

Sample Selection

(GB: Chapter 19.3-19.7)
Isabel Casas

Motivation

- So far we have assumed a random sample from the population.
 - Each individual of the population has the same chance to be selected

Motivation

- So far we have assumed a random sample from the population.
 - Each individual of the population has the same chance to be selected
- Often this is not the case (truncated samples):
 - Self selection: individual choice of whether or not to participate in the activity of interest
 - Sample selection: Those who participate in the activity of interest are oversampled (only participants)

Motivation

We define:

- Specify population from where the sample has to be obtained
- *Random sample* a sample where all individuals have the same probability to be chosen
- The *selected sample* is a nonrandom sample following a selection rule. The nonrandom nature is due to
 - sample design
 - behaviour of the units being sample
 - nonresponse in surveys
 - Selection mechanisms that result into nonrandom samples

Nonrandomly selected samples

Example 19.1 (Saving Function). We have only families whose household head is greater than 45 old

$$saving = \beta_0 + \beta_1 \text{ income} + \beta_2 \text{ age} + \beta_3 \text{ married} + \beta_4 \text{ kid} + \epsilon$$

We are interested in results of every family but we only have data for $age \geq 45 \Rightarrow$ sample selection issue.

- Not such a big problem. We define the population as the group of families greater than 45 and we have a random sample of that subset

Nonrandomly selected samples

Example 19.2 (Truncation based on wealth). We are interested in estimating how a worker eligibility in a certain plan affects the wealth of the family.

$$wealth = \beta_0 + \beta_1 plan + \beta_2 educ + \beta_3 age + \beta_4 income + \epsilon$$

We only have data for people with a net wealth less than \$200000. We are interested in results of every family but we have a selected sample depending of the y variable

- This is a more serious sampling problem

Nonrandomly selected samples

Example 19.3 (Wage offer function). We are interested in estimating the wage offer function for people of working age.

We are interested in results for every working age worker but we only have the ones who have took a job offer.

This is a self-selection example because people self-select into employment. It makes it harder to determine the causation.

More examples

- Self-selection example: sample only individuals who stay in a diet program (the drop outs are not taken into account)
- Sample selection example: we record wages from individuals in the labour force. We are only picking the sample of "industrious" people rather than the population as a whole
- Surviving selection: We are interested in data on assets for 15 years, we only take the ones who have been traded in the market for the last 15 years. We do not take into account the drop outs.

It doesn't matter what is the name we give to the sample selection problem, the result is that it can be difficult to determine the causation because of the bias induced.

Motivation

- How does sample selection affect estimation results?
- Should we do something about it?

Outline

- Selection on the basis of X : exogenous sampling
 - Eg: Segmenting an existing sample by gender or status.
 - Sample selection problem can be ignored
- Selection on the basis of y (the LHS variable)
 - More difficult to resolve
 - ① Deterministic selection rule (known): truncated regression
 - ② Random selection rule (unknown): behavioural selection
- ② If Selection is determined by behaviour
 - Selection rule is unknown
 - Probit models for selection (with and without endogeneity)
 - Tobit models for selection (with and without endogeneity)
 - Structural Tobits with selection

Cases where the sample selection can be ignored

- Selection on the basis of \mathbf{X} (no endogeneity)
 - The selection rule is independent of \mathbf{X} and ϵ or
 - The selection rule is deterministic



OLS yields consistent estimators

- Selection on the basis of \mathbf{X} (with endogeneity)
 - The selection rule is independent of \mathbf{z} (instruments) and ϵ or
 - The selection rule is deterministic



2SLS yields consistent estimators

Selection on the basis of \mathbf{X}

Assume the following model:

$$\mathbf{y} = \beta_0 + \beta_1 \mathbf{x}_1 + \dots + \beta_k \mathbf{x}_k + \epsilon \quad E(\epsilon|\mathbf{X}) = 0$$

- Let $s = 0, 1$ be the selection indicator.
 - Rule: $20 < age < 45$ ($s = 1$), otherwise $s = 0$
 - $age \in \mathbf{X}$
- The crucial assumption is that ϵ (the unobserved part) is randomly distributed in the population.
- Sample selection is a problem if we get a non-random sample of the ϵ 's (e.g. all those with high ϵ 's)

Selection on the basis of \mathbf{X}

- If \mathbf{X} are exogenous variables and
- If $E(\epsilon|\mathbf{X}, s) = 0$, then sample selection can be ignored
- I.e. in particular if:
 - s is independent of (\mathbf{X}, ϵ) : $E(\epsilon|\mathbf{X}, s) = E(\epsilon|\mathbf{X}) = 0$ or
 - s is a deterministic function of \mathbf{X} : $E(\epsilon|\mathbf{X}, s) = E(\epsilon|\mathbf{X}) = 0$
- Intuition:
 - If s is independent of (\mathbf{X}, ϵ) and hence $(\mathbf{X}, y) \Rightarrow$ the selected sample is also a random sample
 - If s is a deterministic function of $\mathbf{X} \Rightarrow$ the selection does not give us those with particularly high (or low) ϵ 's, as \mathbf{X} contain no info about the ϵ 's
- Estimation:
 - OLS yields consistent estimates

Selection on the basis of X

- This result can be extended to the case with endogenous X 's as long as:

$$E(\epsilon | z, s) = 0$$

- Where z is the vector of instruments (including exogenous X 's), i.e. if
 - s is independent of (z, ϵ) , or
 - s is a deterministic function of z
- This is stronger than what we need for random sample 2SLS ($E(z, \epsilon) = 0$)
- Estimation by 2SLS still yields consistent and asymptotically normal estimates.
- But not if s is a function of y !! (we return to this case)

Selection on the basis of \mathbf{X}

Example:

$$wage = \mathbf{X}\beta + \gamma \text{ age} + \epsilon$$

$$E(\epsilon|\mathbf{X}, \text{age}) = 0$$

Selection rules:

- Deterministic: $k_1 < \text{age} < k_2 \Rightarrow$ determined by RHS \Rightarrow OLS is consistent
- Random: $\text{age} > k_1 + r$ where r is a random variable, then OLS still ok if r does not contain info about ϵ :

$$E(\epsilon|\mathbf{X}, \text{age}, r) = 0$$

Selection on the basis of \mathbf{X}

The result can also be extended to non-linear models:

- Probit
- Logit

If selection is ignorable: $E(y|\mathbf{X}, s) = E(y|\mathbf{X})$. This occurs if:

- s is a deterministic function of \mathbf{X} , or
- s is independent of (\mathbf{X}, y)

Exercise 1 (5 minutes)

Discuss the following questions with your colleagues:

- In the linear model: When can selection be ignored? What is the intuition for this?
- In the probit/logit models: When can selection be ignored? What is the intuition for this?
- What if we select only those with:
 - $(\mathbf{x}_1 + \mathbf{x}_2)^2 > 5$?
 - $(\mathbf{x}_1 + \mathbf{x}_2)^2 > Y$?
 - $(\mathbf{x}_1 + \mathbf{x}_2)^2 > \text{average of } Y$?

Look at simulations!!!

Deterministic selection of y : Truncated regression

- Selection on the basis of y
- Selection rule is deterministic and known so it can be accounted for
- Find the log-lik function of the truncated distribution and obtain estimates by ML.
- If $y|\mathbf{X} \sim N(0, \sigma^2)$ then we can use the truncated Tobit model.
- It is less efficient than censored regression
- If there is heteroskedasticity or the distribution is misspecified \Rightarrow inconsistency

Truncated regression

A standard linear model:

$$E(y_j|\mathbf{x}_j) = \mathbf{x}_j\boldsymbol{\beta}$$

Selection rule:

$$s_j = \mathbf{1}[a_1 < y_j < a_2]$$

- Rule is known and observable and it depends on \mathbf{y}
 - Income equations: Individuals in certain income brackets
 - Labour supply: normal time workers
- The value of $\mathbf{X}\boldsymbol{\beta}$ will impose the observed ϵ 's in our sample and vice versa because the value of \mathbf{y} is bounded.
- In our estimation, we need to condition on the rule

Truncated regression

- We know the probability of $y_j < y$ conditional on \mathbf{x}_j :
 $F(y|\mathbf{x}_j; \beta, \gamma)$
- What is the distribution of y_j conditional on \mathbf{x}_j and selection
 $a_1 < y < a_2$?

Truncated regression

- We know the probability of $y_j < y$ conditional on \mathbf{x}_j :
 $F(y|\mathbf{x}_j; \beta, \gamma)$
- What is the distribution of y_j conditional on \mathbf{x}_j and selection
 $a_1 < y < a_2$?

$$\begin{aligned}
 P(y_j < y | \mathbf{x}_j, s = 1) &= \frac{P(y_j < y, s = 1 | \mathbf{x}_j)}{P(s = 1 | \mathbf{x}_j)} \\
 &= \frac{P(y_j < y, a_1 < y_j < a_2 | \mathbf{x}_j)}{P(a_1 < y_j < a_2 | \mathbf{x}_j)} \\
 &= \frac{P(a_1 < y_j < y | \mathbf{x}_j)}{P(a_1 < y_j < a_2 | \mathbf{x}_j)} \\
 &= \frac{F(y | \mathbf{x}_j; \beta, \gamma) - F(a_1 | \mathbf{x}_j; \beta, \gamma)}{F(a_2 | \mathbf{x}_j; \beta, \gamma) - F(a_1 | \mathbf{x}_j; \beta, \gamma)}
 \end{aligned}$$

Truncated regression

Now the density function of y_j given \mathbf{x}_j and $s_j = 1$ can be found by differentiating wrt y :

$$f(y|\mathbf{x}_j, s = 1) = \frac{f(y|\mathbf{x}_j; \beta, \gamma)}{F(a_2|\mathbf{x}_j; \beta, \gamma) - F(a_1|\mathbf{x}_j; \beta, \gamma)}$$

Which is just the unconditional density (the numerator) divided by the probability of being sampled (the denominator) I.e. we are simply rescaling the densities to the new interval: (a_1, a_2)

How to get the log-lik?

- Take logs of the above and
- sum over the selected sample.

Truncated regression

Special case: if $\mathbf{y}|\mathbf{X} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2) \Rightarrow$ use Φ and ϕ instead of F and f :

- The Truncated Tobit Model (or Truncated Normal Regression Model)

Censored regression models (Ch. 16): Only \mathbf{y} was not observed for some individuals but information about \mathbf{X} from these could still be used.

Truncated regression model (now): Neither \mathbf{X} nor \mathbf{y} is observed for some individuals

Same code but censoring is preferable over truncation

Truncated regression

If misspecification:

- Heteroscedasticity, or
- Non-normality

This implies:

- Changes the log-lik
- Inconsistent estimates (same as in Chapter 16)

Exercise 2 (5 minutes)

Discuss with your colleagues:

- Think of an example where selection on y is likely to appear.
- How will that ? intuitively affect your estimations if not corrected for?
- Explain ? using your own words ? what the correction does (to the likelihood function)? Consider, e.g., an observation with an extreme X but a y value between a_1 and a_2 . How will its density be altered?

Selection determined by behaviour

Cases to consider:

- A) Selection can be modelled as probit
- B) ... and with endogenous RHS variables
- D) Selection can be modelled as tobit
- D) ... and with endogenous RHS variables
- E) Structural Tobit with selection

Selection determined by behaviour

Classic example: Wage offer equation (Example 19.5)

- How education affects how much a person could earn in the labour market.
- Your first intuition is to run OLS as if you had a random sample
- But only individuals with wages above their reservation wage have jobs and can thus be sampled
- Reservation wage varies across individuals (otherwise, we could use approach from previous slides)

Other example:

- Non-response in surveys

Selection determined by behaviour

The wage equation example (formally):

- Wage: $\log w_j = \mathbf{x}_{j1}\beta_1 + \epsilon_{1j}$
- Reservation wage:

$$\log w_j^{res} = \mathbf{x}_{j2}\beta_2 + \gamma_2 \underbrace{a_j}_{\text{no wage income}} + \epsilon_{2j}$$
- Selection rule: $\log w_j > \log w_j^{res}$
 - We don't know $w_j^{res} \Rightarrow$ not a truncated or censored regression model
 - Instead, selection depends on unobservables (that are correlated with the error term in the wage equation):

$$\log w_j > \log w_j^{res} \Rightarrow \mathbf{x}_{j1}\beta_1 - \mathbf{x}_{j2}\beta_2 - \gamma_2 a_j + \epsilon_{1j} - \epsilon_{2j} > 0$$

Selection determined by behaviour

Structural equation (what we are interested on studying):

$$y_{1j} = \log w_j = \mathbf{x}_{j1}\boldsymbol{\beta}_1 + \epsilon_{j1}$$

Selection equation:

$$\begin{aligned} y_{2j} &= \mathbf{1}[\mathbf{x}_{j1}\boldsymbol{\beta}_1 - \mathbf{x}_{j2}\boldsymbol{\beta}_2 - \gamma_2 a_j + \epsilon_{1j} - \epsilon_{2j} > 0] \\ &= \mathbf{1}[\mathbf{x}_j\boldsymbol{\delta}_2 + \nu_{2j} > 0] \end{aligned}$$

A) Probit selection

Model (type II Tobit):

$$\begin{aligned} y_1 &= \mathbf{X}_1 \beta_1 + \epsilon_1 \\ y_2 &= \mathbf{1}[\mathbf{X} \delta_2 + \nu_2 > 0] \quad \text{where } \mathbf{X}_1 \in \mathbf{X} \end{aligned}$$

Assumptions:

- (\mathbf{X}, y_2) are always observed
- y_1 only observed when $y_2 = 1$
- (ϵ_1, ν_2) independent of \mathbf{X} (exogenous)
- $\nu_2 \sim N(0, 1)$
- $E(\epsilon_1 | \nu_2) = \gamma_1 \nu_2$ (rule and error correlated)

\mathbf{X}_1 could in principle contain variables not in \mathbf{X} (but undesirable)

A) Probit selection

Model (type II Tobit):

$$\begin{aligned} \mathbf{y}_1 &= \mathbf{X}_1\beta_1 + \epsilon_1 \\ \mathbf{y}_2 &= \mathbf{1}[\mathbf{X}\delta_2 + \nu_2 > 0] \quad \text{where } \mathbf{X}_1 \in \mathbf{X} \end{aligned}$$

We know how to calculate:

$$E(\mathbf{y}_1|\mathbf{X}) =$$

$$P(\mathbf{y}_2 = 1|\mathbf{X}) =$$

A) Probit selection

Model (type II Tobit):

$$\begin{aligned} y_1 &= \mathbf{X}_1\beta_1 + \epsilon_1 \\ y_2 &= \mathbf{1}[\mathbf{X}\delta_2 + \nu_2 > 0] \quad \text{where } \mathbf{X}_1 \in \mathbf{X} \end{aligned}$$

We know how to calculate:

$$\begin{aligned} E(y_1|\mathbf{X}) &= E(y_1|\mathbf{X}_1) = \mathbf{X}_1\beta_1 \\ P(y_2 = 1|\mathbf{X}) &= \Phi(\mathbf{X}\delta_2) \end{aligned}$$

A) Probit selection

Model (type II Tobit):

$$\begin{aligned} y_1 &= \mathbf{X}_1 \beta_1 + \epsilon_1 \\ y_2 &= \mathbf{1}[\mathbf{X} \delta_2 + \nu_2 > 0] \quad \text{where } \mathbf{X}_1 \in \mathbf{X} \end{aligned}$$

We know how to calculate:

$$\begin{aligned} E(y_1 | \mathbf{X}) &= E(y_1 | \mathbf{X}_1) = \mathbf{X}_1 \beta_1 \\ P(y_2 = 1 | \mathbf{X}) &= \Phi(\mathbf{X} \delta_2) \end{aligned}$$

How do we estimate:

$$E(y_1 | \mathbf{X}, y_2 = 1)$$

A) Probit selection

Model (type II Tobit):

$$\begin{aligned} \mathbf{y}_1 &= \mathbf{X}_1\boldsymbol{\beta}_1 + \boldsymbol{\epsilon}_1 \\ \mathbf{y}_2 &= \mathbf{1}[\mathbf{X}\boldsymbol{\delta}_2 + \boldsymbol{\nu}_2 > 0] \quad \text{where } \mathbf{X}_1 \in \mathbf{X} \end{aligned}$$

- How does $E(\mathbf{y}_1|\mathbf{X}, \mathbf{y}_2 = 1)$ depend on $\boldsymbol{\beta}_1$?
- Notice that: $E(\mathbf{y}_1|\mathbf{X}, \boldsymbol{\nu}_2) = \mathbf{X}_1\boldsymbol{\beta}_1 + \gamma_1\boldsymbol{\nu}_2$ (sample bias?)
- Two cases to consider:
 - 1 $\gamma_1 = 0$ and
 - 2 $\gamma_1 \neq 0$

A) Probit selection

Model (type II Tobit):

$$\begin{aligned} y_1 &= \mathbf{X}_1\beta_1 + \epsilon_1 \\ y_2 &= \mathbf{1}[\mathbf{X}\delta_2 + \nu_2 > 0] \quad \text{where } \mathbf{X}_1 \in \mathbf{X} \end{aligned}$$

Case 1:

- $\gamma_1 = 0$ (ϵ_1 and ν_2 are uncorrelated)

$$E(y_1|\mathbf{X}, \nu_2) = E(y_1|\mathbf{X}) = E(y_1|\mathbf{X}_1) = \mathbf{X}_1\beta_1$$

- Hence

$$E(y_1|\mathbf{X}, y_2) = E(y_1|\mathbf{X}_1) = \mathbf{X}_1\beta_1$$

- In this case there is no sample selection problem
- Estimate β consistently with OLS on the selected sample

A) Probit selection

Model (type II Tobit):

$$\begin{aligned} y_1 &= \mathbf{X}_1\beta_1 + \epsilon_1 \\ y_2 &= \mathbf{1}[\mathbf{X}\delta_2 + \nu_2 > 0] \quad \text{where } \mathbf{X}_1 \in \mathbf{X} \end{aligned}$$

Case 2:

- $\gamma_1 \neq 0$ (ϵ_1, ν_2 are correlated):

$$E(y_1|\mathbf{X}, y_2) = \mathbf{X}_1\beta_1 + \gamma_1 E(\nu_2|\mathbf{X}, y_2) = \mathbf{X}_1\beta_1 + \gamma_1 h(\mathbf{X}, y_2)$$

- Conditioning on the selection variable y_2 changes the expected value of $y_1|\mathbf{X}$
- Because y_2 and \mathbf{X} gives us info about ν_2 and hence ϵ_1
- If we knew $h()$ then we could use OLS on the selected sample and obtain consistent estimates of β_1 and γ_1
- Omitting $h()$ \Rightarrow biased or even inconsistent estimate of β_1

A) Probit selection

Case 2 (continued): The selected sample has $y_2 = 1$, $h(X, 1)$

$$\begin{aligned} h(X, 1) &= E(\nu_2 | \mathbf{X}, y_2 = 1) = E(\nu_2 | \mathbf{X}\delta_2 + \nu_2 > 0) \\ &= E(\nu_2 | \nu_2 > -\mathbf{X}\delta_2) \\ &= \lambda(\mathbf{X}\delta_2) = \frac{\phi(\mathbf{X}\delta_2)}{\Phi(\mathbf{X}\delta_2)} \end{aligned}$$

- This is the inverse Mill's ratio (see chapter 17)
- We then have:

$$E(y_1 | \mathbf{X}, y_2 = 1) = \mathbf{X}_1\beta_1 + \gamma_1\lambda(\mathbf{X}\delta_2)$$

- We know everything in this expression but $\delta_2 \Rightarrow$ we must estimate it
- Estimation with 2 step procedure.

A) Probit selection

Model (type II Tobit):

$$\begin{aligned} y_1 &= \mathbf{X}_1\beta_1 + \epsilon_1 \\ y_2 &= \mathbf{1}[\mathbf{X}\delta_2 + \nu_2 > 0] \quad \text{where } \mathbf{X}_1 \in \mathbf{X} \end{aligned}$$

Procedure 19.1 (Heckit):

- 1 Obtain probit estimate of δ_2 from: $P(y_{j2} = 1 | \mathbf{x}_j) = \Phi(\mathbf{x}_j\delta_2)$
 - The selection model using all observations
 - Construct λ from this estimate and \mathbf{X} : $\hat{\lambda}_j = \lambda(\mathbf{x}_j\hat{\delta}_2)$
- 2 With OLS:

$$y_{1j} = x_{1j}\beta_1 + \gamma_1\hat{\lambda}_j + error_j$$

Estimators $\hat{\beta}_1$ and $\hat{\gamma}_1$ are consistent and asymp. normal

A) Probit selection

NB!

- Test for sample selection bias: t-test for $H_0 : \gamma_1 = 0$ (no bias)
 - The asymptotic variance of β_1 and γ_1 are not affected by $\hat{\delta}_2$ under the null
- If the t-test shows that $\gamma_1 \neq 0$
 - Asymptotic variance of the estimator of β_1 is complicated.
 - The robust standard errors is not enough
- Preferably, we should have $\mathbf{X} \neq \mathbf{X}_1$ in the selection model
 - But in principle works without this because of non-linearity in selection equation

Exercise 3 (15 minutes)

Example 19.6 (Wage offer equation). Estimate the log wage for married women.

- mroz.dat
- $n = 753$, but only 428 women work (look at variable *lwage* and *inlf*)
- Model includes on RHS: *educ*, *exper*, *exper*²
- Selection equation includes further:
age, *kidslt6*, *kidsge6*, *nwifeinc*
- Write the selection equation $y_2 = ?$
- What is \mathbf{X} , \mathbf{X}_1 ?
- library *sampleSelection*
- function *heckit*

Exercise 3

In R,

heckit(selection, outcome, data, method = "2step")

or

selection(selection, outcome, data, method = "2step")

- *selection* is the **formula** of y_2 (variable defining the sample selection) and X
- *outcome* is the **formula** of y_1 (what we are interested on)

A) Probit selection – special case

Case 2 (special case) with stronger assumption: ϵ_1 and ν_2 are bivariate normal

$$\begin{pmatrix} \epsilon_1 \\ \nu_2 \end{pmatrix} \sim N \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & 1 \end{pmatrix} \right)$$

- We can then run partial maximum likelihood instead of two-step procedure.
- We cannot run full ML as \mathbf{y}_1 is not always observed
- Advantage of ML: More efficient, and we get standard errors directly from ML estimation
- In R: *heckit(selection, outcome, data)* or *selection(selection, outcome, data)*

A) Probit selection – an extension

Under the assumption that ϵ_1 and ν_2 are bivariate normal,

- The set-up can also be extended to binary response models:

$$\begin{aligned} y_1 &= \mathbf{1}[\mathbf{X}_1\beta_1 + \epsilon_1 > 0] \\ y_2 &= \mathbf{1}[\mathbf{X}\delta_2 + \nu_2 > 0] \quad \mathbf{X}_1 \in X \end{aligned}$$

- Example:
 - y_1 : Participation in on-the-job training
 - y_2 : Labour market participation
- Estimation method: Partial max-lik as above

Exercise 4 (3 minutes)

Discuss with your colleagues:

- When can't we just run OLS on the equation of interest when we have sample selection?
- Why?
- How can we fix it?
- In R?