mathcal{K}, \mathcal{J} \setminus \mathcal{K}}$  only if either (i) the index functions are known, or (ii)  $\mathcal{K} = \{k\}$ , i.e., the set  $\mathcal{K}$  contains only one element.

- However, unless the shifts move the indices for choices in  $\mathcal{K}$  all by the same amount, the shift in  $Z$  will result in movement not only from the set  $\mathcal{J} \setminus \mathcal{K}$  to the set  $\mathcal{K}$  but also cause movement between choices within  $\mathcal{K}$ .
- Thus, one can use shifts in  $Z$  to recover a LATE-type parameter for  $\Delta_{\mathcal{K}, \mathcal{J} \setminus \mathcal{K}}$  only if either (i) the index functions are known, or (ii)  $\mathcal{K} = \{k\}$ , i.e., the set  $\mathcal{K}$  contains only one element.
- Our analysis establishes a fundamental role for choice theory in recovering the indices needed to perform IV analysis.

- Thus far, we have only considered identification of marginal treatment effect parameters for  $\Delta_{\mathcal{K}, \mathcal{J} \setminus \mathcal{K}}$  and not of the more standard treatment parameters ATE and TT for  $\Delta_{\mathcal{K}, \mathcal{J} \setminus \mathcal{K}}$ .

- Thus far, we have only considered identification of marginal treatment effect parameters for  $\Delta_{\mathcal{K}, \mathcal{J} \setminus \mathcal{K}}$  and not of the more standard treatment parameters ATE and TT for  $\Delta_{\mathcal{K}, \mathcal{J} \setminus \mathcal{K}}$ .
- As in the immediately preceding section, we can follow ? in expressing ATE and TT as integrated versions of MTE or show that ATE and TT can be approximated arbitrarily well by LATE parameters.



- Thus far, we have only considered identification of marginal treatment effect parameters for  $\Delta_{\mathcal{K}, \mathcal{J} \setminus \mathcal{K}}$  and not of the more standard treatment parameters ATE and TT for  $\Delta_{\mathcal{K}, \mathcal{J} \setminus \mathcal{K}}$ .
- As in the immediately preceding section, we can follow ? in expressing ATE and TT as integrated versions of MTE or show that ATE and TT can be approximated arbitrarily well by LATE parameters.
- Given appropriate support conditions, we can again identify MTE over the appropriate range or identify the appropriate LATE parameters and thus identify ATE and TT given the required support conditions.

## Identification: Effect of One Fixed Choice Versus Another

- Consider evaluating the effect of fixed option  $j$  versus fixed option  $k$ ,  $\Delta_{j,k}$ , i.e., the effect for the individual of having no choice except to choose state  $j$  versus no choice except to choose state  $k$ .

## Identification: Effect of One Fixed Choice Versus Another

- Consider evaluating the effect of fixed option  $j$  versus fixed option  $k$ ,  $\Delta_{j,k}$ , i.e., the effect for the individual of having no choice except to choose state  $j$  versus no choice except to choose state  $k$ .
- We show that it is possible to identify averages of  $\Delta_{j,k}$  if one has sufficient support conditions.

## Identification: Effect of One Fixed Choice Versus Another

- Consider evaluating the effect of fixed option  $j$  versus fixed option  $k$ ,  $\Delta_{j,k}$ , i.e., the effect for the individual of having no choice except to choose state  $j$  versus no choice except to choose state  $k$ .
- We show that it is possible to identify averages of  $\Delta_{j,k}$  if one has sufficient support conditions.
- These conditions supplement the standard IV conditions developed for the binary case (?) with the conditions more commonly used in semiparametric estimation.

- We start by considering the analysis if one knows the  $\vartheta$  index functions, say from a semiparametric analysis of discrete choice, and then show that knowledge of the  $\vartheta$  index functions is not necessary.

- We start by considering the analysis if one knows the  $\vartheta$  index functions, say from a semiparametric analysis of discrete choice, and then show that knowledge of the  $\vartheta$  index functions is not necessary.
- For notational purposes, for any  $j, k, \in \mathcal{J}$ , define  $U_{j,k} = U_j - U_k$ , and let  $\vartheta_{j,k}(Z) = \vartheta_j(Z_j) - \vartheta_k(Z_k)$ .

- We start by considering the analysis if one knows the  $\vartheta$  index functions, say from a semiparametric analysis of discrete choice, and then show that knowledge of the  $\vartheta$  index functions is not necessary.
- For notational purposes, for any  $j, k, \in \mathcal{J}$ , define  $U_{j,k} = U_j - U_k$ , and let  $\vartheta_{j,k}(Z) = \vartheta_j(Z_j) - \vartheta_k(Z_k)$ .
- One might try to follow our previous strategy to identify treatment parameters for  $\Delta_{j,k}$  if one could shift  $\vartheta_j - \vartheta_k = \vartheta_{j,k}$  while holding constant  $\{\vartheta_{l,m}\}_{(l,m) \in \mathcal{J} \times \mathcal{J} \setminus \{j,k\}}$ , i.e., while holding all other utility contrasts fixed.

- We start by considering the analysis if one knows the  $\vartheta$  index functions, say from a semiparametric analysis of discrete choice, and then show that knowledge of the  $\vartheta$  index functions is not necessary.
- For notational purposes, for any  $j, k, \in \mathcal{J}$ , define  $U_{j,k} = U_j - U_k$ , and let  $\vartheta_{j,k}(Z) = \vartheta_j(Z_j) - \vartheta_k(Z_k)$ .
- One might try to follow our previous strategy to identify treatment parameters for  $\Delta_{j,k}$  if one could shift  $\vartheta_j - \vartheta_k = \vartheta_{j,k}$  while holding constant  $\{\vartheta_{l,m}\}_{(l,m) \in \mathcal{J} \times \mathcal{J} \setminus \{j,k\}}$ , i.e., while holding all other utility contrasts fixed.
- However, given the structure of the latent variable model determining choices, these are incompatible conditions.



- We start by considering the analysis if one knows the  $\vartheta$  index functions, say from a semiparametric analysis of discrete choice, and then show that knowledge of the  $\vartheta$  index functions is not necessary.
- For notational purposes, for any  $j, k \in \mathcal{J}$ , define  $U_{j,k} = U_j - U_k$ , and let  $\vartheta_{j,k}(Z) = \vartheta_j(Z_j) - \vartheta_k(Z_k)$ .
- One might try to follow our previous strategy to identify treatment parameters for  $\Delta_{j,k}$  if one could shift  $\vartheta_j - \vartheta_k = \vartheta_{j,k}$  while holding constant  $\{\vartheta_{l,m}\}_{(l,m) \in \mathcal{J} \times \mathcal{J} \setminus \{j,k\}}$ , i.e., while holding all other utility contrasts fixed.
- However, given the structure of the latent variable model determining choices, these are incompatible conditions.
- To see this, note that  $\vartheta_{j,k} = \vartheta_{l,k} - \vartheta_{l,j}$  for any  $l$ , and thus  $\vartheta_{j,k}$  cannot be shifted while holding  $\vartheta_{l,j}$  and  $\vartheta_{l,k}$  constant.

- We start by considering the analysis if one knows the  $\vartheta$  index functions, say from a semiparametric analysis of discrete choice, and then show that knowledge of the  $\vartheta$  index functions is not necessary.
- For notational purposes, for any  $j, k, \in \mathcal{J}$ , define  $U_{j,k} = U_j - U_k$ , and let  $\vartheta_{j,k}(Z) = \vartheta_j(Z_j) - \vartheta_k(Z_k)$ .
- One might try to follow our previous strategy to identify treatment parameters for  $\Delta_{j,k}$  if one could shift  $\vartheta_j - \vartheta_k = \vartheta_{j,k}$  while holding constant  $\{\vartheta_{l,m}\}_{(l,m) \in \mathcal{J} \times \mathcal{J} \setminus \{j,k\}}$ , i.e., while holding all other utility contrasts fixed.
- However, given the structure of the latent variable model determining choices, these are incompatible conditions.
- To see this, note that  $\vartheta_{j,k} = \vartheta_{l,k} - \vartheta_{l,j}$  for any  $l$ , and thus  $\vartheta_{j,k}$  cannot be shifted while holding  $\vartheta_{l,j}$  and  $\vartheta_{l,k}$  constant.

- To bypass this problem, we develop a limit strategy to make the consequences of shifting  $\vartheta_{j,k}$  negligible.

- To bypass this problem, we develop a limit strategy to make the consequences of shifting  $\vartheta_{j,k}$  negligible.
- Our strategy relies on an identification at infinity argument.

- To bypass this problem, we develop a limit strategy to make the consequences of shifting  $\vartheta_{j,k}$  negligible.
- Our strategy relies on an identification at infinity argument.
- For example, consider the case where  $\mathcal{J} = \{1, 2, 3\}$ , and consider identification of the MTE parameter for option 3 versus option 1.

- To bypass this problem, we develop a limit strategy to make the consequences of shifting  $\vartheta_{j,k}$  negligible.
- Our strategy relies on an identification at infinity argument.
- For example, consider the case where  $\mathcal{J} = \{1, 2, 3\}$ , and consider identification of the MTE parameter for option 3 versus option 1.
- Recall that  $D_{\mathcal{J}\setminus 3,l}$  is an indicator variable for whether option  $l$  would be chosen if option 3 were not available, so that  $D_{\mathcal{J}\setminus 3,l}\Delta_{3,\mathcal{J}\setminus 3} = D_{\mathcal{J}\setminus 3,l}\Delta_{3,l}$ .

- Since 1 and 2 are the only options if 3 is not available, it follows that  $\Delta_{3,\mathcal{J}\setminus 3} = D_{\mathcal{J}\setminus 3,1}\Delta_{3,1} + D_{\mathcal{J}\setminus 3,2}\Delta_{3,2}$ , and we have that

$$\begin{aligned} & E(\Delta_{3,\mathcal{J}\setminus 3} \mid X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}, R_3(Z) = R_{\mathcal{J}\setminus 3}(Z)) \\ &= E(D_{\mathcal{J}\setminus 3,1}\Delta_{3,1} \mid X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}, R_3(Z) = R_{\mathcal{J}\setminus 3}(Z)) \\ &+ E(D_{\mathcal{J}\setminus 3,2}\Delta_{3,2} \mid X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}, R_3(Z) = R_{\mathcal{J}\setminus 3}(Z)). \end{aligned}$$

The smaller  $\vartheta_2$  is (holding  $\vartheta_1$  and  $\vartheta_3$  fixed), the larger the probability that the “next best option” is 1 and not 2.

- Note that  $E(\Delta_{3,1} \mid X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}, R_3(Z) = R_1(Z))$  does not depend on the  $\vartheta_2$  evaluation point given independence assumption (B-1), so that

$$\begin{aligned} E(\Delta_{3,1} \mid X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}, R_3(Z) = R_1(Z)) \\ = E(\Delta_{3,1} \mid X = x, \vartheta_{\mathcal{J} \setminus 2}(Z) = \vartheta_{\mathcal{J} \setminus 2}, R_3(Z) = R_1(Z)). \end{aligned}$$



- Thus, by assumptions (B-1) and (B-3) and the Dominated Convergence Theorem, we have that

$$\begin{aligned} \lim_{\vartheta_2 \rightarrow -\infty} E(D_{\mathcal{J}\setminus 3,1} \Delta_{3,1} \mid X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}, R_3(Z) = R_{\mathcal{J}\setminus 3}(Z)) \\ = E(\Delta_{3,1} \mid X = x, \vartheta_{\mathcal{J}\setminus 2}(Z) = \vartheta_{\mathcal{J}\setminus 2}, R_3(Z) = R_1(Z)) \end{aligned}$$

while

$$\lim_{\vartheta_2 \rightarrow -\infty} E(D_{\mathcal{J}\setminus 3,2} \Delta_{3,2} \mid X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}, R_3(Z) = R_{\mathcal{J}\setminus 3}(Z)) = 0,$$

so that

$$\begin{aligned} \lim_{\vartheta_2 \rightarrow -\infty} E(\Delta_{3,\mathcal{J}\setminus 3} \mid X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}, R_3(Z) = R_{\mathcal{J}\setminus 3}(Z)) \\ = E(\Delta_{3,1} \mid X = x, \vartheta_{\mathcal{J}\setminus 2}(Z) = \vartheta_{\mathcal{J}\setminus 2}, R_3(Z) = R_1(Z)). \end{aligned}$$

- In other words, as the value of option 2 becomes arbitrarily small, the probability of the “next best option” being 1 becomes arbitrarily close to one.

- In other words, as the value of option 2 becomes arbitrarily small, the probability of the “next best option” being 1 becomes arbitrarily close to one.
- Thus the MTE parameter for option 3 versus the next best option becomes arbitrarily close to the MTE parameter for option 3 versus option 1.

- In other words, as the value of option 2 becomes arbitrarily small, the probability of the “next best option” being 1 becomes arbitrarily close to one.
- Thus the MTE parameter for option 3 versus the next best option becomes arbitrarily close to the MTE parameter for option 3 versus option 1.
- We can identify the MTE parameter for option 3 versus the next best option using the LIV estimand as in Theorem 6, and thus conditioning on  $\vartheta_2$  arbitrarily small we have that the LIV estimand is arbitrarily close to the MTE parameter for option 3 versus option 1.

- This analysis requires the appropriate support conditions in order for the limit operations to be well defined.

- This analysis requires the appropriate support conditions in order for the limit operations to be well defined.
- The following Theorem formalizes this idea, and is for the more general case where  $\mathcal{J}$  is a general finite set.

## Theorem 8

Assume (B-1), (B-3)–(B-5), and (B-2b). Assume that, for any  $t \in \mathbb{R}$ ,

$$\Pr(\vartheta_l(Z_l) \leq t | \vartheta_j(Z_j), \vartheta_k(Z_k)) \geq 0 \quad \forall l \in \mathcal{J} \setminus \{j, k\}.$$

Then

$$\lim_{\max_{l \in \mathcal{J} \setminus \{j, k\}} \{\vartheta_l\} \rightarrow -\infty} \tilde{\Delta}_j^{LIV}(x, \vartheta_{\mathcal{J}}) = E(\Delta_{j,k} | X = x, \vartheta_{j,k}(Z) = \vartheta_{j,k}, R_j(Z) = R_k(Z))$$

for any

$$x \in \lim_{t \rightarrow -\infty} \text{Supp}(X | \vartheta_j(Z_j) = \vartheta_j, \vartheta_k(Z_k) = \vartheta_k, \max_{l \in \mathcal{J} \setminus \{j, k\}} \{\vartheta_l(Z)\} \leq t).$$

## Proof.

By a trivial modification to the proof of Theorem 6, we have that  $\tilde{\Delta}_j^{\text{LIV}}(x, \vartheta_{\mathcal{J}}) = E(\Delta_{j, \mathcal{J} \setminus j} | X = x, \vartheta_{\mathcal{J}}(Z) = \vartheta_{\mathcal{J}}, R_j(Z) = R_{\mathcal{J} \setminus j}(Z))$ . The remainder of the proof follows from an immediate extension of the 3-option case just analyzed. □



- Thus, for  $x$  values in the appropriate limit support, we can approximate  $E(\Delta_{j,k} | X = x, \vartheta_{\{j,k\}}(Z) = \vartheta_{\{j,k\}}, R_j(z) = R_k(z))$  arbitrarily well by  $\Delta_j^{\text{LIV}}(x, \vartheta_{\mathcal{J}})$  for an arbitrarily small  $\max_{l \in \mathcal{J} \setminus \{j,k\}} \{\vartheta_l\}$ .

- Thus, for  $x$  values in the appropriate limit support, we can approximate  $E(\Delta_{j,k} | X = x, \vartheta_{\{j,k\}}(Z) = \vartheta_{\{j,k\}}, R_j(z) = R_k(z))$  arbitrarily well by  $\Delta_j^{\text{LIV}}(x, \vartheta_{\mathcal{J}})$  for an arbitrarily small  $\max_{l \in \mathcal{J} \setminus \{j,k\}} \{\vartheta_l\}$ .
- This analysis uses the  $\vartheta$  index functions directly, but the results can be restated without using the  $\vartheta$  functions directly.

- Thus, for  $x$  values in the appropriate limit support, we can approximate  $E(\Delta_{j,k} | X = x, \vartheta_{\{j,k\}}(Z) = \vartheta_{\{j,k\}}, R_j(z) = R_k(z))$  arbitrarily well by  $\Delta_j^{\text{LIV}}(x, \vartheta_{\mathcal{J}})$  for an arbitrarily small  $\max_{l \in \mathcal{J} \setminus \{j,k\}} \{\vartheta_l\}$ .
- This analysis uses the  $\vartheta$  index functions directly, but the results can be restated without using the  $\vartheta$  functions directly.
- Again consider the three-choice example.

- Thus, for  $x$  values in the appropriate limit support, we can approximate  $E(\Delta_{j,k} | X = x, \vartheta_{\{j,k\}}(Z) = \vartheta_{\{j,k\}}, R_j(z) = R_k(z))$  arbitrarily well by  $\Delta_j^{\text{LIV}}(x, \vartheta_{\mathcal{J}})$  for an arbitrarily small  $\max_{l \in \mathcal{J} \setminus \{j,k\}} \{\vartheta_l\}$ .
- This analysis uses the  $\vartheta$  index functions directly, but the results can be restated without using the  $\vartheta$  functions directly.
- Again consider the three-choice example.
- The central aspect of the identification strategy is to “zero-out” the second choice by making  $\vartheta_2$  arbitrarily small, allowing one to then use the LIV estimand to identify the MTE parameter for the first option versus the third as if the second choice were not an option.

- If we do not know the  $v_2$  function, we cannot condition on it.

- If we do not know the  $\vartheta_2$  function, we cannot condition on it.
- However, if we know that  $\vartheta_2$  is decreasing in a particular element of  $Z$ , say  $Z^{[j]}$ , where  $Z^{[j]}$  does not enter the index function for choices 1 and 3 and where  $\vartheta_2(z_2) \rightarrow 0$  as  $z^{[j]} \rightarrow -\infty$ , then we can follow the same strategy as if we knew the  $\vartheta_2$  index except we condition on  $Z^{[j]}$  being small instead of conditioning on  $\vartheta_2$  being small.

- If we do not know the  $\vartheta_2$  function, we cannot condition on it.
- However, if we know that  $\vartheta_2$  is decreasing in a particular element of  $Z$ , say  $Z^{[j]}$ , where  $Z^{[j]}$  does not enter the index function for choices 1 and 3 and where  $\vartheta_2(z_2) \rightarrow 0$  as  $z^{[j]} \rightarrow -\infty$ , then we can follow the same strategy as if we knew the  $\vartheta_2$  index except we condition on  $Z^{[j]}$  being small instead of conditioning on  $\vartheta_2$  being small.
- The idea naturally extends to the case of more than three options.

- We can follow ? in following a two step identification strategy for ATE and TT parameters of  $\Delta_{j,k}$ .



- We can follow ? in following a two step identification strategy for ATE and TT parameters of  $\Delta_{j,k}$ .
- We first identify the appropriate MTE or LATE parameters and then use them to identify ATE and TT given the appropriate support conditions.

- We can follow ? in following a two step identification strategy for ATE and TT parameters of  $\Delta_{j,k}$ .
- We first identify the appropriate MTE or LATE parameters and then use them to identify ATE and TT given the appropriate support conditions.
- Notice that the required support conditions are now stronger than those required for the ATE and TT parameters of  $\Delta_{j,\mathcal{J}\setminus j}$ .

- We can follow ? in following a two step identification strategy for ATE and TT parameters of  $\Delta_{j,k}$ .
- We first identify the appropriate MTE or LATE parameters and then use them to identify ATE and TT given the appropriate support conditions.
- Notice that the required support conditions are now stronger than those required for the ATE and TT parameters of  $\Delta_{j,\mathcal{J}_j}$ .
- For identification of the ATE and TT parameters of  $\Delta_{j,\mathcal{J}_j}$ , we require a large support assumption only on the  $j$ th index.

- We can follow ? in following a two step identification strategy for ATE and TT parameters of  $\Delta_{j,k}$ .
- We first identify the appropriate MTE or LATE parameters and then use them to identify ATE and TT given the appropriate support conditions.
- Notice that the required support conditions are now stronger than those required for the ATE and TT parameters of  $\Delta_{j,\mathcal{J}_j}$ .
- For identification of the ATE and TT parameters of  $\Delta_{j,\mathcal{J}_j}$ , we require a large support assumption only on the  $j$ th index.
- In particular, we require that it be possible to condition on  $Z$  values that make  $\vartheta_j$  arbitrarily small or arbitrarily large while holding the remaining indices fixed.

- In contrast, for identification of the ATE and TT parameters of  $\Delta_{j,k}$ , we require a large support assumption on each index.

- In contrast, for identification of the ATE and TT parameters of  $\Delta_{j,k}$ , we require a large support assumption on each index.
- We require that for each index we can condition on  $Z$  values that make the index arbitrarily small or arbitrarily large while holding the remaining indices fixed.

- In contrast, for identification of the ATE and TT parameters of  $\Delta_{j,k}$ , we require a large support assumption on each index.
- We require that for each index we can condition on  $Z$  values that make the index arbitrarily small or arbitrarily large while holding the remaining indices fixed.
- The reason for this stronger condition is that for identification of  $\Delta_{j,k}$  we need to use an identification at infinity strategy on all but the  $j$  and  $k$  indices to even obtain the marginal parameters.

- In contrast, for identification of the ATE and TT parameters of  $\Delta_{j,k}$ , we require a large support assumption on each index.
- We require that for each index we can condition on  $Z$  values that make the index arbitrarily small or arbitrarily large while holding the remaining indices fixed.
- The reason for this stronger condition is that for identification of  $\Delta_{j,k}$  we need to use an identification at infinity strategy on all but the  $j$  and  $k$  indices to even obtain the marginal parameters.
- We then need an additional identification at infinity step to use the marginal parameters to recover the ATE and TT parameters.



## Summarizing the Results for the Unordered Model

- We have obtained the following results on the unordered choice model in this section:

- $E(\Delta_{j, \mathcal{J} \setminus j} | X = x, Z = z, R_j(z) = R_{\mathcal{J} \setminus j}(z))$  and  $E(\Delta_{j, \mathcal{J} \setminus j} | X = x, Z = z, R_j(\tilde{z}) \geq R_{\mathcal{J} \setminus j}(\tilde{z}) \geq R_j(z))$  can be identified without a limit argument.
- $E(\Delta_{j,k} | X = x, \{\vartheta_k\}_{k \in \mathcal{J}}, R_j(z) = R_k(z))$  and  $E(\Delta_{j,k} | X = x, \{\vartheta_k\}_{k \in \mathcal{J}}, R_j(\tilde{z}) \geq R_k(\tilde{z}) \geq R_j(z))$  can be identified with a limit argument on each index in  $\mathcal{J} \setminus \{j, k\}$ .
- $\Delta_{j, \mathcal{J} \setminus j}^{\text{ATE}}(x, z)$  and  $\Delta_{j, \mathcal{J} \setminus j}^{\text{TT}}(x, z)$  can be identified with a limit argument using the  $\vartheta_j$  index.
- $\Delta_{j,k}^{\text{ATE}}(x, z)$  and  $\Delta_{j,k}^{\text{TT}}(x, z)$  can be identified with a limit argument using each index.

- These results establish the central role of choice theory (via  $\{\vartheta_k\}_{k \in \mathcal{J}}$ ) and identification at infinity in using an IV strategy to identify a variety of treatment parameters and their extensions to a general multiple choice model.

- These results establish the central role of choice theory (via  $\{\vartheta_k\}_{k \in \mathcal{J}}$ ) and identification at infinity in using an IV strategy to identify a variety of treatment parameters and their extensions to a general multiple choice model.
- Our analysis extends the analysis of ordered outcome models developed in the preceding section to a general unordered case.

- These results establish the central role of choice theory (via  $\{\vartheta_k\}_{k \in \mathcal{J}}$ ) and identification at infinity in using an IV strategy to identify a variety of treatment parameters and their extensions to a general multiple choice model.
- Our analysis extends the analysis of ordered outcome models developed in the preceding section to a general unordered case.
- Local instrumental variables identify the marginal treatment effect corresponding to the effect of one option versus the best alternative option without requiring large support assumptions or knowledge of the parameters of the choice model.

- This result preserves the spirit of the ? LATE analysis and the analysis of ??? . More generally, LIV can provide identification of the marginal treatment effect corresponding to the effect of choosing between one choice set versus not having that choice set available.

- This result preserves the spirit of the ? LATE analysis and the analysis of ??? . More generally, LIV can provide identification of the marginal treatment effect corresponding to the effect of choosing between one choice set versus not having that choice set available.
- However, identification of the more general parameters requires knowledge (identification) of the structural, latent index functions of the multinomial choice model.

- This result preserves the spirit of the ? LATE analysis and the analysis of ??? . More generally, LIV can provide identification of the marginal treatment effect corresponding to the effect of choosing between one choice set versus not having that choice set available.
- However, identification of the more general parameters requires knowledge (identification) of the structural, latent index functions of the multinomial choice model.
- LIV can also provide identification of the effect of one specified choice versus another, requiring large support assumptions but not knowledge of the latent index functions.



- In order to identify some treatment parameters, we require identification of the latent index functions generating the multinomial choice model or else having large support assumptions.

- In order to identify some treatment parameters, we require identification of the latent index functions generating the multinomial choice model or else having large support assumptions.
- This connects the LIV analysis in this paper to the more ambitious but demanding identification conditions for the full multinomial selection model developed in ?, ?, and appendix B of Part I.

- In order to identify some treatment parameters, we require identification of the latent index functions generating the multinomial choice model or else having large support assumptions.
- This connects the LIV analysis in this paper to the more ambitious but demanding identification conditions for the full multinomial selection model developed in ?, ?, and appendix B of Part I.
- We next develop the case of the continuum of outcomes.

## Continuous Treatment

- Thus far we have considered the case of a treatment variable taking a finite number of values.

## Continuous Treatment

- Thus far we have considered the case of a treatment variable taking a finite number of values.
- Now consider the case where the treatment variable  $D$  can take a continuum of values.

## Continuous Treatment

- Thus far we have considered the case of a treatment variable taking a finite number of values.
- Now consider the case where the treatment variable  $D$  can take a continuum of values.
- Suppose that

$$\begin{aligned} Y &= \mu(D, X, U) \\ D &= \vartheta(Z, V), \end{aligned}$$

with  $D$  a continuous random variable.

## Continuous Treatment

- Thus far we have considered the case of a treatment variable taking a finite number of values.
- Now consider the case where the treatment variable  $D$  can take a continuum of values.
- Suppose that

$$\begin{aligned} Y &= \mu(D, X, U) \\ D &= \vartheta(Z, V), \end{aligned}$$

with  $D$  a continuous random variable.

- We do not in general need to restrict  $U$  or  $V$  to be scalar random variables.

- We can rewrite this model in potential outcome notation by defining

$$Y_d \equiv \mu_d(X, U) \equiv \mu(d, X, U).$$

For ease of exposition, we will assume that  $X$  is exogenous in addition to  $Z$  being exogenous, so that  $(X, Z) \perp\!\!\!\perp (U, V)$ .



- We can rewrite this model in potential outcome notation by defining

$$Y_d \equiv \mu_d(X, U) \equiv \mu(d, X, U).$$

For ease of exposition, we will assume that  $X$  is exogenous in addition to  $Z$  being exogenous, so that  $(X, Z) \perp\!\!\!\perp (U, V)$ .

- We assume that  $\mu(d, x, u)$  is continuous in its first argument.

- We can rewrite this model in potential outcome notation by defining

$$Y_d \equiv \mu_d(X, U) \equiv \mu(d, X, U).$$

For ease of exposition, we will assume that  $X$  is exogenous in addition to  $Z$  being exogenous, so that  $(X, Z) \perp\!\!\!\perp (U, V)$ .

- We assume that  $\mu(d, x, u)$  is continuous in its first argument.
- Equivalently, we assume that  $\{Y_d\}$  is continuous in  $d$  for any realization.

- We can rewrite this model in potential outcome notation by defining

$$Y_d \equiv \mu_d(X, U) \equiv \mu(d, X, U).$$

For ease of exposition, we will assume that  $X$  is exogenous in addition to  $Z$  being exogenous, so that  $(X, Z) \perp\!\!\!\perp (U, V)$ .

- We assume that  $\mu(d, x, u)$  is continuous in its first argument.
- Equivalently, we assume that  $\{Y_d\}$  is continuous in  $d$  for any realization.
- Implicit in the continuity assumption is an ordering, that two treatments that are close to one another have associated outcomes that are close to one another.

- We can rewrite this model in potential outcome notation by defining

$$Y_d \equiv \mu_d(X, U) \equiv \mu(d, X, U).$$

For ease of exposition, we will assume that  $X$  is exogenous in addition to  $Z$  being exogenous, so that  $(X, Z) \perp\!\!\!\perp (U, V)$ .

- We assume that  $\mu(d, x, u)$  is continuous in its first argument.
- Equivalently, we assume that  $\{Y_d\}$  is continuous in  $d$  for any realization.
- Implicit in the continuity assumption is an ordering, that two treatments that are close to one another have associated outcomes that are close to one another.
- The restriction is qualitatively different from any restriction we have considered thus far.

- In the previous sections, there are no restrictions connecting  $Y_d$  to  $Y_{d'}$ .

- In the previous sections, there are no restrictions connecting  $Y_d$  to  $Y_{d'}$ .
- Equivalently, there are no restrictions connecting  $\mu_d(X, U)$  and  $\mu_{d'}(X, U)$ .

- In the previous sections, there are no restrictions connecting  $Y_d$  to  $Y_{d'}$ .
- Equivalently, there are no restrictions connecting  $\mu_d(X, U)$  and  $\mu_{d'}(X, U)$ .
- In the case of a continuum of treatments, we now tightly link counterfactual values that correspond to treatments that are close to one another.

- In the previous sections, there are no restrictions connecting  $Y_d$  to  $Y_{d'}$ .
- Equivalently, there are no restrictions connecting  $\mu_d(X, U)$  and  $\mu_{d'}(X, U)$ .
- In the case of a continuum of treatments, we now tightly link counterfactual values that correspond to treatments that are close to one another.
- The literature analyzing continuous endogenous regressors often defines the object of interest not as a treatment effect but instead as the “Average Structural Function” (ASF).



- Following ?, the ASF is defined as:

$$\mu(d, x) = E(Y_d | X = x) = \int \mu(d, x, u) dF_U(u)$$

In other words, the ASF is defined as the average value of  $Y$  that would result from assigning treatment  $d$  to all individuals with  $X = x$ .

- Following ?, the ASF is defined as:

$$\mu(d, x) = E(Y_d | X = x) = \int \mu(d, x, u) dF_U(u)$$

In other words, the ASF is defined as the average value of  $Y$  that would result from assigning treatment  $d$  to all individuals with  $X = x$ .

- If  $D$  is endogenous, the ASF does not in general equal the conditional expected value of  $Y$  in the data,  $E(Y_d | X = x) \neq E(Y | D = d, X = x)$ , since  $\int \mu(d, x, u) dF_U(u) \neq \int \mu(d, x, u) dF_{U|X,D}(u|x, d)$ .

- Following ?, the ASF is defined as:

$$\mu(d, x) = E(Y_d | X = x) = \int \mu(d, x, u) dF_U(u)$$

In other words, the ASF is defined as the average value of  $Y$  that would result from assigning treatment  $d$  to all individuals with  $X = x$ .

- If  $D$  is endogenous, the ASF does not in general equal the conditional expected value of  $Y$  in the data,  $E(Y_d | X = x) \neq E(Y | D = d, X = x)$ , since  $\int \mu(d, x, u) dF_U(u) \neq \int \mu(d, x, u) dF_{U|X,D}(u|x, d)$ .
- This is just a version of the distinction between fixing and conditioning introduced in ? and discussed in Part I.

- Instead of working with the ASF, we can follow the lead of ? and define treatment effect parameters for a continuous treatment.

- Instead of working with the ASF, we can follow the lead of ? and define treatment effect parameters for a continuous treatment.
- Suppose that  $\mu(d, x, u)$  is differentiable in  $d$  for any  $(x, u)$ .

- Instead of working with the ASF, we can follow the lead of ? and define treatment effect parameters for a continuous treatment.
- Suppose that  $\mu(d, x, u)$  is differentiable in  $d$  for any  $(x, u)$ .
- We can define the average treatment effect as

$$\Delta_d^{\text{ATE}}(x) = E\left(\frac{\partial}{\partial d} Y_d | X = x\right) = \int \frac{\partial}{\partial d} \mu(d, x, u) dF_U(u),$$

which is the average effect of a marginal increase in in the treatment if individuals were randomly assigned treatment level  $d$ .

- Instead of working with the ASF, we can follow the lead of ? and define treatment effect parameters for a continuous treatment.
- Suppose that  $\mu(d, x, u)$  is differentiable in  $d$  for any  $(x, u)$ .
- We can define the average treatment effect as

$$\Delta_d^{\text{ATE}}(x) = E\left(\frac{\partial}{\partial d} Y_d | X = x\right) = \int \frac{\partial}{\partial d} \mu(d, x, u) dF_U(u),$$

which is the average effect of a marginal increase in in the treatment if individuals were randomly assigned treatment level  $d$ .

- Note that in this expression the average treatment effect depends on the base treatment level,  $d$ , and for any of the continuum of possible base treatment levels we have a different average treatment effect.

- The average treatment effect is the derivative of the Blundell and Powell ASF:

$$\Delta_d^{\text{ATE}}(x) = \frac{\partial}{\partial d} \mu(d, x).$$

? define treatment on the treated as

$$\begin{aligned} \Delta_d^{\text{TT}}(x) &= E\left(\frac{\partial}{\partial d_1} Y_{d_1} \mid D = d_2, X = x\right) \Bigg|_{d=d_1=d_2} \\ &= \int \left[ \frac{\partial}{\partial d_1} \mu(d_1, x, u) \Bigg|_{d=d_1} \right] dF_{U|X,D}(u|x, d). \end{aligned}$$

which is the average effect among those currently choosing treatment level  $d$  of an incremental increase in the treatment while leaving their unobservables fixed.



- Likewise, define the marginal treatment effect as

$$\begin{aligned}\Delta_d^{\text{MTE}}(x, v) &= E\left(\frac{\partial}{\partial d} Y_d \mid V = v, X = x\right) \\ &= \int \frac{\partial}{\partial d} \mu(d, x, u) dF_{U|V}(u|v).\end{aligned}$$

- Likewise, define the marginal treatment effect as

$$\begin{aligned}\Delta_d^{\text{MTE}}(x, v) &= E\left(\frac{\partial}{\partial d} Y_d \mid V = v, X = x\right) \\ &= \int \frac{\partial}{\partial d} \mu(d, x, u) dF_{U|V}(u|v).\end{aligned}$$

- To illustrate these definitions, suppose  $D$  is schooling level measured as a continuous variable, and suppose  $Y$  is wages.

- Likewise, define the marginal treatment effect as

$$\begin{aligned}\Delta_d^{\text{MTE}}(x, v) &= E\left(\frac{\partial}{\partial d} Y_d \mid V = v, X = x\right) \\ &= \int \frac{\partial}{\partial d} \mu(d, x, u) dF_{U|V}(u|v).\end{aligned}$$

- To illustrate these definitions, suppose  $D$  is schooling level measured as a continuous variable, and suppose  $Y$  is wages.
- Then, e.g.,  $Y_{12}$  would be the potential wage corresponding to receiving exactly 12 years of schooling and  $\mu_{12} = E(Y_{12})$  is the average wage if individuals were exogenously assigned exactly 12 years of schooling.

- $\Delta_{12}^{ATE}$  is the average effect on wages of being assigned marginally more than 12 years of schooling versus being assigned exactly 12 years of schooling, and  $\Delta_{12}^{TT}$  would be the average effect of obtaining marginally more schooling for those who self-select to obtain exactly 12 years of schooling.

- $\Delta_{12}^{ATE}$  is the average effect on wages of being assigned marginally more than 12 years of schooling versus being assigned exactly 12 years of schooling, and  $\Delta_{12}^{TT}$  would be the average effect of obtaining marginally more schooling for those who self-select to obtain exactly 12 years of schooling.
- One approach to identification of the treatment parameters is to impose more structure on the outcome equation while allowing the treatment selection equation to be unspecified.

- $\Delta_{12}^{ATE}$  is the average effect on wages of being assigned marginally more than 12 years of schooling versus being assigned exactly 12 years of schooling, and  $\Delta_{12}^{TT}$  would be the average effect of obtaining marginally more schooling for those who self-select to obtain exactly 12 years of schooling.
- One approach to identification of the treatment parameters is to impose more structure on the outcome equation while allowing the treatment selection equation to be unspecified.
- The nonparametric instrumental variable approach of Angriston, Krueger, and Angriston requires that the unobservables in the outcome equation ( $U$ ) be a scalar random variable and that the outcome be an additive function of the unobservables — Angriston surveys this literature.

- Their additivity assumption imposes the restriction of no treatment effect heterogeneity (conditional on  $X$ ), so that all treatment effect parameters coincide.

- Their additivity assumption imposes the restriction of no treatment effect heterogeneity (conditional on  $X$ ), so that all treatment effect parameters coincide.
- In exchange for this restriction on the outcome equation, they do not require any structure on the first stage equation so that  $D$  does not need to be increasing in  $V$  and  $V$  is not required to be a scalar random variable.



- Their additivity assumption imposes the restriction of no treatment effect heterogeneity (conditional on  $X$ ), so that all treatment effect parameters coincide.
- In exchange for this restriction on the outcome equation, they do not require any structure on the first stage equation so that  $D$  does not need to be increasing in  $V$  and  $V$  is not required to be a scalar random variable.
- Furthermore, they only require that  $U$  be mean independent of  $(X, Z)$ , not that  $(U, V)$  be fully independent of  $(X, Z)$ .

- The additive error term assumption is relaxed by ?, who impose the stronger requirement that the outcome is a strictly increasing function of the error term (i.e.,  $\mu(x, d, u)$  strictly increasing in  $u$ ), while strengthening the required independence property to be  $(Z, X) \perp\!\!\!\perp U$ .

- The additive error term assumption is relaxed by ?, who impose the stronger requirement that the outcome is a strictly increasing function of the error term (i.e.,  $\mu(x, d, u)$  strictly increasing in  $u$ ), while strengthening the required independence property to be  $(Z, X) \perp\!\!\!\perp U$ .
- The restriction of a scalar error term with the outcome strictly increasing in this error term is again a strong restriction on the forms of treatment effect heterogeneity that are possible in the model.

- The additive error term assumption is relaxed by ?, who impose the stronger requirement that the outcome is a strictly increasing function of the error term (i.e.,  $\mu(x, d, u)$  strictly increasing in  $u$ ), while strengthening the required independence property to be  $(Z, X) \perp\!\!\!\perp U$ .
- The restriction of a scalar error term with the outcome strictly increasing in this error term is again a strong restriction on the forms of treatment effect heterogeneity that are possible in the model.
- Suppress  $X$  for ease of exposition.

- The additive error term assumption is relaxed by ?, who impose the stronger requirement that the outcome is a strictly increasing function of the error term (i.e.,  $\mu(x, d, u)$  strictly increasing in  $u$ ), while strengthening the required independence property to be  $(Z, X) \perp\!\!\!\perp U$ .
- The restriction of a scalar error term with the outcome strictly increasing in this error term is again a strong restriction on the forms of treatment effect heterogeneity that are possible in the model.
- Suppress  $X$  for ease of exposition.
- Under their restriction, if  $\mu(d, u) > \mu(d, u')$  at some treatment level  $d$ , then  $\mu(\tilde{d}, u) > \mu(\tilde{d}, u')$  for all treatment levels  $\tilde{d}$ .

- In other words, if individual one has a higher potential outcome at some value of the treatment than a second individual, then that first individual has a higher potential outcome for any value of the treatment than the second individual.

- In other words, if individual one has a higher potential outcome at some value of the treatment than a second individual, then that first individual has a higher potential outcome for any value of the treatment than the second individual.
- Under this restriction, treatment cannot change the rank ordering of outcomes across individuals.

- In other words, if individual one has a higher potential outcome at some value of the treatment than a second individual, then that first individual has a higher potential outcome for any value of the treatment than the second individual.
- Under this restriction, treatment cannot change the rank ordering of outcomes across individuals.
- These restrictions are in contrast with the Roy model and generalized Roy model, where one individual may have a higher with-treatment potential outcome but a lower without-treatment potential outcome compared to a second individual.



- In contrast to these approaches, control variate approaches impose more structure on the selection equation, imposing that the unobservables in the treatment selection equation ( $V$ ) be a scalar random variable, and that the treatment is an additive function of the unobservables or more generally a strictly increasing function of the unobservables.

- In contrast to these approaches, control variate approaches impose more structure on the selection equation, imposing that the unobservables in the treatment selection equation ( $V$ ) be a scalar random variable, and that the treatment is an additive function of the unobservables or more generally a strictly increasing function of the unobservables.
- Such approaches thus impose strong restrictions on the heterogeneity in the treatment selection equation.

- In contrast to these approaches, control variate approaches impose more structure on the selection equation, imposing that the unobservables in the treatment selection equation ( $V$ ) be a scalar random variable, and that the treatment is an additive function of the unobservables or more generally a strictly increasing function of the unobservables.
- Such approaches thus impose strong restrictions on the heterogeneity in the treatment selection equation.
- In exchange for these restrictions, such approaches do not require  $Y$  to be increasing in  $U$  and do not require  $U$  to be a scalar random variable.

- ? consider identification and estimation of the average structural function in a nonparametric model using the control variate approach, building on the work of ? and ?.

- ? consider identification and estimation of the average structural function in a nonparametric model using the control variate approach, building on the work of ? and ?.
- Their approach does not impose any further restrictions on the outcome equation, but does require a large support assumption.

- ? consider identification and estimation of the average structural function in a nonparametric model using the control variate approach, building on the work of ? and ?.
- Their approach does not impose any further restrictions on the outcome equation, but does require a large support assumption.
- Another recent contribution to the control function literature is ?, who restrict  $Y$  to be determined by a stochastic polynomial in  $D$  but do not require a large support assumption.

- ? consider identification and estimation of the average structural function in a nonparametric model using the control variate approach, building on the work of ? and ?.
- Their approach does not impose any further restrictions on the outcome equation, but does require a large support assumption.
- Another recent contribution to the control function literature is ?, who restrict  $Y$  to be determined by a stochastic polynomial in  $D$  but do not require a large support assumption.
- We now further discuss both approaches.

- ? consider identification and estimation of the average structural function in a nonparametric model using the control variate approach, building on the work of ? and ?.
- Their approach does not impose any further restrictions on the outcome equation, but does require a large support assumption.
- Another recent contribution to the control function literature is ?, who restrict  $Y$  to be determined by a stochastic polynomial in  $D$  but do not require a large support assumption.
- We now further discuss both approaches.
- ? proceed as follows.



- ? consider identification and estimation of the average structural function in a nonparametric model using the control variate approach, building on the work of ? and ?.
- Their approach does not impose any further restrictions on the outcome equation, but does require a large support assumption.
- Another recent contribution to the control function literature is ?, who restrict  $Y$  to be determined by a stochastic polynomial in  $D$  but do not require a large support assumption.
- We now further discuss both approaches.
- ? proceed as follows.
- They assume that  $\vartheta(z, v)$  is strictly monotonic in  $v$ .

- Suppose that  $(U, V) \perp\!\!\!\perp (X, Z)$ , and without loss of generality normalize  $V$  to be unit uniform.

- Suppose that  $(U, V) \perp\!\!\!\perp (X, Z)$ , and without loss of generality normalize  $V$  to be unit uniform.
- Then  $V$  is immediately identified (up to the normalization) from  $V = F(Y|X, Z)$ .

- Suppose that  $(U, V) \perp\!\!\!\perp (X, Z)$ , and without loss of generality normalize  $V$  to be unit uniform.
- Then  $V$  is immediately identified (up to the normalization) from  $V = F(Y|X, Z)$ .
- Given identification of  $V$ , they can identify  $E(Y|D, X, V)$ .

- Suppose that  $(U, V) \perp\!\!\!\perp (X, Z)$ , and without loss of generality normalize  $V$  to be unit uniform.
- Then  $V$  is immediately identified (up to the normalization) from  $V = F(Y|X, Z)$ .
- Given identification of  $V$ , they can identify  $E(Y|D, X, V)$ .
- Their independence assumptions imply that  $U \perp\!\!\!\perp D \mid (X, V)$ , so that

$$E(Y|D = d, X = x, V = v) = E(Y_d|X = x, V = v).$$

$E(Y_d|X = x, V = v)$  corresponds to the marginal treatment effect except that it is the conditional expectation in level instead of the derivative of the conditional expectation.

- Then, in parallel to the way ? integrate up the MTE to recover the ATE, Imbens and Newey integrate up  $E(Y_d|X = x, V = v)$  to obtain the ASF:

$$E(Y_d|X = x) = \int E(Y_d|X = x, V = v)dF_V(v) = \int E(Y|D = d, X = x, V = v)dF_V(v)$$

Imbens and Newey do not explicitly consider the ATE, TT, or MTE, but we can adapt the ? weighting analysis summarized in Slide 90 to obtain these parameters as a slight modification of the Imbens and Newey analysis.

- Then, in parallel to the way ? integrate up the MTE to recover the ATE, Imbens and Newey integrate up  $E(Y_d|X = x, V = v)$  to obtain the ASF:

$$E(Y_d|X = x) = \int E(Y_d|X = x, V = v)dF_V(v) = \int E(Y|D = d, X = x, V = v)dF_V(v)$$

Imbens and Newey do not explicitly consider the ATE, TT, or MTE, but we can adapt the ? weighting analysis summarized in Slide 90 to obtain these parameters as a slight modification of the Imbens and Newey analysis.

- First consider the MTE.

- Then, in parallel to the way ? integrate up the MTE to recover the ATE, Imbens and Newey integrate up  $E(Y_d|X = x, V = v)$  to obtain the ASF:

$$E(Y_d|X = x) = \int E(Y_d|X = x, V = v)dF_V(v) = \int E(Y|D = d, X = x, V = v)dF_V(v)$$

Imbens and Newey do not explicitly consider the ATE, TT, or MTE, but we can adapt the ? weighting analysis summarized in Slide 90 to obtain these parameters as a slight modification of the Imbens and Newey analysis.

- First consider the MTE.
- We have that

$$\frac{\partial}{\partial d}E(Y|D = d, X = x, V = v) = E\left(\frac{\partial}{\partial d}Y_d|X = x, V = v\right),$$

so that the MTE is identified.



- Integrating up the MTE we obtain ATE,

$$\begin{aligned} E\left(\frac{\partial}{\partial d} Y_d \mid X = x\right) &= \int E\left(\frac{\partial}{\partial d} Y_d \mid X = x, V = v\right) dF_V(v) \\ &= \int \frac{\partial}{\partial d} E(Y \mid D = d, X = x, V = v) dF_V(v) \end{aligned}$$

and TT

$$\begin{aligned} E\left(\frac{\partial}{\partial d_1} Y_{d_1} \mid D = d_2, X = x\right) \Big|_{d=d_1=d_2} \\ &= \int E\left(\frac{\partial}{\partial d} Y_d \mid X = x, V = v\right) dF_{V \mid D=d_2, X}(v \mid x) \\ &= \int \frac{\partial}{\partial d} E(Y \mid D = d, X = x, V = v) dF_{V \mid D=d_2, X}(v \mid x). \end{aligned}$$

- Note the strong connection between the control variate approach and the LIV/MTE approach of ?.

- Note the strong connection between the control variate approach and the LIV/MTE approach of ?.
- They both proceed by identifying an expectation conditional on the first stage error term, and then integrating that expectation up to obtain the parameter of interest.

- Note the strong connection between the control variate approach and the LIV/MTE approach of ?.
- They both proceed by identifying an expectation conditional on the first stage error term, and then integrating that expectation up to obtain the parameter of interest.
- The primary distinction is that, in the control variate approach with a continuous endogenous treatment, it is possible to assume that the treatment is a strictly increasing function of an error term that is independent of the instruments, to identify this error term, and then to explicitly include the identified first-stage error term as a regressor in the second stage regression for the outcome.

- In contrast, with a discrete endogenous treatment, it is not possible to characterize the treatment as a strictly increasing function of an error term that is independent of the instruments.

- In contrast, with a discrete endogenous treatment, it is not possible to characterize the treatment as a strictly increasing function of an error term that is independent of the instruments.
- It is thus not possible to identify the first-stage error term, and thus not possible to explicitly include an identified first-stage error term in the second stage.

- In contrast, with a discrete endogenous treatment, it is not possible to characterize the treatment as a strictly increasing function of an error term that is independent of the instruments.
- It is thus not possible to identify the first-stage error term, and thus not possible to explicitly include an identified first-stage error term in the second stage.
- The LIV strategy is the approach in the discrete case that by-passes the need to explicitly identify the first stage error term.

- In order to be able to integrate  $E(Y|D = d, X = x, V = v) = E(Y_d|X = x, V = v)$  up to obtain the ASF (or to integrate MTE to obtain ATE), it is necessary to evaluate  $E(Y|D = d, X = x, V = v)$  at all values of  $v$  in the support of the distribution of  $V$  conditional on  $X$ .



- In order to be able to integrate  $E(Y|D = d, X = x, V = v) = E(Y_d|X = x, V = v)$  up to obtain the ASF (or to integrate MTE to obtain ATE), it is necessary to evaluate  $E(Y|D = d, X = x, V = v)$  at all values of  $v$  in the support of the distribution of  $V$  conditional on  $X$ .
- This is a nontrivial requirement.

- In order to be able to integrate  $E(Y|D = d, X = x, V = v) = E(Y_d|X = x, V = v)$  up to obtain the ASF (or to integrate MTE to obtain ATE), it is necessary to evaluate  $E(Y|D = d, X = x, V = v)$  at all values of  $v$  in the support of the distribution of  $V$  conditional on  $X$ .
- This is a nontrivial requirement.
- To show this, suppress  $X$  for ease of exposition.

- In order to be able to integrate  $E(Y|D = d, X = x, V = v) = E(Y_d|X = x, V = v)$  up to obtain the ASF (or to integrate MTE to obtain ATE), it is necessary to evaluate  $E(Y|D = d, X = x, V = v)$  at all values of  $v$  in the support of the distribution of  $V$  conditional on  $X$ .
- This is a nontrivial requirement.
- To show this, suppress  $X$  for ease of exposition.
- One can only evaluate  $E(Y|D = d, V = v)$  at values of  $v$  in the support of the distribution of  $V$  conditional on  $D = d$ , so that the requirement is that the support of the distribution of  $V$  conditional on  $D = d$  equal the support of the unconditional distribution.

- This requires, in turn, a large support assumption on an element of  $Z$ .

- This requires, in turn, a large support assumption on an element of  $Z$ .
- For example, suppose that  $\vartheta(Z, V) = P(Z) + V$ , so that  $D = P(Z) + V$ .

- This requires, in turn, a large support assumption on an element of  $Z$ .
- For example, suppose that  $\vartheta(Z, V) = P(Z) + V$ , so that  $D = P(Z) + V$ .
- Let  $\mathcal{P}$  denote the support of the distribution of  $P(Z)$ .

- This requires, in turn, a large support assumption on an element of  $Z$ .
- For example, suppose that  $\vartheta(Z, V) = P(Z) + V$ , so that  $D = P(Z) + V$ .
- Let  $\mathcal{P}$  denote the support of the distribution of  $P(Z)$ .
- Then

$$\begin{aligned}\text{Supp}(V|D = d) &= \text{Supp}(V|P(Z) + V = d) \\ &= \text{Supp}(V|V = d - P(Z)) = \{d - p : p \in \mathcal{P}\}\end{aligned}$$

where the last equality uses  $Z \perp\!\!\!\perp V$ .

- For example, if  $\mathcal{P} = [a, b]$ , then  $\{d - p : p \in [a, b]\} = [d - b, d - a]$  which does not depend on  $d$  if and only if  $a = -\infty$  and  $b = \infty$ , i.e., if and only if  $\mathcal{P} = \mathbb{R}$ .



- For example, if  $\mathcal{P} = [a, b]$ , then  $\{d - p : p \in [a, b]\} = [d - b, d - a]$  which does not depend on  $d$  if and only if  $a = -\infty$  and  $b = \infty$ , i.e., if and only if  $\mathcal{P} = \mathbb{R}$ .
- For standard models, this requirement in turn necessitates a regressor with unbounded support, analogous to the identification at infinity requirement in selection models shown by ?.

- For example, if  $\mathcal{P} = [a, b]$ , then  $\{d - p : p \in [a, b]\} = [d - b, d - a]$  which does not depend on  $d$  if and only if  $a = -\infty$  and  $b = \infty$ , i.e., if and only if  $\mathcal{P} = \mathbb{R}$ .
- For standard models, this requirement in turn necessitates a regressor with unbounded support, analogous to the identification at infinity requirement in selection models shown by ?.
- We have noted the central role played by identification at infinity assumptions in many different settings throughout this Handbook.

- Next consider the analysis of ?.

- Next consider the analysis of ?.
- They assume that  $(U, V) \perp\!\!\!\perp (X, Z)$ .

- Next consider the analysis of ?.
- They assume that  $(U, V) \perp\!\!\!\perp (X, Z)$ .
- They impose additional structure on the outcome equation, in particular that the outcome equation can be expressed by a finite order stochastic polynomial in the treatment variable:

$$Y = \mu(D, X) + \sum_{j=0}^K D^j U_j$$

so that

$$Y_d = \mu_d(X) + \sum_{j=0}^K d^j U_j.$$

- This specification can be seen as a nonparametric extension of the random coefficient models of ? and ??.

- This specification can be seen as a nonparametric extension of the random coefficient models of ? and ??.
- As a consequence of the structure on the outcome equation, ? are able to identify the ATE without requiring the large support assumption of ?.

- This specification can be seen as a nonparametric extension of the random coefficient models of ? and ??.
- As a consequence of the structure on the outcome equation, ? are able to identify the ATE without requiring the large support assumption of ?.
- Instead of a large support assumption, they require measurable separability of  $D$  and  $V$  conditional on  $X$ .



- Measurable separability is the requirement that any function of  $D$  and  $X$  that almost surely equals a function of  $V$  and  $X$  must be a function of  $X$  only.

- Measurable separability is the requirement that any function of  $D$  and  $X$  that almost surely equals a function of  $V$  and  $X$  must be a function of  $X$  only.
- This assumption can be shown to be equivalent to requiring that  $D$  not lie in a subset of its support if and only if  $V$  lies in a subset of its support (conditional on  $X$ ).

- Measurable separability is the requirement that any function of  $D$  and  $X$  that almost surely equals a function of  $V$  and  $X$  must be a function of  $X$  only.
- This assumption can be shown to be equivalent to requiring that  $D$  not lie in a subset of its support if and only if  $V$  lies in a subset of its support (conditional on  $X$ ).
- As shown by ?, measurable separability between  $D$  and  $V$  follows from the independence assumption  $(U, V) \perp\!\!\!\perp (X, Z)$  along with mild regularity conditions.

- Thus the ? approach allows for identification of the average treatment effect with continuous endogenous regressors without requiring large support assumptions in exchange for requiring a finite-order, stochastic polynomial assumption on the outcome equation.

- Thus the ? approach allows for identification of the average treatment effect with continuous endogenous regressors without requiring large support assumptions in exchange for requiring a finite-order, stochastic polynomial assumption on the outcome equation.
- We next consider the method of matching, which is based on the assumption of conditional independence that is assumed to characterize data structures.

## Matching

- The method of matching assumes selection of treatment based on potential outcomes

$$(Y_0, Y_1) \not\perp D,$$

so  $\Pr(D = 1 \mid Y_0, Y_1)$  depends on  $Y_0, Y_1$ .

## Matching

- The method of matching assumes selection of treatment based on potential outcomes

$$(Y_0, Y_1) \not\perp D,$$

so  $\Pr(D = 1 \mid Y_0, Y_1)$  depends on  $Y_0, Y_1$ .

- It assumes access to variables  $Q$  such that conditioning on  $Q$  removes the dependence:

$$(Y_0, Y_1) \perp\!\!\!\perp D \mid Q. \tag{Q-1}$$

## Matching

- The method of matching assumes selection of treatment based on potential outcomes

$$(Y_0, Y_1) \not\perp D,$$

so  $\Pr(D = 1 \mid Y_0, Y_1)$  depends on  $Y_0, Y_1$ .

- It assumes access to variables  $Q$  such that conditioning on  $Q$  removes the dependence:

$$(Y_0, Y_1) \perp\!\!\!\perp D \mid Q. \quad (Q-1)$$

- Thus,

$$\Pr(D = 1 \mid Q, Y_0, Y_1) = \Pr(D = 1 \mid Q).$$



- Comparisons between treated and untreated can be made at all points in the support of  $Q$  such that

$$0 < \Pr(D = 1 \mid Q) < 1. \quad (\text{Q-2})$$

- Comparisons between treated and untreated can be made at all points in the support of  $Q$  such that

$$0 < \Pr(D = 1 \mid Q) < 1. \quad (\text{Q-2})$$

- The method does not explicitly model choices of treatment or the subjective evaluations of participants, nor is there any distinction between the variables in the outcome equations ( $X$ ) and the variables in the choice equations ( $Z$ ) that is central to the IV method and the method of control functions.

- Comparisons between treated and untreated can be made at all points in the support of  $Q$  such that

$$0 < \Pr(D = 1 \mid Q) < 1. \quad (\text{Q-2})$$

- The method does not explicitly model choices of treatment or the subjective evaluations of participants, nor is there any distinction between the variables in the outcome equations ( $X$ ) and the variables in the choice equations ( $Z$ ) that is central to the IV method and the method of control functions.
- In principle, condition (Q-1) can be satisfied using a set of variables  $Q$  distinct from all or some of the components of  $X$  and  $Z$ .

- Comparisons between treated and untreated can be made at all points in the support of  $Q$  such that

$$0 < \Pr(D = 1 \mid Q) < 1. \quad (\text{Q-2})$$

- The method does not explicitly model choices of treatment or the subjective evaluations of participants, nor is there any distinction between the variables in the outcome equations ( $X$ ) and the variables in the choice equations ( $Z$ ) that is central to the IV method and the method of control functions.
- In principle, condition (Q-1) can be satisfied using a set of variables  $Q$  distinct from all or some of the components of  $X$  and  $Z$ .
- The conditioning variables do not have to be exogenous.

- From condition (Q-1), we recover the distributions of  $Y_0$  and  $Y_1$  given  $Q$ ,  $\Pr(Y_0 \leq y_0 \mid Q = q) = F_0(y_0 \mid Q = q)$  and  $\Pr(Y_1 \leq y_1 \mid Q = q) = F_1(y_1 \mid Q = q)$ , but not the joint distribution  $F(y_0, y_1 \mid Q = q)$ , because we do not observe the same persons in the treated and untreated states.

- From condition (Q-1), we recover the distributions of  $Y_0$  and  $Y_1$  given  $Q$ ,  $\Pr(Y_0 \leq y_0 \mid Q = q) = F_0(y_0 \mid Q = q)$  and  $\Pr(Y_1 \leq y_1 \mid Q = q) = F_1(y_1 \mid Q = q)$ , but not the joint distribution  $F(y_0, y_1 \mid Q = q)$ , because we do not observe the same persons in the treated and untreated states.
- This is a standard evaluation problem common to all econometric estimators.

- From condition (Q-1), we recover the distributions of  $Y_0$  and  $Y_1$  given  $Q$ ,  $\Pr(Y_0 \leq y_0 \mid Q = q) = F_0(y_0 \mid Q = q)$  and  $\Pr(Y_1 \leq y_1 \mid Q = q) = F_1(y_1 \mid Q = q)$ , but not the joint distribution  $F(y_0, y_1 \mid Q = q)$ , because we do not observe the same persons in the treated and untreated states.
- This is a standard evaluation problem common to all econometric estimators.
- Methods for determining which variables belong in  $Q$  rely on untested exogeneity assumptions which we discuss in this section.

- OLS is a special case of matching that focuses on the identification of certain conditional means.



- OLS is a special case of matching that focuses on the identification of certain conditional means.
- In OLS, linear functional forms are maintained as exact representations or valid approximations.

- OLS is a special case of matching that focuses on the identification of certain conditional means.
- In OLS, linear functional forms are maintained as exact representations or valid approximations.
- Considering a common coefficient model, OLS writes

$$Y = Q\alpha + D\beta + U, \quad (\text{Q-3})$$

where  $\alpha$  is the treatment effect and

$$E(U | Q, D) = 0. \quad (\text{Q-4})$$

The assumption is made that the variance-covariance matrix of  $(Q, D)$  is of full rank:

$$\text{Var}(Q, D) \text{ full rank.} \quad (\text{Q-5})$$

Under these conditions, we can identify  $\beta$  even though  $D$  and  $U$  are dependent:  $D \not\perp U$ .

- Controlling for the observable  $Q$  eliminates any spurious mean dependence between  $D$  and  $U$ :  $E(U | D) \neq 0$  but  $E(U | D, Q) = 0$ .

- Controlling for the observable  $Q$  eliminates any spurious mean dependence between  $D$  and  $U$ :  $E(U | D) \neq 0$  but  $E(U | D, Q) = 0$ .
- (Q-4) is the linear regression counterpart to (Q-1).

- Controlling for the observable  $Q$  eliminates any spurious mean dependence between  $D$  and  $U$ :  $E(U | D) \neq 0$  but  $E(U | D, Q) = 0$ .
- (Q-4) is the linear regression counterpart to (Q-1).
- (Q-4) is the linear regression counterpart to (Q-2).

- Controlling for the observable  $Q$  eliminates any spurious mean dependence between  $D$  and  $U$ :  $E(U | D) \neq 0$  but  $E(U | D, Q) = 0$ .
- (Q-4) is the linear regression counterpart to (Q-1).
- (Q-4) is the linear regression counterpart to (Q-2).
- Failure of (Q-5) would mean that using a nonparametric estimator, we might perfectly predict  $D$  given  $Q$ , and that  $\Pr(D = 1 | Q = q) = 1$  or  $0$ .



- (Q-5)' : If the goal of the analysis is to identify  $\beta$ , in place of (Q-4), we can get by with

$$(Q - 4)' : E(U|Q, D) = E(U|Q).$$



- (Q-5)' : If the goal of the analysis is to identify  $\beta$ , in place of (Q-4), we can get by with

$$(Q - 4)' : E(U|Q, D) = E(U|Q).$$

Assuming  $\text{Var}(D | Q) > 0$ , we can identify  $\beta$  even if we cannot separate  $\alpha Q$  from  $E(U|Q)$ .

- (Q-5)' : If the goal of the analysis is to identify  $\beta$ , in place of (Q-4), we can get by with

$$(Q - 4)' : E(U|Q, D) = E(U|Q).$$

Assuming  $\text{Var}(D | Q) > 0$ , we can identify  $\beta$  even if we cannot separate  $\alpha Q$  from  $E(U|Q)$ .

- Matching can be implemented as a nonparametric method.

- Matching can be implemented as a nonparametric method.
- When this is done, the procedure does not require specification of the functional form of the outcome equations.

- Matching can be implemented as a nonparametric method.
- When this is done, the procedure does not require specification of the functional form of the outcome equations.
- It enforces the requirement that (Q-2) be satisfied by estimating functions pointwise in the support of  $Q$ .

- Matching can be implemented as a nonparametric method.
- When this is done, the procedure does not require specification of the functional form of the outcome equations.
- It enforces the requirement that (Q-2) be satisfied by estimating functions pointwise in the support of  $Q$ .
- To link our notation in this section to that in the rest of the chapter, we assume that  $Q = (X, Z)$  and that  $X$  and  $Z$  are the same except where otherwise noted.

- Matching can be implemented as a nonparametric method.
- When this is done, the procedure does not require specification of the functional form of the outcome equations.
- It enforces the requirement that (Q-2) be satisfied by estimating functions pointwise in the support of  $Q$ .
- To link our notation in this section to that in the rest of the chapter, we assume that  $Q = (X, Z)$  and that  $X$  and  $Z$  are the same except where otherwise noted.
- Thus we invoke assumptions (M-1) and (M-2) presented in Slide 12, even though in principle we can use a more general conditioning set.

- Assumptions (M-1) and (M-2) introduced in Section 2 or (Q-1) and (Q-2) rule out the possibility that after conditioning on  $X$  (or  $Q$ ), agents possess more information about their choices than econometricians, and that the unobserved information helps to predict the potential outcomes.



- Assumptions (M-1) and (M-2) introduced in Section 2 or (Q-1) and (Q-2) rule out the possibility that after conditioning on  $X$  (or  $Q$ ), agents possess more information about their choices than econometricians, and that the unobserved information helps to predict the potential outcomes.
- Put another way, the method allows for potential outcomes to affect choices but only through the observed variables,  $Q$ , that predict outcomes.

- Assumptions (M-1) and (M-2) introduced in Section 2 or (Q-1) and (Q-2) rule out the possibility that after conditioning on  $X$  (or  $Q$ ), agents possess more information about their choices than econometricians, and that the unobserved information helps to predict the potential outcomes.
- Put another way, the method allows for potential outcomes to affect choices but only through the observed variables,  $Q$ , that predict outcomes.
- This is the reason why ?? call the method selection on observables.

- This section establishes the following points.

- This section establishes the following points.
- (1) Matching assumptions (M-1) and (M-2) generically imply a flat MTE in  $u_D$ , i.e., they assume that  $E(Y_1 - Y_0 \mid X = x, U_D = u_D)$  does not depend on  $u_D$ .

- This section establishes the following points.
- (1) Matching assumptions (M-1) and (M-2) generically imply a flat MTE in  $u_D$ , i.e., they assume that  $E(Y_1 - Y_0 \mid X = x, U_D = u_D)$  does not depend on  $u_D$ .
- Thus the unobservables central to the Roy model and its extensions and the unobservables central to the modern IV literature are assumed to be absent once the analyst conditions on  $X$ .

- This section establishes the following points.
- (1) Matching assumptions (M-1) and (M-2) generically imply a flat MTE in  $u_D$ , i.e., they assume that  $E(Y_1 - Y_0 \mid X = x, U_D = u_D)$  does not depend on  $u_D$ .
- Thus the unobservables central to the Roy model and its extensions and the unobservables central to the modern IV literature are assumed to be absent once the analyst conditions on  $X$ .
- (M-1) implies that all mean treatment parameters are the same.

- This section establishes the following points.
- (1) Matching assumptions (M-1) and (M-2) generically imply a flat MTE in  $u_D$ , i.e., they assume that  $E(Y_1 - Y_0 \mid X = x, U_D = u_D)$  does not depend on  $u_D$ .
- Thus the unobservables central to the Roy model and its extensions and the unobservables central to the modern IV literature are assumed to be absent once the analyst conditions on  $X$ .
- (M-1) implies that all mean treatment parameters are the same.
- (2) Even if we weaken (M-1) and (M-2) to mean independence instead of full independence, generically the MTE is flat in  $u_D$  under the assumptions of the nonparametric generalized Roy model developed in Slide 90, so again all mean treatment parameters are the same.

- (3) We show that IV and matching make distinct identifying assumptions even though they both invoke conditional independence assumptions.



- (3) We show that IV and matching make distinct identifying assumptions even though they both invoke conditional independence assumptions.
- (4) We compare matching with IV and control function (sample selection) methods.

- (3) We show that IV and matching make distinct identifying assumptions even though they both invoke conditional independence assumptions.
- (4) We compare matching with IV and control function (sample selection) methods.
- Matching assumes that conditioning on observables eliminates the dependence between  $(Y_0, Y_1)$  and  $D$ .

- (3) We show that IV and matching make distinct identifying assumptions even though they both invoke conditional independence assumptions.
- (4) We compare matching with IV and control function (sample selection) methods.
- Matching assumes that conditioning on observables eliminates the dependence between  $(Y_0, Y_1)$  and  $D$ .
- The control function principle models the dependence.

- (3) We show that IV and matching make distinct identifying assumptions even though they both invoke conditional independence assumptions.
- (4) We compare matching with IV and control function (sample selection) methods.
- Matching assumes that conditioning on observables eliminates the dependence between  $(Y_0, Y_1)$  and  $D$ .
- The control function principle models the dependence.
- (5) We present some examples that demonstrate that if the assumptions of the method of matching are violated, the method can produce substantially biased estimators of the parameters of interest.

- (6) We show that standard methods for selecting the conditioning variables used in matching assume exogeneity.

- (6) We show that standard methods for selecting the conditioning variables used in matching assume exogeneity.
- This is a property shared with many econometric estimators, as noted in Part I, section 5.2.

- (6) We show that standard methods for selecting the conditioning variables used in matching assume exogeneity.
- This is a property shared with many econometric estimators, as noted in Part I, section 5.2.
- Violations of the exogeneity assumption can produce biased estimators.

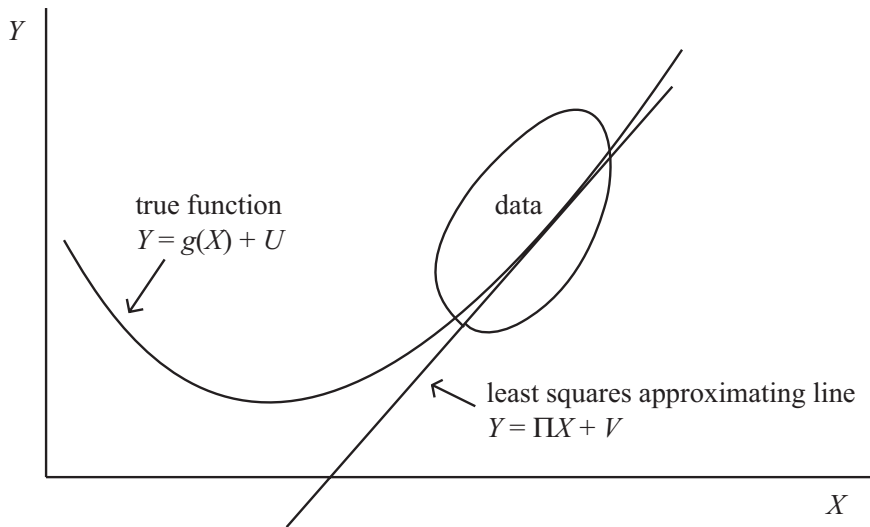
- (6) We show that standard methods for selecting the conditioning variables used in matching assume exogeneity.
- This is a property shared with many econometric estimators, as noted in Part I, section 5.2.
- Violations of the exogeneity assumption can produce biased estimators.
- Nonparametric versions of matching embodying (M-2) avoid the problem of making inferences outside the support of the data.



- (6) We show that standard methods for selecting the conditioning variables used in matching assume exogeneity.
- This is a property shared with many econometric estimators, as noted in Part I, section 5.2.
- Violations of the exogeneity assumption can produce biased estimators.
- Nonparametric versions of matching embodying (M-2) avoid the problem of making inferences outside the support of the data.
- This problem is implicit in any application of least squares.

- (6) We show that standard methods for selecting the conditioning variables used in matching assume exogeneity.
- This is a property shared with many econometric estimators, as noted in Part I, section 5.2.
- Violations of the exogeneity assumption can produce biased estimators.
- Nonparametric versions of matching embodying (M-2) avoid the problem of making inferences outside the support of the data.
- This problem is implicit in any application of least squares.
- Figure 22 shows the support problem that can arise in linear least squares when the linearity of the regression is used to extrapolate estimates determined in one empirical support to new supports.

Figure 22: The Least Squares Extrapolation Problem Avoided by Using Nonparametric Regression or Matching



- Careful attention to support problems is a virtue of any nonparametric method, including, but not unique to, nonparametric matching.

- Careful attention to support problems is a virtue of any nonparametric method, including, but not unique to, nonparametric matching.
- ? show that the bias from neglecting the problem of limited support can be substantial.

- Careful attention to support problems is a virtue of any nonparametric method, including, but not unique to, nonparametric matching.
- ? show that the bias from neglecting the problem of limited support can be substantial.
- See also the discussion in ?.

- We now show that matching implies that conditional on  $X$ , the marginal return is assumed to be the same as the average return (marginal = average).

- We now show that matching implies that conditional on  $X$ , the marginal return is assumed to be the same as the average return (marginal = average).
- This is a strong behavioral assumption implicit in statistical conditional independence assumption (M-1).



- We now show that matching implies that conditional on  $X$ , the marginal return is assumed to be the same as the average return (marginal = average).
- This is a strong behavioral assumption implicit in statistical conditional independence assumption (M-1).
- It says that the marginal participant has the same return as the average participant.

## Matching Assumption (M-1) Implies a Flat MTE

- An immediate consequence of (M-1) is that the MTE does not depend on  $U_D$ .

## Matching Assumption (M-1) Implies a Flat MTE

- An immediate consequence of (M-1) is that the MTE does not depend on  $U_D$ .
- This is so because  $(Y_0, Y_1) \perp\!\!\!\perp D \mid X$  implies that  $(Y_0, Y_1) \perp\!\!\!\perp U_D \mid X$  and hence that

$$\Delta^{\text{MTE}}(x, u_D) = E(Y_1 - Y_0 \mid X = x, U_D = u_D) = E(Y_1 - Y_0 \mid X = x). \quad (60)$$

- This, in turn, implies that  $\Delta^{\text{MTE}}$  conditional on  $X$  is flat in  $u_D$ , so that matching invokes assumption (C-1) invoked in Slide 197.
- Under our assumptions for the generalized Roy model, it assumes that  $E(Y | P(Z) = p)$  is linear in  $p$ . Thus the method of matching assumes that mean marginal returns and average returns are the same and all mean treatment effects are the same given  $X$ .

- However, one can still distinguish marginal from average effects of the observables ( $X$ ) using matching.

- However, one can still distinguish marginal from average effects of the observables ( $X$ ) using matching.
- See ?.

- However, one can still distinguish marginal from average effects of the observables ( $X$ ) using matching.
- See ?.
- It is sometimes said that the matching assumptions are “for free” (See, e.g., ?) because one can always replace unobserved  $F_1(Y_1 | X = x, D = 0)$  with observed  $F_1(Y_1 | X = x, D = 1)$  and unobserved  $F_0(Y_0 | X = x, D = 1)$  with observed  $F_0(Y_0 | X = x, D = 0)$ .

- However, one can still distinguish marginal from average effects of the observables ( $X$ ) using matching.
- See ?.
- It is sometimes said that the matching assumptions are “for free” (See, e.g., ?) because one can always replace unobserved  $F_1(Y_1 | X = x, D = 0)$  with observed  $F_1(Y_1 | X = x, D = 1)$  and unobserved  $F_0(Y_0 | X = x, D = 1)$  with observed  $F_0(Y_0 | X = x, D = 0)$ .
- Such substitutions do not contradict any observed data.



- While the claim is true, it ignores the counterfactual states generated under the matching assumptions.

- While the claim is true, it ignores the counterfactual states generated under the matching assumptions.
- The assumed absence of selection on unobservables is not a “for free” assumption, and produces fundamentally different counterfactual states for the same model under matching and selection assumptions.

- While the claim is true, it ignores the counterfactual states generated under the matching assumptions.
- The assumed absence of selection on unobservables is not a “for free” assumption, and produces fundamentally different counterfactual states for the same model under matching and selection assumptions.
- To explore these issues in depth, consider a nonparametric regression model more general than the linear regression model (Q-3).

- Without assumption (M-1), a nonparametric regression of  $Y$  on  $D$  conditional on  $X$  identifies a nonparametric mean difference:

$$\begin{aligned}\Delta^{\text{OLS}}(X) &= E(Y_1 | X, D = 1) - E(Y_0 | X, D = 0) \\ &= E(Y_1 - Y_0 | X, D = 1) + \{E(Y_0 | X, D = 1) - E(Y_0 | X, D = 0)\} \quad (61)\end{aligned}$$

The term in braces in the second expression arises from selection on pre-treatment levels of the outcome.

- Without assumption (M-1), a nonparametric regression of  $Y$  on  $D$  conditional on  $X$  identifies a nonparametric mean difference:

$$\begin{aligned}\Delta^{\text{OLS}}(X) &= E(Y_1 | X, D = 1) - E(Y_0 | X, D = 0) \\ &= E(Y_1 - Y_0 | X, D = 1) + \{E(Y_0 | X, D = 1) - E(Y_0 | X, D = 0)\} \quad (61)\end{aligned}$$

The term in braces in the second expression arises from selection on pre-treatment levels of the outcome.

- OLS identifies the parameter treatment on the treated (the first term in the second line of (61)) plus a bias term in braces corresponding to selection on the levels.

- The OLS estimator can be represented as a weighted average of  $\Delta^{\text{MTE}}$ .

- The OLS estimator can be represented as a weighted average of  $\Delta^{\text{MTE}}$ .
- The weight is given in table 2B where  $U_1$  and  $U_0$  for the OLS model are defined as deviations from conditional expectations,  $U_1 = Y_1 - E(Y_1 | X)$ ,  $U_0 = Y_0 - E(Y_0 | X)$ .

- The OLS estimator can be represented as a weighted average of  $\Delta^{\text{MTE}}$ .
- The weight is given in table 2B where  $U_1$  and  $U_0$  for the OLS model are defined as deviations from conditional expectations,  $U_1 = Y_1 - E(Y_1 | X)$ ,  $U_0 = Y_0 - E(Y_0 | X)$ .
- Unlike the weights for  $\Delta^{\text{TT}}$  and  $\Delta^{\text{ATE}}$ , the OLS weights do not necessarily integrate to one and they are not necessarily nonnegative.



- The OLS estimator can be represented as a weighted average of  $\Delta^{\text{MTE}}$ .
- The weight is given in table 2B where  $U_1$  and  $U_0$  for the OLS model are defined as deviations from conditional expectations,  $U_1 = Y_1 - E(Y_1 | X)$ ,  $U_0 = Y_0 - E(Y_0 | X)$ .
- Unlike the weights for  $\Delta^{\text{TT}}$  and  $\Delta^{\text{ATE}}$ , the OLS weights do not necessarily integrate to one and they are not necessarily nonnegative.
- Application of IV eliminates the contribution of the second term of equation (61).

- The OLS estimator can be represented as a weighted average of  $\Delta^{\text{MTE}}$ .
- The weight is given in table 2B where  $U_1$  and  $U_0$  for the OLS model are defined as deviations from conditional expectations,  $U_1 = Y_1 - E(Y_1 | X)$ ,  $U_0 = Y_0 - E(Y_0 | X)$ .
- Unlike the weights for  $\Delta^{\text{TT}}$  and  $\Delta^{\text{ATE}}$ , the OLS weights do not necessarily integrate to one and they are not necessarily nonnegative.
- Application of IV eliminates the contribution of the second term of equation (61).
- The weights for the first term are the same as the weights for  $\Delta^{\text{TT}}$  and hence they integrate to one.

- The OLS weights for our generalized Roy model example are plotted in figure 2B.

- The OLS weights for our generalized Roy model example are plotted in figure 2B.
- The negative component of the OLS weight leads to a smaller OLS treatment estimate compared to the other treatment effects in table 3.

- The OLS weights for our generalized Roy model example are plotted in figure 2B.
- The negative component of the OLS weight leads to a smaller OLS treatment estimate compared to the other treatment effects in table 3.
- This table shows the estimated OLS treatment effect for the generalized Roy example.

- The OLS weights for our generalized Roy model example are plotted in figure 2B.
- The negative component of the OLS weight leads to a smaller OLS treatment estimate compared to the other treatment effects in table 3.
- This table shows the estimated OLS treatment effect for the generalized Roy example.
- The large negative selection bias in this example is consistent with comparative advantage as emphasized by ? and detected empirically by ? and ?.

- The OLS weights for our generalized Roy model example are plotted in figure 2B.
- The negative component of the OLS weight leads to a smaller OLS treatment estimate compared to the other treatment effects in table 3.
- This table shows the estimated OLS treatment effect for the generalized Roy example.
- The large negative selection bias in this example is consistent with comparative advantage as emphasized by ? and detected empirically by ? and ?.
- People who are good in sector 1 (i.e., receive treatment) may be very poor in sector 0 (those who receive no treatment).

- Hence the bias in OLS for the parameter treatment on the treated may be negative  
( $E(Y_0 | X, D = 1) - E(Y_0 | X, D = 0) < 0$ ).



- Hence the bias in OLS for the parameter treatment on the treated may be negative  
( $E(Y_0 | X, D = 1) - E(Y_0 | X, D = 0) < 0$ ).
- The differences among the policy relevant treatment effects, the conventional treatment effects and the OLS estimand are illustrated in figure 4A and table 3 for the generalized Roy model example.

- Hence the bias in OLS for the parameter treatment on the treated may be negative ( $E(Y_0 | X, D = 1) - E(Y_0 | X, D = 0) < 0$ ).
- The differences among the policy relevant treatment effects, the conventional treatment effects and the OLS estimand are illustrated in figure 4A and table 3 for the generalized Roy model example.
- As is evident from table 3, it is not at all clear that the instrumental variable estimator, with instruments that satisfy classical properties, performs better than nonparametric OLS in identifying the policy relevant treatment effect in this example.

- Hence the bias in OLS for the parameter treatment on the treated may be negative ( $E(Y_0 | X, D = 1) - E(Y_0 | X, D = 0) < 0$ ).
- The differences among the policy relevant treatment effects, the conventional treatment effects and the OLS estimand are illustrated in figure 4A and table 3 for the generalized Roy model example.
- As is evident from table 3, it is not at all clear that the instrumental variable estimator, with instruments that satisfy classical properties, performs better than nonparametric OLS in identifying the policy relevant treatment effect in this example.
- While IV eliminates the term in braces in (61), it reweights the MTE differently from what might be desired for many policy analyses.

- If there is no selection on unobserved variables conditional on covariates,  $U_D \perp\!\!\!\perp (Y_0, Y_1) \mid X$ , then  
 $E(U_1 \mid X, U_D) = E(U_1 \mid X) = 0$  and  
 $E(U_0 \mid X, U_D) = E(U_0 \mid X) = 0$  so that the OLS weights are unity and OLS identifies both ATE and the parameter treatment on the treated (TT), which are the same under this assumption.

- If there is no selection on unobserved variables conditional on covariates,  $U_D \perp\!\!\!\perp (Y_0, Y_1) \mid X$ , then  
 $E(U_1 \mid X, U_D) = E(U_1 \mid X) = 0$  and  
 $E(U_0 \mid X, U_D) = E(U_0 \mid X) = 0$  so that the OLS weights are unity and OLS identifies both ATE and the parameter treatment on the treated (TT), which are the same under this assumption.
- This condition is an implication of matching condition (M-1).

- If there is no selection on unobserved variables conditional on covariates,  $U_D \perp\!\!\!\perp (Y_0, Y_1) \mid X$ , then  $E(U_1 \mid X, U_D) = E(U_1 \mid X) = 0$  and  $E(U_0 \mid X, U_D) = E(U_0 \mid X) = 0$  so that the OLS weights are unity and OLS identifies both ATE and the parameter treatment on the treated (TT), which are the same under this assumption.
- This condition is an implication of matching condition (M-1).
- Given the assumed conditional independence in terms of  $X$ , we can identify ATE and TT without use of any instrument  $Z$  satisfying assumptions (A-1)–(A-2).

- If there is no selection on unobserved variables conditional on covariates,  $U_D \perp\!\!\!\perp (Y_0, Y_1) \mid X$ , then  $E(U_1 \mid X, U_D) = E(U_1 \mid X) = 0$  and  $E(U_0 \mid X, U_D) = E(U_0 \mid X) = 0$  so that the OLS weights are unity and OLS identifies both ATE and the parameter treatment on the treated (TT), which are the same under this assumption.
- This condition is an implication of matching condition (M-1).
- Given the assumed conditional independence in terms of  $X$ , we can identify ATE and TT without use of any instrument  $Z$  satisfying assumptions (A-1)–(A-2).
- If there is such a  $Z$ , the conditional independence condition implies under (A-1)–(A-5) that  $E(Y \mid X, P(Z) = p)$  is linear in  $p$ .

- The conditional independence assumption invoked in the method of matching has come into widespread use for much the same reason that OLS has come into widespread use.



- The conditional independence assumption invoked in the method of matching has come into widespread use for much the same reason that OLS has come into widespread use.
- It is easy to implement with modern software and makes little demands of the data because it assumes the existence of  $X$  variables that satisfy the conditional independence assumptions.

- The conditional independence assumption invoked in the method of matching has come into widespread use for much the same reason that OLS has come into widespread use.
- It is easy to implement with modern software and makes little demands of the data because it assumes the existence of  $X$  variables that satisfy the conditional independence assumptions.
- The crucial conditional independence assumption is not testable.

- The conditional independence assumption invoked in the method of matching has come into widespread use for much the same reason that OLS has come into widespread use.
- It is easy to implement with modern software and makes little demands of the data because it assumes the existence of  $X$  variables that satisfy the conditional independence assumptions.
- The crucial conditional independence assumption is not testable.
- As we note below, additional assumptions on the  $X$  are required to test the validity of the matching assumptions.

- If the sole interest is to identify treatment on the treated,  $\Delta^{TT}$ , it is apparent from representation (61) that we can weaken (M-1) to

$$(M-1)' \quad Y_0 \perp\!\!\!\perp D \mid X.$$

- If the sole interest is to identify treatment on the treated,  $\Delta^{TT}$ , it is apparent from representation (61) that we can weaken (M-1) to

$$(M-1)' \quad Y_0 \perp\!\!\!\perp D \mid X.$$

- This is possible because  $E(Y_1 \mid X, D = 1)$  is known from data on outcomes of the treated and only need to construct  $E(Y_0 \mid X, D = 1)$ .

- If the sole interest is to identify treatment on the treated,  $\Delta^{TT}$ , it is apparent from representation (61) that we can weaken (M-1) to

$$(M-1)' \quad Y_0 \perp\!\!\!\perp D \mid X.$$

- This is possible because  $E(Y_1 \mid X, D = 1)$  is known from data on outcomes of the treated and only need to construct  $E(Y_0 \mid X, D = 1)$ .
- In this case, MTE is not restricted to be flat in  $u_D$  and all treatment parameters are not the same.

- If the sole interest is to identify treatment on the treated,  $\Delta^{TT}$ , it is apparent from representation (61) that we can weaken (M-1) to

$$(M-1)' \quad Y_0 \perp\!\!\!\perp D \mid X.$$

- This is possible because  $E(Y_1 \mid X, D = 1)$  is known from data on outcomes of the treated and only need to construct  $E(Y_0 \mid X, D = 1)$ .
- In this case, MTE is not restricted to be flat in  $u_D$  and all treatment parameters are not the same.
- A straightforward implication of (M-1)' in the Roy model, where selection is made solely on the gain, is that persons must sort into treatment status positively in terms of levels of  $Y_1$ .

- If the sole interest is to identify treatment on the treated,  $\Delta^{TT}$ , it is apparent from representation (61) that we can weaken (M-1) to

$$(M-1)' \quad Y_0 \perp\!\!\!\perp D \mid X.$$

- This is possible because  $E(Y_1 \mid X, D = 1)$  is known from data on outcomes of the treated and only need to construct  $E(Y_0 \mid X, D = 1)$ .
- In this case, MTE is not restricted to be flat in  $u_D$  and all treatment parameters are not the same.
- A straightforward implication of (M-1)' in the Roy model, where selection is made solely on the gain, is that persons must sort into treatment status positively in terms of levels of  $Y_1$ .
- We now consider more generally the implications of assuming mean independence of the errors rather than full independence.



## Matching and MTE Using Mean Independence Conditions

- To identify all mean treatment parameters, one can weaken the assumption (M-1) to the condition that  $Y_0$  and  $Y_1$  are mean independent of  $D$  conditional on  $X$ .

## Matching and MTE Using Mean Independence Conditions

- To identify all mean treatment parameters, one can weaken the assumption (M-1) to the condition that  $Y_0$  and  $Y_1$  are mean independent of  $D$  conditional on  $X$ .
- However,  $(Y_0, Y_1)$  will be mean independent of  $D$  conditional on  $X$  without  $U_D$  being independent of  $Y_0, Y_1$  conditional on  $X$  only if fortuitous balancing occurs, with regions of positive dependence of  $(Y_0, Y_1)$  on  $U_D$  and regions of negative dependence of  $(Y_0, Y_1)$  on  $U_D$  just exactly offsetting each other.

## Matching and MTE Using Mean Independence Conditions

- To identify all mean treatment parameters, one can weaken the assumption (M-1) to the condition that  $Y_0$  and  $Y_1$  are mean independent of  $D$  conditional on  $X$ .
- However,  $(Y_0, Y_1)$  will be mean independent of  $D$  conditional on  $X$  without  $U_D$  being independent of  $Y_0, Y_1$  conditional on  $X$  only if fortuitous balancing occurs, with regions of positive dependence of  $(Y_0, Y_1)$  on  $U_D$  and regions of negative dependence of  $(Y_0, Y_1)$  on  $U_D$  just exactly offsetting each other.
- Such balancing is not generic in the Roy model and in the generalized Roy model.

- In particular, assume that  $Y_j = \mu_j(X) + U_j$  for  $j = 0, 1$  and further assume that  $D = \mathbf{1}[Y_1 - Y_0 \geq C(Z) + U_C]$ .

- In particular, assume that  $Y_j = \mu_j(X) + U_j$  for  $j = 0, 1$  and further assume that  $D = \mathbf{1}[Y_1 - Y_0 \geq C(Z) + U_C]$ .
- Let  $V = U_C - (U_1 - U_0)$ .

- In particular, assume that  $Y_j = \mu_j(X) + U_j$  for  $j = 0, 1$  and further assume that  $D = \mathbf{1}[Y_1 - Y_0 \geq C(Z) + U_C]$ .
- Let  $V = U_C - (U_1 - U_0)$ .
- Assume  $(U_0, U_1, V) \perp\!\!\!\perp (X, Z)$ .

- In particular, assume that  $Y_j = \mu_j(X) + U_j$  for  $j = 0, 1$  and further assume that  $D = \mathbf{1}[Y_1 - Y_0 \geq C(Z) + U_C]$ .
- Let  $V = U_C - (U_1 - U_0)$ .
- Assume  $(U_0, U_1, V) \perp\!\!\!\perp (X, Z)$ .
- Then if  $V \perp\!\!\!\perp (U_1 - U_0)$ , and  $U_C$  has a log concave density, then  $E(Y_1 - Y_0 | X, V = v)$  is decreasing in  $v$ ,  $\Delta^{TT}(x) > \Delta^{ATE}(x)$ , and the matching conditions do not hold.

- In particular, assume that  $Y_j = \mu_j(X) + U_j$  for  $j = 0, 1$  and further assume that  $D = \mathbf{1}[Y_1 - Y_0 \geq C(Z) + U_C]$ .
- Let  $V = U_C - (U_1 - U_0)$ .
- Assume  $(U_0, U_1, V) \perp\!\!\!\perp (X, Z)$ .
- Then if  $V \perp\!\!\!\perp (U_1 - U_0)$ , and  $U_C$  has a log concave density, then  $E(Y_1 - Y_0 | X, V = v)$  is decreasing in  $v$ ,  $\Delta^{\text{TT}}(x) > \Delta^{\text{ATE}}(x)$ , and the matching conditions do not hold.
- If  $V \perp\!\!\!\perp (U_1 - U_0)$  but  $V$  does not have a log concave density, then it is still the case that  $(U_1 - U_0, V)$  is negative quadrant dependent.



- One can show that  $(U_1 - U_0, V)$  being negative quadrant dependent implies that  $\Delta^{TT}(x) > \Delta^{ATE}(x)$ , and thus again that the matching conditions cannot hold.

- One can show that  $(U_1 - U_0, V)$  being negative quadrant dependent implies that  $\Delta^{TT}(x) > \Delta^{ATE}(x)$ , and thus again that the matching conditions cannot hold.
- We now develop a more general analysis.

- One can show that  $(U_1 - U_0, V)$  being negative quadrant dependent implies that  $\Delta^{\text{TT}}(x) > \Delta^{\text{ATE}}(x)$ , and thus again that the matching conditions cannot hold.
- We now develop a more general analysis.
- Suppose that we assume selection model (7) so that  $D = \mathbf{1}[P(Z) \geq U_D]$ , where  $Z$  is independent of  $(Y_0, Y_1)$  conditional on  $X$ , where  $U_D = F_{V|X}(V)$  and  $P(Z) = F_{V|X}(\mu_D(Z))$ .

- One can show that  $(U_1 - U_0, V)$  being negative quadrant dependent implies that  $\Delta^{\text{TT}}(x) > \Delta^{\text{ATE}}(x)$ , and thus again that the matching conditions cannot hold.
- We now develop a more general analysis.
- Suppose that we assume selection model (7) so that  $D = \mathbf{1}[P(Z) \geq U_D]$ , where  $Z$  is independent of  $(Y_0, Y_1)$  conditional on  $X$ , where  $U_D = F_{V|X}(V)$  and  $P(Z) = F_{V|X}(\mu_D(Z))$ .
- Consider the weaker mean independence assumptions in place of assumption (M-1):

(M-4)

$$E(Y_1|X, D) = E(Y_1|X), \quad E(Y_0|X, D) = E(Y_0|X).$$

- This assumption is all that is needed to identify the mean treatment parameters because under it

$$E(Y|X = x, Z = z, D = 1) = E(Y_1|X = x, Z = z, D = 1) = E(Y_1|X = x)$$

and

$$E(Y|X = x, Z = z, D = 0) = E(Y_0|X = x, Z = z, D = 0) = E(Y_0|X = x).$$

- This assumption is all that is needed to identify the mean treatment parameters because under it

$$E(Y|X = x, Z = z, D = 1) = E(Y_1|X = x, Z = z, D = 1) = E(Y_1|X = x)$$

and

$$E(Y|X = x, Z = z, D = 0) = E(Y_0|X = x, Z = z, D = 0) = E(Y_0|X = x).$$

- Thus we can identify all the mean treatment parameters over the support that satisfies (M-2).

- Recalling that  $\Delta = Y_1 - Y_0$ , (M-3) implies in terms of  $U_D$  that

$$\begin{aligned} E(\Delta|X = x, Z = z, U_D \leq P(z)) &= E(\Delta|X = x) \\ \Leftrightarrow E(\Delta^{\text{MTE}}(X, U_D)|X = x, U_D \leq P(z)) &= E(\Delta|X = x), \end{aligned}$$

and hence

$$E(\Delta^{\text{MTE}}(X, U_D)|X = x, U_D \leq P(z)) = E(\Delta^{\text{MTE}}(X, U_D)|X = x, U_D > P(z)).$$

If the support of  $P(Z)$  is the full unit interval conditional on  $X = x$ , then  $\Delta^{\text{MTE}}(X, U_D) = E(\Delta|X = x)$  for all  $U_D$ .



- Recalling that  $\Delta = Y_1 - Y_0$ , (M-3) implies in terms of  $U_D$  that

$$E(\Delta|X = x, Z = z, U_D \leq P(z)) = E(\Delta|X = x)$$

$$\Leftrightarrow E(\Delta^{\text{MTE}}(X, U_D)|X = x, U_D \leq P(z)) = E(\Delta|X = x),$$

and hence

$$E(\Delta^{\text{MTE}}(X, U_D)|X = x, U_D \leq P(z)) = E(\Delta^{\text{MTE}}(X, U_D)|X = x, U_D > P(z)).$$

If the support of  $P(Z)$  is the full unit interval conditional on  $X = x$ , then  $\Delta^{\text{MTE}}(X, U_D) = E(\Delta|X = x)$  for all  $U_D$ .

- If the support of  $P(Z)$  is a proper subset of the full unit interval, then generically (M-3) will hold only if  $\Delta^{\text{MTE}}(X, U_D) = E(\Delta|X = x)$  for all  $U_D$ , though positive and negative parts could balance out for any particular value of  $X$ .

- To see this, note that

$$\begin{aligned} & E_Z \left( E(\Delta^{\text{MTE}}(X, U_D) | X = x, U_D \leq P(z)) | X = x, D = 1 \right) \\ &= E_Z \left( E(\Delta^{\text{MTE}}(X, U_D) | X = x, U_D > P(z)) | X = x, D = 0 \right). \end{aligned}$$

- To see this, note that

$$\begin{aligned} & E_Z \left( E(\Delta^{\text{MTE}}(X, U_D) | X = x, U_D \leq P(z)) | X = x, D = 1 \right) \\ &= E_Z \left( E(\Delta^{\text{MTE}}(X, U_D) | X = x, U_D > P(z)) | X = x, D = 0 \right). \end{aligned}$$

- Working with  $V = F_{V|X}^{-1}(U_D)$ , suppose that  $D = \mathbf{1}[\mu_D(Z, V) \geq 0]$ .

- To see this, note that

$$\begin{aligned} & E_Z \left( E(\Delta^{\text{MTE}}(X, U_D) | X = x, U_D \leq P(z)) | X = x, D = 1 \right) \\ &= E_Z \left( E(\Delta^{\text{MTE}}(X, U_D) | X = x, U_D > P(z)) | X = x, D = 0 \right). \end{aligned}$$

- Working with  $V = F_{V|X}^{-1}(U_D)$ , suppose that  $D = \mathbf{1}[\mu_D(Z, V) \geq 0]$ .
- Let  $\Omega(z) = \{v : \mu_D(z, v) \geq 0\}$ .

- To see this, note that

$$\begin{aligned} & E_Z \left( E(\Delta^{\text{MTE}}(X, U_D) | X = x, U_D \leq P(z)) | X = x, D = 1 \right) \\ &= E_Z \left( E(\Delta^{\text{MTE}}(X, U_D) | X = x, U_D > P(z)) | X = x, D = 0 \right). \end{aligned}$$

- Working with  $V = F_{V|X}^{-1}(U_D)$ , suppose that  $D = \mathbf{1}[\mu_D(Z, V) \geq 0]$ .
- Let  $\Omega(z) = \{v : \mu_D(z, v) \geq 0\}$ .
- Then (M-3) implies that

$$E(\Delta^{\text{MTE}}(X, V) | X = x, V \in \Omega(z)) = E(\Delta^{\text{MTE}}(X, V) | X = x, V \in (\Omega(z))^c)$$

so we expect that generically under assumption (M-3) we obtain a flat MTE in terms of  $V = F_{V|X}^{-1}(U_D)$ .

- To see this, note that

$$\begin{aligned} & E_Z \left( E(\Delta^{\text{MTE}}(X, U_D) | X = x, U_D \leq P(z)) | X = x, D = 1 \right) \\ &= E_Z \left( E(\Delta^{\text{MTE}}(X, U_D) | X = x, U_D > P(z)) | X = x, D = 0 \right). \end{aligned}$$

- Working with  $V = F_{V|X}^{-1}(U_D)$ , suppose that  $D = \mathbf{1}[\mu_D(Z, V) \geq 0]$ .
- Let  $\Omega(z) = \{v : \mu_D(z, v) \geq 0\}$ .
- Then (M-3) implies that

$$E(\Delta^{\text{MTE}}(X, V) | X = x, V \in \Omega(z)) = E(\Delta^{\text{MTE}}(X, V) | X = x, V \in (\Omega(z))^c)$$

so we expect that generically under assumption (M-3) we obtain a flat MTE in terms of  $V = F_{V|X}^{-1}(U_D)$ .

- We conduct a parallel analysis for the nonseparable choice model in Appendix, Slide 1155, and obtain similar conditions.

- Matching assumes a flat MTE, i.e., that marginal returns conditional on  $X$  and  $V$  do not depend on  $V$  (alternatively, that marginal returns do not depend on  $U_D$  given  $X$ ).

- Matching assumes a flat MTE, i.e., that marginal returns conditional on  $X$  and  $V$  do not depend on  $V$  (alternatively, that marginal returns do not depend on  $U_D$  given  $X$ ).
- We already noted in Slide 12 that IV and matching invoke very different assumptions.



- Matching assumes a flat MTE, i.e., that marginal returns conditional on  $X$  and  $V$  do not depend on  $V$  (alternatively, that marginal returns do not depend on  $U_D$  given  $X$ ).
- We already noted in Slide 12 that IV and matching invoke very different assumptions.
- Matching requires no exclusion restrictions whereas IV is based on the existence of exclusion restrictions.

- Matching assumes a flat MTE, i.e., that marginal returns conditional on  $X$  and  $V$  do not depend on  $V$  (alternatively, that marginal returns do not depend on  $U_D$  given  $X$ ).
- We already noted in Slide 12 that IV and matching invoke very different assumptions.
- Matching requires no exclusion restrictions whereas IV is based on the existence of exclusion restrictions.
- Superficially, we can bridge these literatures by invoking matching with an exclusion condition:  $(Y_0, Y_1) \not\perp\!\!\!\perp D \mid X$  but  $(Y_0, Y_1) \perp\!\!\!\perp D \mid X, Z$ .

- Matching assumes a flat MTE, i.e., that marginal returns conditional on  $X$  and  $V$  do not depend on  $V$  (alternatively, that marginal returns do not depend on  $U_D$  given  $X$ ).
- We already noted in Slide 12 that IV and matching invoke very different assumptions.
- Matching requires no exclusion restrictions whereas IV is based on the existence of exclusion restrictions.
- Superficially, we can bridge these literatures by invoking matching with an exclusion condition:  $(Y_0, Y_1) \not\perp\!\!\!\perp D \mid X$  but  $(Y_0, Y_1) \perp\!\!\!\perp D \mid X, Z$ .
- This looks like an IV condition, but it is not.

- We explore the relationship between matching with exclusion and IV in Appendix, Slide 1163, and demonstrate a fundamental contradiction between the two identifying conditions.

- We explore the relationship between matching with exclusion and IV in Appendix, Slide 1163, and demonstrate a fundamental contradiction between the two identifying conditions.
- For an additively separable representation of the outcome equations  $U_1 = Y_1 - E(Y_1|X)$  and  $U_0 = Y_0 - E(Y_0|X)$ , we establish that if  $(U_0, U_1)$  is mean independent of  $D$  conditional on  $(X, Z)$ , as required by IV, but  $(U_0, U_1)$  is not mean independent of  $D$  conditional on  $X$  alone, then  $U_0$  is dependent on  $Z$  conditional on  $X$ , contrary to all assumptions used to justify instrumental variables.

- We explore the relationship between matching with exclusion and IV in Appendix, Slide 1163, and demonstrate a fundamental contradiction between the two identifying conditions.
- For an additively separable representation of the outcome equations  $U_1 = Y_1 - E(Y_1|X)$  and  $U_0 = Y_0 - E(Y_0|X)$ , we establish that if  $(U_0, U_1)$  is mean independent of  $D$  conditional on  $(X, Z)$ , as required by IV, but  $(U_0, U_1)$  is not mean independent of  $D$  conditional on  $X$  alone, then  $U_0$  is dependent on  $Z$  conditional on  $X$ , contrary to all assumptions used to justify instrumental variables.
- We next consider how to implement matching.

## Implementing the Method of Matching

- We draw on [Dehejais and Todd](#) and [Lechner](#) to describe the mechanics of matching. [Lechner](#) presents a comprehensive treatment of the main issues and a guide to software.

- To operationalize the method of matching, we assume two samples: “ $t$ ” for treatment and “ $c$ ” for comparison group.



- To operationalize the method of matching, we assume two samples: “ $t$ ” for treatment and “ $c$ ” for comparison group.
- Treatment group members have  $D = 1$  and control group members have  $D = 0$ .

- To operationalize the method of matching, we assume two samples: “ $t$ ” for treatment and “ $c$ ” for comparison group.
- Treatment group members have  $D = 1$  and control group members have  $D = 0$ .
- Unless otherwise noted, we assume that observations are statistically independent within and across groups.

- To operationalize the method of matching, we assume two samples: “ $t$ ” for treatment and “ $c$ ” for comparison group.
- Treatment group members have  $D = 1$  and control group members have  $D = 0$ .
- Unless otherwise noted, we assume that observations are statistically independent within and across groups.
- Simple matching methods are based on the following idea.

- To operationalize the method of matching, we assume two samples: “ $t$ ” for treatment and “ $c$ ” for comparison group.
- Treatment group members have  $D = 1$  and control group members have  $D = 0$ .
- Unless otherwise noted, we assume that observations are statistically independent within and across groups.
- Simple matching methods are based on the following idea.
- For each person  $i$  in the treatment group, we find some group of “comparable” persons.

- To operationalize the method of matching, we assume two samples: “ $t$ ” for treatment and “ $c$ ” for comparison group.
- Treatment group members have  $D = 1$  and control group members have  $D = 0$ .
- Unless otherwise noted, we assume that observations are statistically independent within and across groups.
- Simple matching methods are based on the following idea.
- For each person  $i$  in the treatment group, we find some group of “comparable” persons.
- The same individual may be in both treated and control groups if that person is treated at one time and untreated at another.

- We denote outcomes for person  $i$  in the treatment group by  $Y_i^t$  and we match these outcomes to the outcomes of a subsample of persons in the comparison group to estimate a treatment effect.

- We denote outcomes for person  $i$  in the treatment group by  $Y_i^t$  and we match these outcomes to the outcomes of a subsample of persons in the comparison group to estimate a treatment effect.
- In principle, we can use a different subsample as a comparison group for each person.

- We denote outcomes for person  $i$  in the treatment group by  $Y_i^t$  and we match these outcomes to the outcomes of a subsample of persons in the comparison group to estimate a treatment effect.
- In principle, we can use a different subsample as a comparison group for each person.
- In practice, we can construct matches on the basis of a neighborhood  $\xi(X_i)$ , where  $X_i$  is a vector of characteristics for person  $i$ .



- We denote outcomes for person  $i$  in the treatment group by  $Y_i^t$  and we match these outcomes to the outcomes of a subsample of persons in the comparison group to estimate a treatment effect.
- In principle, we can use a different subsample as a comparison group for each person.
- In practice, we can construct matches on the basis of a neighborhood  $\xi(X_i)$ , where  $X_i$  is a vector of characteristics for person  $i$ .
- Neighbors to treated person  $i$  are persons in the comparison sample whose characteristics are in neighborhood  $\xi(X_i)$ .

- We denote outcomes for person  $i$  in the treatment group by  $Y_i^t$  and we match these outcomes to the outcomes of a subsample of persons in the comparison group to estimate a treatment effect.
- In principle, we can use a different subsample as a comparison group for each person.
- In practice, we can construct matches on the basis of a neighborhood  $\xi(X_i)$ , where  $X_i$  is a vector of characteristics for person  $i$ .
- Neighbors to treated person  $i$  are persons in the comparison sample whose characteristics are in neighborhood  $\xi(X_i)$ .
- Suppose that there are  $N_c$  persons in the comparison sample and  $N_t$  in the treatment sample.

- Thus the persons in the comparison sample who are neighbors to  $i$ , are persons  $j$  for whom  $X_j \in \xi(X_i)$ , i.e., the set of persons  $\mathcal{A}_i = \{j \mid X_j \in \xi(X_i)\}$ .

- Thus the persons in the comparison sample who are neighbors to  $i$ , are persons  $j$  for whom  $X_j \in \xi(X_i)$ , i.e., the set of persons  $\mathcal{A}_i = \{j \mid X_j \in \xi(X_i)\}$ .
- Let  $W(i, j)$  be the weight placed on observation  $j$  in forming a comparison with observation  $i$  and further assume that the weights sum to one,  $\sum_{j=1}^{N_c} W(i, j) = 1$ , and that  $0 \leq W(i, j) \leq 1$ .

- Thus the persons in the comparison sample who are neighbors to  $i$ , are persons  $j$  for whom  $X_j \in \xi(X_i)$ , i.e., the set of persons  $\mathcal{A}_i = \{j \mid X_j \in \xi(X_i)\}$ .
- Let  $W(i, j)$  be the weight placed on observation  $j$  in forming a comparison with observation  $i$  and further assume that the weights sum to one,  $\sum_{j=1}^{N_c} W(i, j) = 1$ , and that  $0 \leq W(i, j) \leq 1$ .
- Form a weighted comparison group mean for person  $i$ , given by

$$\bar{Y}_i^c = \sum_{j=1}^{N_c} W(i, j) Y_j^c. \quad (62)$$

- Thus the persons in the comparison sample who are neighbors to  $i$ , are persons  $j$  for whom  $X_j \in \xi(X_i)$ , i.e., the set of persons  $\mathcal{A}_i = \{j \mid X_j \in \xi(X_i)\}$ .
- Let  $W(i, j)$  be the weight placed on observation  $j$  in forming a comparison with observation  $i$  and further assume that the weights sum to one,  $\sum_{j=1}^{N_c} W(i, j) = 1$ , and that  $0 \leq W(i, j) \leq 1$ .
- Form a weighted comparison group mean for person  $i$ , given by

$$\bar{Y}_i^c = \sum_{j=1}^{N_c} W(i, j) Y_j^c. \quad (62)$$

- The estimated treatment effect for person  $i$  is  $Y_i - \bar{Y}_i^c$ .

- This selects a set of comparison group members associated with  $i$  and the mean of their outcomes.

- This selects a set of comparison group members associated with  $i$  and the mean of their outcomes.
- Unlike IV or the control function approach, the method of matching identifies counterfactuals for each treated member.



- This selects a set of comparison group members associated with  $i$  and the mean of their outcomes.
- Unlike IV or the control function approach, the method of matching identifies counterfactuals for each treated member.
- ? and ? survey a variety of alternative matching schemes proposed in the literature.

- This selects a set of comparison group members associated with  $i$  and the mean of their outcomes.
- Unlike IV or the control function approach, the method of matching identifies counterfactuals for each treated member.
- ? and ? survey a variety of alternative matching schemes proposed in the literature.
- ?? provides a comprehensive survey.

- This selects a set of comparison group members associated with  $i$  and the mean of their outcomes.
- Unlike IV or the control function approach, the method of matching identifies counterfactuals for each treated member.
- ? and ? survey a variety of alternative matching schemes proposed in the literature.
- ?? provides a comprehensive survey.
- In this chapter, we briefly consider two widely-used methods.

- The nearest-neighbor matching estimator defines  $\mathcal{A}_i$  such that only one  $j$  is selected so that it is closest to  $X_i$  in some metric:

$$\mathcal{A}_i = \{j \mid \min_{j \in \{1, \dots, N_c\}} \|X_i - X_j\|\},$$

where “ $\| \ \|$ ” is a metric measuring distance in the  $X$  characteristics space.

- The nearest-neighbor matching estimator defines  $\mathcal{A}_i$  such that only one  $j$  is selected so that it is closest to  $X_i$  in some metric:

$$\mathcal{A}_i = \{j \mid \min_{j \in \{1, \dots, N_c\}} \|X_i - X_j\|\},$$

where “ $\| \ \|$ ” is a metric measuring distance in the  $X$  characteristics space.

- The Mahalanobis metric is one widely used metric for implementing the nearest neighbor matching estimator.

- The nearest-neighbor matching estimator defines  $\mathcal{A}_i$  such that only one  $j$  is selected so that it is closest to  $X_i$  in some metric:

$$\mathcal{A}_i = \{j \mid \min_{j \in \{1, \dots, N_c\}} \|X_i - X_j\|\},$$

where “ $\| \ \|$ ” is a metric measuring distance in the  $X$  characteristics space.

- The Mahalanobis metric is one widely used metric for implementing the nearest neighbor matching estimator.
- This metric defines neighborhoods for  $i$  as

$$\| \ \| = (X_i - X_j)' \sum_c^{-1} (X_i - X_j),$$

where  $\sum_c$  is the covariance matrix in the comparison sample.

- The weighting scheme for the nearest neighbor matching estimator is

$$W(i, j) = \begin{cases} 1 & \text{if } j \in \mathcal{A}_i, \\ 0 & \text{otherwise.} \end{cases}$$

The nearest neighbor in the metric “ $\|\cdot\|$ ” is used in the match.

- The weighting scheme for the nearest neighbor matching estimator is

$$W(i, j) = \begin{cases} 1 & \text{if } j \in \mathcal{A}_i, \\ 0 & \text{otherwise.} \end{cases}$$

The nearest neighbor in the metric “ $\|\cdot\|$ ” is used in the match.

- A version of nearest-neighbor matching, called “caliper” matching (?), makes matches to person  $i$  only if

$$\|X_i - X_j\| < \varepsilon,$$

where  $\varepsilon$  is a pre-specified tolerance.



- The weighting scheme for the nearest neighbor matching estimator is

$$W(i, j) = \begin{cases} 1 & \text{if } j \in \mathcal{A}_i, \\ 0 & \text{otherwise.} \end{cases}$$

The nearest neighbor in the metric “ $\|\cdot\|$ ” is used in the match.

- A version of nearest-neighbor matching, called “caliper” matching (?), makes matches to person  $i$  only if

$$\|X_i - X_j\| < \varepsilon,$$

where  $\varepsilon$  is a pre-specified tolerance.

- Otherwise, person  $i$  is bypassed and no match is made to him or her.

- Kernel matching uses the entire comparison sample, so that  $\mathcal{A}_i = \{1, \dots, N_c\}$ , and sets

$$W(i, j) = \frac{K(X_j - X_i)}{\sum_{j=1}^{N_c} K(X_j - X_i)},$$

where  $K$  is a kernel.

- Kernel matching uses the entire comparison sample, so that  $\mathcal{A}_i = \{1, \dots, N_c\}$ , and sets

$$W(i, j) = \frac{K(X_j - X_i)}{\sum_{j=1}^{N_c} K(X_j - X_i)},$$

where  $K$  is a kernel.

- Kernel matching is a smooth method that reuses and weights the comparison group sample observations differently for each person  $i$  in the treatment group with a different  $X_i$ .

- Kernel matching uses the entire comparison sample, so that  $\mathcal{A}_i = \{1, \dots, N_c\}$ , and sets

$$W(i, j) = \frac{K(X_j - X_i)}{\sum_{j=1}^{N_c} K(X_j - X_i)},$$

where  $K$  is a kernel.

- Kernel matching is a smooth method that reuses and weights the comparison group sample observations differently for each person  $i$  in the treatment group with a different  $X_i$ .
- Kernel matching can be defined pointwise at each sample point  $X_i$  or for broader intervals.

- For example, the impact of treatment on the treated can be estimated by forming the mean difference across the  $i$ :

$$\hat{\Delta}^{TT} = \frac{1}{N_t} \sum_{i=1}^{N_t} (Y_i^t - \bar{Y}_i^c) = \frac{1}{N_t} \sum_{i=1}^{N_t} (Y_i^t - \sum_{j=1}^{N_c} W(i,j) Y_j^c). \quad (63)$$

We can define this mean for various subsets of the treatment sample defined in various ways.

- For example, the impact of treatment on the treated can be estimated by forming the mean difference across the  $i$ :

$$\hat{\Delta}^{TT} = \frac{1}{N_t} \sum_{i=1}^{N_t} (Y_i^t - \bar{Y}_i^c) = \frac{1}{N_t} \sum_{i=1}^{N_t} (Y_i^t - \sum_{j=1}^{N_c} W(i,j) Y_j^c). \quad (63)$$

We can define this mean for various subsets of the treatment sample defined in various ways.

- More efficient estimators weight the observations accounting for the variance (?????) .

- Matching assumes that conditioning on  $X$  eliminates selection bias.

- Matching assumes that conditioning on  $X$  eliminates selection bias.
- The method requires no functional form assumptions for outcome equations.



- Matching assumes that conditioning on  $X$  eliminates selection bias.
- The method requires no functional form assumptions for outcome equations.
- If, however, a functional form assumption is maintained, as in the econometric procedure proposed by ?, it is possible to implement the matching assumption using standard regression analysis.

- Suppose, for example, that  $Y_0$  is linearly related to observables  $X$  and an unobservable  $U_0$ , so that

$$E(Y_0 | X, D = 0) = X\alpha + E(U_0 | X, D = 0),$$

and

$$E(U_0 | X, D = 0) = E(U_0 | X)$$

is linear in  $X$  ( $E(U | X) = \varphi X$ ).

- Suppose, for example, that  $Y_0$  is linearly related to observables  $X$  and an unobservable  $U_0$ , so that

$$E(Y_0 | X, D = 0) = X\alpha + E(U_0 | X, D = 0),$$

and

$$E(U_0 | X, D = 0) = E(U_0 | X)$$

is linear in  $X$  ( $E(U | X) = \varphi X$ ).

- Under these assumptions, controlling for  $X$  via linear regression allows one to identify  $E(Y_0 | X, D = 1)$  from the data on nonparticipants.

- Suppose, for example, that  $Y_0$  is linearly related to observables  $X$  and an unobservable  $U_0$ , so that

$$E(Y_0 | X, D = 0) = X\alpha + E(U_0 | X, D = 0),$$

and

$$E(U_0 | X, D = 0) = E(U_0 | X)$$

is linear in  $X$  ( $E(U | X) = \varphi X$ ).

- Under these assumptions, controlling for  $X$  via linear regression allows one to identify  $E(Y_0 | X, D = 1)$  from the data on nonparticipants.
- Under assumption (Q-4)', setting  $X = Q$ , this approach justifies OLS equation (Q-3) for identifying treatment effects.

- Such functional form assumptions are not strictly required to implement the method of matching.

- Such functional form assumptions are not strictly required to implement the method of matching.
- Moreover, in practice, users of the method of ? do not impose the common support condition (M-2) for the distribution of  $X$  when generating estimates of the treatment effect.

- Such functional form assumptions are not strictly required to implement the method of matching.
- Moreover, in practice, users of the method of ? do not impose the common support condition (M-2) for the distribution of  $X$  when generating estimates of the treatment effect.
- The distribution of  $X$  may be very different in the treatment group ( $D = 1$ ) and comparison group ( $D = 0$ ) samples, so that comparability is only achieved by imposing linearity in the parameters and extrapolating over different regions.

- Such functional form assumptions are not strictly required to implement the method of matching.
- Moreover, in practice, users of the method of ? do not impose the common support condition (M-2) for the distribution of  $X$  when generating estimates of the treatment effect.
- The distribution of  $X$  may be very different in the treatment group ( $D = 1$ ) and comparison group ( $D = 0$ ) samples, so that comparability is only achieved by imposing linearity in the parameters and extrapolating over different regions.
- One advantage of the method of ? is that it uses data parsimoniously.



- Such functional form assumptions are not strictly required to implement the method of matching.
- Moreover, in practice, users of the method of ? do not impose the common support condition (M-2) for the distribution of  $X$  when generating estimates of the treatment effect.
- The distribution of  $X$  may be very different in the treatment group ( $D = 1$ ) and comparison group ( $D = 0$ ) samples, so that comparability is only achieved by imposing linearity in the parameters and extrapolating over different regions.
- One advantage of the method of ? is that it uses data parsimoniously.
- If the  $X$  are high dimensional, the number of observations in each cell when matching can get very small.

- Another solution to this problem that reduces the dimension of the matching problem without imposing arbitrary linearity assumptions is based on the probability of participation or the “propensity score,”  $P(X) = \Pr(D = 1 | X)$ .

- Another solution to this problem that reduces the dimension of the matching problem without imposing arbitrary linearity assumptions is based on the probability of participation or the “propensity score,”  $P(X) = \Pr(D = 1 | X)$ .
- ? demonstrate that under assumptions (M-1) and (M-2),

$$(Y_0, Y_1) \perp\!\!\!\perp D | P(X) \text{ for } X \in \chi_c, \quad (64)$$

for some set  $\chi_c$ , where it is assumed that (M-2) holds in the set.

- Another solution to this problem that reduces the dimension of the matching problem without imposing arbitrary linearity assumptions is based on the probability of participation or the “propensity score,”  $P(X) = \Pr(D = 1 | X)$ .
- ? demonstrate that under assumptions (M-1) and (M-2),

$$(Y_0, Y_1) \perp\!\!\!\perp D \mid P(X) \text{ for } X \in \chi_c, \quad (64)$$

for some set  $\chi_c$ , where it is assumed that (M-2) holds in the set.

- Conditioning either on  $P(X)$  or on  $X$  produces conditional independence.

- Conditioning on  $P(X)$  reduces the dimension of the matching problem down to matching on the scalar  $P(X)$ .

- Conditioning on  $P(X)$  reduces the dimension of the matching problem down to matching on the scalar  $P(X)$ .
- The analysis of ? assumes that  $P(X)$  is known rather than estimated.

- Conditioning on  $P(X)$  reduces the dimension of the matching problem down to matching on the scalar  $P(X)$ .
- The analysis of ? assumes that  $P(X)$  is known rather than estimated.
- ?, ?, and ? present the asymptotic distribution theory for the kernel matching estimator in the cases in which  $P(X)$  is known and in which it is estimated both parametrically and nonparametrically.

- Conditioning on  $P$  identifies all treatment parameters but as we have seen, it imposes the assumption of a flat MTE.



- Conditioning on  $P$  identifies all treatment parameters but as we have seen, it imposes the assumption of a flat MTE.
- Marginal returns and average returns are the same.

- Conditioning on  $P$  identifies all treatment parameters but as we have seen, it imposes the assumption of a flat MTE.
- Marginal returns and average returns are the same.
- A consequence of (64) is that

$$\begin{aligned} E(Y_1|D = 0, P(X)) &= E(Y_1|D = 1, P(X)) = E(Y_1|P(X)), \\ E(Y_0|D = 1, P(X)) &= E(Y_0|D = 0, P(X)) = E(Y_0|P(X)). \end{aligned}$$

- Support condition (M-2) has the unattractive feature that if the analyst has too much information about the decision of who takes treatment, so that  $P(X) = 1$  or  $0$ , the method breaks down at such values of  $X$  because people cannot be compared at a common  $X$ .

- Support condition (M-2) has the unattractive feature that if the analyst has too much information about the decision of who takes treatment, so that  $P(X) = 1$  or  $0$ , the method breaks down at such values of  $X$  because people cannot be compared at a common  $X$ .
- The method of matching assumes that, given  $X$ , some unspecified randomization in the economic environment allocates people to treatment.

- Support condition (M-2) has the unattractive feature that if the analyst has too much information about the decision of who takes treatment, so that  $P(X) = 1$  or  $0$ , the method breaks down at such values of  $X$  because people cannot be compared at a common  $X$ .
- The method of matching assumes that, given  $X$ , some unspecified randomization in the economic environment allocates people to treatment.
- This justifies assumption (Q-5) in the OLS example.

- Support condition (M-2) has the unattractive feature that if the analyst has too much information about the decision of who takes treatment, so that  $P(X) = 1$  or  $0$ , the method breaks down at such values of  $X$  because people cannot be compared at a common  $X$ .
- The method of matching assumes that, given  $X$ , some unspecified randomization in the economic environment allocates people to treatment.
- This justifies assumption (Q-5) in the OLS example.
- The fact that the cases  $P(X) = 1$  and  $P(X) = 0$  must be eliminated suggests that methods for choosing  $X$  based on the fit of the model to data on  $D$  are potentially problematic, as we discuss below.

- Offsetting these disadvantages, the method of matching with a known conditioning set that produces condition (M-2) does not require separability of outcome or choice equations, exogeneity of conditioning variables, exclusion restrictions, or adoption of specific functional forms of outcome equations.

- Offsetting these disadvantages, the method of matching with a known conditioning set that produces condition (M-2) does not require separability of outcome or choice equations, exogeneity of conditioning variables, exclusion restrictions, or adoption of specific functional forms of outcome equations.
- Such features are commonly used in conventional selection (control function) methods and conventional applications of IV although as we have demonstrated in Slide 152, recent work in semiparametric estimation relaxes these assumptions.



- Offsetting these disadvantages, the method of matching with a known conditioning set that produces condition (M-2) does not require separability of outcome or choice equations, exogeneity of conditioning variables, exclusion restrictions, or adoption of specific functional forms of outcome equations.
- Such features are commonly used in conventional selection (control function) methods and conventional applications of IV although as we have demonstrated in Slide 152, recent work in semiparametric estimation relaxes these assumptions.
- As noted in Slide 700, the method of matching does not strictly require (M-1).

- Offsetting these disadvantages, the method of matching with a known conditioning set that produces condition (M-2) does not require separability of outcome or choice equations, exogeneity of conditioning variables, exclusion restrictions, or adoption of specific functional forms of outcome equations.
- Such features are commonly used in conventional selection (control function) methods and conventional applications of IV although as we have demonstrated in Slide 152, recent work in semiparametric estimation relaxes these assumptions.
- As noted in Slide 700, the method of matching does not strictly require (M-1).
- One can get by with weaker mean independence assumptions (M-3) in the place of the stronger conditions (M-1).

- However, if (M-3) is invoked, the assumption that one can replace  $X$  by  $P(X)$  does not follow from the analysis of ?, and is an additional new assumption.

- However, if (M-3) is invoked, the assumption that one can replace  $X$  by  $P(X)$  does not follow from the analysis of ?, and is an additional new assumption.
- Methods for implementing matching are provided in ? and are discussed extensively in ?.

- However, if (M-3) is invoked, the assumption that one can replace  $X$  by  $P(X)$  does not follow from the analysis of ?, and is an additional new assumption.
- Methods for implementing matching are provided in ? and are discussed extensively in ?.
- See ??? for software and extensive discussion of the mechanics of matching.

- However, if (M-3) is invoked, the assumption that one can replace  $X$  by  $P(X)$  does not follow from the analysis of ?, and is an additional new assumption.
- Methods for implementing matching are provided in ? and are discussed extensively in ?.
- See ??? for software and extensive discussion of the mechanics of matching.
- We now contrast the identifying assumptions used in the method of control functions with those used in matching.

## Comparing Matching and Control Functions Approaches

- The method of matching eliminates the dependence between  $(Y_0, Y_1)$  and  $D$ ,  $(Y_0, Y_1) \not\perp\!\!\!\perp D$ , by assuming access to conditioning variables  $X$  such that (M-1) is satisfied:  $(Y_0, Y_1) \perp\!\!\!\perp D \mid X$ .

## Comparing Matching and Control Functions Approaches

- The method of matching eliminates the dependence between  $(Y_0, Y_1)$  and  $D$ ,  $(Y_0, Y_1) \not\perp\!\!\!\perp D$ , by assuming access to conditioning variables  $X$  such that (M-1) is satisfied:  $(Y_0, Y_1) \perp\!\!\!\perp D \mid X$ .
- By conditioning on observables, one can identify the distributions of  $Y_0$  and  $Y_1$  over the support of  $X$  satisfying (M-2).



- Other methods model the dependence that gives rise to the spurious relationship and in this way attempt to eliminate it.

- Other methods model the dependence that gives rise to the spurious relationship and in this way attempt to eliminate it.
- IV involves exclusion and a different type of conditional independence,  $(Y_0, Y_1) \perp\!\!\!\perp Z \mid X$ , as well as a rank condition ( $\Pr(D = 1 \mid X, Z)$  depends on  $Z$ ).

- Other methods model the dependence that gives rise to the spurious relationship and in this way attempt to eliminate it.
- IV involves exclusion and a different type of conditional independence,  $(Y_0, Y_1) \perp\!\!\!\perp Z \mid X$ , as well as a rank condition ( $\Pr(D = 1 \mid X, Z)$  depends on  $Z$ ).
- The instrument  $Z$  plays the role of the implicit randomization used in matching by allocating people to treatment status in a way that does not depend on  $(Y_0, Y_1)$ .

- Other methods model the dependence that gives rise to the spurious relationship and in this way attempt to eliminate it.
- IV involves exclusion and a different type of conditional independence,  $(Y_0, Y_1) \perp\!\!\!\perp Z \mid X$ , as well as a rank condition ( $\Pr(D = 1 \mid X, Z)$  depends on  $Z$ ).
- The instrument  $Z$  plays the role of the implicit randomization used in matching by allocating people to treatment status in a way that does not depend on  $(Y_0, Y_1)$ .
- We have already established that matching and IV make very different assumptions.

- Other methods model the dependence that gives rise to the spurious relationship and in this way attempt to eliminate it.
- IV involves exclusion and a different type of conditional independence,  $(Y_0, Y_1) \perp\!\!\!\perp Z \mid X$ , as well as a rank condition ( $\Pr(D = 1 \mid X, Z)$  depends on  $Z$ ).
- The instrument  $Z$  plays the role of the implicit randomization used in matching by allocating people to treatment status in a way that does not depend on  $(Y_0, Y_1)$ .
- We have already established that matching and IV make very different assumptions.
- Thus, in general, a matching assumption that  $(Y_0, Y_1) \perp\!\!\!\perp D \mid X, Z$  neither implies nor is implied by  $(Y_0, Y_1) \perp\!\!\!\perp Z \mid X$ .

- One special case where they are equivalent is when treatment status is assigned by randomization with full compliance (letting  $\xi = 1$  denote assignment to treatment,  $\xi = 1 \Rightarrow A = 1$  and  $\xi = 0 \Rightarrow A = 0$ ) and  $Z = \xi$ , so that the instrument is the assignment mechanism.

- One special case where they are equivalent is when treatment status is assigned by randomization with full compliance (letting  $\xi = 1$  denote assignment to treatment,  $\xi = 1 \Rightarrow A = 1$  and  $\xi = 0 \Rightarrow A = 0$ ) and  $Z = \xi$ , so that the instrument is the assignment mechanism.
- $A = 1$  if the person actually receives treatment, and  $A = 0$  otherwise.

- One special case where they are equivalent is when treatment status is assigned by randomization with full compliance (letting  $\xi = 1$  denote assignment to treatment,  $\xi = 1 \Rightarrow A = 1$  and  $\xi = 0 \Rightarrow A = 0$ ) and  $Z = \xi$ , so that the instrument is the assignment mechanism.
- $A = 1$  if the person actually receives treatment, and  $A = 0$  otherwise.
- The method of control functions explicitly models the dependence between  $(Y_0, Y_1)$  and  $D$  and attempts to eliminate it.



- One special case where they are equivalent is when treatment status is assigned by randomization with full compliance (letting  $\xi = 1$  denote assignment to treatment,  $\xi = 1 \Rightarrow A = 1$  and  $\xi = 0 \Rightarrow A = 0$ ) and  $Z = \xi$ , so that the instrument is the assignment mechanism.
- $A = 1$  if the person actually receives treatment, and  $A = 0$  otherwise.
- The method of control functions explicitly models the dependence between  $(Y_0, Y_1)$  and  $D$  and attempts to eliminate it.
- ? provides a comprehensive review of these methods.

- In Slide 1005, we present a summary of some of the general principles underlying the method of control functions, the method of control variates, replacement functions, and proxy approaches as they apply to the selection problem.

- In Slide 1005, we present a summary of some of the general principles underlying the method of control functions, the method of control variates, replacement functions, and proxy approaches as they apply to the selection problem.
- All of these methods attempt to eliminate the  $\theta$  in (U-1) that produces the dependence captured in (U-2).

- In Slide 1005, we present a summary of some of the general principles underlying the method of control functions, the method of control variates, replacement functions, and proxy approaches as they apply to the selection problem.
- All of these methods attempt to eliminate the  $\theta$  in (U-1) that produces the dependence captured in (U-2).
- In this section, we relate matching to the form of the control function introduced in ? and ??.

- In Slide 1005, we present a summary of some of the general principles underlying the method of control functions, the method of control variates, replacement functions, and proxy approaches as they apply to the selection problem.
- All of these methods attempt to eliminate the  $\theta$  in (U-1) that produces the dependence captured in (U-2).
- In this section, we relate matching to the form of the control function introduced in ? and ??.
- This version was used in our analysis of local instrumental variables (LIV) in Slide 152, where we compare LIV with control function approaches and show that LIV and LATE estimate derivatives of the control functions.

- We analyze conditional means because of their familiarity.

- We analyze conditional means because of their familiarity.
- Using the fact that  $E(\mathbf{1}(Y \leq y) | X) = F(y | X)$ , the analysis applies to marginal distributions as well.

- We analyze conditional means because of their familiarity.
- Using the fact that  $E(\mathbf{1}(Y \leq y) | X) = F(y | X)$ , the analysis applies to marginal distributions as well.
- Thus we work with conditional expectations of  $(Y_0, Y_1)$  given  $(X, Z, D)$ , where  $Z$  is assumed to include at least one variable not in  $X$ .



- We analyze conditional means because of their familiarity.
- Using the fact that  $E(\mathbf{1}(Y \leq y) | X) = F(y | X)$ , the analysis applies to marginal distributions as well.
- Thus we work with conditional expectations of  $(Y_0, Y_1)$  given  $(X, Z, D)$ , where  $Z$  is assumed to include at least one variable not in  $X$ .
- Conventional applications of the control function method assume additive separability, which is not required in matching.

- We analyze conditional means because of their familiarity.
- Using the fact that  $E(\mathbf{1}(Y \leq y) | X) = F(y | X)$ , the analysis applies to marginal distributions as well.
- Thus we work with conditional expectations of  $(Y_0, Y_1)$  given  $(X, Z, D)$ , where  $Z$  is assumed to include at least one variable not in  $X$ .
- Conventional applications of the control function method assume additive separability, which is not required in matching.
- Strictly speaking, additive separability is not required in the application of control functions either.

- We analyze conditional means because of their familiarity.
- Using the fact that  $E(\mathbf{1}(Y \leq y) | X) = F(y | X)$ , the analysis applies to marginal distributions as well.
- Thus we work with conditional expectations of  $(Y_0, Y_1)$  given  $(X, Z, D)$ , where  $Z$  is assumed to include at least one variable not in  $X$ .
- Conventional applications of the control function method assume additive separability, which is not required in matching.
- Strictly speaking, additive separability is not required in the application of control functions either.
- What is required is a model relating the outcome unobservables to the observables and the unobservables in the choice of treatment equation.

- Various assumptions give operational content to (U-1) defined in Slide 12.

- Various assumptions give operational content to (U-1) defined in Slide 12.
- For the additively separable case (2), the control function for mean outcomes models the conditional expectations of  $Y_1$  and  $Y_0$  given  $X$ ,  $Z$ , and  $D$  as

$$E(Y_1|Z, X, D = 1) = \mu_1(X) + E(U_1|Z, X, D = 1)$$

$$E(Y_0|Z, X, D = 0) = \mu_0(X) + E(U_0|Z, X, D = 0).$$

- In the traditional method of control functions, the analyst models  $E(U_1|Z, X, D = 1)$  and  $E(U_0|Z, X, D = 0)$ .

- In the traditional method of control functions, the analyst models  $E(U_1|Z, X, D = 1)$  and  $E(U_0|Z, X, D = 0)$ .
- If these functions can be independently varied against  $\mu_1(X)$  and  $\mu_0(X)$  respectively, one can identify  $\mu_1(X)$  and  $\mu_0(X)$  up to constant terms.

- In the traditional method of control functions, the analyst models  $E(U_1|Z, X, D = 1)$  and  $E(U_0|Z, X, D = 0)$ .
- If these functions can be independently varied against  $\mu_1(X)$  and  $\mu_0(X)$  respectively, one can identify  $\mu_1(X)$  and  $\mu_0(X)$  up to constant terms.
- It is not required that  $X$  or  $Z$  be stochastically independent of  $U_1$  or  $U_0$ , although conventional methods often assume this.



- Assume that  $(U_0, U_1, V) \perp\!\!\!\perp (X, Z)$  and adopt equation (7) as the treatment choice model augmented so that  $X$  and  $Z$  are determinants of treatment choice, using  $V$  as the latent variable that generates  $D$  given  $X, Z$ :  $D = \mathbf{1}(\mu_D(Z) \geq 0)$ .

- Assume that  $(U_0, U_1, V) \perp\!\!\!\perp (X, Z)$  and adopt equation (7) as the treatment choice model augmented so that  $X$  and  $Z$  are determinants of treatment choice, using  $V$  as the latent variable that generates  $D$  given  $X, Z$ :  $D = \mathbf{1}(\mu_D(Z) \geq 0)$ .
- Let  $U_D = F_{V|X}(V)$  and  $P(Z) = F_{V|X}(\mu_D(Z))$ .

- Assume that  $(U_0, U_1, V) \perp\!\!\!\perp (X, Z)$  and adopt equation (7) as the treatment choice model augmented so that  $X$  and  $Z$  are determinants of treatment choice, using  $V$  as the latent variable that generates  $D$  given  $X, Z$ :  $D = \mathbf{1}(\mu_D(Z) \geq 0)$ .
- Let  $U_D = F_{V|X}(V)$  and  $P(Z) = F_{V|X}(\mu_D(Z))$ .
- In this notation, the control functions are

$$\begin{aligned} E(U_1|Z, D=1) &= E(U_1|\mu_D(Z) \geq V) = E(U_1 | P(Z) \geq U_D) = K_1(P(Z)) \text{ and} \\ E(U_0|Z, D=0) &= E(U_0|\mu_D(Z) < V) = E(U_0 | P(Z) < U_D) = K_0(P(Z)), \end{aligned}$$

so the control function only depends on the propensity score  $P(Z)$ .

- Assume that  $(U_0, U_1, V) \perp\!\!\!\perp (X, Z)$  and adopt equation (7) as the treatment choice model augmented so that  $X$  and  $Z$  are determinants of treatment choice, using  $V$  as the latent variable that generates  $D$  given  $X, Z$ :  $D = \mathbf{1}(\mu_D(Z) \geq 0)$ .
- Let  $U_D = F_{V|X}(V)$  and  $P(Z) = F_{V|X}(\mu_D(Z))$ .
- In this notation, the control functions are

$$\begin{aligned} E(U_1|Z, D=1) &= E(U_1|\mu_D(Z) \geq V) = E(U_1 | P(Z) \geq U_D) = K_1(P(Z)) \text{ and} \\ E(U_0|Z, D=0) &= E(U_0|\mu_D(Z) < V) = E(U_0 | P(Z) < U_D) = K_0(P(Z)), \end{aligned}$$

so the control function only depends on the propensity score  $P(Z)$ .

- The key assumption needed to represent the control function solely as a function of  $P(Z)$  is

$$(U_0, U_1, V) \perp\!\!\!\perp X, Z. \quad (\text{CF-1})$$

- This assumption is not strictly required but it is traditional and useful in relating LIV and selection models (as in Slide 152) and selection models and matching (this section).

- This assumption is not strictly required but it is traditional and useful in relating LIV and selection models (as in Slide 152) and selection models and matching (this section).
- Under this condition

$$E(Y_1|Z, X, D = 1) = \mu_1(X) + K_1(P(Z)),$$

$$E(Y_0|Z, X, D = 0) = \mu_0(X) + K_0(P(Z)),$$

with  $\lim_{P \rightarrow 1} K_1(P) = 0$  and  $\lim_{P \rightarrow 0} K_0(P) = 0$ .

- This assumption is not strictly required but it is traditional and useful in relating LIV and selection models (as in Slide 152) and selection models and matching (this section).
- Under this condition

$$E(Y_1|Z, X, D = 1) = \mu_1(X) + K_1(P(Z)),$$

$$E(Y_0|Z, X, D = 0) = \mu_0(X) + K_0(P(Z)),$$

with  $\lim_{P \rightarrow 1} K_1(P) = 0$  and  $\lim_{P \rightarrow 0} K_0(P) = 0$ .

- It is assumed that  $Z$  can be independently varied for all  $X$ , and the limits are obtained by changing  $Z$  while holding  $X$  fixed.

- This assumption is not strictly required but it is traditional and useful in relating LIV and selection models (as in Slide 152) and selection models and matching (this section).
- Under this condition

$$E(Y_1|Z, X, D = 1) = \mu_1(X) + K_1(P(Z)),$$

$$E(Y_0|Z, X, D = 0) = \mu_0(X) + K_0(P(Z)),$$

with  $\lim_{P \rightarrow 1} K_1(P) = 0$  and  $\lim_{P \rightarrow 0} K_0(P) = 0$ .

- It is assumed that  $Z$  can be independently varied for all  $X$ , and the limits are obtained by changing  $Z$  while holding  $X$  fixed.
- These limit results state that when the values of  $X, Z$  are such that the probability of being in a sample ( $D = 1$  or  $D = 0$ , respectively) is 1, there is no selection bias and one can separate out  $\mu_1(X)$  from  $K_1(P(Z))$  and  $\mu_0(X)$  from  $K_0(P(Z))$ .



- This is the same identification at infinity condition that is required to identify ATE and TT in IV for models with heterogeneous responses.

- This is the same identification at infinity condition that is required to identify ATE and TT in IV for models with heterogeneous responses.
- As noted in Slide 152, unlike the method of matching based on (M-1), the method of control functions allows the marginal treatment effect to be different from the average treatment effect and from the conditional effect of treatment on the treated.

- This is the same identification at infinity condition that is required to identify ATE and TT in IV for models with heterogeneous responses.
- As noted in Slide 152, unlike the method of matching based on (M-1), the method of control functions allows the marginal treatment effect to be different from the average treatment effect and from the conditional effect of treatment on the treated.
- Although conventional practice has been to derive the functional forms of  $K_0(P)$  and  $K_1(P)$  by making distributional assumptions about  $(U_0, U_1, V)$  such as normality or other conventional distributional assumptions, this is not an intrinsic feature of the method and there are many nonnormal and semiparametric versions of this method.

- This is the same identification at infinity condition that is required to identify ATE and TT in IV for models with heterogeneous responses.
- As noted in Slide 152, unlike the method of matching based on (M-1), the method of control functions allows the marginal treatment effect to be different from the average treatment effect and from the conditional effect of treatment on the treated.
- Although conventional practice has been to derive the functional forms of  $K_0(P)$  and  $K_1(P)$  by making distributional assumptions about  $(U_0, U_1, V)$  such as normality or other conventional distributional assumptions, this is not an intrinsic feature of the method and there are many nonnormal and semiparametric versions of this method.
- See ? for a survey.

- In its semiparametric implementation, the method of control functions requires an exclusion restriction (a variable in  $Z$  not in  $X$ ) to achieve nonparametric identification.

- In its semiparametric implementation, the method of control functions requires an exclusion restriction (a variable in  $Z$  not in  $X$ ) to achieve nonparametric identification.
- Without any functional-form assumptions one cannot rule out a worst case analysis where, for example, if  $X = Z$ , then  $K_1(P(X)) = \tau\mu(X)$  where  $\tau$  is a scalar.

- In its semiparametric implementation, the method of control functions requires an exclusion restriction (a variable in  $Z$  not in  $X$ ) to achieve nonparametric identification.
- Without any functional-form assumptions one cannot rule out a worst case analysis where, for example, if  $X = Z$ , then  $K_1(P(X)) = \tau\mu(X)$  where  $\tau$  is a scalar.
- In this situation, there is perfect collinearity between the control function and the conditional mean of the outcome equation, and it is impossible to separately identify either.

- In its semiparametric implementation, the method of control functions requires an exclusion restriction (a variable in  $Z$  not in  $X$ ) to achieve nonparametric identification.
- Without any functional-form assumptions one cannot rule out a worst case analysis where, for example, if  $X = Z$ , then  $K_1(P(X)) = \tau\mu(X)$  where  $\tau$  is a scalar.
- In this situation, there is perfect collinearity between the control function and the conditional mean of the outcome equation, and it is impossible to separately identify either.
- Even though this case is not generic, it is possible.



- In its semiparametric implementation, the method of control functions requires an exclusion restriction (a variable in  $Z$  not in  $X$ ) to achieve nonparametric identification.
- Without any functional-form assumptions one cannot rule out a worst case analysis where, for example, if  $X = Z$ , then  $K_1(P(X)) = \tau\mu(X)$  where  $\tau$  is a scalar.
- In this situation, there is perfect collinearity between the control function and the conditional mean of the outcome equation, and it is impossible to separately identify either.
- Even though this case is not generic, it is possible.
- The method of matching does not require an exclusion restriction, but at the cost of ruling out essential heterogeneity.

- In the general case, the method of control functions requires that in certain limit sets of  $Z$ ,  $P(Z) = 1$  and  $P(Z) = 0$  in order to achieve full nonparametric identification.

- In the general case, the method of control functions requires that in certain limit sets of  $Z$ ,  $P(Z) = 1$  and  $P(Z) = 0$  in order to achieve full nonparametric identification.
- The conventional method of matching does not invoke such limit set arguments.

- In the general case, the method of control functions requires that in certain limit sets of  $Z$ ,  $P(Z) = 1$  and  $P(Z) = 0$  in order to achieve full nonparametric identification.
- The conventional method of matching does not invoke such limit set arguments.
- All methods of evaluation, including matching and control functions, require that treatment parameters be defined on a common support that is the intersection of the supports of  $X$  given  $D = 1$  and  $X$  given  $D = 0$ :  
 $\text{Supp}(X|D = 1) \cap \text{Supp}(X|D = 0)$ .

- In the general case, the method of control functions requires that in certain limit sets of  $Z$ ,  $P(Z) = 1$  and  $P(Z) = 0$  in order to achieve full nonparametric identification.
- The conventional method of matching does not invoke such limit set arguments.
- All methods of evaluation, including matching and control functions, require that treatment parameters be defined on a common support that is the intersection of the supports of  $X$  given  $D = 1$  and  $X$  given  $D = 0$ :  
 $\text{Supp}(X|D = 1) \cap \text{Supp}(X|D = 0)$ .
- This is the requirement for any estimator that seeks to identify treatment effects by comparing samples of treated persons with samples of untreated persons.

- In this version of the method of control functions,  $P(Z)$  is a conditioning variable used to predict  $U_1$  conditional on  $D$  and  $U_0$  conditional on  $D$ .

- In this version of the method of control functions,  $P(Z)$  is a conditioning variable used to predict  $U_1$  conditional on  $D$  and  $U_0$  conditional on  $D$ .
- In the method of matching, it is used as a conditioning variable to eliminate the stochastic dependence between  $(U_0, U_1)$  and  $D$ .

- In this version of the method of control functions,  $P(Z)$  is a conditioning variable used to predict  $U_1$  conditional on  $D$  and  $U_0$  conditional on  $D$ .
- In the method of matching, it is used as a conditioning variable to eliminate the stochastic dependence between  $(U_0, U_1)$  and  $D$ .
- In the method of LATE or LIV,  $P(Z)$  is used as an instrument.



- In this version of the method of control functions,  $P(Z)$  is a conditioning variable used to predict  $U_1$  conditional on  $D$  and  $U_0$  conditional on  $D$ .
- In the method of matching, it is used as a conditioning variable to eliminate the stochastic dependence between  $(U_0, U_1)$  and  $D$ .
- In the method of LATE or LIV,  $P(Z)$  is used as an instrument.
- In the method of control functions, as conventionally applied,  $(U_0, U_1) \perp\!\!\!\perp (X, Z)$ , but this assumption is not intrinsic to the method.

- In this version of the method of control functions,  $P(Z)$  is a conditioning variable used to predict  $U_1$  conditional on  $D$  and  $U_0$  conditional on  $D$ .
- In the method of matching, it is used as a conditioning variable to eliminate the stochastic dependence between  $(U_0, U_1)$  and  $D$ .
- In the method of LATE or LIV,  $P(Z)$  is used as an instrument.
- In the method of control functions, as conventionally applied,  $(U_0, U_1) \perp\!\!\!\perp (X, Z)$ , but this assumption is not intrinsic to the method.
- This assumption plays no role in matching if the correct conditioning set is known.

- However, as noted below, exogeneity plays a key role in devising algorithms to select the conditioning variables.

- However, as noted below, exogeneity plays a key role in devising algorithms to select the conditioning variables.
- In addition, as noted in Slide 412, exogeneity is helpful in making out-of-sample forecasts.

- However, as noted below, exogeneity plays a key role in devising algorithms to select the conditioning variables.
- In addition, as noted in Slide 412, exogeneity is helpful in making out-of-sample forecasts.
- The method of control functions does not require that  $(U_0, U_1) \perp\!\!\!\perp D \mid (X, Z)$ , which is a central requirement of matching.

- However, as noted below, exogeneity plays a key role in devising algorithms to select the conditioning variables.
- In addition, as noted in Slide 412, exogeneity is helpful in making out-of-sample forecasts.
- The method of control functions does not require that  $(U_0, U_1) \perp\!\!\!\perp D \mid (X, Z)$ , which is a central requirement of matching.
- Equivalently, the method of control functions does not require

$$(U_0, U_1) \perp\!\!\!\perp V \mid (X, Z), \quad \text{or that } (U_0, U_1) \perp\!\!\!\perp V \mid X$$

whereas matching does and typically equates  $X$  and  $Z$ .

- However, as noted below, exogeneity plays a key role in devising algorithms to select the conditioning variables.
- In addition, as noted in Slide 412, exogeneity is helpful in making out-of-sample forecasts.
- The method of control functions does not require that  $(U_0, U_1) \perp\!\!\!\perp D \mid (X, Z)$ , which is a central requirement of matching.
- Equivalently, the method of control functions does not require

$$(U_0, U_1) \perp\!\!\!\perp V \mid (X, Z), \quad \text{or that } (U_0, U_1) \perp\!\!\!\perp V \mid X$$

whereas matching does and typically equates  $X$  and  $Z$ .

- Thus matching assumes access to a richer set of conditioning variables than is assumed in the method of control functions.

- The method of control functions allows for outcome unobservables to be dependent on  $D$  even after conditioning on  $(X, Z)$ , and it models this dependence.



- The method of control functions allows for outcome unobservables to be dependent on  $D$  even after conditioning on  $(X, Z)$ , and it models this dependence.
- The method of matching assumes no such  $D$  dependence.

- The method of control functions allows for outcome unobservables to be dependent on  $D$  even after conditioning on  $(X, Z)$ , and it models this dependence.
- The method of matching assumes no such  $D$  dependence.
- Thus in this regard, and maintaining all of the assumptions invoked for control functions in this section, matching is a special case of the method of control functions in which under assumptions (M-1) and (M-2),

$$\begin{aligned}E(U_1|X, D = 1) &= E(U_1|X) \\E(U_0|X, D = 0) &= E(U_0|X).\end{aligned}$$

- In the method of control functions, in the case where  $(X, Z) \perp\!\!\!\perp (U_0, U_1, V)$ , where the  $Z$  can include some or all of the elements of  $X$ , the conditional expectation of  $Y$  given  $X, Z, D$  is

$$\begin{aligned}
 E(Y|X, Z, D) &= E(Y_1|X, Z, D=1)D + E(Y_0|X, Z, D=0)(1-D) & (65) \\
 &= \mu_0(X) + [\mu_1(X) - \mu_0(X)]D \\
 &\quad + E(U_1|P(Z), D=1)D + E(U_0|P(Z), D=0)(1-D) \\
 &= \mu_0(X) + K_0(P(Z)) + [\mu_1(X) - \mu_0(X) + K_1(P(Z)) \\
 &\quad - K_0(P(Z))]D.
 \end{aligned}$$

- In the method of control functions, in the case where  $(X, Z) \perp\!\!\!\perp (U_0, U_1, V)$ , where the  $Z$  can include some or all of the elements of  $X$ , the conditional expectation of  $Y$  given  $X, Z, D$  is

$$\begin{aligned}
 E(Y|X, Z, D) &= E(Y_1|X, Z, D=1) D + E(Y_0|X, Z, D=0) (1 - D) & (65) \\
 &= \mu_0(X) + [\mu_1(X) - \mu_0(X)] D \\
 &\quad + E(U_1|P(Z), D=1) D + E(U_0|P(Z), D=0) (1 - D) \\
 &= \mu_0(X) + K_0(P(Z)) + [\mu_1(X) - \mu_0(X) + K_1(P(Z)) \\
 &\quad - K_0(P(Z))] D.
 \end{aligned}$$

- The coefficient on  $D$  in the final equation combines  $\mu_1(X) - \mu_0(X)$  with  $K_1(P(Z)) - K_0(P(Z))$ .

- In the method of control functions, in the case where  $(X, Z) \perp\!\!\!\perp (U_0, U_1, V)$ , where the  $Z$  can include some or all of the elements of  $X$ , the conditional expectation of  $Y$  given  $X, Z, D$  is

$$\begin{aligned}
 E(Y|X, Z, D) &= E(Y_1|X, Z, D=1) D + E(Y_0|X, Z, D=0) (1 - D) & (65) \\
 &= \mu_0(X) + [\mu_1(X) - \mu_0(X)] D \\
 &\quad + E(U_1|P(Z), D=1) D + E(U_0|P(Z), D=0) (1 - D) \\
 &= \mu_0(X) + K_0(P(Z)) + [\mu_1(X) - \mu_0(X) + K_1(P(Z)) \\
 &\quad - K_0(P(Z))] D.
 \end{aligned}$$

- The coefficient on  $D$  in the final equation combines  $\mu_1(X) - \mu_0(X)$  with  $K_1(P(Z)) - K_0(P(Z))$ .
- It does not correspond to any treatment effect.

- In the method of control functions, in the case where  $(X, Z) \perp\!\!\!\perp (U_0, U_1, V)$ , where the  $Z$  can include some or all of the elements of  $X$ , the conditional expectation of  $Y$  given  $X, Z, D$  is

$$\begin{aligned}
 E(Y|X, Z, D) &= E(Y_1|X, Z, D=1)D + E(Y_0|X, Z, D=0)(1-D) & (65) \\
 &= \mu_0(X) + [\mu_1(X) - \mu_0(X)]D \\
 &\quad + E(U_1|P(Z), D=1)D + E(U_0|P(Z), D=0)(1-D) \\
 &= \mu_0(X) + K_0(P(Z)) + [\mu_1(X) - \mu_0(X) + K_1(P(Z)) \\
 &\quad - K_0(P(Z))]D.
 \end{aligned}$$

- The coefficient on  $D$  in the final equation combines  $\mu_1(X) - \mu_0(X)$  with  $K_1(P(Z)) - K_0(P(Z))$ .
- It does not correspond to any treatment effect.
- To identify  $\mu_1(X) - \mu_0(X)$ , one must isolate it from  $K_1(P(Z)) - K_0(P(Z))$ .

- Under assumptions (M-1) and (M-2) of the method of matching, the conditional expectation of  $Y$  conditional on  $P(X)$  and  $D$  is

$$E(Y|P(X), D) = \mu_0(P(X)) + E(U_0|P(X)) + [(\mu_1(P(X)) - \mu_0(P(X))) + E(U_1|P(X)) - E(U_0|P(X))]D. \quad (66)$$

The coefficient on  $D$  in this expression is now interpretable and is the average treatment effect.

- Under assumptions (M-1) and (M-2) of the method of matching, the conditional expectation of  $Y$  conditional on  $P(X)$  and  $D$  is

$$E(Y|P(X), D) = \mu_0(P(X)) + E(U_0|P(X)) + [(\mu_1(P(X)) - \mu_0(P(X))) + E(U_1|P(X)) - E(U_0|P(X))] D. \quad (66)$$

The coefficient on  $D$  in this expression is now interpretable and is the average treatment effect.

- If we assume that  $(U_0, U_1) \perp\!\!\!\perp X$ , which is not strictly required, we reach a more familiar representation

$$E(Y|P(X), D) = \mu_0(P(X)) + [\mu_1(P(X)) - \mu_0(P(X))] D, \quad (67)$$

since  $E(U_1|P(X)) = E(U_0|P(X)) = 0$ .



- Under assumptions (M-1) and (M-2) of the method of matching, the conditional expectation of  $Y$  conditional on  $P(X)$  and  $D$  is

$$E(Y|P(X), D) = \mu_0(P(X)) + E(U_0|P(X)) + [(\mu_1(P(X)) - \mu_0(P(X))) + E(U_1|P(X)) - E(U_0|P(X))] D. \quad (66)$$

The coefficient on  $D$  in this expression is now interpretable and is the average treatment effect.

- If we assume that  $(U_0, U_1) \perp\!\!\!\perp X$ , which is not strictly required, we reach a more familiar representation

$$E(Y|P(X), D) = \mu_0(P(X)) + [\mu_1(P(X)) - \mu_0(P(X))] D, \quad (67)$$

since  $E(U_1|P(X)) = E(U_0|P(X)) = 0$ .

- A parallel derivation can be made conditioning on  $X$  instead of  $P(X)$ .

- Under the assumptions that justify matching, treatment effects ATE or TT (conditional on  $P(X)$ ) are identified from the coefficient on  $D$  in either (66) or (67).

- Under the assumptions that justify matching, treatment effects ATE or TT (conditional on  $P(X)$ ) are identified from the coefficient on  $D$  in either (66) or (67).
- Condition (M-2) guarantees that  $D$  is not perfectly predictable by  $X$  (or  $P(X)$ ), so the variation in  $D$  identifies the treatment parameter.

- The coefficient on  $D$  in equation (65) for the more general control function model does not correspond to any treatment parameter, whereas the coefficients on  $D$  in equations (66) and (67) correspond to treatment parameters under the assumptions of the matching model.

- The coefficient on  $D$  in equation (65) for the more general control function model does not correspond to any treatment parameter, whereas the coefficients on  $D$  in equations (66) and (67) correspond to treatment parameters under the assumptions of the matching model.
- Under assumption (CF-1),  $\mu_1(P(X)) - \mu_0(P(X)) = \text{ATE}$  and  $\text{ATE} = \text{TT} = \text{MTE}$ , so the method of matching identifies all of the (conditional on  $P(X)$ ) mean treatment parameters.

- Under the assumptions justifying matching, when means of  $Y_1$  and  $Y_0$  are the parameters of interest, and  $X$  satisfies (M-1) and (M-2), the bias terms vanish.

- Under the assumptions justifying matching, when means of  $Y_1$  and  $Y_0$  are the parameters of interest, and  $X$  satisfies (M-1) and (M-2), the bias terms vanish.
- They do not vanish in the more general case considered by the method of control functions.

- Under the assumptions justifying matching, when means of  $Y_1$  and  $Y_0$  are the parameters of interest, and  $X$  satisfies (M-1) and (M-2), the bias terms vanish.
- They do not vanish in the more general case considered by the method of control functions.
- This is the mathematical counterpart of the randomization implicit in matching: conditional on  $X$  or  $P(X)$ ,  $(U_0, U_1)$  are random with respect to  $D$ .



- Under the assumptions justifying matching, when means of  $Y_1$  and  $Y_0$  are the parameters of interest, and  $X$  satisfies (M-1) and (M-2), the bias terms vanish.
- They do not vanish in the more general case considered by the method of control functions.
- This is the mathematical counterpart of the randomization implicit in matching: conditional on  $X$  or  $P(X)$ ,  $(U_0, U_1)$  are random with respect to  $D$ .
- The method of control functions allows these error terms to be nonrandom with respect to  $D$  and models the dependence.

- Under the assumptions justifying matching, when means of  $Y_1$  and  $Y_0$  are the parameters of interest, and  $X$  satisfies (M-1) and (M-2), the bias terms vanish.
- They do not vanish in the more general case considered by the method of control functions.
- This is the mathematical counterpart of the randomization implicit in matching: conditional on  $X$  or  $P(X)$ ,  $(U_0, U_1)$  are random with respect to  $D$ .
- The method of control functions allows these error terms to be nonrandom with respect to  $D$  and models the dependence.
- In the absence of functional form assumptions, it requires an exclusion restriction (a variable in  $Z$  not in  $X$ ) to separate out  $K_0(P(Z))$  from the coefficient on  $D$ .

- Matching produces identification without exclusion restrictions whereas identification with exclusion restrictions is a central feature of the control function method in the absence of functional form assumptions.

- Matching produces identification without exclusion restrictions whereas identification with exclusion restrictions is a central feature of the control function method in the absence of functional form assumptions.
- The fact that the control function approach allows for more general dependencies among the unobservables and the conditioning variables than the matching approach allows is implicitly recognized in the work of ? and ?.

- Matching produces identification without exclusion restrictions whereas identification with exclusion restrictions is a central feature of the control function method in the absence of functional form assumptions.
- The fact that the control function approach allows for more general dependencies among the unobservables and the conditioning variables than the matching approach allows is implicitly recognized in the work of ? and ?.
- Their “sensitivity analyses” for matching when there are unobserved conditioning variables are, in their essence, sensitivity analyses using control functions.

- $\theta$ ,  $\gamma$  and  $\delta$  explicitly model the relationship between matching and selection models using factor structure models, treating the omitted conditioning variables as unobserved factors and estimating their distribution.

- $\theta$ ,  $\gamma$  and  $\delta$  explicitly model the relationship between matching and selection models using factor structure models, treating the omitted conditioning variables as unobserved factors and estimating their distribution.
- Abbring and Heckman discuss this work in Part III.

## Comparing Matching and Classical Control Function Methods for a Generalized Roy Model

- Figure 10, developed in connection with our discussion of instrumental variables, shows the contrast between the shape of the MTE and the OLS matching estimand as a function of  $p$  for the extended Roy model developed in Slide 152.



## Comparing Matching and Classical Control Function Methods for a Generalized Roy Model

- Figure 10, developed in connection with our discussion of instrumental variables, shows the contrast between the shape of the MTE and the OLS matching estimand as a function of  $p$  for the extended Roy model developed in Slide 152.
- The  $MTE(p)$  shows its typical declining shape associated with diminishing returns, and the assumptions justifying matching are violated.

## Comparing Matching and Classical Control Function Methods for a Generalized Roy Model

- Figure 10, developed in connection with our discussion of instrumental variables, shows the contrast between the shape of the MTE and the OLS matching estimand as a function of  $p$  for the extended Roy model developed in Slide 152.
- The  $MTE(p)$  shows its typical declining shape associated with diminishing returns, and the assumptions justifying matching are violated.
- Matching attempts to impose a flat  $MTE(p)$  and therefore flattens the estimated  $MTE(p)$  compared to its true value.

- It understates marginal returns at low levels of  $p$  (associated with unobservables that make it likely to participate in treatment) and overstates marginal returns at high levels of  $p$ .

- It understates marginal returns at low levels of  $p$  (associated with unobservables that make it likely to participate in treatment) and overstates marginal returns at high levels of  $p$ .
- To further illustrate the bias in matching and how the control function eliminates it, we perform sensitivity analyses under different assumptions about the parameters of the underlying selection model.

- In particular, we assume that the data are generated by the model of equations 5 and 6, where  $\mu_D(Z) = Z\gamma$ ,  $\mu_0(X) = \mu_0$ ,  $\mu_1(X) = \mu_1$ , and

$$\begin{aligned}(U_0, U_1, V)' &\sim N(0, \Sigma) \\ \text{corr}(U_j, V) &= \rho_{jV} \\ \text{Var}(U_j) &= \sigma_j^2; \quad j = \{0, 1\}.\end{aligned}$$

We assume in this section that  $D = \mathbf{1}[\mu_D(Z) + V \geq 0]$ , in conformity with the examples presented in ?, from which we draw.

- In particular, we assume that the data are generated by the model of equations 5 and 6, where  $\mu_D(Z) = Z\gamma$ ,  $\mu_0(X) = \mu_0$ ,  $\mu_1(X) = \mu_1$ , and

$$\begin{aligned} (U_0, U_1, V)' &\sim N(0, \Sigma) \\ \text{corr}(U_j, V) &= \rho_{jV} \\ \text{Var}(U_j) &= \sigma_j^2; \quad j = \{0, 1\}. \end{aligned}$$

We assume in this section that  $D = \mathbf{1}[\mu_D(Z) + V \geq 0]$ , in conformity with the examples presented in ?, from which we draw.

- This reformulation of choice model (7) simply entails a change in the sign of  $V$ .

- In particular, we assume that the data are generated by the model of equations 5 and 6, where  $\mu_D(Z) = Z\gamma$ ,  $\mu_0(X) = \mu_0$ ,  $\mu_1(X) = \mu_1$ , and

$$\begin{aligned} (U_0, U_1, V)' &\sim N(0, \Sigma) \\ \text{corr}(U_j, V) &= \rho_{jV} \\ \text{Var}(U_j) &= \sigma_j^2; \quad j = \{0, 1\}. \end{aligned}$$

We assume in this section that  $D = \mathbf{1}[\mu_D(Z) + V \geq 0]$ , in conformity with the examples presented in ?, from which we draw.

- This reformulation of choice model (7) simply entails a change in the sign of  $V$ .
- We assume that  $Z \perp\!\!\!\perp (U_0, U_1, V)$ .

- Using the selection formulae derived in Appendix, Slide 1184, we can write the biases conditional on  $P(Z) = p$  using propensity score matching in a generalized Roy model as

$$\begin{aligned} \text{Bias } TT(Z = z) &= \text{Bias } TT(P(Z) = p) = \sigma_0 \rho_{0V} M(p) \\ \text{Bias } ATE(Z = z) &= \text{Bias } ATE(P(Z) = p) = M(p) [\sigma_1 \rho_{1V} (1 - p) + \sigma_0 \rho_{0V} p], \end{aligned}$$

where  $M(p) = \frac{\phi(\Phi^{-1}(1-p))}{p(1-p)}$ ,  $\phi(\cdot)$  and  $\Phi(\cdot)$  are the pdf and cdf of a standard normal random variable and the propensity score  $P(z)$  is evaluated at  $P(z) = p$ .



- Using the selection formulae derived in Appendix, Slide 1184, we can write the biases conditional on  $P(Z) = p$  using propensity score matching in a generalized Roy model as

$$\begin{aligned} \text{Bias } TT(Z = z) &= \text{Bias } TT(P(Z) = p) = \sigma_0 \rho_{0V} M(p) \\ \text{Bias } ATE(Z = z) &= \text{Bias } ATE(P(Z) = p) = M(p) [\sigma_1 \rho_{1V} (1 - p) + \sigma_0 \rho_{0V} p], \end{aligned}$$

where  $M(p) = \frac{\phi(\Phi^{-1}(1-p))}{p(1-p)}$ ,  $\phi(\cdot)$  and  $\Phi(\cdot)$  are the pdf and cdf of a standard normal random variable and the propensity score  $P(z)$  is evaluated at  $P(z) = p$ .

- We assume that  $\mu_1 = \mu_0$  so that the true average treatment effect is zero.

- We simulate the mean bias for TT (table 10) and ATE (table 11) for different values of the  $\rho_{jV}$  and  $\sigma_j$ .

- We simulate the mean bias for TT (table 10) and ATE (table 11) for different values of the  $\rho_{jV}$  and  $\sigma_j$ .
- The results in the tables show that, as we let the variances of the outcome equations grow, the value of the mean bias that we obtain can become substantial.

- We simulate the mean bias for TT (table 10) and ATE (table 11) for different values of the  $\rho_{jV}$  and  $\sigma_j$ .
- The results in the tables show that, as we let the variances of the outcome equations grow, the value of the mean bias that we obtain can become substantial.
- With larger correlations between the outcomes and the unobservables generating choices, come larger biases.

- We simulate the mean bias for TT (table 10) and ATE (table 11) for different values of the  $\rho_{jV}$  and  $\sigma_j$ .
- The results in the tables show that, as we let the variances of the outcome equations grow, the value of the mean bias that we obtain can become substantial.
- With larger correlations between the outcomes and the unobservables generating choices, come larger biases.
- These tables demonstrate the greater generality of the control function approach, which models the bias rather than assuming it away by conditioning.

Table 10: Mean Bias for Treatment on the Treated

$\rho_{0V}$	Average Bias ( $\sigma_0 = 1$ )	Average Bias ( $\sigma_0 = 2$ )
-1.00	-1.7920	-3.5839
-0.75	-1.3440	-2.6879
-0.50	-0.8960	-1.7920
-0.25	-0.4480	-0.8960
0.00	0.0000	0.0000
0.25	0.4480	0.8960
0.50	0.8960	1.7920
0.75	1.3440	2.6879
1.00	1.7920	3.5839

$$\text{Bias TT} = \rho_{0V} * \sigma_0 * M(p)$$

$$M(p) = \frac{\phi(\Phi^{-1}(1-p))}{p*(1-p)}$$

Source: Heckman and Navarro (2004)

Table 11: Mean Bias for Average Treatment Effect

$(\sigma_0 = 1)$									
$\rho_{0V}$	-1.00	-0.75	-0.50	-0.25	0	0.25	0.50	0.75	1.00
$\rho_{1V}(\sigma_1 = 1)$									
-1.00	-1.7920	-1.5680	-1.3440	-1.1200	-0.8960	-0.6720	-0.4480	-0.2240	0
-0.75	-1.5680	-1.3440	-1.1200	-0.8960	-0.6720	-0.4480	-0.2240	0	0.2240
-0.50	-1.3440	-1.1200	-0.8960	-0.6720	-0.4480	-0.2240	0	0.2240	0.4480
-0.25	-1.1200	-0.8960	-0.6720	-0.4480	-0.2240	0	0.2240	0.4480	0.6720
0	-0.8960	-0.6720	-0.4480	-0.2240	0	0.2240	0.4480	0.6720	0.8960
0.25	-0.6720	-0.4480	-0.2240	0	0.2240	0.4480	0.6720	0.8960	1.1200
0.50	-0.4480	-0.2240	0	0.2240	0.4480	0.6720	0.8960	1.1200	1.3440
0.75	-0.2240	0	0.2240	0.4480	0.6720	0.8960	1.1200	1.3440	1.5680
1.00	0	0.2240	0.4480	0.6720	0.8960	1.1200	1.3440	1.5680	1.7920
$\rho_{1V}(\sigma_1 = 2)$									
-1.00	-2.6879	-2.2399	-1.7920	-1.3440	-0.8960	-0.4480	0	0.4480	0.8960
-0.75	-2.4639	-2.0159	-1.5680	-1.1200	-0.6720	-0.2240	0.2240	0.6720	1.1200
-0.50	-2.2399	-1.7920	-1.3440	-0.8960	-0.4480	0	0.4480	0.8960	1.3440
-0.25	-2.0159	-1.5680	-1.1200	-0.6720	-0.2240	0.2240	0.6720	1.1200	1.5680
0	-1.7920	-1.3440	-0.8960	-0.4480	0	0.4480	0.8960	1.3440	1.7920
0.25	-1.5680	-1.1200	-0.6720	-0.2240	0.2240	0.6720	1.1200	1.5680	2.0159
0.50	-1.3440	-0.8960	-0.4480	0	0.4480	0.8960	1.3440	1.7920	2.2399
0.75	-1.1200	-0.6720	-0.2240	0.2240	0.6720	1.1200	1.5680	2.0159	2.4639
1.00	-0.8960	-0.4480	0	0.4480	0.8960	1.3440	1.7920	2.2399	2.6879

$$\text{Bias ATE} = \rho_{1V} * \sigma_1 * M_1(p) - \rho_{0V} * \sigma_0 * M_0(p)$$

$$M_1(p) = \frac{\phi(\Phi^{-1}(p))}{1-p}$$

$$M_0(p) = \frac{-\phi(\Phi^{-1}(1-p))}{[1-p]}$$

Source: Heckman and Navarro (2004)

- Even if the correlation between the observables and the unobservables ( $\rho_{jV}$ ) is small, so that one might think that selection on unobservables is relatively unimportant, we still obtain substantial biases if we do not control for relevant omitted conditioning variables.



- Even if the correlation between the observables and the unobservables ( $\rho_{jV}$ ) is small, so that one might think that selection on unobservables is relatively unimportant, we still obtain substantial biases if we do not control for relevant omitted conditioning variables.
- Only for special values of the parameters do we avoid bias by matching.

- Even if the correlation between the observables and the unobservables ( $\rho_{jV}$ ) is small, so that one might think that selection on unobservables is relatively unimportant, we still obtain substantial biases if we do not control for relevant omitted conditioning variables.
- Only for special values of the parameters do we avoid bias by matching.
- These examples also demonstrate that sensitivity analyses can be conducted for analysis based on control function methods even when they are not fully identified.

- Even if the correlation between the observables and the unobservables ( $\rho_{jV}$ ) is small, so that one might think that selection on unobservables is relatively unimportant, we still obtain substantial biases if we do not control for relevant omitted conditioning variables.
- Only for special values of the parameters do we avoid bias by matching.
- These examples also demonstrate that sensitivity analyses can be conducted for analysis based on control function methods even when they are not fully identified.
- ? provides an example.

## The Informational Requirements of Matching and the Bias When They Are Not Satisfied

- In this section, we present some examples of when matching “works” and when it breaks down.

## The Informational Requirements of Matching and the Bias When They Are Not Satisfied

- In this section, we present some examples of when matching “works” and when it breaks down.
- This section is based on ?.

## The Informational Requirements of Matching and the Bias When They Are Not Satisfied

- In this section, we present some examples of when matching “works” and when it breaks down.
- This section is based on ?.
- In particular, we show how matching on some of the relevant information but not all can make the bias using matching worse for standard treatment parameters.

## The Informational Requirements of Matching and the Bias When They Are Not Satisfied

- In this section, we present some examples of when matching “works” and when it breaks down.
- This section is based on ?.
- In particular, we show how matching on some of the relevant information but not all can make the bias using matching worse for standard treatment parameters.
- These examples also introduce factor models that play a key role in the analysis of Abbring and Heckman in Part III.

- Slide 12 of this chapter discussed informational asymmetries between the econometrician and the agents whose behavior they are analyzing.



- Slide 12 of this chapter discussed informational asymmetries between the econometrician and the agents whose behavior they are analyzing.
- The method of matching assumes that the econometrician has access to and uses all of the relevant information in the precise sense defined there.

- Slide 12 of this chapter discussed informational asymmetries between the econometrician and the agents whose behavior they are analyzing.
- The method of matching assumes that the econometrician has access to and uses all of the relevant information in the precise sense defined there.
- That means that the  $X$  that guarantees conditional independence (M-1) is available and is used.

- Slide 12 of this chapter discussed informational asymmetries between the econometrician and the agents whose behavior they are analyzing.
- The method of matching assumes that the econometrician has access to and uses all of the relevant information in the precise sense defined there.
- That means that the  $X$  that guarantees conditional independence (M-1) is available and is used.
- The concept of relevant information is a delicate one and it is difficult to find the true conditioning set.

- Assume that the economic model generating the data is a generalized Roy model of the form

$$\begin{aligned}
 D^* &= Z\gamma + V && \text{where} \\
 Z &\perp\!\!\!\perp V && \text{and} \\
 V &= \alpha_{V1}f_1 + \alpha_{V2}f_2 + \varepsilon_V \\
 D &= \begin{cases} 1 & \text{if } D^* \geq 0 \\ 0 & \text{otherwise} \end{cases} ,
 \end{aligned}$$

and

$$\begin{aligned}
 Y_1 &= \mu_1 + U_1 && \text{where } U_1 = \alpha_{11}f_1 + \alpha_{12}f_2 + \varepsilon_1, \\
 Y_0 &= \mu_0 + U_0 && \text{where } U_0 = \alpha_{01}f_1 + \alpha_{02}f_2 + \varepsilon_0.
 \end{aligned}$$

- Assume that the economic model generating the data is a generalized Roy model of the form

$$\begin{aligned}
 D^* &= Z\gamma + V && \text{where} \\
 Z &\perp\!\!\!\perp V && \text{and} \\
 V &= \alpha_{V1}f_1 + \alpha_{V2}f_2 + \varepsilon_V \\
 D &= \begin{cases} 1 & \text{if } D^* \geq 0 \\ 0 & \text{otherwise} \end{cases} ,
 \end{aligned}$$

and

$$\begin{aligned}
 Y_1 &= \mu_1 + U_1 && \text{where } U_1 = \alpha_{11}f_1 + \alpha_{12}f_2 + \varepsilon_1, \\
 Y_0 &= \mu_0 + U_0 && \text{where } U_0 = \alpha_{01}f_1 + \alpha_{02}f_2 + \varepsilon_0.
 \end{aligned}$$

- We remind the reader that contrary to the analysis throughout the rest of this chapter we add  $V$  and do not subtract it in the decision equation.

- This is the familiar representation.

- This is the familiar representation.
- By a change in sign in  $V$ , we can go back and forth between the specification used in this section and the specification used in other sections of the chapter.

- This is the familiar representation.
- By a change in sign in  $V$ , we can go back and forth between the specification used in this section and the specification used in other sections of the chapter.
- In this specification,  $(f_1, f_2, \varepsilon_V, \varepsilon_1, \varepsilon_0)$  are assumed to be mean zero random variables that are mutually independent of each other and  $Z$  so that all the correlation among the elements of  $(U_0, U_1, V)$  is captured by  $f = (f_1, f_2)$ .



- This is the familiar representation.
- By a change in sign in  $V$ , we can go back and forth between the specification used in this section and the specification used in other sections of the chapter.
- In this specification,  $(f_1, f_2, \varepsilon_V, \varepsilon_1, \varepsilon_0)$  are assumed to be mean zero random variables that are mutually independent of each other and  $Z$  so that all the correlation among the elements of  $(U_0, U_1, V)$  is captured by  $f = (f_1, f_2)$ .
- Models that take this form are known as factor models and have been applied in the context of selection models by ?, ??, and ?, among others.

- This is the familiar representation.
- By a change in sign in  $V$ , we can go back and forth between the specification used in this section and the specification used in other sections of the chapter.
- In this specification,  $(f_1, f_2, \varepsilon_V, \varepsilon_1, \varepsilon_0)$  are assumed to be mean zero random variables that are mutually independent of each other and  $Z$  so that all the correlation among the elements of  $(U_0, U_1, V)$  is captured by  $f = (f_1, f_2)$ .
- Models that take this form are known as factor models and have been applied in the context of selection models by ?, ??, and ?, among others.
- We keep implicit any dependence on  $X$  which may be general.

- Generically, the minimal relevant information for this model when the factor loadings are not zero ( $\alpha_{ij} \neq 0$ ) is, for general values of the factor loadings,

$$I_R = \{f_1, f_2\}.$$

Recall that we assume independence between  $Z$  and all error terms.

- Generically, the minimal relevant information for this model when the factor loadings are not zero ( $\alpha_{ij} \neq 0$ ) is, for general values of the factor loadings,

$$I_R = \{f_1, f_2\}.$$

Recall that we assume independence between  $Z$  and all error terms.

- If the econometrician has access to  $I_R$  and uses it, (M-1) is satisfied conditional on  $I_R$ .

- Generically, the minimal relevant information for this model when the factor loadings are not zero ( $\alpha_{ij} \neq 0$ ) is, for general values of the factor loadings,

$$I_R = \{f_1, f_2\}.$$

Recall that we assume independence between  $Z$  and all error terms.

- If the econometrician has access to  $I_R$  and uses it, (M-1) is satisfied conditional on  $I_R$ .
- Note that  $I_R$  plays the role of  $\theta$  in (U-1).

- Generically, the minimal relevant information for this model when the factor loadings are not zero ( $\alpha_{ij} \neq 0$ ) is, for general values of the factor loadings,

$$I_R = \{f_1, f_2\}.$$

Recall that we assume independence between  $Z$  and all error terms.

- If the econometrician has access to  $I_R$  and uses it, (M-1) is satisfied conditional on  $I_R$ .
- Note that  $I_R$  plays the role of  $\theta$  in (U-1).
- In the case where the economist knows  $I_R$ , the economist's information set  $\sigma(I_E)$  contains the relevant information ( $\sigma(I_E) \supseteq \sigma(I_R)$ ).

- The agent's information set may include different variables.

- The agent's information set may include different variables.
- If we assume that  $\varepsilon_0, \varepsilon_1$  are shocks to outcomes not known to the agent at the time treatment decisions are made, but the agent knows all other aspects of the model, the agent's information is

$$I_A = \{f_1, f_2, Z, \varepsilon_V\}.$$

Under perfect certainty, the agent's information set includes  $\varepsilon_1$  and  $\varepsilon_0$ :

$$I_A = \{f_1, f_2, Z, \varepsilon_V, \varepsilon_1, \varepsilon_0\}.$$

In either case, all of the information available to the agent is not required to satisfy conditional independence (M-1).



- The agent's information set may include different variables.
- If we assume that  $\varepsilon_0, \varepsilon_1$  are shocks to outcomes not known to the agent at the time treatment decisions are made, but the agent knows all other aspects of the model, the agent's information is

$$I_A = \{f_1, f_2, Z, \varepsilon_V\}.$$

Under perfect certainty, the agent's information set includes  $\varepsilon_1$  and  $\varepsilon_0$ :

$$I_A = \{f_1, f_2, Z, \varepsilon_V, \varepsilon_1, \varepsilon_0\}.$$

In either case, all of the information available to the agent is not required to satisfy conditional independence (M-1).

- All three information sets guarantee conditional independence, but only the first is minimal relevant.

- In the notation of Slide 12, the observing economist may know some variables not in  $I_A$ ,  $I_{R^*}$  or  $I_R$  but may not know all of the variables in  $I_R$ .

- In the notation of Slide 12, the observing economist may know some variables not in  $I_A$ ,  $I_{R^*}$  or  $I_R$  but may not know all of the variables in  $I_R$ .
- In the following subsections, we study what happens when the matching assumption that  $\sigma(I_E) \supseteq \sigma(I_R)$  does not hold.

- In the notation of Slide 12, the observing economist may know some variables not in  $I_A$ ,  $I_{R^*}$  or  $I_R$  but may not know all of the variables in  $I_R$ .
- In the following subsections, we study what happens when the matching assumption that  $\sigma(I_E) \supseteq \sigma(I_R)$  does not hold.
- That is, we analyze what happens to the bias from matching as the amount of information used by the econometrician is changed.

- In the notation of Slide 12, the observing economist may know some variables not in  $I_A$ ,  $I_{R^*}$  or  $I_R$  but may not know all of the variables in  $I_R$ .
- In the following subsections, we study what happens when the matching assumption that  $\sigma(I_E) \supseteq \sigma(I_R)$  does not hold.
- That is, we analyze what happens to the bias from matching as the amount of information used by the econometrician is changed.
- In order to get closed form expressions for the biases of the treatment parameters, we make the additional assumption that

$$(f_1, f_2, \varepsilon_V, \varepsilon_1, \varepsilon_0) \sim N(0, \Sigma),$$

where  $\Sigma$  is a matrix with  $(\sigma_{f_1}^2, \sigma_{f_2}^2, \sigma_{\varepsilon_V}^2, \sigma_{\varepsilon_1}^2, \sigma_{\varepsilon_0}^2)$  on the diagonal and zero in all the non-diagonal elements.

- This assumption links matching models to conventional normal selection models of the sort developed in Part I and further analyzed in Slide 12 of this Part.

- This assumption links matching models to conventional normal selection models of the sort developed in Part I and further analyzed in Slide 12 of this Part.
- However, the examples based on this specification illustrate more general principles.

- This assumption links matching models to conventional normal selection models of the sort developed in Part I and further analyzed in Slide 12 of this Part.
- However, the examples based on this specification illustrate more general principles.
- We now analyze various commonly encountered cases.



## The Economist Uses the Minimal Relevant Information:

$$\sigma(I_R) \subseteq \sigma(I_E)$$

- We begin by analyzing the case in which the information used by the economist is  $I_E = \{Z, f_1, f_2\}$ , so that the econometrician has access to a relevant information set and it is larger than the minimal relevant information set.

## The Economist Uses the Minimal Relevant Information:

$$\sigma(I_R) \subseteq \sigma(I_E)$$

- We begin by analyzing the case in which the information used by the economist is  $I_E = \{Z, f_1, f_2\}$ , so that the econometrician has access to a relevant information set and it is larger than the minimal relevant information set.
- In this case, it is straightforward to show that matching identifies all of the mean treatment parameters with no bias.

## The Economist Uses the Minimal Relevant Information:

$$\sigma(I_R) \subseteq \sigma(I_E)$$

- We begin by analyzing the case in which the information used by the economist is  $I_E = \{Z, f_1, f_2\}$ , so that the econometrician has access to a relevant information set and it is larger than the minimal relevant information set.
- In this case, it is straightforward to show that matching identifies all of the mean treatment parameters with no bias.
- The matching estimator has population mean

$$E(Y_1|D=1, I_E) - E(Y_0|D=0, I_E) = \mu_1 - \mu_0 + (\alpha_{11} - \alpha_{01}) f_1 + (\alpha_{12} - \alpha_{02}) f_2,$$

and all of the mean treatment parameters collapse to this same expression since, conditional on knowing  $f_1$  and  $f_2$ , there is no selection because  $(\varepsilon_0, \varepsilon_1) \perp\!\!\!\perp V$ .

- Recall that, for arbitrary choices of  $\alpha_{11}, \alpha_{01}, \alpha_{12}, \alpha_{02}$ ,  $I_R = \{f_1, f_2\}$  and the economist needs less information to achieve (M-1) than is contained in  $I_E$ .

- Recall that, for arbitrary choices of  $\alpha_{11}, \alpha_{01}, \alpha_{12}, \alpha_{02}$ ,  $I_R = \{f_1, f_2\}$  and the economist needs less information to achieve (M-1) than is contained in  $I_E$ .
- In this case, the analysis of ? tells us that knowledge of  $(Z, f_1, f_2)$  and knowledge of  $P(Z, f_1, f_2)$  are equally useful in identifying all of the treatment parameters conditional on  $P$ .

- Recall that, for arbitrary choices of  $\alpha_{11}, \alpha_{01}, \alpha_{12}, \alpha_{02}$ ,  $I_R = \{f_1, f_2\}$  and the economist needs less information to achieve (M-1) than is contained in  $I_E$ .
- In this case, the analysis of ? tells us that knowledge of  $(Z, f_1, f_2)$  and knowledge of  $P(Z, f_1, f_2)$  are equally useful in identifying all of the treatment parameters conditional on  $P$ .
- If we write the propensity score as

$$P(I_E) = \Pr \left( \frac{\varepsilon_V}{\sigma_{\varepsilon_V}} > \frac{-Z\gamma - \alpha_{V1}f_1 - \alpha_{V2}f_2}{\sigma_{\varepsilon_V}} \right) = 1 - \Phi \left( \frac{-Z\gamma - \alpha_{V1}f_1 - \alpha_{V2}f_2}{\sigma_{\varepsilon_V}} \right) = p,$$

the event  $\left( D^* \begin{matrix} \leq \\ \geq \end{matrix} 0, \text{ given } f = \tilde{f} \text{ and } Z = z \right)$  can be written as  $\frac{\varepsilon_V}{\sigma_{\varepsilon_V}} \begin{matrix} \leq \\ \geq \end{matrix} \Phi^{-1} \left( 1 - P(z, \tilde{f}) \right)$ , where  $\Phi$  is the cdf of a standard normal random variable and  $f = (f_1, f_2)$ .

- We abuse notation slightly by using  $z$  as the realized fixed value of  $Z$  and  $\tilde{f}$  as the realized value of  $f$ .

- We abuse notation slightly by using  $z$  as the realized fixed value of  $Z$  and  $\tilde{f}$  as the realized value of  $f$ .
- The population matching condition (M-1) implies that

$$\begin{aligned} & E\left(Y_1|D=1, P(I_E) = P(z, \tilde{f})\right) - E\left(Y_0|D=0, P(I_E) = P(z, \tilde{f})\right) \\ &= \mu_1 - \mu_0 + E\left(U_1|D=1, P(I_E) = P(z, \tilde{f})\right) - E\left(U_0|D=0, P(I_E) = P(z, \tilde{f})\right) \\ &= \mu_1 - \mu_0 + E\left(U_1 \mid \frac{\varepsilon_V}{\sigma_{\varepsilon_V}} > \Phi^{-1}\left(1 - P(z, \tilde{f})\right)\right) - E\left(U_0 \mid \frac{\varepsilon_V}{\sigma_{\varepsilon_V}} \leq \Phi^{-1}\left(1 - P(z, \tilde{f})\right)\right) \\ &= \mu_1 - \mu_0. \end{aligned}$$



- This expression is equal to all of the treatment parameters discussed in this chapter, since

$$E\left(U_1 \mid \frac{\varepsilon_V}{\sigma_{\varepsilon_V}} > \Phi^{-1}\left(1 - P(z, \tilde{f})\right)\right) = \frac{\text{Cov}(U_1, \varepsilon_V)}{\sigma_{\varepsilon_V}} M_1\left(P(z, \tilde{f})\right)$$

and

$$E\left(U_0 \mid \frac{\varepsilon_V}{\sigma_{\varepsilon_V}} \leq \Phi^{-1}\left(1 - P(z, \tilde{f})\right)\right) = \frac{\text{Cov}(U_0, \varepsilon_V)}{\sigma_{\varepsilon_V}} M_0\left(P(z, \tilde{f})\right),$$

where

$$M_1(P(z, \tilde{f})) = \frac{\phi(\Phi^{-1}(1 - P(z, \tilde{f})))}{P(z, \tilde{f})}$$

$$M_0(P(z, \tilde{f})) = -\frac{\phi(\Phi^{-1}(1 - P(z, \tilde{f})))}{1 - P(z, \tilde{f})},$$

where  $\phi$  is the density of a standard normal random variable.

- As a consequence of the assumptions about mutual independence of the errors

$$\text{Cov}(U_i, \varepsilon_V) = \text{Cov}(\alpha_{i1}f_1 + \alpha_{i2}f_2 + \varepsilon_i, \varepsilon_V) = 0, \quad i = 0, 1.$$

- As a consequence of the assumptions about mutual independence of the errors

$$\text{Cov}(U_i, \varepsilon_V) = \text{Cov}(\alpha_{i1}f_1 + \alpha_{i2}f_2 + \varepsilon_i, \varepsilon_V) = 0, \quad i = 0, 1.$$

- In the context of the generalized Roy model, the case considered in this subsection is the one matching is designed to solve.

- As a consequence of the assumptions about mutual independence of the errors

$$\text{Cov}(U_i, \varepsilon_V) = \text{Cov}(\alpha_{i1}f_1 + \alpha_{i2}f_2 + \varepsilon_i, \varepsilon_V) = 0, \quad i = 0, 1.$$

- In the context of the generalized Roy model, the case considered in this subsection is the one matching is designed to solve.
- Even though a selection model generates the data, the fact that the information used by the econometrician includes the minimal relevant information makes matching a correct solution to the selection problem.

- As a consequence of the assumptions about mutual independence of the errors

$$\text{Cov}(U_i, \varepsilon_V) = \text{Cov}(\alpha_{i1}f_1 + \alpha_{i2}f_2 + \varepsilon_i, \varepsilon_V) = 0, \quad i = 0, 1.$$

- In the context of the generalized Roy model, the case considered in this subsection is the one matching is designed to solve.
- Even though a selection model generates the data, the fact that the information used by the econometrician includes the minimal relevant information makes matching a correct solution to the selection problem.
- We can estimate the treatment parameters with no bias since, as a consequence of our assumptions,  $(U_0, U_1) \perp\!\!\!\perp D \mid (f, Z)$ , which is exactly what matching requires.

- As a consequence of the assumptions about mutual independence of the errors

$$\text{Cov}(U_i, \varepsilon_V) = \text{Cov}(\alpha_{i1}f_1 + \alpha_{i2}f_2 + \varepsilon_i, \varepsilon_V) = 0, \quad i = 0, 1.$$

- In the context of the generalized Roy model, the case considered in this subsection is the one matching is designed to solve.
- Even though a selection model generates the data, the fact that the information used by the econometrician includes the minimal relevant information makes matching a correct solution to the selection problem.
- We can estimate the treatment parameters with no bias since, as a consequence of our assumptions,  $(U_0, U_1) \perp\!\!\!\perp D \mid (f, Z)$ , which is exactly what matching requires.
- The minimal relevant information set is even smaller.

- For arbitrary factor loadings, we only need to know  $(f_1, f_2)$  to secure conditional independence.

- For arbitrary factor loadings, we only need to know  $(f_1, f_2)$  to secure conditional independence.
- We can define the propensity score solely in terms of  $f_1$  and  $f_2$ , and the Rosenbaum-Rubin result still goes through.



- For arbitrary factor loadings, we only need to know  $(f_1, f_2)$  to secure conditional independence.
- We can define the propensity score solely in terms of  $f_1$  and  $f_2$ , and the Rosenbaum-Rubin result still goes through.
- Our analysis in this section focuses on treatment parameters conditional on particular values of  $P(Z, f) = P(z, \tilde{f})$ , i.e., for fixed values of  $p$ , but we could condition more finely.

- For arbitrary factor loadings, we only need to know  $(f_1, f_2)$  to secure conditional independence.
- We can define the propensity score solely in terms of  $f_1$  and  $f_2$ , and the Rosenbaum-Rubin result still goes through.
- Our analysis in this section focuses on treatment parameters conditional on particular values of  $P(Z, f) = P(z, \tilde{f})$ , i.e., for fixed values of  $p$ , but we could condition more finely.
- Conditioning on  $P(z, \tilde{f})$  defines the treatment parameters more coarsely.

- For arbitrary factor loadings, we only need to know  $(f_1, f_2)$  to secure conditional independence.
- We can define the propensity score solely in terms of  $f_1$  and  $f_2$ , and the Rosenbaum-Rubin result still goes through.
- Our analysis in this section focuses on treatment parameters conditional on particular values of  $P(Z, f) = P(z, \tilde{f})$ , i.e., for fixed values of  $p$ , but we could condition more finely.
- Conditioning on  $P(z, \tilde{f})$  defines the treatment parameters more coarsely.
- We can use either fine or coarse conditioning to construct the unconditional treatment effects.

- In this example, using more information than what is in the relevant information set (i.e., using  $Z$ ) is harmless.

- In this example, using more information than what is in the relevant information set (i.e., using  $Z$ ) is harmless.
- But this is not generally true.

- In this example, using more information than what is in the relevant information set (i.e., using  $Z$ ) is harmless.
- But this is not generally true.
- If  $Z \not\perp\!\!\!\perp (U_0, U_1, V)$ , adding  $Z$  to the conditioning set can violate conditional independence assumption (M-1):

$$(Y_0, Y_1) \perp\!\!\!\perp D \mid (f_1, f_2),$$

but

$$(Y_0, Y_1) \not\perp\!\!\!\perp D \mid (f_1, f_2, Z).$$

Adding extra variables can destroy the crucial conditional independence property of matching.

- In this example, using more information than what is in the relevant information set (i.e., using  $Z$ ) is harmless.
- But this is not generally true.
- If  $Z \not\perp\!\!\!\perp (U_0, U_1, V)$ , adding  $Z$  to the conditioning set can violate conditional independence assumption (M-1):

$$(Y_0, Y_1) \perp\!\!\!\perp D \mid (f_1, f_2),$$

but

$$(Y_0, Y_1) \not\perp\!\!\!\perp D \mid (f_1, f_2, Z).$$

Adding extra variables can destroy the crucial conditional independence property of matching.

- We present an example of this point below.

- We first consider a case where  $Z \perp\!\!\!\perp (U_0, U_1, V)$  but the analyst conditions on  $Z$  and not  $(f_1, f_2)$ .



- We first consider a case where  $Z \perp\!\!\!\perp (U_0, U_1, V)$  but the analyst conditions on  $Z$  and not  $(f_1, f_2)$ .
- In this case, there is selection on the unobservables that are not conditioned on.

## The Economist does not Use All of the Minimal Relevant Information

- Next, suppose that the information used by the econometrician is

$$I_E = \{Z\},$$

and there is selection on the unobservable (to the analyst)  $f_1$  and  $f_2$ , i.e., the factor loadings  $\alpha_{ij}$  are all non zero.

## The Economist does not Use All of the Minimal Relevant Information

- Next, suppose that the information used by the econometrician is

$$I_E = \{Z\},$$

and there is selection on the unobservable (to the analyst)  $f_1$  and  $f_2$ , i.e., the factor loadings  $\alpha_{ij}$  are all non zero.

- Recall that we assume that  $Z$  and the  $f$  are independent.

## The Economist does not Use All of the Minimal Relevant Information

- Next, suppose that the information used by the econometrician is

$$I_E = \{Z\},$$

and there is selection on the unobservable (to the analyst)  $f_1$  and  $f_2$ , i.e., the factor loadings  $\alpha_{ij}$  are all non zero.

- Recall that we assume that  $Z$  and the  $f$  are independent.
- In this case, the event  $\left(D^* \begin{matrix} \leq \\ > \end{matrix} 0, Z = z\right)$  is characterized by

$$\frac{\alpha_{V1}f_1 + \alpha_{V2}f_2 + \varepsilon_V}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \alpha_{V2}^2\sigma_{f_2}^2 + \sigma_{\varepsilon_V}^2}} \begin{matrix} \leq \\ > \end{matrix} \Phi^{-1}(1 - P(z)).$$

- Using the analysis presented in Appendix, Slide 1184, the bias for the different treatment parameters is given by

$$\text{Bias } TT(Z = z) = \text{Bias } TT(P(Z) = P(z)) = \eta_0 M(P(z)), \quad (68)$$

where  $M(P(z)) = M_1(P(z)) - M_0(P(z))$ .

$$\text{Bias } ATE(Z = z) = \text{Bias } ATE(P(Z) = P(z)) = M(P(z))\{\eta_1[1 - P(z)] + \eta_0 P(z)\}, \quad (69)$$

where

$$\eta_1 = \frac{\alpha_{V1}\alpha_{11}\sigma_{f_1}^2 + \alpha_{V2}\alpha_{12}\sigma_{f_2}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \alpha_{V2}^2\sigma_{f_2}^2 + \sigma_{\varepsilon V}^2}}$$

$$\eta_0 = \frac{\alpha_{V1}\alpha_{01}\sigma_{f_1}^2 + \alpha_{V2}\alpha_{02}\sigma_{f_2}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \alpha_{V2}^2\sigma_{f_2}^2 + \sigma_{\varepsilon V}^2}}.$$

- It is not surprising that matching on sets of variables that exclude the relevant conditioning variables produces bias for the conditional (on  $P(z)$ ) treatment parameters.

- It is not surprising that matching on sets of variables that exclude the relevant conditioning variables produces bias for the conditional (on  $P(z)$ ) treatment parameters.
- The advantage of working with a closed form expression for the bias is that it allows us to answer questions about the *magnitude* of this bias under different assumptions about the information available to the analyst, and to present some simple examples.

- It is not surprising that matching on sets of variables that exclude the relevant conditioning variables produces bias for the conditional (on  $P(z)$ ) treatment parameters.
- The advantage of working with a closed form expression for the bias is that it allows us to answer questions about the *magnitude* of this bias under different assumptions about the information available to the analyst, and to present some simple examples.
- We next use expressions (68) and (69) as benchmarks against which to compare the relative size of the bias when we enlarge the econometrician's information set beyond  $Z$ .



Adding Information to the Econometrician's Information Set  $I_E$ :  
Using Some but not All the Information from the Minimal Relevant  
Information Set  $I_R$

- Suppose that the econometrician uses more information but not all of the information in the minimal relevant information set.

Adding Information to the Econometrician's Information Set  $I_E$ :  
Using Some but not All the Information from the Minimal Relevant  
Information Set  $I_R$

- Suppose that the econometrician uses more information but not all of the information in the minimal relevant information set.
- He still reports values of the parameters conditional on specific  $p$  values but now the model for  $p$  has different conditioning variables.

## Adding Information to the Econometrician's Information Set $I_E$ : Using Some but not All the Information from the Minimal Relevant Information Set $I_R$

- Suppose that the econometrician uses more information but not all of the information in the minimal relevant information set.
- He still reports values of the parameters conditional on specific  $p$  values but now the model for  $p$  has different conditioning variables.
- For example, the data set assumed in the preceding section might be augmented or else the econometrician decides to use information previously available.

- In particular, assume that the econometrician's information set is

$$I'_E = \{Z, f_2\},$$

and that he uses this information set.

- In particular, assume that the econometrician's information set is

$$I'_E = \{Z, f_2\},$$

and that he uses this information set.

- Under conditions 1 and 2 presented below, the biases for the treatment parameters conditional on values of  $P = p$  are reduced in absolute value relative to their values in Slide 773 by changing the conditioning set in this way.

- In particular, assume that the econometrician's information set is

$$I'_E = \{Z, f_2\},$$

and that he uses this information set.

- Under conditions 1 and 2 presented below, the biases for the treatment parameters conditional on values of  $P = p$  are reduced in absolute value relative to their values in Slide 773 by changing the conditioning set in this way.
- But these conditions are not generally satisfied, so that adding extra information does not necessarily reduce bias and may actually increase it.

- To show how this happens in our model, we define expressions comparable to  $\eta_1$  and  $\eta_0$  for this case:

$$\eta'_1 = \frac{\alpha_{V1}\alpha_{11}\sigma_{f_1}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \sigma_{\varepsilon_V}^2}}$$
$$\eta'_0 = \frac{\alpha_{V1}\alpha_{01}\sigma_{f_1}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \sigma_{\varepsilon_V}^2}}.$$

- To show how this happens in our model, we define expressions comparable to  $\eta_1$  and  $\eta_0$  for this case:

$$\eta'_1 = \frac{\alpha_{V1}\alpha_{11}\sigma_{f_1}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \sigma_{\varepsilon_V}^2}}$$

$$\eta'_0 = \frac{\alpha_{V1}\alpha_{01}\sigma_{f_1}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \sigma_{\varepsilon_V}^2}}.$$

- We compare the biases under the two cases using formulae (68)–(69), suitably modified, but keeping  $p$  fixed at a specific value even though this implies different conditioning sets in terms of  $(z, \tilde{f})$ .



**Condition 1** *The bias produced by using matching to estimate  $\tau$  is smaller in absolute value for any given  $p$  when the new information set  $\sigma(I_E)$  is used if*

$$|\eta_0| > |\eta'_0|.$$

There is a similar result for ATE:

**Condition 2** *The bias produced by using matching to estimate ATE is smaller in absolute value for any given  $p$  when the new information set  $\sigma(I_E)$  is used if*

$$|\eta_1(1-p) + \eta_0 p| > |\eta'_1(1-p) + \eta'_0 p|.$$

## Proof.

These conditions are a direct consequence of formulae (68) and (69), modified to allow for the different covariance structure produced by the information structure assumed in this section (replacing  $\eta_0$  with  $\eta'_0$ ,  $\eta_1$  with  $\eta'_1$ ). □

- It is important to notice that we condition on the same value of  $p$  in deriving these expressions although the variables in  $P$  are different across different specifications of the model.

- It is important to notice that we condition on the same value of  $p$  in deriving these expressions although the variables in  $P$  are different across different specifications of the model.
- Propensity-score matching defines them conditional on  $P = p$ , so we are being faithful to that method.

- These conditions do not always hold.

- These conditions do not always hold.
- In general, whether or not the bias will be reduced by adding additional conditioning variables depends on the relative importance of the additional information in both the outcome equations and on the signs of the terms inside the absolute value.

- Consider whether Condition (1) is satisfied in general.



- Consider whether Condition (1) is satisfied in general.
- Assume  $\eta_0 > 0$  for all  $\alpha_{02}, \alpha_{12}$ .

- Consider whether Condition (1) is satisfied in general.
- Assume  $\eta_0 > 0$  for all  $\alpha_{02}, \alpha_{V2}$ .
- Then  $\eta_0 > \eta'_0$  if

$$\eta_0 = \frac{\alpha_{V1}\alpha_{01}\sigma_{f_1}^2 + (\alpha_{V2}^2) \left(\frac{\alpha_{02}}{\alpha_{V2}}\right) \sigma_{f_2}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \alpha_{V2}^2\sigma_{f_2}^2 + \sigma_{\varepsilon_V}^2}} > \frac{\alpha_{V1}\alpha_{11}\sigma_{f_1}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \sigma_{\varepsilon_V}^2}} = \eta'_0.$$

- Consider whether Condition (1) is satisfied in general.
- Assume  $\eta_0 > 0$  for all  $\alpha_{02}, \alpha_{V2}$ .
- Then  $\eta_0 > \eta'_0$  if

$$\eta_0 = \frac{\alpha_{V1}\alpha_{01}\sigma_{f_1}^2 + (\alpha_{V2}^2) \left( \frac{\alpha_{02}}{\alpha_{V2}} \right) \sigma_{f_2}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \alpha_{V2}^2\sigma_{f_2}^2 + \sigma_{\varepsilon V}^2}} > \frac{\alpha_{V1}\alpha_{11}\sigma_{f_1}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \sigma_{\varepsilon V}^2}} = \eta'_0.$$

- When  $\frac{\alpha_{02}}{\alpha_{V2}} = 0$ , clearly  $\eta_0 < \eta'_0$ .

- Consider whether Condition (1) is satisfied in general.
- Assume  $\eta_0 > 0$  for all  $\alpha_{02}, \alpha_{V2}$ .
- Then  $\eta_0 > \eta'_0$  if

$$\eta_0 = \frac{\alpha_{V1}\alpha_{01}\sigma_{f_1}^2 + (\alpha_{V2}^2) \left(\frac{\alpha_{02}}{\alpha_{V2}}\right) \sigma_{f_2}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \alpha_{V2}^2\sigma_{f_2}^2 + \sigma_{\varepsilon_V}^2}} > \frac{\alpha_{V1}\alpha_{11}\sigma_{f_1}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \sigma_{\varepsilon_V}^2}} = \eta'_0.$$

- When  $\frac{\alpha_{02}}{\alpha_{V2}} = 0$ , clearly  $\eta_0 < \eta'_0$ .
- Adding information to the conditioning set increases bias.

- Consider whether Condition (1) is satisfied in general.
- Assume  $\eta_0 > 0$  for all  $\alpha_{02}, \alpha_{V2}$ .
- Then  $\eta_0 > \eta'_0$  if

$$\eta_0 = \frac{\alpha_{V1}\alpha_{01}\sigma_{f_1}^2 + (\alpha_{V2}^2) \left(\frac{\alpha_{02}}{\alpha_{V2}}\right) \sigma_{f_2}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \alpha_{V2}^2\sigma_{f_2}^2 + \sigma_{\varepsilon_V}^2}} > \frac{\alpha_{V1}\alpha_{11}\sigma_{f_1}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \sigma_{\varepsilon_V}^2}} = \eta'_0.$$

- When  $\frac{\alpha_{02}}{\alpha_{V2}} = 0$ , clearly  $\eta_0 < \eta'_0$ .
- Adding information to the conditioning set increases bias.
- We can vary  $\left(\frac{\alpha_{02}}{\alpha_{V2}}\right)$  holding all of the other parameters constant and hence can make the left hand side arbitrarily large.

- Consider whether Condition (1) is satisfied in general.
- Assume  $\eta_0 > 0$  for all  $\alpha_{02}, \alpha_{V2}$ .
- Then  $\eta_0 > \eta'_0$  if

$$\eta_0 = \frac{\alpha_{V1}\alpha_{01}\sigma_{f_1}^2 + (\alpha_{V2}^2) \left(\frac{\alpha_{02}}{\alpha_{V2}}\right) \sigma_{f_2}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \alpha_{V2}^2\sigma_{f_2}^2 + \sigma_{\varepsilon_V}^2}} > \frac{\alpha_{V1}\alpha_{11}\sigma_{f_1}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \sigma_{\varepsilon_V}^2}} = \eta'_0.$$

- When  $\frac{\alpha_{02}}{\alpha_{V2}} = 0$ , clearly  $\eta_0 < \eta'_0$ .
- Adding information to the conditioning set increases bias.
- We can vary  $\left(\frac{\alpha_{02}}{\alpha_{V2}}\right)$  holding all of the other parameters constant and hence can make the left hand side arbitrarily large.
- As  $\alpha_{02}$  increases, there is some critical value  $\alpha_{02}^*$  beyond which  $\eta_0 > \eta'_0$ .

- If we assumed that  $\eta_0 < 0$ , however, the opposite conclusion would hold, and the conditions for reduction in bias would be harder to meet, as the relative importance of the new information is increased.

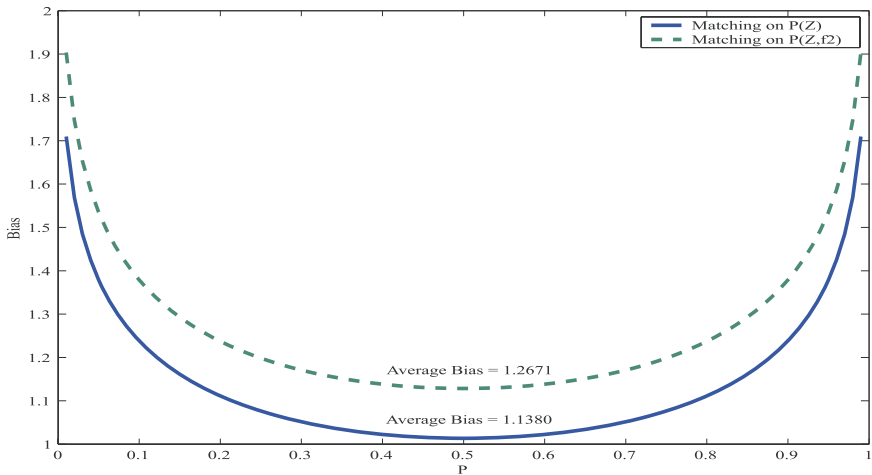
- If we assumed that  $\eta_0 < 0$ , however, the opposite conclusion would hold, and the conditions for reduction in bias would be harder to meet, as the relative importance of the new information is increased.
- Similar expressions can be derived for ATE and MTE, in which the direction of the effect depends on the signs of the terms in the absolute value.



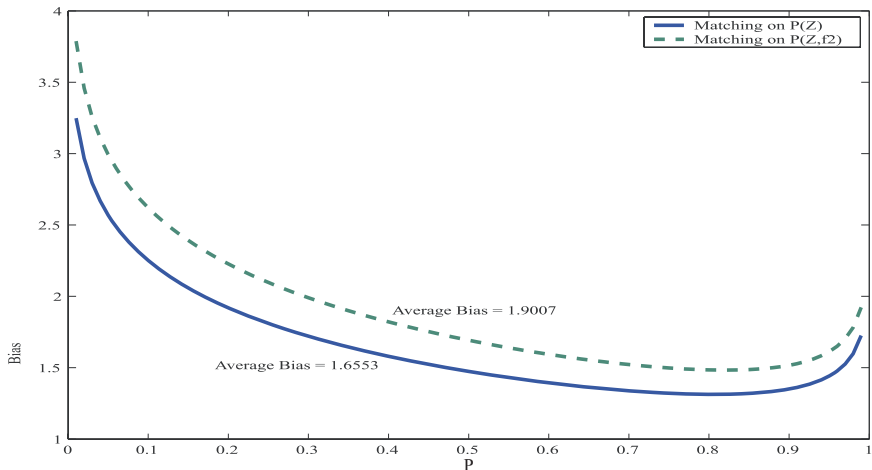
- If we assumed that  $\eta_0 < 0$ , however, the opposite conclusion would hold, and the conditions for reduction in bias would be harder to meet, as the relative importance of the new information is increased.
- Similar expressions can be derived for ATE and MTE, in which the direction of the effect depends on the signs of the terms in the absolute value.
- Figures 23A and 23B illustrate the point that adding some but not all information from the minimal relevant set might increase the point-wise bias and the unconditional or average bias for ATE and TT, respectively.

- If we assumed that  $\eta_0 < 0$ , however, the opposite conclusion would hold, and the conditions for reduction in bias would be harder to meet, as the relative importance of the new information is increased.
- Similar expressions can be derived for ATE and MTE, in which the direction of the effect depends on the signs of the terms in the absolute value.
- Figures 23A and 23B illustrate the point that adding some but not all information from the minimal relevant set might increase the point-wise bias and the unconditional or average bias for ATE and TT, respectively.
- Values of the parameters of the model are presented at the base of the figures.

Figure 23: A. Bias for Treatment on the Treated



## B. Bias for Average Treatment Effect



Note: Using proxy  $\tilde{Z}$  for  $f_2$  increases the bias. Correlation  $(\tilde{Z}, f_2) = 0.5$ .

Model:

$$\begin{aligned} V &= Z + f_1 + f_2 + \varepsilon_V; & Y_1 &= 2f_1 + 0.1f_2 + \varepsilon_1; & Y_0 &= f_1 + 0.1f_2 + \varepsilon_0 \\ \varepsilon_V &\sim N(0, 1); & \varepsilon_1 &\sim N(0, 1); & \varepsilon_0 &\sim N(0, 1) \\ f_1 &\sim N(0, 1); & f_2 &\sim N(0, 1) \end{aligned}$$

Source: Heckman and Navarro (2005)

- In these figures, we compare conditioning on  $P(z)$ , which in general is not guaranteed to eliminate bias, with conditioning on  $P(z)$  and  $f_2$  but not  $f_1$ .

- In these figures, we compare conditioning on  $P(z)$ , which in general is not guaranteed to eliminate bias, with conditioning on  $P(z)$  and  $f_2$  but not  $f_1$ .
- Adding  $f_2$  to the conditioning increases bias.

- In these figures, we compare conditioning on  $P(z)$ , which in general is not guaranteed to eliminate bias, with conditioning on  $P(z)$  and  $f_2$  but not  $f_1$ .
- Adding  $f_2$  to the conditioning increases bias.
- The fact that the point-wise (and overall) bias might increase when adding some but not all information from  $I_R$  is a feature that is not shared by the method of control functions.



- In these figures, we compare conditioning on  $P(z)$ , which in general is not guaranteed to eliminate bias, with conditioning on  $P(z)$  and  $f_2$  but not  $f_1$ .
- Adding  $f_2$  to the conditioning increases bias.
- The fact that the point-wise (and overall) bias might increase when adding some but not all information from  $I_R$  is a feature that is not shared by the method of control functions.
- Because the method of control functions models the stochastic dependence of the unobservables in the outcome equations on the observables, changing the variables observed by the econometrician to include  $f_2$  does not generate bias.

- In these figures, we compare conditioning on  $P(z)$ , which in general is not guaranteed to eliminate bias, with conditioning on  $P(z)$  and  $f_2$  but not  $f_1$ .
- Adding  $f_2$  to the conditioning increases bias.
- The fact that the point-wise (and overall) bias might increase when adding some but not all information from  $I_R$  is a feature that is not shared by the method of control functions.
- Because the method of control functions models the stochastic dependence of the unobservables in the outcome equations on the observables, changing the variables observed by the econometrician to include  $f_2$  does not generate bias.
- It only changes the control function used.

- That is, by adding  $f_2$  we change the control function from

$$K_1(P(Z) = P(z)) = \eta_1 M_1(P(z))$$

$$K_0(P(Z) = P(z)) = \eta_0 M_0(P(z))$$

to

$$K'_1(P(Z, f_2) = P(z, \tilde{f}_2)) = \eta'_1 M_1(P(z, \tilde{f}_2))$$

$$K'_0(P(Z, f_2) = P(z, \tilde{f}_2)) = \eta'_0 M_0(P(z, \tilde{f}_2))$$

but do not generate any bias in using the control function estimator.

- That is, by adding  $f_2$  we change the control function from

$$K_1(P(Z) = P(z)) = \eta_1 M_1(P(z))$$

$$K_0(P(Z) = P(z)) = \eta_0 M_0(P(z))$$

to

$$K'_1(P(Z, f_2) = P(z, \tilde{f}_2)) = \eta'_1 M_1(P(z, \tilde{f}_2))$$

$$K'_0(P(Z, f_2) = P(z, \tilde{f}_2)) = \eta'_0 M_0(P(z, \tilde{f}_2))$$

but do not generate any bias in using the control function estimator.

- This is a major advantage of this method.

- It controls for the bias of the omitted conditioning variables by modeling it.

- It controls for the bias of the omitted conditioning variables by modeling it.
- Of course, if the model for the bias term is not valid, neither is the correction for the bias.

- It controls for the bias of the omitted conditioning variables by modeling it.
- Of course, if the model for the bias term is not valid, neither is the correction for the bias.
- Semiparametric selection estimators are designed to protect the analyst against model misspecification.

- It controls for the bias of the omitted conditioning variables by modeling it.
- Of course, if the model for the bias term is not valid, neither is the correction for the bias.
- Semiparametric selection estimators are designed to protect the analyst against model misspecification.
- (See, e.g., ?).



- It controls for the bias of the omitted conditioning variables by modeling it.
- Of course, if the model for the bias term is not valid, neither is the correction for the bias.
- Semiparametric selection estimators are designed to protect the analyst against model misspecification.
- (See, e.g., ?).
- Matching evades this problem by assuming that the analyst always knows the correct conditioning variables and that they satisfy (M-1).

- It controls for the bias of the omitted conditioning variables by modeling it.
- Of course, if the model for the bias term is not valid, neither is the correction for the bias.
- Semiparametric selection estimators are designed to protect the analyst against model misspecification.
- (See, e.g., ?).
- Matching evades this problem by assuming that the analyst always knows the correct conditioning variables and that they satisfy (M-1).
- In actual empirical settings, agents rarely know the relevant information set.

- It controls for the bias of the omitted conditioning variables by modeling it.
- Of course, if the model for the bias term is not valid, neither is the correction for the bias.
- Semiparametric selection estimators are designed to protect the analyst against model misspecification.
- (See, e.g., ?).
- Matching evades this problem by assuming that the analyst always knows the correct conditioning variables and that they satisfy (M-1).
- In actual empirical settings, agents rarely know the relevant information set.
- Instead they use proxies.

## Adding Information to the Econometrician's Information Set: Using Proxies for the Relevant Information

- Suppose that instead of knowing some part of the minimal relevant information set, such as  $f_2$ , the analyst has access to a proxy for it.

## Adding Information to the Econometrician's Information Set: Using Proxies for the Relevant Information

- Suppose that instead of knowing some part of the minimal relevant information set, such as  $f_2$ , the analyst has access to a proxy for it.
- In particular, assume that he has access to a variable  $\tilde{Z}$  that is correlated with  $f_2$  but that is not the full minimal relevant information set.

## Adding Information to the Econometrician's Information Set: Using Proxies for the Relevant Information

- Suppose that instead of knowing some part of the minimal relevant information set, such as  $f_2$ , the analyst has access to a proxy for it.
- In particular, assume that he has access to a variable  $\tilde{Z}$  that is correlated with  $f_2$  but that is not the full minimal relevant information set.
- That is, define the econometrician's information to be

$$\tilde{I}_E = \{Z, \tilde{Z}\},$$

and suppose that he uses it so  $I_E = \tilde{I}_E$ .

- In order to obtain closed-form expressions for the biases we assume that

$$\begin{aligned}\tilde{Z} &\sim N(0, \sigma_{\tilde{Z}}^2) \\ \text{corr}(\tilde{Z}, f_2) &= \rho, \text{ and } \tilde{Z} \perp\!\!\!\perp (\varepsilon_0, \varepsilon_1, \varepsilon_V, f_1).\end{aligned}$$

We define expressions comparable to  $\eta$  and  $\eta'$  :

$$\begin{aligned}\tilde{\eta}_1 &= \frac{\alpha_{11}\alpha_{V1}\sigma_{f_1}^2 + \alpha_{12}\alpha_{V2}(1 - \rho^2)\sigma_{f_2}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \alpha_{V2}^2\sigma_{f_2}^2(1 - \rho^2) + \sigma_{\varepsilon_V}^2}} \\ \tilde{\eta}_0 &= \frac{\alpha_{01}\alpha_{V1}\sigma_{f_1}^2 + \alpha_{02}\alpha_{V2}(1 - \rho^2)\sigma_{f_2}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \alpha_{V2}^2\sigma_{f_2}^2(1 - \rho^2) + \sigma_{\varepsilon_V}^2}}.\end{aligned}$$

- By substituting for  $l'_E$  by  $\tilde{l}_E$  and  $\eta'_j$  by  $\tilde{\eta}_j$  ( $j = 0, 1$ ) in Conditions (1) and (2) of Slide 776, we can obtain results for the bias in this case.



- By substituting for  $I'_E$  by  $\tilde{I}_E$  and  $\eta'_j$  by  $\tilde{\eta}_j$  ( $j = 0, 1$ ) in Conditions (1) and (2) of Slide 776, we can obtain results for the bias in this case.
- Whether  $\tilde{I}_E$  will be bias-reducing depends on how well it spans  $I_R$  and on the signs of the terms in the absolute values in those conditions in Slide 776.

- In this case, however, there is another parameter to consider: the correlation  $\rho$  between  $\tilde{Z}$  and  $f_2$ ,  $\rho$ .

- In this case, however, there is another parameter to consider: the correlation  $\rho$  between  $\tilde{Z}$  and  $f_2$ ,  $\rho$ .
- If  $|\rho| = 1$  we are back to the case of  $\tilde{l}_E = l_E$  because  $\tilde{Z}$  is a perfect proxy for  $f_2$ . If  $\rho = 0$ , we are essentially back to the case analyzed in Slide 776.

- In this case, however, there is another parameter to consider: the correlation  $\rho$  between  $\tilde{Z}$  and  $f_2$ ,  $\rho$ .
- If  $|\rho| = 1$  we are back to the case of  $\tilde{l}_E = l_E$  because  $\tilde{Z}$  is a perfect proxy for  $f_2$ . If  $\rho = 0$ , we are essentially back to the case analyzed in Slide 776.
- Because we know that the bias at a particular value of  $\rho$  might either increase or decrease when  $f_2$  is used as a conditioning variable but  $f_1$  is not, we know that it is not possible to determine whether the bias increases or decreases as we change the correlation between  $f_2$  and  $\tilde{Z}$ .

- In this case, however, there is another parameter to consider: the correlation  $\rho$  between  $\tilde{Z}$  and  $f_2$ ,  $\rho$ .
- If  $|\rho| = 1$  we are back to the case of  $\tilde{l}_E = l_E$  because  $\tilde{Z}$  is a perfect proxy for  $f_2$ . If  $\rho = 0$ , we are essentially back to the case analyzed in Slide 776.
- Because we know that the bias at a particular value of  $\rho$  might either increase or decrease when  $f_2$  is used as a conditioning variable but  $f_1$  is not, we know that it is not possible to determine whether the bias increases or decreases as we change the correlation between  $f_2$  and  $\tilde{Z}$ .
- That is, we know that going from  $\rho = 0$  to  $|\rho| = 1$  might change the bias in any direction.

- In this case, however, there is another parameter to consider: the correlation  $\rho$  between  $\tilde{Z}$  and  $f_2$ ,  $\rho$ .
- If  $|\rho| = 1$  we are back to the case of  $\tilde{l}_E = l_E$  because  $\tilde{Z}$  is a perfect proxy for  $f_2$ . If  $\rho = 0$ , we are essentially back to the case analyzed in Slide 776.
- Because we know that the bias at a particular value of  $\rho$  might either increase or decrease when  $f_2$  is used as a conditioning variable but  $f_1$  is not, we know that it is not possible to determine whether the bias increases or decreases as we change the correlation between  $f_2$  and  $\tilde{Z}$ .
- That is, we know that going from  $\rho = 0$  to  $|\rho| = 1$  might change the bias in any direction.
- Use of a better proxy in this correlational sense may produce a *more* biased estimate.

- From the analysis of Slide 776, it is straightforward to derive conditions under which the bias generated when the econometrician's information is  $\tilde{I}_E$  is smaller than when it is  $I'_E$ .

- From the analysis of Slide 776, it is straightforward to derive conditions under which the bias generated when the econometrician's information is  $\tilde{I}_E$  is smaller than when it is  $I'_E$ .
- That is, it can be the case that knowing *the proxy* variable  $\tilde{Z}$  is *better* than knowing the actual variable  $f_2$ .

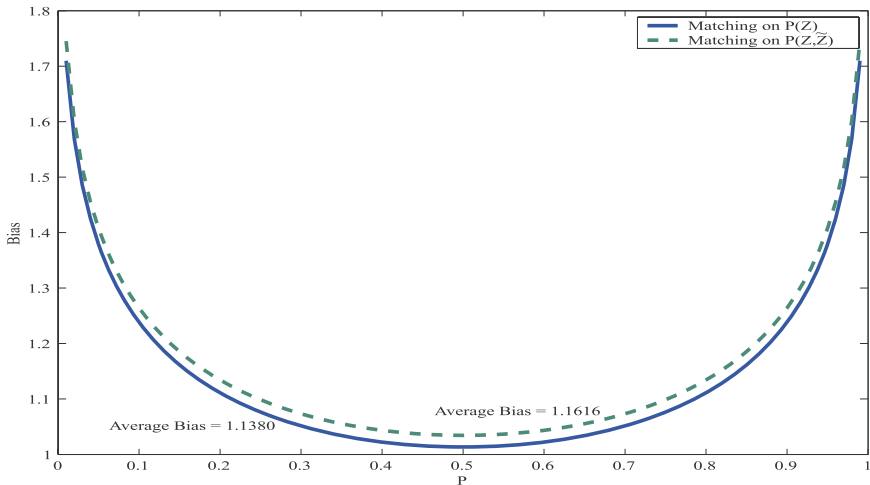


- From the analysis of Slide 776, it is straightforward to derive conditions under which the bias generated when the econometrician's information is  $\tilde{I}_E$  is smaller than when it is  $I_E$ .
- That is, it can be the case that knowing *the proxy* variable  $\tilde{Z}$  is *better* than knowing the actual variable  $f_2$ .
- Returning to the analysis of treatment on the treated as an example (i.e., Condition (1)), the bias in absolute value (at a fixed value of  $\rho$ ) is reduced when  $\tilde{Z}$  is used instead of  $f_2$  if

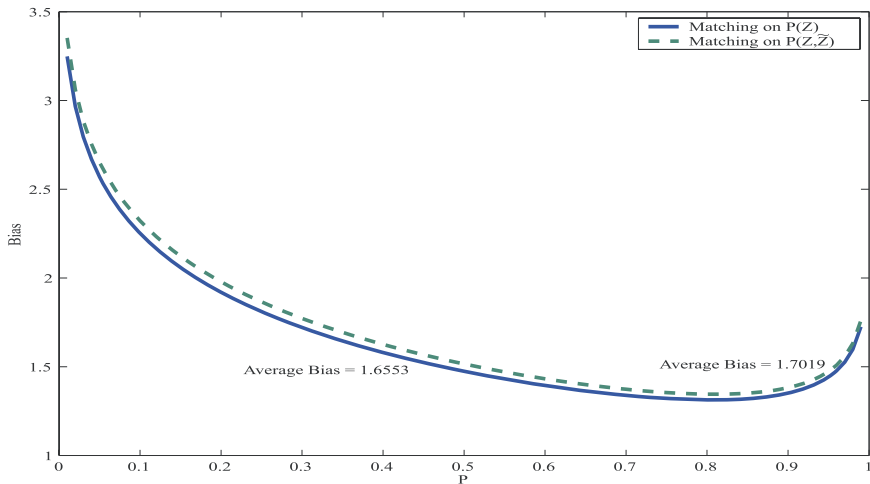
$$\left| \frac{\alpha_{01}\alpha_{V1}\sigma_{f_1}^2 + \alpha_{02}\alpha_{V2}(1-\rho^2)\sigma_{f_2}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \alpha_{V2}^2\sigma_{f_2}^2(1-\rho^2) + \sigma_{\varepsilon_V}^2}} \right| < \left| \frac{\alpha_{01}\alpha_{V1}\sigma_{f_1}^2}{\sqrt{\alpha_{V1}^2\sigma_{f_1}^2 + \sigma_{\varepsilon_V}^2}} \right|.$$

Figures 24A and 24B, use the same true model as used in the previous section to illustrate the two points being made here.

Figure 24: A. Bias for Treatment on the Treated



## B. Bias for Average Treatment Effect



Note: Using proxy  $\tilde{Z}$  for  $f_2$  increases the bias. Correlation  $(\tilde{Z}, f_2) = 0.5$ .

Model:

$$\begin{aligned} V &= Z + f_1 + f_2 + \varepsilon_V; & Y_1 &= 2f_1 + 0.1f_2 + \varepsilon_1; & Y_0 &= f_1 + 0.1f_2 + \varepsilon_0 \\ \varepsilon_V &\sim N(0, 1); & \varepsilon_1 &\sim N(0, 1); & \varepsilon_0 &\sim N(0, 1) \\ f_1 &\sim N(0, 1); & f_2 &\sim N(0, 1) \end{aligned}$$

Source: Heckman and Navarro (2005)

- Namely, *using* a *proxy* for an unobserved relevant variable *might increase the bias*.

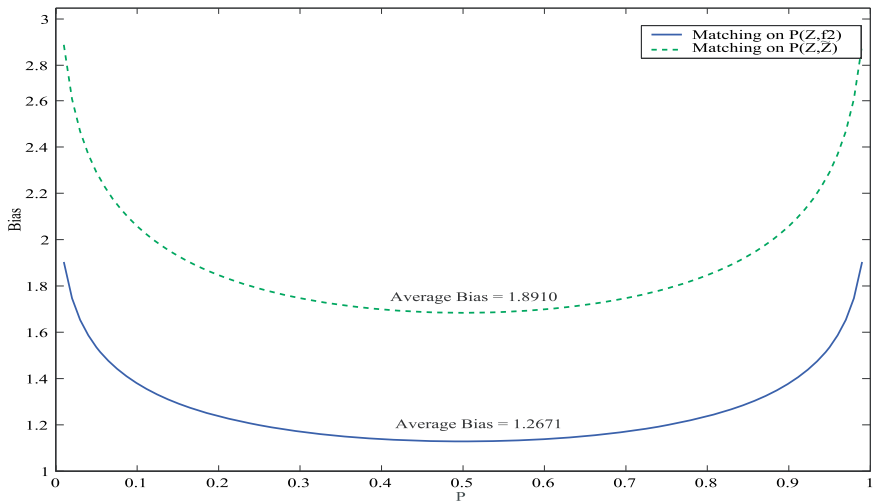
- Namely, *using* a *proxy* for an unobserved relevant variable *might increase the bias*.
- On the other hand, it *might be better* in terms of bias to use a *proxy* than to use the actual variable,  $f_2$ .

- Namely, *using* a *proxy* for an unobserved relevant variable *might increase the bias*.
- On the other hand, it *might be better* in terms of bias to use a *proxy* than to use the actual variable,  $f_2$ .
- However, as Figures 25A and 25B show, by changing  $\alpha_{02}$  from 0.1 to 1, using a proxy might increase the bias versus using the actual variable  $f_2$ .

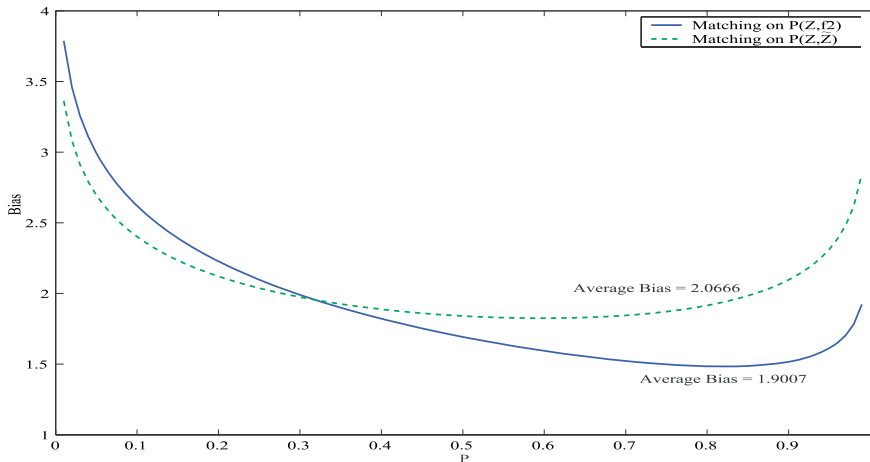
- Namely, *using* a *proxy* for an unobserved relevant variable *might increase the bias*.
- On the other hand, it *might be better* in terms of bias to use a *proxy* than to use the actual variable,  $f_2$ .
- However, as Figures 25A and 25B show, by changing  $\alpha_{02}$  from 0.1 to 1, using a proxy might increase the bias versus using the actual variable  $f_2$ .
- Notice that the bias need not be universally negative or positive but depends on  $p$ .



Figure 25: A. Bias for Treatment on the Treated



## B. Bias for Average Treatment Effect



Note: Using proxy  $\tilde{Z}$  for  $f_2$  increases the bias. Correlation  $(\tilde{Z}, f_2) = 0.5$ .

Model:

$$\begin{aligned} V &= Z + f_1 + f_2 + \varepsilon_V; & Y_1 &= 2f_1 + 0.1f_2 + \varepsilon_1; & Y_0 &= f_1 + f_2 + \varepsilon_0 \\ \varepsilon_V &\sim N(0, 1); & \varepsilon_1 &\sim N(0, 1); & \varepsilon_0 &\sim N(0, 1) \\ f_1 &\sim N(0, 1); & f_2 &\sim N(0, 1) \end{aligned}$$

Source: Heckman and Navarro (2005)

- The point of these examples is that matching makes very knife-edge assumptions.

- The point of these examples is that matching makes very knife-edge assumptions.
- If the analyst gets the right conditioning set, (M-1) is satisfied and there is no bias.

- The point of these examples is that matching makes very knife-edge assumptions.
- If the analyst gets the right conditioning set, (M-1) is satisfied and there is no bias.
- But determining the correct information set is not a trivial task, as we note in Slide 806.

- The point of these examples is that matching makes very knife-edge assumptions.
- If the analyst gets the right conditioning set, (M-1) is satisfied and there is no bias.
- But determining the correct information set is not a trivial task, as we note in Slide 806.
- Having good proxies in the standard usage of that term can create substantial bias in estimating treatment effects.

- The point of these examples is that matching makes very knife-edge assumptions.
- If the analyst gets the right conditioning set, (M-1) is satisfied and there is no bias.
- But determining the correct information set is not a trivial task, as we note in Slide 806.
- Having good proxies in the standard usage of that term can create substantial bias in estimating treatment effects.
- Half a loaf may be worse than none.



## The Case of a Discrete Outcome Variable

- ? construct parallel examples for cases including discrete dependent variables.

## The Case of a Discrete Outcome Variable

- ? construct parallel examples for cases including discrete dependent variables.
- In particular, they consider nonnormal, nonseparable equations for odds ratios and probabilities.

## The Case of a Discrete Outcome Variable

- ? construct parallel examples for cases including discrete dependent variables.
- In particular, they consider nonnormal, nonseparable equations for odds ratios and probabilities.
- The proposition that matching identifies the correct treatment parameter if the econometrician's information set includes all the minimal relevant information is true more generally, provided that any additional extraneous information used is exogenous in a sense to be defined precisely in the next section.

## On the Use of Model Selection Criteria to Choose Matching Variables

- We have already shown by way of example that adding more variables from a minimal relevant information set, but not all variables in it, may increase bias.

## On the Use of Model Selection Criteria to Choose Matching Variables

- We have already shown by way of example that adding more variables from a minimal relevant information set, but not all variables in it, may increase bias.
- By a parallel argument, adding additional variables to the relevant conditioning set may make the bias worse.

## On the Use of Model Selection Criteria to Choose Matching Variables

- We have already shown by way of example that adding more variables from a minimal relevant information set, but not all variables in it, may increase bias.
- By a parallel argument, adding additional variables to the relevant conditioning set may make the bias worse.
- Although we have used our prototypical Roy model as our point of departure, the point is more general.

- There is no rigorous rule for choosing the conditioning variables that produce (M-1).

- There is no rigorous rule for choosing the conditioning variables that produce (M-1).
- Adding variables that are statistically significant in the treatment choice equation is not guaranteed to select a set of conditioning variables that satisfies condition (M-1).



- There is no rigorous rule for choosing the conditioning variables that produce (M-1).
- Adding variables that are statistically significant in the treatment choice equation is not guaranteed to select a set of conditioning variables that satisfies condition (M-1).
- This is demonstrated by the analysis of Slide 776 that shows that adding  $f_2$  when it determines  $D$  may increase bias at any selected value of  $p$ .

- The existing literature (e.g., ?) proposes criteria based on selecting a set of conditioning variables based on a goodness of fit criterion ( $\lambda$ ), where a higher  $\lambda$  means a better fit in the equation predicting  $D$ .

- The existing literature (e.g., ?) proposes criteria based on selecting a set of conditioning variables based on a goodness of fit criterion ( $\lambda$ ), where a higher  $\lambda$  means a better fit in the equation predicting  $D$ .
- The intuition behind such criteria is that by using some measure of goodness of fit as a guiding principle one is using information relevant to the decision process.

- The existing literature (e.g., ?) proposes criteria based on selecting a set of conditioning variables based on a goodness of fit criterion ( $\lambda$ ), where a higher  $\lambda$  means a better fit in the equation predicting  $D$ .
- The intuition behind such criteria is that by using some measure of goodness of fit as a guiding principle one is using information relevant to the decision process.
- In the example of Slide 776, using  $f_2$  improves goodness of fit of the model for  $D$ , but increases bias for the parameters.

- The existing literature (e.g., ?) proposes criteria based on selecting a set of conditioning variables based on a goodness of fit criterion ( $\lambda$ ), where a higher  $\lambda$  means a better fit in the equation predicting  $D$ .
- The intuition behind such criteria is that by using some measure of goodness of fit as a guiding principle one is using information relevant to the decision process.
- In the example of Slide 776, using  $f_2$  improves goodness of fit of the model for  $D$ , but increases bias for the parameters.
- In general, such a rule is deficient if  $f_1$  is not known or is not used.

- An implicit assumption underlying such procedures is that the added conditioning variables  $\mathcal{X}$  are exogenous in the following sense:



$$(Y_0, Y_1) \perp\!\!\!\perp D | I_{\text{int}}, \mathcal{X} \quad (\text{E-1})$$

where  $I_{\text{int}}$  is interpreted as the variables initially used as conditioning variables before  $\mathcal{X}$  is added.

- $$(Y_0, Y_1) \perp\!\!\!\perp D \mid I_{\text{int}}, \mathcal{X} \quad (\text{E-1})$$

where  $I_{\text{int}}$  is interpreted as the variables initially used as conditioning variables before  $\mathcal{X}$  is added.

- Failure of exogeneity is a failure of (M-1) for the augmented conditioning set, and matching estimators based on the augmented information set  $(I_{\text{int}}, \mathcal{X})$  are biased when the condition is not satisfied.



- Exogeneity assumption (E-1) is not usually invoked in the matching literature, which largely focuses on problem P-1, evaluating a program in place, rather than extrapolating to new environments (P-2).

- Exogeneity assumption (E-1) is not usually invoked in the matching literature, which largely focuses on problem P-1, evaluating a program in place, rather than extrapolating to new environments (P-2).
- Indeed, the robustness of matching to such exogeneity assumptions is trumpeted as one of the virtues of the method.

- Exogeneity assumption (E-1) is not usually invoked in the matching literature, which largely focuses on problem P-1, evaluating a program in place, rather than extrapolating to new environments (P-2).
- Indeed, the robustness of matching to such exogeneity assumptions is trumpeted as one of the virtues of the method.
- In this section, we show some examples that illustrate the general point that standard model selection criteria fail to produce correctly specified conditioning sets unless some version of exogeneity condition (E-1) is satisfied.

- In the literature, the use of model selection criteria is justified in two different ways.

- In the literature, the use of model selection criteria is justified in two different ways.
- Sometimes it is claimed that they provide a *relative* guide.

- In the literature, the use of model selection criteria is justified in two different ways.
- Sometimes it is claimed that they provide a *relative* guide.
- Sets of variables with better goodness of fit in predicting  $D$  (a higher  $\lambda$  in the notation of table 12) are alleged to be better than sets of variables with lower  $\lambda$  in the sense that they generate lower biases.

- In the literature, the use of model selection criteria is justified in two different ways.
- Sometimes it is claimed that they provide a *relative* guide.
- Sets of variables with better goodness of fit in predicting  $D$  (a higher  $\lambda$  in the notation of table 12) are alleged to be better than sets of variables with lower  $\lambda$  in the sense that they generate lower biases.
- However, we have already shown that this is not true.

- In the literature, the use of model selection criteria is justified in two different ways.
- Sometimes it is claimed that they provide a *relative* guide.
- Sets of variables with better goodness of fit in predicting  $D$  (a higher  $\lambda$  in the notation of table 12) are alleged to be better than sets of variables with lower  $\lambda$  in the sense that they generate lower biases.
- However, we have already shown that this is not true.
- We know that enlarging the analyst's information from  $I_{\text{int}} = \{Z\}$  to  $I'_{\text{int}} = \{Z, f_2\}$  will improve fit since  $f_2$  is also in  $I_A$  and  $I_R$ .



- In the literature, the use of model selection criteria is justified in two different ways.
- Sometimes it is claimed that they provide a *relative* guide.
- Sets of variables with better goodness of fit in predicting  $D$  (a higher  $\lambda$  in the notation of table 12) are alleged to be better than sets of variables with lower  $\lambda$  in the sense that they generate lower biases.
- However, we have already shown that this is not true.
- We know that enlarging the analyst's information from  $I_{\text{int}} = \{Z\}$  to  $I'_{\text{int}} = \{Z, f_2\}$  will improve fit since  $f_2$  is also in  $I_A$  and  $I_R$ .
- But, going from  $I_{\text{int}}$  to  $I'_{\text{int}}$  might increase the bias.

Table 12: Goodness of fit statistics  $\lambda$

Variables in Probit	Goodness of fit statistics $\lambda$		Average Bias	
	Correct in-sample prediction rate	Pseudo $R^2$	TT	ATE
$Z$	66.88%	0.1284	1.1380	1.6553
$Z, f_2$	75.02%	0.2791	1.2671	1.9007
$Z, f_1, f_2$	83.45%	0.4844	0.0000	0.0000
$Z, S_1$	77.38%	0.3282	0.9612	1.3981
$Z, S_2$	92.25%	0.7498	0.9997	1.4541

Source: Heckman and Navarro (2004)

- So it is not true that combinations of variables that increase some measure of fit  $\lambda$  necessarily reduce the bias.

- So it is not true that combinations of variables that increase some measure of fit  $\lambda$  necessarily reduce the bias.
- Table 12 illustrates this point using our normal example.

- So it is not true that combinations of variables that increase some measure of fit  $\lambda$  necessarily reduce the bias.
- Table 12 illustrates this point using our normal example.
- Going from row 1 to row 2 (adding  $f_2$ ) improves goodness of fit and increases the unconditional or overall bias for all three treatment parameters, because (E-1) is violated.

- The following rule of thumb argument is sometimes invoked as an absolute standard against which to compare alternative models.

- The following rule of thumb argument is sometimes invoked as an absolute standard against which to compare alternative models.
- In versions of the argument, the analyst asserts that there is a combination of variables  $I''$  that satisfy (M-1) and hence produces zero bias and a value of  $\lambda = \lambda''$  larger than that of any other  $I$ .

- The following rule of thumb argument is sometimes invoked as an absolute standard against which to compare alternative models.
- In versions of the argument, the analyst asserts that there is a combination of variables  $I''$  that satisfy (M-1) and hence produces zero bias and a value of  $\lambda = \lambda''$  larger than that of any other  $I$ .
- In our examples, conditioning on  $\{Z, f_1, f_2\}$  generates zero bias.



- The following rule of thumb argument is sometimes invoked as an absolute standard against which to compare alternative models.
- In versions of the argument, the analyst asserts that there is a combination of variables  $I''$  that satisfy (M-1) and hence produces zero bias and a value of  $\lambda = \lambda''$  larger than that of any other  $I$ .
- In our examples, conditioning on  $\{Z, f_1, f_2\}$  generates zero bias.
- We can exclude  $Z$  and still obtain zero bias.

- The following rule of thumb argument is sometimes invoked as an absolute standard against which to compare alternative models.
- In versions of the argument, the analyst asserts that there is a combination of variables  $I''$  that satisfy (M-1) and hence produces zero bias and a value of  $\lambda = \lambda''$  larger than that of any other  $I$ .
- In our examples, conditioning on  $\{Z, f_1, f_2\}$  generates zero bias.
- We can exclude  $Z$  and still obtain zero bias.
- Because  $Z$  is a determinant of  $D$ , this shows immediately that the best fitting model does not necessarily identify the minimal relevant information set.

- In this example including  $Z$  is innocuous because there is still zero bias and the added conditioning variables satisfy (E-1) where  $I_{\text{int}} = (f_1, f_2)$ .

- In this example including  $Z$  is innocuous because there is still zero bias and the added conditioning variables satisfy (E-1) where  $I_{\text{int}} = (f_1, f_2)$ .
- In general, such a rule is not innocuous if  $Z$  is not exogenous.

- In this example including  $Z$  is innocuous because there is still zero bias and the added conditioning variables satisfy (E-1) where  $I_{\text{int}} = (f_1, f_2)$ .
- In general, such a rule is not innocuous if  $Z$  is not exogenous.
- If goodness of fit is used as a rule to choose variables on which to match, there is no guarantee it produces a desirable conditioning set.

- In this example including  $Z$  is innocuous because there is still zero bias and the added conditioning variables satisfy (E-1) where  $I_{\text{int}} = (f_1, f_2)$ .
- In general, such a rule is not innocuous if  $Z$  is not exogenous.
- If goodness of fit is used as a rule to choose variables on which to match, there is no guarantee it produces a desirable conditioning set.
- If we include in the conditioning set variables  $\mathcal{X}$  that violate (E-1), they may improve the fit of predicted probabilities but worsen the bias.

- ? produce a series of examples that have the following feature.

- ? produce a series of examples that have the following feature.
- Variables  $S$  (shown at the base of table 12) are added to the information set that improve the prediction of  $D$  but are correlated with  $(U_0, U_1)$ .



- ? produce a series of examples that have the following feature.
- Variables  $S$  (shown at the base of table 12) are added to the information set that improve the prediction of  $D$  but are correlated with  $(U_0, U_1)$ .
- Their particular examples use imperfect proxies  $(S_1, S_2)$  for  $(f_1, f_2)$ .

- ? produce a series of examples that have the following feature.
- Variables  $S$  (shown at the base of table 12) are added to the information set that improve the prediction of  $D$  but are correlated with  $(U_0, U_1)$ .
- Their particular examples use imperfect proxies  $(S_1, S_2)$  for  $(f_1, f_2)$ .
- The point is more general.

- ? produce a series of examples that have the following feature.
- Variables  $S$  (shown at the base of table 12) are added to the information set that improve the prediction of  $D$  but are correlated with  $(U_0, U_1)$ .
- Their particular examples use imperfect proxies  $(S_1, S_2)$  for  $(f_1, f_2)$ .
- The point is more general.
- The  $S$  variables fail exogeneity and produce greater bias for TT and ATE but they improve the prediction of  $D$  as measured by the correct in-sample prediction rate and the pseudo- $R^2$ .

- ? produce a series of examples that have the following feature.
- Variables  $S$  (shown at the base of table 12) are added to the information set that improve the prediction of  $D$  but are correlated with  $(U_0, U_1)$ .
- Their particular examples use imperfect proxies  $(S_1, S_2)$  for  $(f_1, f_2)$ .
- The point is more general.
- The  $S$  variables fail exogeneity and produce greater bias for TT and ATE but they improve the prediction of  $D$  as measured by the correct in-sample prediction rate and the pseudo- $R^2$ .
- See the bottom two rows of table 12.

- We next turn to the method of randomization, which is frequently held up to be an ideal approach for evaluating social programs.

- We next turn to the method of randomization, which is frequently held up to be an ideal approach for evaluating social programs.
- Randomization attempts to use a random assignment to achieve the conditional independence assumed in matching.

## Randomized Evaluations

- This section analyzes randomized social experiments as tools for evaluating social programs.

## Randomized Evaluations

- This section analyzes randomized social experiments as tools for evaluating social programs.
- In the introduction to this chapter, we discussed an ideal randomization where treatment status is randomly assigned.



## Randomized Evaluations

- This section analyzes randomized social experiments as tools for evaluating social programs.
- In the introduction to this chapter, we discussed an ideal randomization where treatment status is randomly assigned.
- In this section, we discuss actual social experiments, where self-selection decisions often intrude on the randomization decisions of experimenters.

- Two cases have been made for the application of social experimentation.

- Two cases have been made for the application of social experimentation.
- One case is a classical argument in experimental design.

- Two cases have been made for the application of social experimentation.
- One case is a classical argument in experimental design.
- Inducing variation in regressors increases precision of estimates and the power of tests.

- Two cases have been made for the application of social experimentation.
- One case is a classical argument in experimental design.
- Inducing variation in regressors increases precision of estimates and the power of tests.
- The other case focuses on solving endogeneity and self-selection problems.

- Two cases have been made for the application of social experimentation.
- One case is a classical argument in experimental design.
- Inducing variation in regressors increases precision of estimates and the power of tests.
- The other case focuses on solving endogeneity and self-selection problems.
- Randomization is an instrumental variable.

- Two cases have been made for the application of social experimentation.
- One case is a classical argument in experimental design.
- Inducing variation in regressors increases precision of estimates and the power of tests.
- The other case focuses on solving endogeneity and self-selection problems.
- Randomization is an instrumental variable.
- The two cases are mutually compatible, but involve different emphases.

- Both cases can be motivated within a linear regression model for outcome  $Y$  with treatment indicator  $D$  and covariates  $X$ :

$$Y = X\alpha + D\beta + U, \quad (70)$$

where  $U$  is an unobservable.



- Both cases can be motivated within a linear regression model for outcome  $Y$  with treatment indicator  $D$  and covariates  $X$ :

$$Y = X\alpha + D\beta + U, \quad (70)$$

where  $U$  is an unobservable.

- $\beta$  may be the same for all observations (conditional on  $X$ ) as in the common coefficient setup, or it may be a variable coefficient of the type extensively in this chapter.

- Both cases can be motivated within a linear regression model for outcome  $Y$  with treatment indicator  $D$  and covariates  $X$ :

$$Y = X\alpha + D\beta + U, \quad (70)$$

where  $U$  is an unobservable.

- $\beta$  may be the same for all observations (conditional on  $X$ ) as in the common coefficient setup, or it may be a variable coefficient of the type extensively in this chapter.
- $D$  (and the  $X$ ) may be statistically dependent on  $U$ .

- Both cases can be motivated within a linear regression model for outcome  $Y$  with treatment indicator  $D$  and covariates  $X$ :

$$Y = X\alpha + D\beta + U, \quad (70)$$

where  $U$  is an unobservable.

- $\beta$  may be the same for all observations (conditional on  $X$ ) as in the common coefficient setup, or it may be a variable coefficient of the type extensively in this chapter.
- $D$  (and the  $X$ ) may be statistically dependent on  $U$ .
- We also entertain the possibility that when  $\beta$  is random it is dependent on  $D$ , as in the generalized Roy model.

- Both cases for social experimentation seek to secure identification of some parameters of (70) or parameters that can be generated from (70).

- Both cases for social experimentation seek to secure identification of some parameters of (70) or parameters that can be generated from (70).
- Analysts advocating the first case for experimentation typically assume a common coefficient model for  $\alpha$  and  $\beta$ .

- Both cases for social experimentation seek to secure identification of some parameters of (70) or parameters that can be generated from (70).
- Analysts advocating the first case for experimentation typically assume a common coefficient model for  $\alpha$  and  $\beta$ .
- They address the problem that variation in  $(X, D)$  may be insufficient to identify or precisely estimate  $(\alpha, \beta)$ .

- Both cases for social experimentation seek to secure identification of some parameters of (70) or parameters that can be generated from (70).
- Analysts advocating the first case for experimentation typically assume a common coefficient model for  $\alpha$  and  $\beta$ .
- They address the problem that variation in  $(X, D)$  may be insufficient to identify or precisely estimate  $(\alpha, \beta)$ .
- Manipulating  $(X, D)$  through randomization, or more generally, through controlled variation, can secure identification.

- Both cases for social experimentation seek to secure identification of some parameters of (70) or parameters that can be generated from (70).
- Analysts advocating the first case for experimentation typically assume a common coefficient model for  $\alpha$  and  $\beta$ .
- They address the problem that variation in  $(X, D)$  may be insufficient to identify or precisely estimate  $(\alpha, \beta)$ .
- Manipulating  $(X, D)$  through randomization, or more generally, through controlled variation, can secure identification.
- It is typically assumed that  $(X, D)$  is independent of  $U$  or at least mean independent.



- Both cases for social experimentation seek to secure identification of some parameters of (70) or parameters that can be generated from (70).
- Analysts advocating the first case for experimentation typically assume a common coefficient model for  $\alpha$  and  $\beta$ .
- They address the problem that variation in  $(X, D)$  may be insufficient to identify or precisely estimate  $(\alpha, \beta)$ .
- Manipulating  $(X, D)$  through randomization, or more generally, through controlled variation, can secure identification.
- It is typically assumed that  $(X, D)$  is independent of  $U$  or at least mean independent.
- This is the traditional case analyzed in a large literature on experimental design in statistics.

- Good examples in economics of experimentation designed to increase the variation in the regressors are studies by ?, ?, and ???.

- Good examples in economics of experimentation designed to increase the variation in the regressors are studies by ?, ?, and ???.
- The papers by Conlisk show how experimental manipulation can solve a multicollinearity problem.

- Good examples in economics of experimentation designed to increase the variation in the regressors are studies by ?, ?, and ???.
- The papers by Conlisk show how experimental manipulation can solve a multicollinearity problem.
- In analyzing the effects of taxes on labor supply, it is necessary to isolate the effect of wages (the substitution effect) from the effect of pure asset income (the income effect) on labor supply.

- Good examples in economics of experimentation designed to increase the variation in the regressors are studies by ?, ?, and ???.
- The papers by Conlisk show how experimental manipulation can solve a multicollinearity problem.
- In analyzing the effects of taxes on labor supply, it is necessary to isolate the effect of wages (the substitution effect) from the effect of pure asset income (the income effect) on labor supply.
- In observational data, empirical measures of wages and asset income are highly intercorrelated.

- Good examples in economics of experimentation designed to increase the variation in the regressors are studies by ?, ?, and ???.
- The papers by Conlisk show how experimental manipulation can solve a multicollinearity problem.
- In analyzing the effects of taxes on labor supply, it is necessary to isolate the effect of wages (the substitution effect) from the effect of pure asset income (the income effect) on labor supply.
- In observational data, empirical measures of wages and asset income are highly intercorrelated.
- In addition, asset income is often poorly measured.

- By experimentally assigning these variables as in the negative income tax experiments, it is possible to identify both income and substitution effects in labor supply equations (see ?).

- By experimentally assigning these variables as in the negative income tax experiments, it is possible to identify both income and substitution effects in labor supply equations (see ?).
- ? shows how variation in the prices paid for electricity across the day can identify price effects that cannot be identified in regimes with uniform prices across all hours of the day.



- Random assignment is not essential to this approach.

- Random assignment is not essential to this approach.
- Any regressor assignment rule based on variables  $Q$  that are stochastically independent of  $U$  will suffice, although the efficiency of the estimates will depend on the choice of  $Q$  and care must be taken to avoid inducing multicollinearity by the choice of an assignment rule.

- The second case for social experiments and the one that receives the most attention in applied work in economics and in this chapter focuses on the dependence between  $(X, D)$  and  $U$  that invalidates least squares as an estimator of the causal effect of  $X$  and  $D$  on  $Y$ .

- The second case for social experiments and the one that receives the most attention in applied work in economics and in this chapter focuses on the dependence between  $(X, D)$  and  $U$  that invalidates least squares as an estimator of the causal effect of  $X$  and  $D$  on  $Y$ .
- This is the problem of least squares bias raised by ? and extensively discussed in Part I.

- The second case for social experiments and the one that receives the most attention in applied work in economics and in this chapter focuses on the dependence between  $(X, D)$  and  $U$  that invalidates least squares as an estimator of the causal effect of  $X$  and  $D$  on  $Y$ .
- This is the problem of least squares bias raised by ? and extensively discussed in Part I.
- In the second case, experimental variation in  $(X, D)$  is sought to make it “exogenous” or “external” to  $U$ .

- The second case for social experiments and the one that receives the most attention in applied work in economics and in this chapter focuses on the dependence between  $(X, D)$  and  $U$  that invalidates least squares as an estimator of the causal effect of  $X$  and  $D$  on  $Y$ .
- This is the problem of least squares bias raised by ? and extensively discussed in Part I.
- In the second case, experimental variation in  $(X, D)$  is sought to make it “exogenous” or “external” to  $U$ .
- A popular argument in favor of experiments is that they produce simple, transparent estimates of the effects of the programs being evaluated in the presence of such biases.

- A quotation from ? is apt:

The beauty of randomized evaluations is that the results are what they are: we compare the outcome in the treatment *group* with the outcome in the control group, see whether they are different, and if so by how much. Interpreting quasi-experiments sometimes requires statistical legerdemain, which makes them less attractive . . .

- This argument assumes that interesting evaluation questions can be answered by the marginal distributions produced from experiments.



- This argument assumes that interesting evaluation questions can be answered by the marginal distributions produced from experiments.
- It also assumes that no economic model is needed to interpret evidence, contrary to a main theme of this chapter.

- Randomization is an instrument.

- Randomization is an instrument.
- As such, it shares all of the assets and liabilities of IV already discussed.

- Randomization is an instrument.
- As such, it shares all of the assets and liabilities of IV already discussed.
- In particular, randomization applied to a correlated random coefficient (or a model of essential heterogeneity) raises the same issues about the multiplicity of parameters identified by different randomizations as were discussed there in connection with the multiplicity of parameters identified by different instruments.

- The two popular arguments for social experimentation are closely related.

- The two popular arguments for social experimentation are closely related.
- Exogenous variation in  $(X, D)$  can, if judiciously administered, solve collinearity, precision, and endogeneity problems.

- The two popular arguments for social experimentation are closely related.
- Exogenous variation in  $(X, D)$  can, if judiciously administered, solve collinearity, precision, and endogeneity problems.
- Applying the terminology of Part I to the analysis of model (70), randomization can identify a model that can solve all three policy evaluation problems: (P-1), the problem of internal validity; (P-2), the problem of extrapolation to new environments (by virtue of the linearity of (70)); and (P-3), the problem of forecasting new policies that can be described by identifiable functions of  $(X, D)$  and any external variables.

- As noted in the concluding section of Part I, the modern literature tends to reject functional form assumptions such as those embodied in equation (70).



- As noted in the concluding section of Part I, the modern literature tends to reject functional form assumptions such as those embodied in equation (70).
- It has evolved towards a more focused attempt to solve problem P-1 to protect against endogeneity of  $D$  with respect to  $U$ .

- As noted in the concluding section of Part I, the modern literature tends to reject functional form assumptions such as those embodied in equation (70).
- It has evolved towards a more focused attempt to solve problem P-1 to protect against endogeneity of  $D$  with respect to  $U$ .
- Sometimes the parameter being identified is not clearly specified.

- As noted in the concluding section of Part I, the modern literature tends to reject functional form assumptions such as those embodied in equation (70).
- It has evolved towards a more focused attempt to solve problem P-1 to protect against endogeneity of  $D$  with respect to  $U$ .
- Sometimes the parameter being identified is not clearly specified.
- When it is, this focus implements Marschak's Maxim of doing one thing well, as discussed in Part I.

- Common to the literature on IV estimation, proponents of randomization often ignore the consequences of heterogeneity in  $\beta$  and dependence of  $\beta$  on  $D$ —the problem of essential heterogeneity.

- Common to the literature on IV estimation, proponents of randomization often ignore the consequences of heterogeneity in  $\beta$  and dependence of  $\beta$  on  $D$ —the problem of essential heterogeneity.
- Our discussion in the previous sections applies with full force to randomization as an instrument.

- Common to the literature on IV estimation, proponents of randomization often ignore the consequences of heterogeneity in  $\beta$  and dependence of  $\beta$  on  $D$ —the problem of essential heterogeneity.
- Our discussion in the previous sections applies with full force to randomization as an instrument.
- Only if the randomization (instrument) corresponds exactly to the policy that is sought to be evaluated will the IV (randomization) identify the parameters of economic interest.

- Common to the literature on IV estimation, proponents of randomization often ignore the consequences of heterogeneity in  $\beta$  and dependence of  $\beta$  on  $D$ —the problem of essential heterogeneity.
- Our discussion in the previous sections applies with full force to randomization as an instrument.
- Only if the randomization (instrument) corresponds exactly to the policy that is sought to be evaluated will the IV (randomization) identify the parameters of economic interest.
- This section considers the case for randomization as an instrumental variable to solve endogeneity problems.

## Randomization as an Instrumental Variable

- The argument justifying randomization as an instrument assumes that randomization (or more generally the treatment assignment rule) does not alter subjective or objective potential outcomes.



## Randomization as an Instrumental Variable

- The argument justifying randomization as an instrument assumes that randomization (or more generally the treatment assignment rule) does not alter subjective or objective potential outcomes.
- This is covered by assumption (PI-3) presented in Part I.

## Randomization as an Instrumental Variable

- The argument justifying randomization as an instrument assumes that randomization (or more generally the treatment assignment rule) does not alter subjective or objective potential outcomes.
- This is covered by assumption (PI-3) presented in Part I.
- We also maintain absence of general equilibrium effects (PI-4) throughout this section.

## Randomization as an Instrumental Variable

- The argument justifying randomization as an instrument assumes that randomization (or more generally the treatment assignment rule) does not alter subjective or objective potential outcomes.
- This is covered by assumption (PI-3) presented in Part I.
- We also maintain absence of general equilibrium effects (PI-4) throughout this section.
- We discuss violations of (PI-3) when we discuss randomization bias.

- To be explicit about particular randomization mechanisms, we return to our touchstone generalized Roy model.

- To be explicit about particular randomization mechanisms, we return to our touchstone generalized Roy model.
- Potential outcomes are  $(Y_0, Y_1)$  and cost of participation is  $C$ .

- To be explicit about particular randomization mechanisms, we return to our touchstone generalized Roy model.
- Potential outcomes are  $(Y_0, Y_1)$  and cost of participation is  $C$ .
- Assume perfect certainty in the absence of randomization.

- To be explicit about particular randomization mechanisms, we return to our touchstone generalized Roy model.
- Potential outcomes are  $(Y_0, Y_1)$  and cost of participation is  $C$ .
- Assume perfect certainty in the absence of randomization.
- Under self-selection, the treatment choice is governed by

$$D = \mathbf{1}(Y_1 - Y_0 - C \geq 0).$$

This model of program participation abstracts from the important practical feature of many social programs that multiple agents contribute to decisions about program participation.

- To be explicit about particular randomization mechanisms, we return to our touchstone generalized Roy model.
- Potential outcomes are  $(Y_0, Y_1)$  and cost of participation is  $C$ .
- Assume perfect certainty in the absence of randomization.
- Under self-selection, the treatment choice is governed by

$$D = \mathbf{1}(Y_1 - Y_0 - C \geq 0).$$

This model of program participation abstracts from the important practical feature of many social programs that multiple agents contribute to decisions about program participation.

- We consider a more general framework in Slide 882.



- We assume additive separability between the observables  $(X, W)$  and the unobservables  $(U_0, U_1, U_C)$  for convenience:

$$\begin{aligned} Y_1 &= \mu_1(X) + U_1, & Y_0 &= \mu_0(X) + U_0, \\ C &= \mu_C(W) + U_C, & V &= U_1 - U_0 - U_C, \\ \mu_I(X, W) &= \mu_1(X) - \mu_0(X) - \mu_C(W), & Z &= (X, W). \end{aligned}$$

Only some components of  $X$  and/or  $W$  may be randomized.

- We assume additive separability between the observables  $(X, W)$  and the unobservables  $(U_0, U_1, U_C)$  for convenience:

$$\begin{aligned} Y_1 &= \mu_1(X) + U_1, & Y_0 &= \mu_0(X) + U_0, \\ C &= \mu_C(W) + U_C, & V &= U_1 - U_0 - U_C, \\ \mu_I(X, W) &= \mu_1(X) - \mu_0(X) - \mu_C(W), & Z &= (X, W). \end{aligned}$$

Only some components of  $X$  and/or  $W$  may be randomized.

- Randomization can be performed unconditionally or conditional on strata,  $Q$ , where the strata may or may not include components of  $(X, W)$  that are not randomized.

- We assume additive separability between the observables  $(X, W)$  and the unobservables  $(U_0, U_1, U_C)$  for convenience:

$$\begin{aligned} Y_1 &= \mu_1(X) + U_1, & Y_0 &= \mu_0(X) + U_0, \\ C &= \mu_C(W) + U_C, & V &= U_1 - U_0 - U_C, \\ \mu_I(X, W) &= \mu_1(X) - \mu_0(X) - \mu_C(W), & Z &= (X, W). \end{aligned}$$

Only some components of  $X$  and/or  $W$  may be randomized.

- Randomization can be performed unconditionally or conditional on strata,  $Q$ , where the strata may or may not include components of  $(X, W)$  that are not randomized.
- Specifically, it can be performed conditional on  $X$ , just as in our analysis of IV.

- We assume additive separability between the observables  $(X, W)$  and the unobservables  $(U_0, U_1, U_C)$  for convenience:

$$\begin{aligned}
 Y_1 &= \mu_1(X) + U_1, & Y_0 &= \mu_0(X) + U_0, \\
 C &= \mu_C(W) + U_C, & V &= U_1 - U_0 - U_C, \\
 \mu_I(X, W) &= \mu_1(X) - \mu_0(X) - \mu_C(W), & Z &= (X, W).
 \end{aligned}$$

Only some components of  $X$  and/or  $W$  may be randomized.

- Randomization can be performed unconditionally or conditional on strata,  $Q$ , where the strata may or may not include components of  $(X, W)$  that are not randomized.
- Specifically, it can be performed conditional on  $X$ , just as in our analysis of IV.
- Parameters can be defined conditional on  $X$ .

- Examples of treatments randomly assigned include the tax/benefit plans of the negative income tax programs; the price of electricity over the course of the day; variable tolls and bonuses; textbooks to pupils; reducing class size.

- Examples of treatments randomly assigned include the tax/benefit plans of the negative income tax programs; the price of electricity over the course of the day; variable tolls and bonuses; textbooks to pupils; reducing class size.
- Under invariance condition (PI-3), the functions  $\mu_0(X)$ ,  $\mu_1(X)$ ,  $\mu_C(W)$  (and hence  $\mu_I(X, W)$ ) are invariant to such modifications.

- Examples of treatments randomly assigned include the tax/benefit plans of the negative income tax programs; the price of electricity over the course of the day; variable tolls and bonuses; textbooks to pupils; reducing class size.
- Under invariance condition (PI-3), the functions  $\mu_0(X)$ ,  $\mu_1(X)$ ,  $\mu_C(W)$  (and hence  $\mu_I(X, W)$ ) are invariant to such modifications.
- The intervention is assumed to change the arguments of functions without shifting the functions themselves.

- Examples of treatments randomly assigned include the tax/benefit plans of the negative income tax programs; the price of electricity over the course of the day; variable tolls and bonuses; textbooks to pupils; reducing class size.
- Under invariance condition (PI-3), the functions  $\mu_0(X)$ ,  $\mu_1(X)$ ,  $\mu_C(W)$  (and hence  $\mu_I(X, W)$ ) are invariant to such modifications.
- The intervention is assumed to change the arguments of functions without shifting the functions themselves.
- Thus for the intervention of randomization, the functions are assumed to be structural in the sense of ?.



- Examples of treatments randomly assigned include the tax/benefit plans of the negative income tax programs; the price of electricity over the course of the day; variable tolls and bonuses; textbooks to pupils; reducing class size.
- Under invariance condition (PI-3), the functions  $\mu_0(X)$ ,  $\mu_1(X)$ ,  $\mu_C(W)$  (and hence  $\mu_I(X, W)$ ) are invariant to such modifications.
- The intervention is assumed to change the arguments of functions without shifting the functions themselves.
- Thus for the intervention of randomization, the functions are assumed to be structural in the sense of ?.
- The distributions of  $(U_0, U_1, U_C)$  conditional on  $X$ , and hence the distribution of  $V$  conditional on  $X$ , are also invariant.

- Under full compliance, the manipulated  $Z$  are the same as the  $Z$  facing the agent.

- Under full compliance, the manipulated  $Z$  are the same as the  $Z$  facing the agent.
- We formalize this assumption:

(R-4)

*The  $Z$  assigned agent  $\omega$  conditional on  $X$  are the  $Z$  realized and acted on by the agent conditional on  $X$ .*

- In terms of the generalized Roy model, this assumption states that the  $Z$  assigned  $\omega$  given  $X$  is the  $W$  that appears in the cost function and the derived decision rule.

- Some randomizations alter the environments facing agents in a more fundamental way by introducing new random variables into the model instead of modifying the variables that would be present in a pre-experimental environment.

- Some randomizations alter the environments facing agents in a more fundamental way by introducing new random variables into the model instead of modifying the variables that would be present in a pre-experimental environment.
- Comparisons of these randomizations involve an implicit dynamics, better expositied using the dynamic models presented in Part III.

- Some randomizations alter the environments facing agents in a more fundamental way by introducing new random variables into the model instead of modifying the variables that would be present in a pre-experimental environment.
- Comparisons of these randomizations involve an implicit dynamics, better expositied using the dynamic models presented in Part III.
- For simplicity and to present some main ideas, we initially invoke an implicit dynamics suitable to the generalized Roy model.



- Some randomizations alter the environments facing agents in a more fundamental way by introducing new random variables into the model instead of modifying the variables that would be present in a pre-experimental environment.
- Comparisons of these randomizations involve an implicit dynamics, better expositied using the dynamic models presented in Part III.
- For simplicity and to present some main ideas, we initially invoke an implicit dynamics suitable to the generalized Roy model.
- We develop a more explicit dynamic model of randomized evaluation in Slide 882.

- The most commonly used randomizations restrict eligibility either in advance of agent decisions about participation in a program or after agent decisions are made, but before actual participation begins.

- The most commonly used randomizations restrict eligibility either in advance of agent decisions about participation in a program or after agent decisions are made, but before actual participation begins.
- Unlike statistical discussions of randomization, we build agent choice front and center into our analysis.

- The most commonly used randomizations restrict eligibility either in advance of agent decisions about participation in a program or after agent decisions are made, but before actual participation begins.
- Unlike statistical discussions of randomization, we build agent choice front and center into our analysis.
- Agents choose and experimenters can only manipulate choice sets.

- Let  $\xi = 1$  if an agent is eligible to participate in the program;  
 $\xi = 0$  otherwise.

- Let  $\xi = 1$  if an agent is eligible to participate in the program;  $\xi = 0$  otherwise.
- $\tilde{\xi} = \{0, 1\}$  is the set of possible values of  $\xi$ .

- Let  $\xi = 1$  if an agent is eligible to participate in the program;  $\xi = 0$  otherwise.
- $\tilde{\xi} = \{0, 1\}$  is the set of possible values of  $\xi$ .
- Let  $D$  indicate participation under ordinary conditions.

- Let  $\xi = 1$  if an agent is eligible to participate in the program;  $\xi = 0$  otherwise.
- $\tilde{\xi} = \{0, 1\}$  is the set of possible values of  $\xi$ .
- Let  $D$  indicate participation under ordinary conditions.
- In the absence of randomization,  $D$  is an indicator of whether the agent actually participates in the program.



- Let  $\xi = 1$  if an agent is eligible to participate in the program;  $\xi = 0$  otherwise.
- $\tilde{\xi} = \{0, 1\}$  is the set of possible values of  $\xi$ .
- Let  $D$  indicate participation under ordinary conditions.
- In the absence of randomization,  $D$  is an indicator of whether the agent actually participates in the program.
- Let actual participation be  $A$ .

- Let  $\xi = 1$  if an agent is eligible to participate in the program;  $\xi = 0$  otherwise.
- $\tilde{\xi} = \{0, 1\}$  is the set of possible values of  $\xi$ .
- Let  $D$  indicate participation under ordinary conditions.
- In the absence of randomization,  $D$  is an indicator of whether the agent actually participates in the program.
- Let actual participation be  $A$ .
- By construction, under invariance condition (PI-3) presented in Part I,

$$A = D\xi. \quad (71)$$

This assumes that eligibility is strictly enforced.

- There is a distinction between desired participation by the agent ( $D$ ) and actual participation ( $A$ ).

- There is a distinction between desired participation by the agent ( $D$ ) and actual participation ( $A$ ).
- This distinction is conceptually distinct from the *ex-ante*, *ex-post* distinction.

- There is a distinction between desired participation by the agent ( $D$ ) and actual participation ( $A$ ).
- This distinction is conceptually distinct from the *ex-ante*, *ex-post* distinction.
- At all stages of the application and enrollment process, agents may be perfectly informed about their value of  $\xi$  and desire to participate ( $D$ ), but may not be allowed to participate.

- There is a distinction between desired participation by the agent ( $D$ ) and actual participation ( $A$ ).
- This distinction is conceptually distinct from the *ex-ante*, *ex-post* distinction.
- At all stages of the application and enrollment process, agents may be perfectly informed about their value of  $\xi$  and desire to participate ( $D$ ), but may not be allowed to participate.
- On the other hand, the agent may be surprised by  $\xi$  after applying to the program.

- There is a distinction between desired participation by the agent ( $D$ ) and actual participation ( $A$ ).
- This distinction is conceptually distinct from the *ex-ante*, *ex-post* distinction.
- At all stages of the application and enrollment process, agents may be perfectly informed about their value of  $\xi$  and desire to participate ( $D$ ), but may not be allowed to participate.
- On the other hand, the agent may be surprised by  $\xi$  after applying to the program.
- In this case, there is revelation of information and there is a distinction between *ex ante* expectations and *ex post* realizations.

- There is a distinction between desired participation by the agent ( $D$ ) and actual participation ( $A$ ).
- This distinction is conceptually distinct from the *ex-ante*, *ex-post* distinction.
- At all stages of the application and enrollment process, agents may be perfectly informed about their value of  $\xi$  and desire to participate ( $D$ ), but may not be allowed to participate.
- On the other hand, the agent may be surprised by  $\xi$  after applying to the program.
- In this case, there is revelation of information and there is a distinction between *ex ante* expectations and *ex post* realizations.
- Our analysis covers both cases.



We consider two types of randomization of eligibility.

Randomization of Type 1. *A random mechanism (possibly conditional on  $(X, Z)$ ) is used to determine  $\xi$ . The probability of eligibility is  $\Pr(\xi = 1 \mid X, Z)$ .*

- For this type of randomization, in the context of the generalized Roy model, it is assumed that

- (e-1a)  $\xi \perp\!\!\!\perp (U_0, U_1, U_C) \mid X, Z$  (**Randomization of Eligibility**)

and

- (e-1b)  $\Pr(A = 1 \mid X, Z, \xi)$  depends on  $\xi$ .

- This randomization affects the eligibility of the agent for the program but because agents still self-select, there is no assurance that eligible agents will participate in the program.

- This randomization affects the eligibility of the agent for the program but because agents still self-select, there is no assurance that eligible agents will participate in the program.
- This condition does not impose exogeneity on  $X, Z$ .

- This randomization affects the eligibility of the agent for the program but because agents still self-select, there is no assurance that eligible agents will participate in the program.
- This condition does not impose exogeneity on  $X, Z$ .
- Thus  $Z$  can fail as an instrument but  $\xi$  remains a valid instrument.

- This randomization affects the eligibility of the agent for the program but because agents still self-select, there is no assurance that eligible agents will participate in the program.
- This condition does not impose exogeneity on  $X, Z$ .
- Thus  $Z$  can fail as an instrument but  $\xi$  remains a valid instrument.
- Alternatively, (e-1a) and (e-1b) may be formulated according to the notation of ?.

- This randomization affects the eligibility of the agent for the program but because agents still self-select, there is no assurance that eligible agents will participate in the program.
- This condition does not impose exogeneity on  $X, Z$ .
- Thus  $Z$  can fail as an instrument but  $\xi$  remains a valid instrument.
- Alternatively, (e-1a) and (e-1b) may be formulated according to the notation of ?.
- Define  $A(z, e)$  to be the value of  $A$  when we set  $Z = z$  and  $\xi = e$ .



- This randomization affects the eligibility of the agent for the program but because agents still self-select, there is no assurance that eligible agents will participate in the program.
- This condition does not impose exogeneity on  $X, Z$ .
- Thus  $Z$  can fail as an instrument but  $\xi$  remains a valid instrument.
- Alternatively, (e-1a) and (e-1b) may be formulated according to the notation of ?.
- Define  $A(z, e)$  to be the value of  $A$  when we set  $Z = z$  and  $\xi = e$ .
- Define  $\mathcal{Z}$  as the set of admissible  $Z$  and  $\tilde{\xi}$  as the set of admissible  $\xi$ .

- This randomization affects the eligibility of the agent for the program but because agents still self-select, there is no assurance that eligible agents will participate in the program.
- This condition does not impose exogeneity on  $X, Z$ .
- Thus  $Z$  can fail as an instrument but  $\xi$  remains a valid instrument.
- Alternatively, (e-1a) and (e-1b) may be formulated according to the notation of ?.
- Define  $A(z, e)$  to be the value of  $A$  when we set  $Z = z$  and  $\xi = e$ .
- Define  $\mathcal{Z}$  as the set of admissible  $Z$  and  $\tilde{\xi}$  as the set of admissible  $\xi$ .
- In this notation, we may rewrite assumptions (e-1a) and (e-1b) as

- (e-1a)'  $\xi \perp\!\!\!\perp \left( Y_0, Y_1, \{A(z, e)\}_{(z,e) \in \mathcal{Z} \times \tilde{\xi}} \right) \mid X, Z$

and

- (e-1b)'  $\Pr(A = 1 \mid X, Z, \xi)$  depends on  $\xi$ .

- A second type of randomization conditions on individuals manifesting a desire to participate through their decision to apply to the program.

- A second type of randomization conditions on individuals manifesting a desire to participate through their decision to apply to the program.
- This type of randomization is widely used.

Randomization of Type 2: *Eligibility may be a function of  $D$  (conditionally on some or all components of  $X, Z, Q$  or unconditionally ). It is common to deny entry into programs among people who applied and were accepted into the program ( $D = 1$  ) so the probability of eligibility is  $\Pr(\xi = 1 \mid X, Z, Q, D = 1)$ . This assumes (PI-3) stated in Part I.*

- For this type of randomization of eligibility, it is assumed that

and

- (e-2b)  $\Pr(A = 1 \mid X, Z, D = 1, \xi = 1) = 1;$   
 $\Pr(A = 1 \mid X, Z, D = 1, \xi = 0) = 0.$

- For this type of randomization of eligibility, it is assumed that
- (e-2a)  $\xi \perp\!\!\!\perp (U_0, U_1) \mid X, Z, Q, D = 1$

and

- (e-2b)  $\Pr(A = 1 \mid X, Z, D = 1, \xi = 1) = 1;$   
 $\Pr(A = 1 \mid X, Z, D = 1, \xi = 0) = 0.$



- Agent failure to comply with the eligibility rules or protocols of experiments can lead to violations of (e-1) and/or (e-2).

- Agent failure to comply with the eligibility rules or protocols of experiments can lead to violations of (e-1) and/or (e-2).
- An equivalent way to formulate (e-2a) and (e-2b) uses the Imbens-Angrist notation for IV:

- (e-2a)'  $\xi \perp\!\!\!\perp (Y_0, Y_1) \mid X, Z, Q, D = 1$

and

- (e-2b)'  $\Pr(A = 1 \mid X, Z, D = 1, \xi = 1) = 1;$   
 $\Pr(A = 1 \mid X, Z, D = 1, \xi = 0) = 0.$

- Both randomizations are instruments as defined in Slide 152.

- Both randomizations are instruments as defined in Slide 152.
- Under the stated conditions, both satisfy (IV-1) and (IV-2), suitably redefined for eligibility randomizations, replacing  $D$  by  $A$ .

- A variety of conditioning variables is permitted by these definitions.

- A variety of conditioning variables is permitted by these definitions.
- Thus, (e-1) and (e-2) allow for the possibility that the conventional instruments  $Z$  fail (IV-1) and (IV-2), but nonetheless the randomization generates a valid instrument  $\xi$ .

- A variety of conditioning variables is permitted by these definitions.
- Thus, (e-1) and (e-2) allow for the possibility that the conventional instruments  $Z$  fail (IV-1) and (IV-2), but nonetheless the randomization generates a valid instrument  $\xi$ .
- The simplest randomizations do not condition on any variables.



- A variety of conditioning variables is permitted by these definitions.
- Thus, (e-1) and (e-2) allow for the possibility that the conventional instruments  $Z$  fail (IV-1) and (IV-2), but nonetheless the randomization generates a valid instrument  $\xi$ .
- The simplest randomizations do not condition on any variables.
- We next consider what these instruments identify.

## What Does Randomization Identify?

- Under invariance assumption (PI-3) and under one set of randomization assumptions just presented, IV is an instrument that identifies some treatment effect for an ongoing program.

## What Does Randomization Identify?

- Under invariance assumption (PI-3) and under one set of randomization assumptions just presented, IV is an instrument that identifies some treatment effect for an ongoing program.
- The question is: which treatment effect?

## What Does Randomization Identify?

- Under invariance assumption (PI-3) and under one set of randomization assumptions just presented, IV is an instrument that identifies some treatment effect for an ongoing program.
- The question is: which treatment effect?
- Following our discussion of IV with essential heterogeneity presented in Slide 152, different randomizations (or instruments) identify different parameters unless there is a common coefficient model ( $Y_1 - Y_0 = \beta(X)$  is the same for everyone given  $X$ ) or unless there is no dependence between the treatment effect ( $Y_1 - Y_0$ ) and the indicator  $D$  of the agents' desire to participate in the treatment.

## What Does Randomization Identify?

- Under invariance assumption (PI-3) and under one set of randomization assumptions just presented, IV is an instrument that identifies some treatment effect for an ongoing program.
- The question is: which treatment effect?
- Following our discussion of IV with essential heterogeneity presented in Slide 152, different randomizations (or instruments) identify different parameters unless there is a common coefficient model ( $Y_1 - Y_0 = \beta(X)$  is the same for everyone given  $X$ ) or unless there is no dependence between the treatment effect ( $Y_1 - Y_0$ ) and the indicator  $D$  of the agents' desire to participate in the treatment.
- In these two special cases, all mean treatment parameters are the same.

- Using IV, we can identify the marginal distributions  $F_0(y_0 | X)$  and  $F_1(y_1 | X)$ .

- Using IV, we can identify the marginal distributions  $F_0(y_0 | X)$  and  $F_1(y_1 | X)$ .
- In a model with essential heterogeneity, the instruments generated by randomization can identify parameters that are far from the parameters of economic interest.

- Using IV, we can identify the marginal distributions  $F_0(y_0 | X)$  and  $F_1(y_1 | X)$ .
- In a model with essential heterogeneity, the instruments generated by randomization can identify parameters that are far from the parameters of economic interest.
- Randomization of components of  $W$  (or  $Z$  given  $X$ ) under (R-4) and conditions (IV-1) and (IV-2) from Slide 12 produces instruments with the same problems and possibilities as analyzed in our discussion of instrumental variables.



- Using IV, we can identify the marginal distributions  $F_0(y_0 | X)$  and  $F_1(y_1 | X)$ .
- In a model with essential heterogeneity, the instruments generated by randomization can identify parameters that are far from the parameters of economic interest.
- Randomization of components of  $W$  (or  $Z$  given  $X$ ) under (R-4) and conditions (IV-1) and (IV-2) from Slide 12 produces instruments with the same problems and possibilities as analyzed in our discussion of instrumental variables.
- Using  $W$  as an instrument may lead to negative weights on the underlying LATEs or MTEs.

- Thus, unless we condition on the other instruments, the IV defined by randomization can be negative even if all of the underlying treatment effects or LATEs and MTEs generating choice behavior are positive.

- Thus, unless we condition on the other instruments, the IV defined by randomization can be negative even if all of the underlying treatment effects or LATEs and MTEs generating choice behavior are positive.
- The weighted average of the MTE generated by the instrument may be far from the policy relevant treatment effect.

- Thus, unless we condition on the other instruments, the IV defined by randomization can be negative even if all of the underlying treatment effects or LATEs and MTEs generating choice behavior are positive.
- The weighted average of the MTE generated by the instrument may be far from the policy relevant treatment effect.
- Under (PI-3) and (e-1), or equivalently (e-1)', the first type of eligibility randomization identifies  $\Pr(D = 1 | X, Z)$  (the choice probability) and hence relative subjective evaluations, and the marginal outcome distributions  $F_0(y_0 | X, D = 0)$  and  $F_1(y_1 | X, D = 1)$  for the eligible population ( $\xi = 1$ ).

- Thus, unless we condition on the other instruments, the IV defined by randomization can be negative even if all of the underlying treatment effects or LATEs and MTEs generating choice behavior are positive.
- The weighted average of the MTE generated by the instrument may be far from the policy relevant treatment effect.
- Under (PI-3) and (e-1), or equivalently (e-1)', the first type of eligibility randomization identifies  $\Pr(D = 1 | X, Z)$  (the choice probability) and hence relative subjective evaluations, and the marginal outcome distributions  $F_0(y_0 | X, D = 0)$  and  $F_1(y_1 | X, D = 1)$  for the eligible population ( $\xi = 1$ ).
- Agents made eligible for the program self-select as usual.

- For those deemed ineligible ( $\xi = 0$ ), under our assumptions, we would identify the distribution of  $Y_0$ , which can be partitioned into components for those who would have participated in the program had it not been for the randomization and components for those who would not have participated if offered the opportunity to do so:

$$F_0(y_0 | X) = F_0(y_0 | X, D = 0) \Pr(D = 0 | X) + F_0(y_0 | X, D = 1) \Pr(D = 1 | X).$$

- For those deemed ineligible ( $\xi = 0$ ), under our assumptions, we would identify the distribution of  $Y_0$ , which can be partitioned into components for those who would have participated in the program had it not been for the randomization and components for those who would not have participated if offered the opportunity to do so:

$$F_0(y_0 | X) = F_0(y_0 | X, D = 0) \Pr(D = 0 | X) + F_0(y_0 | X, D = 1) \Pr(D = 1 | X).$$

- Since we know  $F_0(y_0 | X, D = 0)$  and  $\Pr(D = 1 | X)$  from the eligible population, we can identify  $F_0(y_0 | X, D = 1)$ .

- For those deemed ineligible ( $\xi = 0$ ), under our assumptions, we would identify the distribution of  $Y_0$ , which can be partitioned into components for those who would have participated in the program had it not been for the randomization and components for those who would not have participated if offered the opportunity to do so:

$$F_0(y_0 | X) = F_0(y_0 | X, D = 0) \Pr(D = 0 | X) + F_0(y_0 | X, D = 1) \Pr(D = 1 | X).$$

- Since we know  $F_0(y_0 | X, D = 0)$  and  $\Pr(D = 1 | X)$  from the eligible population, we can identify  $F_0(y_0 | X, D = 1)$ .
- This is the new piece of information produced by the randomization compared to what can be obtained from standard observational data.



- In particular, we can identify the parameter  $\tau_T$ ,  $E(Y_1 - Y_0 | X, D = 1)$ , but without further assumptions, we cannot identify the other treatment parameters ATE ( $= E(Y_1 - Y_0 | X)$ ) or the joint distributions  $F(y_0, y_1 | X)$  or  $F(y_0, y_1 | X, D = 1)$ .

- In particular, we can identify the parameter  $\tau$ ,  $E(Y_1 - Y_0 | X, D = 1)$ , but without further assumptions, we cannot identify the other treatment parameters ATE ( $= E(Y_1 - Y_0 | X)$ ) or the joint distributions  $F(y_0, y_1 | X)$  or  $F(y_0, y_1 | X, D = 1)$ .
- To show that  $\xi$  is a valid instrument for  $A$ , form the Wald estimand,

$$IV_{(e-1)} = \frac{E(Y | \xi = 1, Z, X) - E(Y | \xi = 0, Z, X)}{\Pr(A = 1 | \xi = 1, Z, X) - \Pr(A = 1 | \xi = 0, Z, X)}. \quad (72)$$

Under invariance assumption (PI-3),  $\Pr(D = 1 | Z, X)$  is the same in the presence or absence of randomization.

- Assuming full compliance so that agents randomized to ineligibility do not show up in the program,

$$\Pr(A = 1 \mid \xi = 0, Z, X) = 0,$$

and

$$\begin{aligned} E(Y \mid \xi = 0, Z, X) &= E(Y_0 \mid Z, X) \\ &= E(Y_0 \mid D = 1, X, Z) \Pr(D = 1 \mid X, Z) \\ &\quad + E(Y_0 \mid D = 0, X, Z) \Pr(D = 0 \mid X, Z). \end{aligned}$$

If  $Z$  also satisfies the requirement (IV-1) that it is an instrument, then  $E(Y_0 \mid Z, X) = E(Y_0 \mid X)$ .

- Assuming full compliance so that agents randomized to ineligibility do not show up in the program,

$$\Pr(A = 1 \mid \xi = 0, Z, X) = 0,$$

and

$$\begin{aligned} E(Y \mid \xi = 0, Z, X) &= E(Y_0 \mid Z, X) \\ &= E(Y_0 \mid D = 1, X, Z) \Pr(D = 1 \mid X, Z) \\ &\quad + E(Y_0 \mid D = 0, X, Z) \Pr(D = 0 \mid X, Z). \end{aligned}$$

If  $Z$  also satisfies the requirement (IV-1) that it is an instrument, then  $E(Y_0 \mid Z, X) = E(Y_0 \mid X)$ .

- Under (e-1) or (e-1)' we do not have to assume that  $Z$  is a valid instrument.

- Using (e-1) and assumption (PI-3), the first term in the numerator of (72) can be written as

$$E(Y | \xi = 1, Z, X) = E(Y_1 | D = 1, Z, X) \Pr(D = 1 | Z, X) + E(Y_0 | D = 0, Z, X) \Pr(D = 0 | Z, X).$$

- Using (e-1) and assumption (PI-3), the first term in the numerator of (72) can be written as

$$E(Y | \xi = 1, Z, X) = E(Y_1 | D = 1, Z, X) \Pr(D = 1 | Z, X) + E(Y_0 | D = 0, Z, X) \Pr(D = 0 | Z, X).$$

- Substituting this expression into the numerator of equation (72) and collecting terms,  $IV_{(e-1)}$  identifies the parameter treatment on the treated:

$$IV_{(e-1)} = E(Y_1 - Y_0 | D = 1, Z, X).$$

- Using (e-1) and assumption (PI-3), the first term in the numerator of (72) can be written as

$$E(Y | \xi = 1, Z, X) = E(Y_1 | D = 1, Z, X) \Pr(D = 1 | Z, X) + E(Y_0 | D = 0, Z, X) \Pr(D = 0 | Z, X).$$

- Substituting this expression into the numerator of equation (72) and collecting terms,  $IV_{(e-1)}$  identifies the parameter treatment on the treated:

$$IV_{(e-1)} = E(Y_1 - Y_0 | D = 1, Z, X).$$

- It does not identify the other mean treatment effects, such as LATE or the average treatment effect ATE, unless the common coefficient model governs the data or  $(Y_1 - Y_0)$  is mean independent of  $D$ .

- Using the result that  $F(y | X) = E(\mathbf{1}(Y \leq y) | X)$ ,  $IV_{(e-1)}$  also identifies  $F_0(y_0 | X, D = 1)$ , since we can compute conditional means of  $\mathbf{1}(Y \leq y)$  for all  $y$ .



- Using the result that  $F(y | X) = E(\mathbf{1}(Y \leq y) | X)$ ,  $IV_{(e=1)}$  also identifies  $F_0(y_0 | X, D = 1)$ , since we can compute conditional means of  $\mathbf{1}(Y \leq y)$  for all  $y$ .
- The distribution  $F_1(y_1 | X, D = 1)$  can be identified from observational data.

- Using the result that  $F(y | X) = E(\mathbf{1}(Y \leq y) | X)$ ,  $IV_{(e=1)}$  also identifies  $F_0(y_0 | X, D = 1)$ , since we can compute conditional means of  $\mathbf{1}(Y \leq y)$  for all  $y$ .
- The distribution  $F_1(y_1 | X, D = 1)$  can be identified from observational data.
- Thus we can identify the outcome distributions for  $Y_0$  and for  $Y_1$  separately, conditional on  $D = 1, X, Z$ , but without additional assumptions we cannot identify the joint distribution of outcomes or the other treatment parameters.

- Randomization not conditional on  $(X, Z)$  ( $\xi \perp\!\!\!\perp (X, Z)$ ) creates an instrument  $\xi$  that satisfies the monotonicity or uniformity conditions.

- Randomization not conditional on  $(X, Z)$  ( $\xi \perp\!\!\!\perp (X, Z)$ ) creates an instrument  $\xi$  that satisfies the monotonicity or uniformity conditions.
- If the randomization is performed on  $(X, Z)$  strata, then the IV must be used conditional on the strata variables to ensure monotonicity is satisfied.

- The second type of eligibility randomization proceeds conditionally on  $D = 1$ .

- The second type of eligibility randomization proceeds conditionally on  $D = 1$ .
- Accordingly, data generated from such experiments do not identify choice probabilities ( $\Pr(D = 1 | X, Z)$ ) and hence do not identify the subjective evaluations of agents (??).

- The second type of eligibility randomization proceeds conditionally on  $D = 1$ .
- Accordingly, data generated from such experiments do not identify choice probabilities ( $\Pr(D = 1 | X, Z)$ ) and hence do not identify the subjective evaluations of agents (??).
- Under (PI-3) and (e-2) (or equivalent conditions (e-2)') randomization identifies  $F_0(y_0 | D = 1, X, Z)$  from the data on the randomized-out participants.

- The second type of eligibility randomization proceeds conditionally on  $D = 1$ .
- Accordingly, data generated from such experiments do not identify choice probabilities ( $\Pr(D = 1 | X, Z)$ ) and hence do not identify the subjective evaluations of agents (??).
- Under (PI-3) and (e-2) (or equivalent conditions (e-2)') randomization identifies  $F_0(y_0 | D = 1, X, Z)$  from the data on the randomized-out participants.
- This conditional distribution cannot be constructed from ordinary observational data unless additional assumptions are invoked.



- The second type of eligibility randomization proceeds conditionally on  $D = 1$ .
- Accordingly, data generated from such experiments do not identify choice probabilities ( $\Pr(D = 1 | X, Z)$ ) and hence do not identify the subjective evaluations of agents (??).
- Under (PI-3) and (e-2) (or equivalent conditions (e-2)') randomization identifies  $F_0(y_0 | D = 1, X, Z)$  from the data on the randomized-out participants.
- This conditional distribution cannot be constructed from ordinary observational data unless additional assumptions are invoked.
- From the data for the eligible ( $\xi = 1$ ) population, we identify  $F_1(y_1 | D = 1, X, Z)$ .

- The Wald estimator for mean outcomes in this case is

$$IV_{(e-2)} = \frac{E(Y | D = 1, \xi = 1, X, Z) - E(Y | D = 1, \xi = 0, X, Z)}{\Pr(A = 1 | D = 1, \xi = 1, X, Z) - \Pr(A = 1 | D = 1, \xi = 0, X, Z)}.$$

- The Wald estimator for mean outcomes in this case is

$$IV_{(e-2)} = \frac{E(Y | D = 1, \xi = 1, X, Z) - E(Y | D = 1, \xi = 0, X, Z)}{\Pr(A = 1 | D = 1, \xi = 1, X, Z) - \Pr(A = 1 | D = 1, \xi = 0, X, Z)}.$$

- Under (e-2)/(e-2)',

$$\Pr(A = 1 | D = 1, \xi = 1, X, Z) = 1,$$

$$\Pr(A = 1 | D = 1, \xi = 0, X, Z) = 0,$$

$$E(Y | A = 0, D = 1, \xi = 0, X, Z) = E(Y_0 | D = 1, X, Z), \text{ and}$$

$$E(Y | A = 1, D = 1, \xi = 1, X, Z) = E(Y_1 | D = 1, X, Z).$$

- The Wald estimator for mean outcomes in this case is

$$IV_{(e-2)} = \frac{E(Y | D = 1, \xi = 1, X, Z) - E(Y | D = 1, \xi = 0, X, Z)}{\Pr(A = 1 | D = 1, \xi = 1, X, Z) - \Pr(A = 1 | D = 1, \xi = 0, X, Z)}.$$

- Under (e-2)/(e-2)',

$$\Pr(A = 1 | D = 1, \xi = 1, X, Z) = 1,$$

$$\Pr(A = 1 | D = 1, \xi = 0, X, Z) = 0,$$

$$E(Y | A = 0, D = 1, \xi = 0, X, Z) = E(Y_0 | D = 1, X, Z), \text{ and}$$

$$E(Y | A = 1, D = 1, \xi = 1, X, Z) = E(Y_1 | D = 1, X, Z).$$

- Thus,

$$IV_{(e-2)} = E(Y_1 - Y_0 | D = 1, X, Z).$$

- In the general model with essential heterogeneity, randomized trials with full compliance that do not disturb the activity being evaluated answer a limited set of questions, and do not in general identify the policy relevant treatment effect (PRTE).

- In the general model with essential heterogeneity, randomized trials with full compliance that do not disturb the activity being evaluated answer a limited set of questions, and do not in general identify the policy relevant treatment effect (PRTE).
- Randomizations have to be carefully chosen to make sure that they answer interesting economic questions.

- In the general model with essential heterogeneity, randomized trials with full compliance that do not disturb the activity being evaluated answer a limited set of questions, and do not in general identify the policy relevant treatment effect (PRTE).
- Randomizations have to be carefully chosen to make sure that they answer interesting economic questions.
- Their analysis has to be supplemented with the methods previously analyzed to answer the full range of policy questions addressed there.

- Thus far we have assumed that the randomizations do not violate the invariance assumption (PI-3).



- Thus far we have assumed that the randomizations do not violate the invariance assumption (PI-3).
- Yet many randomizations alter the environment they are studying and inject what may be unwelcome sources of uncertainty into agent decision making.

- Thus far we have assumed that the randomizations do not violate the invariance assumption (PI-3).
- Yet many randomizations alter the environment they are studying and inject what may be unwelcome sources of uncertainty into agent decision making.
- We now examine the consequences of violations of invariance.

## Randomization Bias

- If randomization alters the program being evaluated, the outcomes of a randomized trial may bear little resemblance to the outcomes generated by an ongoing version of the program that has not been subject to randomization.

## Randomization Bias

- If randomization alters the program being evaluated, the outcomes of a randomized trial may bear little resemblance to the outcomes generated by an ongoing version of the program that has not been subject to randomization.
- In this case, assumption (PI-3) is violated.

## Randomization Bias

- If randomization alters the program being evaluated, the outcomes of a randomized trial may bear little resemblance to the outcomes generated by an ongoing version of the program that has not been subject to randomization.
- In this case, assumption (PI-3) is violated.
- Such violations are termed “Hawthorne effects” and are called “Randomization Bias” in the economics literature.

## Randomization Bias

- If randomization alters the program being evaluated, the outcomes of a randomized trial may bear little resemblance to the outcomes generated by an ongoing version of the program that has not been subject to randomization.
- In this case, assumption (PI-3) is violated.
- Such violations are termed “Hawthorne effects” and are called “Randomization Bias” in the economics literature.
- The process of randomization may affect objective outcomes, subjective outcomes or both.

- Even if (PI-3) is violated, randomization may still be a valid instrument for the altered program.

- Even if (PI-3) is violated, randomization may still be a valid instrument for the altered program.
- Although the program studied may be changed, under the assumptions made in Slide 856, randomization can produce “internally valid” treatment effects for the altered program.



- Even if (PI-3) is violated, randomization may still be a valid instrument for the altered program.
- Although the program studied may be changed, under the assumptions made in Slide 856, randomization can produce “internally valid” treatment effects for the altered program.
- Thus randomization can answer policy question (P-1) for a program changed by randomization, but not for the program as it would operate in the absence of randomization.

- As noted repeatedly, a distinctive feature of the econometric approach to social program evaluation is its emphasis on choice and agent subjective evaluations of programs.

- As noted repeatedly, a distinctive feature of the econometric approach to social program evaluation is its emphasis on choice and agent subjective evaluations of programs.
- This feature accounts for the distinction between the statistician's invariance assumption (PI-1) and the economist's invariance assumption (PI-3).

- As noted repeatedly, a distinctive feature of the econometric approach to social program evaluation is its emphasis on choice and agent subjective evaluations of programs.
- This feature accounts for the distinction between the statistician's invariance assumption (PI-1) and the economist's invariance assumption (PI-3).
- (These are presented in Part I.) It is instructive to consider the case where assumption (PI-1) is valid but assumption (PI-3) is not.

- As noted repeatedly, a distinctive feature of the econometric approach to social program evaluation is its emphasis on choice and agent subjective evaluations of programs.
- This feature accounts for the distinction between the statistician's invariance assumption (PI-1) and the economist's invariance assumption (PI-3).
- (These are presented in Part I.) It is instructive to consider the case where assumption (PI-1) is valid but assumption (PI-3) is not.
- This case might arise when randomization alters risk-averse agent decision behavior but has no effects on potential outcomes.

- As noted repeatedly, a distinctive feature of the econometric approach to social program evaluation is its emphasis on choice and agent subjective evaluations of programs.
- This feature accounts for the distinction between the statistician's invariance assumption (PI-1) and the economist's invariance assumption (PI-3).
- (These are presented in Part I.) It is instructive to consider the case where assumption (PI-1) is valid but assumption (PI-3) is not.
- This case might arise when randomization alters risk-averse agent decision behavior but has no effects on potential outcomes.
- Thus the  $R(s, \omega)$  are affected, but not the  $Y(s, \omega)$ .

- In this case, the parameter  $ATE(X) = E(Y_1 - Y_0 | X)$  is the same in the ongoing program as in the population generated by the randomized trial.

- In this case, the parameter  $ATE(X) = E(Y_1 - Y_0 | X)$  is the same in the ongoing program as in the population generated by the randomized trial.
- However, treatment parameters conditional on choices such as  $TT(X) = E(Y_1 - Y_0 | X, D = 1)$ ,  $TUT(X) = E(Y_1 - Y_0 | X, D = 0)$  are not, in general, invariant.



- In this case, the parameter  $ATE(X) = E(Y_1 - Y_0 | X)$  is the same in the ongoing program as in the population generated by the randomized trial.
- However, treatment parameters conditional on choices such as  $TT(X) = E(Y_1 - Y_0 | X, D = 1)$ ,  $TUT(X) = E(Y_1 - Y_0 | X, D = 0)$  are not, in general, invariant.
- If the subjective valuations are altered, so are the parameters based on choices produced by the subjective valuations.

- In this case, the parameter  $ATE(X) = E(Y_1 - Y_0 | X)$  is the same in the ongoing program as in the population generated by the randomized trial.
- However, treatment parameters conditional on choices such as  $TT(X) = E(Y_1 - Y_0 | X, D = 1)$ ,  $TUT(X) = E(Y_1 - Y_0 | X, D = 0)$  are not, in general, invariant.
- If the subjective valuations are altered, so are the parameters based on choices produced by the subjective valuations.
- Different random variables generate the conditioning sets in the randomized and nonrandomized regimes and, in general, they will have a different dependence structure with the outcomes  $Y(s, \omega)$ .

- This arises because randomization alters the composition of participants in the conditioning set that defines the treatment parameter.

- This arises because randomization alters the composition of participants in the conditioning set that defines the treatment parameter.
- This analysis applies with full force to LATE.

- This arises because randomization alters the composition of participants in the conditioning set that defines the treatment parameter.
- This analysis applies with full force to LATE.
- LATE based on  $P(Z)$  for two distinct values of  $Z$  ( $Z = z$  and  $Z = z'$ ) is  $E(Y_1 - Y_0 \mid X, P(z') \leq U_D \leq P(z))$ .

- This arises because randomization alters the composition of participants in the conditioning set that defines the treatment parameter.
- This analysis applies with full force to LATE.
- LATE based on  $P(Z)$  for two distinct values of  $Z$  ( $Z = z$  and  $Z = z'$ ) is  $E(Y_1 - Y_0 | X, P(z') \leq U_D \leq P(z))$ .
- In the randomized trial, violation of (PI-3) because of lack of invariance of  $R(s, \omega)$  changes  $U_D$  and the values of  $P(Z)$  for the same  $Z = z$ .

- This arises because randomization alters the composition of participants in the conditioning set that defines the treatment parameter.
- This analysis applies with full force to LATE.
- LATE based on  $P(Z)$  for two distinct values of  $Z$  ( $Z = z$  and  $Z = z'$ ) is  $E(Y_1 - Y_0 | X, P(z') \leq U_D \leq P(z))$ .
- In the randomized trial, violation of (PI-3) because of lack of invariance of  $R(s, \omega)$  changes  $U_D$  and the values of  $P(Z)$  for the same  $Z = z$ .
- In general, this alters LATE.

- The case where (PI-1) holds, but (PI-3) does not, generates invariant conditional (on choice) parameters if there is no treatment effect heterogeneity or if there is such heterogeneity that is independent of  $D$ .



- The case where (PI-1) holds, but (PI-3) does not, generates invariant conditional (on choice) parameters if there is no treatment effect heterogeneity or if there is such heterogeneity that is independent of  $D$ .
- These are the familiar conditions: (a)  $Y_1 - Y_0$  is the same for all people with the same  $X = x$  or (b)  $Y_1 - Y_0$  is (mean) independent of  $D$  given  $X = x$ .

- The case where (PI-1) holds, but (PI-3) does not, generates invariant conditional (on choice) parameters if there is no treatment effect heterogeneity or if there is such heterogeneity that is independent of  $D$ .
- These are the familiar conditions: (a)  $Y_1 - Y_0$  is the same for all people with the same  $X = x$  or (b)  $Y_1 - Y_0$  is (mean) independent of  $D$  given  $X = x$ .
- In these cases, the MTE is flat in  $U_D$ .

- In general, in a model with essential heterogeneity, even if the Rubin invariance conditions (PI-1) and (PI-2) are satisfied, but conditions (PI-3) and (PI-4) are not, treatment parameters defined conditional on choices are not invariant to the choice of randomization.

- In general, in a model with essential heterogeneity, even if the Rubin invariance conditions (PI-1) and (PI-2) are satisfied, but conditions (PI-3) and (PI-4) are not, treatment parameters defined conditional on choices are not invariant to the choice of randomization.
- This insight shows the gain in clarity in interpreting what experiments identify from adopting a choice-theoretic, econometric approach to the evaluation of social programs, as opposed to the conventional approach adopted by statisticians.

- In general, in a model with essential heterogeneity, even if the Rubin invariance conditions (PI-1) and (PI-2) are satisfied, but conditions (PI-3) and (PI-4) are not, treatment parameters defined conditional on choices are not invariant to the choice of randomization.
- This insight shows the gain in clarity in interpreting what experiments identify from adopting a choice-theoretic, econometric approach to the evaluation of social programs, as opposed to the conventional approach adopted by statisticians.
- We now show another advantage of the economic approach in an analysis of noncompliance and its implications for interpreting experimental evidence.

## Compliance

- The statistical treatment effect literature extensively analyzes the problem of noncompliance.

## Compliance

- The statistical treatment effect literature extensively analyzes the problem of noncompliance.
- Persons assigned to a treatment may not accept it.

## Compliance

- The statistical treatment effect literature extensively analyzes the problem of noncompliance.
- Persons assigned to a treatment may not accept it.
- In the notation of equation (72), let  $\xi = 1$  if a person is assigned to treatment,  $\xi = 0$  otherwise.



## Compliance

- The statistical treatment effect literature extensively analyzes the problem of noncompliance.
- Persons assigned to a treatment may not accept it.
- In the notation of equation (72), let  $\xi = 1$  if a person is assigned to treatment,  $\xi = 0$  otherwise.
- Compliance is said to be perfect when  $\xi = 1 \Rightarrow A = 1$  and  $\xi = 0 \Rightarrow A = 0$ .

## Compliance

- The statistical treatment effect literature extensively analyzes the problem of noncompliance.
- Persons assigned to a treatment may not accept it.
- In the notation of equation (72), let  $\xi = 1$  if a person is assigned to treatment,  $\xi = 0$  otherwise.
- Compliance is said to be perfect when  $\xi = 1 \Rightarrow A = 1$  and  $\xi = 0 \Rightarrow A = 0$ .
- In the presence of self selection by agents, these conditions do not, in general, hold.

## Compliance

- The statistical treatment effect literature extensively analyzes the problem of noncompliance.
- Persons assigned to a treatment may not accept it.
- In the notation of equation (72), let  $\xi = 1$  if a person is assigned to treatment,  $\xi = 0$  otherwise.
- Compliance is said to be perfect when  $\xi = 1 \Rightarrow A = 1$  and  $\xi = 0 \Rightarrow A = 0$ .
- In the presence of self selection by agents, these conditions do not, in general, hold.
- People assigned to treatment may not comply ( $\xi = 1$  but  $D = 0$ ).

- This is also called the “dropout” problem ( ?; ?).

- This is also called the “dropout” problem ( ?; ?).
- In its formulation of this problem, the literature assumes that outcomes are measured for each participant but that outcomes realized are not always those intended by the randomizers.

- This is also called the “dropout” problem ( ?; ?).
- In its formulation of this problem, the literature assumes that outcomes are measured for each participant but that outcomes realized are not always those intended by the randomizers.
- In addition, people denied treatment may find substitutes for the treatment outside of the program.

- This is also called the “dropout” problem ( ?; ?).
- In its formulation of this problem, the literature assumes that outcomes are measured for each participant but that outcomes realized are not always those intended by the randomizers.
- In addition, people denied treatment may find substitutes for the treatment outside of the program.
- This is the problem of substitution bias.

- This is also called the “dropout” problem ( ?; ?).
- In its formulation of this problem, the literature assumes that outcomes are measured for each participant but that outcomes realized are not always those intended by the randomizers.
- In addition, people denied treatment may find substitutes for the treatment outside of the program.
- This is the problem of substitution bias.
- Since self-selection is an integral part of choice models, noncompliance, as the term is used by the statisticians, is a feature of most social experiments.



- The econometric approach builds in the possibility of self-selection as an integral part of model specification.

- The econometric approach builds in the possibility of self-selection as an integral part of model specification.
- As emphasized in the econometric literature since the work of Heckman, Ichimura, and Todd, agent decisions to participate are informative about their subjective evaluations of the program.

- The econometric approach builds in the possibility of self-selection as an integral part of model specification.
- As emphasized in the econometric literature since the work of Heckman, Ichimura, and Todd, agent decisions to participate are informative about their subjective evaluations of the program.
- In the dynamic setting discussed in section 3 of Part III, agent decisions to attrite from a program are informative about their update of information about the program ( Heckman; Ichimura; Todd; and Todd).

- The econometric approach builds in the possibility of self-selection as an integral part of model specification.
- As emphasized in the econometric literature since the work of Heckman, Ichimura, and Todd, agent decisions to participate are informative about their subjective evaluations of the program.
- In the dynamic setting discussed in section 3 of Part III, agent decisions to attrite from a program are informative about their update of information about the program ( Heckman; Heckman; Heckman; and Heckman ).
- Noncompliance is a source of information about subjective evaluations of programs.

- Noncompliance is a problem if the goal of the social experiment is to estimate  $ATE(X) = E(Y_1 - Y_0 | X)$  without using the econometric methods previously discussed.

- Noncompliance is a problem if the goal of the social experiment is to estimate  $ATE(X) = E(Y_1 - Y_0 | X)$  without using the econometric methods previously discussed.
- We established in Slide 869 that in the presence of self-selection, in a general case with essential heterogeneity, experiments under assumptions (PI-3) and (PI-4) and (e-1) or (e-2) identify  $E(Y_1 - Y_0 | X, D = 1)$  instead of  $ATE(X)$ .

- Concerns about noncompliance often arise from adoption of the Neyman-Cox-Rubin “causal model” discussed in Part I, section 4.4.

- Concerns about noncompliance often arise from adoption of the Neyman-Cox-Rubin “causal model” discussed in Part I, section 4.4.
- Experiments are conceived as tools for direct allocation of agricultural treatments.



- Concerns about noncompliance often arise from adoption of the Neyman-Cox-Rubin “causal model” discussed in Part I, section 4.4.
- Experiments are conceived as tools for direct allocation of agricultural treatments.
- For that reason, that literature elevates ATE to pre-eminence as the parameter of interest because it is thought that this parameter can be produced by experiments.

- Concerns about noncompliance often arise from adoption of the Neyman-Cox-Rubin “causal model” discussed in Part I, section 4.4.
- Experiments are conceived as tools for direct allocation of agricultural treatments.
- For that reason, that literature elevates ATE to pre-eminence as the parameter of interest because it is thought that this parameter can be produced by experiments.
- In social experiments, it is rare that the experimenter can force anyone to do anything.

- Concerns about noncompliance often arise from adoption of the Neyman-Cox-Rubin “causal model” discussed in Part I, section 4.4.
- Experiments are conceived as tools for direct allocation of agricultural treatments.
- For that reason, that literature elevates ATE to pre-eminence as the parameter of interest because it is thought that this parameter can be produced by experiments.
- In social experiments, it is rare that the experimenter can force anyone to do anything.
- As the old adage goes, “you can lead a horse to water but you cannot make it drink.” Agent choice behavior intervenes.

- Thus it is no accident that if they are not compromised, the two randomizations most commonly implemented directly identify parameters conditional on choices.

- Thus it is no accident that if they are not compromised, the two randomizations most commonly implemented directly identify parameters conditional on choices.
- There is a more general version of the noncompliance problem which requires a dynamic formulation.

- Thus it is no accident that if they are not compromised, the two randomizations most commonly implemented directly identify parameters conditional on choices.
- There is a more general version of the noncompliance problem which requires a dynamic formulation.
- Agents are assigned to treatment ( $\xi = 1$ ) and some accept ( $D = 1$ ) but drop out of the program at a later stage.

- Thus it is no accident that if they are not compromised, the two randomizations most commonly implemented directly identify parameters conditional on choices.
- There is a more general version of the noncompliance problem which requires a dynamic formulation.
- Agents are assigned to treatment ( $\xi = 1$ ) and some accept ( $D = 1$ ) but drop out of the program at a later stage.
- We need to modify the formulation in this section to cover this case.

- Thus it is no accident that if they are not compromised, the two randomizations most commonly implemented directly identify parameters conditional on choices.
- There is a more general version of the noncompliance problem which requires a dynamic formulation.
- Agents are assigned to treatment ( $\xi = 1$ ) and some accept ( $D = 1$ ) but drop out of the program at a later stage.
- We need to modify the formulation in this section to cover this case.
- We now turn to that modification.



## The Dynamics of Dropout and Program Participation

- Actual programs are more dynamic in character than the stylized program just analyzed.

## The Dynamics of Dropout and Program Participation

- Actual programs are more dynamic in character than the stylized program just analyzed.
- Multiple actors are involved, such as the agents being studied and the groups administering the programs.

## The Dynamics of Dropout and Program Participation

- Actual programs are more dynamic in character than the stylized program just analyzed.
- Multiple actors are involved, such as the agents being studied and the groups administering the programs.
- People apply, are accepted, enroll, and complete the program.

## The Dynamics of Dropout and Program Participation

- Actual programs are more dynamic in character than the stylized program just analyzed.
- Multiple actors are involved, such as the agents being studied and the groups administering the programs.
- People apply, are accepted, enroll, and complete the program.
- A fully dynamic analysis, along the lines of the models developed in Part III of our contribution to this Handbook, analyzes each of these decisions, accounting for the updating of agent and program administrators' information.

## The Dynamics of Dropout and Program Participation

- Actual programs are more dynamic in character than the stylized program just analyzed.
- Multiple actors are involved, such as the agents being studied and the groups administering the programs.
- People apply, are accepted, enroll, and complete the program.
- A fully dynamic analysis, along the lines of the models developed in Part III of our contribution to this Handbook, analyzes each of these decisions, accounting for the updating of agent and program administrators' information.
- This section briefly discusses some new issues that arise in a more dynamic formulation of the dropout problem.

- ?, ?, ?, and ?? discuss these issues in greater depth.

- $\tau$ ,  $\tau$ ,  $\tau$ , and  $\tau\tau$  discuss these issues in greater depth.
- In this subsection, we analyze the effects of dropouts on inferences from social experiments and assume no attrition.

- $\chi^2$ ,  $\chi^2$ ,  $\chi^2$ , and  $\chi^2$  discuss these issues in greater depth.
- In this subsection, we analyze the effects of dropouts on inferences from social experiments and assume no attrition.
- Our analysis of this case is of interest both in its own right and as a demonstration of the power of our approach.



- Consider a stylized multiple stage program.

- Consider a stylized multiple stage program.
- In stage “0”, the agent (possibly in conjunction with program officials) decides to participate or not to participate in the program.

- Consider a stylized multiple stage program.
- In stage “0”, the agent (possibly in conjunction with program officials) decides to participate or not to participate in the program.
- This is an enrollment phase prior to treatment.

- Consider a stylized multiple stage program.
- In stage “0”, the agent (possibly in conjunction with program officials) decides to participate or not to participate in the program.
- This is an enrollment phase prior to treatment.
- Let  $D_0 = 1$  denote that the agent does not choose to participate.

- Consider a stylized multiple stage program.
- In stage “0”, the agent (possibly in conjunction with program officials) decides to participate or not to participate in the program.
- This is an enrollment phase prior to treatment.
- Let  $D_0 = 1$  denote that the agent does not choose to participate.
- $D_0 = 0$  denotes that the agent participates and receives some treatment among  $J$  possible program levels beyond the no treatment state.

- Consider a stylized multiple stage program.
- In stage “0”, the agent (possibly in conjunction with program officials) decides to participate or not to participate in the program.
- This is an enrollment phase prior to treatment.
- Let  $D_0 = 1$  denote that the agent does not choose to participate.
- $D_0 = 0$  denotes that the agent participates and receives some treatment among  $J$  possible program levels beyond the no treatment state.
- The outcome associated with state “0” is  $Y_0$ .

- Consider a stylized multiple stage program.
- In stage “0”, the agent (possibly in conjunction with program officials) decides to participate or not to participate in the program.
- This is an enrollment phase prior to treatment.
- Let  $D_0 = 1$  denote that the agent does not choose to participate.
- $D_0 = 0$  denotes that the agent participates and receives some treatment among  $J$  possible program levels beyond the no treatment state.
- The outcome associated with state “0” is  $Y_0$ .
- This assumes that acts of inquiry about a program or registration in it have no effect on outcomes.

- One could disaggregate stage “0” into recruitment, application, and acceptance stages, but for expositional simplicity we collapse these into one stage.



- One could disaggregate stage “0” into recruitment, application, and acceptance stages, but for expositional simplicity we collapse these into one stage.
- If the  $J$  possible treatment stages are ordered, say, by the intensity of treatment, then “1” is the least amount of treatment and “ $J$ ” is the greatest amount.

- One could disaggregate stage “0” into recruitment, application, and acceptance stages, but for expositional simplicity we collapse these into one stage.
- If the  $J$  possible treatment stages are ordered, say, by the intensity of treatment, then “1” is the least amount of treatment and “ $J$ ” is the greatest amount.
- A more general model would allow people to transit to stage  $j$  but not complete it.

- One could disaggregate stage “0” into recruitment, application, and acceptance stages, but for expositional simplicity we collapse these into one stage.
- If the  $J$  possible treatment stages are ordered, say, by the intensity of treatment, then “1” is the least amount of treatment and “ $J$ ” is the greatest amount.
- A more general model would allow people to transit to stage  $j$  but not complete it.
- The  $J$  distinct stages can be interpreted quite generally.

- One could disaggregate stage “0” into recruitment, application, and acceptance stages, but for expositional simplicity we collapse these into one stage.
- If the  $J$  possible treatment stages are ordered, say, by the intensity of treatment, then “1” is the least amount of treatment and “ $J$ ” is the greatest amount.
- A more general model would allow people to transit to stage  $j$  but not complete it.
- The  $J$  distinct stages can be interpreted quite generally.
- If a person no longer participates in the program after stage  $j$ ,  $j = 1, \dots, J$ , we set indicator  $D_j = 1$ .

- One could disaggregate stage “0” into recruitment, application, and acceptance stages, but for expositional simplicity we collapse these into one stage.
- If the  $J$  possible treatment stages are ordered, say, by the intensity of treatment, then “1” is the least amount of treatment and “ $J$ ” is the greatest amount.
- A more general model would allow people to transit to stage  $j$  but not complete it.
- The  $J$  distinct stages can be interpreted quite generally.
- If a person no longer participates in the program after stage  $j$ ,  $j = 1, \dots, J$ , we set indicator  $D_j = 1$ .
- The person is assumed to receive stage  $j$  treatment.

- $D_J = 1$  corresponds to completion of the program in all  $J$  stages of its treatment phase.

- $D_J = 1$  corresponds to completion of the program in all  $J$  stages of its treatment phase.
- Note that, by construction,  $\sum_{j=0}^J D_j = 1$ .

- $D_J = 1$  corresponds to completion of the program in all  $J$  stages of its treatment phase.
- Note that, by construction,  $\sum_{j=0}^J D_j = 1$ .
- The sequential updating model developed below in Part III can be used to formalize these decisions and their associated outcomes.



- $D_J = 1$  corresponds to completion of the program in all  $J$  stages of its treatment phase.
- Note that, by construction,  $\sum_{j=0}^J D_j = 1$ .
- The sequential updating model developed below in Part III can be used to formalize these decisions and their associated outcomes.
- We can also use the simple multinomial choice model developed and analyzed in appendix B of Part I.

- Let  $\{D_j(z)\}_{z \in \mathcal{Z}}$  be the set of potential treatment choices for choice  $j$  associated with setting  $Z = z$ .

- Let  $\{D_j(z)\}_{z \in \mathcal{Z}}$  be the set of potential treatment choices for choice  $j$  associated with setting  $Z = z$ .
- For each  $Z = z$ ,  $\sum_{j=0}^J D_j(z) = 1$ .

- Let  $\{D_j(z)\}_{z \in \mathcal{Z}}$  be the set of potential treatment choices for choice  $j$  associated with setting  $Z = z$ .
- For each  $Z = z$ ,  $\sum_{j=0}^J D_j(z) = 1$ .
- Using the methods exposted in Part III, we could update the information sets at each stage.

- Let  $\{D_j(z)\}_{z \in \mathcal{Z}}$  be the set of potential treatment choices for choice  $j$  associated with setting  $Z = z$ .
- For each  $Z = z$ ,  $\sum_{j=0}^J D_j(z) = 1$ .
- Using the methods exposted in Part III, we could update the information sets at each stage.
- We keep this updating implicit.

- Let  $\{D_j(z)\}_{z \in \mathcal{Z}}$  be the set of potential treatment choices for choice  $j$  associated with setting  $Z = z$ .
- For each  $Z = z$ ,  $\sum_{j=0}^J D_j(z) = 1$ .
- Using the methods exposted in Part III, we could update the information sets at each stage.
- We keep this updating implicit.
- Different components of  $Z$  may determine different choice indicators.

- Array the collections of choice indicators evaluated at each  $Z = z$  into a vector

$$D(z) = (\{D_1(z)\}_{z \in \mathcal{Z}}, \dots, \{D_J(z)\}_{z \in \mathcal{Z}}).$$

The potential outcomes associated with each of the  $J + 1$  states are

$$Y_j = \mu_j(X, U_j), \quad j = 0, \dots, J.$$

$Y_0$  is the no treatment state, and the  $Y_j, j \geq 1$ , correspond to outcomes associated with dropping out at various stages of the program.

- In the absence of randomization, the observed  $Y$  is

$$Y = \sum_{j=0}^J D_j Y_j,$$

the Roy-Quandt switching regime model.



- In the absence of randomization, the observed  $Y$  is

$$Y = \sum_{j=0}^J D_j Y_j,$$

the Roy-Quandt switching regime model.

- Let  $\tilde{Y} = (Y_0, \dots, Y_J)$  denote the vector of potential outcomes associated with all phases of the program.

- In the absence of randomization, the observed  $Y$  is

$$Y = \sum_{j=0}^J D_j Y_j,$$

the Roy-Quandt switching regime model.

- Let  $\tilde{Y} = (Y_0, \dots, Y_J)$  denote the vector of potential outcomes associated with all phases of the program.
- Through selection, the  $Y_j$  for persons with  $D_j = 1$  may be different from the  $Y_j$  for persons with  $D_j = 0$ .

- Appendix B of Part I gives conditions under which the distributions of the  $Y_j$  and the subjective evaluations  $R_j$ ,  $j = 0, \dots, J$ , that generate choices  $D_j$  are identified.

- Appendix B of Part I gives conditions under which the distributions of the  $Y_j$  and the subjective evaluations  $R_j$ ,  $j = 0, \dots, J$ , that generate choices  $D_j$  are identified.
- Using the tools for multiple outcome models developed in Slide 471, we can use IV and our extensions of IV to identify the treatment parameters discussed there.

- In this subsection, we consider what randomizations at various stages identify.

- In this subsection, we consider what randomizations at various stages identify.
- We assume that the randomizations do not disturb the program.

- In this subsection, we consider what randomizations at various stages identify.
- We assume that the randomizations do not disturb the program.
- Thus we invoke assumption (PI-3).

- In this subsection, we consider what randomizations at various stages identify.
- We assume that the randomizations do not disturb the program.
- Thus we invoke assumption (PI-3).
- Recall that we also assume absence of general equilibrium effects (PI-4).



- In this subsection, we consider what randomizations at various stages identify.
- We assume that the randomizations do not disturb the program.
- Thus we invoke assumption (PI-3).
- Recall that we also assume absence of general equilibrium effects (PI-4).
- Let  $\xi_j = 1$  denote whether the person is eligible to move beyond stage  $j$ .

- In this subsection, we consider what randomizations at various stages identify.
- We assume that the randomizations do not disturb the program.
- Thus we invoke assumption (PI-3).
- Recall that we also assume absence of general equilibrium effects (PI-4).
- Let  $\xi_j = 1$  denote whether the person is eligible to move beyond stage  $j$ .
- $\xi_j = 0$  means the person is randomized out of the program after completing stage  $j$ .

- A randomization at stage  $j$  with  $\xi_j = 1$  means the person is allowed to continue on to stage  $j + 1$ , although the agent may still choose not to.

- A randomization at stage  $j$  with  $\xi_j = 1$  means the person is allowed to continue on to stage  $j + 1$ , although the agent may still choose not to.
- We set  $\xi_J \equiv 1$  to simplify the notation.

- A randomization at stage  $j$  with  $\xi_j = 1$  means the person is allowed to continue on to stage  $j + 1$ , although the agent may still choose not to.
- We set  $\xi_J \equiv 1$  to simplify the notation.
- The  $\xi_j$  are ordered in a natural way:  $\xi_j = 1$  only if  $\xi_\ell = 1$ ,  $\ell = 0, \dots, j - 1$ .

- A randomization at stage  $j$  with  $\xi_j = 1$  means the person is allowed to continue on to stage  $j + 1$ , although the agent may still choose not to.
- We set  $\xi_J \equiv 1$  to simplify the notation.
- The  $\xi_j$  are ordered in a natural way:  $\xi_j = 1$  only if  $\xi_\ell = 1$ ,  $\ell = 0, \dots, j - 1$ .
- Array the  $\xi_j$  into a vector  $\xi$  and denote its support by  $\tilde{\xi}$ .

- Because of agent self-selection, a person who does not choose to participate at stage  $j$  cannot be forced to do so.

- Because of agent self-selection, a person who does not choose to participate at stage  $j$  cannot be forced to do so.
- For a person who would choose  $k$  ( $D_k = 1$ ) in a nonexperimental environment,  $Y_k$  is observed if  $\prod_{\ell=0}^{k-1} \xi_{\ell} = 1$ .



- Because of agent self-selection, a person who does not choose to participate at stage  $j$  cannot be forced to do so.
- For a person who would choose  $k$  ( $D_k = 1$ ) in a nonexperimental environment,  $Y_k$  is observed if  $\prod_{\ell=0}^{k-1} \xi_{\ell} = 1$ .
- Otherwise, if  $\xi_{k-1} = 0$  but, say,  $\prod_{\ell=0}^{k'-1} \xi_{\ell} = 1$  and  $\prod_{\ell=0}^{k'} \xi_{\ell} = 0$  for  $k' < k$ , we observe  $Y_{k'}$  for the agent.

- Because of agent self-selection, a person who does not choose to participate at stage  $j$  cannot be forced to do so.
- For a person who would choose  $k$  ( $D_k = 1$ ) in a nonexperimental environment,  $Y_k$  is observed if  $\prod_{\ell=0}^{k-1} \xi_\ell = 1$ .
- Otherwise, if  $\xi_{k-1} = 0$  but, say,  $\prod_{\ell=0}^{k'-1} \xi_\ell = 1$  and  $\prod_{\ell=0}^{k'} \xi_\ell = 0$  for  $k' < k$ , we observe  $Y_{k'}$  for the agent.
- From an experiment with randomization administered at different stages, we observe

$$Y = \sum_{j=0}^J D_j \left( \sum_{k=0}^j \left( \prod_{\ell=0}^{k-1} \xi_\ell \right) (1 - \xi_k) Y_k \right).$$

To understand this formula, consider a program with three stages ( $J = 3$ ) after the initial participation stage.

- For a person who would like to complete the program ( $D_3 = 1$ ), but is stopped by randomization after stage 2, we observe  $Y_2$  instead of  $Y_3$ .

- For a person who would like to complete the program ( $D_3 = 1$ ), but is stopped by randomization after stage 2, we observe  $Y_2$  instead of  $Y_3$ .
- If the person is randomized out after stage 1, we observe  $Y_1$  instead of  $Y_3$ .

- For a person who would like to complete the program ( $D_3 = 1$ ), but is stopped by randomization after stage 2, we observe  $Y_2$  instead of  $Y_3$ .
- If the person is randomized out after stage 1, we observe  $Y_1$  instead of  $Y_3$ .
- Let  $A_k$  be the indicator that we observe the agent with a stage  $k$  outcome.

- For a person who would like to complete the program ( $D_3 = 1$ ), but is stopped by randomization after stage 2, we observe  $Y_2$  instead of  $Y_3$ .
- If the person is randomized out after stage 1, we observe  $Y_1$  instead of  $Y_3$ .
- Let  $A_k$  be the indicator that we observe the agent with a stage  $k$  outcome.
- This can happen if a person would have chosen to stop at stage  $k$  ( $D_k = 1$ ) and survives randomization through  $k$  ( $\prod_{\ell=0}^{k-1} \xi_\ell = 1$ ), or if a person would have chosen to stop at a stage later than  $k$  but was thwarted from doing so by the randomization and settles for the best attainable state given the constraint imposed by the randomization.

- We can express  $A_k$  as

$$A_k = D_k \prod_{\ell=0}^{k-1} \xi_\ell + \sum_{j \geq k} D_j \left( \prod_{\ell=0}^{k-1} \xi_\ell \right) (1 - \xi_k), \quad k = 1, \dots, J.$$

If a person who chooses  $D_k = 1$  survives all stages of randomization through  $k - 1$  and hence is allowed to transit to  $k$ , we observe  $Y_k$  for that person.

- We can express  $A_k$  as

$$A_k = D_k \prod_{\ell=0}^{k-1} \xi_\ell + \sum_{j \geq k} D_j \left( \prod_{\ell=0}^{k-1} \xi_\ell \right) (1 - \xi_k), \quad k = 1, \dots, J.$$

If a person who chooses  $D_k = 1$  survives all stages of randomization through  $k - 1$  and hence is allowed to transit to  $k$ , we observe  $Y_k$  for that person.

- For persons who would choose  $D_j = 1$ ,  $j > k$ , but get randomized out at  $k$ , i.e.,  $\left( \prod_{\ell=0}^{k-1} \xi_\ell \right) (1 - \xi_k) = 1$ , we also observe  $Y_k$ .



- We now state the conditions under which sequential randomizations are instrumental variables for the  $A_j$ .

- We now state the conditions under which sequential randomizations are instrumental variables for the  $A_j$ .
- Let  $A_i(z, e_i)$  be the value of  $A_i$  when  $Z = z$  and  $\xi_i = e_i$ .

- We now state the conditions under which sequential randomizations are instrumental variables for the  $A_j$ .
- Let  $A_i(z, e_i)$  be the value of  $A_i$  when  $Z = z$  and  $\xi_i = e_i$ .
- Array the  $A_i$ ,  $i = 1, \dots, J$ , into a vector

$$A(z, e) = (A_1(z, e_1), A_2(z, e_2), \dots, A_J(z, e_J)).$$

- We now state the conditions under which sequential randomizations are instrumental variables for the  $A_j$ .
- Let  $A_i(z, e_i)$  be the value of  $A_i$  when  $Z = z$  and  $\xi_i = e_i$ .
- Array the  $A_i$ ,  $i = 1, \dots, J$ , into a vector

$$A(z, e) = (A_1(z, e_1), A_2(z, e_2), \dots, A_J(z, e_J)).$$

- A variety of randomization mechanisms are possible in which randomization depends on the information known to the randomizer at each stage of the program.

- IV conditions for  $\xi$  are satisfied under the following sequential randomization assumptions.

and

- (e-3b):  $\Pr\left(A_i = 1 \mid X, Z, D_\ell = 1 \text{ for } \ell < i, \xi_i, \prod_{\ell=0}^{i-1} \xi_\ell = 1\right)$  depends on  $\xi_i$ , for  $i = 1, \dots, J$ .

- IV conditions for  $\xi$  are satisfied under the following sequential randomization assumptions.
- They parallel the sequential randomization conditions developed in the dynamic models analyzed in Part III:

and

- (e-3b):  $\Pr\left(A_i = 1 \mid X, Z, D_\ell = 1 \text{ for } \ell < i, \xi_i, \prod_{\ell=0}^{i-1} \xi_\ell = 1\right)$  depends on  $\xi_i$ , for  $i = 1, \dots, J$ .

- IV conditions for  $\xi$  are satisfied under the following sequential randomization assumptions.
- They parallel the sequential randomization conditions developed in the dynamic models analyzed in Part III:
- (e-3a):  $\xi_i \perp\!\!\!\perp$   
 $\left( \tilde{Y}, \{A(z, e)\}_{(z, e) \in \mathcal{Z} \times \tilde{\xi}} \mid X, Z, D_\ell = 1 \text{ for } \ell < i, \prod_{\ell=0}^{i-1} \xi_\ell = 1 \right)$ ,  
 for  $i = 1, \dots, J$ ,

and

- (e-3b):  $\Pr \left( A_i = 1 \mid X, Z, D_\ell = 1 \text{ for } \ell < i, \xi_i, \prod_{\ell=0}^{i-1} \xi_\ell = 1 \right)$   
 depends on  $\xi_i$ , for  $i = 1, \dots, J$ .

- These expressions assume that the components of  $\tilde{Y} = (Y_0, \dots, Y_J)$  that are realized are known to the randomizer after the dropout decision is made, and thus cannot enter the conditioning set for the sequential randomizations.



- To fix ideas, consider a randomization of eligibility  $\xi_0$ , setting  $\xi_1 = \dots = \xi_J = 1$ .

- To fix ideas, consider a randomization of eligibility  $\xi_0$ , setting  $\xi_1 = \dots = \xi_J = 1$ .
- This is a randomization that makes people eligible for participation at all stages of the program.

- To fix ideas, consider a randomization of eligibility  $\xi_0$ , setting  $\xi_1 = \dots = \xi_J = 1$ .
- This is a randomization that makes people eligible for participation at all stages of the program.
- We investigate what this randomization identifies, assuming invariance conditions (PI-3) and (PI-4) hold.

- For those declared eligible,

$$E(Y | \xi_0 = 1) = \sum_{j=0}^J E(Y_j | D_j = 1) \Pr(D_j = 1). \quad (73)$$

For those declared ineligible,

$$E(Y | \xi_0 = 0) = \sum_{j=0}^J E(Y_0 | D_j = 1) \Pr(D_j = 1), \quad (74)$$

since agents cannot participate in any stage of the program and are all in the state “0” with outcome  $Y_0$ .

- For those declared eligible,

$$E(Y | \xi_0 = 1) = \sum_{j=0}^J E(Y_j | D_j = 1) \Pr(D_j = 1). \quad (73)$$

For those declared ineligible,

$$E(Y | \xi_0 = 0) = \sum_{j=0}^J E(Y_0 | D_j = 1) \Pr(D_j = 1), \quad (74)$$

since agents cannot participate in any stage of the program and are all in the state “0” with outcome  $Y_0$ .

- From observed choice behavior, we can identify each of the components of (73).

- We observe  $\Pr(D_j = 1)$  from observed choices of treatment, and we observe  $E(Y_j | D_j = 1)$  from observed outcomes for each treatment choice.

- We observe  $\Pr(D_j = 1)$  from observed choices of treatment, and we observe  $E(Y_j | D_j = 1)$  from observed outcomes for each treatment choice.
- Except for the choice probabilities ( $\Pr(D_j = 1), j = 0, \dots, J$ ) and  $E(Y_0 | D_0 = 1)$ , we cannot identify individual components of (74) for  $J > 1$ .

- We observe  $\Pr(D_j = 1)$  from observed choices of treatment, and we observe  $E(Y_j | D_j = 1)$  from observed outcomes for each treatment choice.
- Except for the choice probabilities ( $\Pr(D_j = 1), j = 0, \dots, J$ ) and  $E(Y_0 | D_0 = 1)$ , we cannot identify individual components of (74) for  $J > 1$ .
- When  $J = 1$ , we can identify all of the components of (74).



- We observe  $\Pr(D_j = 1)$  from observed choices of treatment, and we observe  $E(Y_j | D_j = 1)$  from observed outcomes for each treatment choice.
- Except for the choice probabilities ( $\Pr(D_j = 1), j = 0, \dots, J$ ) and  $E(Y_0 | D_0 = 1)$ , we cannot identify individual components of (74) for  $J > 1$ .
- When  $J = 1$ , we can identify all of the components of (74).
- The individual components of (74) cannot, without further assumptions, be identified by the experiment, although the sum can be.

- We observe  $\Pr(D_j = 1)$  from observed choices of treatment, and we observe  $E(Y_j | D_j = 1)$  from observed outcomes for each treatment choice.
- Except for the choice probabilities ( $\Pr(D_j = 1), j = 0, \dots, J$ ) and  $E(Y_0 | D_0 = 1)$ , we cannot identify individual components of (74) for  $J > 1$ .
- When  $J = 1$ , we can identify all of the components of (74).
- The individual components of (74) cannot, without further assumptions, be identified by the experiment, although the sum can be.
- Comparing the treatment group with the control group, we obtain the “intention to treat” parameter with respect to the randomization of  $\xi_0$  alone, setting  $\xi_1 = \dots = \xi_J = 1$  for anyone for whom  $\xi_0 = 1$ .

$$E(Y | \xi_0 = 1) - E(Y | \xi_0 = 0) = \sum_{j=1}^J E(Y_j - Y_0 | D_j = 1) \Pr(D_j = 1). \quad (75)$$



$$E(Y | \xi_0 = 1) - E(Y | \xi_0 = 0) = \sum_{j=1}^J E(Y_j - Y_0 | D_j = 1) \Pr(D_j = 1). \quad (75)$$

- For  $J > 1$ , this simple experimental estimator does not identify the effect of full participation in the program for those who participate ( $E(Y_J - Y_0 | D_J = 1)$ ) unless additional assumptions are invoked, such as the assumption that partial participation has the same mean effect as full participation for persons who drop out at the early stages, i.e.,  $E(Y_j - Y_0 | D_j = 1) = E(Y_J - Y_0 | D_J = 1)$  for all  $j$ .



$$E(Y | \xi_0 = 1) - E(Y | \xi_0 = 0) = \sum_{j=1}^J E(Y_j - Y_0 | D_j = 1) \Pr(D_j = 1). \quad (75)$$

- For  $J > 1$ , this simple experimental estimator does not identify the effect of full participation in the program for those who participate ( $E(Y_J - Y_0 | D_J = 1)$ ) unless additional assumptions are invoked, such as the assumption that partial participation has the same mean effect as full participation for persons who drop out at the early stages, i.e.,  $E(Y_j - Y_0 | D_j = 1) = E(Y_J - Y_0 | D_J = 1)$  for all  $j$ .
- This assumption might be appropriate if just getting into the program is all that matters—a pure signalling effect.

- A second set of conditions for identification of this parameter is that  $E(Y_j - Y_0 \mid D_j = 1) = 0$  for all  $j < J$ .

- A second set of conditions for identification of this parameter is that  $E(Y_j - Y_0 | D_j = 1) = 0$  for all  $j < J$ .
- Under those conditions, if we divide the mean difference by  $\Pr(D_J = 1)$ , we obtain the “Bloom” estimator ( ?, ?)

$$IV_{\text{Bloom}} = \frac{E(Y | \xi_0 = 1) - E(Y | \xi_0 = 0)}{\Pr(D_J = 1)},$$

assuming  $\Pr(D_J = 1) \neq 0$ .

- A second set of conditions for identification of this parameter is that  $E(Y_j - Y_0 | D_j = 1) = 0$  for all  $j < J$ .
- Under those conditions, if we divide the mean difference by  $\Pr(D_J = 1)$ , we obtain the “Bloom” estimator ( ?, ?)

$$IV_{\text{Bloom}} = \frac{E(Y | \xi_0 = 1) - E(Y | \xi_0 = 0)}{\Pr(D_J = 1)},$$

assuming  $\Pr(D_J = 1) \neq 0$ .

- This is an IV estimator using  $\xi_0$  as the instrument for  $A_J$ .



- A second set of conditions for identification of this parameter is that  $E(Y_j - Y_0 | D_j = 1) = 0$  for all  $j < J$ .
- Under those conditions, if we divide the mean difference by  $\Pr(D_J = 1)$ , we obtain the “Bloom” estimator ( ?, ?)

$$IV_{\text{Bloom}} = \frac{E(Y | \xi_0 = 1) - E(Y | \xi_0 = 0)}{\Pr(D_J = 1)},$$

assuming  $\Pr(D_J = 1) \neq 0$ .

- This is an IV estimator using  $\xi_0$  as the instrument for  $A_J$ .
- In general, the mean difference between the treated and the controlled identifies only the composite term shown in (75).

- A second set of conditions for identification of this parameter is that  $E(Y_j - Y_0 | D_j = 1) = 0$  for all  $j < J$ .
- Under those conditions, if we divide the mean difference by  $\Pr(D_J = 1)$ , we obtain the “Bloom” estimator ( ?, ?)

$$IV_{\text{Bloom}} = \frac{E(Y | \xi_0 = 1) - E(Y | \xi_0 = 0)}{\Pr(D_J = 1)},$$

assuming  $\Pr(D_J = 1) \neq 0$ .

- This is an IV estimator using  $\xi_0$  as the instrument for  $A_J$ .
- In general, the mean difference between the treated and the controlled identifies only the composite term shown in (75).
- In this case, the simple randomization estimator identifies a not-so-simple or easily interpreted parameter.

- More generally, if we randomize persons out after completing stage  $k$  ( $(\prod_{\ell=0}^{k-1} \xi_\ell) (1 - \xi_k) = 1$ ) and for another group establish full eligibility at all stages ( $\prod_{\ell=0}^J \xi_\ell = 1$ ), we obtain

$$\begin{aligned}
 & E \left[ Y \left| \prod_{\ell=0}^J \xi_\ell = 1 \right. \right] - E \left[ Y \left| \left( \prod_{\ell=0}^{k-1} \xi_\ell \right) (1 - \xi_k) = 1 \right. \right] \\
 &= \sum_{j=k}^J E(Y_j - Y_k \mid D_j = 1) \Pr(D_j = 1),
 \end{aligned}$$

- Hence, since we know  $E(Y_k | D_k = 1)$  and  $\Pr(D_k = 1)$  from observational data, we can identify the combination of parameters

$$\sum_{j=k+1}^J E(Y_k | D_j = 1) \Pr(D_j = 1), \quad (76)$$

for each randomization that stops persons from advancing beyond level  $k$ ,  $k = 0, \dots, J - 1$ .

- Observe that a randomization of eligibility that prevents people from going to stage  $J - 1$  but not to stage  $J$  ( $(\prod_{\ell=0}^{J-2} \xi_{\ell}) (1 - \xi_{J-1}) = 1$ ) identifies  $E(Y_J - Y_{J-1} | D_J = 1)$ :

$$\begin{aligned}
 & E(Y | \xi_0 = 1, \dots, \xi_{J-2} = 1, \xi_{J-1} = 0) \\
 = & \left[ \sum_{j=0}^{J-1} E(Y_j | D_j = 1) \Pr(D_j = 1) \right] + E(Y_{J-1} | D_J = 1) \Pr(D_J = 1).
 \end{aligned}$$

- Observe that a randomization of eligibility that prevents people from going to stage  $J - 1$  but not to stage  $J$  ( $[\prod_{\ell=0}^{J-2} \xi_{\ell}] (1 - \xi_{J-1}) = 1$ ) identifies  $E(Y_J - Y_{J-1} | D_J = 1)$ :

$$\begin{aligned}
 & E(Y | \xi_0 = 1, \dots, \xi_{J-2} = 1, \xi_{J-1} = 0) \\
 = & \left[ \sum_{j=0}^{J-1} E(Y_j | D_j = 1) \Pr(D_j = 1) \right] + E(Y_{J-1} | D_J = 1) \Pr(D_J = 1).
 \end{aligned}$$

- Thus,

$$\begin{aligned}
 & E(Y | \xi_0 = 1, \dots, \xi_J = 1) - E(Y | \xi_0 = 1, \dots, \xi_{J-1} = 1, \xi_J = 0) \\
 = & E(Y_J - Y_{J-1} | D_J = 1) \Pr(D_J = 1).
 \end{aligned}$$

Since  $\Pr(D_J = 1)$  is observed from choice data, as is  $E(Y_J | D_J = 1)$ , we can identify  $E(Y_{J-1} | D_J = 1)$  from the experiment.

- In the general case under assumptions (PI-3) and (PI-4), a randomization that prevents agents from moving beyond stage  $\ell$  ( $\xi_0 = 1, \dots, \xi_{\ell-1} = 1, \xi_\ell = 0$ ) directly identifies

$$\begin{aligned}
 & E(Y \mid \xi_0 = 1, \dots, \xi_{\ell-1} = 1, \xi_\ell = 0) = \\
 & = \underbrace{\sum_{j=0}^{\ell} E(Y_j \mid D_j = 1) \Pr(D_j = 1)}_{\text{all components known from observational data}} \\
 & + \underbrace{\sum_{j=\ell+1}^J E(Y_\ell \mid D_j = 1) \Pr(D_j = 1)}_{\text{sum and probability weights known, but not individual } E(Y_\ell \mid D_j=1)}.
 \end{aligned}$$

- In the general case under assumptions (PI-3) and (PI-4), a randomization that prevents agents from moving beyond stage  $\ell$  ( $\xi_0 = 1, \dots, \xi_{\ell-1} = 1, \xi_\ell = 0$ ) directly identifies

$$\begin{aligned}
 & E(Y \mid \xi_0 = 1, \dots, \xi_{\ell-1} = 1, \xi_\ell = 0) = \\
 & = \underbrace{\sum_{j=0}^{\ell} E(Y_j \mid D_j = 1) \Pr(D_j = 1)}_{\text{all components known from observational data}} \\
 & + \underbrace{\sum_{j=\ell+1}^J E(Y_\ell \mid D_j = 1) \Pr(D_j = 1)}_{\text{sum and probability weights known, but not individual } E(Y_\ell \mid D_j=1)}.
 \end{aligned}$$

- All of the components of the first set of terms on the right-hand side are known from observational data.



- The probabilities in the second set of terms are known, but the individual conditional expectations  $E(Y_\ell | D_j = 1)$ ,  $j = \ell + 1, \dots, J$ , are not known without further assumptions.

- The probabilities in the second set of terms are known, but the individual conditional expectations  $E(Y_\ell | D_j = 1)$ ,  $j = \ell + 1, \dots, J$ , are not known without further assumptions.
- Randomization at stage  $\ell$  is an IV.

- The probabilities in the second set of terms are known, but the individual conditional expectations  $E(Y_\ell | D_j = 1)$ ,  $j = \ell + 1, \dots, J$ , are not known without further assumptions.
- Randomization at stage  $\ell$  is an IV.
- To show this, decompose the observed outcome  $Y$  into components associated with each value of  $A_j$ , the indicator associated with observing a stage  $j$  outcome:

$$Y = \sum_{j=0}^J A_j Y_j.$$

- The probabilities in the second set of terms are known, but the individual conditional expectations  $E(Y_\ell | D_j = 1)$ ,  $j = \ell + 1, \dots, J$ , are not known without further assumptions.
- Randomization at stage  $\ell$  is an IV.
- To show this, decompose the observed outcome  $Y$  into components associated with each value of  $A_j$ , the indicator associated with observing a stage  $j$  outcome:

$$Y = \sum_{j=0}^J A_j Y_j.$$

- We can interpret  $\xi_\ell$  as an instrument for  $A_\ell$ .

- Keeping the conditioning on  $X, Z$  implicit, we obtain

$$\begin{aligned}
 IV_{\xi_\ell} &= \frac{E[Y | \xi_\ell = 0] - E[Y | \xi_\ell = 1]}{\Pr(A_\ell = 1 | \xi_\ell = 0) - \Pr(A_\ell = 1 | \xi_\ell = 1)} \\
 &= \frac{\sum_{j=\ell+1}^J E[Y_\ell - Y_j | D_j = 1] \Pr(D_j = 1)}{\sum_{j=\ell+1}^J \Pr(D_j = 1)}, \ell = 0, \dots, J-1.
 \end{aligned}$$

- Keeping the conditioning on  $X, Z$  implicit, we obtain

$$\begin{aligned} \text{IV}_{\xi_\ell} &= \frac{E[Y | \xi_\ell = 0] - E[Y | \xi_\ell = 1]}{\Pr(A_\ell = 1 | \xi_\ell = 0) - \Pr(A_\ell = 1 | \xi_\ell = 1)} \\ &= \frac{\sum_{j=\ell+1}^J E[Y_\ell - Y_j | D_j = 1] \Pr(D_j = 1)}{\sum_{j=\ell+1}^J \Pr(D_j = 1)}, \ell = 0, \dots, J-1. \end{aligned}$$

- By the preceding analysis, we know the numerator term but not the individual components.

- Keeping the conditioning on  $X, Z$  implicit, we obtain

$$\begin{aligned} \text{IV}_{\xi_\ell} &= \frac{E[Y | \xi_\ell = 0] - E[Y | \xi_\ell = 1]}{\Pr(A_\ell = 1 | \xi_\ell = 0) - \Pr(A_\ell = 1 | \xi_\ell = 1)} \\ &= \frac{\sum_{j=\ell+1}^J E[Y_\ell - Y_j | D_j = 1] \Pr(D_j = 1)}{\sum_{j=\ell+1}^J \Pr(D_j = 1)}, \ell = 0, \dots, J-1. \end{aligned}$$

- By the preceding analysis, we know the numerator term but not the individual components.
- We know the denominator from choices measured in observational data and invariance assumption (PI-3).

- Keeping the conditioning on  $X, Z$  implicit, we obtain

$$\begin{aligned} \text{IV}_{\xi_\ell} &= \frac{E[Y | \xi_\ell = 0] - E[Y | \xi_\ell = 1]}{\Pr(A_\ell = 1 | \xi_\ell = 0) - \Pr(A_\ell = 1 | \xi_\ell = 1)} \\ &= \frac{\sum_{j=\ell+1}^J E[Y_\ell - Y_j | D_j = 1] \Pr(D_j = 1)}{\sum_{j=\ell+1}^J \Pr(D_j = 1)}, \ell = 0, \dots, J-1. \end{aligned}$$

- By the preceding analysis, we know the numerator term but not the individual components.
- We know the denominator from choices measured in observational data and invariance assumption (PI-3).
- The IV formalism is less helpful in the general case.



- Table 13 summarizes the parameters or combinations of parameters that can be identified from randomizations performed at different stages.

- Table 13 summarizes the parameters or combinations of parameters that can be identified from randomizations performed at different stages.
- It presents the array of factual and counterfactual conditional mean outcomes  $E(Y_j | D_\ell = 1)$ ,  $j = 0, \dots, J$  and  $\ell = 0, \dots, J$ .

- Table 13 summarizes the parameters or combinations of parameters that can be identified from randomizations performed at different stages.
- It presents the array of factual and counterfactual conditional mean outcomes  $E(Y_j | D_\ell = 1)$ ,  $j = 0, \dots, J$  and  $\ell = 0, \dots, J$ .
- The conditional mean outcomes obtained from observational data are on the diagonal of the table ( $E(Y_j | D_j = 1), j = 0, \dots, J$ ).

- Table 13 summarizes the parameters or combinations of parameters that can be identified from randomizations performed at different stages.
- It presents the array of factual and counterfactual conditional mean outcomes  $E(Y_j | D_\ell = 1)$ ,  $j = 0, \dots, J$  and  $\ell = 0, \dots, J$ .
- The conditional mean outcomes obtained from observational data are on the diagonal of the table ( $E(Y_j | D_j = 1), j = 0, \dots, J$ ).
- Because of choices of agents, experiments do not directly identify the elements in the table that are above the diagonal.

- Table 13 summarizes the parameters or combinations of parameters that can be identified from randomizations performed at different stages.
- It presents the array of factual and counterfactual conditional mean outcomes  $E(Y_j | D_\ell = 1)$ ,  $j = 0, \dots, J$  and  $\ell = 0, \dots, J$ .
- The conditional mean outcomes obtained from observational data are on the diagonal of the table ( $E(Y_j | D_j = 1), j = 0, \dots, J$ ).
- Because of choices of agents, experiments do not directly identify the elements in the table that are above the diagonal.
- Under assumptions (PI-3) and (PI-4), experiments described at the base of the table identify the *combinations* of the parameters below the diagonal.

**Table 13:** Parameters and Combinations of Parameters That Can be Identified by Different Randomizations

Choice Probabilities (known)		Outcome						
		$Y_0$	$Y_1$	...	$Y_j$	...	$Y_{j-1}$	$Y_j$
$\Pr(D_0 = 1)$	$D_0$	$E(Y_0   D_0 = 1)$	$E(Y_1   D_0 = 1)$	...	$E(Y_j   D_0 = 1)$	...	$E(Y_{j-1}   D_0 = 1)$	$E(Y_j   D_0 = 1)$
$\Pr(D_1 = 1)$	$D_1$	$E(Y_0   D_1 = 1)$	$E(Y_1   D_1 = 1)$	...	$E(Y_j   D_1 = 1)$	...	$E(Y_{j-1}   D_1 = 1)$	$E(Y_j   D_1 = 1)$
$\Pr(D_2 = 1)$	$D_2$	$E(Y_0   D_2 = 1)$	$E(Y_1   D_2 = 1)$	...	$E(Y_j   D_2 = 1)$	...	$E(Y_{j-1}   D_2 = 1)$	$E(Y_j   D_2 = 1)$
$\vdots$	$\vdots$	$\vdots$	$\vdots$		$\vdots$		$\vdots$	$\vdots$
		$\vdots$	$\vdots$		$\vdots$		$\vdots$	$\vdots$
		$\vdots$	$\vdots$		$\vdots$		$\vdots$	$\vdots$
$\Pr(D_j = 1)$	$D_j$	$E(Y_0   D_j = 1)$	$E(Y_1   D_j = 1)$	...	$E(Y_j   D_j = 1)$	...	$E(Y_{j-1}   D_j = 1)$	$E(Y_j   D_j = 1)$
$\vdots$	$\vdots$	$\vdots$	$\vdots$		$\vdots$		$\vdots$	$\vdots$
$\Pr(D_{j-1} = 1)$	$D_{j-1}$	$E(Y_0   D_{j-1} = 1)$	$E(Y_1   D_{j-1} = 1)$	...	$E(Y_j   D_{j-1} = 1)$	...	$E(Y_{j-1}   D_{j-1} = 1)$	$E(Y_j   D_{j-1} = 1)$
$\Pr(D_j = 1)$	$D_j$	$E(Y_0   D_j = 1)$	$E(Y_1   D_j = 1)$	...	$E(Y_j   D_j = 1)$	...	$E(Y_{j-1}   D_j = 1)$	$E(Y_j   D_j = 1)$
Randomization		$\xi_0 = 0$	$\xi_1 = 0$	...	$\xi_j = 0$	...	$\xi_{j-1} = 0$	$\xi_j = 0$
New Identified Combinations of Parameters		$\sum_{\ell=1}^j \{E(Y_0   D_\ell = 1) \times \Pr(D_\ell = 1)\}$	$\sum_{\ell=2}^j \{E(Y_1   D_\ell = 1) \times \Pr(D_\ell = 1)\}$	...	$\sum_{\ell=j+1}^j \{E(Y_j   D_\ell = 1) \times \Pr(D_\ell = 1)\}$	...	$E(Y_{j-1}   D_j = 1)$	

- Recall that if  $\xi_\ell = 0$ , the agent cannot advance beyond stage  $\ell$ .

- Recall that if  $\xi_\ell = 0$ , the agent cannot advance beyond stage  $\ell$ .
- If we randomly deny eligibility to move to  $J$  ( $\xi_{J-1} = 0$ ), we point identify  $E(Y_{J-1} | D_J = 1)$ , as well as the parameters that can be obtained from observational data.



- Recall that if  $\xi_\ell = 0$ , the agent cannot advance beyond stage  $\ell$ .
- If we randomly deny eligibility to move to  $J$  ( $\xi_{J-1} = 0$ ), we point identify  $E(Y_{J-1} | D_J = 1)$ , as well as the parameters that can be obtained from observational data.
- In general, we can only identify the combinations of parameters shown at the base of the table.

- Recall that if  $\xi_\ell = 0$ , the agent cannot advance beyond stage  $\ell$ .
- If we randomly deny eligibility to move to  $J$  ( $\xi_{J-1} = 0$ ), we point identify  $E(Y_{J-1} | D_J = 1)$ , as well as the parameters that can be obtained from observational data.
- In general, we can only identify the combinations of parameters shown at the base of the table.
- Following ?, ???? , and ? , we can use the identified combinations from different randomizations to bound the admissible values of counterfactuals below the diagonal of the table 13.

- ? present a test for a strengthened version of the identifying assumptions made by Bloom.

- ? present a test for a strengthened version of the identifying assumptions made by Bloom.
- They perform a sensitivity analysis to analyze departures from the assumption that dropouts have the same outcomes as nonparticipants.

- ? present a test for a strengthened version of the identifying assumptions made by Bloom.
- They perform a sensitivity analysis to analyze departures from the assumption that dropouts have the same outcomes as nonparticipants.
- ? apply the Manski bounds in carefully executed empirical examples and show the difficulties involved in using the Bloom estimator in experiments with multiple outcomes.

- ? present a test for a strengthened version of the identifying assumptions made by Bloom.
- They perform a sensitivity analysis to analyze departures from the assumption that dropouts have the same outcomes as nonparticipants.
- ? apply the Manski bounds in carefully executed empirical examples and show the difficulties involved in using the Bloom estimator in experiments with multiple outcomes.
- We next turn to some evidence on the importance of randomization bias.

## Evidence on Randomization Bias

- Violations of assumption (PI-3) in the general case with essential heterogeneity affect the interpretation of the outputs of social experiments.

## Evidence on Randomization Bias

- Violations of assumption (PI-3) in the general case with essential heterogeneity affect the interpretation of the outputs of social experiments.
- They are manifestations of a more general problem termed “Hawthorne effects” that arise from observing any population (see ??) . How important is this theoretical possibility in practice?



## Evidence on Randomization Bias

- Violations of assumption (PI-3) in the general case with essential heterogeneity affect the interpretation of the outputs of social experiments.
- They are manifestations of a more general problem termed “Hawthorne effects” that arise from observing any population (see ??) . How important is this theoretical possibility in practice?
- Surprisingly, very little is known about the answer to this question for the social experiments conducted in economics.

## Evidence on Randomization Bias

- Violations of assumption (PI-3) in the general case with essential heterogeneity affect the interpretation of the outputs of social experiments.
- They are manifestations of a more general problem termed “Hawthorne effects” that arise from observing any population (see ??) . How important is this theoretical possibility in practice?
- Surprisingly, very little is known about the answer to this question for the social experiments conducted in economics.
- This is so because randomized social experimentation has usually only been implemented on “pilot projects” or “demonstration projects” designed to evaluate new programs never previously estimated.

- Disruption by randomization cannot be confirmed or denied using data from these experiments.

- Disruption by randomization cannot be confirmed or denied using data from these experiments.
- In one ongoing program evaluated by randomization by the Manpower Demonstration Research Corporation (MDRC), participation was compulsory for the target population (?).

- Disruption by randomization cannot be confirmed or denied using data from these experiments.
- In one ongoing program evaluated by randomization by the Manpower Demonstration Research Corporation (MDRC), participation was compulsory for the target population (?).
- Hence randomization did not affect applicant pools or assessments of applicant eligibility by program administrators.

- There is some information on the importance of randomization, although it is indirect.

- There is some information on the importance of randomization, although it is indirect.
- In the 1980s, the U.S. Department of Labor financed a large-scale experimental evaluation of the ongoing, large-scale manpower training program authorized under the Job Training Partnership Act (JTPA).

- There is some information on the importance of randomization, although it is indirect.
- In the 1980s, the U.S. Department of Labor financed a large-scale experimental evaluation of the ongoing, large-scale manpower training program authorized under the Job Training Partnership Act (JTPA).
- A study by ? gives some indirect information from which it is possible to determine whether randomization bias was present in an ongoing program.



- There is some information on the importance of randomization, although it is indirect.
- In the 1980s, the U.S. Department of Labor financed a large-scale experimental evaluation of the ongoing, large-scale manpower training program authorized under the Job Training Partnership Act (JTPA).
- A study by ? gives some indirect information from which it is possible to determine whether randomization bias was present in an ongoing program.
- Job training in the United States is organized through geographically decentralized centers.

- These centers receive incentive payments for placing unemployed persons and persons on welfare in “high-paying” jobs.

- These centers receive incentive payments for placing unemployed persons and persons on welfare in “high-paying” jobs.
- The participation of centers in the experiment was not compulsory.

- These centers receive incentive payments for placing unemployed persons and persons on welfare in “high-paying” jobs.
- The participation of centers in the experiment was not compulsory.
- Funds were set aside to compensate job centers for the administrative costs of participating in the experiment.

- These centers receive incentive payments for placing unemployed persons and persons on welfare in “high-paying” jobs.
- The participation of centers in the experiment was not compulsory.
- Funds were set aside to compensate job centers for the administrative costs of participating in the experiment.
- The funds set aside range from 5 percent to 10 percent of the total operating costs of the centers.

- In attempting to enroll geographically dispersed sites, MDRC experienced a training center refusal rate in excess of 90 percent.

- In attempting to enroll geographically dispersed sites, MDRC experienced a training center refusal rate in excess of 90 percent.
- The reasons for refusal to participate are given in table 14.

- In attempting to enroll geographically dispersed sites, MDRC experienced a training center refusal rate in excess of 90 percent.
- The reasons for refusal to participate are given in table 14.
- (The reasons stated there are not mutually exclusive.)



- In attempting to enroll geographically dispersed sites, MDRC experienced a training center refusal rate in excess of 90 percent.
- The reasons for refusal to participate are given in table 14.
- (The reasons stated there are not mutually exclusive.)
- Leading the list are ethical and public relations objections to randomization.

- In attempting to enroll geographically dispersed sites, MDRC experienced a training center refusal rate in excess of 90 percent.
- The reasons for refusal to participate are given in table 14.
- (The reasons stated there are not mutually exclusive.)
- Leading the list are ethical and public relations objections to randomization.
- Major fears (items 2 and 3) were expressed about the effects of randomization on the quality of applicant pool, which would impede the profitability of the training centers.

- In attempting to enroll geographically dispersed sites, MDRC experienced a training center refusal rate in excess of 90 percent.
- The reasons for refusal to participate are given in table 14.
- (The reasons stated there are not mutually exclusive.)
- Leading the list are ethical and public relations objections to randomization.
- Major fears (items 2 and 3) were expressed about the effects of randomization on the quality of applicant pool, which would impede the profitability of the training centers.
- By randomizing, the centers had to widen the available pool of persons deemed eligible, and there was great concern about the effects of this widening on applicant quality—precisely the behavior ruled out by assumptions (PI-3) and (PI-4).

**Table 14:** Percentage of Local JTPA Agencies Citing Specific Concerns About Participating in the Experiment

Concern	Percentage of Training Centers Citing the Concern
1. Ethical and public relations implications of:	
a. Random assignment in social programs	61.8
b. Denial of services to controls	54.4
2. Potential negative effect of creation of a control group on achievement of client recruitment goals	47.8
3. Potential negative impact on performance standards	25.4
4. Implementation of the study when service providers do intake	21.1
5. Objections of service providers to the study	17.5
6. Potential staff administrative burden	16.2
7. Possible lack of support by elected officials	15.8
8. Legality of random assignment and possible grievances	14.5
9. Procedures for providing controls with referrals to other services	14.0
10. Special recruitment problems for out-of-school youth	10.5
Sample size	228

Notes: Concerns noted by fewer than 5 percent of the training centers are not listed. Percentages add up to more than 100.0 because training centers could raise more than one concern.

Source: Based on responses of 228 local JTPA agencies contacted about possible participation in the National JTPA Study.

Source: Heckman (1992), based on Doolittle and Traeger (1990).

- In attempting to entice centers to participate, MDRC had to reduce the randomized rejection probability from  $\frac{1}{2}$  to as low as  $\frac{1}{6}$  for certain centers.

- In attempting to entice centers to participate, MDRC had to reduce the randomized rejection probability from  $\frac{1}{2}$  to as low as  $\frac{1}{6}$  for certain centers.
- The resulting reduction in the size of the control group impairs the power of statistical tests designed to test the null hypothesis of no program effect.

- In attempting to entice centers to participate, MDRC had to reduce the randomized rejection probability from  $\frac{1}{2}$  to as low as  $\frac{1}{6}$  for certain centers.
- The resulting reduction in the size of the control group impairs the power of statistical tests designed to test the null hypothesis of no program effect.
- Compensation for participation was expanded sevenfold in order to get any centers to participate in the experiment.



- The MDRC analysts conclude:

Implementing a complex random assignment research design in an ongoing program providing a variety of services does inevitably change its operation in some ways. The most likely difference arising from a random assignment field study of program impacts is a change in the mix of clients served. Expanded recruitment efforts, needed to generate the control group, draw in additional applicants who are not identical to the people previously served. A second likely change is that the treatment categories may somewhat restrict program staff's flexibility to change service recommendations

(?, p. 121).

These authors go on to note that

some [training centers] because of severe recruitment problems or up-front services cannot implement the type of random assignment model needed to answer the various impact questions without major changes in procedures

(?, p. 123).

- This indirect evidence is hardly decisive even about the JTPA experiment, much less all experiments.

- This indirect evidence is hardly decisive even about the JTPA experiment, much less all experiments.
- Training centers may offer these arguments only as a means of avoiding administrative scrutiny, and there may be no “real” effect of randomization.

- This indirect evidence is hardly decisive even about the JTPA experiment, much less all experiments.
- Training centers may offer these arguments only as a means of avoiding administrative scrutiny, and there may be no “real” effect of randomization.
- During the JTPA experiment conducted at Corpus Christi, Texas, center administrators successfully petitioned the government of Texas for a waiver of its performance standards on the ground that the experiment disrupted center operations.

- This indirect evidence is hardly decisive even about the JTPA experiment, much less all experiments.
- Training centers may offer these arguments only as a means of avoiding administrative scrutiny, and there may be no “real” effect of randomization.
- During the JTPA experiment conducted at Corpus Christi, Texas, center administrators successfully petitioned the government of Texas for a waiver of its performance standards on the ground that the experiment disrupted center operations.
- Self-selection likely guarantees that participant sites are the least likely sites to suffer disruption.

- Such selective participation in the experiment calls into question the validity of experimental estimates as a statement about the JTPA system as a whole, as it clearly poses a threat to external validity — problem (P-2) as defined in Part I.

- Such selective participation in the experiment calls into question the validity of experimental estimates as a statement about the JTPA system as a whole, as it clearly poses a threat to external validity — problem (P-2) as defined in Part I.
- ? report similar problems in a randomized evaluation of a job training program in Norway.



- ? note that subjects in drug trials were less likely to participate in randomized trials than in nonexperimental studies.

- ? note that subjects in drug trials were less likely to participate in randomized trials than in nonexperimental studies.
- They discuss one study of drugs administered to children afflicted with a disease.

- ? note that subjects in drug trials were less likely to participate in randomized trials than in nonexperimental studies.
- They discuss one study of drugs administered to children afflicted with a disease.
- The study had two components.

- ? note that subjects in drug trials were less likely to participate in randomized trials than in nonexperimental studies.
- They discuss one study of drugs administered to children afflicted with a disease.
- The study had two components.
- The nonexperimental phase of the study had a 4 percent refusal rate, while 34 percent of a subsample of the same parents refused to participate in a randomized subtrial, although the treatments were equally nonthreatening.

- These authors cite further evidence suggesting that refusal to participate in randomization schemes is selective.

- These authors cite further evidence suggesting that refusal to participate in randomization schemes is selective.
- In a study of treatment of adults with cirrhosis, no effect of the treatment was found for participants in a randomized trial.

- These authors cite further evidence suggesting that refusal to participate in randomization schemes is selective.
- In a study of treatment of adults with cirrhosis, no effect of the treatment was found for participants in a randomized trial.
- But the death rates for those randomized out of the treatment were substantially lower than among those individuals who refused to participate in the experiment, despite the fact that both groups were administered the same alternative treatment.

- These authors cite further evidence suggesting that refusal to participate in randomization schemes is selective.
- In a study of treatment of adults with cirrhosis, no effect of the treatment was found for participants in a randomized trial.
- But the death rates for those randomized out of the treatment were substantially lower than among those individuals who refused to participate in the experiment, despite the fact that both groups were administered the same alternative treatment.
- Part of any convincing identification strategy by randomization requires that the agent document the absence of randomization bias.



- These authors cite further evidence suggesting that refusal to participate in randomization schemes is selective.
- In a study of treatment of adults with cirrhosis, no effect of the treatment was found for participants in a randomized trial.
- But the death rates for those randomized out of the treatment were substantially lower than among those individuals who refused to participate in the experiment, despite the fact that both groups were administered the same alternative treatment.
- Part of any convincing identification strategy by randomization requires that the agent document the absence of randomization bias.
- We next consider some evidence on the importance of dropping out and noncompliance with experimental protocols.

## Evidence on Dropping Out and Substitution Bias

- Dropouts are a feature of all social programs.

## Evidence on Dropping Out and Substitution Bias

- Dropouts are a feature of all social programs.
- Randomization may raise dropout rates, but the evidence for such effects is weak.

## Evidence on Dropping Out and Substitution Bias

- Dropouts are a feature of all social programs.
- Randomization may raise dropout rates, but the evidence for such effects is weak.
- In addition, most social programs have good substitutes, so that the estimated effect of a program as typically estimated has to be defined relative to the full range of substitute activities in which non-participants engage.

## Evidence on Dropping Out and Substitution Bias

- Dropouts are a feature of all social programs.
- Randomization may raise dropout rates, but the evidence for such effects is weak.
- In addition, most social programs have good substitutes, so that the estimated effect of a program as typically estimated has to be defined relative to the full range of substitute activities in which non-participants engage.
- Experiments exacerbate this problem by creating a pool of persons who attempt to take training who then flock to substitute programs when they are placed in an experimental control group ( $\xi = 0$  in the simple randomization analyzed in Slides 833–876).

- Table 15 (reproduced from ?) demonstrates the practical importance of both dropout and substitution bias in experimental evaluations.

- Table 15 (reproduced from ?) demonstrates the practical importance of both dropout and substitution bias in experimental evaluations.
- It reports the rates of treatment group dropout and control group substitution from a variety of social experiments.

- Table 15 (reproduced from ?) demonstrates the practical importance of both dropout and substitution bias in experimental evaluations.
- It reports the rates of treatment group dropout and control group substitution from a variety of social experiments.
- It reveals that the fraction of treatment group members receiving program services is often less than 0.7, and sometimes less than 0.5.



- Table 15 (reproduced from ?) demonstrates the practical importance of both dropout and substitution bias in experimental evaluations.
- It reports the rates of treatment group dropout and control group substitution from a variety of social experiments.
- It reveals that the fraction of treatment group members receiving program services is often less than 0.7, and sometimes less than 0.5.
- Furthermore, the observed characteristics of the treatment group members who drop out often differ from those who remain and receive the program services.

**Table 15:** Fraction of Experimental Treatment and Control Groups Receiving Services in Experimental Evaluations of Employment and Training Programs

Study	Authors/time period	Target group(s)	Fraction of treatments receiving services	Fraction of controls receiving services
1. NSW	Hollister, et al. (1984) (9 months after RA)	Long-term AFDC women	0.95	0.11
		Ex-addicts	NA	0.03
		17-20 year old high school dropouts	NA	0.04
2. SWIM	Friedlander and Hamilton (1993) (Time period not reported)	AFDC women: applicants and recipients		
		a. Job search assistance	0.54	0.01
		b. Work experience	0.21	0.01
		c. Classroom training/OJT	0.39	0.21
		d. Any activity	0.69	0.30
		AFDC-U unemployed fathers		
		a. Job search assistance	0.60	0.01
		b. Work experience	0.21	0.01
3. JOBSTART	Cave, et al. (1993) (12 months after RA)	c. Classroom training/OJT	0.34	0.22
		d. Any activity	0.70	0.23
		Youth high school dropouts		
		Classroom training/OJT	0.90	0.26
4. Project Independence	Kemple, et al. (1995) (24 months after RA)	AFDC women: applicants and recipients		
		a. Job search assistance	0.43	0.19
		b. Classroom training/OJT	0.42	0.31
		c. Any activity	0.64	0.40

Table 15 [Continued]

Study	Authors/time period	Target group(s)	Fraction of treatments receiving services	Fraction of controls receiving services
5. New Chance	Quint, et al. (1994) (18 months after RA)	Teenage single mothers		
		Any education services	0.82	0.48
		Any training services	0.26	0.15
6. National JTPA Study	Heckman and Smith (1998) (18 months after RA)	Any education or training	0.87	0.55
		Self-reported from survey data		
		Adult males	0.38	0.24
		Adult females	0.51	0.33
		Male youth	0.50	0.32
		Female youth	0.81	0.42
		Combined Administrative Survey Data		
		Adult males	0.74	0.25
Adult females	0.78	0.34		
Male youth	0.81	0.34		
Female youth	0.81	0.42		

Notes: RA = random assignment. H.S. = high school. AFDC = Aid to Families with Dependent Children. OJI = On the Job Training. Service receipt includes any employment and training services. The services received by the controls in the NSW study are CETA and WIN jobs. For the Long Term AFDC Women, this measure also includes regular public sector employment during the period.

Sources for data: Maynard and Brown (1980), p. 169, Table A14; Masters and Maynard (1981), p. 148, Table A.15; Friedlander and Hamilton (1993), p. 22, Table 3.1; Cave, et al. (1993), p. 95, Table 4.1; Quint, et al. (1994), p. 110, Table 4.9; and Kemple, et al. (1995), p. 58, Table 3.5; Heckman and Smith (1998) and calculations by the authors.

Source: Heckman, LaLonde and Smith (1999).

- With regard to substitution bias, table 15 shows that as many as 40% of the controls in some experiments received substitute services elsewhere.

- With regard to substitution bias, table 15 shows that as many as 40% of the controls in some experiments received substitute services elsewhere.
- In a simple one treatment experiment with full compliance ( $\xi = 1 \Rightarrow A = 1$  and  $\xi = 0 \Rightarrow A = 0$ ), all individuals assigned to the treatment group receive the treatment and there is no control group substitution, so that the difference between the fraction of treatments and controls that receive the treatment equals 1.0.

- With regard to substitution bias, table 15 shows that as many as 40% of the controls in some experiments received substitute services elsewhere.
- In a simple one treatment experiment with full compliance ( $\xi = 1 \Rightarrow A = 1$  and  $\xi = 0 \Rightarrow A = 0$ ), all individuals assigned to the treatment group receive the treatment and there is no control group substitution, so that the difference between the fraction of treatments and controls that receive the treatment equals 1.0.
- In practice, this difference is often well below 1.0.

- With regard to substitution bias, table 15 shows that as many as 40% of the controls in some experiments received substitute services elsewhere.
- In a simple one treatment experiment with full compliance ( $\xi = 1 \Rightarrow A = 1$  and  $\xi = 0 \Rightarrow A = 0$ ), all individuals assigned to the treatment group receive the treatment and there is no control group substitution, so that the difference between the fraction of treatments and controls that receive the treatment equals 1.0.
- In practice, this difference is often well below 1.0.
- Randomization reduced and delayed receipt of training in the experimental control group but by no means eliminated it.



- With regard to substitution bias, table 15 shows that as many as 40% of the controls in some experiments received substitute services elsewhere.
- In a simple one treatment experiment with full compliance ( $\xi = 1 \Rightarrow A = 1$  and  $\xi = 0 \Rightarrow A = 0$ ), all individuals assigned to the treatment group receive the treatment and there is no control group substitution, so that the difference between the fraction of treatments and controls that receive the treatment equals 1.0.
- In practice, this difference is often well below 1.0.
- Randomization reduced and delayed receipt of training in the experimental control group but by no means eliminated it.
- Many of the treatment group members received no treatment.

- The extent of both substitution and dropout depends on the characteristics of the treatment being evaluated and the local program environment.

- The extent of both substitution and dropout depends on the characteristics of the treatment being evaluated and the local program environment.
- In the NSW study, where the treatment was relatively unique and of high enough quality to be clearly perceived as valuable by participants, dropout and substitution rates were low enough to approximate the ideal case.

- The extent of both substitution and dropout depends on the characteristics of the treatment being evaluated and the local program environment.
- In the NSW study, where the treatment was relatively unique and of high enough quality to be clearly perceived as valuable by participants, dropout and substitution rates were low enough to approximate the ideal case.
- In contrast, for the NJS and for other programs that provide low cost services widely available from other sources, substitution and dropout rates are high.

- The extent of both substitution and dropout depends on the characteristics of the treatment being evaluated and the local program environment.
- In the NSW study, where the treatment was relatively unique and of high enough quality to be clearly perceived as valuable by participants, dropout and substitution rates were low enough to approximate the ideal case.
- In contrast, for the NJS and for other programs that provide low cost services widely available from other sources, substitution and dropout rates are high.
- In the NJS, the substitution problem is accentuated by the fact that the program relied on outside vendors to provide most of its training.

- Many of these vendors, such as community colleges, provided the same training to the general public, often with subsidies from other government programs such as Pell Grants.

- Many of these vendors, such as community colleges, provided the same training to the general public, often with subsidies from other government programs such as Pell Grants.
- In addition, in order to help in recruiting sites to participate in the NJS, evaluators allowed them to provide control group members with a list of alternative training providers in the community.

- Many of these vendors, such as community colleges, provided the same training to the general public, often with subsidies from other government programs such as Pell Grants.
- In addition, in order to help in recruiting sites to participate in the NJS, evaluators allowed them to provide control group members with a list of alternative training providers in the community.
- Of the 16 sites in the NJS, 14 took advantage of this opportunity to alert control group members to substitute training opportunities.



- There are counterpart findings in the application of randomized clinical trials.

- There are counterpart findings in the application of randomized clinical trials.
- For example, ? notes that AIDS patients denied potentially life-saving drugs took steps to undo random assignment.

- There are counterpart findings in the application of randomized clinical trials.
- For example, ? notes that AIDS patients denied potentially life-saving drugs took steps to undo random assignment.
- Patients had the pills they were taking tested to see if they were getting a placebo or an unsatisfactory treatment, and were likely to drop out of the experiment in either case or to seek more effective medication, or both.

- There are counterpart findings in the application of randomized clinical trials.
- For example, ? notes that AIDS patients denied potentially life-saving drugs took steps to undo random assignment.
- Patients had the pills they were taking tested to see if they were getting a placebo or an unsatisfactory treatment, and were likely to drop out of the experiment in either case or to seek more effective medication, or both.
- In the MDRC experiment, in some sites qualified trainees found alternative avenues for securing exactly the same training presented by the same subcontractors by using other methods of financial support.

- ? discuss a variety of other problems that sometimes plague social experiments.

- ? discuss a variety of other problems that sometimes plague social experiments.
- Our discussion up to this point has focused on point identification of parameters over the empirical supports.

- ? discuss a variety of other problems that sometimes plague social experiments.
- Our discussion up to this point has focused on point identification of parameters over the empirical supports.
- A large and emerging literature produces bounds on the parameters and distributions when point identification is not possible.

- ? discuss a variety of other problems that sometimes plague social experiments.
- Our discussion up to this point has focused on point identification of parameters over the empirical supports.
- A large and emerging literature produces bounds on the parameters and distributions when point identification is not possible.
- We now consider bounds on the parameters within the framework of economic models of choice and the MTE.



## Bounding and Sensitivity Analysis

- Thus far we have assumed full support conditions and have presented conditions for identification over those supports.

## Bounding and Sensitivity Analysis

- Thus far we have assumed full support conditions and have presented conditions for identification over those supports.
- We now consider partial identification in the context of the MTE framework.

## Bounding and Sensitivity Analysis

- Thus far we have assumed full support conditions and have presented conditions for identification over those supports.
- We now consider partial identification in the context of the MTE framework.
- We return to the two-outcome model to develop the basic approach in a simpler setting.

- The central evaluation problem is that we observe the distribution of  $(Y, D, X, Z) = (DY_1 + (1 - D)Y_0, D, X, Z)$ , but do not observe the distribution of all of the components that comprise it  $(Y_1, Y_0, D, X, Z)$ .

- The central evaluation problem is that we observe the distribution of  $(Y, D, X, Z) = (DY_1 + (1 - D)Y_0, D, X, Z)$ , but do not observe the distribution of all of the components that comprise it  $(Y_1, Y_0, D, X, Z)$ .
- Let  $\eta$  denote a distribution for  $(Y_1, Y_0, D, X, Z)$ , and let it be known that  $\eta$  belongs to the class  $\mathcal{H} \subset \mathcal{F}$ , where  $\mathcal{F}$  is the space of all probability distributions on  $(Y_1, Y_0, D, X, Z)$ .

- The central evaluation problem is that we observe the distribution of  $(Y, D, X, Z) = (DY_1 + (1 - D)Y_0, D, X, Z)$ , but do not observe the distribution of all of the components that comprise it  $(Y_1, Y_0, D, X, Z)$ .
- Let  $\eta$  denote a distribution for  $(Y_1, Y_0, D, X, Z)$ , and let it be known that  $\eta$  belongs to the class  $\mathcal{H} \subset \mathcal{F}$ , where  $\mathcal{F}$  is the space of all probability distributions on  $(Y_1, Y_0, D, X, Z)$ .
- Let  $P_\eta$  denote the resulting distribution of  $(DY_1 + (1 - D)Y_0, D, X, Z)$  if  $\eta$  is the distribution for  $(Y_1, Y_0, D, X, Z)$ .

- Let  $\eta^0$  and  $P_{\eta^0}$  denote the corresponding true distributions.

- Let  $\eta^0$  and  $P_{\eta^0}$  denote the corresponding true distributions.
- Knowledge of the distribution of  $(DY_1 + (1 - D)Y_0, D, X, Z)$  allows us to infer that  $\eta$  lies in the set  $\{\eta \in \mathcal{H} : P_\eta = P_{\eta^0}\}$ .



- Let  $\eta^0$  and  $P_{\eta^0}$  denote the corresponding true distributions.
- Knowledge of the distribution of  $(DY_1 + (1 - D)Y_0, D, X, Z)$  allows us to infer that  $\eta$  lies in the set  $\{\eta \in \mathcal{H} : P_\eta = P_{\eta^0}\}$ .
- All elements of  $\{\eta \in \mathcal{H} : P_\eta = P_{\eta^0}\}$  are consistent with the true distribution of the observed data.

- Let  $\mathcal{H}^0 = \{\eta \in \mathcal{H} : P_\eta = P_{\eta^0}\}$ .

- Let  $\mathcal{H}^0 = \{\eta \in \mathcal{H} : P_\eta = P_{\eta^0}\}$ .
- Let  $E_\eta$  denote expectation with respect to the measure  $\eta$ , i.e.,  $E_\eta(A) = \int A d\eta$ , so that  $E(A) = E_{\eta^0}(A)$ .

- Let  $\mathcal{H}^0 = \{\eta \in \mathcal{H} : P_\eta = P_{\eta^0}\}$ .
- Let  $E_\eta$  denote expectation with respect to the measure  $\eta$ , i.e.,  $E_\eta(A) = \int A d\eta$ , so that  $E(A) = E_{\eta^0}(A)$ .
- Consider inference for ATE,  $E(Y_1 - Y_0)$ .

- Let  $\mathcal{H}^0 = \{\eta \in \mathcal{H} : P_\eta = P_{\eta^0}\}$ .
- Let  $E_\eta$  denote expectation with respect to the measure  $\eta$ , i.e.,  $E_\eta(A) = \int A d\eta$ , so that  $E(A) = E_{\eta^0}(A)$ .
- Consider inference for ATE,  $E(Y_1 - Y_0)$ .
- Knowledge of the distribution of the observed variables allows us to infer that

$$E(Y_1 - Y_0) \in \{E_\eta(Y_1 - Y_0) : \eta \in \mathcal{H}^0\}.$$

- The identification analyses of the previous sections proceed by imposing sufficient restrictions on  $\mathcal{H}$  such that  $\{E_{\eta}(Y_1 - Y_0) : \eta \in \mathcal{H}^0\}$  contains only one element and thus  $E(Y_1 - Y_0)$  is point identified.

- The identification analyses of the previous sections proceed by imposing sufficient restrictions on  $\mathcal{H}$  such that  $\{E_{\eta}(Y_1 - Y_0) : \eta \in \mathcal{H}^0\}$  contains only one element and thus  $E(Y_1 - Y_0)$  is point identified.
- Bounding analysis proceeds by finding a set  $\mathcal{B}$  such that  $\mathcal{B} \supseteq \{E_{\eta}(Y_1 - Y_0) : \eta \in \mathcal{H}^0\}$ .

- The identification analyses of the previous sections proceed by imposing sufficient restrictions on  $\mathcal{H}$  such that  $\{E_{\eta}(Y_1 - Y_0) : \eta \in \mathcal{H}^0\}$  contains only one element and thus  $E(Y_1 - Y_0)$  is point identified.
- Bounding analysis proceeds by finding a set  $\mathcal{B}$  such that  $\mathcal{B} \supseteq \{E_{\eta}(Y_1 - Y_0) : \eta \in \mathcal{H}^0\}$ .
- One goal of bounding analysis is to construct  $\mathcal{B}$  such that  $\mathcal{B} = \{E_{\eta}(Y_1 - Y_0) : \eta \in \mathcal{H}^0\}$  in which case the bounds are said to be *sharp*.



- If the bounds are sharp, then the bounds exploit all information and no smaller bounds can be constructed without imposing additional structure.

- If the bounds are sharp, then the bounds exploit all information and no smaller bounds can be constructed without imposing additional structure.
- In contrast, if  $\{E_\eta(Y_1 - Y_0) : \eta \in \mathcal{H}^0\}$  is a proper subset of  $\mathcal{B}$ , then smaller bounds can be constructed.

- If the bounds are sharp, then the bounds exploit all information and no smaller bounds can be constructed without imposing additional structure.
- In contrast, if  $\{E_\eta(Y_1 - Y_0) : \eta \in \mathcal{H}^0\}$  is a proper subset of  $\mathcal{B}$ , then smaller bounds can be constructed.
- In every example we consider, the set  $\{E_\eta(Y_1 - Y_0) : \eta \in \mathcal{H}^0\}$  is a closed interval, so that
$$\{E_\eta(Y_1 - Y_0) : \eta \in \mathcal{H}^0\} = \left[ \inf_{\eta \in \mathcal{H}^0} E_\eta(Y_1 - Y_0), \sup_{\eta \in \mathcal{H}^0} E_\eta(Y_1 - Y_0) \right].$$

- *Sensitivity* analysis is a commonly used procedure.

- *Sensitivity* analysis is a commonly used procedure.
- It varies the parameters fixed in a model and determines the sensitivity of estimates to the perturbations of the parameter.

- *Sensitivity* analysis is a commonly used procedure.
- It varies the parameters fixed in a model and determines the sensitivity of estimates to the perturbations of the parameter.
- Sensitivity analysis is formally equivalent to bounding.

- *Sensitivity* analysis is a commonly used procedure.
- It varies the parameters fixed in a model and determines the sensitivity of estimates to the perturbations of the parameter.
- Sensitivity analysis is formally equivalent to bounding.
- In particular, in sensitivity analysis, one parameterizes  $\eta$  and then constructs bounds based on letting the parameters vary over some set.

- Parameterize  $\eta$  as  $\eta(\theta)$  for some parameter vector  $\theta \in \Theta$ , and let  $\theta^0$  be the “true” parameter value so that  $\eta^0 = \eta(\theta^0)$ .



- Parameterize  $\eta$  as  $\eta(\theta)$  for some parameter vector  $\theta \in \Theta$ , and let  $\theta^0$  be the “true” parameter value so that  $\eta^0 = \eta(\theta^0)$ .
- $\theta$  is typically finite dimensional, though it need not be.

- Parameterize  $\eta$  as  $\eta(\theta)$  for some parameter vector  $\theta \in \Theta$ , and let  $\theta^0$  be the “true” parameter value so that  $\eta^0 = \eta(\theta^0)$ .
- $\theta$  is typically finite dimensional, though it need not be.
- Let  $\Theta^0 = \{\theta \in \Theta : P_{\eta(\theta)} = P_{\eta(\theta^0)}\}$ .

- Parameterize  $\eta$  as  $\eta(\theta)$  for some parameter vector  $\theta \in \Theta$ , and let  $\theta^0$  be the “true” parameter value so that  $\eta^0 = \eta(\theta^0)$ .
- $\theta$  is typically finite dimensional, though it need not be.
- Let  $\Theta^0 = \{\theta \in \Theta : P_{\eta(\theta)} = P_{\eta(\theta^0)}\}$ .
- If  $\theta$  is point identified given the observed variables, then  $\Theta^0$  will contain only one element, but if not all parameters are identified given the observed data then  $\Theta$  will contain more than one element.

- Consider

$$\{E_{\eta(\theta)}(Y_1 - Y_0) : \theta \in \Theta^0\}.$$

- Consider

$$\{E_{\eta(\theta)}(Y_1 - Y_0) : \theta \in \Theta^0\}.$$

- This can trivially be seen as a special case of bounding analysis by taking  $\mathcal{H} = \{\eta(\theta) : \theta \in \Theta\}$  and  $\mathcal{H}^0 = \{\eta(\theta) : \theta \in \Theta^0\}$ .

- Consider

$$\{E_{\eta(\theta)}(Y_1 - Y_0) : \theta \in \Theta^0\}.$$

- This can trivially be seen as a special case of bounding analysis by taking  $\mathcal{H} = \{\eta(\theta) : \theta \in \Theta\}$  and  $\mathcal{H}^0 = \{\eta(\theta) : \theta \in \Theta^0\}$ .
- Likewise, by taking a proper parameterization, any bounding analysis can be seen as a special case of sensitivity analysis.

- We consider bounds on ATE.

- We consider bounds on ATE.
- The corresponding bounds on treatment on the treated follow with trivial modifications.



- We consider bounds on ATE.
- The corresponding bounds on treatment on the treated follow with trivial modifications.
- We focus on bounds that exploit instrumental variable type assumptions or latent index assumptions, and we do not attempt to survey the entire literature on bounds.

- We consider bounds on ATE.
- The corresponding bounds on treatment on the treated follow with trivial modifications.
- We focus on bounds that exploit instrumental variable type assumptions or latent index assumptions, and we do not attempt to survey the entire literature on bounds.
- We begin by describing the bounds that only assume that the outcome variables are bounded.

- We consider bounds on ATE.
- The corresponding bounds on treatment on the treated follow with trivial modifications.
- We focus on bounds that exploit instrumental variable type assumptions or latent index assumptions, and we do not attempt to survey the entire literature on bounds.
- We begin by describing the bounds that only assume that the outcome variables are bounded.
- We then consider imposing additional assumptions.

- We consider imposing the assumption of comparative advantage in the decision rule, then consider instead imposing an instrumental variables type assumption, and conclude by considering the combination of comparative advantage and instrumental variables assumptions.

- We consider imposing the assumption of comparative advantage in the decision rule, then consider instead imposing an instrumental variables type assumption, and conclude by considering the combination of comparative advantage and instrumental variables assumptions.
- We examine the relative power of these alternative assumptions to narrow the very wide bounds that result from only imposing that the outcome variables are bounded.

## Outcome is Bounded

- We first consider bounds on  $E(Y_1 - Y_0)$  that only assume that the outcomes be bounded.

## Outcome is Bounded

- We first consider bounds on  $E(Y_1 - Y_0)$  that only assume that the outcomes be bounded.
- We consider this case as a point of contrast for the later bounds that exploit instrumental variable conditions, and also for the pedagogical purpose of showing the bounding methodology in a simple context.

We impose that the outcomes are bounded with probability 1,  
**Assumption B: Outcome is Bounded** For  $j = 0, 1$ ,



- In our notation this corresponds to

$$\mathcal{H} = \{\eta \in \mathcal{F} : \eta[y^l \leq Y_1 \leq y^u] = 1, \eta[y^l \leq Y_0 \leq y^u] = 1\}.$$

- In our notation this corresponds to

$$\mathcal{H} = \{\eta \in \mathcal{F} : \eta[y^l \leq Y_1 \leq y^u] = 1, \eta[y^l \leq Y_0 \leq y^u] = 1\}.$$

- For example, if  $Y$  is an indicator variable, then the bounds are  $y^l = 0$  and  $y^u = 1$ .

- Following ? and ?, use the law of iterated expectations to obtain,

$$E(Y_1) = \Pr[D = 1]E(Y_1|D = 1) + (1 - \Pr[D = 1])E(Y_1|D = 0)$$

$$E(Y_0) = \Pr[D = 1]E(Y_0|D = 1) + (1 - \Pr[D = 1])E(Y_0|D = 0).$$

- $\Pr[D = 1]$ ,  $E(Y_1|D = 1)$ , and  $E(Y_0|D = 0)$  are identified, while  $E(Y_0|D = 1)$  and  $E(Y_1|D = 0)$  are bounded by  $y^l$  and  $y^u$ , so that

$$\Pr[D = 1]E(Y_1 | D = 1) + (1 - \Pr[D = 1])y^l$$

$$\leq E(Y_1) \leq \Pr[D = 1]E(Y_1 | D = 1) + (1 - \Pr[D = 1])y^u,$$

$$\Pr[D = 1]y^l + (1 - \Pr[D = 1])E(Y_0 | D = 0)$$

$$\leq E(Y_0) \leq \Pr[D = 1]y^u + (1 - \Pr[D = 1])E(Y_0 | D = 0)$$

- Thus

$$\mathcal{B} = [B^L, B^U],$$

with

$$B^L = (\Pr[D = 1]E(Y | D = 1) + (1 - \Pr[D = 1])y^l) \\ - (\Pr[D = 1]y^u + (1 - \Pr[D = 1])E(Y | D = 0)),$$

$$B^U = (\Pr[D = 1]E(Y | D = 1) + (1 - \Pr[D = 1])y^u) \\ - (\Pr[D = 1]y^l + (1 - \Pr[D = 1])E(Y | D = 0))$$

with the width of these bounds given by

$$B^U - B^L = y^u - y^l.$$

- Thus

$$B = [B^L, B^U],$$

with

$$B^L = (\Pr[D = 1]E(Y | D = 1) + (1 - \Pr[D = 1])y^l) \\ - (\Pr[D = 1]y^u + (1 - \Pr[D = 1])E(Y | D = 0)),$$

$$B^U = (\Pr[D = 1]E(Y | D = 1) + (1 - \Pr[D = 1])y^u) \\ - (\Pr[D = 1]y^l + (1 - \Pr[D = 1])E(Y | D = 0))$$

with the width of these bounds given by

$$B^U - B^L = y^u - y^l.$$

- For example, if  $Y = 0, 1$ , then the width of the bounds equals 1,  $B^U - B^L = 1$ .

- These bounds are sharp.

- These bounds are sharp.
- To show this, for any  $M \in [B^L, B^U]$ , one can trivially construct a distribution  $\eta$  of  $(Y_0, Y_1, D)$  which is consistent with the observed data, consistent with the restriction that the outcomes are bounded, and for which  $E_\eta(Y_1 - Y_0) = M$ , thus showing that  $M \in [B^L, B^U]$ .



- These bounds are sharp.
- To show this, for any  $M \in [B^L, B^U]$ , one can trivially construct a distribution  $\eta$  of  $(Y_0, Y_1, D)$  which is consistent with the observed data, consistent with the restriction that the outcomes are bounded, and for which  $E_\eta(Y_1 - Y_0) = M$ , thus showing that  $M \in [B^L, B^U]$ .
- Since this is true for any  $M \in [B^L, B^U]$ , it follows that  $[B^L, B^U] \subseteq \{E_\eta(Y_1 - Y_0) : \eta \in \mathcal{H}_0\}$ .

- These bounds are sharp.
- To show this, for any  $M \in [B^L, B^U]$ , one can trivially construct a distribution  $\eta$  of  $(Y_0, Y_1, D)$  which is consistent with the observed data, consistent with the restriction that the outcomes are bounded, and for which  $E_\eta(Y_1 - Y_0) = M$ , thus showing that  $M \in [B^L, B^U]$ .
- Since this is true for any  $M \in [B^L, B^U]$ , it follows that  $[B^L, B^U] \subseteq \{E_\eta(Y_1 - Y_0) : \eta \in \mathcal{H}_0\}$ .
- Since we have already shown that  $[B^L, B^U]$  are valid bounds,  $[B^L, B^U] \supseteq \{E_\eta(Y_1 - Y_0) : \eta \in \mathcal{H}_0\}$ , we conclude that  $[B^L, B^U] = \{E_\eta(Y_1 - Y_0) : \eta \in \mathcal{H}_0\}$  and thus that the bounds are sharp.

- This illustrates a common technique towards the construction of sharp bounds: in a first step, construct a natural set of bounds, and in a second step, use a proof by construction to show that the bounds are sharp.

- This illustrates a common technique towards the construction of sharp bounds: in a first step, construct a natural set of bounds, and in a second step, use a proof by construction to show that the bounds are sharp.
- Note the following features of these bounds.

- This illustrates a common technique towards the construction of sharp bounds: in a first step, construct a natural set of bounds, and in a second step, use a proof by construction to show that the bounds are sharp.
- Note the following features of these bounds.
- First, as noted by ?, these bounds always include zero.

- Thus, bounds that only exploit that the outcomes are bounded can never reject the null of zero average treatment effect.

- Thus, bounds that only exploit that the outcomes are bounded can never reject the null of zero average treatment effect.
- The bounds themselves depend on the data, but the width of the bounds,  $B^U - B^L = y^u - y^l$ , is completely driven by the assumed bounds on  $Y_1, Y_0$ .

- Thus, bounds that only exploit that the outcomes are bounded can never reject the null of zero average treatment effect.
- The bounds themselves depend on the data, but the width of the bounds,  $B^U - B^L = y^u - y^l$ , is completely driven by the assumed bounds on  $Y_1, Y_0$ .
- For example, if  $Y_1$  and  $Y_0$  are binary, the width of the bounds is always 1.



## Latent Index Model: Roy Model

- The bounds that only impose that the outcomes are bounded are typically very wide, never provide point identification, and can never reject the null of zero average treatment effect.

## Latent Index Model: Roy Model

- The bounds that only impose that the outcomes are bounded are typically very wide, never provide point identification, and can never reject the null of zero average treatment effect.
- This lack of identifying power raises the question of whether one can impose additional structure to narrow the bounds.

## Latent Index Model: Roy Model

- The bounds that only impose that the outcomes are bounded are typically very wide, never provide point identification, and can never reject the null of zero average treatment effect.
- This lack of identifying power raises the question of whether one can impose additional structure to narrow the bounds.
- The central issue with bounding analysis is to explore the trade-off between assumptions and width of the resulting bounds.

- In this section, we discuss bounds that follow from maintaining Assumption B, that the outcomes are bounded, but also add the assumption of a Roy model for selection into treatment.

- In this section, we discuss bounds that follow from maintaining Assumption B, that the outcomes are bounded, but also add the assumption of a Roy model for selection into treatment.
- Such an assumption substantially narrows the width of the bounds compared to only imposing that the outcomes themselves are bounded, but does not provide point identification.

- Again impose Assumption B: the outcomes are bounded.

## Assumption RM: Roy Model

$$D = \mathbf{1}[Y_1 \geq Y_0]. \quad (77)$$

- Again impose Assumption B: the outcomes are bounded.
- In addition, assume a model of comparative advantage, in particular,

Assumption RM: Roy Model

$$D = \mathbf{1}[Y_1 \geq Y_0]. \quad (77)$$

- Restriction RM imposes a special case of a latent index model,  $D = \mathbf{1}[Y^* \geq 0]$  with  $Y^* = Y_1 - Y_0$ .



- Restriction RM imposes a special case of a latent index model,  $D = \mathbf{1}[Y^* \geq 0]$  with  $Y^* = Y_1 - Y_0$ .
- Using the assumption of a Roy model while maintaining the assumption that the outcomes are bounded, we can narrow the bounds compared to the case where we only imposed that the outcomes are bounded.

- Restriction RM imposes a special case of a latent index model,  $D = \mathbf{1}[Y^* \geq 0]$  with  $Y^* = Y_1 - Y_0$ .
- Using the assumption of a Roy model while maintaining the assumption that the outcomes are bounded, we can narrow the bounds compared to the case where we only imposed that the outcomes are bounded.
- ? constructs the sharp bounds for the competing risks model, which is formally equivalent to a Roy model.

- Restriction RM imposes a special case of a latent index model,  $D = \mathbf{1}[Y^* \geq 0]$  with  $Y^* = Y_1 - Y_0$ .
- Using the assumption of a Roy model while maintaining the assumption that the outcomes are bounded, we can narrow the bounds compared to the case where we only imposed that the outcomes are bounded.
- ? constructs the sharp bounds for the competing risks model, which is formally equivalent to a Roy model.
- ? constructs the same bounds for the Roy model.

- Following ? and ?, we have that

$$\begin{aligned} E[Y_1|D = 1] &= E[Y_1|Y_0 \leq Y_1] \\ &\geq E[Y_0|Y_0 \leq Y_1] \\ &= E[Y_0|D = 1] \end{aligned}$$

and by a parallel argument,  $E[Y_0|D = 0] \geq E[Y_1|D = 0]$ .

- Following ? and ?, we have that

$$\begin{aligned} E[Y_1|D = 1] &= E[Y_1|Y_0 \leq Y_1] \\ &\geq E[Y_0|Y_0 \leq Y_1] \\ &= E[Y_0|D = 1] \end{aligned}$$

and by a parallel argument,  $E[Y_0|D = 0] \geq E[Y_1|D = 0]$ .

- We thus have upper bounds on  $E(Y_0|D = 1)$  and  $E(Y_1|D = 0)$ .

- Following ? and ?, we have that

$$\begin{aligned} E[Y_1|D = 1] &= E[Y_1|Y_0 \leq Y_1] \\ &\geq E[Y_0|Y_0 \leq Y_1] \\ &= E[Y_0|D = 1] \end{aligned}$$

and by a parallel argument,  $E[Y_0|D = 0] \geq E[Y_1|D = 0]$ .

- We thus have upper bounds on  $E(Y_0|D = 1)$  and  $E(Y_1|D = 0)$ .
- The lower bounds on  $E[Y_0|D = 1]$  and  $E[Y_1|D = 0]$  are the same as for the bounds that only imposed that the outcomes are bounded.

We then have:

$$E(Y_1 - Y_0) \in \mathcal{B} \equiv [B^L, B^U],$$

with

$$B^L = \left( \Pr[D = 1]E(Y|D = 1) + (1 - \Pr[D = 1])y' \right) \\ - \left( \Pr[D = 1]E(Y|D = 1) + (1 - \Pr[D = 1])E(Y|D = 0) \right)$$

$$B^U = \left( \Pr[D = 1]E(Y|D = 1) + (1 - \Pr[D = 1])E(Y|D = 0) \right) \\ - \left( \Pr[D = 1]y' + (1 - \Pr[D = 1])E(Y|D = 0) \right),$$

and we can rewrite these bounds as

$$B^L = \left( 1 - \Pr[D = 1] \right) \left( y' - E(Y|D = 0) \right) \\ B^U = \Pr[D = 1] \left( E(Y|D = 1) - y' \right),$$

with the width of the bounds given by

$$B^U - B^L = E(Y) - y'.$$

- For example, if  $Y = 0, 1$ , then the width of the bounds is given by  $B^U - B^L = \Pr(Y = 1)$ .



- For example, if  $Y = 0, 1$ , then the width of the bounds is given by  $B^U - B^L = \Pr(Y = 1)$ .
- Following an argument similar to that presented in the previous section, one can show that these bounds are sharp.

- For example, if  $Y = 0, 1$ , then the width of the bounds is given by  $B^U - B^L = \Pr(Y = 1)$ .
- Following an argument similar to that presented in the previous section, one can show that these bounds are sharp.
- Note the following features of these bounds.

- For example, if  $Y = 0, 1$ , then the width of the bounds is given by  $B^U - B^L = \Pr(Y = 1)$ .
- Following an argument similar to that presented in the previous section, one can show that these bounds are sharp.
- Note the following features of these bounds.
- First, the bounds do not involve  $y^u$ , and actually the same bounds will hold if we were to weaken the maintained assumption that  $\Pr[y^l \leq Y_j \leq y^u] = 1$  for  $j = 0, 1$ , to instead only require that  $\Pr[y^l \leq Y_j] = 1$ .

- For example, if  $Y = 0, 1$ , then the width of the bounds is given by  $B^U - B^L = \Pr(Y = 1)$ .
- Following an argument similar to that presented in the previous section, one can show that these bounds are sharp.
- Note the following features of these bounds.
- First, the bounds do not involve  $y^u$ , and actually the same bounds will hold if we were to weaken the maintained assumption that  $\Pr[y^l \leq Y_j \leq y^u] = 1$  for  $j = 0, 1$ , to instead only require that  $\Pr[y^l \leq Y_j] = 1$ .
- The width of the bounds imposing comparative advantage are  $E(Y) - y^l$ , so that the bounds will never provide point identification (as long as  $E(Y) > y^l$ ).

- For example, if  $Y$  is binary, the width of the bounds is  $\Pr[Y = 1]$ , the bounds will not provide point identification unless all individuals have  $Y = 0$ .

- For example, if  $Y$  is binary, the width of the bounds is  $\Pr[Y = 1]$ , the bounds will not provide point identification unless all individuals have  $Y = 0$ .
- However, the bounds will always improve upon the bounds that impose only that the outcome is bounded – imposing comparative advantage shrinks the width of the bounds from  $y^u - y^l$  to  $E(Y) - y^l$ , thus shrinking the bounds by an amount equal to  $y^u - E(Y)$ .

- For example, if  $Y$  is binary, the width of the bounds is  $\Pr[Y = 1]$ , the bounds will not provide point identification unless all individuals have  $Y = 0$ .
- However, the bounds will always improve upon the bounds that impose only that the outcome is bounded – imposing comparative advantage shrinks the width of the bounds from  $y^u - y^l$  to  $E(Y) - y^l$ , thus shrinking the bounds by an amount equal to  $y^u - E(Y)$ .
- For example, if  $Y$  is binary, then imposing the bounds shrinks the width of the bounds from 1 to  $\Pr[Y = 1]$ .

- For example, if  $Y$  is binary, the width of the bounds is  $\Pr[Y = 1]$ , the bounds will not provide point identification unless all individuals have  $Y = 0$ .
- However, the bounds will always improve upon the bounds that impose only that the outcome is bounded – imposing comparative advantage shrinks the width of the bounds from  $y^u - y^l$  to  $E(Y) - y^l$ , thus shrinking the bounds by an amount equal to  $y^u - E(Y)$ .
- For example, if  $Y$  is binary, then imposing the bounds shrinks the width of the bounds from 1 to  $\Pr[Y = 1]$ .
- Finally, note that the bounds will always include zero, so that imposing comparative advantage does not by itself allow one to ever reject the null of zero average treatment effect.



## Bounds that Exploit an Instrument

- The previous section considered bounds that exploit knowledge of the selection process, in particular that selection is determined by a Roy model.

## Bounds that Exploit an Instrument

- The previous section considered bounds that exploit knowledge of the selection process, in particular that selection is determined by a Roy model.
- An alternative way to narrow the bounds over simply imposing that the outcome is bounded is to assume access to an instrument.

## Bounds that Exploit an Instrument

- The previous section considered bounds that exploit knowledge of the selection process, in particular that selection is determined by a Roy model.
- An alternative way to narrow the bounds over simply imposing that the outcome is bounded is to assume access to an instrument.
- We now discuss bounds with various types of instrumental variables assumptions.

- We begin with the ? analysis for bounds that exploit a mean-independence condition, then consider the ? analysis for bounds that exploit a full statistical independence condition, and finally conclude with a discussion of ? who combine an instrumental variable assumption with a nonparametric selection model.

## Instrumental Variables: Mean Independence Condition

- Again impose Assumption B so that the outcomes are bounded.

## Instrumental Variables: Mean Independence Condition

- Again impose Assumption B so that the outcomes are bounded.
- In addition, following ?, impose a mean-independence assumption:

## Assumption IV:

$$E(Y_1|Z = z) = E(Y_1)$$

$$E(Y_0|Z = z) = E(Y_0)$$

for  $z \in \mathcal{Z}$  where  $\mathcal{Z}$  denotes the support of the distribution of  $Z$ .

- For any  $z \in \mathcal{Z}$ , following the exact same series of steps as for the bounds that only imposed Assumption [B], we have that

$$E(DY|Z = z) + (1 - P(z))y^l \leq E(Y_1|Z = z) \leq E(DY|Z = z) + (1 - P(z))y^u.$$



- For any  $z \in \mathcal{Z}$ , following the exact same series of steps as for the bounds that only imposed Assumption [B], we have that

$$E(DY|Z = z) + (1 - P(z))y^l \leq E(Y_1|Z = z) \leq E(DY|Z = z) + (1 - P(z))y^u.$$

- By the IV assumption, we have  $E(Y_1|Z = z) = E(Y_1)$ .

- For any  $z \in \mathcal{Z}$ , following the exact same series of steps as for the bounds that only imposed Assumption [B], we have that

$$E(DY|Z = z) + (1 - P(z))y^l \leq E(Y_1|Z = z) \leq E(DY|Z = z) + (1 - P(z))y^u.$$

- By the IV assumption, we have  $E(Y_1|Z = z) = E(Y_1)$ .
- Since these bounds hold for any  $z \in \mathcal{Z}$ , we have

$$\sup_{z \in \mathcal{Z}} \{E(DY|Z = z) + (1 - P(z))y^l\} \leq E(Y_1) \leq \inf_{z \in \mathcal{Z}} \{E(DY|Z = z) + (1 - P(z))y^u\}.$$

- Applying the same analysis for  $E(Y_0)$ , we have

$$E(Y_1 - Y_0) \in \mathcal{B} = [B^L, B^U],$$

with  $B^L = \sup_{z \in \mathcal{Z}} \{E(DY | Z = z) + (1 - P(z))y^l\} - \inf_{z \in \mathcal{Z}} \{(E((1 - D)Y | Z = z) + P(z)y^u)\}$ ,

$B^U = \inf_{z \in \mathcal{Z}} \{E(DY | Z = z) + (1 - P(z))y^u\} - \sup_{z \in \mathcal{Z}} \{(E((1 - D)Y | Z = z) + P(z)y^l)\}$ .

- As discussed by ?, these bounds are sharp under the mean-independence condition.

- As discussed by Angriston and Pischke (2009), these bounds are sharp under the mean-independence condition.
- As noted by Angriston and Pischke (2009), these bounds do not necessarily include zero, so that it may be possible to use the bounds to test the null of zero average treatment effect.

- As discussed by ?, these bounds are sharp under the mean-independence condition.
- As noted by ?, these bounds do not necessarily include zero, so that it may be possible to use the bounds to test the null of zero average treatment effect.
- Let  $p^u = \sup_{z \in \mathcal{Z}} \Pr[D = 1 \mid Z = z]$ ,  
 $p^l = \inf_{z \in \mathcal{Z}} \Pr[D = 1 \mid Z = z]$ .

- As discussed by ? , these bounds are sharp under the mean-independence condition.
- As noted by ? , these bounds do not necessarily include zero, so that it may be possible to use the bounds to test the null of zero average treatment effect.
- Let  $p^u = \sup_{z \in \mathcal{Z}} \Pr[D = 1 \mid Z = z]$ ,  
 $p^l = \inf_{z \in \mathcal{Z}} \Pr[D = 1 \mid Z = z]$ .
- A trivial modification to Corollary 1 and Corollary 2 of Proposition 6 of ? shows that

- 1  $p^u \geq \frac{1}{2}$  and  $p^l \geq \frac{1}{2}$  is a necessary condition for  $B^L = B^U$ , i.e., for point identification from the mean independence condition.
- 2 If  $Y_1, Y_0$  are independent of  $D$ , then the width of the IV-bounds is  $((1 - p^u) + p^l)(y^u - y^l)$ . Thus, if  $Y_1, Y_0$  are independent of  $D$ , the bounds will collapse to point identification if and only if  $p^u = 1, p^l = 0$ .



- Note that it is neither necessary nor sufficient for  $P(z)$  to be a nontrivial function of  $z$  for these bounds to improve upon the bounds that only imposed that the outcome is bounded.

- Note that it is neither necessary nor sufficient for  $P(z)$  to be a nontrivial function of  $z$  for these bounds to improve upon the bounds that only imposed that the outcome is bounded.
- Likewise, comparing these bounds to the comparative advantage bounds shows that neither set of bounds will in general be narrower than the other.

- Note that it is neither necessary nor sufficient for  $P(z)$  to be a nontrivial function of  $z$  for these bounds to improve upon the bounds that only imposed that the outcome is bounded.
- Likewise, comparing these bounds to the comparative advantage bounds shows that neither set of bounds will in general be narrower than the other.
- Finally, note that these bounds are relatively complicated, and to evaluate the bounds and the width of the bounds requires use of  $P(z)$ ,  $E(YD \mid Z = z)$ , and  $E(Y(1 - D) \mid Z = z)$  for all  $z \in \mathcal{Z}$ .

## Instrumental Variables: Statistical Independence Condition

- While Manski constructs sharp bounds for mean-independence conditions, ? construct sharp bounds for the statistical independence condition for the case where  $Y$  and  $Z$  are binary.

## Instrumental Variables: Statistical Independence Condition

- While Manski constructs sharp bounds for mean-independence conditions, ? construct sharp bounds for the statistical independence condition for the case where  $Y$  and  $Z$  are binary.
- Balke and Pearl impose the same independence condition as the ? LATE independence condition.

## Instrumental Variables: Statistical Independence Condition

- While Manski constructs sharp bounds for mean-independence conditions, ? construct sharp bounds for the statistical independence condition for the case where  $Y$  and  $Z$  are binary.
- Balke and Pearl impose the same independence condition as the ? LATE independence condition.
- In particular, let  $D_0, D_1$  denote the counterfactual choices that would have been made had  $Z$  been set exogenously to 0 and 1 respectively, and impose the following assumption:

## Instrumental Variables: Statistical Independence Condition

- While Manski constructs sharp bounds for mean-independence conditions, ? construct sharp bounds for the statistical independence condition for the case where  $Y$  and  $Z$  are binary.
- Balke and Pearl impose the same independence condition as the ? LATE independence condition.
- In particular, let  $D_0, D_1$  denote the counterfactual choices that would have been made had  $Z$  been set exogenously to 0 and 1 respectively, and impose the following assumption:

Assumption: IV-BP

$$(Y_0, Y_1, D_0, D_1) \perp\!\!\!\perp Z$$

- Note that this strengthens the Manski conditions not only in imposing that potential outcomes are statistically independent of  $Z$  instead of mean-independent of  $Z$ , but also imposing that the counterfactual choices are independent of  $Z$ .



- For the case of  $Z$  and  $Y$  binary, Balke and Pearl manage to transform the problem of constructing sharp bounds into a linear programming problem.

- For the case of  $Z$  and  $Y$  binary, Balke and Pearl manage to transform the problem of constructing sharp bounds into a linear programming problem.
- Assuming that the identified set is a closed interval, the sharp bounds are by definition  $[B^L, B^U]$  with

$$B^L = \inf_{\eta \in \mathcal{H}^0} E_{\eta}(Y_1 - Y_0)$$

$$B^U = \sup_{\eta \in \mathcal{H}^0} E_{\eta}(Y_1 - Y_0).$$

- For the case of  $Z$  and  $Y$  binary, Balke and Pearl manage to transform the problem of constructing sharp bounds into a linear programming problem.
- Assuming that the identified set is a closed interval, the sharp bounds are by definition  $[B^L, B^U]$  with

$$B^L = \inf_{\eta \in \mathcal{H}^0} E_{\eta}(Y_1 - Y_0)$$

$$B^U = \sup_{\eta \in \mathcal{H}^0} E_{\eta}(Y_1 - Y_0).$$

- In general, the constrained set of distributions,  $\eta \in \mathcal{H}^0$ , may be high dimensional and non-convex.

- Using the assumption that  $Z$  and  $Y$  are binary, they transform the problem into the minimization of a linear function over a finite dimensional vector space subject to a set of linear constraints.

- Using the assumption that  $Z$  and  $Y$  are binary, they transform the problem into the minimization of a linear function over a finite dimensional vector space subject to a set of linear constraints.
- The resulting bounds are somewhat complex.

- Using the assumption that  $Z$  and  $Y$  are binary, they transform the problem into the minimization of a linear function over a finite dimensional vector space subject to a set of linear constraints.
- The resulting bounds are somewhat complex.
- For some distributions of the observed data, they will coincide with the Manski mean-independence bounds, but for other distributions of the observed data they will be narrower than the Manski mean-independence bounds.

- Using the assumption that  $Z$  and  $Y$  are binary, they transform the problem into the minimization of a linear function over a finite dimensional vector space subject to a set of linear constraints.
- The resulting bounds are somewhat complex.
- For some distributions of the observed data, they will coincide with the Manski mean-independence bounds, but for other distributions of the observed data they will be narrower than the Manski mean-independence bounds.
- Thus, imposing statistical independence does narrow the bounds over the mean independence bounds.

- It is not immediately clear how to generalize the Balke and Pearl analysis to distributions with continuous  $Z$  or  $Y$ , or how to construct sharp bounds under the statistical independence condition for  $Z$  or  $Y$  continuous.



- It is not immediately clear how to generalize the Balke and Pearl analysis to distributions with continuous  $Z$  or  $Y$ , or how to construct sharp bounds under the statistical independence condition for  $Z$  or  $Y$  continuous.
- The appropriate generalization of Balke and Pearl's analysis to a more general setting remains an open question.

## Instrumental Variables: Nonparametric Selection Model/LATE Conditions

- We started with the mean independence version of the instrumental variables condition, and then discussed strengthening the instrumental variables condition to full independence in the special case where  $Y$  and  $Z$  are binary.

## Instrumental Variables: Nonparametric Selection Model/LATE Conditions

- We started with the mean independence version of the instrumental variables condition, and then discussed strengthening the instrumental variables condition to full independence in the special case where  $Y$  and  $Z$  are binary.
- The result of shifting from mean independence to full independence is to sometimes reduce the width of the resulting bounds but also to have an even more complicated form for the bounds.

- We now consider further strengthening the instrumental variables either by imposing a nonparametric selection model for the first stage as in ? or by imposing instrumental variable conditions of the form considered by ?.

- We now consider further strengthening the instrumental variables either by imposing a nonparametric selection model for the first stage as in ? or by imposing instrumental variable conditions of the form considered by ?.
- The sharp bounds corresponding to these strengthened versions of instrumental variables do not reduce the bounds compared to imposing a weaker form of the instrumental variables assumption but produces a much simpler form for the bounds.

- Let  $D(z)$  denote the counterfactual choices that would have been made had  $Z$  been set exogenously to  $z$ .

- Let  $D(z)$  denote the counterfactual choices that would have been made had  $Z$  been set exogenously to  $z$ .
- Consider the LATE independence, rank, and monotonicity conditions (IV-1), (IV-2), (IV-3) respectively of ? presented in Slides 12 and 152.

- Note that the LATE monotonicity assumption (IV-3) strengthens assumption [IV-BP]. The LATE independence assumption (IV-1) is exactly the same as assumption [IV-BP] except that the assumption is stated here without requiring  $Z$  to be binary.



- Note that the LATE monotonicity assumption (IV-3) strengthens assumption [IV-BP]. The LATE independence assumption (IV-1) is exactly the same as assumption [IV-BP] except that the assumption is stated here without requiring  $Z$  to be binary.
- In their context of binary  $Z$  and  $Y$ , Balke and Pearl discuss the LATE monotonicity condition and show that the LATE monotonicity condition imposes constraints on the observed data which imply that the ? bounds and the Manski mean-independence bounds will coincide.

- Consider the nonparametric selection model of ?:

**Nonparametric Selection Model S:**  $D = \mathbf{1}[\mu(Z) \geq U]$  and  $Z \perp\!\!\!\perp (Y_0, Y_1, U)$ . This is a consequence of equations (7) and assumptions (A-1)–(A-5) presented in Slide 152.

- From ?, we have that the Imbens and Angrist conditions (IV-1)–(IV-3) are equivalent to imposing a nonparametric selection model of the form  $S$ .

- From ?, we have that the Imbens and Angrist conditions (IV-1)–(IV-3) are equivalent to imposing a nonparametric selection model of the form  $S$ .
- Thus, the bounds derived under one set of assumptions will be valid under the alternative set of assumptions, and bounds that are sharp under one set will be sharp under the alternative set of assumptions.

- From ?, we have that the Imbens and Angrist conditions (IV-1)–(IV-3) are equivalent to imposing a nonparametric selection model of the form  $S$ .
- Thus, the bounds derived under one set of assumptions will be valid under the alternative set of assumptions, and bounds that are sharp under one set will be sharp under the alternative set of assumptions.
- This equivalence implies that the Balke and Pearl result also holds for the selection model: if  $Z$  and  $Y$  are binary, then the sharp bounds under the nonparametric selection model coincide with the sharp bounds under mean independence IV.

- We now consider the more general case where neither  $Z$  nor  $Y$  need be binary.

- We now consider the more general case where neither  $Z$  nor  $Y$  need be binary.
- ? derived bounds on the average treatment effect under the assumptions that the outcomes are generated from a bounded outcome nonparametric selection model for treatment without requiring that  $Z$  or  $Y$  be binary or any other restrictions on the support of the distributions of  $Z$  and  $Y$  beyond the assumption that the outcomes are bounded (Assumption [B]).

- In particular, they derived the following bounds on the average treatment effect:

$$B^L \leq E(Y_1 - Y_0) \leq B^U,$$

with

$$B^U = E(DY | P(Z) = p^u) + (1 - p^u)y^u - E((1 - D)Y | P(Z) = p^l) - p^l y^l$$

$$B^L = E(DY | P(Z) = p^u) + (1 - p^u)y^l - E((1 - D)Y | P(Z) = p^l) - p^l y^u.$$



- Note that these bounds do not necessarily include zero.

- Note that these bounds do not necessarily include zero.
- The width of the bounds is

$$B^U - B^L = ((1 - p^u) + p^l)(y^u - y^l).$$

- Note that these bounds do not necessarily include zero.
- The width of the bounds is

$$B^U - B^L = ((1 - p^u) + p^l)(y^u - y^l).$$

- For example, if  $Y$  is binary then the width of the bounds is simply  $B^U - B^L = ((1 - p^u) + p^l)$ .

- Note that these bounds do not necessarily include zero.
- The width of the bounds is

$$B^U - B^L = ((1 - p^u) + p^l)(y^u - y^l).$$

- For example, if  $Y$  is binary then the width of the bounds is simply  $B^U - B^L = ((1 - p^u) + p^l)$ .
- Trivially,  $p^u = 1$  and  $p^l = 0$  is necessary and sufficient for the bounds to collapse to point identification, with the width of the bounds linearly related to the distance between  $p^u$  and 1 and the distance between  $p^l$  and 0.

- Note that it is necessary and sufficient for  $P(z)$  to be a nontrivial function of  $z$  for these bounds to improve upon the bounds that only imposed that the outcomes are bounded.

- Note that it is necessary and sufficient for  $P(z)$  to be a nontrivial function of  $z$  for these bounds to improve upon the bounds that only imposed that the outcomes are bounded.
- Evaluating the width of the bounds only requires  $p^u, p^l$ .

- Note that it is necessary and sufficient for  $P(z)$  to be a nontrivial function of  $z$  for these bounds to improve upon the bounds that only imposed that the outcomes are bounded.
- Evaluating the width of the bounds only requires  $p^u, p^l$ .
- The only additional information required to evaluate the bounds themselves is  $E(DY | P(Z) = p^u)$  and  $E((1 - D)Y | P(Z) = p^l)$ .

- ? analyze how these bounds compare to the ? mean independence bounds, and analyze whether these bounds are sharp.



- ? analyze how these bounds compare to the ? mean independence bounds, and analyze whether these bounds are sharp.
- They show that the selection model imposes restrictions on the observed data such that the ? mean independence bounds collapse to the simpler ? bounds.

- In particular, given assumption S, they show that

$$\inf_{z \in \mathcal{Z}} \{E(DY | Z = z) + (1 - P(z))y^u\} = E(DY | P(Z) = p^u) + (1 - p^u)y^u$$

$$\sup_{z \in \mathcal{Z}} \{E((1 - D)Y | Z = z) + P(z)y^l\} = E((1 - D)Y | P(Z) = p^l) - p^l y^l$$

and thus the ? upper bound collapses to the ? upper bound under assumption S.

- In particular, given assumption S, they show that

$$\inf_{z \in \mathcal{Z}} \{E(DY | Z = z) + (1 - P(z))y^u\} = E(DY | P(Z) = p^u) + (1 - p^u)y^u$$

$$\sup_{z \in \mathcal{Z}} \{E((1 - D)Y | Z = z) + P(z)y^l\} = E((1 - D)Y | P(Z) = p^l) - p^l y^l$$

and thus the ? upper bound collapses to the ? upper bound under assumption S.

- The parallel result holds for the lower bounds.

- In particular, given assumption S, they show that

$$\inf_{z \in \mathcal{Z}} \{E(DY | Z = z) + (1 - P(z))y^u\} = E(DY | P(Z) = p^u) + (1 - p^u)y^u$$

$$\sup_{z \in \mathcal{Z}} \{E((1 - D)Y | Z = z) + P(z)y^l\} = E((1 - D)Y | P(Z) = p^l) - p^l y^l$$

and thus the ? upper bound collapses to the ? upper bound under assumption S.

- The parallel result holds for the lower bounds.
- Furthermore, ? establish that the ? bounds are sharp given assumptions [B] and S.

- Thus, somewhat surprisingly, imposing the stronger assumption of the existence of an instrument in a nonparametric selection model does not narrow the bounds compared to the case of imposing only the weaker assumption of mean independence, but does impose structure on the data which substantially simplifies the form of the the mean-independence bounds.

- Thus, somewhat surprisingly, imposing the stronger assumption of the existence of an instrument in a nonparametric selection model does not narrow the bounds compared to the case of imposing only the weaker assumption of mean independence, but does impose structure on the data which substantially simplifies the form of the the mean-independence bounds.
- By the ? equivalence result, the same conclusion holds for the LATE assumptions – imposing the LATE assumptions does not narrow the bounds compared to only imposing the weaker assumption of mean independence, but does impose restrictions on the data that substantially simplify the form of the bounds.

- Thus, somewhat surprisingly, imposing the stronger assumption of the existence of an instrument in a nonparametric selection model does not narrow the bounds compared to the case of imposing only the weaker assumption of mean independence, but does impose structure on the data which substantially simplifies the form of the the mean-independence bounds.
- By the ? equivalence result, the same conclusion holds for the LATE assumptions – imposing the LATE assumptions does not narrow the bounds compared to only imposing the weaker assumption of mean independence, but does impose restrictions on the data that substantially simplify the form of the bounds.
- ? extend these bounds.

## Combining Comparative Advantage and Instrumental Variables

- We have thus far examined bounds that impose a comparative advantage model, and bounds that exploit an instrumental variables assumption.



## Combining Comparative Advantage and Instrumental Variables

- We have thus far examined bounds that impose a comparative advantage model, and bounds that exploit an instrumental variables assumption.
- In general, neither restriction has more identifying power than the other.

## Combining Comparative Advantage and Instrumental Variables

- We have thus far examined bounds that impose a comparative advantage model, and bounds that exploit an instrumental variables assumption.
- In general, neither restriction has more identifying power than the other.
- We now consider combining both types of assumptions.

- Assume  $D = \mathbf{1}[Y_1 - Y_0 \geq C(Z)]$ , with  $Z$  observed and  $Z \perp\!\!\!\perp (Y_0, Y_1)$ .

- Assume  $D = \mathbf{1}[Y_1 - Y_0 \geq C(Z)]$ , with  $Z$  observed and  $Z \perp\!\!\!\perp (Y_0, Y_1)$ .
- This is a Roy model with a cost  $C(Z)$  of treatment, with the cost of treatment a function of an “instrument”  $Z$ .

- Assume  $D = \mathbf{1}[Y_1 - Y_0 \geq C(Z)]$ , with  $Z$  observed and  $Z \perp\!\!\!\perp (Y_0, Y_1)$ .
- This is a Roy model with a cost  $C(Z)$  of treatment, with the cost of treatment a function of an “instrument”  $Z$ .
- For ease of exposition, assume that  $Z$  is a continuous scalar random variable and that  $(Y_0, Y_1)$  are continuous random variables.

- Assume  $D = \mathbf{1}[Y_1 - Y_0 \geq C(Z)]$ , with  $Z$  observed and  $Z \perp\!\!\!\perp (Y_0, Y_1)$ .
- This is a Roy model with a cost  $C(Z)$  of treatment, with the cost of treatment a function of an “instrument”  $Z$ .
- For ease of exposition, assume that  $Z$  is a continuous scalar random variable and that  $(Y_0, Y_1)$  are continuous random variables.
- Also for ease of exposition, assume that  $\mathcal{Z}$  (the support of the distribution  $Z$ ) is compact and that  $C(\cdot)$  is a continuous function.

- Assume  $D = \mathbf{1}[Y_1 - Y_0 \geq C(Z)]$ , with  $Z$  observed and  $Z \perp\!\!\!\perp (Y_0, Y_1)$ .
- This is a Roy model with a cost  $C(Z)$  of treatment, with the cost of treatment a function of an “instrument”  $Z$ .
- For ease of exposition, assume that  $Z$  is a continuous scalar random variable and that  $(Y_0, Y_1)$  are continuous random variables.
- Also for ease of exposition, assume that  $\mathcal{Z}$  (the support of the distribution  $Z$ ) is compact and that  $C(\cdot)$  is a continuous function.
- These assumptions are only imposed for ease of exposition.

- The model is a special case of the nonparametric selection model considered by ?, but with more structure that we can now exploit.



- The model is a special case of the nonparametric selection model considered by ?, but with more structure that we can now exploit.
- Begin by following steps similar to ?.

- The model is a special case of the nonparametric selection model considered by ?, but with more structure that we can now exploit.
- Begin by following steps similar to ?.
- Using the fact that  $D = \mathbf{1}[Y_1 - Y_0 \geq C(Z)]$  and that  $Z \perp\!\!\!\perp (Y_0, Y_1)$ , we have

$$P(Z) = 1 - F_{Y_1 - Y_0}(C(Z))$$

where  $F_{Y_1 - Y_0}$  is the distribution function of  $Y_1 - Y_0$ .

- Given our assumptions, we have that there will exist  $z^u$  and  $z^l$  such that

$$C(z^u) = \sup\{C(z) : z \in \mathcal{Z}\}, \quad P(z^u) = 1 - F_{Y_1 - Y_0}(C(z^u)) = \inf\{P(Z) : z \in \mathcal{Z}\}$$

$$C(z^l) = \inf\{C(z) : z \in \mathcal{Z}\}, \quad P(z^l) = 1 - F_{Y_1 - Y_0}(C(z^l)) = \sup\{P(Z) : z \in \mathcal{Z}\}$$

- Given our assumptions, we have that there will exist  $z^u$  and  $z^l$  such that

$$C(z^u) = \sup\{C(z) : z \in \mathcal{Z}\}, \quad P(z^u) = 1 - F_{Y_1 - Y_0}(C(z^u)) = \inf\{P(Z) : z \in \mathcal{Z}\}$$

$$C(z^l) = \inf\{C(z) : z \in \mathcal{Z}\}, \quad P(z^l) = 1 - F_{Y_1 - Y_0}(C(z^l)) = \sup\{P(Z) : z \in \mathcal{Z}\}$$

- In other words,  $Z = z^u$  is associated with the highest possible cost of treatment and thus the lowest possible conditional probability of  $D = 1$ , while  $Z = z^l$  is associated with the lowest possible cost of treatment and thus the highest possible conditional probability of  $D = 1$ .

- Given our assumptions, we have that there will exist  $z^u$  and  $z^l$  such that

$$C(z^u) = \sup\{C(z) : z \in \mathcal{Z}\}, \quad P(z^u) = 1 - F_{Y_1 - Y_0}(C(z^u)) = \inf\{P(Z) : z \in \mathcal{Z}\}$$

$$C(z^l) = \inf\{C(z) : z \in \mathcal{Z}\}, \quad P(z^l) = 1 - F_{Y_1 - Y_0}(C(z^l)) = \sup\{P(Z) : z \in \mathcal{Z}\}$$

- In other words,  $Z = z^u$  is associated with the highest possible cost of treatment and thus the lowest possible conditional probability of  $D = 1$ , while  $Z = z^l$  is associated with the lowest possible cost of treatment and thus the highest possible conditional probability of  $D = 1$ .
- Since  $P(\cdot)$  for  $z \in \mathcal{Z}$  is identified, we have that  $z^u$  and  $z^l$  are identified.

- Consider identification of  $C(z)$ .

- Consider identification of  $C(z)$ .
- Using the model and the independence assumptions, we have

$$\begin{aligned}
 \frac{\partial}{\partial z} E(Y|Z = z) &= \frac{\partial}{\partial z} E(YD|Z = z) + \frac{\partial}{\partial z} E(Y(1 - D)|Z = z) \\
 &= \frac{\partial}{\partial z} \int_{C(z)}^{\infty} E(Y_1|Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) \\
 &\quad + \frac{\partial}{\partial z} \int_{-\infty}^{C(z)} E(Y_0|Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) \\
 &= - \left[ E(Y_1|Y_1 - Y_0 = C(z)) - E(Y_0|Y_1 - Y_0 = C(z)) \right] \\
 &\quad \times f_{Y_1 - Y_0}(C(z)) C'(z) \\
 &= - C(z) C'(z) f_{Y_1 - Y_0}(C(z))
 \end{aligned}$$

and

$$\begin{aligned}
 \frac{\partial}{\partial z} P(z) &= \frac{\partial}{\partial z} \int_{C(z)}^{\infty} dF_{Y_1 - Y_0}(t) \\
 &= -C'(z) f_{Y_1 - Y_0}(C(z))
 \end{aligned}$$

- Thus

$$\left[ \frac{\partial}{\partial z} E(Y|Z=z) / \frac{\partial}{\partial z} P(z) \right] = C(z)$$

for any  $z \in \mathcal{Z}$  such that  $\frac{\partial}{\partial z} P(z) \neq 0$ , i.e for any  $z \in \mathcal{Z}$  such that  $C'(z) \neq 0$  and  $F_{Y_1-Y_0}(C(z)) \neq 0$ .



- Thus

$$\left[ \frac{\partial}{\partial z} E(Y|Z=z) / \frac{\partial}{\partial z} P(z) \right] = C(z)$$

for any  $z \in \mathcal{Z}$  such that  $\frac{\partial}{\partial z} P(z) \neq 0$ , i.e for any  $z \in \mathcal{Z}$  such that  $C'(z) \neq 0$  and  $F_{Y_1-Y_0}(C(z)) \neq 0$ .

- We thus conclude that  $C(z)$  is identified for  $z \in \mathcal{Z}$ .

- Our goal is to identify  $E(Y_1 - Y_0)$ .

- Our goal is to identify  $E(Y_1 - Y_0)$ .
- For any  $z \in \mathcal{Z}$ , we have by the law of iterated expectations that

$$\begin{aligned} E(Y_j) &= \int E(Y_j | Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) \\ &= \int_{-\infty}^{C(z)} E(Y_j | Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) + \int_{C(z)}^{\infty} E(Y_j | Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) \end{aligned}$$

for  $j = 0, 1$ .

- Using the model for  $D$  and the assumption that  $Z \perp\!\!\!\perp (Y_0, Y_1)$ , we have

$$\int_{C(z)}^{\infty} E(Y_1 | Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) = E(DY | Z = z) \quad (78)$$

$$\int_{-\infty}^{C(z)} E(Y_0 | Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) = E((1 - D)Y | Z = z). \quad (79)$$

- Using the model for  $D$  and the assumption that  $Z \perp\!\!\!\perp (Y_0, Y_1)$ , we have

$$\int_{C(z)}^{\infty} E(Y_1 | Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) = E(DY | Z = z) \quad (78)$$

$$\int_{-\infty}^{C(z)} E(Y_0 | Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) = E((1 - D)Y | Z = z). \quad (79)$$

- We identify the right hand sides of these equations for any  $z \in \mathcal{Z}$ , and thus identify the left hand sides for any  $z \in \mathcal{Z}$ .

- Using the model for  $D$  and the assumption that  $Z \perp\!\!\!\perp (Y_0, Y_1)$ , we have

$$\int_{C(z)}^{\infty} E(Y_1 | Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) = E(DY | Z = z) \quad (78)$$

$$\int_{-\infty}^{C(z)} E(Y_0 | Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) = E((1 - D)Y | Z = z). \quad (79)$$

- We identify the right hand sides of these equations for any  $z \in \mathcal{Z}$ , and thus identify the left hand sides for any  $z \in \mathcal{Z}$ .
- In particular, consider evaluating equation (78) at  $z = z^l$  and equation (79) at  $z = z^u$ .

- Using the model for  $D$  and the assumption that  $Z \perp\!\!\!\perp (Y_0, Y_1)$ , we have

$$\int_{C(z)}^{\infty} E(Y_1 | Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) = E(DY | Z = z) \quad (78)$$

$$\int_{-\infty}^{C(z)} E(Y_0 | Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) = E((1 - D)Y | Z = z). \quad (79)$$

- We identify the right hand sides of these equations for any  $z \in \mathcal{Z}$ , and thus identify the left hand sides for any  $z \in \mathcal{Z}$ .
- In particular, consider evaluating equation (78) at  $z = z^l$  and equation (79) at  $z = z^u$ .
- Then, to bound  $E(Y_1 - Y_0)$ , we need to bound

$$\int_{-\infty}^{C(z^l)} E(Y_1 | Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) \text{ and}$$

$$\int_{C(z^u)}^{\infty} E(Y_0 | Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t).$$

- We have

$$\begin{aligned}
 & \int_{-\infty}^{C(z')} E(Y_1 | Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) \\
 &= (1 - P(z')) E[Y_1 | Z = z', Y_1 \leq Y_0 + C(z')] \\
 &\leq (1 - P(z')) E[Y_0 + C(z') | Z = z', Y_1 \leq Y_0 + C(z')] \\
 &= E[(1 - D)Y | Z = z'] + (1 - P(z'))C(z') \\
 &= E[(1 - D)Y | Z = z'] - \left[ \frac{\partial}{\partial z} E(Y | Z = z) / \frac{\partial}{\partial z} \ln(1 - P(z)) \right] \Big|_{z=z'},
 \end{aligned}$$

where the inequality arises from the conditioning  $Y_1 \leq Y_0 + C(z')$ .



- We have

$$\begin{aligned}
 & \int_{-\infty}^{C(z')} E(Y_1 | Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) \\
 &= (1 - P(z')) E[Y_1 | Z = z', Y_1 \leq Y_0 + C(z')] \\
 &\leq (1 - P(z')) E[Y_0 + C(z') | Z = z', Y_1 \leq Y_0 + C(z')] \\
 &= E[(1 - D)Y | Z = z'] + (1 - P(z'))C(z') \\
 &= E[(1 - D)Y | Z = z'] - \left[ \frac{\partial}{\partial z} E(Y | Z = z) / \frac{\partial}{\partial z} \ln(1 - P(z)) \right] \Big|_{z=z'}
 \end{aligned}$$

where the inequality arises from the conditioning  $Y_1 \leq Y_0 + C(z')$ .

- The final expression follows from our derivation of  $C(z)$ .

- Since  $\Pr[y^l \leq Y_1 \leq y^u] = 1$  by assumption, we have

$$\begin{aligned}
 (1 - P(z^l))y^l &\leq \int_{-\infty}^{c(z^l)} E(Y_1 | Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) \\
 &\leq E[(1 - D)Y | Z = z^l] - \left[ \frac{\partial}{\partial z} E(Y | Z = z) / \frac{\partial}{\partial z} \ln(1 - P(z)) \right] \Big|_{z=z^l}.
 \end{aligned}$$

- Since  $\Pr[y^l \leq Y_1 \leq y^u] = 1$  by assumption, we have

$$\begin{aligned} (1 - P(z^l))y^l &\leq \int_{-\infty}^{C(z^l)} E(Y_1 | Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) \\ &\leq E[(1 - D)Y | Z = z^l] - \left[ \frac{\partial}{\partial z} E(Y | Z = z) / \frac{\partial}{\partial z} \ln(1 - P(z)) \right] \Bigg|_{z=z^l}. \end{aligned}$$

- By a parallel argument, we have

$$\begin{aligned} P(z^u)y^l &\leq \int_{C(z^u)}^{\infty} E(Y_0 | Y_1 - Y_0 = t) dF_{Y_1 - Y_0}(t) \\ &\leq E[DY | Z = z^u] + \left[ \frac{\partial}{\partial z} E(Y | Z = z) / \frac{\partial}{\partial z} \ln P(z) \right] \Bigg|_{z=z^u}. \end{aligned}$$

- We thus have the bounds:

$$B^L \leq E(Y_1 - Y_0) \leq B^U,$$

with

$$\begin{aligned} B^U = & E(Y | Z = z') \\ & - \left[ \frac{\partial}{\partial z} E(Y | Z = z) \bigg/ \frac{\partial}{\partial z} \ln(1 - P(z)) \right] \bigg|_{z=z'} \\ & - E((1 - D)Y | Z = z'') - P(z'')y' \end{aligned}$$

$$\begin{aligned} B^L = & E(DY | Z = z') + [1 - P(z')]y' - E(Y | Z = z'') \\ & - \left[ \frac{\partial}{\partial z} E(Y | Z = z) \bigg/ \frac{\partial}{\partial z} \ln P(z) \right] \bigg|_{z=z''}. \end{aligned}$$

- The last two terms in  $B^U$  come from the lower bound for  $E(Y_0)$  and the first two terms come from the upper bound for  $E(Y_1)$  just derived.

- The last two terms in  $B^U$  come from the lower bound for  $E(Y_0)$  and the first two terms come from the upper bound for  $E(Y_1)$  just derived.
- The terms for  $B^L$  are decomposed in an analogous fashion, reversing the roles of the upper and lower bounds for  $E(Y_1)$  and  $E(Y_0)$ .

- The last two terms in  $B^U$  come from the lower bound for  $E(Y_0)$  and the first two terms come from the upper bound for  $E(Y_1)$  just derived.
- The terms for  $B^L$  are decomposed in an analogous fashion, reversing the roles of the upper and lower bounds for  $E(Y_1)$  and  $E(Y_0)$ .
- These bounds improve over the bounds that only impose a nonparametric selection model (Assumption S) without imposing the Roy model structure.

- The last two terms in  $B^U$  come from the lower bound for  $E(Y_0)$  and the first two terms come from the upper bound for  $E(Y_1)$  just derived.
- The terms for  $B^L$  are decomposed in an analogous fashion, reversing the roles of the upper and lower bounds for  $E(Y_1)$  and  $E(Y_0)$ .
- These bounds improve over the bounds that only impose a nonparametric selection model (Assumption S) without imposing the Roy model structure.
- We next consider some alternative approaches to the solution of selection and hence evaluation problems developed in the literature using replacement functions, proxy functions, and other conditions.



## Control Functions, Replacement Functions, and Proxy Variables

- This chapter analyzes the main tools used to evaluate social programs in the presence of selection bias in observational data.

## Control Functions, Replacement Functions, and Proxy Variables

- This chapter analyzes the main tools used to evaluate social programs in the presence of selection bias in observational data.
- Yet many other tools have not been analyzed.

## Control Functions, Replacement Functions, and Proxy Variables

- This chapter analyzes the main tools used to evaluate social programs in the presence of selection bias in observational data.
- Yet many other tools have not been analyzed.
- We briefly summarize these approaches.

## Control Functions, Replacement Functions, and Proxy Variables

- This chapter analyzes the main tools used to evaluate social programs in the presence of selection bias in observational data.
- Yet many other tools have not been analyzed.
- We briefly summarize these approaches.
- ? establishes conditions under which some of the methods we discuss produce identification of econometric models.

## Control Functions, Replacement Functions, and Proxy Variables

- This chapter analyzes the main tools used to evaluate social programs in the presence of selection bias in observational data.
- Yet many other tools have not been analyzed.
- We briefly summarize these approaches.
- ? establishes conditions under which some of the methods we discuss produce identification of econometric models.
- We use some of these tools in Part III.

The methods of replacement functions and proxy variables all start from characterizations (U-1) and (U-2) which we repeat for convenience:

(U-1)

$$(Y_0, Y_1) \perp\!\!\!\perp D \mid X, Z, \theta,$$

but

(U-2)

$(Y_0, Y_1) \not\perp D \mid X, Z.$

where  $\theta$  is not observed by the analyst and  $(Y_0, Y_1)$  are not observed directly but  $Y$  is observed as are the  $X, Z$ :

$$Y = DY_1 + (1 - D) Y_0.$$

- Missing variables  $\theta$  produce selection bias which creates a problem with using observational data to evaluate social programs.



- Missing variables  $\theta$  produce selection bias which creates a problem with using observational data to evaluate social programs.
- From (U-1), if we condition on  $\theta$ , we would satisfy the condition (M-1) for matching, and hence could identify the parameters and distributions that can be identified if the conditions required for matching are satisfied.

- The most direct approach to controlling for  $\theta$  is to assume access to a function  $\tau(X, Z, Q)$  that perfectly proxies  $\theta$ :

$$\theta = \tau(X, Z, Q). \quad (80)$$

- The most direct approach to controlling for  $\theta$  is to assume access to a function  $\tau(X, Z, Q)$  that perfectly proxies  $\theta$ :

$$\theta = \tau(X, Z, Q). \quad (80)$$

- This approach based on a perfect proxy is called the **method of replacement functions** by ?.

- The most direct approach to controlling for  $\theta$  is to assume access to a function  $\tau(X, Z, Q)$  that perfectly proxies  $\theta$ :

$$\theta = \tau(X, Z, Q). \quad (80)$$

- This approach based on a perfect proxy is called the **method of replacement functions** by ?.
- In (U-1), we can substitute for  $\theta$  in terms of observables  $(X, Z, Q)$  . Then

$$(Y_0, Y_1) \perp\!\!\!\perp D \mid X, Z, Q.$$

- The most direct approach to controlling for  $\theta$  is to assume access to a function  $\tau(X, Z, Q)$  that perfectly proxies  $\theta$ :

$$\theta = \tau(X, Z, Q). \quad (80)$$

- This approach based on a perfect proxy is called the **method of replacement functions** by ?.
- In (U-1), we can substitute for  $\theta$  in terms of observables  $(X, Z, Q)$  . Then

$$(Y_0, Y_1) \perp\!\!\!\perp D \mid X, Z, Q.$$

- We can condition nonparametrically on  $(X, Z, Q)$  and do not have to know the exact functional form of  $\tau$  although knowledge of  $\tau$  might reduce the dimensionality of the matching problem.

- $\theta$  can be a vector and  $\tau$  can be a vector of functions.

- $\theta$  can be a vector and  $\tau$  can be a vector of functions.
- This method has been used in the economics of education for decades (see the references in ?).

- $\theta$  can be a vector and  $\tau$  can be a vector of functions.
- This method has been used in the economics of education for decades (see the references in ?).
- If  $\theta$  is ability and  $\tau$  is a test score, it is sometimes assumed that the test score is a perfect proxy (or replacement function) for  $\theta$  and  $\tau$  is entered into the regressions of earnings on schooling to escape the problem of ability bias, typically assuming a linear relationship between  $\tau$  and  $\theta$ .



- ? discuss the literature that uses replacement functions in this way.

- ? discuss the literature that uses replacement functions in this way.
- ? apply this method and consider nonparametric identification of the  $\tau$  function.

- ? discuss the literature that uses replacement functions in this way.
- ? apply this method and consider nonparametric identification of the  $\tau$  function.
- ? provides a rigorous proof of identification for this approach in a general nonparametric setting.

- The method of replacement functions assumes that (80) is a perfect proxy.

- The method of replacement functions assumes that (80) is a perfect proxy.
- In many applications, this assumption is far too strong.

- The method of replacement functions assumes that (80) is a perfect proxy.
- In many applications, this assumption is far too strong.
- More often, we measure  $\theta$  with error.

- The method of replacement functions assumes that (80) is a perfect proxy.
- In many applications, this assumption is far too strong.
- More often, we measure  $\theta$  with error.
- This produces a factor model or measurement error model (?).

- ? surveys this method.



- ? surveys this method.
- We can represent the factor model in a general way by a system of equations:

$$Y_j = g_j(X, Z, Q, \theta, \varepsilon_j), \quad j = 1, \dots, J. \quad (81)$$

- ? surveys this method.
- We can represent the factor model in a general way by a system of equations:

$$Y_j = g_j(X, Z, Q, \theta, \varepsilon_j), \quad j = 1, \dots, J. \quad (81)$$

- A linear factor model separable in the unobservables writes

$$Y_j = g_j(X, Z, Q) + \lambda_j \theta + \varepsilon_j, \quad j = 1, \dots, J, \quad (82)$$

where

$$(X, Z, Q) \perp\!\!\!\perp (\theta, \varepsilon_j), \varepsilon_j \perp\!\!\!\perp \theta, \quad j = 1, \dots, J, \quad (83)$$

and the  $\varepsilon_j$  are mutually independent.

- Observe that under (81) and (82),  $Y_j$  controlling for  $X, Z, Q$  only imperfectly proxies  $\theta$  because of the presence of  $\varepsilon_j$ .

- Observe that under (81) and (82),  $Y_j$  controlling for  $X, Z, Q$  only imperfectly proxies  $\theta$  because of the presence of  $\varepsilon_j$ .
- The  $\theta$  are called factors,  $\lambda_j$  factor loadings and the  $\varepsilon_j$  “uniquenesses” (see, e.g., ?).

- A large literature, partially reviewed in Part III, section 1, and in ?, shows how to establish identification of econometric models under factor structure assumptions.

- A large literature, partially reviewed in Part III, section 1, and in ?, shows how to establish identification of econometric models under factor structure assumptions.
- ?, ? and ? establish identification in nonlinear models of the form (81).

- A large literature, partially reviewed in Part III, section 1, and in ?, shows how to establish identification of econometric models under factor structure assumptions.
- ?, ? and ? establish identification in nonlinear models of the form (81).
- The key to identification is multiple, but imperfect (because of  $\varepsilon_j$ ), measurements on  $\theta$  from the  $Y_j$ ,  $j = 1, \dots, J$  and  $X, Z, Q$ , and possibly other measurement systems that depend on  $\theta$ .

- $\tau$ ,  $\tau^*$  and  $\tau^{\#}$  apply and develop these methods.



- ?, ?? and ?? apply and develop these methods.
- Under assumption (83), they show how to nonparametrically identify the econometric model and the distributions of the unobservables  $F_\theta(\theta)$  and  $F_{\varepsilon_j}(\varepsilon_j)$ .

- ?, ?? and ?? apply and develop these methods.
- Under assumption (83), they show how to nonparametrically identify the econometric model and the distributions of the unobservables  $F_\theta(\theta)$  and  $F_{\varepsilon_j}(\varepsilon_j)$ .
- In the context of classical simultaneous equations models, identification is secured by using covariance restrictions across equations exploiting the low dimensionality of vector  $\theta$  compared to the high dimensional vector of (imperfect) measurements on it.

- $?$ ,  $??$  and  $??$  apply and develop these methods.
- Under assumption (83), they show how to nonparametrically identify the econometric model and the distributions of the unobservables  $F_\theta(\theta)$  and  $F_{\varepsilon_j}(\varepsilon_j)$ .
- In the context of classical simultaneous equations models, identification is secured by using covariance restrictions across equations exploiting the low dimensionality of vector  $\theta$  compared to the high dimensional vector of (imperfect) measurements on it.
- The recent literature ( $???$ ) extends the linear model to a nonlinear setting.

- The recent econometric literature applies in special cases the idea of the **control function principle** introduced in ?.

- The recent econometric literature applies in special cases the idea of the **control function principle** introduced in ?.
- This principle, versions of which can be traced back to ?, partitions  $\theta$  in (U-1) into two or more components,  $\theta = (\theta_1, \theta_2)$ , where only one component of  $\theta$  is the source of bias.

Thus it is assumed that (U-1) is true, and (U-1)' is also true:

(U-1)'

$$(Y_0, Y_1) \perp\!\!\!\perp D \mid X, Z, \theta_1,$$

and (U-2) holds.

- For example, in the normal selection model analyzed in Part I, in Slide 819, we broke  $U_1$ , the error term associated with  $Y_1$ , into two components:

$$U_1 = E(U_1 | V) + \varepsilon,$$

where  $V$  plays the role of  $\theta_1$  and arises from the choice equation.

- For example, in the normal selection model analyzed in Part I, in Slide 819, we broke  $U_1$ , the error term associated with  $Y_1$ , into two components:

$$U_1 = E(U_1 | V) + \varepsilon,$$

where  $V$  plays the role of  $\theta_1$  and arises from the choice equation.

- Under normality,  $\varepsilon$  is independent of  $E(U_1 | V)$ .



- Further,

$$E(U_1 | V) = \frac{\text{Cov}(U_1, V)}{\text{Var}(V)} V, \quad (84)$$

assuming  $E(U_1) = 0$  and  $E(V) = 0$ .

- Further,

$$E(U_1 | V) = \frac{\text{Cov}(U_1, V)}{\text{Var}(V)} V, \quad (84)$$

assuming  $E(U_1) = 0$  and  $E(V) = 0$ .

- In that section, we show how to construct a control function in the context of the choice model

$$D = \mathbf{1} [\mu_D(Z) \geq V].$$

- Further,

$$E(U_1 | V) = \frac{\text{Cov}(U_1, V)}{\text{Var}(V)} V, \quad (84)$$

assuming  $E(U_1) = 0$  and  $E(V) = 0$ .

- In that section, we show how to construct a control function in the context of the choice model

$$D = \mathbf{1} [\mu_D(Z) \geq V].$$

- Controlling for  $V$  controls for the component of  $\theta_1$  in (U-1)' that gives rise to the spurious dependence.

- The ?? application of the control function principle assumes functional form (84 ) but assumes that  $V$  can be perfectly proxied by a first stage equation.

- The ?? application of the control function principle assumes functional form (84 ) but assumes that  $V$  can be perfectly proxied by a first stage equation.
- Thus they use a replacement function in their first stage.

- The ?? application of the control function principle assumes functional form (84 ) but assumes that  $V$  can be perfectly proxied by a first stage equation.
- Thus they use a replacement function in their first stage.
- Their method does not work when one can only condition on  $D$  rather than on  $D^* = \mu_D(Z) - V$ .

- The ?? application of the control function principle assumes functional form (84 ) but assumes that  $V$  can be perfectly proxied by a first stage equation.
- Thus they use a replacement function in their first stage.
- Their method does not work when one can only condition on  $D$  rather than on  $D^* = \mu_D(Z) - V$ .
- In the sample selection model, it is not necessary to use  $V$ .

- As developed in Part I and reviewed in Slide 338 and Slide 727 of this chapter, under additive separability for the outcome equation for  $Y_1$ , we can write

$$E(Y_1 | X, Z, D = 1) = \mu_1(X) + \underbrace{E(U_1 | \mu_D(Z) \geq V)}_{\text{control function}}$$

so we “expect out” rather than solve out the effect of the component of  $V$  on  $U_1$  and thus control for selection bias under our maintained assumptions.



- In terms of the propensity score, under the conditions specified in Part I, we may write the preceding expression in terms of  $P(Z)$ :

$$E(Y_1 | X, Z, D = 1) = \mu_1(X) + K_1(P(Z)),$$

where  $K_1(P(Z)) = E(U_1 | X, Z, D = 1)$ .

- In terms of the propensity score, under the conditions specified in Part I, we may write the preceding expression in terms of  $P(Z)$ :

$$E(Y_1 | X, Z, D = 1) = \mu_1(X) + K_1(P(Z)),$$

where  $K_1(P(Z)) = E(U_1 | X, Z, D = 1)$ .

- It is not necessary to know  $V$  or be able to estimate it.

- In terms of the propensity score, under the conditions specified in Part I, we may write the preceding expression in terms of  $P(Z)$ :

$$E(Y_1 | X, Z, D = 1) = \mu_1(X) + K_1(P(Z)),$$

where  $K_1(P(Z)) = E(U_1 | X, Z, D = 1)$ .

- It is not necessary to know  $V$  or be able to estimate it.
- The ?? application of the control function principle assumes that the analyst can condition on and estimate  $V$ .

- The Blundell-Powell method and the method of ? build heavily on (84) and implicitly make strong distributional and functional form assumptions that are not intrinsic to the method of control functions.

- The Blundell-Powell method and the method of ? build heavily on (84) and implicitly make strong distributional and functional form assumptions that are not intrinsic to the method of control functions.
- As just noted, their method uses a replacement function to obtain  $E(U_1 | V)$  in the first step of their procedures.

- The Blundell-Powell method and the method of ? build heavily on (84) and implicitly make strong distributional and functional form assumptions that are not intrinsic to the method of control functions.
- As just noted, their method uses a replacement function to obtain  $E(U_1 | V)$  in the first step of their procedures.
- The general control function method does not require a replacement function approach.

- The Blundell-Powell method and the method of ? build heavily on (84) and implicitly make strong distributional and functional form assumptions that are not intrinsic to the method of control functions.
- As just noted, their method uses a replacement function to obtain  $E(U_1 | V)$  in the first step of their procedures.
- The general control function method does not require a replacement function approach.
- The literature has begun to distinguish between the more general control function approach and the *control variate* approach that uses a first stage replacement function.

- ? develops the method of unobservable instruments which is a version of the replacement function approach applied to nonlinear models.



- ? develops the method of unobservable instruments which is a version of the replacement function approach applied to nonlinear models.
- Her unobservable instruments play the role of covariance restrictions used to identify classical simultaneous equations models (see ?).

- ? develops the method of unobservable instruments which is a version of the replacement function approach applied to nonlinear models.
- Her unobservable instruments play the role of covariance restrictions used to identify classical simultaneous equations models (see ?).
- Her approach is distinct from and therefore complementary with linear factor models.

- ? develops the method of unobservable instruments which is a version of the replacement function approach applied to nonlinear models.
- Her unobservable instruments play the role of covariance restrictions used to identify classical simultaneous equations models (see ?).
- Her approach is distinct from and therefore complementary with linear factor models.
- Instead of assuming  $(X, Z, Q) \perp\!\!\!\perp (\theta, \varepsilon_j)$ , she assumes in a two equation system that  $(\theta, \varepsilon_1) \perp\!\!\!\perp Y_2 \mid Y_1, X, Z$ .

- ? develops the method of unobservable instruments which is a version of the replacement function approach applied to nonlinear models.
- Her unobservable instruments play the role of covariance restrictions used to identify classical simultaneous equations models (see ?).
- Her approach is distinct from and therefore complementary with linear factor models.
- Instead of assuming  $(X, Z, Q) \perp\!\!\!\perp (\theta, \varepsilon_j)$ , she assumes in a two equation system that  $(\theta, \varepsilon_1) \perp\!\!\!\perp Y_2 \mid Y_1, X, Z$ .
- See the discussion in ?.

- We have not discussed panel data methods in this chapter.

- We have not discussed panel data methods in this chapter.
- The most commonly used panel data method is difference-in-differences as discussed in ?, ?, ?, and ?, to cite only a few key papers.

- We have not discussed panel data methods in this chapter.
- The most commonly used panel data method is difference-in-differences as discussed in ?, ?, ?, and ?, to cite only a few key papers.
- Most of the estimators we have discussed can be adapted to a panel data setting.

- ? develop difference-in-differences matching estimators.



- ? develop difference-in-differences matching estimators.
- ? extends this work.

- ? develop difference-in-differences matching estimators.
- ? extends this work.
- Separability between errors and observables is a key feature of the panel data approach in its standard application.

- ? develop difference-in-differences matching estimators.
- ? extends this work.
- Separability between errors and observables is a key feature of the panel data approach in its standard application.
- ? and ? present analyses of nonseparable panel data methods.

## Summary

- This chapter summarizes the main methods used to identify mean treatment effect parameters under semiparametric and nonparametric assumptions.

## Summary

- This chapter summarizes the main methods used to identify mean treatment effect parameters under semiparametric and nonparametric assumptions.
- We have used the marginal treatment effect as the unifying parameter to straddle a diverse econometric literature summarized in table 1 of this chapter.

## Summary

- This chapter summarizes the main methods used to identify mean treatment effect parameters under semiparametric and nonparametric assumptions.
- We have used the marginal treatment effect as the unifying parameter to straddle a diverse econometric literature summarized in table 1 of this chapter.
- For each estimator, we establish what it identifies, the economic content of the estimand and the identifying assumptions of the method.

# Appendices

## Relationships Among Parameters Using the Index Structure

- Given the index structure, a simple relationship exists among the parameters.



## Relationships Among Parameters Using the Index Structure

- Given the index structure, a simple relationship exists among the parameters.
- It is immediate from the definitions  $D = \mathbf{1}(U_D \leq P(z))$  and  $\Delta = Y_1 - Y_0$  that

$$\Delta^{\text{TT}}(x, P(z)) = E(\Delta | X = x, U_D \leq P(z)). \quad (85)$$

## Relationships Among Parameters Using the Index Structure

- Given the index structure, a simple relationship exists among the parameters.
- It is immediate from the definitions  $D = \mathbf{1}(U_D \leq P(z))$  and  $\Delta = Y_1 - Y_0$  that

$$\Delta^{\text{TT}}(x, P(z)) = E(\Delta | X = x, U_D \leq P(z)). \quad (85)$$

- Next consider  $\Delta^{\text{LATE}}(x, P(z), P(z'))$ .

• Note that  $E(Y|X = x, P(Z) = P(z))$

$$\begin{aligned}
 &= P(z) \left[ E(Y_1|X = x, P(Z) = P(z), D = 1) \right] \\
 &\quad + (1 - P(z)) \left[ E(Y_0|X = x, P(Z) = P(z), D = 0) \right] \\
 &= \int_0^{P(z)} E(Y_1|X = x, U_D = u_D) du_D + \int_{P(z)}^1 E(Y_0|X = x, U_D = u_D) du_D,
 \end{aligned}$$

so that

$$\begin{aligned}
 &E(Y|X = x, P(Z) = P(z)) - E(Y|X = x, P(Z) = P(z')) \\
 &= \int_{P(z')}^{P(z)} E(Y_1|X = x, U_D = u_D) du_D - \int_{P(z')}^{P(z)} E(Y_0|X = x, U_D = u_D) du_D,
 \end{aligned}$$

and thus

$$\Delta^{\text{LATE}}(x, P(z), P(z')) = E(\Delta|X = x, P(z') \leq U_D \leq P(z)).$$

- Notice that this expression could be taken as an alternative definition of LATE.

- Notice that this expression could be taken as an alternative definition of LATE.
- Note that, in this expression, we could replace  $P(z)$  and  $P(z')$  with  $u_D$  and  $u'_D$ .

- Notice that this expression could be taken as an alternative definition of LATE.
- Note that, in this expression, we could replace  $P(z)$  and  $P(z')$  with  $u_D$  and  $u'_D$ .
- No instrument needs to be available to define LATE.

- We can rewrite these relationships in succinct form in the following way:

$$\begin{aligned}
 \Delta^{\text{MTE}}(x, u_D) &= E(\Delta | X = x, U_D = u_D) \\
 \Delta^{\text{ATE}}(x) &= \int_0^1 E(\Delta | X = x, U_D = u_D) du_D \\
 P(z)[\Delta^{\text{TT}}(x, P(z))] &= \int_0^{P(z)} E(\Delta | X = x, U_D = u_D) du_D \\
 (P(z) - P(z'))[\Delta^{\text{LATE}}(x, P(z), P(z'))] &= \int_{P(z')}^{P(z)} E(\Delta | X = x, U_D = u_D) du_D.
 \end{aligned} \tag{86}$$

- We can rewrite these relationships in succinct form in the following way:

$$\begin{aligned}
 \Delta^{\text{MTE}}(x, u_D) &= E(\Delta | X = x, U_D = u_D) \\
 \Delta^{\text{ATE}}(x) &= \int_0^1 E(\Delta | X = x, U_D = u_D) du_D \\
 P(z)[\Delta^{\text{TT}}(x, P(z))] &= \int_0^{P(z)} E(\Delta | X = x, U_D = u_D) du_D \\
 (P(z) - P(z'))[\Delta^{\text{LATE}}(x, P(z), P(z'))] &= \int_{P(z')}^{P(z)} E(\Delta | X = x, U_D = u_D) du_D.
 \end{aligned} \tag{86}$$

- We stress that everywhere in these expressions we can replace  $P(z)$  with  $u_D$  and  $P(z')$  with  $u'_D$ .



- We can rewrite these relationships in succinct form in the following way:

$$\begin{aligned}
 \Delta^{\text{MTE}}(x, u_D) &= E(\Delta | X = x, U_D = u_D) \\
 \Delta^{\text{ATE}}(x) &= \int_0^1 E(\Delta | X = x, U_D = u_D) du_D \\
 P(z)[\Delta^{\text{TT}}(x, P(z))] &= \int_0^{P(z)} E(\Delta | X = x, U_D = u_D) du_D \\
 (P(z) - P(z'))[\Delta^{\text{LATE}}(x, P(z), P(z'))] &= \int_{P(z')}^{P(z)} E(\Delta | X = x, U_D = u_D) du_D.
 \end{aligned} \tag{86}$$

- We stress that everywhere in these expressions we can replace  $P(z)$  with  $u_D$  and  $P(z')$  with  $u'_D$ .
- Each parameter is an average value of MTE,  $E(\Delta | X = x, U_D = u_D)$ , but for values of  $U_D$  lying in different intervals and with different weighting functions.

- MTE defines the treatment effect more finely than do LATE, ATE, or TT.

- MTE defines the treatment effect more finely than do LATE, ATE, or TT.
- The relationship between MTE and LATE or TT conditional on  $P(z)$  is analogous to the relationship between a probability density function and a cumulative distribution function.

- MTE defines the treatment effect more finely than do LATE, ATE, or TT.
- The relationship between MTE and LATE or TT conditional on  $P(z)$  is analogous to the relationship between a probability density function and a cumulative distribution function.
- The probability density function and the cumulative distribution function represent the same information, but for some purposes the density function is more easily interpreted.

- Likewise, knowledge of TT for all  $P(z)$  evaluation points is equivalent to knowledge of the MTE for all  $u_D$  evaluation points, so it is not the case that knowledge of one provides more information than knowledge of the other.

- Likewise, knowledge of TT for all  $P(z)$  evaluation points is equivalent to knowledge of the MTE for all  $u_D$  evaluation points, so it is not the case that knowledge of one provides more information than knowledge of the other.
- However, in many choice-theoretic contexts it is often easier to interpret MTE than the TT or LATE parameters.

- Likewise, knowledge of TT for all  $P(z)$  evaluation points is equivalent to knowledge of the MTE for all  $u_D$  evaluation points, so it is not the case that knowledge of one provides more information than knowledge of the other.
- However, in many choice-theoretic contexts it is often easier to interpret MTE than the TT or LATE parameters.
- It has the interpretation as a measure of willingness to pay on the part of people on a specified margin of participation in the program.

- $\Delta^{\text{MTE}}(x, u_D)$  is the average effect for people who are just indifferent between participation in the program ( $D = 1$ ) or not ( $D = 0$ ) if the instrument is externally set so that  $P(Z) = u_D$ .



- $\Delta^{\text{MTE}}(x, u_D)$  is the average effect for people who are just indifferent between participation in the program ( $D = 1$ ) or not ( $D = 0$ ) if the instrument is externally set so that  $P(Z) = u_D$ .
- For values of  $u_D$  close to zero,  $\Delta^{\text{MTE}}(x, u_D)$  is the average effect for individuals with unobservable characteristics that make them the most inclined to participate in the program ( $D = 1$ ), and for values of  $u_D$  close to one it is the average treatment effect for individuals with unobserved (by the econometrician) characteristics that make them the least inclined to participate.

- $\Delta^{\text{MTE}}(x, u_D)$  is the average effect for people who are just indifferent between participation in the program ( $D = 1$ ) or not ( $D = 0$ ) if the instrument is externally set so that  $P(Z) = u_D$ .
- For values of  $u_D$  close to zero,  $\Delta^{\text{MTE}}(x, u_D)$  is the average effect for individuals with unobservable characteristics that make them the most inclined to participate in the program ( $D = 1$ ), and for values of  $u_D$  close to one it is the average treatment effect for individuals with unobserved (by the econometrician) characteristics that make them the least inclined to participate.
- ATE integrates  $\Delta^{\text{MTE}}(x, u_D)$  over the entire support of  $U_D$  (from  $u_D = 0$  to  $u_D = 1$ ).

- It is the average effect for an individual chosen at random from the entire population.

- It is the average effect for an individual chosen at random from the entire population.
- $\Delta^{TT}(x, P(z))$  is the average treatment effect for persons who chose to participate at the given value of  $P(Z) = P(z)$ ; it integrates  $\Delta^{MTE}(x, u_D)$  up to  $u_D = P(z)$ .

- It is the average effect for an individual chosen at random from the entire population.
- $\Delta^{TT}(x, P(z))$  is the average treatment effect for persons who chose to participate at the given value of  $P(Z) = P(z)$ ; it integrates  $\Delta^{MTE}(x, u_D)$  up to  $u_D = P(z)$ .
- As a result, it is primarily determined by the MTE parameter for individuals whose unobserved characteristics make them the most inclined to participate in the program.

- LATE is the average treatment effect for someone who would not participate if  $P(Z) \leq P(z)$  and would participate if  $P(Z) \geq P(z)$ .

- LATE is the average treatment effect for someone who would not participate if  $P(Z) \leq P(z')$  and would participate if  $P(Z) \geq P(z)$ .
- The parameter  $\Delta^{\text{LATE}}(x, P(z), P(z'))$  integrates  $\Delta^{\text{MTE}}(x, u_D)$  from  $u_D = P(z')$  to  $u_D = P(z)$ .

- Using the third expression in equation (86) to substitute into equation (85), we obtain an alternative expression for the TT parameter as a weighted average of MTE parameters:

$$\Delta^{TT}(x) = \int_0^1 \frac{1}{p} \left[ \int_0^p E(\Delta | X = x, U_D = u_D) du_D \right] dF_{P(Z)|X,D}(p|x, D = 1).$$



- Using the third expression in equation (86) to substitute into equation (85), we obtain an alternative expression for the TT parameter as a weighted average of MTE parameters:

$$\Delta^{TT}(x) = \int_0^1 \frac{1}{p} \left[ \int_0^p E(\Delta|X=x, U_D = u_D) du_D \right] dF_{P(Z)|X,D}(p|x, D=1).$$

- Using Bayes' rule, it follows that

$$dF_{P(Z)|X,D}(p|x, 1) = \frac{\Pr(D=1|X=x, P(Z)=p)}{\Pr(D=1|X=x)} dF_{P(Z)|X}(p|x).$$

- Since  $\Pr(D = 1|X = x, P(Z) = p) = p$ , it follows that

$$\Delta^{TT}(x) = \frac{1}{\Pr(D = 1|X = x)} \int_0^1 \left( \int_0^p E(\Delta|X = x, U_D = u_D) du_D \right) dF_{P(Z)|X}(p|x). \quad (87)$$

- Note further that since

$\Pr(D = 1|X = x) = E(P(Z)|X = x) = \int_0^1 (1 - F_{P(Z)|X}(t|x)) dt$ ,

we can reinterpret (87) as a weighted average of local IV parameters where the weighting is similar to that obtained from a length-biased, size-biased, or  $P$ -biased sample:

$$\begin{aligned} \Delta^{\text{IT}}(x) &= \frac{1}{\Pr(D = 1 | X = x)} \\ &\quad \times \int_0^1 \left( \int_0^1 \mathbf{1}(u_D \leq p) E(\Delta | X = x, U_D = u_D) du_D \right) dF_{P(Z)|X}(p|x) \\ &= \frac{1}{\int (1 - F_{P(Z)|X}(t|x)) dt} \\ &\quad \times \int_0^1 \left( \int_0^1 E(\Delta | X = x, U_D = u_D) \mathbf{1}(u_D \leq p) dF_{P(Z)|X}(p|x) \right) du_D \\ &= \int_0^1 E(\Delta | X = x, U_D = u_D) \left( \frac{1 - F_{P(Z)|X}(u_D|x)}{\int (1 - F_{P(Z)|X}(t|x)) dt} \right) du_D \\ &= \int_0^1 E(\Delta | X = x, U_D = u_D) g_x(u_D) du_D, \end{aligned}$$

where  $g_x(u_D) = \frac{1 - F_{P(Z)|X}(u_D|x)}{\int (1 - F_{P(Z)|X}(t|x)) dt}$ .

- Thus  $g_x(u_D)$  is a *weighted distribution* (?).

- Thus  $g_x(u_D)$  is a *weighted distribution* (?).
- Since  $g_x(u_D)$  is a nonincreasing function of  $u_D$ , we have that drawings from  $g_x(u_D)$  oversample persons with low values of  $U_D$ , i.e., values of unobserved characteristics that make them the most likely to participate in the program no matter what their value of  $P(Z)$ .

- Since

$$\Delta^{\text{MTE}}(x, u_D) = E(\Delta | X = x, U_D = u_D)$$

it follows that

$$\Delta^{\text{TT}}(x) = \int_0^1 \Delta^{\text{MTE}}(x, u_D) g_x(u_D) du_D.$$

The TT parameter is thus a weighted version of MTE, where  $\Delta^{\text{MTE}}(x, u_D)$  is given the largest weight for low  $u_D$  values and is given zero weight for  $u_D \geq p_x^{\text{max}}$ , where  $p_x^{\text{max}}$  is the maximum value in the support of  $P(Z)$  conditional on  $X = x$ .

- Figure A-1 graphs the relationship between  $\Delta^{\text{MTE}}(u_D)$ ,  $\Delta^{\text{ATE}}$  and  $\Delta^{\text{TT}}(P(z))$ , assuming that the gains are the greatest for those with the lowest  $U_D$  values and that the gains decline as  $U_D$  increases.

- Figure A-1 graphs the relationship between  $\Delta^{\text{MTE}}(u_D)$ ,  $\Delta^{\text{ATE}}$  and  $\Delta^{\text{TT}}(P(z))$ , assuming that the gains are the greatest for those with the lowest  $U_D$  values and that the gains decline as  $U_D$  increases.
- The curve is the MTE parameter as a function of  $u_D$ , and is drawn for the special case where the outcome variable is binary so that MTE parameter is bounded between  $-1$  and  $1$ .

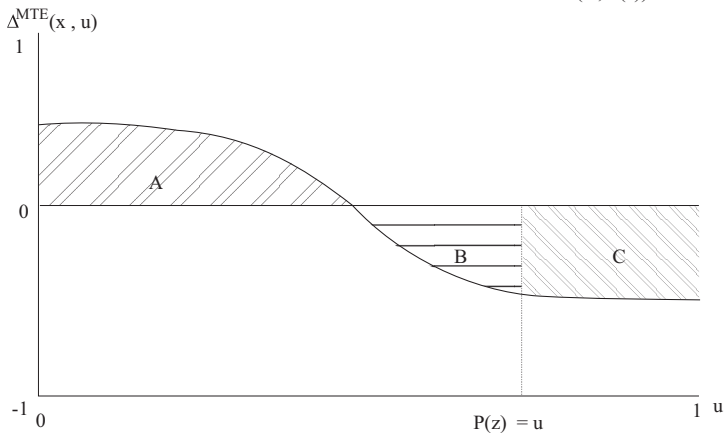


- Figure A-1 graphs the relationship between  $\Delta^{\text{MTE}}(u_D)$ ,  $\Delta^{\text{ATE}}$  and  $\Delta^{\text{TT}}(P(z))$ , assuming that the gains are the greatest for those with the lowest  $U_D$  values and that the gains decline as  $U_D$  increases.
- The curve is the MTE parameter as a function of  $u_D$ , and is drawn for the special case where the outcome variable is binary so that MTE parameter is bounded between  $-1$  and  $1$ .
- The ATE parameter averages  $\Delta^{\text{MTE}}(u_D)$  over the full unit interval (i.e., is the area under A minus the area under B and C in the figure).

Figure A-1: MTE Integrates to ATE and TT Under Full Support (for dichotomous outcome)

$$\Delta^{\text{ATE}}(x) = A - (B + C)$$

$$\Delta^{\text{TT}}(x, P(z)) = A - B$$



Source: Heckman and Vytlačil (2000).

- $\Delta^{TT}(P(z))$  averages  $\Delta^{MTE}(u_D)$  up to the point  $P(z)$  (is the area under A minus the area under B in the figure).

- $\Delta^{TT}(P(z))$  averages  $\Delta^{MTE}(u_D)$  up to the point  $P(z)$  (is the area under A minus the area under B in the figure).
- Because  $\Delta^{MTE}(u_D)$  is assumed to be declining in  $u_D$ , the TT parameter for any given  $P(z)$  evaluation point is larger than the ATE parameter.

- Equation (86) relates each of the other parameters to the MTE parameter.

- Equation (86) relates each of the other parameters to the MTE parameter.
- One can also relate each of the other parameters to the LATE parameter.

- Equation (86) relates each of the other parameters to the MTE parameter.
- One can also relate each of the other parameters to the LATE parameter.
- This relationship turns out to be useful later on in this chapter when we encounter conditions where LATE can be identified but MTE cannot.

- Equation (86) relates each of the other parameters to the MTE parameter.
- One can also relate each of the other parameters to the LATE parameter.
- This relationship turns out to be useful later on in this chapter when we encounter conditions where LATE can be identified but MTE cannot.
- MTE is the limit form of LATE:

$$\Delta^{\text{MTE}}(x, p) = \lim_{p' \rightarrow p} \Delta^{\text{LATE}}(x, p, p').$$



- Equation (86) relates each of the other parameters to the MTE parameter.
- One can also relate each of the other parameters to the LATE parameter.
- This relationship turns out to be useful later on in this chapter when we encounter conditions where LATE can be identified but MTE cannot.
- MTE is the limit form of LATE:

$$\Delta^{\text{MTE}}(x, p) = \lim_{p' \rightarrow p} \Delta^{\text{LATE}}(x, p, p').$$

- Direct relationships between LATE and the other parameters are easily derived.

- The relationship between LATE and ATE is immediate:

$$\Delta^{\text{ATE}}(x) = \Delta^{\text{LATE}}(x, 0, 1).$$

- The relationship between LATE and ATE is immediate:

$$\Delta^{\text{ATE}}(x) = \Delta^{\text{LATE}}(x, 0, 1).$$

- Using Bayes' rule, the relationship between LATE and TT is

$$\Delta^{\text{TT}}(x) = \int_0^1 \Delta^{\text{LATE}}(x, 0, p) \frac{p}{\Pr(D = 1|X = x)} dF_{P(Z)|X}(p|x). \quad (88)$$

## Relaxing Additive Separability and Independence

- There are two central assumptions that underlie the latent index representation used in this chapter: that  $V$  is independent of  $Z$ , and that  $V$  and  $Z$  are additively separable in the index.

## Relaxing Additive Separability and Independence

- There are two central assumptions that underlie the latent index representation used in this chapter: that  $V$  is independent of  $Z$ , and that  $V$  and  $Z$  are additively separable in the index.
- The latent index model with these two restrictions implies the independence and monotonicity assumptions of ? and the latent index model implied by those assumptions implies a latent index model with a representation that satisfies both the independence and the monotonicity assumptions.

## Relaxing Additive Separability and Independence

- There are two central assumptions that underlie the latent index representation used in this chapter: that  $V$  is independent of  $Z$ , and that  $V$  and  $Z$  are additively separable in the index.
- The latent index model with these two restrictions implies the independence and monotonicity assumptions of ? and the latent index model implied by those assumptions implies a latent index model with a representation that satisfies both the independence and the monotonicity assumptions.
- In this appendix, we consider the sensitivity of the analysis presented in the text to relaxation of either of these assumptions.

- First, consider allowing  $V$  and  $Z$  to be nonseparable in the treatment index:

$$D^* = \mu_D(Z, V)$$
$$D = \begin{cases} 1 & \text{if } D^* \geq 0 \\ 0 & \text{otherwise} \end{cases} ,$$

while maintaining the assumption that  $Z$  is independent of  $(V, U_1, U_0)$ . We do not impose any restrictions on the cross partials of  $\mu_D$ .

- First, consider allowing  $V$  and  $Z$  to be nonseparable in the treatment index:

$$D^* = \mu_D(Z, V)$$

$$D = \begin{cases} 1 & \text{if } D^* \geq 0 \\ 0 & \text{otherwise} \end{cases} ,$$

while maintaining the assumption that  $Z$  is independent of  $(V, U_1, U_0)$ . We do not impose any restrictions on the cross partials of  $\mu_D$ .

- The monotonicity condition of ? is that for any  $(z, z')$  pair,  $\mu_D(z, v) \geq \mu_D(z', v)$  for all  $v$ , or  $\mu_D(z, v) \leq \mu_D(z', v)$  for all  $v$ .



- First, consider allowing  $V$  and  $Z$  to be nonseparable in the treatment index:

$$D^* = \mu_D(Z, V)$$

$$D = \begin{cases} 1 & \text{if } D^* \geq 0 \\ 0 & \text{otherwise} \end{cases},$$

while maintaining the assumption that  $Z$  is independent of  $(V, U_1, U_0)$ . We do not impose any restrictions on the cross partials of  $\mu_D$ .

- The monotonicity condition of ? is that for any  $(z, z')$  pair,  $\mu_D(z, v) \geq \mu_D(z', v)$  for all  $v$ , or  $\mu_D(z, v) \leq \mu_D(z', v)$  for all  $v$ .
- ? shows that monotonicity always implies one representation of  $\mu_D$  as  $\mu_D(z, v) = \mu_D(z) + v$ .

- We now reconsider the analysis in the text without imposing the monotonicity condition by considering the latent index model without additive separability.

- We now reconsider the analysis in the text without imposing the monotonicity condition by considering the latent index model without additive separability.
- Since we have imposed no structure on the  $\mu_D(z, v)$  index, one can easily show that this model is equivalent to imposing the independence condition of ? without imposing their monotonicity condition.

- We now reconsider the analysis in the text without imposing the monotonicity condition by considering the latent index model without additive separability.
- Since we have imposed no structure on the  $\mu_D(z, v)$  index, one can easily show that this model is equivalent to imposing the independence condition of ? without imposing their monotonicity condition.
- A random coefficient discrete choice model with  $\mu_D = Z\gamma + \varepsilon$  where  $\gamma$  and  $\varepsilon$  are random, and  $\gamma$  can assume positive or negative values is an example of this case, i.e.,  $V = (\gamma, \varepsilon)$ .

- We impose the regularity condition that, for any  $z \in \text{Supp}(Z)$ ,  $\mu_D(z, V)$  is absolutely continuous with respect to Lebesgue measure.

- We impose the regularity condition that, for any  $z \in \text{Supp}(Z)$ ,  $\mu_D(z, V)$  is absolutely continuous with respect to Lebesgue measure.
- Let

$$\Omega(z) = \{v : \mu_D(z, v) \geq 0\},$$

so that

$$P(z) \equiv \Pr(D = 1 | Z = z) = \Pr(V \in \Omega(z)).$$

- We impose the regularity condition that, for any  $z \in \text{Supp}(Z)$ ,  $\mu_D(z, V)$  is absolutely continuous with respect to Lebesgue measure.
- Let

$$\Omega(z) = \{v : \mu_D(z, v) \geq 0\},$$

so that

$$P(z) \equiv \Pr(D = 1 | Z = z) = \Pr(V \in \Omega(z)).$$

- Under additive separability,  $P(z) = P(z') \Leftrightarrow \Omega(z) = \Omega(z')$ .

- This equivalence enables us to define the parameters in terms of the  $P(z)$  index instead of the full  $z$  vector.



- This equivalence enables us to define the parameters in terms of the  $P(z)$  index instead of the full  $z$  vector.
- In the more general case without additive separability, it is possible to have  $(z, z')$  such that  $P(z) = P(z')$  and  $\Omega(z) \neq \Omega(z')$ .

- This equivalence enables us to define the parameters in terms of the  $P(z)$  index instead of the full  $z$  vector.
- In the more general case without additive separability, it is possible to have  $(z, z')$  such that  $P(z) = P(z')$  and  $\Omega(z) \neq \Omega(z')$ .
- We present a random coefficient choice model example of this case in Slide 381 in the text.

- This equivalence enables us to define the parameters in terms of the  $P(z)$  index instead of the full  $z$  vector.
- In the more general case without additive separability, it is possible to have  $(z, z')$  such that  $P(z) = P(z')$  and  $\Omega(z) \neq \Omega(z')$ .
- We present a random coefficient choice model example of this case in Slide 381 in the text.
- In this case, we can no longer replace  $Z = z$  with  $P(Z) = P(z)$  in the conditioning sets.

- Define, using  $\Delta = Y_1 - Y_0$ ,

$$\Delta^{\text{MTE}}(x, v) = E(\Delta | X = x, V = v).$$

- Define, using  $\Delta = Y_1 - Y_0$ ,

$$\Delta^{\text{MTE}}(x, v) = E(\Delta | X = x, V = v).$$

- For ATE, we obtain the same expression as before:

$$\Delta^{\text{ATE}}(x) = \int_{-\infty}^{\infty} E(\Delta | X = x, V = v) dF_{V|X}(v).$$

- Define, using  $\Delta = Y_1 - Y_0$ ,

$$\Delta^{\text{MTE}}(x, v) = E(\Delta | X = x, V = v).$$

- For ATE, we obtain the same expression as before:

$$\Delta^{\text{ATE}}(x) = \int_{-\infty}^{\infty} E(\Delta | X = x, V = v) dF_{V|X}(v).$$

- For TT, we obtain a similar but slightly more complicated expression:

$$\begin{aligned} \Delta^{\text{TT}}(x, z) &\equiv E(\Delta | X = x, Z = z, D = 1) \\ &= E(\Delta | X = x, V \in \Omega(z)) \\ &= \frac{1}{P(z)} \int_{\Omega(z)} E(\Delta | X = x, V = v) dF_{V|X}(v). \end{aligned}$$

- Because it is no longer the case that we can define the parameter solely in terms of  $P(z)$  instead of  $z$ , it is possible to have  $(z, z')$  such that  $P(z) = P(z')$  but  $\Delta^{TT}(x, z) \neq \Delta^{TT}(x, z')$ .

- Following the same derivation as used in the text for the TT parameter not conditional on  $Z$ ,

$$\begin{aligned}
 \Delta^{\text{TT}}(x) &\equiv E(\Delta|X=x, D=1) \\
 &= \int E(\Delta|X=x, Z=z, D=1) dF_{Z|X, D}(z|x, 1) \\
 &= \frac{1}{\Pr(D=1|X=x)} \\
 &\quad \times \int \left[ \int_{-\infty}^{\infty} \mathbf{1}[v \in \Omega(z)] E(\Delta|X=x, V=v) dF_{V|X}(v) \right] dF_{Z|X}(z|x) \\
 &= \frac{1}{\Pr(D=1|X=x)} \\
 &\quad \times \int_{-\infty}^{\infty} \left[ \int \mathbf{1}[v \in \Omega(z)] E(\Delta|X=x, V=v) dF_{Z|X}(z|x) \right] dF_{V|X}(v) \\
 &= \int_{-\infty}^{\infty} E(\Delta|X=x, V=v) g_x(v) dv
 \end{aligned}$$

where

$$g_x(v) = \frac{\int \mathbf{1}[v \in \Omega(z)] dF_{Z|X}(z|x)}{\Pr(D=1|X=x)} = \frac{\Pr(D=1|V=v, X=x)}{\Pr(D=1|X=x)}.$$



- Thus the definitions of the parameters and the relationships among them that are developed in the main text of this chapter generalize in a straightforward way to the nonseparable case.

- Thus the definitions of the parameters and the relationships among them that are developed in the main text of this chapter generalize in a straightforward way to the nonseparable case.
- Separability allows us to define the parameters in terms of  $P(z)$  instead of  $z$  and allows for slightly simpler expressions, but is not crucial for the definition of parameters or the relationship among them.

- Separability is, however, crucial to the form of LATE when we allow  $V$  and  $Z$  to be additively nonseparable in the treatment index.

- Separability is, however, crucial to the form of LATE when we allow  $V$  and  $Z$  to be additively nonseparable in the treatment index.
- For simplicity, we will keep the conditioning on  $X$  implicit.

- Separability is, however, crucial to the form of LATE when we allow  $V$  and  $Z$  to be additively nonseparable in the treatment index.
- For simplicity, we will keep the conditioning on  $X$  implicit.
- Define the following sets

$$A(z, z') = \{v : \mu_D(z, v) \geq 0, \mu_D(z', v) \geq 0\}$$

$$B(z, z') = \{v : \mu_D(z, v) \geq 0, \mu_D(z', v) < 0\}$$

$$C(z, z') = \{v : \mu_D(z, v) < 0, \mu_D(z', v) < 0\}$$

$$D(z, z') = \{v : \mu_D(z, v) < 0, \mu_D(z', v) \geq 0\}.$$

- Separability is, however, crucial to the form of LATE when we allow  $V$  and  $Z$  to be additively nonseparable in the treatment index.
- For simplicity, we will keep the conditioning on  $X$  implicit.
- Define the following sets

$$A(z, z') = \{v : \mu_D(z, v) \geq 0, \mu_D(z', v) \geq 0\}$$

$$B(z, z') = \{v : \mu_D(z, v) \geq 0, \mu_D(z', v) < 0\}$$

$$C(z, z') = \{v : \mu_D(z, v) < 0, \mu_D(z', v) < 0\}$$

$$D(z, z') = \{v : \mu_D(z, v) < 0, \mu_D(z', v) \geq 0\}.$$

- Monotonicity implies that either  $B(z, z')$  or  $D(z, z')$  is empty.

- Suppressing the  $z, z'$  arguments, we have:

$$\begin{aligned} E(Y|Z = z) &= \Pr(A \cup B)E(Y_1|A \cup B) + \Pr(C \cup D)E(Y_0|C \cup D) \\ E(Y|Z = z') &= \Pr(A \cup D)E(Y_1|A \cup D) + \Pr(B \cup C)E(Y_0|B \cup C) \end{aligned}$$

so that

$$\begin{aligned} \frac{E(Y|Z = z) - E(Y|Z = z')}{\Pr(D = 1|Z = z) - \Pr(D = 1|Z = z')} &= \frac{E(Y|Z = z) - E(Y|Z = z')}{\Pr(A \cup B) - \Pr(A \cup D)} \\ &= \frac{\Pr(B)E(Y_1 - Y_0|B) - \Pr(D)E(Y_1 - Y_0|D)}{\Pr(B) - \Pr(D)} \\ &= w_B E(\Delta|B) - w_D E(\Delta|D) \end{aligned}$$

with

$$\begin{aligned} w_B &= \frac{\Pr(B|B \cup D)}{\Pr(B|B \cup D) - \Pr(D|B \cup D)} \\ w_D &= \frac{\Pr(D|B \cup D)}{\Pr(B|B \cup D) - \Pr(D|B \cup D)}. \end{aligned}$$

- Under monotonicity, either  $\Pr(B) = 0$  and LATE identifies  $E(\Delta|D)$  or  $\Pr(D) = 0$  and LATE identifies  $E(\Delta|B)$ .



- Under monotonicity, either  $\Pr(B) = 0$  and LATE identifies  $E(\Delta|D)$  or  $\Pr(D) = 0$  and LATE identifies  $E(\Delta|B)$ .
- Without monotonicity, the IV estimator used as the sample analogue to LATE converges to the above weighted difference in the two terms, and the relationship between LATE and the other treatment parameters presented in the text no longer holds.

- Consider what would happen if we could condition on a given  $v$ .

- Consider what would happen if we could condition on a given  $v$ .
- For  $v \in A \cup C$ , the denominator is zero and the parameter is not well defined.

- Consider what would happen if we could condition on a given  $v$ .
- For  $v \in A \cup C$ , the denominator is zero and the parameter is not well defined.
- For  $v \in B$ , the parameter is  $E(\Delta|V = v)$ , for  $v \in D$ , the parameter is  $E(\Delta|V = v)$ .

- Consider what would happen if we could condition on a given  $v$ .
- For  $v \in A \cup C$ , the denominator is zero and the parameter is not well defined.
- For  $v \in B$ , the parameter is  $E(\Delta|V = v)$ , for  $v \in D$ , the parameter is  $E(\Delta|V = v)$ .
- If we could restrict conditioning to  $v \in B$  (or  $v \in D$ ), we would obtain monotonicity within the restricted sample.

- Now consider LIV.

- Now consider LIV.
- For simplicity, assume  $z$  is a scalar.

- Now consider LIV.
- For simplicity, assume  $z$  is a scalar.
- Assume  $\mu_D(z, v)$  is continuously differentiable in  $(z, v)$ , with  $\mu^j(z, v)$  denoting the partial derivative with respect to the  $j$ th argument.



- Now consider LIV.
- For simplicity, assume  $z$  is a scalar.
- Assume  $\mu_D(z, v)$  is continuously differentiable in  $(z, v)$ , with  $\mu^j(z, v)$  denoting the partial derivative with respect to the  $j$ th argument.
- Assume that  $\mu_D(Z, V)$  is absolutely continuous with respect to Lebesgue measure.

- Now consider LIV.
- For simplicity, assume  $z$  is a scalar.
- Assume  $\mu_D(z, v)$  is continuously differentiable in  $(z, v)$ , with  $\mu^j(z, v)$  denoting the partial derivative with respect to the  $j$ th argument.
- Assume that  $\mu_D(Z, V)$  is absolutely continuous with respect to Lebesgue measure.
- Fix some evaluation point,  $z_0$ .

- Now consider LIV.
- For simplicity, assume  $z$  is a scalar.
- Assume  $\mu_D(z, v)$  is continuously differentiable in  $(z, v)$ , with  $\mu^j(z, v)$  denoting the partial derivative with respect to the  $j$ th argument.
- Assume that  $\mu_D(Z, V)$  is absolutely continuous with respect to Lebesgue measure.
- Fix some evaluation point,  $z_0$ .
- One can show that there may be at most a countable number of  $v$  points such that  $\mu_D(z_0, v) = 0$ .

- Let  $j \in \mathcal{J} = \{1, \dots, L\}$  index the set of  $v$  evaluation points such that  $\mu_D(z_0, v) = 0$ , where  $L$  may be infinity, and thus write:  $\mu_D(z_0, v_j) = 0$  for all  $j \in \mathcal{J}$ .

- Let  $j \in \mathcal{J} = \{1, \dots, L\}$  index the set of  $v$  evaluation points such that  $\mu_D(z_0, v) = 0$ , where  $L$  may be infinity, and thus write:  $\mu_D(z_0, v_j) = 0$  for all  $j \in \mathcal{J}$ .
- (Both the number of such evaluation points and the evaluation points themselves depend on the evaluation point,  $z_0$ , but we suppress this dependence for notational convenience.) Assume that there exists  $\{B_k\}_{k \in \mathcal{J}}$ ,  $\sum_{k \in \mathcal{J}} B_k < \infty$  such that
 
$$\left| \frac{\mu^1(z, v_k)}{\mu^2(z, v_k)} \right| \leq B_k \text{ for } k \in \mathcal{J} \text{ and all } z \text{ in some neighborhood of } z_0.$$

- One can show that

$$\frac{\partial}{\partial z} \left[ E(Y|Z = z) \right] \Big|_{z=z_0} = \sum_{k=1}^L \frac{\mu^1(z_0, v_k)}{|\mu^2(z_0, v_k)|} E(\Delta | V = v_k)$$

and

$$\frac{\partial}{\partial z} [Pr(D = 1|Z = z)] \Big|_{z=z_0} = \sum_{k=1}^L \frac{\mu^1(z_0, v_k)}{|\mu^2(z_0, v_k)|} \Big|_{z=z_0} .$$

- One can show that

$$\frac{\partial}{\partial z} \left[ E(Y|Z = z) \right] \Big|_{z=z_0} = \sum_{k=1}^L \frac{\mu^1(z_0, v_k)}{|\mu^2(z_0, v_k)|} E(\Delta | V = v_k)$$

and

$$\frac{\partial}{\partial z} [Pr(D = 1|Z = z)] \Big|_{z=z_0} = \sum_{k=1}^L \frac{\mu^1(z_0, v_k)}{|\mu^2(z_0, v_k)|} \Big|_{z=z_0} .$$

- LIV is the ratio of these two terms, and does not in general equal the MTE.

- One can show that

$$\frac{\partial}{\partial z} \left[ E(Y|Z = z) \right] \Big|_{z=z_0} = \sum_{k=1}^L \frac{\mu^1(z_0, v_k)}{|\mu^2(z_0, v_k)|} E(\Delta | V = v_k)$$

and

$$\frac{\partial}{\partial z} [Pr(D = 1|Z = z)] \Big|_{z=z_0} = \sum_{k=1}^L \frac{\mu^1(z_0, v_k)}{|\mu^2(z_0, v_k)|} \Big|_{z=z_0} .$$

- LIV is the ratio of these two terms, and does not in general equal the MTE.
- Thus, the relationship between LIV and MTE breaks down in the nonseparable case.



- As an example, take the case where  $L$  is finite and  $\frac{\mu^1(z_0, v_k)}{|\mu^2(z_0, v_k)|}$  does not vary with  $k$ .

- As an example, take the case where  $L$  is finite and  $\frac{\mu^1(z_0, v_k)}{|\mu^2(z_0, v_k)|}$  does not vary with  $k$ .
- For this case,

$$\begin{aligned} \Delta^{\text{LIV}}(z_0) &= \Pr(\mu^1(z_0, V) > 0 | \mu(z_0, V) = 0) \cdot E\left(\Delta \mid \mu_D(z_0, V) = 0, \mu^1(z_0, V) > 0\right) \\ &\quad - \Pr(\mu^1(z_0, V) < 0 | \mu(z_0, V) = 0) E\left(\Delta \mid \mu_D(z_0, V) = 0, \mu^1(z_0, V) < 0\right). \end{aligned}$$

Thus, while the definition of the parameters and the relationship among them does not depend crucially on the additive separability assumption, the connection between the LATE or LIV estimators and the underlying parameters crucially depends on the additive separability assumption.

- Next consider the assumption that  $V$  and  $Z$  are separable in the treatment index while allowing them to be stochastically dependent:

$$D^* = \mu_D(Z) - V$$
$$D = \begin{cases} 1 & \text{if } D^* \geq 0 \\ 0 & \text{otherwise} \end{cases},$$

with  $Z$  independent of  $(U_0, U_1)$ , but allowing  $Z$  and  $V$  to be stochastically dependent.

- Next consider the assumption that  $V$  and  $Z$  are separable in the treatment index while allowing them to be stochastically dependent:

$$D^* = \mu_D(Z) - V$$

$$D = \begin{cases} 1 & \text{if } D^* \geq 0 \\ 0 & \text{otherwise} \end{cases},$$

with  $Z$  independent of  $(U_0, U_1)$ , but allowing  $Z$  and  $V$  to be stochastically dependent.

- The analysis of ? can be easily adapted to show that the latent index model with separability but without imposing independence is equivalent to imposing the monotonicity assumption of Imbens and Angrist without imposing their independence assumption.

- We have

$$\Omega(z) = \{v : \mu_D(z) \geq v\}$$

and

$$P(z) \equiv \Pr(D = 1|Z = z) = \Pr(V \in \Omega(z)|Z = z).$$

- We have

$$\Omega(z) = \{v : \mu_D(z) \geq v\}$$

and

$$P(z) \equiv \Pr(D = 1|Z = z) = \Pr(V \in \Omega(z)|Z = z).$$

- Note that  $\Omega(z) = \Omega(z') \Rightarrow \mu_D(z) = \mu_D(z')$ , but  $\Omega(z) = \Omega(z')$  does not imply  $P(z) = P(z')$  since the distribution of  $V$  conditional on  $Z = z$  need not equal the distribution of  $V$  conditional on  $Z = z'$ .

- We have

$$\Omega(z) = \{v : \mu_D(z) \geq v\}$$

and

$$P(z) \equiv \Pr(D = 1|Z = z) = \Pr(V \in \Omega(z)|Z = z).$$

- Note that  $\Omega(z) = \Omega(z') \Rightarrow \mu_D(z) = \mu_D(z')$ , but  $\Omega(z) = \Omega(z')$  does not imply  $P(z) = P(z')$  since the distribution of  $V$  conditional on  $Z = z$  need not equal the distribution of  $V$  conditional on  $Z = z'$ .
- Likewise,  $P(z) = P(z')$  does not imply  $\Omega(z) = \Omega(z')$ .

- We have

$$\Omega(z) = \{v : \mu_D(z) \geq v\}$$

and

$$P(z) \equiv \Pr(D = 1|Z = z) = \Pr(V \in \Omega(z)|Z = z).$$

- Note that  $\Omega(z) = \Omega(z') \Rightarrow \mu_D(z) = \mu_D(z')$ , but  $\Omega(z) = \Omega(z')$  does not imply  $P(z) = P(z')$  since the distribution of  $V$  conditional on  $Z = z$  need not equal the distribution of  $V$  conditional on  $Z = z'$ .
- Likewise,  $P(z) = P(z')$  does not imply  $\Omega(z) = \Omega(z')$ .
- As occurred in the nonseparable case, we can no longer replace  $Z = z$  with  $P(Z) = P(z)$  in the conditioning sets.



- Consider the definition of the parameters and the relationship among them.

- Consider the definition of the parameters and the relationship among them.
- The definition of MTE and ATE in no way involves  $Z$ , nor does the relationship between them, so that both their definition and their relationship remains unchanged by allowing  $Z$  and  $V$  to be dependent.

- Now consider the TT parameter where now we make the dependence of  $X$  explicit:

$$\begin{aligned}
 \Delta^{\text{TT}}(x, z) &= E(\Delta | X = x, Z = z, V \leq \mu_D(z)) \\
 &= \frac{1}{P(z)} \int_{-\infty}^{\mu_D(z)} E(\Delta | X = x, V = v) dF_{V|Z, X}(v | z, x) \\
 &= \frac{1}{P(z)} \int_{-\infty}^{\mu_D(z)} E(\Delta | X = x, V = v) \frac{f_{Z|V, X}(z|v, x)}{f_{Z|X}(z|x)} dF_{V|X}(v | x)
 \end{aligned}$$

where  $f_{Z|X}$  and  $f_{Z|V, X}$  denote the densities corresponding to  $F_{Z|X}$  and  $F_{Z|V, X}$  with respect to the appropriate dominating measure.

- We thus obtain

$$\begin{aligned}
 & \Delta^{\text{TT}}(x) \\
 &= E(\Delta | X = x, V \leq \mu_D(Z)) \\
 &= \frac{1}{\Pr(D = 1 | X = x)} \\
 & \quad \times \int \left[ \int_{-\infty}^{\mu_D(z)} E(\Delta | X = x, V = v) \left[ \frac{f_{Z|U,X}(z|v,x)}{f_{Z|X}(z|x)} dF_{V|X}(v|x) \right] dF_{Z|X}(z|x) \right] \\
 &= \frac{1}{\Pr(D = 1 | X = x)} \\
 & \quad \times \int_{-\infty}^{\infty} \left[ \int \mathbf{1}[v \leq \mu_D(z)] E(\Delta | X = x, V = v) \left[ \frac{f_{Z|U,X}(z|v,x)}{f_{Z|X}(z|x)} dF_{Z|X}(z|x) \right] dF_{V|X}(v|x) \right] \\
 &= \frac{1}{\Pr(D = 1 | X = x)} \\
 & \quad \times \int_{-\infty}^{\infty} \left[ \int \mathbf{1}[v \leq \mu_D(z)] E(\Delta | X = x, V = v) dF_{Z|V,X}(z|v,x) \right] dF_{V|X}(v|x) \\
 &= \int_{-\infty}^{\infty} E(\Delta | X = x, V = v) g_x(v) dv
 \end{aligned}$$

where

$$g_x(v) = \frac{\Pr(D = 1 | V = v, X = x)}{\Pr(D = 1 | X = x)}.$$

- Thus the definitions of parameters and the relationships among the parameters that are developed in the text generalize naturally to the case where  $Z$  and  $V$  are stochastically dependent.

- Thus the definitions of parameters and the relationships among the parameters that are developed in the text generalize naturally to the case where  $Z$  and  $V$  are stochastically dependent.
- Independence (combined with the additive separability assumption) allows us to define the parameters in terms of  $P(z)$  instead of  $z$  and allows for slightly simpler expressions, but is not crucial for the definition of parameters or the relationship among them.

- We next investigate LATE when we allow  $V$  and  $Z$  to be stochastically dependent.

- We next investigate LATE when we allow  $V$  and  $Z$  to be stochastically dependent.
- We have

$$\begin{aligned}
 E(Y|X=x, Z=z) &= P(z) \left[ E(Y_1|X=x, Z=z, D=1) \right] \\
 &\quad + (1 - P(z)) \left[ E(Y_0|X=x, Z=z, D=0) \right] \\
 &= \int_{-\infty}^{\mu_D(z)} E(Y_1|X=x, V=v) dF_{V|X,Z}(v|x, z) \\
 &\quad + \int_{\mu_D(z)}^{\infty} E(Y_0|X=x, V=v) dF_{V|X,Z}(v|x, z),
 \end{aligned}$$



- For simplicity, take the case where  $\mu_D(z) > \mu_D(z')$ .

- For simplicity, take the case where  $\mu_D(z) > \mu_D(z')$ .
- Then

$$\begin{aligned}
 & E(Y|X = x, Z = z) - E(Y|X = x, Z = z') \\
 = & \left[ \int_{\mu_D(z')}^{\mu_D(z)} E(Y_1|X = x, V = v) dF_{V|X,Z}(v|x, z) \right. \\
 & \left. - \int_{\mu_D(z')}^{\mu_D(z)} E(Y_0|X = x, V = v) dF_{V|X,Z}(v|x, z') \right] \\
 & + \int_{-\infty}^{\mu_D(z')} E(Y_1|X = x, V = v) \left( dF_{V|X,Z}(v|x, z) - dF_{V|X,Z}(v|x, z') \right) \\
 & + \int_{\mu_D(z)}^{\infty} E(Y_0|X = x, V = v) \left( dF_{V|X,Z}(v|x, z) - dF_{V|X,Z}(v|x, z') \right)
 \end{aligned}$$

## • Thus

$$\begin{aligned} & \Delta^{\text{LATE}}(x, z, z') \\ &= \delta_0(z)E(Y_1|X = x, Z = z, \mu_D(z') \leq V \leq \mu_D(z)) \\ & \quad - \delta_0(z')E(Y_0|X = x, Z = z', \mu_D(z') \leq V \leq \mu_D(z)) \\ & \quad + \left[ \delta_1(z)E(Y_1|X = x, Z = z, V \leq \mu_D(z')) - \delta_1(z')E(Y_1|X = x, Z = z', V \leq \mu_D(z')) \right] \\ & \quad + \left[ \delta_2(z)E(Y_0|X = x, Z = z, V > \mu_D(z)) - \delta_2(z')E(Y_0|X = x, Z = z', V > \mu_D(z)) \right], \end{aligned}$$

with

$$\delta_0(t) = \frac{\Pr(\mu_D(z') \leq V \leq \mu_D(z) | Z = t)}{\Pr(V \leq \mu_D(z) | Z = z, X = x) - \Pr(V \leq \mu_D(z') | Z = z', X = x)}$$

$$\delta_1(t) = \frac{\Pr(V \leq \mu_D(z') | Z = t)}{\Pr(V \leq \mu_D(z) | Z = z, X = x) - \Pr(V \leq \mu_D(z') | Z = z', X = x)}$$

$$\delta_2(t) = \frac{\Pr(V > \mu_D(z) | Z = t)}{\Pr(V \leq \mu_D(z) | Z = z, X = x) - \Pr(V \leq \mu_D(z') | Z = z', X = x)}.$$

- Note that  $\delta_0(z) = \delta_0(z') = 1$  and the two terms in brackets are zero in the case where  $Z$  and  $V$  are independent.

- Note that  $\delta_0(z) = \delta_0(z') = 1$  and the two terms in brackets are zero in the case where  $Z$  and  $V$  are independent.
- In the more general case,  $\delta_0$  may be bigger or smaller than 1, and the terms in brackets are of unknown sign.

- Note that  $\delta_0(z) = \delta_0(z') = 1$  and the two terms in brackets are zero in the case where  $Z$  and  $V$  are independent.
- In the more general case,  $\delta_0$  may be bigger or smaller than 1, and the terms in brackets are of unknown sign.
- In general, LATE may be negative even when  $\Delta$  is positive for all individuals.

- Now consider LIV.

- Now consider LIV.
- For simplicity, take the case where  $Z$  is a continuous scalar r.v.



- Now consider LIV.
- For simplicity, take the case where  $Z$  is a continuous scalar r.v.
- Let  $f_{V|Z}(v|z)$  denote the density of  $V$  conditional on  $Z = z$ , and assume that this density is differentiable in  $z$ .

- Then we obtain

$$\begin{aligned} \frac{\partial E(Y|X=x, Z=z)}{\partial z} &= E(\Delta|X=x, V=\mu_D(z))\mu_D'(z)f_{V|Z,X}(v|x, \mu_D(z)) \\ &+ \left[ \int_{-\infty}^{\mu_D(z)} E(Y_1|X=x, V=v) \frac{\partial f_{V|Z,X}(v|z, x)}{\partial z} dv \right. \\ &\left. + \int_{\mu_D(z)}^{\infty} E(Y_0|X=x, V=v) \frac{\partial f_{V|Z,X}(v|z, x)}{\partial z} dv \right], \end{aligned}$$

and

$$\frac{\partial \Pr(D=1|Z=z)}{\partial z} = f_{V|Z,X}(v|x, \mu_D(z))\mu_D'(z) + \int_{-\infty}^{\mu_D(z)} \frac{\partial f_{V|Z,X}(v|z, x)}{\partial z} dv.$$

- Then we obtain

$$\begin{aligned} \frac{\partial E(Y|X = x, Z = z)}{\partial z} &= E(\Delta|X = x, V = \mu_D(z))\mu'_D(z)f_{V|Z,X}(v|x, \mu_D(z)) \\ &+ \left[ \int_{-\infty}^{\mu_D(z)} E(Y_1|X = x, V = v) \frac{\partial f_{V|Z,X}(v|z, x)}{\partial z} dv \right. \\ &\left. + \int_{\mu_D(z)}^{\infty} E(Y_0|X = x, V = v) \frac{\partial f_{V|Z,X}(v|z, x)}{\partial z} dv \right], \end{aligned}$$

and

$$\frac{\partial \Pr(D = 1|Z = z)}{\partial z} = f_{V|Z,X}(v|x, \mu_D(z))\mu'_D(z) + \int_{-\infty}^{\mu_D(z)} \frac{\partial f_{V|Z,X}(v|z, x)}{\partial z} dv.$$

- LIV is the ratio of the two terms.

- Then we obtain

$$\begin{aligned} \frac{\partial E(Y|X = x, Z = z)}{\partial z} &= E(\Delta|X = x, V = \mu_D(z))\mu'_D(z)f_{V|Z,X}(v|x, \mu_D(z)) \\ &+ \left[ \int_{-\infty}^{\mu_D(z)} E(Y_1|X = x, V = v) \frac{\partial f_{V|Z,X}(v|z, x)}{\partial z} dv \right. \\ &\left. + \int_{\mu_D(z)}^{\infty} E(Y_0|X = x, V = v) \frac{\partial f_{V|Z,X}(v|z, x)}{\partial z} dv \right], \end{aligned}$$

and

$$\frac{\partial \Pr(D = 1|Z = z)}{\partial z} = f_{V|Z,X}(v|x, \mu_D(z))\mu'_D(z) + \int_{-\infty}^{\mu_D(z)} \frac{\partial f_{V|Z,X}(v|z, x)}{\partial z} dv.$$

- LIV is the ratio of the two terms.
- Thus, without the independence condition, the relationship between LIV and the MTE breaks down.

## Proof.

(Equation (30))

$$\begin{aligned}
 E(Y_p | X) &= \int E(Y_p | X, V = v, Z_p = z) dF_{V, Z_p | X}(v, z) \\
 &= \int (\mathbf{1}_\Omega(z) E(Y_1 | X, V = v, Z_p = z) \\
 &\quad + \mathbf{1}_{\Omega^c}(z) E(Y_0 | X, V = v, Z_p = z)) dF_{V, Z_p | X}(v, z) \\
 &= \int (\mathbf{1}_\Omega(z) E(Y_1 | X, V = v) + \mathbf{1}_{\Omega^c}(z) E(Y_0 | X, V = v)) dF_{V, Z_p | X}(v, z) \\
 &= \int \left[ \int (\mathbf{1}_\Omega(z) E(Y_1 | X, V = v) + \mathbf{1}_{\Omega^c}(z) E(Y_0 | X, V = v)) dF_{Z_p | X}(z) \right] dF_{V | X}(v) \\
 &= \int [\Pr[Z_p \in \Omega | X] E(Y_1 | X, V = v) \\
 &\quad + (1 - \Pr[Z_p \in \Omega(z) | X]) E(Y_0 | X, V = v)] dF_{V | X}(v)
 \end{aligned}$$

Q.E.D.

## Proof.

where  $\Omega^c(z)$  denotes the complement of  $\Omega(z)$  and where the first equality follows from the law of iterated expectations; the second equality follows by plugging in our threshold crossing model for  $D$ ; the third equality follows from independence  $Z \perp\!\!\!\perp (Y_1, Y_0, V) \mid X$ ; the fourth and fifth equalities follow by an application of Fubini's Theorem and a rearrangement of terms. Fubini's Theorem may be applied by assumption (A-4). Thus comparing policy  $p$  to policy  $p'$ , we obtain (30).

$$E(Y_p \mid X) - E(Y_{p'} \mid X) = \int E(\Delta \mid X, V = v)(\Pr[Z_p \in \Omega \mid X] - \Pr[Z_{p'} \in \Omega \mid X]) dF_{V|X}(v).$$



Proof.

(Equation (32))

$$\begin{aligned}
 E(Y_p | X) &= \int E(Y_p | X, V = v, Z_p = z) dF_{V, Z_p | X}(v, z) \\
 &= \int \left[ \mathbf{1}_{[-\infty, \mu_D(z)]}(v) E(Y_1 | X, Z = z, V = v) + \mathbf{1}_{(\mu_D(z), \infty]}(v) E(Y_0 | X, Z = z, V = v) \right] dF_{V, Z_p | X}(v, z) \\
 &= \int \left[ \mathbf{1}_{[-\infty, \mu_D(z)]}(v) E(Y_1 | X, V = v) + \mathbf{1}_{(\mu_D(z), \infty]}(v) E(Y_0 | X, V = v) \right] dF_{V, Z_p | X}(v, z) \\
 &= \int \left[ \int \left( \mathbf{1}_{[-\infty, \mu_D(z)]}(v) E(Y_1 | X, V = v) + \mathbf{1}_{(\mu_D(z), \infty]}(v) E(Y_0 | X, V = v) \right) dF_{Z_p | V}(z | v) \right] dF_{V | X}(v) \\
 &= \int \left[ (1 - \Pr[\mu_D(Z_p) < v | V = v]) E(Y_1 | X, V = v) + \Pr[\mu_D(Z_p) < v | V = v] E(Y_0 | X, V = v) \right] dF_{V | X}(v),
 \end{aligned}$$

Q.E.D.

## Proof.

where the first equality follows from the law of iterated expectations; the second equality follows by plugging in our model for  $D$ ; the third equality follows from independence  $Z \perp\!\!\!\perp (Y_1, Y_0) \mid X, V$ ; the fourth equality follows by an application of Fubini's Theorem; and the final equality follows immediately. Thus comparing policy  $p$  to policy  $p'$ , we obtain (32) in the text.  $\square$



## Derivation of PRTE and Implications of Noninvariance for PRTE

## Proof.

**(Equation (10))** To simplify the notation, assume that  $\Upsilon(Y) = Y$ . Modifications required for the more general case are obvious. Define  $\mathbf{1}_{\mathcal{P}}(t)$  to be the indicator function for the event  $t \in \mathcal{P}$ . Then  $E(Y_p | X)$

$$\begin{aligned}
 &= \int_0^1 E(Y_p | X, P_p(Z_p) = t) dF_{P_p|X}(t) \\
 &= \int_0^1 \left[ \int_0^1 [\mathbf{1}_{[0,t]}(u_D) E(Y_{1,p} | X, U_D = u_D) + \mathbf{1}_{(t,1]}(u_D) E(Y_{0,p} | X, U_D = u_D)] du_D \right] dF_{P_p|X}(t) \\
 &= \int_0^1 \left[ \int_0^1 [\mathbf{1}_{[u_D,1]}(t) E(Y_{1,p} | X, U_D = u_D) + \mathbf{1}_{(0,u_D]}(t) E(Y_{0,p} | X, U_D = u_D)] dF_{P_p|X}(t) \right] du_D \\
 &= \int_0^1 [(1 - F_{P_p|X}(u_D)) E(Y_{1,p} | X, U_D = u_D) + F_{P_p|X}(u_D) E(Y_{0,p} | X, U_D = u_D)] du_D.
 \end{aligned}$$

This derivation involves changing the order of integration. *Q.E.D.*

## Proof.

Note that from (A-4),

$$E|\mathbf{1}_{[0,t]}(u_D)E(Y_{1,p} | X, U_D = u_D) + \mathbf{1}_{(t,1]}(u_D)E(Y_{0,p} | X, U_D = u_D)| \leq E(|Y_1| + |Y_0|) < \infty,$$

so the change in the order of integration is valid by Fubini's theorem. Comparing policy  $p$  to policy  $p'$ ,

$$E(Y_p | X) - E(Y_{p'} | X) = \int_0^1 E(\Delta | X, U_D = u_D)(F_{P_{p'}|X}(u_D) - F_{P_p|X}(u_D)) du_D,$$

which gives the required weights. (Recall  $\Delta = Y_1 - Y_0$  and from (A-7) we can drop the  $p, p'$  subscripts on outcomes and errors.)  $\square$

- **Relaxing (A-7): Implications of Noninvariance for PRTE.**

- Suppose that all of the assumptions invoked up through Slide 139 are satisfied, including additive separability in the latent index choice equation (7) (equivalently, the monotonicity or uniformity condition).

- Suppose that all of the assumptions invoked up through Slide 139 are satisfied, including additive separability in the latent index choice equation (7) (equivalently, the monotonicity or uniformity condition).
- Impose the normalization that the distribution of  $U_D$  is unit uniform ( $U_D = F_{V|X}(V | X)$ ).

- Suppose that all of the assumptions invoked up through Slide 139 are satisfied, including additive separability in the latent index choice equation (7) (equivalently, the monotonicity or uniformity condition).
- Impose the normalization that the distribution of  $U_D$  is unit uniform ( $U_D = F_{V|X}(V | X)$ ).
- Suppose however, contrary to (A-7), that the distribution of  $(Y_1, Y_0, U_D, X)$  is different under the two regimes  $p$  and  $p'$ .

- Suppose that all of the assumptions invoked up through Slide 139 are satisfied, including additive separability in the latent index choice equation (7) (equivalently, the monotonicity or uniformity condition).
- Impose the normalization that the distribution of  $U_D$  is unit uniform ( $U_D = F_{V|X}(V | X)$ ).
- Suppose however, contrary to (A-7), that the distribution of  $(Y_1, Y_0, U_D, X)$  is different under the two regimes  $p$  and  $p'$ .
- Thus, let  $(Y_{1,p}, Y_{0,p}, U_{D,p}, X_p)$  and  $(Y_{1,p'}, Y_{0,p'}, U_{D,p'}, X_{p'})$  denote the random vectors under regimes  $p$  and  $p'$ , respectively.



- Following the same analysis as used to derive equation (10), the PRTE conditional on  $X$  is given by

$$E(Y_p | X_p = x) - E(Y_{p'} | X_{p'} = x) \\ = \int_0^1 E(Y_{1,p} - Y_{0,p} | X_p = x, U_{D,p} = u) [F_{P_{p'}|X_{p'}}(u | x) - F_{P_p|X_p}(u | x)] du \quad (I)$$

$$+ \int_0^1 [E(Y_{0,p} | X_p = x, U_{D,p} = u) - E(Y_{0,p'} | X_{p'} = x, U_{D,p'} = u)] du \quad (II)$$

$$+ \int_0^1 \left[ (1 - F_{P_{p'}|X_{p'}}(u | x))(E(Y_{1,p} - Y_{0,p} | X_p = x, U_{D,p} = u) - E(Y_{1,p'} - Y_{0,p'} | X_{p'} = x, U_{D,p'} = u)) \right] du . \quad (III)$$

- Thus, when the policy affects the distribution of  $(Y_1, Y_0, U_D, X)$ , the PRTE is given by the sum of three terms: (I) the value of PRTE if the policy did not affect  $(Y_1, Y_0, X, U_D)$ ; (II) the weighted effect of the policy change on  $E(Y_0 | X, U_D)$ ; and (III) the weighted effect of the policy change on MTE.

- Thus, when the policy affects the distribution of  $(Y_1, Y_0, U_D, X)$ , the PRTE is given by the sum of three terms: (I) the value of PRTE if the policy did not affect  $(Y_1, Y_0, X, U_D)$ ; (II) the weighted effect of the policy change on  $E(Y_0 | X, U_D)$ ; and (III) the weighted effect of the policy change on MTE.
- Evaluating the PRTE requires knowledge of the MTE function in both regimes, knowledge of  $E(Y_0 | X = x, U_D = u)$  in both regimes, as well as knowledge of the distribution of  $P(Z)$  in both regimes.

- Note, however, that if we assume that the distribution of  $(Y_{1,p}, Y_{0,p}, U_{D,p})$  conditional on  $X_p = x$  equals the distribution of  $(Y_{1,p'}, Y_{0,p'}, U_{D,p'})$  conditional on  $X_{p'} = x$ , then
 
$$E(Y_{1,p} \mid U_{D,p} = u, X_p = x) = E(Y_{1,p'} \mid U_{D,p'} = u, X_{p'} = x),$$

$$E(Y_{0,p} \mid U_{D,p} = u, X_p = x) = E(Y_{0,p'} \mid U_{D,p'} = u, X_{p'} = x),$$
 and thus the last two terms vanish and the expression for PRTE simplifies to the expression of equation (10).

## Deriving the IV Weights on MTE

- We consider instrumental variables conditional on  $X = x$  using a general function of  $Z$  as an instrument.

## Deriving the IV Weights on MTE

- We consider instrumental variables conditional on  $X = x$  using a general function of  $Z$  as an instrument.
- To simplify the notation, we keep the conditioning on  $X$  implicit.

## Deriving the IV Weights on MTE

- We consider instrumental variables conditional on  $X = x$  using a general function of  $Z$  as an instrument.
- To simplify the notation, we keep the conditioning on  $X$  implicit.
- Let  $J(Z)$  be any function of  $Z$  such that  $\text{Cov}(J(Z), D) \neq 0$ .

## Deriving the IV Weights on MTE

- We consider instrumental variables conditional on  $X = x$  using a general function of  $Z$  as an instrument.
- To simplify the notation, we keep the conditioning on  $X$  implicit.
- Let  $J(Z)$  be any function of  $Z$  such that  $\text{Cov}(J(Z), D) \neq 0$ .
- Consider the population analogue of the IV estimator,

$$\frac{\text{Cov}(J(Z), Y)}{\text{Cov}(J(Z), D)}$$



- First consider the numerator of this expression,

$$\begin{aligned}\text{Cov}(J(Z), Y) &= E([J(Z) - E(J(Z))] Y) \\ &= E((J(Z) - E(J(Z))) (Y_0 + D(Y_1 - Y_0))) \\ &= E((J(Z) - E(J(Z))) D(Y_1 - Y_0))\end{aligned}$$

where the second equality comes from substituting in the definition of  $Y$  and the third equality follows from conditional independence assumption (A-1).

- First consider the numerator of this expression,

$$\begin{aligned}
 \text{Cov}(J(Z), Y) &= E([J(Z) - E(J(Z))] Y) \\
 &= E((J(Z) - E(J(Z))) (Y_0 + D(Y_1 - Y_0))) \\
 &= E((J(Z) - E(J(Z))) D(Y_1 - Y_0))
 \end{aligned}$$

where the second equality comes from substituting in the definition of  $Y$  and the third equality follows from conditional independence assumption (A-1).

- Define  $\tilde{J}(Z) \equiv J(Z) - E(J(Z))$ .

- Then

$$\begin{aligned}
 \text{Cov}(J(Z), Y) &= E\left(\tilde{J}(Z) \mathbf{1}[U_D \leq P(Z)] (Y_1 - Y_0)\right) \\
 &= E\left(\tilde{J}(Z) \mathbf{1}[U_D \leq P(Z)] E(Y_1 - Y_0 \mid Z, U_D)\right) \\
 &= E\left(\tilde{J}(Z) \mathbf{1}[U_D \leq P(Z)] E(Y_1 - Y_0 \mid U_D)\right) \\
 &= E_{U_D}\left(E_Z\left[\tilde{J}(Z) \mathbf{1}[U_D \leq P(Z)] \mid U_D\right] E(Y_1 - Y_0 \mid U_D)\right) \\
 &= \int_0^1 \left\{ E(\tilde{J}(Z) \mid P(Z) \geq u_D) \Pr(P(Z) \geq u_D) E(Y_1 - Y_0 \mid U_D = u_D) \right\} du_D \\
 &= \int_0^1 \Delta^{\text{MTE}}(x, u_D) E(\tilde{J}(Z) \mid P(Z) \geq u_D) \Pr(P(Z) \geq u_D) du_D
 \end{aligned}$$

- The first equality follows from plugging in the model for  $D$ ; the second equality follows from the law of iterated expectations with the inside expectation conditional on  $(Z, U_D)$ ; the third equality follows from conditional independence assumption (A-1); the fourth equality follows from Fubini's Theorem and the law of iterated expectations with the inside expectation conditional on  $(U_D = u_D)$ ; (and implicitly on  $X$ ); this allows to reverse the order of integration in a multiple integral; the fifth equality follows from the normalization that  $U_D$  is distributed unit uniform conditional on  $X$ ; and the final equality follows from plugging in the definition of  $\Delta^{\text{MTE}}$ .

- Next consider the denominator of the IV estimand.

- Next consider the denominator of the IV estimand.
- Observe that by iterated expectations

$$\text{Cov}(J(Z), D) = \text{Cov}(J(Z), P(Z)).$$

- Next consider the denominator of the IV estimand.
- Observe that by iterated expectations

$$\text{Cov}(J(Z), D) = \text{Cov}(J(Z), P(Z)).$$

- Thus, the population analogue of the IV estimator is given by

$$\int_0^1 \Delta^{\text{MTE}}(u_D) \omega(u_D) du_D \quad (89)$$

where

$$\omega(u_D) = \frac{E(\tilde{J}(Z) \mid P(Z) \geq u_D) \Pr(P(Z) \geq u_D)}{\text{Cov}(J(Z), P(Z))} \quad (90)$$

where by assumption  $\text{Cov}(J(Z), P(Z)) \neq 0$ .

- If  $J(Z)$  and  $P(Z)$  are continuous random variables, then an interpretation of the weight can be derived from (90) by noting that

$$\begin{aligned} & \int (j - E(J(Z))) \int_{u_D}^1 f_{P,J}(t, j) dt dj \\ &= \int (j - E(J(Z))) f_J(j) \int_{u_D}^1 f_{P|J}(t | J(Z) = j) dt dj. \end{aligned}$$



- If  $J(Z)$  and  $P(Z)$  are continuous random variables, then an interpretation of the weight can be derived from (90) by noting that

$$\begin{aligned} \int (j - E(J(Z))) \int_{u_D}^1 f_{P,J}(t, j) dt dj \\ = \int (j - E(J(Z))) f_J(j) \int_{u_D}^1 f_{P|J}(t | J(Z) = j) dt dj. \end{aligned}$$

- Write

$$\begin{aligned} \int_{u_D}^1 f_{P|J}(t | J(Z) = j) dt &= 1 - F_{P|J}(u_D | J(Z) = j) \\ &= S_{P|J}(u_D | J(Z) = j) \end{aligned}$$

where  $S_{P|J}(u_D | J(Z) = j)$  is the probability of  $(P(Z) \geq u_D)$  given  $J(Z) = j$  (and implicitly  $X = x$ ).

- Likewise,  $\Pr[P(Z) > U_D \mid J(Z)] = S_{P|J}(U_D \mid J(Z))$ .

- Likewise,  $\Pr[P(Z) > U_D | J(Z)] = S_{P|J}(U_D | J(Z))$ .
- Using these results, we may write the weight as

$$\omega(u_D) = \frac{\text{Cov}(J(Z), S_{P|J}(u_D | J(Z)))}{\text{Cov}(J(Z), S_{P|J}(U_D | J(Z)))}$$

- Likewise,  $\Pr[P(Z) > U_D | J(Z)] = S_{P|J}(U_D | J(Z))$ .
- Using these results, we may write the weight as

$$\omega(u_D) = \frac{\text{Cov}(J(Z), S_{P|J}(u_D | J(Z)))}{\text{Cov}(J(Z), S_{P|J}(U_D | J(Z)))}$$

- For fixed  $u_D$  and  $x$  evaluation points,  $S_{P|J}(u_D | J(Z))$  is a function of the random variable  $J(Z)$ .

- Likewise,  $\Pr[P(Z) > U_D | J(Z)] = S_{P|J}(U_D | J(Z))$ .
- Using these results, we may write the weight as

$$\omega(u_D) = \frac{\text{Cov}(J(Z), S_{P|J}(u_D | J(Z)))}{\text{Cov}(J(Z), S_{P|J}(U_D | J(Z)))}$$

- For fixed  $u_D$  and  $x$  evaluation points,  $S_{P|J}(u_D | J(Z))$  is a function of the random variable  $J(Z)$ .
- The numerator of the preceding expression is the covariance between  $J(Z)$  and the probability that the random variable  $P(Z)$  is greater than the evaluation point  $u_D$  conditional on  $J(Z)$ .

- $S_{P|J}(U_D | J(Z))$  is a function of the random variables  $U_D$  and  $J(Z)$ .

- $S_{P|J}(U_D | J(Z))$  is a function of the random variables  $U_D$  and  $J(Z)$ .
- The denominator of the above expression is the covariance between  $J(Z)$  and the probability that the random variable  $P(Z)$  is greater than the random variable  $U_D$  conditional on  $J(Z)$ .

- $S_{P|J}(U_D | J(Z))$  is a function of the random variables  $U_D$  and  $J(Z)$ .
- The denominator of the above expression is the covariance between  $J(Z)$  and the probability that the random variable  $P(Z)$  is greater than the random variable  $U_D$  conditional on  $J(Z)$ .
- Thus, it is clear that if the covariance between  $J(Z)$  and the conditional probability that  $(P(Z) > u_D)$  given  $J(Z)$  is positive for all  $u_D$ , then the weights are positive.



- $S_{P|J}(U_D | J(Z))$  is a function of the random variables  $U_D$  and  $J(Z)$ .
- The denominator of the above expression is the covariance between  $J(Z)$  and the probability that the random variable  $P(Z)$  is greater than the random variable  $U_D$  conditional on  $J(Z)$ .
- Thus, it is clear that if the covariance between  $J(Z)$  and the conditional probability that  $(P(Z) > u_D)$  given  $J(Z)$  is positive for all  $u_D$ , then the weights are positive.
- The conditioning is trivially satisfied if  $J(Z) = P(Z)$ , so the weights are positive and IV estimates a gross treatment effect.

- $S_{P|J}(U_D | J(Z))$  is a function of the random variables  $U_D$  and  $J(Z)$ .
- The denominator of the above expression is the covariance between  $J(Z)$  and the probability that the random variable  $P(Z)$  is greater than the random variable  $U_D$  conditional on  $J(Z)$ .
- Thus, it is clear that if the covariance between  $J(Z)$  and the conditional probability that  $(P(Z) > u_D)$  given  $J(Z)$  is positive for all  $u_D$ , then the weights are positive.
- The conditioning is trivially satisfied if  $J(Z) = P(Z)$ , so the weights are positive and IV estimates a gross treatment effect.
- If the  $J(Z)$  and  $P(Z)$  are discrete valued, we obtain expressions and (25) and (26) in the text.

## Yitzhaki's Theorem and the IV Weights (?)

- Assume  $(Y, X)$  i.i.d.,  $E(|Y|) < \infty$ ,  $E(|X|) < \infty$ ,  
 $g(X) = E(Y | X)$ ,  $g'(X)$  exists and  $E(|g'(x)|) < \infty$ .

## Yitzhaki's Theorem and the IV Weights (?)

- Assume  $(Y, X)$  i.i.d.,  $E(|Y|) < \infty$ ,  $E(|X|) < \infty$ ,  
 $g(X) = E(Y | X)$ ,  $g'(X)$  exists and  $E(|g'(x)|) < \infty$ .
- Let  $\mu_Y = E(Y)$  and  $\mu_X = E(X)$ .

## Yitzhaki's Theorem and the IV Weights (?)

- Assume  $(Y, X)$  i.i.d.,  $E(|Y|) < \infty$ ,  $E(|X|) < \infty$ ,  
 $g(X) = E(Y | X)$ ,  $g'(X)$  exists and  $E(|g'(x)|) < \infty$ .
- Let  $\mu_Y = E(Y)$  and  $\mu_X = E(X)$ .
- Then,

$$\frac{\text{Cov}(Y, X)}{\text{Var}(X)} = \int_{-\infty}^{\infty} g'(t) \omega(t) dt,$$

where

$$\begin{aligned} \omega(t) &= \frac{1}{\text{Var}(X)} \int_t^{\infty} (x - \mu_X) f_X(x) dx \\ &= \frac{1}{\text{Var}(X)} E(X - \mu_X | X > t) \Pr(X > t). \end{aligned}$$

## Proof.

$$\begin{aligned}\text{Cov}(Y, X) &= \text{Cov}(E(Y | X), X) = \text{Cov}(g(X), X) \\ &= \int_{-\infty}^{\infty} g(t)(t - \mu_X) f_X(t) dt.\end{aligned}$$

Integration by parts implies that

$$\begin{aligned}&= g(t) \int_{-\infty}^t (x - \mu_X) f_X(x) dx \Big|_{-\infty}^{\infty} \\ &\quad - \int_{-\infty}^{\infty} g'(t) \int_{-\infty}^t (x - \mu_X) f_X(x) dx dt \\ &= \int_{-\infty}^{\infty} g'(t) \int_t^{\infty} (x - \mu_X) f_X(x) dx dt,\end{aligned}$$

since  $E(X - \mu_X) = 0$  and the first term in the first expression vanishes.

- Therefore,

$$\text{Cov}(Y, X) = \int_{-\infty}^{\infty} g'(t) E(X - \mu_X | X > t) \Pr(X > t) dt,$$

so

$$\omega(t) = \frac{1}{\text{Var}(X)} E(X - \mu_X | X > t) \Pr(X > t). \quad \blacksquare$$

Notice that:

- (i) The weights are non-negative ( $\omega(t) \geq 0$ ).
- (ii) They integrate to one (use an integration by parts formula)
- (iii)  $\omega(t) \rightarrow 0$  when  $t \rightarrow -\infty$ , and  $\omega(t) \rightarrow 0$  when  $t \rightarrow \infty$ .



- We get the formula in the text when we use  $P(Z)$ , with a suitably defined domain, in place of  $X$ .

- We get the formula in the text when we use  $P(Z)$ , with a suitably defined domain, in place of  $X$ .
- We apply Yitzhaki's result to the treatment effect model:

$$Y = \alpha + \beta D + \varepsilon,$$

$$\begin{aligned} E(Y | P(Z)) &= \alpha + E(\beta | D = 1, P(Z)) P(Z) \\ &= \alpha + E(\beta | P(Z) > u_D, P(Z)) P(Z) \\ &= g(P(Z)). \end{aligned}$$

- We get the formula in the text when we use  $P(Z)$ , with a suitably defined domain, in place of  $X$ .
- We apply Yitzhaki's result to the treatment effect model:

$$Y = \alpha + \beta D + \varepsilon,$$

$$\begin{aligned} E(Y | P(Z)) &= \alpha + E(\beta | D = 1, P(Z)) P(Z) \\ &= \alpha + E(\beta | P(Z) > u_D, P(Z)) P(Z) \\ &= g(P(Z)). \end{aligned}$$

- By the law of iterated expectations, we eliminate the conditioning on  $D = 0$ .

- Using our previous results for OLS,

$$IV = \frac{\text{Cov}(Y, P(Z))}{\text{Cov}(D, P(Z))} = \int g'(t) \omega(t) dt,$$

$$g'(t) = \left. \frac{\partial [E(\beta \mid D = 1, P(Z))] P(Z)}{\partial P(Z)} \right|_{P(Z)=t},$$

$$\omega(t) = \frac{\int_t^1 [\varphi - E(P(Z))] f_P(\varphi) d\varphi}{\text{Cov}(P(Z), D)}.$$

- Using our previous results for OLS,

$$IV = \frac{\text{Cov}(Y, P(Z))}{\text{Cov}(D, P(Z))} = \int g'(t) \omega(t) dt,$$

$$g'(t) = \left. \frac{\partial [E(\beta | D = 1, P(Z))] P(Z)}{\partial P(Z)} \right|_{P(Z)=t},$$

$$\omega(t) = \frac{\int_t^1 [\varphi - E(P(Z))] f_P(\varphi) d\varphi}{\text{Cov}(P(Z), D)}.$$

- Under (A-1) to (A-5) and separability,  $g'(t) = \Delta^{\text{MTE}}(t)$  but  $g'(t) = \text{LIV}$ , for  $P(Z)$  as an instrument.

## Relationship of our Weights to the Yitzhaki Weights

- Under our assumptions the Yitzhaki weights and ours are equivalent.

## Relationship of our Weights to the Yitzhaki Weights

- Under our assumptions the Yitzhaki weights and ours are equivalent.
- Using (22),

$$\begin{aligned}\text{Cov}(J(Z), Y) &= E(Y \cdot \tilde{J}) = E(E(Y|Z) \cdot \tilde{J}(Z)) \\ &= E(E(Y|P(Z)) \cdot \tilde{J}(Z)) = E(g(P(Z)) \cdot \tilde{J}(Z)).\end{aligned}$$

## Relationship of our Weights to the Yitzhaki Weights

- Under our assumptions the Yitzhaki weights and ours are equivalent.
- Using (22),

$$\begin{aligned} \text{Cov}(J(Z), Y) &= E(Y \cdot \tilde{J}) = E(E(Y | Z) \cdot \tilde{J}(Z)) \\ &= E(E(Y | P(Z)) \cdot \tilde{J}(Z)) = E(g(P(Z)) \cdot \tilde{J}(Z)). \end{aligned}$$

- The third equality follows from index sufficiency and  $\tilde{J} = J(Z) - E(J(Z) | P(Z) \geq u_D)$ , where  $E(Y | P(Z)) = g(P(Z))$ .



- Writing out the expectation and assuming that  $J(Z)$  and  $P(Z)$  are continuous random variables with joint density  $f_{P,J}$  and that  $J(Z)$  has support  $[\underline{J}, \bar{J}]$ ,

$$\begin{aligned} \text{Cov}(J(Z), Y) &= \int_0^1 \int_{\underline{J}}^{\bar{J}} g(u_D) \tilde{j} f_{P,J}(u_D, j) \, dj du_D \\ &= \int_0^1 g(u_D) \int_{\underline{J}}^{\bar{J}} \tilde{j} f_{P,J}(u_D, j) \, dj du_D . \end{aligned}$$

- Using an integration by parts argument as in ? and as summarized in ?, we obtain

$$\begin{aligned}
 \text{Cov}(J(Z), Y) &= g(u_D) \int_0^{u_D} \int_{\underline{J}}^{\bar{J}} \tilde{j} f_{P,J}(p, j) \, dj dp \Big|_0^1 \\
 &\quad - \int_0^1 g'(u_D) \int_0^{u_D} \int_{\underline{J}}^{\bar{J}} \tilde{j} f_{P,J}(p, j) \, dj dp du_D \\
 &= \int_0^1 g'(u_D) \int_{u_D}^1 \int_{\underline{J}}^{\bar{J}} \tilde{j} f_{P,J}(p, j) \, dj dp du_D \\
 &= \int_0^1 g'(u_D) E(\tilde{J}(Z) | P(Z) \geq u_D) \Pr(P(Z) \geq u_D) \, du_D,
 \end{aligned}$$

which is then exactly the expression given in (22), where

$$g'(u_D) = \frac{\partial E(Y | P(Z) = p)}{\partial P(Z)} \Big|_{p=u_D} = \Delta^{\text{MTE}}(u_D).$$

## Derivation of the Weights for the Mixture of Normals Example

- Writing  $E_1$  as the expectation for group 1, letting  $\mu_1$  be the mean of  $Z$  for population 1 and  $\mu_{11}$  be the mean of the first component of  $Z$ ,

$$\begin{aligned}
 E_1(Z_1 | \gamma'Z > v) &= \mu_{11} + \frac{\gamma' \Sigma_1^1}{\gamma' \Sigma_1 \gamma} E_1(Z_1 - \mu_1 | \gamma'Z > v) \\
 &= \mu_{11} + \frac{\gamma' \Sigma_1^1}{(\gamma' \Sigma_1 \gamma)^{1/2}} E_1 \left( \frac{\gamma' (Z - \mu_1)}{(\gamma' \Sigma_1 \gamma)^{1/2}} \mid \frac{\gamma' (Z - \mu_1)}{(\gamma' \Sigma_1 \gamma)^{1/2}} > \frac{(v - \gamma' \mu_1)}{(\gamma' \Sigma_1 \gamma)^{1/2}} \right) \\
 &= \mu_{11} + \frac{\gamma' \Sigma_1^1}{(\gamma' \Sigma_1 \gamma)^{1/2}} \lambda \left( \frac{(v - \gamma' \mu_1)}{(\gamma' \Sigma_1 \gamma)^{1/2}} \right),
 \end{aligned}$$

where

$$\lambda(c) = \frac{1}{\sqrt{2\pi}} \frac{e^{-c^2/2}}{\Phi(-c)},$$

where  $\Phi(\cdot)$  is the unit normal cumulative distribution function.

- By the same logic, in the second group:

$$E_2(Z_1 | \gamma'Z > v) = \mu_{21} + \frac{\gamma'\Sigma_2^1}{(\gamma'\Sigma_2\gamma)^{1/2}} \lambda \left( \frac{(v - \gamma'\mu_2)}{(\gamma'\Sigma_2\gamma)^{1/2}} \right).$$

- By the same logic, in the second group:

$$E_2(Z_1 \mid \gamma' Z > v) = \mu_{21} + \frac{\gamma' \Sigma_2^1}{(\gamma' \Sigma_2 \gamma)^{1/2}} \lambda \left( \frac{(v - \gamma' \mu_2)}{(\gamma' \Sigma_2 \gamma)^{1/2}} \right).$$

- Therefore for the overall population we obtain

$$\begin{aligned} & E(Z_1 - E(Z_1) \mid \gamma' Z > v) \Pr(\gamma' Z > v) \\ &= (P_1 \mu_{11} + P_2 \mu_{21}) \Pr(\gamma' Z > v) + \frac{P_1 \gamma' \Sigma_1^1}{(\gamma' \Sigma_1 \gamma)^{1/2} \sqrt{2\pi}} \exp \left[ -\frac{1}{2} \left( \frac{v - \gamma' \mu_1}{(\gamma' \Sigma_1 \gamma)^{1/2}} \right)^2 \right] \\ &+ \frac{P_2 \gamma' \Sigma_2^1}{(\gamma' \Sigma_2 \gamma)^{1/2} \sqrt{2\pi}} \exp \left[ -\frac{1}{2} \left( \frac{v - \gamma' \mu_2}{(\gamma' \Sigma_2 \gamma)^{1/2}} \right)^2 \right] - (P_1 \mu_{11} + P_2 \mu_{21}) \Pr(\gamma' Z > v) \\ &= \frac{P_1 \gamma' \Sigma_1^1}{(\gamma' \Sigma_1 \gamma)^{1/2} \sqrt{2\pi}} \exp \left[ -\frac{1}{2} \left( \frac{v - \gamma' \mu_1}{(\gamma' \Sigma_1 \gamma)^{1/2}} \right)^2 \right] \\ &+ \frac{P_2 \gamma' \Sigma_2^1}{(\gamma' \Sigma_2 \gamma)^{1/2} \sqrt{2\pi}} \exp \left[ -\frac{1}{2} \left( \frac{v - \gamma' \mu_2}{(\gamma' \Sigma_2 \gamma)^{1/2}} \right)^2 \right]. \end{aligned}$$

- We need  $\text{Cov}(D, Z_1)$ .

- We need  $\text{Cov}(D, Z_1)$ .
- To obtain it, observe that

$$\begin{aligned} D &= \mathbf{1}[\gamma'Z - V > 0] \\ E(Z_1 D) &= E(Z_1 \mathbf{1}(\gamma'Z - V \geq 0)). \end{aligned}$$

- We need  $\text{Cov}(D, Z_1)$ .
- To obtain it, observe that

$$D = \mathbf{1}[\gamma'Z - V > 0]$$
$$E(Z_1 D) = E(Z_1 \mathbf{1}(\gamma'Z - V \geq 0)).$$

- Let  $E_1$  denote the expectation for Group 1, and let  $E_2$  denote the expectation for Group 2.



$$\begin{aligned}
 E(Z_1 D) &= \left\{ P_1 \left[ \mu_{11} + \frac{\gamma' \Sigma_1^1}{\gamma' \Sigma_1 \gamma + \sigma_V^2} E_1(Z_1 - \mu_{11} | \gamma' Z - V \geq 0) \right] \right. \\
 &\quad \left. + P_2 \left[ \mu_{21} + \frac{\gamma' \Sigma_2^1}{\gamma' \Sigma_2 \gamma + \sigma_V^2} E_2(Z_1 - \mu_{21} | \gamma' Z - V \geq 0) \right] \right\} \Pr[(\gamma' Z - V) > 0] \\
 &= (P_1 \mu_{11} + P_2 \mu_{21}) \Pr(\gamma' Z - V \geq 0) \\
 &\quad + \frac{P_1 \gamma' \Sigma_1^1}{(\gamma' \Sigma_1 \gamma + \sigma_V^2)^{1/2} \sqrt{2\pi}} \exp \left[ - \left( \frac{-\gamma' \mu_1}{(\gamma' \Sigma_1 \gamma + \sigma_V^2)^{1/2}} \right)^2 \right] \\
 &\quad + \frac{P_2 \gamma' \Sigma_2^1}{(\gamma' \Sigma_2 \gamma + \sigma_V^2)^{1/2} \sqrt{2\pi}} \exp \left[ - \left( \frac{-\gamma' \mu_2}{(\gamma' \Sigma_2 \gamma + \sigma_V^2)^{1/2}} \right)^2 \right].
 \end{aligned}$$

- Because

$$E(D)E(Z_1) = \Pr(\gamma'Z - V \geq 0) (P_1\mu_{11} + P_2\mu_{21})$$

and

$$\text{Cov}(D, Z_1) = E(Z_1D) - E(Z_1)E(D)$$

$$\begin{aligned} \therefore \text{Cov}(D, Z_1) &= \frac{P_1\gamma'\Sigma_1^1}{(\gamma'\Sigma_1\gamma + \sigma_V^2)^{1/2} \sqrt{2\pi}} \exp \left[ - \left( \frac{-\gamma'\mu_1}{(\gamma'\Sigma_1\gamma + \sigma_V^2)^{1/2}} \right)^2 \right] \\ &+ \frac{P_2\gamma'\Sigma_2^1}{(\gamma'\Sigma_2\gamma + \sigma_V^2)^{1/2} \sqrt{2\pi}} \exp \left[ - \left( \frac{-\gamma'\mu_2}{(\gamma'\Sigma_2\gamma + \sigma_V^2)^{1/2}} \right)^2 \right]. \end{aligned}$$

- Thus the IV weights for this set-up are:

$$\begin{aligned} \tilde{\omega}_{IV}(v) = & \left\{ \left[ \frac{P_1 \gamma' \Sigma_1^1}{(\gamma' \Sigma_1 \gamma)^{1/2}} \exp \left[ -\frac{1}{2} \left( \frac{v - \gamma' \mu_1}{(\gamma' \Sigma_1 \gamma)^{1/2}} \right)^2 \right] \right. \right. \\ & \left. \left. + \frac{P_2 \gamma' \Sigma_2^1}{(\gamma' \Sigma_2 \gamma)^{1/2}} \exp \left[ -\frac{1}{2} \left( \frac{v - \gamma' \mu_2}{(\gamma' \Sigma_2 \gamma)^{1/2}} \right)^2 \right] \right] f_V(v) \right\} \\ & \times \left\{ \frac{P_1 \gamma' \Sigma_1^1}{(\gamma' \Sigma_1 \gamma + \sigma_V^2)^{1/2}} \exp \left[ -\left( \frac{-\gamma' \mu_1}{(\gamma' \Sigma_1 \gamma + \sigma_V^2)^{1/2}} \right)^2 \right] \right. \\ & \left. + \frac{P_2 \gamma' \Sigma_2^1}{(\gamma' \Sigma_2 \gamma + \sigma_V^2)^{1/2}} \exp \left[ -\left( \frac{-\gamma' \mu_2}{(\gamma' \Sigma_2 \gamma + \sigma_V^2)^{1/2}} \right)^2 \right] \right\}^{-1}, \end{aligned}$$

where  $\sigma_V^2$  represents the variance of  $V$ .

- Clearly,  $\tilde{\omega}_{IV}(-\infty) = 0$ ,  $\tilde{\omega}_{IV}(\infty) = 0$  and the weights integrate to one over the support of  $V = (-\infty, \infty)$ .

- Clearly,  $\tilde{\omega}_{IV}(-\infty) = 0$ ,  $\tilde{\omega}_{IV}(\infty) = 0$  and the weights integrate to one over the support of  $V = (-\infty, \infty)$ .
- Observe that the weights must be positive if  $P_2 = 0$ .

- Clearly,  $\tilde{\omega}_{IV}(-\infty) = 0$ ,  $\tilde{\omega}_{IV}(\infty) = 0$  and the weights integrate to one over the support of  $V = (-\infty, \infty)$ .
- Observe that the weights must be positive if  $P_2 = 0$ .
- Thus the structure of the covariances of the instrument with the choice index  $\gamma'Z$  is a key determinant of the positivity of the weights for any instrument.

- Clearly,  $\tilde{\omega}_{IV}(-\infty) = 0$ ,  $\tilde{\omega}_{IV}(\infty) = 0$  and the weights integrate to one over the support of  $V = (-\infty, \infty)$ .
- Observe that the weights must be positive if  $P_2 = 0$ .
- Thus the structure of the covariances of the instrument with the choice index  $\gamma'Z$  is a key determinant of the positivity of the weights for any instrument.
- It has nothing to do with the *ceteris paribus* effect of  $Z_1$  on  $\gamma'Z$  or  $P(Z)$  in the general case.

- A necessary condition for  $\omega_{IV} < 0$  over some values of  $v$  is that  $\text{sign}(\gamma' \Sigma_1^1) = -\text{sign}(\gamma' \Sigma_2^1)$ , i.e., that the covariance between  $Z_1$  and  $\gamma' Z$  be of opposite signs in the two subpopulations so  $Z_1$  and  $P(Z)$  have different relationships in the two component populations.



- A necessary condition for  $\omega_{IV} < 0$  over some values of  $v$  is that  $\text{sign}(\gamma' \Sigma_1^1) = -\text{sign}(\gamma' \Sigma_2^1)$ , i.e., that the covariance between  $Z_1$  and  $\gamma' Z$  be of opposite signs in the two subpopulations so  $Z_1$  and  $P(Z)$  have different relationships in the two component populations.
- Without loss of generality, assume that  $\gamma' \Sigma_1^1 > 0$ .

- A necessary condition for  $\omega_{IV} < 0$  over some values of  $v$  is that  $\text{sign}(\gamma' \Sigma_1^1) = -\text{sign}(\gamma' \Sigma_2^1)$ , i.e., that the covariance between  $Z_1$  and  $\gamma' Z$  be of opposite signs in the two subpopulations so  $Z_1$  and  $P(Z)$  have different relationships in the two component populations.
- Without loss of generality, assume that  $\gamma' \Sigma_1^1 > 0$ .
- If it equals zero, we fail the rank condition in the first population and we are back to a one subpopulation model with positive weights.

- The numerator of the expression for  $\omega_{IV}(v)$  switches signs if for some values of  $v$ ,

$$\frac{P_1 \gamma' \Sigma_1^{-1}}{(\gamma' \Sigma_1 \gamma)^{1/2}} \exp \left[ -\frac{1}{2} \left( \frac{v - \gamma' \mu_1}{(\gamma' \Sigma_1 \gamma)^{1/2}} \right)^2 \right] < -\frac{P_2 \gamma' \Sigma_2^{-1}}{(\gamma' \Sigma_2 \gamma)^{1/2}} \exp \left[ -\frac{1}{2} \left( \frac{v - \gamma' \mu_2}{(\gamma' \Sigma_2 \gamma)^{1/2}} \right)^2 \right]$$

while for other values the inequality is reversed.

- The numerator of the expression for  $\omega_{IV}(v)$  switches signs if for some values of  $v$ ,

$$\frac{P_1 \gamma' \Sigma_1^{-1}}{(\gamma' \Sigma_1 \gamma)^{1/2}} \exp \left[ -\frac{1}{2} \left( \frac{v - \gamma' \mu_1}{(\gamma' \Sigma_1 \gamma)^{1/2}} \right)^2 \right] < -\frac{P_2 \gamma' \Sigma_2^{-1}}{(\gamma' \Sigma_2 \gamma)^{1/2}} \exp \left[ -\frac{1}{2} \left( \frac{v - \gamma' \mu_2}{(\gamma' \Sigma_2 \gamma)^{1/2}} \right)^2 \right]$$

while for other values the inequality is reversed.

- Observe that the denominator is a constant.

- The numerator of the expression for  $\omega_{IV}(v)$  switches signs if for some values of  $v$ ,

$$\frac{P_1 \gamma' \Sigma_1^{-1}}{(\gamma' \Sigma_1 \gamma)^{1/2}} \exp \left[ -\frac{1}{2} \left( \frac{v - \gamma' \mu_1}{(\gamma' \Sigma_1 \gamma)^{1/2}} \right)^2 \right] < -\frac{P_2 \gamma' \Sigma_2^{-1}}{(\gamma' \Sigma_2 \gamma)^{1/2}} \exp \left[ -\frac{1}{2} \left( \frac{v - \gamma' \mu_2}{(\gamma' \Sigma_2 \gamma)^{1/2}} \right)^2 \right]$$

while for other values the inequality is reversed.

- Observe that the denominator is a constant.
- Rewriting and taking logarithms, we obtain under the assumption that  $\text{sign}(\gamma' \Sigma_1^{-1}) = -\text{sign}(\gamma' \Sigma_2^{-1})$ , the following expression:

$$\frac{1}{2} \left[ \frac{(v - \gamma' \mu_2)^2}{\gamma' \Sigma_2 \gamma} - \frac{(v - \gamma' \mu_1)^2}{\gamma' \Sigma_1 \gamma} \right] < \ln \left( \frac{1 - P_1}{P_1} \right) + \ln \left[ \frac{-\gamma' \Sigma_2^{-1}}{\gamma' \Sigma_1^{-1}} \right] + \ln \left[ \frac{\gamma' \Sigma_1 \gamma}{\gamma' \Sigma_2 \gamma} \right],$$

where we assume  $0 < P_1 < 1$ .

- Observe that  $\frac{1-P_1}{P_1}$  can be made as large or as small a non-negative number as we like by varying  $P_1$ .

- Observe that  $\frac{1-P_1}{P_1}$  can be made as large or as small a non-negative number as we like by varying  $P_1$ .
- Varying  $(\mu_1, \mu_2)$  does not affect the right hand side.

- Observe that  $\frac{1-P_1}{P_1}$  can be made as large or as small a non-negative number as we like by varying  $P_1$ .
- Varying  $(\mu_1, \mu_2)$  does not affect the right hand side.
- For  $\mu_1 = \mu_2 = 0$ , the inequality becomes

$$\frac{1}{2}v^2 \left[ \frac{1}{\gamma'\Sigma_2\gamma} - \frac{1}{\gamma'\Sigma_1\gamma} \right] < \ln \left( \frac{1-P_1}{P_1} \right) + \ln \left[ \frac{-\gamma'\Sigma_2^1}{\gamma'\Sigma_1^1} \right] + \ln \left[ \frac{\gamma'\Sigma_1\gamma}{\gamma'\Sigma_2\gamma} \right].$$



- Suppose that  $\gamma' \Sigma_2 \gamma < \gamma' \Sigma_1 \gamma$ .

- Suppose that  $\gamma' \Sigma_2 \gamma < \gamma' \Sigma_1 \gamma$ .
- Then the left hand side is positive except when  $v = 0$ .

- Suppose that  $\gamma' \Sigma_2 \gamma < \gamma' \Sigma_1 \gamma$ .
- Then the left hand side is positive except when  $v = 0$ .
- For any fixed  $\gamma, \Sigma_1, \Sigma_2$  we can find a value of  $P_1$  sufficiently small so that right hand side of the equation is positive and for any such value of  $P_1$  there will be a  $v$  sufficiently small for the inequality to be satisfied.

- Suppose that  $\gamma' \Sigma_2 \gamma < \gamma' \Sigma_1 \gamma$ .
- Then the left hand side is positive except when  $v = 0$ .
- For any fixed  $\gamma, \Sigma_1, \Sigma_2$  we can find a value of  $P_1$  sufficiently small so that right hand side of the equation is positive and for any such value of  $P_1$  there will be a  $v$  sufficiently small for the inequality to be satisfied.
- There is also a value of  $v$  that reverses the inequality.

- The inequality is satisfied for some  $v^* \geq 0$ .

- The inequality is satisfied for some  $v^* \geq 0$ .
- But with  $v$  arbitrarily large, the inequality can be reversed so that the weight will switch signs at some value of  $v$ .

- The inequality is satisfied for some  $v^* \geq 0$ .
- But with  $v$  arbitrarily large, the inequality can be reversed so that the weight will switch signs at some value of  $v$ .
- The key necessary condition is that  $\text{Cov}(Z_1, \gamma'Z)$  be of opposite signs in the two subpopulations.

- The inequality is satisfied for some  $v^* \geq 0$ .
- But with  $v$  arbitrarily large, the inequality can be reversed so that the weight will switch signs at some value of  $v$ .
- The key necessary condition is that  $\text{Cov}(Z_1, \gamma' Z)$  be of opposite signs in the two subpopulations.
- Using  $Z_1$  as an IV, but not conditioning or controlling for the other components of  $Z$ , produces sometimes negative and sometimes positive movements in the components of  $Z_2, \dots, Z_k$  which can offset the *ceteris paribus* ( $Z_2 = z_2, \dots, Z_k = z_k$ ) movements of  $Z_1$ .



## Local Instrumental Variables for the Random Coefficient Model

- Consider the model:

$$D = \mathbf{1}[Z\gamma \geq 0],$$

where  $\gamma$  is a random variable.

## Local Instrumental Variables for the Random Coefficient Model

- Consider the model:

$$D = \mathbf{1}[Z\gamma \geq 0],$$

where  $\gamma$  is a random variable.

- For ease of exposition, we leave implicit the conditioning on  $X$  covariates.

## Local Instrumental Variables for the Random Coefficient Model

- Consider the model:

$$D = \mathbf{1}[Z\gamma \geq 0],$$

where  $\gamma$  is a random variable.

- For ease of exposition, we leave implicit the conditioning on  $X$  covariates.
- Assume that  $(Y_0, Y_1, \gamma) \perp\!\!\!\perp Z$ .

## Local Instrumental Variables for the Random Coefficient Model

- Consider the model:

$$D = \mathbf{1}[Z\gamma \geq 0],$$

where  $\gamma$  is a random variable.

- For ease of exposition, we leave implicit the conditioning on  $X$  covariates.
- Assume that  $(Y_0, Y_1, \gamma) \perp\!\!\!\perp Z$ .
- Assume that  $\gamma$  has a density that is absolutely continuous with respect to Lebesgue measure on  $\mathbb{R}^K$ .

- We have

$$E(Y | Z = z) = E(DY_1 | Z = z) + E((1 - D)Y_0 | Z = z).$$

- We have

$$E(Y | Z = z) = E(DY_1 | Z = z) + E((1 - D)Y_0 | Z = z).$$

- To simplify the exposition, consider the first term,  
 $E(DY_1 | Z = z)$ .

- We have

$$E(Y | Z = z) = E(DY_1 | Z = z) + E((1 - D)Y_0 | Z = z).$$

- To simplify the exposition, consider the first term,  $E(DY_1 | Z = z)$ .
- In this proof, let  $Z^{[K]}$  denote the  $K$ th element of  $Z$  and  $Z^{[-K]}$  denote all other elements of  $Z$ , and write  $Z = (Z^{[-K]}, Z^{[K]})$ .

- Using the model, the independence assumption, and the law of iterated expectations, we have

$$\begin{aligned} E(DY | Z = z) &= E\left(\mathbf{1}[z\gamma \geq 0]Y_1\right) = E\left(\mathbf{1}[z\gamma \geq 0]E(Y_1 | \gamma)\right) \\ &= E\left(\mathbf{1}\left\{z^{[K]}\gamma^{[K]} \geq -z^{[-K]}\gamma^{[-K]}\right\} E(Y_1 | \gamma)\right), \end{aligned}$$

where the final outer expectation is over  $\gamma$ .



- Using the model, the independence assumption, and the law of iterated expectations, we have

$$\begin{aligned} E(DY | Z = z) &= E\left(\mathbf{1}[z\gamma \geq 0]Y_1\right) = E\left(\mathbf{1}[z\gamma \geq 0]E(Y_1 | \gamma)\right) \\ &= E\left(\mathbf{1}\left\{z^{[K]}\gamma^{[K]} \geq -z^{[-K]}\gamma^{[-K]}\right\} E(Y_1 | \gamma)\right), \end{aligned}$$

where the final outer expectation is over  $\gamma$ .

- Consider taking the derivative with respect to the  $K$ th element of  $Z$  assumed to be continuous.

- Using the model, the independence assumption, and the law of iterated expectations, we have

$$\begin{aligned} E(DY | Z = z) &= E(\mathbf{1}[z\gamma \geq 0] Y_1) = E(\mathbf{1}[z\gamma \geq 0] E(Y_1 | \gamma)) \\ &= E(\mathbf{1}\left\{z^{[K]}\gamma^{[K]} \geq -z^{[-K]}\gamma^{[-K]}\right\} E(Y_1 | \gamma)), \end{aligned}$$

where the final outer expectation is over  $\gamma$ .

- Consider taking the derivative with respect to the  $K$ th element of  $Z$  assumed to be continuous.
- Partition  $z$ ,  $\gamma$ , and  $g$  as  $z = (z^{[-K]}, z^{[K]})$ ,  $\gamma = (\gamma^{[-K]}, \gamma^{[K]})$ , and  $g = (g^{[-K]}, g^{[K]})$ , where  $z$  is a realization of  $Z$  and  $g$  is a realization of  $\gamma$ .

- For simplicity, suppose that the  $K$ th element of  $z$  is positive,  $z^{[K]} > 0$ .

- For simplicity, suppose that the  $K$ th element of  $z$  is positive,  $z^{[K]} > 0$ .
- We obtain

$$\begin{aligned} E(DY | Z = z) &= E\left[E\left(\mathbf{1}\left\{z^{[K]}\gamma^{[K]} \geq -z^{[-K]}\gamma^{[-K]}\right\} E(Y_1 | \gamma) \mid \gamma^{[-K]}\right)\right] \\ &= E\left[E\left(\mathbf{1}\left\{\gamma^{[K]} \geq \frac{-z^{[-K]}\gamma^{[-K]}}{z^{[K]}}\right\} E(Y_1 | \gamma) \mid \gamma^{[-K]}\right)\right], \end{aligned}$$

where the inside expectation is over  $\gamma^{[K]}$  conditional on  $\gamma^{[-K]}$ , i.e., is over the  $K$ th element of  $\gamma$  conditional on all other components of  $\gamma$ .

- Computing the derivative with respect to  $z^{[K]}$ , we obtain

$$\frac{\partial}{\partial z^{[K]}} E(DY | Z = z) = \int E(Y_1 | \gamma = M(g^{[-K]})) \tilde{w}(g^{[-K]}) dg^{[-K]},$$

where

$$M(g^{[-K]}) = ((g^{[-K]})', \frac{-z^{[-K]} g^{[-K]}}{z^{[K]}})',$$

$$\text{and } \tilde{w}(g^{[-K]}) = \frac{z^{[-K]} g^{[-K]}}{(z^{[K]})^2} f\left(g^{[-K]}, \frac{-z^{[-K]} g^{[-K]}}{z^{[K]}}\right),$$

with  $f(\cdot)$  the density of  $\gamma$  (with respect to Lebesgue measure), and where for notational simplicity we suppress the dependence of the function  $M(\cdot)$  and the weights  $\tilde{w}(\cdot)$  on the  $z$  evaluation point.

- Computing the derivative with respect to  $z^{[K]}$ , we obtain

$$\frac{\partial}{\partial z^{[K]}} E(DY | Z = z) = \int E(Y_1 | \gamma = M(g^{[-K]})) \tilde{w}(g^{[-K]}) dg^{[-K]},$$

where

$$M(g^{[-K]}) = ((g^{[-K]})', \frac{-z^{[-K]} g^{[-K]}}{z^{[K]}})',$$

$$\text{and } \tilde{w}(g^{[-K]}) = \frac{z^{[-K]} g^{[-K]}}{(z^{[K]})^2} f\left(g^{[-K]}, \frac{-z^{[-K]} g^{[-K]}}{z^{[K]}}\right),$$

with  $f(\cdot)$  the density of  $\gamma$  (with respect to Lebesgue measure), and where for notational simplicity we suppress the dependence of the function  $M(\cdot)$  and the weights  $\tilde{w}(\cdot)$  on the  $z$  evaluation point.

- In this expression, we are averaging over  $E(Y_1 | \gamma = g)$ , but only over  $g$  evaluation points such that  $zg = 0$ .

- In particular, the expression averages over the  $K - 1$  space of  $g^{[-K]}$ , while for each potential realization of  $g^{[-K]}$  it is filling in the value of  $g^{[K]}$  such that  $z^{[K]}g^{[K]} = -z^{[-K]}g^{[-K]}$  so that  $z^{[K]}g^{[K]} + z^{[-K]}g^{[-K]} = 0$ .

- In particular, the expression averages over the  $K - 1$  space of  $g^{[-K]}$ , while for each potential realization of  $g^{[-K]}$  it is filling in the value of  $g^{[K]}$  such that  $z^{[K]}g^{[K]} = -z^{[-K]}g^{[-K]}$  so that  $z^{[K]}g^{[K]} + z^{[-K]}g^{[-K]} = 0$ .
- Note that the weights  $\tilde{w}(g^{[-K]})$  will be zero for any  $g^{[-K]}$  such that  $f(g^{[-K]}, \frac{-z^{[-K]}g^{[-K]}}{z^{[K]}}) = 0$ , i.e., the weights will be zero for any  $g^{[-K]}$  such that there does not exist  $g^{[K]}$  in the conditional support of  $\gamma^{[K]}$  with  $z^{[K]}g^{[K]} = -z^{[-K]}g^{[-K]}$ .



- Following the same logic for  $E((1 - D)Y_0 | Z = z)$ , we obtain

$$\frac{\partial}{\partial z^{[K]}} E((1 - D)Y | Z = z) = - \int E(Y_0 | \gamma = M(g^{[-K]})) \tilde{w}(g^{[-K]}) dg^{[-K]}$$

and likewise have

$$\frac{\partial}{\partial z^{[K]}} \Pr(D = 1 | Z = z) = \int \tilde{w}(g^{[-K]}) dg^{[-K]}$$

so that

$$\frac{\frac{\partial}{\partial z^{[K]}} E(Y | Z = z)}{\frac{\partial}{\partial z^{[K]}} \Pr(D = 1 | Z = z)} = \int E(Y_1 - Y_0 | \gamma = M(g^{[-K]})) w(g^{[-K]}) dg^{[-K]},$$

where

$$w(g^{[-K]}) = \tilde{w}(g^{[-K]}) / \int \tilde{w}(g^{[-K]}) dg^{[-K]}.$$

- Now consider the question of whether this expression will have both positive and negative weights.

- Now consider the question of whether this expression will have both positive and negative weights.
- Recall that  $\tilde{w}(g^{[-K]}) = \frac{z^{[-K]}g^{[-K]}}{(z^{[K]})^2} f(g^{[-K]}, \frac{-z^{[-K]}g^{[-K]}}{z^{[K]}})$ .

- Now consider the question of whether this expression will have both positive and negative weights.
- Recall that  $\tilde{w}(g^{[-K]}) = \frac{z^{[-K]}g^{[-K]}}{(z^{[K]})^2} f(g^{[-K]}, \frac{-z^{[-K]}g^{[-K]}}{z^{[K]}})$ .
- Thus,

$$\tilde{w}(g^{[-K]}) \geq 0 \quad \text{if } z^{[-K]}g^{[-K]} > 0, \quad \tilde{w}(g^{[-K]}) \leq 0 \quad \text{if } z^{[-K]}g^{[-K]} < 0,$$

and will be nonzero if  $z^{[-K]}g^{[-K]} \neq 0$  and there exists  $g^{[K]}$  in the conditional support of  $\gamma^{[K]}$  with  $z^{[K]}g^{[K]} = z^{[-K]}g^{[-K]}$ , i.e., with  $zg = 0$ .

- Now consider the question of whether this expression will have both positive and negative weights.

- Recall that  $\tilde{w}(g^{[-K]}) = \frac{z^{[-K]}g^{[-K]}}{(z^{[K]})^2} f(g^{[-K]}, \frac{-z^{[-K]}g^{[-K]}}{z^{[K]}})$ .

- Thus,

$$\tilde{w}(g^{[-K]}) \geq 0 \quad \text{if } z^{[-K]}g^{[-K]} > 0, \quad \tilde{w}(g^{[-K]}) \leq 0 \quad \text{if } z^{[-K]}g^{[-K]} < 0,$$

and will be nonzero if  $z^{[-K]}g^{[-K]} \neq 0$  and there exists  $g^{[K]}$  in the conditional support of  $\gamma^{[K]}$  with  $z^{[K]}g^{[K]} = z^{[-K]}g^{[-K]}$ , i.e., with  $zg = 0$ .

- We thus have that there will be both positive and negative weights on the MTE if there exist values of  $g$  in the support of  $\gamma$  with both  $z^{[-K]}g^{[-K]} > 0$  and  $zg = 0$ , and there exist other values of  $g$  in the support of  $\gamma$  with  $z^{[-K]}g^{[-K]} < 0$  and  $zg = 0$ .

## Generalized Ordered Choice Model with Stochastic Thresholds

- The ordered choice model presented in the text with parameterized, but nonstochastic, thresholds is analyzed in ? who establish its nonparametric identifiability under the conditions they specify.

## Generalized Ordered Choice Model with Stochastic Thresholds

- The ordered choice model presented in the text with parameterized, but nonstochastic, thresholds is analyzed in ? who establish its nonparametric identifiability under the conditions they specify.
- Treating the  $W_s$  (or components of it) as unobservables, we obtain the generalized ordered choice model analyzed in ? and ?.

## Generalized Ordered Choice Model with Stochastic Thresholds

- The ordered choice model presented in the text with parameterized, but nonstochastic, thresholds is analyzed in ? who establish its nonparametric identifiability under the conditions they specify.
- Treating the  $W_s$  (or components of it) as unobservables, we obtain the generalized ordered choice model analyzed in ? and ?.
- In this appendix, we present the main properties of this more general model.



- The thresholds are now written as  $Q_s + C_s(W_s)$  in place of  $C_s(W_s)$ , where  $Q_s$  is a random variable.

- The thresholds are now written as  $Q_s + C_s(W_s)$  in place of  $C_s(W_s)$ , where  $Q_s$  is a random variable.
- In addition to the order on the  $C_s(W_s)$  in the text, we impose the order  $Q_s + C_s(W_s) \geq Q_{s-1} + C_{s-1}(W_{s-1})$ ,  $s = 2, \dots, \bar{S} - 1$ .

- The thresholds are now written as  $Q_s + C_s(W_s)$  in place of  $C_s(W_s)$ , where  $Q_s$  is a random variable.
- In addition to the order on the  $C_s(W_s)$  in the text, we impose the order  $Q_s + C_s(W_s) \geq Q_{s-1} + C_{s-1}(W_{s-1})$ ,  $s = 2, \dots, \bar{S} - 1$ .
- We impose the requirement that  $Q_{\bar{S}} = \infty$  and  $Q_0 = -\infty$ .

- The thresholds are now written as  $Q_s + C_s(W_s)$  in place of  $C_s(W_s)$ , where  $Q_s$  is a random variable.
- In addition to the order on the  $C_s(W_s)$  in the text, we impose the order  $Q_s + C_s(W_s) \geq Q_{s-1} + C_{s-1}(W_{s-1})$ ,  $s = 2, \dots, \bar{S} - 1$ .
- We impose the requirement that  $Q_{\bar{S}} = \infty$  and  $Q_0 = -\infty$ .
- The latent index  $D_s^*$  is as defined in the text, but now

$$\begin{aligned} D_s &= \mathbf{1}[C_{s-1}(W_{s-1}) + Q_{s-1} < \mu_D(Z) - V \leq C_s(W_s) + Q_s] \\ &= \mathbf{1}[l_{s-1}(Z, W_{s-1}) - Q_{s-1} > V \geq l_s(Z, W_s) - Q_s], \end{aligned}$$

where  $l_s = \mu_D(Z) - C_s(W_s)$ .

- Using the fact that  $I_s(Z, W_s) - Q_s < I_{s-1}(Z, W_{s-1}) - Q_{s-1}$ , we obtain

$$\begin{aligned} \mathbf{1}[I_{s-1}(Z, W_{s-1}) - Q_{s-1} > V \geq I_s(Z, W_s) - Q_s] &= \\ &= \mathbf{1}[V + Q_{s-1} < I_{s-1}(Z, W_{s-1})] - \mathbf{1}[V + Q_s \leq I_s(Z, W_s)]. \end{aligned}$$

The nonparametric identifiability of this choice model is established in ? and ?.

- We retain assumptions (OC-2) – (OC-6), but alter (OC-1) to

- We retain assumptions (OC-2) – (OC-6), but alter (OC-1) to (OC-1)'  $(Q_s, U_s, V) \perp\!\!\!\perp (Z, W) \mid X, \quad s = 1, \dots, \bar{S}$ .

- ? shows that this model with no transition specific instruments (with  $W_s$  degenerate for each  $s$ ) implies and is implied by the independence and monotonicity conditions of ? for an ordered model.



- ? shows that this model with no transition specific instruments (with  $W_s$  degenerate for each  $s$ ) implies and is implied by the independence and monotonicity conditions of ? for an ordered model.
- Define  $Q = (Q_1, \dots, Q_{\bar{S}})$ .

- ? shows that this model with no transition specific instruments (with  $W_s$  degenerate for each  $s$ ) implies and is implied by the independence and monotonicity conditions of ? for an ordered model.
- Define  $Q = (Q_1, \dots, Q_{\bar{S}})$ .
- Redefine  $\pi_s(Z, W_s) = F_{V+Q_s}(\mu_D(Z) + C_s(W_s))$  and define  $\pi(Z, W) = [\pi_1(Z, W_1), \dots, \pi_{\bar{S}-1}(Z, W_{\bar{S}-1})]$ .

- ? shows that this model with no transition specific instruments (with  $W_s$  degenerate for each  $s$ ) implies and is implied by the independence and monotonicity conditions of ? for an ordered model.
- Define  $Q = (Q_1, \dots, Q_{\bar{S}})$ .
- Redefine  $\pi_s(Z, W_s) = F_{V+Q_s}(\mu_D(Z) + C_s(W_s))$  and define  $\pi(Z, W) = [\pi_1(Z, W_1), \dots, \pi_{\bar{S}-1}(Z, W_{\bar{S}-1})]$ .
- Redefine  $U_{D,s} = F_{V+Q_s}(V + Q_s)$ .

- We have that

$$\begin{aligned}
 E(Y | Z, W) &= \\
 &= E \left( \sum_{s=1}^{\bar{S}} \mathbf{1}[I_{s-1}(Z, W_{s-1}) - Q_{s-1} > V \geq I_s(Z, W_s) - Q_s] Y_s \mid Z, W \right) \\
 &= \sum_{s=1}^{\bar{S}} \left( E(\mathbf{1}[V + Q_{s-1} < I_{s-1}(Z, W_{s-1})] Y_s \mid Z, W) \right. \\
 &\quad \left. - E(\mathbf{1}[V + Q_s \leq I_s(Z, W_s)] Y_s \mid Z, W) \right) \\
 &= \sum_{s=1}^{\bar{S}} \left( \int_{-\infty}^{I_{s-1}(Z, W_{s-1})} E(Y_s \mid V + Q_{s-1} = t) dF_{V+Q_{s-1}}(t) \right. \\
 &\quad \left. - \int_{-\infty}^{I_s(Z, W_s)} E(Y_s \mid V + Q_s = t) dF_{V+Q_s}(t) \right) \\
 &= \sum_{s=1}^{\bar{S}} \left( \int_0^{\pi_{s-1}(Z, W_{s-1})} E(Y_s \mid U_{D,s-1} = t) dt - \int_0^{\pi_s(Z, W_s)} E(Y_s \mid U_{D,s} = t) dt \right).
 \end{aligned}$$

- We thus have the index sufficiency restriction that  $E(Y | Z, W) = E(Y | \pi(Z, W))$ , and in the general case  $\frac{\partial}{\partial \pi_s} E(Y | \pi(Z, W) = \pi) = E(Y_{s+1} - Y_s | U_{D,s} = \pi_s)$ .

- We thus have the index sufficiency restriction that  $E(Y | Z, W) = E(Y | \pi(Z, W))$ , and in the general case  $\frac{\partial}{\partial \pi_s} E(Y | \pi(Z, W) = \pi) = E(Y_{s+1} - Y_s | U_{D,s} = \pi_s)$ .
- Also, notice that we have the restriction that  $\frac{\partial^2}{\partial \pi_s \partial \pi_{s'}} E(Y | \pi(Z, W) = \pi) = 0$  if  $|s - s'| > 1$ .

- We thus have the index sufficiency restriction that  $E(Y | Z, W) = E(Y | \pi(Z, W))$ , and in the general case  $\frac{\partial}{\partial \pi_s} E(Y | \pi(Z, W) = \pi) = E(Y_{s+1} - Y_s | U_{D,s} = \pi_s)$ .
- Also, notice that we have the restriction that  $\frac{\partial^2}{\partial \pi_s \partial \pi_{s'}} E(Y | \pi(Z, W) = \pi) = 0$  if  $|s - s'| > 1$ .
- Under full independence between  $U_s$  and  $V + Q_s$ ,  $s = 1, \dots, \bar{S}$ , we can test full independence for the more general choice model by testing for linearity of  $E(Y | \pi(Z, W) = \pi)$  in  $\pi$ .

- Define

$$\Delta_{s+1,s}^{\text{MTE}}(x, u) = E(Y_{s+1} - Y_s \mid X = x, U_{D,s} = u),$$

so that our result above can be rewritten as

$$\frac{\partial}{\partial \pi_s} E(Y \mid \pi(Z, W) = \pi) = \Delta_{s+1,s}^{\text{MTE}}(x, \pi_s).$$



- Define

$$\Delta_{s+1,s}^{\text{MTE}}(x, u) = E(Y_{s+1} - Y_s \mid X = x, U_{D,s} = u),$$

so that our result above can be rewritten as

$$\frac{\partial}{\partial \pi_s} E(Y \mid \pi(Z, W) = \pi) = \Delta_{s+1,s}^{\text{MTE}}(x, \pi_s).$$

- Since  $\pi_s(Z, W_s)$  can be nonparametrically identified from

$$\pi_s(Z, W_s) = \Pr \left( \sum_{j=s+1}^{\bar{s}} D_j = 1 \mid Z, W_s \right),$$

we have identification of MTE for all evaluation points within the appropriate support.

- The policy relevant treatment effect is defined analogously.

- The policy relevant treatment effect is defined analogously.
- $H_s^p$  is defined as the cumulative distribution function of  $\mu_D(Z) - C_s(W_s)$ .

- We have that

$$\begin{aligned}
 E_p(Y_p) &= \\
 &= E_p(E(Y | V, Q, Z, W)) \\
 &= E_p\left(\sum_{s=1}^{\bar{s}} \mathbf{1}[I_{s-1}(Z, W_{s-1}) - Q_{s-1} > V \geq I_s(Z, W_s) - Q_s] E(Y_s | V, Q, Z, W)\right) \\
 &= E_p\left(\sum_{s=1}^{\bar{s}} \mathbf{1}[I_{s-1}(Z, W_{s-1}) - Q_{s-1} > V \geq I_s(Z, W_s) - Q_s] E(Y_s | V, Q)\right) \\
 &= \sum_{s=1}^{\bar{s}} E_p\left(E(Y_s | V, Q) \{H_s^p(V + Q_s) - H_{s-1}^p(V + Q_{s-1})\}\right) \\
 &= \sum_{s=1}^{\bar{s}} \int \left(E(Y_s | V = v, Q = q) \{H_s^p(v + q_s) - H_{s-1}^p(v + q_{s-1})\}\right) dF_{V,Q}(v, q) \\
 &= \sum_{s=1}^{\bar{s}} \left(\int E(Y_s | V + Q_s = t) H_s^p(t) dF_{V+Q_s}(t)\right) \\
 &= - \int E(Y_s | V + Q_{s-1} = t) H_{s-1}^p(t) dF_{V+Q_{s-1}}(t)
 \end{aligned}$$

- $V$ ,  $Q_s$  enter additively, and

$$\begin{aligned}\Delta_{p,p'}^{\text{PRTE}} &= E_{p'}(Y) - E_p(Y) \\ &= \sum_{s=1}^{\bar{S}-1} \int \left( E(Y_{s+1} - Y_s \mid V + Q_s = t) \{ H_s^p(t) - H_s^{p'}(t) \} \right) dF_{V+Q_s}(t).\end{aligned}$$

- Alternatively, we can express this result in terms of MTE,

$$E_p(Y_p) = \sum_{s=1}^{\bar{s}} \left( \int E(Y_s | U_{D,s} = t) \tilde{H}_s^p(t) dt - \int E(Y_s | U_{D,s-1} = t) \tilde{H}_{s-1}^p(t) dt \right)$$

so that

$$\begin{aligned} \Delta_{p,p'}^{\text{PRTE}} &= E_{p'}(Y) - E_p(Y) \\ &= \sum_{s=1}^{\bar{s}-1} \int \left( E(Y_{s+1} - Y_s | U_{D,s} = t) \{ \tilde{H}_s^p(t) - \tilde{H}_s^{p'}(t) \} \right) dt \end{aligned}$$

where  $\tilde{H}_s^p$  is the cumulative distribution function of the random variable  $F_{U_{D,s}}(\mu_D(Z) - C_s(W_s))$ .

## Derivation of PRTE Weights for the Ordered Choice Model

- To derive the  $\omega_{p,p'}$  weights used in expression (7.5), let  $I_s(Z, W_s) = \mu_D(Z) - C_s(W_s)$ , and let  $H_s^p(\cdot)$  denote the cumulative distribution function of  $I_s(Z, W_s)$  under regime  $p$ ,  $H_s^p(t) = \int \mathbf{1}[\mu_D(z) - C_s(w_s) \leq t] dF_{Z,W}^p(z, w)$ .

## Derivation of PRTE Weights for the Ordered Choice Model

- To derive the  $\omega_{p,p'}$  weights used in expression (7.5), let  $l_s(Z, W_s) = \mu_D(Z) - C_s(W_s)$ , and let  $H_s^p(\cdot)$  denote the cumulative distribution function of  $l_s(Z, W_s)$  under regime  $p$ ,  $H_s^p(t) = \int \mathbf{1}[\mu_D(z) - C_s(w_s) \leq t] dF_{Z,W}^p(z, w)$ .
- Because  $C_0(W_0) = -\infty$  and  $C_{\bar{5}}(W_{\bar{5}}) = \infty$ ,  $l_0(Z, W_0) = \infty$  and  $l_{\bar{5}}(Z, W_{\bar{5}}) = -\infty$ ,  $H_0^p(t) = 0$  and  $H_{\bar{5}}^p(t) = 1$  for any policy  $p$  and for all evaluation points.



## Derivation of PRTE Weights for the Ordered Choice Model

- To derive the  $\omega_{p,p'}$  weights used in expression (7.5), let  $I_s(Z, W_s) = \mu_D(Z) - C_s(W_s)$ , and let  $H_s^p(\cdot)$  denote the cumulative distribution function of  $I_s(Z, W_s)$  under regime  $p$ ,  $H_s^p(t) = \int \mathbf{1}[\mu_D(z) - C_s(w_s) \leq t] dF_{Z,W}^p(z, w)$ .
- Because  $C_0(W_0) = -\infty$  and  $C_{\bar{5}}(W_{\bar{5}}) = \infty$ ,  $I_0(Z, W_0) = \infty$  and  $I_{\bar{5}}(Z, W_{\bar{5}}) = -\infty$ ,  $H_0^p(t) = 0$  and  $H_{\bar{5}}^p(t) = 1$  for any policy  $p$  and for all evaluation points.
- Since  $I_{s-1}(Z, W_{s-1})$  is always larger than  $I_s(Z, W_s)$ , we obtain

$$\mathbf{1}[I_s(Z, W_s) \leq V < I_{s-1}(Z, W_{s-1})] = \mathbf{1}[V < I_{s-1}(Z, W_{s-1})] - \mathbf{1}[V \leq I_s(Z, W_s)],$$

so that under assumption (OC-1),

$$E_p(\mathbf{1}[I_s(Z, W_s) \leq V \leq I_{s-1}(Z, W_{s-1})] | V) = H_s^p(V) - H_{s-1}^p(V).$$

- Collecting these results we obtain:

$$\begin{aligned} E_p(Y) &= E_p[E(Y | V, Z, W)] = \\ &= \sum_{s=1}^{\bar{s}} \int \left[ E(Y_s | V = v) \{H_s^p(v) - H_{s-1}^p(v)\} \right] f_V(v) dv. \end{aligned}$$

- Comparing two policies under  $p$  and  $p'$ , the policy relevant treatment effect is  $\Delta_{p,p'}^{\text{PRTE}} = E_{p'}(Y) - E_p(Y) = \sum_{s=1}^{\bar{S}-1} \int E(Y_{s+1} - Y_s | V = v) [H_s^p(v) - H_s^{p'}(v)] f_V(v) dv.$

- Comparing two policies under  $p$  and  $p'$ , the policy relevant treatment effect is  $\Delta_{p,p'}^{\text{PRTE}} = E_{p'}(Y) - E_p(Y) = \sum_{s=1}^{\bar{S}-1} \int E(Y_{s+1} - Y_s | V = v) [H_s^p(v) - H_s^{p'}(v)] f_V(v) dv$ .
- Alternatively, we can express this in terms of  $\Delta^{\text{MTE}}$ :  
 $\Delta_{p,p'}^{\text{PRTE}} = \sum_{s=1}^{\bar{S}-1} \int \Delta_{s,s+1}^{\text{MTE}}(u) [\tilde{H}_s^p(u) - \tilde{H}_s^{p'}(u)] du$  where  $\tilde{H}_s^p(t)$  is the cumulative distribution function of  $F_V(\mu_D(Z) - C_s(W_s))$  under policy  $p$ ,  
 $\tilde{H}_s^p(t) = \int \mathbf{1}[F_V(\mu_D(z) - C_s(w_s)) \leq t] dF_{Z,W_s}^p(z, w_s)$ .

## Derivation of the Weights for IV in the Ordered Choice Model

- We first derive  $\text{Cov}(J(Z, W), Y)$ .

## Derivation of the Weights for IV in the Ordered Choice Model

- We first derive  $\text{Cov}(J(Z, W), Y)$ .
- Its derivation is typical of the other terms needed to form (47) in the text.

- Defining  $\tilde{J}(Z, W) = J(Z, W) - E(J(Z, W))$ , we obtain, since  $\text{Cov}(J(Z, W), Y) = E(\tilde{J}(Z, W) Y)$ ,

$$\begin{aligned} E(\tilde{J}(Z, W) Y) &= E \left[ \tilde{J}(Z, W) \sum_{s=1}^{\bar{S}} \mathbf{1}[I_s(Z, W_s) \leq V < I_{s-1}(Z, W_{s-1})] E(Y_s | V, Z, W) \right] \\ &= \sum_{s=1}^{\bar{S}} E \left[ \tilde{J}(Z, W) \mathbf{1}[I_s(Z, W_s) \leq V < I_{s-1}(Z, W_{s-1})] E(Y_s | V) \right] \end{aligned}$$

where the first equality comes from the definition of  $Y$  and the law of iterated expectations, and the second equality follows from linearity of expectations and independence assumption (OC-1).

- Let  $H_s(\cdot)$  equal  $H_s^p(\cdot)$  for  $p$  equal to the policy that characterizes the observed data, i.e.,  $H_s(\cdot)$  is the cumulative distribution function of  $I_s(Z, W_s)$ ,

$$H_s^p(t) = \Pr(I_s(Z, W_s) \leq t) = \Pr(\mu_D(Z) - C_s(W_s) \leq t).$$



- Using the law of iterated expectations, we obtain

$$\begin{aligned}
 E(\tilde{J}(Z, W)Y) &= \\
 &= \sum_{s=1}^{\bar{S}} E \left[ E \left( \tilde{J}(Z, W) \left\{ \mathbf{1}[V < I_{s-1}(Z, W_{s-1})] - \mathbf{1}[V \leq I_s(Z, W_s)] \right\} \mid V \right) E(Y_s \mid V) \right] \\
 &= \sum_{s=1}^{\bar{S}} \int [E(Y_s \mid V = v) \{K_{s-1}(v) - K_s(v)\}] f_V(v) dv \\
 &= \sum_{s=1}^{\bar{S}-1} \int [E(Y_{s+1} - Y_s \mid V = v) K_s(v)] f_V(v) dv
 \end{aligned}$$

where  $K_s(v) = E(\tilde{J}(Z, W) \mid I_s(Z, W_s) > v) (1 - H_s(v))$  and we use the fact that  $K_{\bar{S}}(v) = K_0(v) = 0$ .

- Now consider the denominator of the IV estimand,

$$\begin{aligned}
 E(\tilde{S}J(Z, W)) &= \\
 &= E \left[ \tilde{J}(Z, W) \sum_{s=1}^{\bar{S}} s \mathbf{1}[l_s(Z, W_s) \leq V < l_{s-1}(Z, W_{s-1})] \right] \\
 &= \sum_{s=1}^{\bar{S}} s E \left[ \tilde{J}(Z, W) \mathbf{1}[l_s(Z, W_s) \leq V < l_{s-1}(Z, W_{s-1})] \right] \\
 &= \sum_{s=1}^{\bar{S}} s E_V \left[ E \left( \tilde{J}(Z, W) \left\{ \mathbf{1}[V < l_{s-1}(Z, W_{s-1})] - \mathbf{1}[V \leq l_s(Z, W_s)] \right\} \middle| V \right) \right] \\
 &= \sum_{s=1}^{\bar{S}} s \int [K_{s-1}(v) - K_s(v)] f_V(v) dv = \sum_{s=1}^{\bar{S}-1} \int K_s(v) f_V(v) dv.
 \end{aligned}$$

- Collecting results, we obtain an expression for the IV estimand (47):

$$\frac{\text{Cov}(J, Y)}{\text{Cov}(J, S)} = \sum_{s=1}^{\bar{S}-1} \int E(Y_{s+1} - Y_s \mid V = v) \omega(s, v) f_V(v) dv$$

where

$$\omega(s, v) = \frac{K_s(v)}{\sum_{s=1}^{\bar{S}} s \int [K_{s-1}(v) - K_s(v)] f_V(v) dv} = \frac{K_s(v)}{\sum_{s=1}^{\bar{S}-1} \int K_s(v) f_V(v) dv}$$

and clearly

$$\sum_{s=1}^{\bar{S}-1} \int \omega(s, v) f_V(v) dv = 1, \quad \omega(0, v) = 0, \quad \text{and} \quad \omega(\bar{S}, v) = 0.$$

## Proof of Theorem 6

- We now prove Theorem 6.

## Proof

- The basic idea is that we can bring the model back to a two choice set up of  $j$  versus the “next best” option.

## Proof

- The basic idea is that we can bring the model back to a two choice set up of  $j$  versus the “next best” option.
- We prove the result for the second assertion, that  $\Delta_j^{\text{LIV}}(x, z)$  recovers the marginal treatment effect parameter.

## Proof

- The basic idea is that we can bring the model back to a two choice set up of  $j$  versus the “next best” option.
- We prove the result for the second assertion, that  $\Delta_j^{\text{LIV}}(x, z)$  recovers the marginal treatment effect parameter.
- The first assertion, that  $\Delta_j^{\text{Wald}}(x, z^{[-j]}, z^{[j]}, \tilde{z}^{[j]})$  recovers a LATE parameter, follows from a trivial modification to the same proof strategy.

## Proof

- The basic idea is that we can bring the model back to a two choice set up of  $j$  versus the “next best” option.
- We prove the result for the second assertion, that  $\Delta_j^{\text{LIV}}(x, z)$  recovers the marginal treatment effect parameter.
- The first assertion, that  $\Delta_j^{\text{Wald}}(x, z^{[-j]}, z^{[j]}, \tilde{z}^{[j]})$  recovers a LATE parameter, follows from a trivial modification to the same proof strategy.
- Recall that  $R_{\mathcal{J} \setminus j}(z) = \max_{i \in \mathcal{J} \setminus j} \{R_i(z)\}$  and that  $I_{\mathcal{J} \setminus j} = \operatorname{argmax}_{i \in \mathcal{J} \setminus j} (R_i(Z))$ .



## Proof

- The basic idea is that we can bring the model back to a two choice set up of  $j$  versus the “next best” option.
- We prove the result for the second assertion, that  $\Delta_j^{\text{LIV}}(x, z)$  recovers the marginal treatment effect parameter.
- The first assertion, that  $\Delta_j^{\text{Wald}}(x, z^{[-j]}, z^{[j]}, \tilde{z}^{[j]})$  recovers a LATE parameter, follows from a trivial modification to the same proof strategy.
- Recall that  $R_{\mathcal{J} \setminus j}(z) = \max_{i \in \mathcal{J} \setminus j} \{R_i(z)\}$  and that  $I_{\mathcal{J} \setminus j} = \operatorname{argmax}_{i \in \mathcal{J} \setminus j} (R_i(Z))$ .
- We may write  $Y = Y_{I_{\mathcal{J} \setminus j}} + D_{\mathcal{J}j}(Y_j - Y_{I_{\mathcal{J} \setminus j}})$ .

- We have

$$\begin{aligned}\Pr(D_{\mathcal{J},j} = 1 \mid X = x, Z = z) &= \Pr(R_j(z_j) \geq R_{\mathcal{J} \setminus j}(z) \mid X = x, Z = z) \\ &= \Pr(\vartheta_j(z_j) \geq R_{\mathcal{J} \setminus j}(z) + V_j \mid X = x, Z = z).\end{aligned}$$

- We have

$$\begin{aligned}\Pr(D_{\mathcal{J},j} = 1 \mid X = x, Z = z) &= \Pr(R_j(z_j) \geq R_{\mathcal{J} \setminus j}(z) \mid X = x, Z = z) \\ &= \Pr(\vartheta_j(z_j) \geq R_{\mathcal{J} \setminus j}(z) + V_j \mid X = x, Z = z).\end{aligned}$$

- Using independence assumption (B-1),  $R_{\mathcal{J} \setminus j}(z) - V_j$  is independent of  $Z$  conditional on  $X$ , so that

$$\Pr(D_{\mathcal{J},j} = 1 \mid X = x, Z = z) = \Pr(\vartheta_j(z_j) \geq R_{\mathcal{J} \setminus j}(z) + V_j \mid X = x).$$

- We have

$$\begin{aligned}\Pr(D_{\mathcal{J},j} = 1 \mid X = x, Z = z) &= \Pr(R_j(z_j) \geq R_{\mathcal{J}\setminus j}(z) \mid X = x, Z = z) \\ &= \Pr(\vartheta_j(z_j) \geq R_{\mathcal{J}\setminus j}(z) + V_j \mid X = x, Z = z).\end{aligned}$$

- Using independence assumption (B-1),  $R_{\mathcal{J}\setminus j}(z) - V_j$  is independent of  $Z$  conditional on  $X$ , so that

$$\Pr(D_{\mathcal{J},j} = 1 \mid X = x, Z = z) = \Pr(\vartheta_j(z_j) \geq R_{\mathcal{J}\setminus j}(z) + V_j \mid X = x).$$

- $\vartheta_k(\cdot)$  does not depend on  $z^{[j]}$  for  $k \neq j$  by assumption (B-2b), and thus  $R_{\mathcal{J}\setminus j}(z)$  does not depend on  $z^{[j]}$ , and we will therefore (with an abuse of notation) write  $R_{\mathcal{J}\setminus j}(z^{[-j]})$  for  $R_{\mathcal{J}\setminus j}(z)$ .

- We have

$$\begin{aligned}\Pr(D_{\mathcal{J},j} = 1 \mid X = x, Z = z) &= \Pr(R_j(z_j) \geq R_{\mathcal{J}\setminus j}(z) \mid X = x, Z = z) \\ &= \Pr(\vartheta_j(z_j) \geq R_{\mathcal{J}\setminus j}(z) + V_j \mid X = x, Z = z).\end{aligned}$$

- Using independence assumption (B-1),  $R_{\mathcal{J}\setminus j}(z) - V_j$  is independent of  $Z$  conditional on  $X$ , so that

$$\Pr(D_{\mathcal{J},j} = 1 \mid X = x, Z = z) = \Pr(\vartheta_j(z_j) \geq R_{\mathcal{J}\setminus j}(z) + V_j \mid X = x).$$

- $\vartheta_k(\cdot)$  does not depend on  $z^{[j]}$  for  $k \neq j$  by assumption (B-2b), and thus  $R_{\mathcal{J}\setminus j}(z)$  does not depend on  $z^{[j]}$ , and we will therefore (with an abuse of notation) write  $R_{\mathcal{J}\setminus j}(z^{[-j]})$  for  $R_{\mathcal{J}\setminus j}(z)$ .
- Write  $F(\cdot; X = x, Z^{[-j]} = z^{[-j]})$  for the distribution function of  $R_{\mathcal{J}\setminus j}(z^{[-j]}) + V_j$  conditional on  $X = x$ .

- Then

$$\Pr(D_{\mathcal{J},j} = 1 \mid X = x, Z = z) = F(\vartheta_j(z_j); x, z^{[-j]}),$$

and

$$\frac{\partial}{\partial z^{[j]}} \Pr(D_{\mathcal{J},j} = 1 \mid X = x, Z = z) = \left[ \frac{\partial}{\partial z^{[j]}} \vartheta_j(z_j) \right] f_{X|Z^{[-j]}}(\vartheta_j(z_j); X = x, Z^{[-j]} = z^{[-j]}),$$

where  $f_{X|Z^{[-j]}}(\cdot; X = x, Z^{[-j]} = z^{[-j]})$  is the density of  $R_{\mathcal{J}\setminus j}(z^{[-j]}) - V_j$  conditional on  $X = x$ .

- Then

$$\Pr(D_{\mathcal{J},j} = 1 \mid X = x, Z = z) = F(\vartheta_j(z_j); x, z^{[-j]}),$$

and

$$\frac{\partial}{\partial z^{[j]}} \Pr(D_{\mathcal{J},j} = 1 \mid X = x, Z = z) = \left[ \frac{\partial}{\partial z^{[j]}} \vartheta_j(z_j) \right] f_{X|Z^{[-j]}}(\vartheta_j(z_j); X = x, Z^{[-j]} = z^{[-j]}),$$

where  $f_{X|Z^{[-j]}}(\cdot; X = x, Z^{[-j]} = z^{[-j]})$  is the density of  $R_{\mathcal{J}\setminus j}(z^{[-j]}) - V_j$  conditional on  $X = x$ .

- Consider

$$\begin{aligned} E(Y \mid X = x, Z = z) &= E(Y_{I_{\mathcal{J}\setminus j}} \mid X = x, Z = z) \\ &\quad + E(D_{\mathcal{J},j}(Y_j - Y_{I_{\mathcal{J}\setminus j}}) \mid X = x, Z = z). \end{aligned}$$

- As a consequence of (B-1), (B-3)–(B-5), and (B-2b), we have that  $E\left(Y_{I_{j \setminus j}} \mid X = x, Z = z\right)$  does not depend on  $z^{[j]}$ .



- As a consequence of (B-1), (B-3)–(B-5), and (B-2b), we have that  $E\left(Y_{I_{\mathcal{J}\setminus j}} \mid X = x, Z = z\right)$  does not depend on  $z^{[j]}$ .
- Using the assumptions and the law of iterated expectations, we may write

$$\begin{aligned}
 E\left(D_{\mathcal{J},j}(Y_j - Y_{I_{\mathcal{J}\setminus j}}) \mid X = x, Z = z\right) &= \int_{-\infty}^{\vartheta_j(z)} E(Y_j - Y_{I_{\mathcal{J}\setminus j}} \mid X = x, Z = z, R_{\mathcal{J}\setminus j}(z^{[-j]} \\
 &\quad + V_j = t) f_{X|Z^{[-j]}}(t; X = x, Z^{[-j]} = z^{[-j]}) dt \\
 &= \int_{-\infty}^{\vartheta_j(z)} E(Y_j - Y_{I_{\mathcal{J}\setminus j}} \mid X = x, Z^{[-j]} = z^{[-j]}, R_{\mathcal{J}\setminus j}(z^{[-j]} \\
 &\quad - V_j = t) f_{X|Z^{[-j]}}(t; X = x, Z^{[-j]} = z^{[-j]}) dt.
 \end{aligned}$$

- Thus,

$$\begin{aligned} & \frac{\partial}{\partial z^{[j]}} E(Y \mid X = x, Z = z) \\ &= E\left(Y_j - Y_{I_{\mathcal{J} \setminus j}} \mid X = x, Z^{[-j]} = z^{[-j]}, R_j(z) = R_{\mathcal{J} \setminus j}(z)\right) \\ & \quad \times \left[ \frac{\partial}{\partial z^{[j]}} \vartheta_j(z_j) \right] f_{X|Z^{[-j]}}(\vartheta_j(z_j) \mid X = x, Z^{[-j]} = z^{[-j]}). \end{aligned}$$

- Thus,

$$\begin{aligned} & \frac{\partial}{\partial z^j} E(Y | X = x, Z = z) \\ &= E\left(Y_j - Y_{I_{\mathcal{J}^c}} | X = x, Z^{[-j]} = z^{[-j]}, R_j(z) = R_{\mathcal{J}^c}(z)\right) \\ & \quad \times \left[ \frac{\partial}{\partial z^j} \vartheta_j(z_j) \right] f_{X|Z^{[-j]}}(\vartheta_j(z_j) | X = x, Z^{[-j]} = z^{[-j]}). \end{aligned}$$

- Combining results, we have

$$\begin{aligned} & \frac{\partial}{\partial z^j} E(Y | X = x, Z = z) / \frac{\partial}{\partial z^j} \Pr(D_{\mathcal{J}^c} = 1 | X = x, Z = z) \\ &= E\left(Y_j - Y_{I_{\mathcal{J}^c}} | X = x, Z^{[-j]} = z^{[-j]}, R_j(z) = R_{\mathcal{J}^c}(z)\right). \end{aligned}$$

- Finally, noting that

$$\begin{aligned} E\left(Y_j - Y_{I_{\mathcal{J}^c}} \mid X = x, Z^{[-j]} = z^{[-j]}, R_j(z) = R_{\mathcal{J}^c}(z)\right) \\ = E\left(Y_j - Y_{I_{\mathcal{J}^c}} \mid X = x, Z = z, R_j(z) = R_{\mathcal{J}^c}(z)\right) \end{aligned}$$

provides the stated result.

- Finally, noting that

$$\begin{aligned} E\left(Y_j - Y_{I_{\mathcal{J}^c}} \mid X = x, Z^{[-j]} = z^{[-j]}, R_j(z) = R_{\mathcal{J}^c}(z)\right) \\ = E\left(Y_j - Y_{I_{\mathcal{J}^c}} \mid X = x, Z = z, R_j(z) = R_{\mathcal{J}^c}(z)\right) \end{aligned}$$

provides the stated result.

- The proof for the LATE result follows from the parallel argument using discrete changes in the instrument.

## Flat MTE within a General Nonseparable Matching Framework

- The result in the text that conditional mean independence of  $Y_0$  and  $Y_1$  in terms of  $D$  given  $X$  implies a flat MTE holds in a more general nonseparable model.

## Flat MTE within a General Nonseparable Matching Framework

- The result in the text that conditional mean independence of  $Y_0$  and  $Y_1$  in terms of  $D$  given  $X$  implies a flat MTE holds in a more general nonseparable model.
- We establish this claim and also establish some additional restrictions implied by an IV assumption.

- Assume a nonseparable selection model,  
 $D = \mathbf{1}[\mu_D(X, Z, V) \geq 0]$ , with  $Z$  independent of  $(Y_0, Y_1, V)$  conditional on  $X$ .



- Assume a nonseparable selection model,  
 $D = \mathbf{1}[\mu_D(X, Z, V) \geq 0]$ , with  $Z$  independent of  $(Y_0, Y_1, V)$  conditional on  $X$ .
- Let  $\Omega(x, z) = \{v : \mu_D(x, z, v) \geq 0\}$ .

- Assume a nonseparable selection model,  
 $D = \mathbf{1}[\mu_D(X, Z, V) \geq 0]$ , with  $Z$  independent of  $(Y_0, Y_1, V)$  conditional on  $X$ .
- Let  $\Omega(x, z) = \{v : \mu_D(x, z, v) \geq 0\}$ .
- Let  $\Omega(x, z)^c$  denote the complement of  $\Omega(x, z)$ .

- Consider the mean independence assumption

$$(M-3) \quad E(Y_1|X, D) = E(Y_1|X), \quad E(Y_0|X, D) = E(Y_0|X).$$

(M-3) implies that for  $\Delta = Y_1 - Y_0$

$$E(\Delta|X = x, V \in \Omega(X, Z)) = E(\Delta|X = x, V \in \Omega(X, Z)^c),$$

where  $c$  here denotes “complement”. Thus,

$$\begin{aligned} E_{Z|X}(E(\Delta^{\text{MTE}}(x, V)|X = x, V \in \Omega(x, Z))|X = x) \\ = E_{Z|X}(E(\Delta^{\text{MTE}}(x, V)|X = x, V \in \Omega(x, Z)^c) | X = x) \end{aligned}$$

for all  $x$  in the support of  $X$ .

- Consider the mean independence assumption

$$(M-3) \quad E(Y_1|X, D) = E(Y_1|X), \quad E(Y_0|X, D) = E(Y_0|X).$$

(M-3) implies that for  $\Delta = Y_1 - Y_0$

$$E(\Delta|X = x, V \in \Omega(X, Z)) = E(\Delta|X = x, V \in \Omega(X, Z)^c),$$

where  $c$  here denotes “complement”. Thus,

$$\begin{aligned} E_{Z|X}(E(\Delta^{\text{MTE}}(x, V)|X = x, V \in \Omega(x, Z))|X = x) \\ = E_{Z|X}(E(\Delta^{\text{MTE}}(x, V)|X = x, V \in \Omega(x, Z)^c) | X = x) \end{aligned}$$

for all  $x$  in the support of  $X$ .

- (We assume  $0 < \Pr(D = 1|X) < 1$  .)

- Consider the mean independence assumption

$$(M-3) \quad E(Y_1|X, D) = E(Y_1|X), \quad E(Y_0|X, D) = E(Y_0|X).$$

(M-3) implies that for  $\Delta = Y_1 - Y_0$

$$E(\Delta|X = x, V \in \Omega(X, Z)) = E(\Delta|X = x, V \in \Omega(X, Z)^c),$$

where  $c$  here denotes “complement”. Thus,

$$\begin{aligned} E_{Z|X}(E(\Delta^{\text{MTE}}(x, V)|X = x, V \in \Omega(x, Z))|X = x) \\ = E_{Z|X}(E(\Delta^{\text{MTE}}(x, V)|X = x, V \in \Omega(x, Z)^c) | X = x) \end{aligned}$$

for all  $x$  in the support of  $X$ .

- (We assume  $0 < \Pr(D = 1|X) < 1$  .)
- This establishes that the MTE is flat.

Now suppose that (M-3) holds, but suppose that there is an instrument  $Z$  such that

$$(M-3)' \quad E(Y_1|X, Z, D) \neq E(Y_1|X), \quad E(Y_0|X, Z, D) \neq E(Y_0|X).$$

- (Note:  $E(Y_j|X, Z) = E(Y_j|X)$  by assumption).

- (Note:  $E(Y_j|X, Z) = E(Y_j|X)$  by assumption).
- In this case, (M-3) implies that

$$\begin{aligned} E_{Z|X}(E(\Delta^{\text{MTE}}(X, V)|X = x, V \in \Omega(x, Z)) | X = x) \\ = E_{Z|X}(E(\Delta^{\text{MTE}}(X, V)|X = x, V \in (\Omega(x, Z))^c | X = x), \end{aligned}$$

but (M-3)' implies that there exists  $z$  in the support of  $Z$  conditional on  $X$  such that

$$E(\Delta^{\text{MTE}}(X, V)|X = x, V \in \Omega(x, z)) \neq E(\Delta^{\text{MTE}}(X, V)|X = x)$$

and

$$E(\Delta^{\text{MTE}}(X, V)|X = x, V \in \Omega(x, z)^c) \neq E(\Delta^{\text{MTE}}(X, V)|X = x)$$

so that  $\Delta^{\text{MTE}}(X, V)$  is not constant in  $V$ .



- Note that, if  $E(Y_1|X, Z = z, D = 1) \neq E(Y_1|X, Z = z', D = 1)$  for any  $z, z'$  evaluation points in the support of  $Z$  conditional on  $X$ , then  $E(Y_1|X, Z, D) \neq E(Y_1|X)$ .

- Note that, if  $E(Y_1|X, Z = z, D = 1) \neq E(Y_1|X, Z = z', D = 1)$  for any  $z, z'$  evaluation points in the support of  $Z$  conditional on  $X$ , then  $E(Y_1|X, Z, D) \neq E(Y_1|X)$ .
- Thus, (M-3)' is testable, given the maintained assumption that  $Z$  is a proper exclusion restriction.

- Note that, if  $E(Y_1|X, Z = z, D = 1) \neq E(Y_1|X, Z = z', D = 1)$  for any  $z, z'$  evaluation points in the support of  $Z$  conditional on  $X$ , then  $E(Y_1|X, Z, D) \neq E(Y_1|X)$ .
- Thus, (M-3)' is testable, given the maintained assumption that  $Z$  is a proper exclusion restriction.
- Note that (M-3)' implies (M-3), so it is a stronger condition.

Now assume

$$(M-1)' \quad E(Y_1|X, Z, D) = E(Y_1|X), \quad E(Y_0|X, Z, D) = E(Y_0|X).$$

- In this case, we get a stronger restriction on MTE than is produced from (M-3).

- In this case, we get a stronger restriction on MTE than is produced from (M-3).
- We obtain

$$E(\Delta^{\text{MTE}}(X, V)|X = x, V \in \Omega(x, z)) = E(\Delta^{\text{MTE}}(X, V)|X = x)$$

and

$$E(\Delta^{\text{MTE}}(X, V)|X = x, V \in \Omega(x, z)^c) = E(\Delta^{\text{MTE}}(X, V)|X = x)$$

for all  $(x, z)$  in the proper support.

- In this case, we get a stronger restriction on MTE than is produced from (M-3).
- We obtain

$$E(\Delta^{\text{MTE}}(X, V)|X = x, V \in \Omega(x, z)) = E(\Delta^{\text{MTE}}(X, V)|X = x)$$

and

$$E(\Delta^{\text{MTE}}(X, V)|X = x, V \in \Omega(x, z)^c) = E(\Delta^{\text{MTE}}(X, V)|X = x)$$

for all  $(x, z)$  in the proper support.

- Again, the MTE is not flat.

## The Relationship Between Exclusion Conditions in IV and Exclusion Conditions in Matching

- We now investigate the relationship between IV and matching identification conditions.



## The Relationship Between Exclusion Conditions in IV and Exclusion Conditions in Matching

- We now investigate the relationship between IV and matching identification conditions.
- They are very distinct.

## The Relationship Between Exclusion Conditions in IV and Exclusion Conditions in Matching

- We now investigate the relationship between IV and matching identification conditions.
- They are very distinct.
- We analyze mean treatment parameters.

## The Relationship Between Exclusion Conditions in IV and Exclusion Conditions in Matching

- We now investigate the relationship between IV and matching identification conditions.
- They are very distinct.
- We analyze mean treatment parameters.
- We define  $(U_0, U_1)$  by  $U_0 = Y_0 - E(Y_0|X)$  and  $U_1 = Y_1 - E(Y_1|X)$ .

## The Relationship Between Exclusion Conditions in IV and Exclusion Conditions in Matching

- We now investigate the relationship between IV and matching identification conditions.
- They are very distinct.
- We analyze mean treatment parameters.
- We define  $(U_0, U_1)$  by  $U_0 = Y_0 - E(Y_0|X)$  and  $U_1 = Y_1 - E(Y_1|X)$ .
- We consider standard IV as a form of matching where matching does not hold conditional on  $X$  but does hold conditional on  $(X, Z)$ , where  $Z$  is the instrument.

- Consider the following two matching conditions based on an exclusion restriction  $Z$ :

(M-4)  $(U_0, U_1)$  are mean independent of  $D$  conditional on  $(X, Z)$ . ( $E(U_0 | X, Z, D) = E(U_0 | X, Z)$  and  $E(U_1 | X, Z, D) = E(U_1 | X, Z)$ .)

(M-5)  $(U_0, U_1)$  are not mean independent of  $D$  conditional on  $X$ . ( $E(U_0 | X, D) \neq E(U_0 | X)$  and  $E(U_1 | X, D) \neq E(U_1 | X)$ .)

- (M-4) says that the matching conditions hold conditional on  $(X, Z)$ .

- (M-4) says that the matching conditions hold conditional on  $(X, Z)$ .
- However, (M-5) says that the matching conditions do not hold if one only conditions on  $X$ .

- (M-4) says that the matching conditions hold conditional on  $(X, Z)$ .
- However, (M-5) says that the matching conditions do not hold if one only conditions on  $X$ .
- By the definitions of  $U_0, U_1$ , these conditions are equivalent to stating that  $Y_0, Y_1$  are mean independent of  $D$  conditional on  $(X, Z)$  but not mean independent of  $D$  conditional on  $X$ .



- (M-4) says that the matching conditions hold conditional on  $(X, Z)$ .
- However, (M-5) says that the matching conditions do not hold if one only conditions on  $X$ .
- By the definitions of  $U_0, U_1$ , these conditions are equivalent to stating that  $Y_0, Y_1$  are mean independent of  $D$  conditional on  $(X, Z)$  but not mean independent of  $D$  conditional on  $X$ .
- These look like instrumental variable conditions.

- (M-4) says that the matching conditions hold conditional on  $(X, Z)$ .
- However, (M-5) says that the matching conditions do not hold if one only conditions on  $X$ .
- By the definitions of  $U_0, U_1$ , these conditions are equivalent to stating that  $Y_0, Y_1$  are mean independent of  $D$  conditional on  $(X, Z)$  but not mean independent of  $D$  conditional on  $X$ .
- These look like instrumental variable conditions.
- We now consider whether these assumptions are compatible with standard IV conditions as used by ?? and ? to use IV to identify treatment parameters when responses are heterogenous (the model of essential heterogeneity).

For ATE, they show that standard IV identifies ATE if:

(ATE-1)  $U_0$  is mean independent of  $Z$  conditional on  $X$ .

(ATE-2)  $D(U_1 - U_0)$  is mean independent of  $Z$  conditional on  $X$ .

They show that standard IV identifies TT if:

(TT-1)  $U_0$  is mean independent of  $Z$  conditional on  $X$ .

(TT-2)  $U_1 - U_0$  is mean independent of  $Z$  conditional on  $D = 1$  and on  $X$ .

The conventional assumption in means is that:

(IV-1)'  $(U_0, U_1)$  are mean independent of  $Z$  conditional on  $X$ .

(IV-2) Rank condition (IV-2) is still required:

$\Pr(D = 1 | Z, X)$  is a nondegenerate function of  $Z$ .

- Condition (IV-1)' is a commonly invoked instrumental variable condition, even though ? and ? show it is neither necessary nor sufficient to identify ATE or TT by linear IV.

- Condition (IV-1)' is a commonly invoked instrumental variable condition, even though ? and ? show it is neither necessary nor sufficient to identify ATE or TT by linear IV.
- In Slide 152, we used the stronger condition (IV-1):  $(U_0, U_1) \perp\!\!\!\perp Z|X$  along with the rank conditions.

- Condition (IV-1)' is a commonly invoked instrumental variable condition, even though ? and ? show it is neither necessary nor sufficient to identify ATE or TT by linear IV.
- In Slide 152, we used the stronger condition (IV-1):  $(U_0, U_1) \perp\!\!\!\perp Z|X$  along with the rank conditions.
- Clearly, (IV-1) implies (IV-1)'.



- We now show that assumptions (M-4) and (M-5) are inconsistent with any of the sets of IV assumptions.

- We now show that assumptions (M-4) and (M-5) are inconsistent with any of the sets of IV assumptions.
- In particular, we show that assuming (M-4) and that  $U_0$  is mean independent of  $Z$  conditional on  $X$  jointly imply that  $U_0$  is mean independent of  $D$  conditional on  $X$ .

- We now show that assumptions (M-4) and (M-5) are inconsistent with any of the sets of IV assumptions.
- In particular, we show that assuming (M-4) and that  $U_0$  is mean independent of  $Z$  conditional on  $X$  jointly imply that  $U_0$  is mean independent of  $D$  conditional on  $X$ .
- If (M-4) and (M-5) hold, then  $Z$  cannot satisfy condition (IV-1)' (or stronger condition (IV-1)), (ATE-1) or (TT-1).

- We now show that assumptions (M-4) and (M-5) are inconsistent with any of the sets of IV assumptions.
- In particular, we show that assuming (M-4) and that  $U_0$  is mean independent of  $Z$  conditional on  $X$  jointly imply that  $U_0$  is mean independent of  $D$  conditional on  $X$ .
- If (M-4) and (M-5) hold, then  $Z$  cannot satisfy condition (IV-1)' (or stronger condition (IV-1)), (ATE-1) or (TT-1).
- Thus matching based on an exclusion restriction and IV are distinct conditions.

We show this by establishing a series of claims:

### Claim

Condition (M-4) and (IV-1)' jointly imply  $U_0$  is mean independent of  $D$  conditional on  $X$ . Thus, (M-4) and [(IV-1)' or (ATE-1) or (TT-1)] jointly imply that (M-5) cannot hold.

## Proof.

Assume (M-4) and (IV-1)'. We have:

$$\begin{aligned} E(U_0|D, X, Z) &= E(U_0|X, Z) \\ &= E(U_0|X) \end{aligned}$$

where the first equality follows from (M-4) and the second equality follows from (IV-1)'. Thus,

$$\begin{aligned} E(U_0|D, X) &= E_Z[E(U_0|D, X, Z)|D, X] \\ &= E_Z[E(U_0|X)|D, X] \\ &= E(U_0|X) \end{aligned}$$



- Thus (M-4) and (M-5) are inconsistent with any of the sets of IV assumptions that we have considered.

- Thus (M-4) and (M-5) are inconsistent with any of the sets of IV assumptions that we have considered.
- However, this analysis raises the question of whether it is still possible to invoke (M-5) and the assumption that  $U_1$  is not mean independent of  $D$  conditional on  $X$ .



- Thus (M-4) and (M-5) are inconsistent with any of the sets of IV assumptions that we have considered.
- However, this analysis raises the question of whether it is still possible to invoke (M-5) and the assumption that  $U_1$  is not mean independent of  $D$  conditional on  $X$ .
- The following results show that it is not possible.

## Claim

(M-4) and (IV-1)' imply  $U_1$  is mean independent of  $D$  conditional on  $X$ .

Proof.

Follows with trivial modification from the proof to Claim 1.

- A similar claim can be shown for (TT-1) and (TT-2).

## Claim

(M-4) and (TT-1), (TT-2) imply  $U_1$  is mean independent of  $D$  conditional on  $X$ .

## Proof.

Assume (M-4) and (TT-1), (TT-2). We have:

(N-1)

$$E(U_0|X, Z, D) = E(U_0|X, Z) = E(U_0|X)$$

where the first equality follows from (M-4) and the second equality follows from (TT-1). Using the result from the proof of Claim 1, we obtain

(N-2)

$$E(U_0|X, Z, D) = E(U_0|X, D)$$



## Proof.

By (TT-2), we have

$$\begin{aligned} E(U_1|X, Z, D = 1) - E(U_1|X, D = 1) \\ = E(U_0|X, Z, D = 1) - E(U_0|X, D = 1). \end{aligned}$$

By equation (N-2), the right hand side of the preceding expression is zero, and we thus have

(N-3)

$$E(U_1|X, Z, D = 1) = E(U_1|X, D = 1).$$

By (M-4), we have

(N-4)

$$E(U_1|X, Z, D = 1) = E(U_1|X, Z).$$



## Proof.

Combining equations (N-3) and (N-4), we obtain

$$E(U_1|X, Z) = E(U_1|X, D = 1).$$

Integrating both sides of this expression against the distribution of  $Z$  conditional on  $X$ , we obtain

$$E(U_1|X) = E(U_1|X, D = 1).$$





- It is straightforward to show that (M-4) and (ATE-1), (ATE-2) jointly imply that  $U_1$  is mean independent of  $D$  conditional on  $X$ .

- In summary,  $(U_0, U_1)$  mean independent of  $D$  conditional on  $(X, Z)$  but not conditional on  $X$  implies that  $U_0$  is dependent on  $Z$  conditional on  $X$  in contradiction to all of the assumptions used to justify instrumental variables.

- In summary,  $(U_0, U_1)$  mean independent of  $D$  conditional on  $(X, Z)$  but not conditional on  $X$  implies that  $U_0$  is dependent on  $Z$  conditional on  $X$  in contradiction to all of the assumptions used to justify instrumental variables.
- Thus  $(U_0, U_1)$  mean independent of  $D$  conditional on  $(X, Z)$  but not conditional on  $X$  implies that none of the three sets of IV conditions will hold.

- In summary,  $(U_0, U_1)$  mean independent of  $D$  conditional on  $(X, Z)$  but not conditional on  $X$  implies that  $U_0$  is dependent on  $Z$  conditional on  $X$  in contradiction to all of the assumptions used to justify instrumental variables.
- Thus  $(U_0, U_1)$  mean independent of  $D$  conditional on  $(X, Z)$  but not conditional on  $X$  implies that none of the three sets of IV conditions will hold.
- In addition, if we weaken these conditions to only consider  $U_1$ , so that we assume that  $U_1$  is mean independent of  $D$  conditional on  $(X, Z)$  but not conditional on  $X$ , we obtain that  $U_1$  is dependent on  $Z$  conditional on  $X$ .

- We have shown that this implies that (IV-1) does not hold, and implies that (TT-1, TT-2) will not hold.

- We have shown that this implies that (IV-1) does not hold, and implies that (TT-1,TT-2) will not hold.
- A similar line of argument shows that (ATE-1,ATE-2) will not hold.

- We have shown that this implies that (IV-1) does not hold, and implies that (TT-1,TT-2) will not hold.
- A similar line of argument shows that (ATE-1,ATE-2) will not hold.
- Thus, the exclusion conditioning in matching is not the same as the exclusion conditioning in IV.

## Selection Formulae for the Matching Examples

- Consider a generalized Roy model of the form  $Y_1 = \mu_1 + U_1$ ;  $Y_0 = \mu_0 + U_0$ ;  $D^* = \mu_D(Z) + V$ ;  $D = 1$  if  $D^* \geq 0$ ,  $= 0$  otherwise; and  $Y = DY_1 + (1 - D)Y_0$ , where

$$\begin{aligned} (U_0, U_1, V)' &\sim N(0, \Sigma); \text{Var}(U_i) = \sigma_i^2 && i = 0, 1 \\ \text{Var}(V) &= \sigma_V^2; \text{Cov}(U_1, U_0) = \sigma_{10} \\ \text{Cov}(U_1, V) &= \sigma_{1V}; \text{Cov}(U_0, V) = \sigma_{0V}. \end{aligned}$$



## Selection Formulae for the Matching Examples

- Consider a generalized Roy model of the form  $Y_1 = \mu_1 + U_1$ ;  $Y_0 = \mu_0 + U_0$ ;  $D^* = \mu_D(Z) + V$ ;  $D = 1$  if  $D^* \geq 0$ ,  $= 0$  otherwise; and  $Y = DY_1 + (1 - D)Y_0$ , where

$$\begin{aligned} (U_0, U_1, V)' &\sim N(0, \Sigma); \text{Var}(U_i) = \sigma_i^2 & i = 0, 1 \\ \text{Var}(V) &= \sigma_V^2; \text{Cov}(U_1, U_0) = \sigma_{10} \\ \text{Cov}(U_1, V) &= \sigma_{1V}; \text{Cov}(U_0, V) = \sigma_{0V}. \end{aligned}$$

- Assume  $Z \perp\!\!\!\perp (U_0, U_1, V)$ .

## Selection Formulae for the Matching Examples

- Consider a generalized Roy model of the form  $Y_1 = \mu_1 + U_1$ ;  $Y_0 = \mu_0 + U_0$ ;  $D^* = \mu_D(Z) + V$ ;  $D = 1$  if  $D^* \geq 0$ ,  $= 0$  otherwise; and  $Y = DY_1 + (1 - D)Y_0$ , where

$$\begin{aligned} (U_0, U_1, V)' &\sim N(0, \Sigma); \text{Var}(U_i) = \sigma_i^2 & i = 0, 1 \\ \text{Var}(V) &= \sigma_V^2; \text{Cov}(U_1, U_0) = \sigma_{10} \\ \text{Cov}(U_1, V) &= \sigma_{1V}; \text{Cov}(U_0, V) = \sigma_{0V}. \end{aligned}$$

- Assume  $Z \perp\!\!\!\perp (U_0, U_1, V)$ .
- Let  $\phi(\cdot)$  and  $\Phi(\cdot)$  be the pdf and the cdf of a standard normal random variable.

- Then, the propensity score for this model for  $Z = z$  is given by:

$$\Pr(D^* > 0 | Z = z) = \Pr(V > -\mu_D(z)) = P(z) = \Phi\left(\frac{\mu_D(z)}{\sigma_V}\right).$$

Thus  $\frac{\mu_D(z)}{\sigma_V} = \Phi^{-1}(P(z))$ , and

$$\frac{-\mu_D(z)}{\sigma_V} = \Phi^{-1}(1 - P(z)).$$

- Then, the propensity score for this model for  $Z = z$  is given by:

$$\Pr(D^* > 0 | Z = z) = \Pr(V > -\mu_D(z)) = P(z) = \Phi\left(\frac{\mu_D(z)}{\sigma_V}\right).$$

Thus  $\frac{\mu_D(z)}{\sigma_V} = \Phi^{-1}(P(z))$ , and

$$\frac{-\mu_D(z)}{\sigma_V} = \Phi^{-1}(1 - P(z)).$$

- The event  $(V \leq 0, Z = z)$  can be written as

$$\frac{V}{\sigma_V} \leq -\frac{\mu_D(z)}{\sigma_V} \Leftrightarrow \frac{V}{\sigma_V} \leq \Phi^{-1}(1 - P(z)).$$

- Then, the propensity score for this model for  $Z = z$  is given by:

$$\Pr(D^* > 0 | Z = z) = \Pr(V > -\mu_D(z)) = P(z) = \Phi\left(\frac{\mu_D(z)}{\sigma_V}\right).$$

Thus  $\frac{\mu_D(z)}{\sigma_V} = \Phi^{-1}(P(z))$ , and

$$\frac{-\mu_D(z)}{\sigma_V} = \Phi^{-1}(1 - P(z)).$$

- The event  $\left(V \begin{smallmatrix} \leq \\ \geq \end{smallmatrix} 0, Z = z\right)$  can be written as

$$\frac{V}{\sigma_V} \begin{smallmatrix} \leq \\ \geq \end{smallmatrix} -\frac{\mu_D(z)}{\sigma_V} \Leftrightarrow \frac{V}{\sigma_V} \begin{smallmatrix} \leq \\ \geq \end{smallmatrix} \Phi^{-1}(1 - P(z)).$$

- We can write the conditional expectations required to get the biases for the treatment parameters as a function of  $P(z) = p$ .

- For  $U_1$  :

$$\begin{aligned}
 E(U_1 | D^* \geq 0, Z = z) &= \frac{\sigma_{1V}}{\sigma_V} E\left(\frac{V}{\sigma_V} \mid \frac{V}{\sigma_V} \geq \frac{-\mu_D(z)}{\sigma_V}\right) \\
 &= \frac{\sigma_{1V}}{\sigma_V} E\left(\frac{V}{\sigma_V} \mid \frac{V}{\sigma_V} \geq \Phi^{-1}(1 - P(z))\right) \\
 &= \eta_1 M_1(P(z))
 \end{aligned}$$

where

$$\eta_1 = \frac{\sigma_{1V}}{\sigma_V}.$$

- Similarly for  $U_0$  :

$$E(U_0 | D^* > 0, Z = z) = \eta_0 M_1(P(z))$$

$$E(U_0 | D^* < 0, Z = z) = \eta_0 M_0(P(z)),$$

where  $\eta_0 = \frac{\sigma_{0V}}{\sigma_V}$  and  $M_1(P(z)) = \frac{\phi(\Phi^{-1}(1-P(z)))}{P(z)}$  and

$M_0(P(z)) = -\frac{\phi(\Phi^{-1}(1-P(z)))}{1-P(z)}$  are inverse Mills ratio terms.

- Substituting these into the expressions for the biases for the treatment parameters conditional on  $z$  we obtain

$$\begin{aligned}\text{Bias } TT(P(z)) &= \eta_0 M_1(P(z)) - \eta_0 M_0(P(z)) \\ &= \eta_0 M(P(z)),\end{aligned}$$

$$\begin{aligned}\text{Bias } ATE(P(z)) &= \eta_1 M_1(P(z)) - \eta_0 M_0(P(z)) \\ &= M(P(z)) (\eta_1 (1 - P(z)) + \eta_0 P(z)).\end{aligned}$$