

# LECTURE 7

## Introduction to Econometrics

### Endogeneity

Dali Laxton

December 2, 2022

## A LITTLE REVISION: OLS CLASSICAL ASSUMPTIONS

1. Linearity: the regression model is linear in the parameters (coefficients)
2. Random sampling: the data is a random sample drawn from the population and each data point follows the population equation
3. No perfect collinearity: the values of explanatory variables are not all the same and no explanatory variable is a perfect linear function of any other explanatory variable(s)
4. Zero conditional mean: values of explanatory variables must contain no information about the mean of the unobserved factors - explanatory variables are uncorrelated with the error term
5. Homoskedasticity: the error term has a constant variance
6. Normality of the error term: the error term is normally distributed

## ON PREVIOUS LECTURES

- ▶ We discussed what happens if some of the assumptions are violated
- ▶ Linearity of coefficients and no perfect multicollinearity are essential for the definition of OLS estimator
- ▶ Zero mean of the error term is always ensured by the inclusion of intercept
- ▶ Normality of the error term is needed for statistical inference, but it can be shown that if the number of observations is sufficiently high, the OLS estimate will have asymptotically normal distribution even if the stochastic error term is not normal
- ▶ Heteroskedasticity leads to incorrect statistical inference, but we have studied tests to detect it and techniques to overcome this problem

## ON TODAY'S LECTURE

- ▶ The assumption of no correlation between explanatory variables and the error term is crucial
- ▶ Variables that are correlated with the error term are called *endogenous variables* (as opposed to *exogenous variables*)
- ▶ We will show that the estimated coefficients of endogenous variables are inconsistent and biased
- ▶ We will explain in which situations we may encounter endogenous variables
- ▶ We will define the concept of instrumental variables
- ▶ We will derive the 2SLS technique to deal with endogeneity

# ENDOGENOUS VARIABLES

▶ Notation:  $Cov(x_i, \varepsilon_i) \neq 0$  or  $E[X'\varepsilon] \neq 0$

▶ Intuition behind the bias:

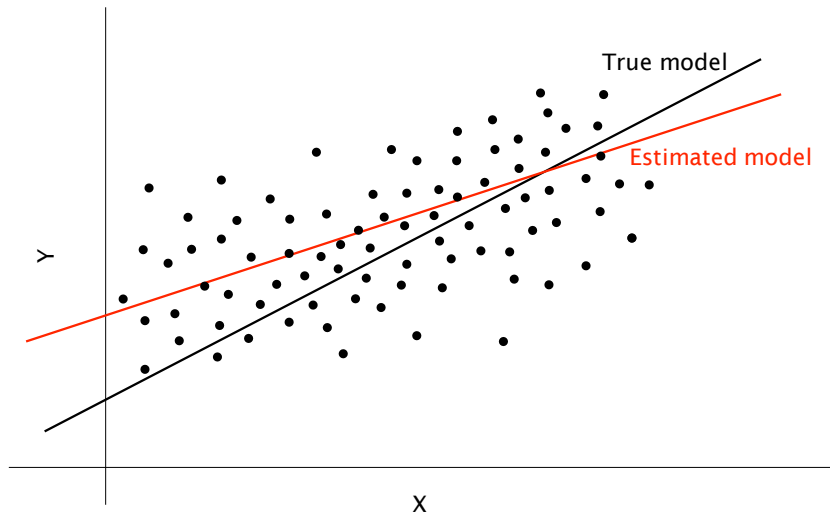
If an explanatory variable  $x$  and the error term  $\varepsilon$  are correlated with each other, the OLS estimate attributes to  $x$  some of the variation in  $y$  that actually came from the error term  $\varepsilon$

▶ Example: Analysis of household consumption patterns

Households with lower income may indicate higher consumption (because of shame)

▶ Leads to inconsistent estimates

# GRAPHICAL REPRESENTATION



## INCONSISTENCY OF ESTIMATES

- ▶ We can express

$$\begin{aligned}\widehat{\beta} &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} = \beta + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\varepsilon \\ &= \beta + \left(\frac{1}{n}\mathbf{X}'\mathbf{X}\right)^{-1} \frac{1}{n}\mathbf{X}'\varepsilon\end{aligned}$$

- ▶ We assume that there exists a finite matrix  $\mathbf{Q}$  so that

$$\frac{1}{n}\mathbf{X}'\mathbf{X} \xrightarrow{n \rightarrow \infty} \mathbf{Q}$$

- ▶ It can be shown that  $\frac{1}{n}\mathbf{X}'\varepsilon \xrightarrow{n \rightarrow \infty} E[\mathbf{X}'\varepsilon] \stackrel{\text{endogeneity}}{\neq} \mathbf{0}$

- ▶ This implies:

$$\widehat{\beta} \xrightarrow{n \rightarrow \infty} \beta + \mathbf{Q}^{-1} \cdot E[\mathbf{X}'\varepsilon] = \beta + \textit{bias}$$

# TYPICAL CASES OF ENDOGENEITY

## 1. Omitted variable bias

An explanatory variable is omitted from the equation and makes part of the error term

## 2. Selection bias

An unobservable characteristic has influence on both dependent and explanatory variables

## 3. Simultaneity

The causal relationship between the dependent variable and the explanatory variable goes in both directions

## 4. Measurement error

Some of the variables are measured with error

- ▶ In all 4 cases, the sign of the bias is given by the sign of  $Cov(\epsilon_i, x_i)$



## OMITTED VARIABLE BIAS

- ▶ Studied on lecture 5
- ▶ True model:  $y_i = \beta x_i + \gamma z_i + u_i$
- ▶ Model as it looks when we omit variable  $z$ :

$$y_i = \beta x_i + \tilde{u}_i \quad \text{implying} \quad \tilde{u}_i = \gamma z_i + u_i$$

- ▶ This gives

$$\text{Cov}(\tilde{u}_i, x_i) = \text{Cov}(\gamma z_i + u_i, x_i) = \gamma \text{Cov}(z_i, x_i) \neq 0$$

- ▶ It can be remedied by including the variable in question, but sometimes we do not have data for it
- ▶ We can include some proxies for such variable, but this may not reduce the bias completely and some endogeneity remains in the equation

## SELECTION BIAS

- ▶ **Very similar to omitted variable bias**
- ▶ **We suppose there is some unobservable characteristic that influences both the level of the dependent variable  $y$  and of the explanatory variable  $x$**
- ▶ **This unobservable characteristic forms part of the error term  $\varepsilon$ , causing  $Cov(\varepsilon, x) \neq 0$  (in the same manner as an omitted variable)**
- ▶ **Example: surveying only non-smoking mothers when inferring the impact of the number of prenatal visits on the birth weight of children.**

# SIMULTANEITY

- ▶ Occurs in models where variables are jointly determined

$$y_{1i} = \alpha_0 + \alpha_1 y_{2i} + \varepsilon_{1i}$$

$$y_{2i} = \beta_0 + \beta_1 y_{1i} + \varepsilon_{2i}$$

- ▶ Intuitively: change in  $y_{1i}$  will cause a change in  $y_{2i}$ , which will in turn cause  $y_{1i}$  to change again
- ▶ Technically:

$$\begin{aligned} \text{Cov}(\varepsilon_{1i}, y_{2i}) &= \text{Cov}(\varepsilon_{1i}, \beta_0 + \beta_1 y_{1i} + \varepsilon_{2i}) \\ &= \beta_1 \text{Cov}(\varepsilon_{1i}, y_{1i}) \\ &= \beta_1 \text{Cov}(\varepsilon_{1i}, \alpha_0 + \alpha_1 y_{2i} + \varepsilon_{1i}) \\ &= \beta_1 (\alpha_1 \text{Cov}(\varepsilon_{1i}, y_{2i}) + \text{Var}(\varepsilon_{1i})) \end{aligned}$$

$$\text{Cov}(\varepsilon_{1i}, y_{2i}) = \frac{\beta_1}{1 - \alpha_1 \beta_1} \text{Var}(\varepsilon_{1i}) \neq 0$$

# SIMULTANEITY

- ▶ Example:

$$Q_{Di} = \alpha_0 + \alpha_1 P_i + \alpha_2 I_i + \varepsilon_{1i}$$

$$Q_{Si} = \beta_0 + \beta_1 P_i + \varepsilon_{2i}$$

$$Q_{Di} = Q_{Si}$$

where

$Q_D$  ... quantity demanded

$Q_S$  ... quantity supplied

$P$  ... price

$I$  ... income

- ▶ Endogeneity of price: it is determined from the interaction of supply and demand

## MEASUREMENT ERROR I

- ▶ **Measurement error in the dependent variable**
- ▶ **Measurement error is correlated with an explanatory variable**

$$y_i^* = y_i + v_i \quad \text{where} \quad \text{Cov}(v_i, x_i) \neq 0$$

- ▶ True regression model:  $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$
- ▶ Estimated regression:  $y_i^* = \beta_0 + \beta_1 x_i + u_i$  where  
 $u_i = \varepsilon_i + v_i$  and so

$$\text{Cov}(x_i, u_i) = \text{Cov}(x_i, \varepsilon_i + v_i) = \text{Cov}(v_i, x_i) \neq 0$$

- ▶ **Example: analysis of household consumption patterns (above)**

## MEASUREMENT ERROR II

- ▶ Classical measurement error in the explanatory variable

$$x_i^* = x_i + v_i \quad \text{where} \quad \text{Cov}(v_i, x_i) = 0$$

- ▶ True regression model:  $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$

- ▶ **Estimated regression:**  $y_i = \beta_0 + \beta_1 x_i^* + u_i$  where

$$u_i = \varepsilon_i - \beta_1 v_i \text{ and so}$$

$$\text{Cov}(x_i^*, u_i) = \text{Cov}(x_i + v_i, \varepsilon_i - \beta_1 v_i) = -\beta_1 \text{Var}(v_i) \neq 0$$

- ▶ **Causes attenuation bias (estimated coefficient is smaller in absolute value than the true one)**

## INSTRUMENTAL VARIABLES (IV)

- ▶ **Answer to the situation when  $Cov(x, \varepsilon) \neq 0$**
- ▶ **Instrumental variable (or instrument) should be a variable  $z$  such that**
  - 1.  $z$  is uncorrelated with the error term:  $Cov(z, \varepsilon) = 0$**
  - 2.  $z$  is correlated with the explanatory variable  $x$ :  $Cov(x, z) \neq 0$**
- ▶ **Intuition behind instrumental variables approach:**
  - project the endogenous variable  $x$  on the instrument  $z$**
  - this projection is uncorrelated with the error term and can be used as an explanatory variable instead of  $x$**

# INSTRUMENTAL VARIABLES

- ▶ Suppose the equation we want to estimate is:

$$y = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\eta}$$

We can have several instruments for several endogenous variables - we will use the matrix notation  $\mathbf{Z}$  and  $\mathbf{X}$

$\mathbf{X}$  denotes endogenous variable(s)

$\mathbf{Z}$  denotes instrumental variable(s)

Assume that we have at least as many instruments as endogenous variables



## TWO STAGE LEAST SQUARES

- ▶ 2SLS is a method of implementing instrumental variables approach
- ▶ Consists of two steps:
  1. Regress the endogenous variables on the instruments

$$\mathbf{X} = \mathbf{Z}\boldsymbol{\delta} + \mathbf{v} ,$$

get predicted values

$$\hat{\mathbf{X}} = \mathbf{Z}\hat{\boldsymbol{\delta}} = \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{X} ,$$

2. Use these predicted values instead of  $\mathbf{X}$  in the original equation:

$$\mathbf{y} = \hat{\mathbf{X}}\boldsymbol{\beta} + \boldsymbol{\eta}$$

## TWO STAGE LEAST SQUARES

- ▶ The estimate is

$$\begin{aligned}\hat{\beta}^{2SLS} &= (\hat{X}'\hat{X})^{-1} \hat{X}'y \\ &= (X'Z (Z'Z)^{-1} Z'X)^{-1} X'Z (Z'Z)^{-1} Z'y\end{aligned}$$

- ▶ This estimate is consistent, but it has higher variance than OLS (it is not efficient)
- ▶ Intuitively:

Only part of the variation in  $X$  that is uncorrelated with the error term is used for the estimation.

This ensures consistency ( $\hat{X}$  that is uncorrelated with error term).

**But it makes the estimate less precise (higher variance of  $\hat{\beta}$  because not all variation in  $X$  is used.**

## EXAMPLE

- ▶ Estimating the impact of education on the number of children for a sample of women in Botswana
- ▶ OLS:

Source	SS	df	MS	Number of obs = 4361		
Model	12243.0295	3	4081.00985	F( 3, 4357) =	1915.20	
Residual	9284.14679	4357	2.13085765	Prob > F	= 0.0000	
Total	21527.1763	4360	4.93742577	R-squared	= 0.5687	
				Adj R-squared	= 0.5684	
				Root MSE	= 1.4597	

children	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
educ	-.0905755	.0059207	-15.30	0.000	-.102183	-.0789679
age	.3324486	.0165495	20.09	0.000	.3000032	.364894
agesq	-.0026308	.0002726	-9.65	0.000	-.0031652	-.0020964
_cons	-4.138307	.2405942	-17.20	0.000	-4.609994	-3.66662

## EXAMPLE

- ▶ Education may be endogenous - both education and number of children may be influenced by some unobserved socioeconomic factors

Omitted variable bias: family background is an unobserved factor that influences both the number of children and years of education

- ▶ Finding possible instrument:

Something that explains education

But is not correlated with the family background

- ▶ A dummy variable

$$frsthalf = \begin{cases} 1 & \text{if the woman was born in the first} \\ & \text{six months of a year} \\ 0 & \text{otherwise} \end{cases}$$

## EXAMPLE

- ▶ Intuition behind the instrument:
- ▶ The first condition - instrument explains education:
  - School year in Botswana starts in January
    - ⇒ Thus, women born in the first half of the year start school when they are at least six and a half.
  - Schooling is compulsory till the age of 15
    - ⇒ Thus, women born in the first half of the year get less education if they leave school at the age of 15.
- ▶ The second condition - instrument is uncorrelated with the error term:
  - Being born in the first half of the year is uncorrelated with the unobserved socioeconomic factors that influence education and number of children (family background etc.)

# EXAMPLE

## First-stage regressions

---

Number of obs = 4361  
F( 3, 4357) = 175.21  
Prob > F = 0.0000  
R-squared = 0.1077  
Adj R-squared = 0.1070  
Root MSE = 3.7110

educ	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
age	-.1079504	.0420402	-2.57	0.010	-.1903706	-.0255302
agesq	-.0005056	.0006929	-0.73	0.466	-.0018641	.0008529
frsthalf	-.8522854	.1128296	-7.55	0.000	-1.073489	-.6310821
_cons	9.692864	.5980686	16.21	0.000	8.520346	10.86538

# EXAMPLE

Instrumental variables (2SLS) regression

Number of obs = **4361**  
Wald chi2(3) = **5300.22**  
Prob > chi2 = **0.0000**  
R-squared = **0.5502**  
Root MSE = **1.49**

children	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
educ	<b>-.1714989</b>	<b>.0531553</b>	<b>-3.23</b>	<b>0.001</b>	<b>-.2756813 - .0673165</b>
age	<b>.3236052</b>	<b>.0178514</b>	<b>18.13</b>	<b>0.000</b>	<b>.2886171 .3585934</b>
agesq	<b>-.0026723</b>	<b>.0002796</b>	<b>-9.56</b>	<b>0.000</b>	<b>-.0032202 -.0021244</b>
_cons	<b>-3.387805</b>	<b>.5478988</b>	<b>-6.18</b>	<b>0.000</b>	<b>-4.461667 -2.313943</b>

Instrumented: educ

Instruments: age agesq frsthalf

## EXAMPLE

- ▶ Compare the estimates from OLS and 2SLS:
- ▶ OLS:

children	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
educ	<b>-.0905755</b>	<b>.0059207</b>	<b>-15.30</b>	<b>0.000</b>	<b>-.102183 - .0789679</b>

- ▶ 2SLS:

children	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
educ	<b>-.1714989</b>	<b>.0531553</b>	<b>-3.23</b>	<b>0.001</b>	<b>-.2756813 - .0673165</b>

- ▶ Is the bias reduced by IV?
- ▶ Are these results statistically different?



# SUMMARY

- ▶ We showed that the estimated coefficients of endogenous variables are inconsistent and biased
- ▶ In which situations we may encounter endogenous variables
  - Omitted variable (omitting important variable which is correlated to independent variable)
  - Selection bias (unobserved factors influencing both dependent and independent variable)
  - Simultaneity (causality goes both ways)
  - Measurement error (in either dependent or independent variable)
- ▶ We can deal with endogeneity by using instrumental variables (2SLS technique)