

# Definition of Samples

James J. Heckman  
Center for the Economics of Human Development  
University of Chicago

Econ 312, Spring 2023

- General definitions of
  - ① Random Sampling
  - ② Censored Sampling
  - ③ Truncated Sample
  - ④ Choice Based Sample
  - ⑤ Also consider truncated and censored random variables.

(1) Random sampling: (Really *simple* random sampling)

- iid. random variables with density  $f(X)$ . Random sampling in general is derivation of a sample by a *calculatable* rule.

$$\text{Prob. of sample} = f(X_1)f(X_2)f(X_3) \dots f(X_N)$$

- Problem of getting an  $X$  is  $f(X)$ . Thus in a population, the probability of getting into the sample is  $f(X)$ . This is simple random sampling.

(2) Truncated Sample

$$f(X) \quad : \quad \text{density of a random variable}$$
$$a < X < b : b, a \text{ may be infinite}$$

- We observe  $X$  if  $X < R$  (right truncation)
- or if  $X > L$  (left truncation).
- Key property: latent variable  $X$  in the population — we know

$$X^* = X \quad \text{when } L < X < R$$

- (Assume simple random sampling of a larger population). We only observe  $X^*$  and we do not know the number of observations in (larger) random sample for which  $X$  is outside the interval. We only know the reduced sample if density in population (untruncated) is  $f(X)$ , then density of  $X^*$  is

$$\frac{f(X^*)}{\int_L^R f(z) dz} \quad L \leq X^* \leq R$$

- (Note further that there are an infinity of underidentified distributions consistent with the truncated one.)
  - 3 Censored Sample: We observe  $X^*$  as before but *we know the number of observations outside interval*. We encounter two types of censoring:
    - a Type one censoring : we only observe a variable if it lies in a range, number of values of  $Y$  outside the range is known.
    - b Type Two Censoring: *Fixed proportion* of the sample is censored in advance (e.g. stop observing light bulb burnout when we have a proportion - say  $m$ ).

- (4) If we have that in both (3) and (2),  $X$  is a *truncated random variable* (the range of the random variable is truncated).
- (5) New term: coined in recent econometric work — *censored random variable*. It is inherently a bivariate concept. Joint *pdf*  $-f(Y_1, Y_2)$ . Then we have that we observe  $Y_1$  only if  $Y_2$  exceeds some value or lies in some range, e.g.

$$L < Y_2 < R \quad (1)$$

Prob. of this event is

$$\int_L^R f_2(Y_2) dY_2$$

- The random variable  $Y_1$  is *not truncated*. We observe  $Y_1$  only if the condition on  $Y_2$  is satisfied.

- The *sample* may or may not be truncated. Thus, it is the case that if we observe  $Y_1$ , given selection criterion (\*), but we do not know the number of observations in the larger random sample variable for which the  $Y_2$  restriction is violated, we have a truncated sample and a censored random variable. Now clearly we may put a restriction on  $Y_1$  e.g. we observe  $Y_1$  only if  $L_2 < Y_2 < R_2$  and  $L_1 < Y_1 < R_1$ . Thus define  $Y_1^* = Y_1$  for  $L_1 < Y_1 < R_1$ .

$$g(Y_1^*) = \frac{\int_{L_2}^{R_2} f(Y_1^*, Y_2) dY_2}{\int_{L_1}^{R_1} \int_{L_2}^{R_2} f(Y_1, Y_2) dY_1 dY_2}$$

- (6) New term in discrete choice literature – *choice based sampling*. Consider the random variable  $Y$  to be discrete.  $Z$  are exogenous explanatory variables. The theory produces a  $g(Y | Z, \theta)$ : discrete choice model  $h(Z)$  in the distribution of the population exogenous variables.

$$Y_j \in \{1, \dots, J\}$$

elements of choice set.

Exogenous Sampling: we pick  $Z$ , then observe  $Y$ . Sample  $Z$  according to the density  $k(Z)$  and observe the value of  $Y$ , the choice. Likelihood of an observation  $(Y, Z)$  is

$$g(Y | Z, \theta)k(Z)$$

when  $k(Z) = h(Z)$ , we have random sampling. Otherwise we have *stratified* sampling.



# Choice Based Samples

- Pick  $Y$  first (e.g. travel mode). Probability of selecting  $Y$  is  $C(Y)$ .
- $f(Y, Z)$  is the joint density of  $Y$  and  $Z$  in the population.

$$f(Y, Z | \theta) = g(Y | Z, \theta)h(Z) = \varphi(Z | Y)f(Y | \theta)$$

$$f(Y | \theta) = \int g(Y | Z, \theta)h(Z)dZ$$

- Given  $Y$  we observe  $Z$  (the implicit assumption is that we are sampling only on  $Y$ , not on  $Y$  and  $Z$ ). Probability of *sampled*  $Z, Y$  is  $\varphi(Z | Y)C(Y)$ .
- A fact we use later is

$$\begin{aligned}\varphi(Z | Y)C(Y) &= \left\{ \frac{g(Y | Z)h(Z)}{f(Y)} \right\} C(Y) \\ &= \frac{g(Y | Z)h(Z)C(Y)}{\left[ \int g(Y | Z)h(Z)dZ \right]}.\end{aligned}$$

When  $C(Y) = f(Y) = \int g(Y | Z)h(Z)dZ$ , choice based sampling is random sampling.

- Note, the likelihood function in an exogenous sampling scheme is

$$\mathcal{L} = \prod_{i=1}^I f(Y_i, Z_i) = \prod_{i=1}^I f(Y_i | Z_i, \theta)h(Z_i)$$
$$\ln \mathcal{L} = \sum_{i=1}^I \ln f(Y_i | Z_i) + \sum \ln h(Z_i).$$

- By exogeneity, we get the lack of dependence of distribution of  $Z$  on  $\theta$ .

- Likelihood function for a choice-based sampling scheme is

$$\ln \mathcal{L} = \sum_{i=1}^I [\ln g(Y_i | Z_i) + \ln h(Z_i) - \ln f(Y_i) + \ln C(Y_i)].$$

- In several,  $f(Y)$  depends on parameters  $\theta$ .  $\therefore$  Max with  $\theta$ .

$$\frac{\partial \ln \mathcal{L}}{\partial \theta} = \sum_{i=1}^I \frac{\partial \ln g(Y_i | Z_i)}{\partial \theta} - \sum_{i=1}^I \frac{\partial \ln f(Y_i)}{\partial \theta}.$$

- We neglect the second term in forming the usual estimators using only the first term. That is the source of the inconsistency.

## Choice Based Sample:

- An example in discrete choice.
- (c) Draw  $d$  by  $\varphi(d)$ .
- (d) Draw  $X$  by  $f(X | d = 1)$ .
- Joint density of data:

$$\begin{aligned} & \varphi(d = 1)f(X | d = 1, \theta_0) \\ = & \varphi(d = 1) \left[ \frac{\Pr(d = 1 | X, \theta_0)f(X)}{\Pr(d = 1 | \theta_0)} \right] \end{aligned}$$

- Now in a choice-based sample

$$\Pr^*(d = 1 | X) = \frac{f(X | d = 1, \theta_0)\varphi(d = 1)}{g^*(X)}$$

where  $g^*(X)$  is the sampled  $X$  data. Joint density of *data*  $X$  is given by:

$$g^*(X) = f(X | d = 1, \theta)\varphi(d = 1) + f(X | d = 0, \theta)\varphi(d = 1)$$

and

$$\Pr(d = 1 | X) = \frac{f(X | d = 1) \Pr(d = 1)}{f(X)}$$

- Assume  $f(X) > 0$ . Using Bayes' theorem for  $Y$  write:

- $$\Pr^*(d = 1 | X) = \frac{\frac{\Pr(d = 1 | X, \theta)f(X)}{\Pr(d = 1 | \theta)}\varphi(d = 1)}{\frac{\Pr(d = 1 | X, \theta)f(X)}{\Pr(d = 1 | \theta)}\varphi(d = 1) + \frac{\Pr(d = 0 | X, \theta)f(X)}{\Pr(d = 0 | \theta)}\varphi(d = 0)}$$

$$= \frac{\Pr(d = 1 | X, \theta)\varphi(d = 1) / \Pr(d = 1 | \theta)}{\Pr(d = 1 | X, \theta)\frac{\varphi(d = 1)}{\Pr(d = 1 | \theta)} + \Pr(d = 0 | X, \theta)\frac{\varphi(d = 0)}{\Pr(d = 0 | \theta)}}.$$



- Now we missample the population with density  $f(X | d = 1)$  in a choice based sample:

$$\begin{aligned}
 \Pr^*(d = 1 | X) &= \frac{f(X | d = 1, \theta_0)\varphi(d = 1)}{f(X | d = 1, \theta_0)\varphi(d = 1) + f(X | d = 0, \theta_0)\varphi(d = 0)} \\
 &= \frac{\frac{f(X)\Pr(d = 1 | X)}{\Pr(d = 1)}\varphi(d = 1)}{\frac{f(X)\Pr(d = 1 | X)}{\Pr(d = 1)}\varphi(d = 1) + \frac{f(X)\Pr(d = 0 | X)}{\Pr(d = 0)}\varphi(d = 0)} \\
 &= \frac{\Pr(d = 1 | X)}{\Pr(d = 1 | X) + \Pr(d = 0 | X) \frac{\varphi(d = 0)}{\varphi(d = 1)} \cdot \frac{\Pr(d = 1)}{\Pr(d = 0)}} \\
 &= \frac{1}{1 + \left[ \frac{\Pr(d = 0 | X)}{\Pr(d = 1 | X)} \right] \cdot \frac{\varphi(d = 0)}{\varphi(d = 1)} \cdot \frac{\Pr(d = 1)}{\Pr(d = 0)}}
 \end{aligned}$$

- With logit we get

$$\Pr^*(d = 1 | X) = \frac{1}{1 + e^{-(\alpha_0 + X\beta) + \ln \left[ \frac{\varphi(d = 0) \Pr(d = 1)}{\varphi(d = 1) \Pr(d = 0)} \right]}}.$$

This goes into an intercept term:

$$\begin{aligned} &= \frac{e^{\alpha^* + X\beta}}{1 + e^{\alpha^* + X\beta}} \\ \alpha^* &= \alpha_0 - \ln \left[ \frac{\varphi(d = 0) \Pr(d = 1)}{\varphi(d = 1) \Pr(d = 0)} \right]. \end{aligned}$$

- How to solve problem: Reweight data by relative frequency in population.
- (Idea due to C.R. Rao, 1965, 1986.)
- Joint density of the data is

$$f(X | d = 1)\varphi(d = 1).$$

Use Bayes' rule to obtain

$$\frac{P(d = 1 | X)f(X)}{P(d = 1)}\varphi(d = 1).$$

- Now weight by

$$\frac{P(d = 1)}{\varphi(d = 1)}.$$

- Solution: Reweight the data to form the following weighted likelihood:

$$\frac{1}{N} \sum_{i=1}^N \left[ \frac{\Pr(d_i = 1)}{\varphi(d_i = 1)} (d_i^*) \ln \Pr(d_i = 1 | X, \theta) + \frac{\Pr(d_i = 0)}{\varphi(d_i = 0)} (1 - d_i^*) \ln \Pr(d_i = 0 | X, \theta) \right]$$

$$P \int \{ [\Pr(d = 1 | X, \theta_0) f(X | \theta_0)] \ln \Pr(d = 1 | X, \theta) + [\Pr(d = 0 | X, \theta_0) f(X | \theta_0)] \ln \Pr(d = 0 | X, \theta) \} f(X | d) dX$$

- This step uses the result that reweighting the data gives us the true density.
- Better way to see what is giving on:

$$\frac{f(X | d = 1)\varphi(d = 1)}{g^*(X)} = \frac{\Pr(d = 1 | X)f(X)}{g^*(X)} \frac{\varphi(d = 1)}{\Pr(d = 1)}.$$

- Reweight the data: when we reweight the data,  $g^*$  is restored to  $f$ .

$$f(X) = f(X | d = 1)\varphi(d = 1) \left[ \frac{P(d = 1)}{\varphi(d = 1)} \right] + f(X | d = 0)\varphi(d = 0) \frac{\Pr(d = 0)}{\varphi(d = 0)}$$