

# How To Correct for Sampling Biases

James J. Heckman  
University of Chicago

Econ 312, Spring 2022

## References:

- Amemiya, Ch. 10
- Different types of sampling
  - a random sampling
  - b censored sampling
  - c truncated sampling
  - d other non-random (exogenous stratified, choice-based)

## Standard Tobit Model (Tobin, 1958) “Type I Tobit”

$$y_i^* = x_i\beta + u_i$$

- Observe, i.e.,

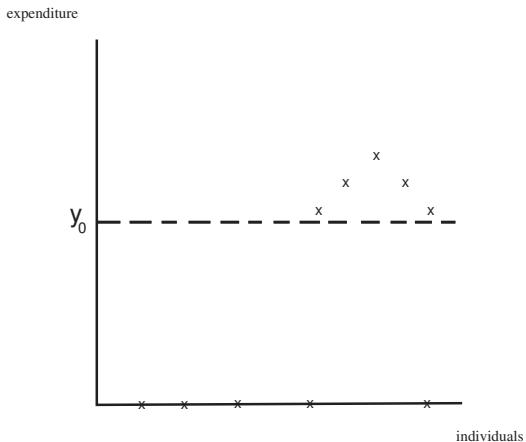
$$y_i = y_i^* \quad \text{if } y_i^* \geq y_0 \text{ or } y_i = 1 (y_i^* \geq y_0) y_i^*$$

$$y_i = 0 \quad \text{if } y_i^* < y_0$$

$$y_i = 1 \quad (y_i^* < y_0) y_i^*$$

- Tobin’s example—expenditure on a durable good only observed if good is purchased

Figure 1



Note: Censored observations might have bought the good if price had been lower.

- Estimator. Assume  $y_i^*/x_i \sim N(x_i\beta, \sigma_u^2)$

## Density of Latent Variables

$$g(y^*) = f(y_i^* | y_i^* < y_0) \Pr(y_i^* < y_0) + f(y_i^* | y_i^* \geq y_0) \cdot \Pr(y_i^* \geq y_0)$$

$$\Pr(y_i^* < y_0) = \Pr(x_i \beta + u_i < y_0) = \Pr\left(\frac{u_i}{\sigma_u} < \frac{y_0 - x_i \beta}{\sigma_u}\right) = \Phi\left(\frac{y_0 - x_i \beta}{\sigma_u}\right)$$

$$f(y_i^* | y_i^* \geq y_0) = \frac{\frac{1}{\sigma_u} \phi\left(\frac{y_i^* - x_i \beta}{\sigma_u}\right)}{1 - \Phi\left(\frac{y_0 - x_i \beta}{\sigma_u}\right)}$$

- **Question: Why?**

$$\begin{aligned} & \Pr(y^* = y_i^* | y_0 \leq y^*) \\ &= \Pr(x\beta + u = y_i^* | y_0 \leq x\beta + u) \\ & \Pr\left(\frac{u}{\sigma_u} = \frac{y_i^* - x\beta}{\sigma_u} \mid \frac{u}{\sigma_u} \geq \frac{y_0 - x\beta}{\sigma_u}\right) \end{aligned}$$

- Note that likelihood can be written as:

$$\mathcal{L} = \underbrace{\Pi_0 \Phi \left( \frac{y_0 - x_i \beta}{\sigma_u} \right) \Pi_1 \left( 1 - \Phi \left( \frac{y_0 - x_i \beta}{\sigma_u} \right) \right)}_{\text{This part you would set with just a simple probit}} \underbrace{\Pi_1 \frac{\frac{1}{\sigma_u} \phi \left( \frac{y_i^* - x_i \beta}{\sigma_u} \right)}{\left\{ 1 - \Phi \left( \frac{y_0 - x_i \beta}{\sigma_u} \right) \right\}}}_{\text{Additional information}}$$

- You could estimate  $\beta$  up to scale using only the information on whether  $y_i \begin{matrix} \geq \\ \leq \end{matrix} y_0$ , but will get more efficient estimate using additional information.
- \* if you know  $y_0$ , you can estimate  $\sigma_u$ .

## Truncated Version of Type I Tobit

Observe  $y_i = y_i^*$  if  $y_i^* > 0$   
( observe nothing for truncated observations )  
( example: only observe wages for workers )

$$\text{Likelihood: } \mathcal{L} = \prod_1 \frac{\frac{1}{\sigma_u} \phi\left(\frac{y_i^* - x_i \beta}{\sigma_u}\right)}{\Phi\left(\frac{x_i \beta}{\sigma_u}\right)}$$
$$\begin{aligned} \Pr(y_i^* > 0) &= \Pr(x\beta + u > 0) \\ &= \Pr\left(\frac{u}{\sigma_u} > \frac{-x\beta}{\sigma_u}\right) \\ &= \Pr\left(u < \frac{x\beta}{\sigma_u}\right) \end{aligned}$$

## Different Ways of Estimating Tobit

- a if censored, could obtain estimates of  $\frac{\beta}{\sigma_u}$  by simple probit
- b run OLS on observations for which  $y_i^*$  is observed

$$E(y_i | x_i\beta + u_i \geq 0) = x_i\beta + \sigma_u E\left(\frac{u_i}{\sigma_u} \mid \frac{u_i}{\sigma_u} > \frac{-x_i\beta}{\sigma_u}\right) \quad (y_0 = 0)$$

- where  $E(y_i | x_i\beta + u_i \geq 0)$  is the conditional mean for truncated normal r.v and

$$\sigma_u E\left(\frac{u_i}{\sigma_u} \mid \frac{u_i}{\sigma_u} > \frac{-x_i\beta}{\sigma_u}\right) \longrightarrow \lambda\left(\frac{x_i\beta}{\sigma_u}\right) = \frac{\phi\left(\frac{-x_i\beta}{\sigma_u}\right)}{\Phi\left(\frac{x_i\beta}{\sigma_u}\right)}$$

- $\lambda\left(\frac{x_i\beta}{\sigma_u}\right)$  known as “Mill’s ratio” ; bias due to censoring, can be viewed as an omitted variables problem



## Heckman Two-Step procedure

- Step 1: estimate  $\frac{\beta}{\sigma_u}$  by probit
- Step 2:

$$\text{form } \hat{\lambda} \left( \frac{x_i \hat{\beta}}{\sigma} \right)$$

regress

$$y_i = x_i \beta + \sigma \hat{\lambda} \left( \frac{x_i \beta}{\sigma} \right) + v + \varepsilon$$

$$v = \sigma \left\{ \lambda \left( \frac{x_i \beta}{\sigma} \right) - \hat{\lambda} \left( \frac{x_i \beta}{\sigma} \right) \right\}$$

$$\varepsilon = u_i - E(u_i | u_i > x_i \beta)$$

- Note: errors ( $v + \varepsilon$ ) will be heteroskedastic;
- need to account for fact that  $\lambda$  is estimated (Durbin problem)
- Two ways of doing this:
  - a Delta method
  - b GMM (Newey, Economic Letters, 1984)
  - c Suppose you run OLS using all the data

$$\begin{aligned}
 E(y_i) &= \Pr(y_i^* \leq 0) \cdot 0 + \Pr(y_i^* > 0) \left[ x_i\beta + \sigma_u E\left(\frac{u_i}{\sigma_u} \mid \frac{u_i}{\sigma_u} > \frac{-x_i\beta}{\sigma}\right) \right] \\
 &= \Phi\left(\frac{x_i\beta}{\sigma}\right) \left[ x_i\beta + \sigma_u \lambda\left(\frac{x_i\beta}{\sigma}\right) \right]
 \end{aligned}$$

- Could estimate model by replacing  $\Phi$  with  $\hat{\phi}$  and  $\lambda$  with  $\hat{\lambda}$ .
- For both (b) and (c), errors are heteroskedastic, meaning that you could use weights to improve efficiency.
- Also need to adjust for estimated regressor.
  - (d) Estimate model by Tobit maximum likelihood directly.

$$\begin{aligned}y_{1i}^* &= x_{1i}\beta + u_{1i} \\y_{2i}^* &= x_{2i}\beta + u_{2i} \\y_{2i} &= y_{2i}^* \quad \text{if } y_{1i}^* \geq 0 \\ &= 0 \quad \text{else}\end{aligned}$$

- Example
  - $y_{2i}$  student test scores
  - $y_{1i}^*$  index representing parents propensity to enroll students in school
  - Test scores only observed for population enrolled

$$\mathcal{L} = \Pi_1 [\Pr (y_{1i}^* > 0) f (y_{2i}|y_{1i}^* > 0)] \Pi_0 [\Pr (y_{1i}^* \leq 0)]$$

$$\begin{aligned} f (y_{2i}^*|y_{1i}^* \geq 0) &= \frac{\int_0^\infty f (y_{1i}^*, y_{2i}^*) dy_{1i}^*}{\int_0^\infty f (y_{1i}^*) dy_{1i}^*} \\ &= \frac{f (y_{2i}) \int_0^\infty f (y_{1i}^*|y_{2i}^*) dy_{1i}^*}{\int_0^\infty f (y_{1i}^*) dy_{1i}^*} \\ &= \frac{1}{\sigma^2} \phi \left( \frac{y_{2i}^* - x_{2i}\beta_2}{\sigma^2} \right) \cdot \frac{\int_0^\infty f (y_{1i}^*|y_{2i}^*) dy_{1i}^*}{\Pr (y_{1i}^* > 0)} \end{aligned}$$

$$y_{1i} \sim N (x_{1i}\beta_1, \sigma^2)$$

$$y_{2i} \sim N (x_{2i}\beta_2, \sigma^2)$$

$$y_{1i}^* | y_{2i}^* \sim N \left( x_{1i}\beta_1 + \frac{\sigma_{12}}{\sigma_2^2} (y_{2i} - x_{2i}\beta_2), \sigma_1^2 - \frac{\sigma_{12}^2}{\sigma_2^2} \right)$$

$$E(y_{1i}^* | u_{2i} = y_{2i}^* - x_{2i}\beta) = x_{1i}\beta_1 + E(u_{1i} | u_{2i} = y_{2i}^* - x_{2i}\beta)$$

$$L = \prod_0 \left[ 1 - \Phi \left( \frac{x_{1i}\beta}{\sigma_1} \right) \right] \prod_1 \frac{1}{\sigma_2} \cdot \phi \left( \frac{y_{2i}^* - x_{2i}\beta_2}{\sigma_2} \right) \cdot \left\{ 1 - \Phi \left( \frac{- \left\{ x_{1i}\beta_1 + \frac{\sigma_{12}}{\sigma_2^2} (y_{2i} - x_{2i}\beta_2) \right\}}{\sigma^x} \right) \right\}$$

## Estimation by Two-Step Approach

- Using data on  $y_{2i}$  for which  $y_{1i} > 0$

$$\begin{aligned} E(y_{2i} | y_{1i} > 0) &= x_{2i}\beta + E(u_{2i} | x_i\beta + u_{1i} > 0) \\ &= x_{2i}\beta + \sigma_2 E\left(\frac{u_{2i}}{\sigma_2} \mid \frac{u_{1i}}{\sigma_1} > \frac{-x_{1i}\beta_1}{\sigma_1}\right) \\ &= x_{2i}\beta + \frac{\sigma_{12}}{\sigma_1\sigma_2} E\left(\frac{u_{1i}}{\sigma_1} \mid \frac{u_{1i}}{\sigma_1} > \frac{-x_{1i}\beta_1}{\sigma_1}\right) \\ &= x_{2i}\beta_2 + \frac{\sigma_{12}}{\sigma_1} \lambda\left(\frac{-x_{1i}\beta_1}{\sigma_1}\right) \end{aligned}$$

## Example: Female labor supply model

$$\begin{aligned} & \max \quad u(L, x) \\ \text{s.t.} \quad & x = wH + v \quad H = 1 - L \\ \text{where } & H : \text{ hours worked} \\ & v : \text{ asset income} \\ & w \text{ given} \\ & P_x = 1 \\ & L : \text{ time spent at home for child care} \end{aligned}$$

$$\frac{\frac{\partial u}{\partial L}}{\frac{\partial u}{\partial x}} = w \quad \text{when } L < 1$$

$$\text{reservation wage} = MRS \big|_{H=0} = w_R$$





## Example: Female labor supply model

- We don't observe  $w_R$  directly.

Model

$$w^0 = x\beta + u \quad (\text{wage person would earn if they worked})$$
$$w^R = z\gamma + v$$
$$w_i = w_i^0 \quad \text{if } w_i^R < w_i^0$$
$$= 0 \quad \text{else}$$

- Fits within previous Tobit framework if we set

$$y_{1i}^* = x\beta - z\gamma + u - v = w^0 - w^R$$
$$y_{2i} = w_i$$

$$\begin{aligned}w^0 &= x_{2i}\beta_2 + u_{2i} \quad \text{given} \\MRS &= \frac{\frac{\partial u}{\partial L}}{\frac{\partial u}{\partial x}} = \gamma H_i + z_i'\alpha + v_i\end{aligned}$$

(Assume functional form for utility function that yields this)

$$\begin{aligned}
w^r (H_i = 0) &= z_i' \alpha + v_i \\
\text{work if } w^0 &= x_{2i} \beta_2 + u_{2i} > z_i \alpha + v_i \\
\text{if work, then } w_i^0 &= MRS \implies x_{2i} \beta_2 + u_{2i} = \alpha H_i + z_i \alpha + v_i \\
&\implies H_i = \frac{x_{2i} \beta_2 - z_i' \alpha + u_{2i} - v_i}{\gamma} \\
&= x_{1i} \beta_1 + u_{1i} \\
\text{where } x_{1i} \beta_1 &= (x_{2i} \beta_2 - z_i \alpha) \gamma^{-1} \\
u_{1i} &= u_{2i} - v_i
\end{aligned}$$

## Type 3 Tobit Model

$$y_{1i}^* = x_{1i}\beta_1 + u_{1i} \leftarrow \text{hours}$$

$$y_{2i}^* = x_{2i}\beta_1 + u_{2i} \leftarrow \text{wage}$$

$$\begin{aligned} y_{1i} &= y_{1i}^* && \text{if } y_{1i}^* > 0 \\ &= 0 && \text{if } y_{1i}^* \leq 0 \end{aligned}$$

$$\begin{aligned} y_{2i} &= y_{2i}^* && \text{if } y_{2i}^* > 0 \\ &= 0 && \text{if } y_{2i}^* \leq 0 \end{aligned}$$

$$\begin{aligned} \text{Here } H_i &= H_i^* && \text{if } H_i^* > 0 \\ &= 0 && \text{if } H_i^* \leq 0 \end{aligned}$$

$$\begin{aligned} w_i &= w_i^0 && \text{if } H_i^* > 0 \\ &= 0 && \text{if } H_i^* \leq 0 \end{aligned}$$

- Note: Type IV Tobit simply adds

$$\begin{aligned} y_{3i} &= y_{3i}^* && \text{if } y_{1i}^* > 0 \\ &= 0 && \text{if } y_{1i}^* \leq 0 \end{aligned}$$

- Can estimate by
  - ① maximum likelihood
  - ② Two-step method

$$E(w_i^0 | H_i > 0) = \gamma H_i + z_i \alpha + E(v_i | H_i > 0)$$

## Type V Tobit Model of Heckman (1978)

$$y_{1i}^* = \gamma y_{2i} + x_{1i}\beta + \delta_2 w_i + u_{1i}$$

$$y_{2i} = \gamma_2 y_{1i}^* + x_{2i}\beta_2 + \delta_2 w_i + u_{2i}$$

- Analysis of an antidiscrimination law on average income of African Americans in  $i$ th state.
- Observe  $x_{1i}$ ,  $x_{2i}$ ,  $y_{2i}$  and  $w_i$

$$w_i = 1 \quad \text{if } y_{1i}^* > 0$$

$$w_i = 0 \quad \text{if } y_{1i}^* \leq 0$$

- $y_{2i}$  = average income of African Americans in the state
- $y_{1i}^*$  = unobservable sentiment towards African Americans
- $w_i$  = if law is in effect

- Adoption of Law is endogenous
- Require restriction  $\gamma\delta_2 + \delta_1 = 0$  so that we can solve for  $y_{1j}^*$  as a function that does not depend on  $w_j$ .
- This class of models known as “dummy endogenous variable” models.

### **Coherency Problem (Suppose Restriction Does Not Bind?)**

- See notes on “Dummy Endogenous Variables in simultaneous equations.”



### References:

- Heckman (AER, 1990) “Varieties of Selection Bias”
- Heckman (1980), “Addendum to Sample Selection Bias as Specification Error”
- Heckman and Robb (1985, 1986)

$$y_1^* = x\beta + u$$

$$y_2^* = z\gamma + v$$

$$y_1 = y_1^* \quad \text{if } y_2^* > 0$$

$$\begin{aligned} E(y_1^* \mid \text{observed}) &= x\beta + E(u \mid x, z\gamma + u > 0) \\ &\quad + [u - E(u \mid x, z\gamma + u > 0)] \\ &\quad \frac{\int_{-\infty}^{\infty} \int_{-\infty}^{-z\gamma} uf(u, v \mid x, z) dv du}{\int_{-\infty}^{\infty} \int_{-\infty}^{-z\gamma} f(uv \mid x, z) dv du} \end{aligned}$$

- Note:

$$\Pr(y_2^* > 0 \mid z) = \Pr(z\gamma + u > 0 \mid z) = P(Z) = 1 - F_v(-z\gamma)$$

$$\Rightarrow F_v(-z\gamma) = 1 - P(Z)$$

$$\Rightarrow -z\gamma = F_v^{-1}(1 - P(Z)) \quad \text{if } F_v$$

- Can replace  $-z\gamma$  in integrals in integrals by  $F_v^{-1}(1 - P(Z))$  if in addition  $f(u, v | x, z) = f(u, v | z\gamma)$  (index sufficiency)
- Then

$$E(y_1^* | y_2 > 0) = x\beta + g(P(z)) + \varepsilon \text{ where } g(P(Z))$$

is bias or “control function.”

- Semiparametric selection model-Approximate bias function by Taylor series in  $P(z\gamma)$ , truncated power series.