#### **Microeconometrics**

Missing data and selection models

Alex Armand alex.armand@novasbe.pt

#### Lecture summary

- Introduction
- When can sample selection be ignored?
- Selection on the response variable: truncated regression
- Incidental truncation

### Example: Normal distribution (0,1)



### Example: Normal distribution censored



### Example: Normal distribution truncated



### Sample selection

Sample selection  $\Rightarrow$  the sample we obtain is not representative of the population of interest.

• Example: wealth equation for all families in a country

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

- $\textbf{0} \quad \textbf{Random sample} \rightarrow \textbf{OLS}$
- **2** Only people less than 65 years old were sampled  $\rightarrow$  selection on x

• What if we use OLS on the selected sample?

- $\textcircled{Only families with wealth greater than zero are sampled}{\rightarrow} selection on y$ 
  - What if we use OLS on the selected sample?





































## Types of sample selection

#### **9** Selection determined by the explanatory variables

- Missing completely at random
- Ø Missing at random
- Missing not at random

#### **2** Selection on the dependent variable

- Truncated regression
- 2 Incidental truncation: Heckman selection model

### Linear model with missing data

A population is represented by the random vector (x, y, z)

• Consider the identified IV case (generalize OLS):

$$y = x\beta + u$$
$$E(z'u) = 0$$

- Random sample
  - rank  $E(z'x) = K \Rightarrow 2SLS$  to consistently estimate  $\beta$
- Selected sample
  - Conditions are not usually enough to consistently estimate  $\beta$

#### Linear model with missing data

• Each observation *i* is supplemented by a selection indicator *s<sub>i</sub>* 

•  $s_i = 1 \Rightarrow$  observation *i* is USED in the estimation

- $s_i = 0 \Rightarrow$  observation *i* is NOT USED in the estimation
- Our sample consists of  $\{(x_i, y_i, z_i, s_i) : i = 1, ..., N\}$
- The IV estimator using the selected sample can be written as

$$\hat{\beta}_{IV} = \left( N^{-1} \sum_{i=1}^{N} s_i z'_i x_i \right)^{-1} \left( N^{-1} \sum_{i=1}^{N} s_i z'_i y_i \right)$$

•  $\hat{\beta}_{IV}$  is called a complete case estimator

#### Linear model with missing data

• With a large sample:

$$plim_{N\to\infty}(\hat{\beta}_{IV}) = \beta + [E(sz'x)]^{-1}E(sz'u)$$

- Assumptions for consistency:
  - Rank condition

rank 
$$E(z'x|s=1) = K$$

Exogeneity or orthogonality condition

$$E(sz'u) = P(s=1)E(z'x|s=1) = 0$$

• When do we achieve these conditions?

#### Missing completely at random (MCAR)

**MCAR**: *s* is independent of (x, y, z)

• Under MCAR

$$E(sz'u) = E(s)E(z'u) = \rho \cdot 0 = 0$$

• Selection does not affect identifying assumptions  $\Rightarrow$  OLS using the  $s_i = 1$  observations is consistent for  $\beta$ 



### Missing at random (MAR)

**MAR**:  $E(y|\mathbf{x}, s) = \mathbf{x}\beta$ 

• With MAR a sufficient condition is

$$E(y|\mathbf{x},s) = E(y|\mathbf{x}) = \mathbf{x}\beta$$

• s can be an arbitrary function of the exogenous variables.



## Missing not at random (MNAR)

MNAR: previous conditions are not valid

- Using MCAR and MAR assumption allow us using the complete case estimator
- Using OLS for the MNAR would lead to biased estimates



## Missing not at random (MNAR)

MNAR: previous conditions are not valid

- Using MCAR and MAR assumption allow us using the complete case estimator
- Using OLS for the MNAR would lead to biased estimates



## Missing not at random (MNAR)

MNAR: previous conditions are not valid

- Using MCAR and MAR assumption allow us using the complete case estimator
- Using OLS for the MNAR would lead to biased estimates



### Dealing with MNAR: s correlated with u?

Suppose the population model is

$$y = x\beta + u$$
  
 $E(u|x) = 0$   
 $Corr(s, u) \neq 0$ 

Suppose

- **1** s is a deterministic function of (x, v) for some variable v
- 2 (u, v) is independent of x
- The conditional mean of y is

$$E(y|\mathbf{x}, \mathbf{v}) = \mathbf{x}\beta + E(u|\mathbf{x}, \mathbf{v}) = \mathbf{x}\beta + E(u|\mathbf{v})$$

### Dealing with MNAR: s correlated with u?

We need to make additional assumptions to identify  $\beta$ !

- Suppose v has zero mean and  $E(u|v) = \gamma v$
- The conditional mean of y can be written as

$$E(y|\mathbf{x}, \mathbf{v}) = \mathbf{x}\beta + \gamma \mathbf{v}$$

- Use OLS of  $y_i$  on  $x_i$ ,  $v_i$  using the selected sample ( $s_i = 1$ ) to consistently estimate  $\beta$  and  $\gamma$ 
  - Notice that all variables only need to be observed when  $s_i = 1$
  - Intuition: selection is like an omitted variable

### Selection on y: truncated regression

The rule for observing a data point depends in a known deterministic way on the dependent variable

- Setting:
  - Random draw (x<sub>i</sub>, y<sub>i</sub>)
  - 2 We only observe the data point if  $s_i = 1$  for known constants  $a_1$  and  $a_2$

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

- Allow for the cases a<sub>1</sub> = −∞ and a<sub>2</sub> = +∞ to have truncation only in one side.
- OLS on the selected sample  $\Rightarrow \beta$  is inconsistent because selection is a function of  $y_i$

### Density conditional on s = 1

Assume that the population conditional density is  $f(y|x; \beta)$  where  $\gamma$  is another set of parameters.

• The density conditional on s = 1 is

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

 $\bullet~\mbox{MLE} \rightarrow \mbox{log-likelihood}$  function for the selected subpopulation

$$\sum_{i=1}^{N} \{\log[f(y_i|\mathsf{x}_i;\beta)] - \log[F(a_2|\mathsf{x}_i;\beta) - F(a_1|\mathsf{x}_i;\beta)]\}$$

### Truncated normal regression model

**Truncated normal regression (or Truncated Tobit) model** makes the following distributional assumption:

$$D(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{x}\beta,\sigma^2)$$

• As with censoring, truncating the sample is costly

- We are interested in  $E(y|x) = x\beta$  in the entire population
- We need to specify all of D(y|x)!
- Differs from the censored normal regression model
  - No information on units not in the subpopulation with  $a_1 < y < a_2$

#### **Example:** $a_1 = -\infty$

Compare the likelihood in the truncation versus censoring cases

Truncated case

$$\left\{\frac{\sigma^{-1}\phi[(y_i - x_i\beta)/\sigma]}{\Phi[(a_2 - x_i\beta)/\sigma]}\right\}^{s}$$

- Completely drop all units with  $s_i = 0$
- Censored case

$$\{\sigma^{-1}\phi[(\mathbf{y}_i-\mathbf{x}_i\beta)/\sigma]\}^{\mathbf{s}_i}\{1-\Phi[(\mathbf{a}_2-\mathbf{x}_i\beta)/\sigma]\}^{1-\mathbf{s}_i}$$

- Uses additional information from the binary selection indicator
- If you have a choice, you should use censored regression

### Incidental truncation: self-selection

Sample selection is **not a deterministic function** of x or y, but it may be related to them.

- Self-selection: *y* is observed only when a certain event is true.
  - The event is often a choice of the individual we observe.
- Example:
  - $y = \log(w^{o})$ , where  $w^{o}$  is the "wage offer"
  - A person works if the wage offer is larger than reservation wage  $(w_i^r)$

$$w_i^o > w_i^r$$

- We observe wage<sup>o</sup> only if the person decides to enter the work force.
- We do not observe *wage<sup>o</sup>* otherwise.

#### Wage offers and incidental truncation

Assume we can model the wage offer and reservation wages as

$$w_i^o = \exp(x_{i1}\beta_1 + u_{i1})$$
  
 $w_i^r = \exp(x_{i2}\beta_2 + \gamma_2 a_i + u_{i2})$ 

• Observe 
$$w_i^o$$
 if  $\log(w_i^o) - \log(w_i^r) > 0$  or  
 $x_{i1}\beta_1 + u_{i1} - x_{i2}\beta_2 - \gamma_2 a_i - u_{i2} > 0$  or

$$x_i\delta_2 + v_{i2} > 0$$

where  $x_i$  includes all nonredundant elements of  $x_{i1}$  and  $x_{i2}$  and  $a_i$ 

• Simplest example of a structural (econometrics) model

## General model: Type II Tobit Model

A general population model:

$$y_1 = x_1\beta_1 + u_1 y_2 = 1[x\delta_2 + v_2 > 0]$$

- y<sub>1</sub> is the **response variable** (only partially observed)
- y<sub>2</sub> is the selection indicator (what we called s before)

#### Assumptions

- **(** $(x, y_2)$ ) are always observed,  $y_1$  is observed only when  $y_2 = 1$
- 2  $(u_1, v_2)$  is independent of x with  $u_1$  having mean zero
- 3  $v_2 \sim \mathcal{N}(0,1)$
- $(u_1|v_2) = \gamma_1 v_2$

#### An estimating equation for $\beta_1$

Under the previous assumptions

$$E(y_1|x, v_2) = x_1\beta_1 + E(u_1|x, v_2) = x_1\beta_1 + E(u_1|v_2) = x_1\beta_1 + \gamma_1v_2$$

• **Problem**: we observe  $y_2$ , not  $v_2$ !

$$E(y_1|x, y_2) = E[E(y_1|x, v_2)|x, y_2] = x_1\beta_1 + \gamma_1 E(v_2|x, y_2)$$

• How can we recover 
$$E(v_2|x, y_2)$$
?

#### Generalized residual: $E(v_2|x, y_2)$

• Decompose the term using the two values of  $y_2$ 

$$E(v_2|x, y_2) = y_2 E(v_2|x, y_2 = 1) + (1 - y_2) E(v_2|x, y_2 = 0)$$

• For the selected sample, notice that

$$\begin{array}{lll} {\it E}(v_2|{\sf x},y_2=1) & = & {\it E}(v_2|{\sf x},v_2>-{\sf x}\delta_2) \\ & = & \frac{\phi(-{\sf x}\delta_2)}{1-\Phi(-{\sf x}\delta_2)} \end{array}$$

 <u>Reminder</u>: step 2 comes from this property of Normal distribution (same rule as tobit):

$$E(z|a < z < b) = \mu + \frac{\phi(\frac{a-\mu}{\sigma}) - \phi(\frac{b-\mu}{\sigma})}{\Phi(\frac{b-\mu}{\sigma}) - \Phi(\frac{a-\mu}{\sigma})}\sigma$$

### Heckit method (Heckman 1976)

If we restrict to the selected sample  $(y_2 = 1)$  we have

$$E(y_1|x, y_2 = 1) = x_1\beta_1 + \gamma_1 E(v_2|x, y_2 = 1) = x_1\beta_1 + \gamma_1 \lambda(x\delta_2)$$

- Two-step estimation method
  - Selection equation: probit of y<sub>i2</sub> on x<sub>i</sub> using all of the data

$$\hat{\lambda}_{i2} = \lambda(\mathsf{x}_i\hat{\delta}_2)$$

**2** Regression equation: OLS of  $y_{i1}$  on  $x_{i1}$ ,  $\hat{\lambda}_{i2}$  in the selected sample

- $H_0: \gamma_1 = 0$  tests for no sample selection problem
- Notice we don't need an exclusion restriction (like in IV)

#### APPLICATION: wage offer for married women

- . use mroz
- . des lwage inlf nwifeinc

variable name	storage type	display format	y value label	variable label
lwage	float	%9.0g		log(wage)
inlf	byte	%9.0g		=1 if in lab frce, 1975
nwifeinc	float	%9.0g		(faminc - wage*hours)/1000
. sum lwage i	nlf educ	kidslt6	nwifeinc	

Max	Min	Std. Dev.	Mean	Obs	Variable
3.218876	-2.054164	.7231978	1.190173	428	lwage
1	0	.4956295	.5683931	753	inlf
17	5	2.280246	12.28685	753	educ
3	0	.523959	.2377158	753	kidslt6
96	0290575	11.6348	20.12896	753	nwifeinc

### APPLICATION: wage equation using OLS

$$log(w_i) = \beta_0 + \beta_1 \cdot educ_i + \beta_2 \cdot exper_i + \beta_3 \cdot expersq_i + u$$

. reg lwage educ exper expersq

	Source	SS	df	MS		Number of obs	=	428
	+-					F(3, 424)	= 2	6.29
	Model	35.0222967	3	11.6740989		Prob > F	= 0.	0000
	Residual	188.305144	424	.444115906		R-squared	= 0.	1568
	+-					Adj R-squared	= 0.	1509
	Total	223.327441	427	.523015084		Root MSE	= .6	6642
-	lwage	Coef.	Std. E	t	P> t	[95% Conf.	Inter	val]
-	lwage   +- educ	Coef. 	Std. E .01414	crr. t 65 7.60	P> t  0.000	[95% Conf. .0796837	Inter 	 val]  2956
-	lwage    educ   exper	Coef. .1074896 .0415665	Std. E .01414 .01317	frr. t 65 7.60 52 3.15	P> t  0.000 0.002	[95% Conf. .0796837 .0156697	Inter .135 .067	val]  2956 4633
	lwage   educ   exper   expersg	Coef. .1074896 .0415665 0008112	Std. E .01414 .01317 .00039	rr. t 65 7.60 52 3.15 32 -2.06	P> t  0.000 0.002 0.040	[95% Conf. .0796837 .0156697 0015841	Inter .135 .067	val]  2956 4633 0382
	lwage   educ   exper   expersq   _cons	Coef. .1074896 .0415665 0008112 5220406	Std. E .01414 .01317 .00039 .19863	rr. t 65 7.60 52 3.15 32 -2.06 21 -2.63	P> t  0.000 0.002 0.040 0.009	[95% Conf. .0796837 .0156697 0015841 9124667	Inter .135 .067 000 131	val] 2956 4633 0382 6144

### **APPLICATION: Heckit regression equation**

Heckman se	election	model	two-step	estimates	Number of	obs	=	753
(regressio	on model	with samp	ole select	ion)	Censored	obs	=	325
					Uncensore	d obs	=	428

Wald	chi2(6)	=	180.10
Prob	> chi2	=	0.0000

 Coef.
 Std. Err.
 z
 P>|z|
 [95% Conf. Interval]

 lwage
 |

 educ
 .1090655
 .015523
 7.03
 0.000
 .0786411
 .13949

 exper
 .0438873
 .0162611
 2.70
 0.007
 .0120163
 .0757584

 expersq
 -.0008591
 .0004389
 -1.96
 0.050
 -.0017194
 1.15e-06

 \_cons
 -.5781032
 .3050062
 -1.90
 0.058
 -1.175904
 .019698

\_\_\_\_\_\_

# **APPLICATION: Heckit selection equation**

inlf							
educ	1	.1309047	.0252542	5.18	0.000	.0814074	.180402
exper	1	.1233476	.0187164	6.59	0.000	.0866641	.1600311
expersq	1	0018871	.0006	-3.15	0.002	003063	0007111
nwifeinc	1	0120237	.0048398	-2.48	0.013	0215096	0025378
age		0528527	.0084772	-6.23	0.000	0694678	0362376
kidslt6		8683285	.1185223	-7.33	0.000	-1.100628	636029
kidsge6		.036005	.0434768	0.83	0.408	049208	.1212179
_cons		.2700768	.508593	0.53	0.595	7267473	1.266901
mills							
lambda	1	.0322619	.1336246	0.24	0.809	2296376	.2941613
rho	-+-	0.04861					
sigma	L	.66362875					
lambda	L	.03226186	.1336246				

# **APPLICATION: Heckit selection equation**

inlf	1						
educ	1	.1309047	.0252542	5.18	0.000	.0814074	.180402
exper	1	.1233476	.0187164	6.59	0.000	.0866641	.1600311
expersq	1	0018871	.0006	-3.15	0.002	003063	0007111
nwifeinc	1	0120237	.0048398	-2.48	0.013	0215096	0025378
age	1	0528527	.0084772	-6.23	0.000	0694678	0362376
kidslt6	1	8683285	.1185223	-7.33	0.000	-1.100628	636029
kidsge6	1	.036005	.0434768	0.83	0.408	049208	.1212179
_cons	1	.2700768	.508593	0.53	0.595	7267473	1.266901
mills	ī						
lambda	1	.0322619	.1336246	0.24	0.809	2296376	.2941613
rho	1	0.04861		Coeffic	ient on t	the selection	n term
sigma	T.	.66362875		in the s	econd s	stage	
lambda	1	.03226186	.1336246				

# **APPLICATION: Heckit selection equation**

\_\_\_\_\_

inlf	Ι.						
educ	L	.1309047	.0252542	5.18	0.000	.0814074	.180402
exper	L	.1233476	.0187164	6.59	0.000	.0866641	.1600311
expersq	L	0018871	.0006	-3.15	0.002	003063	0007111
nwifeinc	L	0120237	.0048398	-2.48	0.013	0215096	0025378
age	L	0528527	.0084772	-6.23	0.000	0694678	0362376
kidslt6	L	8683285	.1185223	-7.33	0.000	-1.100628	636029
kidsge6	L	.036005	.0434768	0.83	0.408	049208	.1212179
_cons	1	.2700768	.508593	0.53	0.595	7267473	1.266901
mills							
lambda	L	.0322619	.1336246	0.24	0.809	2296376	.2941613
rho	+- 	0.04861		Coeffic	ient on t	the selection	n term
sigma	1	.66362875		in the s	econd s	tage is not	significant
lambda	L	.03226186	.1336246	in the s	ocona c	lage to not	- ginnount