

# Econometrics II

Fabian Waldinger (LMU Munich)

# Lecture Structure

- ① Recap from last lecture
- ② Violation of GM3: Measurement Error
- ③ Violation of GM3: Simultaneity

# Recap from Last Lecture

- Dummy Variables: dummy variable trap (example of violation of GM2)
- Violation of GM3: omitted variable bias:  
For a model with two  $X$  variables, one of the excluded from the model, we can derive the omitted variable bias as:

$$E(\tilde{\beta}_2|\mathbf{X}) = \beta_2 + \beta_3 \frac{\text{Cov}(x_2, x_3)}{\text{Var}(x_2)}$$

there is bias if  $\beta_3 \neq 0$  and  $\text{Cov}(x_2, x_3) \neq 0$

The sign of the bias depends on the signs of  $\beta_3$  and  $\text{Cov}(x_2, x_3)$

# Measurement Error in $X$

- Measurement error in one or more  $x$  variables can also lead to a violation of GM3 and, hence, biased OLS estimates
- Here we study measurement error in the simple regression model
- Suppose the true model is:

$$y = \beta_1 + \beta_2 Z + v$$

- But  $Z$  is subject to measurement error  $w$ , we denote the measured explanatory variable  $X$

$$X = Z + w$$

- We also assume that  $E(w) = 0$  (classical measurement error)

# Measurement Error in $X$

- If we substitute the true model for  $Z$  we get:

$$y = \beta_1 + \beta_2 Z + v$$

$$= \beta_1 + \beta_2(X - w) + v$$

$$= \beta_1 + \beta_2 X + u$$

- where:  $u = v - \beta_2 w$
- Hence, the explanatory variable  $X$  is correlated with the error term  $u$

# Measurement Error in X

- To demonstrate that the OLS estimator is inconsistent we start with the OLS estimator

$$\hat{\beta}_2 = \frac{\sum(X_i - \bar{X})(y_i - \bar{y})}{\sum(X_i - \bar{X})^2}$$

- Now substitute the model for  $y$ :

$$\hat{\beta}_2 = \frac{\sum(X_i - \bar{X})([\beta_1 + \beta_2 X_i + u_i] - [\beta_1 + \beta_2 \bar{X} + \bar{u}])}{\sum(X_i - \bar{X})^2}$$

$$\hat{\beta}_2 = \frac{\sum(X_i - \bar{X})(\beta_2[X_i - \bar{X}] + u_i - \bar{u})}{\sum(X_i - \bar{X})^2}$$

$$\hat{\beta}_2 = \beta_2 + \frac{\sum(X_i - \bar{X})(u_i - \bar{u})}{\sum(X_i - \bar{X})^2}$$

# Is OLS Consistent?

- We would like to show whether  $\hat{\beta}_2$  is biased. Which would mean taking expectations of  $\hat{\beta}_2$
- Here both  $X$  and  $u$  depend on  $w$  and hence both denominator and numerator are functions of  $w \rightarrow$  we can therefore not simplify the expression above using expectations
- We therefore investigate the *plim* of this expression (i.e. what happens to  $\hat{\beta}_2$  if the sample size  $\rightarrow \infty$ )

$$plim \hat{\beta}_2 = \beta_2 + plim \left( \frac{\sum (X_i - \bar{X})(u_i - \bar{u})}{\sum (X_i - \bar{X})^2} \right)$$

- The *plim* of the last expression does not exist (the denominator increases indefinitely as the sample size increases, the numerator has no particular limit)

# Is OLS Consistent?

- We therefore divide both denominator and numerator by  $N$

$$\begin{aligned} plim \hat{\beta}_2 &= \beta_2 + plim \left( \frac{\frac{1}{N} \sum (X_i - \bar{X})(u_i - \bar{u})}{\frac{1}{N} \sum (X_i - \bar{X})^2} \right) \\ &= \beta_2 + \frac{plim[\frac{1}{N} \sum (X_i - \bar{X})(u_i - \bar{u})]}{plim[\frac{1}{N} \sum (X_i - \bar{X})^2]} \\ &= \beta_2 + \frac{Cov(X, u)}{Var(X)} \end{aligned}$$



# Is OLS Consistent?

- We can now use variance and covariance rules to simplify the second part of this expression:
- Numerator:

$$\begin{aligned} \text{Cov}(X, u) &= \text{Cov}([Z + w], [v - \beta_2 w]) \\ &= \text{Cov}(Z, v) + \text{Cov}(w, v) + \text{Cov}(Z, -\beta_2 w) + \text{Cov}(w, -\beta_2 w) \\ &= 0 + 0 + 0 - \beta_2 \sigma_w^2 \end{aligned}$$

- Denominator:

$$\begin{aligned} \text{Var}(X) &= \text{Var}(Z + w) \\ &= \text{Var}(Z) + \text{Var}(w) + 2\text{Cov}(Z, w) \\ &= \sigma_Z^2 + \sigma_w^2 + 0 \end{aligned}$$

# Is OLS Consistent?

- Hence the plim of  $\hat{\beta}_2$  is:

$$\begin{aligned} \text{plim} \hat{\beta}_2 &= \beta_2 + \frac{-\beta_2 \sigma_w^2}{\sigma_Z^2 + \sigma_w^2} = \beta_2 - \beta_2 \frac{\sigma_w^2}{\sigma_Z^2 + \sigma_w^2} \\ &= \beta_2 \left( 1 - \frac{\sigma_w^2}{\sigma_Z^2 + \sigma_w^2} \right) = \beta_2 \left( \frac{\sigma_Z^2}{\sigma_Z^2 + \sigma_w^2} \right) \end{aligned}$$

- The last term in parentheses is  $< 1$
- Thus in large samples  $\hat{\beta}_2$  is biased towards 0 (attenuation bias) and the size of the bias depends on the relative sizes of  $\sigma_Z^2$  and  $\sigma_w^2$
- The standard errors will also be biased (derivation beyond the level of this course)
- If the measurement error is not classical (i.e.  $E(w) \neq 0$ ) OLS will also be biased
- If the model contains more than one variable all coefficients will be inconsistent, even if only one variable is measured with error (the variable measured with error will still be attenuated, the sign of the other large sample biases depend on the correlations of the  $X$ s

# Simulation Measurement Error

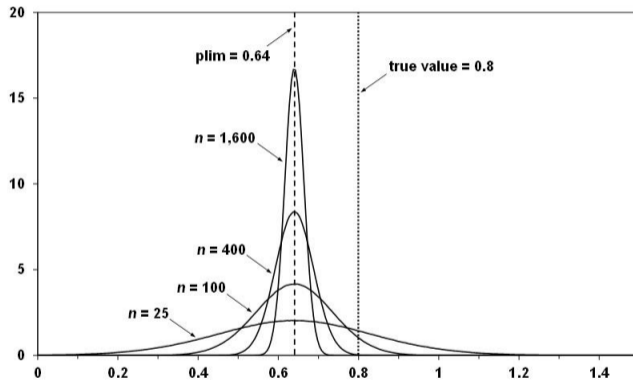
- Chris Dougherty (LSE) simulates how measurement error affects estimates:
- True model:

$$Y = 2.0 + 0.8Z + u$$

- With the values of  $Y$  drawn from a normal distribution with mean 10 and variance 4
- He then creates  $X$  where  $X = Z + w$  and  $w$  is drawn from a normal distribution with mean 0 and variance 1
- What would be the plim of  $\hat{\beta}_2$  if we regress of  $Y$  on  $X$ ?

$$\text{plim}\hat{\beta}_2 = 0.8 - 0.8 \frac{1}{4 + 1} = 0.64$$

# Example Measurement Error



# Simultaneity

- Another leading cause for a violation of GM3 is simultaneity, i.e.
  - $x$  causes  $y$  but also
  - $y$  causes  $x$
- One example is the effect of police on crime:
  - more police on the roads reduces crime
  - more crime often increases the number of policemen on the road
- Another example is the effect of institutions on growth:
  - Better institutions are good for growth
  - Higher growth allows countries to build better institutions

# Simultaneity Bias in OLS

- Consider the two-equation structural model

$$y_1 = \alpha_1 y_2 + \beta_1 z_1 + u_1 \quad (1)$$

$$y_2 = \alpha_2 y_1 + \beta_2 z_2 + u_2 \quad (2)$$

- for simplicity we suppress the intercept in each equation
- The variables  $z_1$  and  $z_2$  are exogenous, so that each is uncorrelated with  $u_1$  and  $u_2$
- Suppose we are interested in estimating the effect of  $y_2$  on  $y_1$

# Simultaneity Bias in OLS

- We obtain the reduced form equation of  $y_2$  (the equation that only depends on exogenous variables) by plugging (1) into (2):

$$y_2 = \alpha_2 y_1 + \beta_2 z_2 + u_2$$

$$= \alpha_2 [\alpha_1 y_2 + \beta_1 z_1 + u_1] + \beta_2 z_2 + u_2$$

- This can be rewritten as:

$$(1 - \alpha_2 \alpha_1) y_2 = \alpha_2 \beta_1 z_1 + \beta_2 z_2 + u_2 + \alpha_2 u_1$$

- If  $\alpha_2 \alpha_1 \neq 1$  we can rewrite this as:

$$y_2 = \frac{\alpha_2 \beta_1}{(1 - \alpha_2 \alpha_1)} z_1 + \frac{\beta_2}{(1 - \alpha_2 \alpha_1)} z_2 + \frac{u_2 + \alpha_2 u_1}{(1 - \alpha_2 \alpha_1)}$$

$$= \pi_{12} z_1 + \pi_{22} z_2 + \nu_2$$

- The parameters  $\pi_{12}$  and  $\pi_{22}$  are called the reduced form parameters, they are nonlinear functions of the structural parameters which appear in the structural equations (1) and (2):
  - $\pi_{12} = \frac{\alpha_2\beta_1}{(1-\alpha_2\alpha_1)}$
  - $\pi_{22} = \frac{\beta_2}{(1-\alpha_2\alpha_1)}$
- The reduced form error  $\nu_2$  is a linear function of the structural error terms  $u_1$  and  $u_2$ :
  - $\nu_2 = \frac{u_2 + \alpha_2 u_1}{(1-\alpha_2\alpha_1)}$
- Thus  $y_2$  is correlated both with  $u_1$  and  $u_2$



# Simultaneity Leads to Violation of GM3

- Hence if we were to estimate equation (1) by OLS:

$$y_1 = \alpha_1 y_2 + \beta_1 z_1 + u_1 \quad (1)$$

- One of the regressors,  $y_2$ , would be correlated with the error term  $u_1$
- This would be a violation of GM3 and hence OLS would be biased
- If you obtain the reduced form equation for  $y_1$  the error term is also a linear function of the structural error terms  $u_1$  and  $u_2$  and hence estimating equation (2) by OLS would also violate GM3

# Simultaneity Bias

- What is the sign of the simultaneity bias in this case?
- As before, we start with the OLS estimator for equation (1)
- To simplify the derivation, we omit  $z_1$  and add a constant to equation to both equations:

$$y_1 = \gamma_1 + \alpha_1 y_2 + u_1$$

$$y_2 = \gamma_2 + \alpha_2 y_1 + \beta_2 z_2 + u_2$$

- In that case, the reduced form equation for  $y_2$  is (make sure that you can derive the reduced form equation at home):

$$y_2 = \frac{\gamma_2 + \alpha_2 \gamma_1}{(1 - \alpha_2 \alpha_1)} + \frac{\beta_2}{(1 - \alpha_2 \alpha_1)} z_2 + \frac{u_2 + \alpha_2 u_1}{(1 - \alpha_2 \alpha_1)}$$

- What is the OLS estimator for  $\alpha_1$  in that case?

$$\hat{\alpha}_1 = \frac{\sum (y_{2i} - \bar{y}_2)(y_{1i} - \bar{y}_1)}{\sum (y_{2i} - \bar{y}_2)^2}$$

# Derivation of Simultaneity Bias

- Plugging in the true model for  $y_1$ :

$$\begin{aligned}\hat{\alpha}_1 &= \frac{\sum(y_{2i} - \bar{y}_2)[(\gamma_1 + \alpha_1 y_{2i} + u_{1i}) - (\gamma_1 + \alpha_1 \bar{y}_2 + \bar{u}_1)]}{\sum(y_{2i} - \bar{y}_2)^2} \\ &= \frac{\sum(y_{2i} - \bar{y}_2)\alpha_1(y_{2i} - \bar{y}_2) + \sum(y_{2i} - \bar{y}_2)(u_{1i} - \bar{u}_1)}{\sum(y_{2i} - \bar{y}_2)^2} \\ &= \alpha_1 + \frac{\sum(y_{2i} - \bar{y}_2)(u_{1i} - \bar{u}_1)}{\sum(y_{2i} - \bar{y}_2)^2}\end{aligned}$$

- Here both  $y_2$  and  $u_1$  are functions of  $u_1$  (see reduced form equation)  $\rightarrow$  we can therefore not simplify the expression above using expectations
- We therefore investigate the *plim* of this expression (i.e. what happens to  $\hat{\alpha}_1$  if the sample size  $\rightarrow \infty$ )

# Is OLS Consistent?

$$plim \hat{\alpha}_1 = \alpha_1 + plim \left( \frac{\sum (y_{2i} - \bar{y}_2)(u_{1i} - \bar{u}_1)}{\sum (y_{2i} - \bar{y}_2)^2} \right)$$

- As before, we multiply both numerator and denominator by  $\frac{1}{N}$

$$plim \hat{\alpha}_1 = \alpha_1 + plim \left( \frac{\frac{1}{N} \sum (y_{2i} - \bar{y}_2)(u_{1i} - \bar{u}_1)}{\frac{1}{N} \sum (y_{2i} - \bar{y}_2)^2} \right)$$

- This allows us to simplify as follows (if  $plim(\frac{1}{N} \sum (y_{2i} - \bar{y}_2)^2) \neq 0$ )

$$\begin{aligned} plim \hat{\alpha}_1 &= \alpha_1 + \frac{plim \frac{1}{N} \sum (y_{2i} - \bar{y}_2)(u_{1i} - \bar{u}_1)}{plim \frac{1}{N} \sum (y_{2i} - \bar{y}_2)^2} \\ &= \alpha_1 + \frac{Cov(y_2, u_1)}{Var(y_2)} \end{aligned}$$

# Sign of Large Sample Bias

- To determine the sign of this bias we need to determine the sign of  $Cov(y_2, u_1)$  (because the sign of  $Var(y_2)$  is always  $> 0$ )
- We therefore plug in the reduced form for  $y_2$ :

$$\begin{aligned}Cov(y_2, u_1) &= Cov\left(\left[\frac{\gamma_2 + \alpha_2\gamma_1}{(1 - \alpha_2\alpha_1)} + \frac{\beta_2}{(1 - \alpha_2\alpha_1)}z_2 + \frac{u_2 + \alpha_2u_1}{(1 - \alpha_2\alpha_1)}\right], u_1\right) \\&= \frac{1}{(1 - \alpha_2\alpha_1)} [Cov(\gamma_2 + \alpha_2\gamma_1 + \beta_2z_2 + u_2 + \alpha_2u_1, u_1)] \\&= \frac{1}{(1 - \alpha_2\alpha_1)} [Cov(\gamma_2, u_1) + Cov(\alpha_2\gamma_1, u_1) + Cov(\beta_2z_2, u_1) \\&\quad + Cov(u_2, u_1) + Cov(\alpha_2u_1, u_1)] \\&= \frac{1}{(1 - \alpha_2\alpha_1)} [0 + 0 + 0 + 0 + \alpha_2\sigma_{u_1}^2] = \frac{\alpha_2\sigma_{u_1}^2}{(1 - \alpha_2\alpha_1)}\end{aligned}$$

- The sign of the bias therefore depends on the sign of  $\alpha_2$  and whether  $\alpha_2\alpha_1$  is smaller or greater than 1

# How Do we Overcome Violations of GM3?

- How can we overcome violations of GM3?
- Solutions for improving basic OLS regressions
  - add control variables (overcomes omitted variable bias)
  - get variables without measurement error
- Design identification strategies that directly address the violation of GM3
  - ① Randomize the variable of interest
  - ② Find quasi-random variation
- In the next few lectures we will study such approaches in a lot of detail
- The notation will be similar across these approaches and follows the notation that was introduced by Joachim Winter and that is also in Mostly Harmless Econometrics (Angrist and Pischke)

# Introductory Example and Notation

- Suppose we wanted to learn about the causal effect of health insurance on health outcomes
- Why would be a simple comparison of people with and without health insurance be problematic?