

Sample selection

- In many economic problems sample we use for estimation is a nonrandomly *selected sample*
- We will analyze three selection mechanisms
 - selection based on selection indicator
 - selection based on response variable
 - selection based on separate selection model (incidental truncation)
- It is always important to correctly define the population in question

Example 1. (*Wooldridge*) *Saving function*

$$saving = \beta_0 + \beta_1 income + \beta_2 age + \beta_3 married + \beta_4 kids + u$$

but we only observe households for which $age > 45$

Example 2. (Wooldridge) *Estimate effect of eligibility in a pension plan on wealth.*

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

but we only have observations for people with $wealth < \$200,000$.

Example 3. (Wooldridge) *Wage offer function. By definition it should represent the (potential) wages of all people - not only the ones who are working. We only have data on wages for people who are actually working - the selection is based on another variable, labor force participation.*

Ignorability of selection

- Define selection indicator s ($s = 1$ if observation is in the sample)
- Selection problem can be ignored if in model

$$y = \beta_1 + \beta_2 x_2 + \dots + \beta_K x_K + \varepsilon$$

$$E(u | z, s) = 0$$

where z is the vector of instruments. The consistent estimate of β can be obtained with standard *2SLS*.

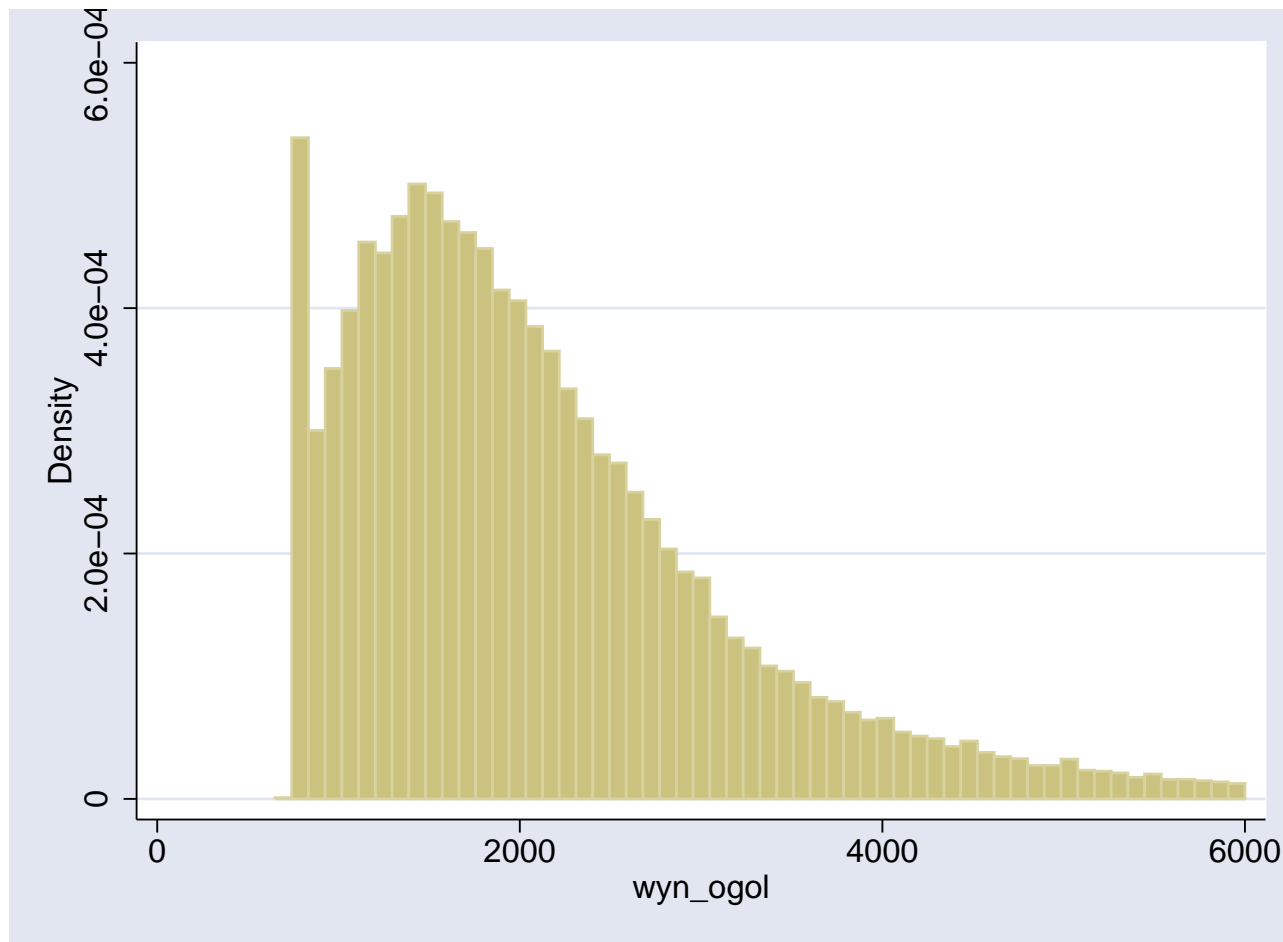
- **Interpretation:** selection can be ignored if it is possible to find instruments z such that given z selection is not correlated with random term u .

- If selection is independent on z, u or based on deterministic function of $s = f(z)$ then if $E(u|z) = 0$ then $E(u|z,s) = E(u|z) = 0$
- The simplest case: $z = x$ (Classical Regression) and selection based on x . In this case we can use *OLS* (Saving function example)

Selection based on response variable

- Hausman and Wise study (1977) - determinants of earnings, sample only for the people participating in negative income experiment, no data available for the people with income higher than same threshold
- Wealth example: we only observe the people whose wealth is smaller than \$200,000.

Example 4. *Gross wages in Poland, sample based on firm survey. .*



Only wages higher than minimal wage are declared!

- Difference with top coding - in the case of top coding we do not observe the response variable but we observe explanatory variable (see top coding for wealth example)
- In the case of sample selection (or truncated sample) - we observe nothing about the units not selected to the sample
- Assume that observation is selected to the sample if $s_i = 1$ and

$$s_i = 1 [a_1 < y_i < a_2]$$

- *cdf* of observed y_i is given by

$$\Pr(y_i \leq Y | \mathbf{x}_i, s_i = 1) = \frac{\Pr(y_i \leq Y, s_i = 1 | \mathbf{x}_i)}{\Pr(s_i = 1 | \mathbf{x}_i)}$$

where

- But $\Pr (s_i = 1 | \mathbf{x}_i) = \Pr (a_1 < y_i < a_2) = F (a_2 | \mathbf{x}_i; \boldsymbol{\beta}, \gamma) - F (a_1 | \mathbf{x}_i; \boldsymbol{\beta}, \gamma)$
- $\Pr (y_i \leq Y, s_i = 1 | \mathbf{x}_i) = \Pr (a_1 < y_i \leq Y | \mathbf{x}_i) = F (Y | \mathbf{x}_i; \boldsymbol{\beta}, \gamma) - F (a_1 | \mathbf{x}_i; \boldsymbol{\beta}, \gamma)$
- Density function is then

$$\frac{\partial \Pr (y_i \leq Y | \mathbf{x}_i, s_i = 1)}{\partial Y} = \frac{f (Y | \mathbf{x}_i; \boldsymbol{\beta}, \gamma)}{F (a_2 | \mathbf{x}_i; \boldsymbol{\beta}, \gamma) - F (a_1 | \mathbf{x}_i; \boldsymbol{\beta}, \gamma)}$$

- If we assume that the distribution $F (\cdot)$ is normal distribution we will obtain the *truncated tobit* (called too truncated normal regression)

Incidental truncation

- In this case we have problem of self-selection of the units in the sample
- Selection is based on decision of the unit partially based on unobserved characteristics

Example 5. (Wooldridge) *Labor force participation and the wage offer. We are interested in $E(w_i^o | x_i)$ and collect the sample for the people in working age. But only for the people actually working we obtain data on wage offer w_i^o . The utility maximization problem for individual is as follows:*

$$\max util(q, h) \quad \text{s.t. } 0 \leq h \leq 168$$

where h_i are the hours worked and $q_i = w_i^o h_i + a_i$ is income and a_i is nonlabor income. We assume that marginal utility w.r.t income is positive and with

respect to hours worked negative. Marginal utility of work is given by

$$\begin{aligned}\frac{\partial \text{util}(q, h_i)}{\partial h_i} &= \frac{\partial \text{util}(q_i, h_i)}{\partial q_i} \frac{\partial q_i}{\partial h_i} + \frac{\partial \text{util}(q, h_i)}{\partial h_i} \\ &= \frac{\partial \text{util}(q_i, h_i)}{\partial q_i} w_i^o + \frac{\partial \text{util}(q, h_i)}{\partial h_i}\end{aligned}$$

First order condition is given by equality

$$\frac{\partial \text{util}(q, h_i)}{\partial h_i} = 0$$

Individual will only work ($h_i > 0$) if for $h_i = 0$ his/her marginal utility of work is positive:

$$\frac{\partial \text{util}(a_i, 0)}{\partial q_i} w_i^o + \frac{\partial \text{util}(a_i, 0)}{\partial h_i} \geq 0$$

and so

$$w_i^o \geq - \frac{\partial \text{util}(a_i, 0)}{\partial h_i} / \frac{\partial \text{util}(a_i, 0)}{\partial q_i} = w_i^r$$

where w_i^r is called reservation wage.

If wage offer is given by $w_i^o = \exp(\mathbf{x}_{i1}\boldsymbol{\beta} + u_{i1})$ and reservation wage by $w_i^r = \exp(\mathbf{x}_{i2}\boldsymbol{\beta} + u_{i2})$ then we have the following two equations in the model:

1. wage equation

$$\log w_i^o = \mathbf{x}_{i1}\boldsymbol{\beta} + u_{i1}$$

2. participation equation

$$\log w_i^o - \log w_i^r = \mathbf{x}_{i1}\boldsymbol{\beta} + u_{i1} - \mathbf{x}_{i2}\boldsymbol{\beta} - u_{i2} = \mathbf{x}_i\boldsymbol{\delta}_2 + v_{i2}$$

If individual is participating ($\log w_i^o - \log w_i^r > 0$) we observe w_i^o . We do not observe w_i^r but only \mathbf{x}_i and the decision to participate or not.

Remarks:

- *as $\mathbf{x}_i = (\mathbf{x}_{i1}, \mathbf{x}_{i2})$ all the variables included in wage equation should also be included in participation equation*
- *as $v_{i2} = u_{i1} - u_{i2}$ we expect positive correlation between u_{i1} and v_{i2}*
- *Difference between the incidental truncation and top coding: we do not know the $\log w_i^r$, so we do not know exactly when the data is censored.*
- *$\log w_i^o - \log w_i^r = \mathbf{x}_i \boldsymbol{\delta}_2 + v_{i2}$ is a reduced form equation (wage is not included in this equation), the parameters of this equation can not be interpreted as parameters of the labor supply function.*

- The basic model

$$y_1 = x_1\beta_1 + u_1$$

$$y_2 = 1 [x\delta_2 + v_2 > 0]$$

- This model is called Heckit or Tobit type II model

- Assumptions

- (x, y_2) are always observed, y_1 is only observed when $y_2 = 1$
- (u_1, v_2) have zero mean and are independent of x
- $v_2 \sim Normal(0, 1)$
- $E(u_1 | v_2) = \gamma_1 v_2$ (which is true e.g. if u_1, v_2 are bivariate normal - in this case $\gamma_1 > 0$ if u_1 and v_2 are positively correlated)

- The conditional expected value of y_1 given v_2 is equal to

$$E(y_1 | \mathbf{x}, v_2) = \mathbf{x}_1 \boldsymbol{\beta}_1 + E(u_1 | v_2) = \mathbf{x}_1 \boldsymbol{\beta}_1 + \gamma_1 v_2$$

- If u_1 and v_2 are not correlated ($E(u_1 | v_2) = 0$ and $\gamma_1 = 0$) then $E(y_1 | \mathbf{x}, v_2) = \mathbf{x}_1 \boldsymbol{\beta}_1$ and as y_2 is a function of v_2 , $E(y_1 | \mathbf{x}, y_2) = E(y_1 | \mathbf{x}, v_2) = \mathbf{x}_1 \boldsymbol{\beta}_1$. In this case $\boldsymbol{\beta}$ can be consistently estimated with *OLS* (no sample selection bias - expected value of y_1 does not depend on unit being selected or not)

- However, if $\gamma_1 \neq 0$ then

$$E(y_1 | \mathbf{x}, y_2) = \mathbf{x}_1 \boldsymbol{\beta}_1 + E(u_1 | \mathbf{x}, y_2) = \mathbf{x}_1 \boldsymbol{\beta}_1 + \gamma_1 h(\mathbf{x}, y_2)$$

where $h(\mathbf{x}, y_2) = E(u_1 | \mathbf{x}, y_2)$.

- $h(\mathbf{x}, 1) = E(u_1 | v_2 > -\mathbf{x}\boldsymbol{\delta}_2) = \lambda(\mathbf{x}\boldsymbol{\delta}_2)$ where $\lambda(\cdot) = \frac{\phi(\cdot)}{\Phi(\cdot)}$ is inverse Mills ratio
- As we only observe y_1 for $y_2 = 1$ we are interested in conditional expectation of y_1 given that individual was selected:

$$\begin{aligned} E(y_1 | \mathbf{x}, y_2 = 1) &= \mathbf{x}_1\boldsymbol{\beta}_1 + \gamma_1 h(\mathbf{x}, 1) \\ &= \mathbf{x}_1\boldsymbol{\beta}_1 + \gamma_1 \lambda(\mathbf{x}\boldsymbol{\delta}_2) \end{aligned}$$

- Using that formula we can analyze the selection bias problem as omitted variable problem. Making *OLS* regression of y_1 on x_1 for selected units we omit $\lambda(\mathbf{x}\boldsymbol{\delta}_2)$ which is correlated with x_1 , which cause the omitted variable bias.
- This equation also suggest how to estimate consistently $\boldsymbol{\beta}_1$ with *OLS*: we should add to the explanatory variables element $\lambda(\mathbf{x}\boldsymbol{\delta}_2)$.

- δ_2 is unknown but it can be consistently estimated with probit.
- Two stage Heckman procedure
 1. Estimate probit regression of y_{i2} on x_i . Calculate $\hat{\lambda}_{i2} = \lambda(x_i \hat{\delta}_2)$
 2. Regress y_{i1} on x_{i1} and $\hat{\lambda}_{i2}$
- Test of the existence of the selection bias ($H_0 : \gamma_1 = 0$) can be based on standard t -statistics
- If $\gamma_1 \neq 0$, the variance matrix should be adjusted to take into account that $\hat{\delta}_2$ is estimated
- Even if $x = x_1$ this model is identified. However in this case identification relies on nonlinearity of $\lambda(\cdot)$. As $\lambda(\cdot)$ may be often well approximated with linear function of x_1 , in such situation we have severe collinearity

problem. Then in order to have precise estimate of β_1 we should have some explanatory variables which belongs to x but does not belong to x_1

Example 6. *(Wooldridge) Wage offer equation for married women. For 753 women in the sample, 428 are working. The labor force participation equation contains variables in wage equation plus income, age, number of young children, number of older children. Results*

	<i>OLS</i>	<i>Heckit</i>
<i>educ</i>	.108 (.014)	.109 (.016)
<i>exper</i>	.042 (.012)	.044 (.016)
<i>exper</i> ²	-.00081 (.00039)	-.00086 (.00044)
<i>constant</i>	-.522 (.199)	-.578 (.307)
$\hat{\lambda}_2$	—	.032 (.134)
<i>Sample size</i>	428	428
<i>R-squared</i>	.157	.157

Differences in estimates in this case are not great, inverse Mills ratio is not significant - no evidence for sample selection bias problem. If we estimate this heckit without exclusion restrictions the estimates become very imprecise

(e.g. *se* for *educ* becomes .119)

- Heckit model can also be estimated with *ML*. This procedure is efficient if u_1 and v_2 are bivariate normal.
- Sometimes we do not only observe the decision to participate but also y_2 (e.g. number of hours worked)
- In this case the selection equation has the form of the tobit

$$y_1 = \mathbf{x}_1\boldsymbol{\beta}_1 + u_1$$

$$y_2 = \max [0, \mathbf{x}\boldsymbol{\delta}_2 + v_2]$$

- This model is called type III tobit model

- Similarly as in heckit model

$$E(y_1 | \mathbf{x}, v_2, y_2 > 0) = \mathbf{x}_1 \boldsymbol{\beta} + \gamma_1 v_2$$

- But in this case v_2 can be consistently estimated as residuals from tobit model.
- Two step procedure:
 1. estimate tobit regression of y_{i2} on \mathbf{x}_i . Calculate $\hat{v}_{i2} = y_{i2} - \mathbf{x}_i \hat{\boldsymbol{\delta}}_2$
 2. estimate *OLS* regression of y_1 on \mathbf{x}_1 and \hat{v}_{i2}