

Misspecification.

Omitted variable bias

Problem (UoL Exam).. (a) Explain what you understand by omitted variable bias using regression model without intercept.

Consider two equations

$$y_i = \beta_1 x_{1i} + \beta_2 x_{2i} + u_i$$

(1) TRUE.

$$y_i = \beta_1 x_{1i} + v_i$$

(2) EST.

(b) Let (1) be a **true model** while (2) be a **false model**. A researcher, using ordinary least squares (OLS), estimates β_1 from the **false model**. Examine the properties of OLS estimator of β_1 in equation (2).

biased, inconsistent

(c) Let now (1) be a **false model** while (2) be a **true model**. A researcher, using ordinary least squares (OLS), estimates β_1 from the **false model** (1). Examine the properties for this OLS estimate of β_1 .

$$b) \text{ in mod (2)} \quad \hat{\beta}_1 = \frac{\sum x_{1i} y_i}{\sum x_{1i}^2} = \frac{\sum x_{1i} (\beta_1 x_{1i} + \beta_2 x_{2i} + u_i)}{\sum x_{1i}^2} =$$

$$= \beta_1 \frac{\sum x_{1i}^2}{\sum x_{1i}^2} + \beta_2 \cdot \frac{\sum x_{1i} \cdot x_{2i}}{\sum x_{1i}^2} + \frac{\sum x_{1i} \cdot u_i}{\sum x_{1i}^2}$$

(*) Deterministic Regressors

$$E(\hat{\beta}_1) = \beta_1 + \underbrace{\beta_2 \cdot \frac{\sum x_{1i} \cdot x_{2i}}{\sum x_{1i}^2}}_{= \text{bias}} + \frac{\sum x_{1i} E(u_i)}{\sum x_{1i}^2} = 0$$

→ sign of bias $\beta_2 > 0$; $\sum x_{1i} x_{2i} > 0 \Rightarrow \text{bias} > 0$

→ no bias $\beta_2 = 0$; $X_1^\top X_2 = \sum x_{1i} \cdot x_{2i} = 0$
when x_1, x_2 orthogonal

$$c) \quad \begin{array}{l} y_i = \beta_1 x_{1i} + \beta_2 x_{2i} + u_i \\ y_i = \beta_1 x_{1i} + v_i \end{array} \quad \begin{array}{l} (1) \text{ EST} \\ (2) \text{ TRUE} \end{array}$$

$$\hat{\beta}_1$$

unbiased, consistent, ineffective

Problem (ICEF Exam). □ Explain why omitting explanatory variable from the regression equation can lead to the violation of Gauss-Markov conditions so standard errors and all tests become invalid.

□ What are consequences of omitting explanatory variable from the regression equation for the estimation regression equations.

$$(True) \quad Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + u_i$$

· X_3 - relevant variable

$$\cdot \text{Cov}(X_2, X_3) \neq 0$$

$$Y_i = \beta_1 + \beta_2 \cdot X_{2i} + u_i^* \quad u_i^* = \beta_3 \cdot X_{3i} + u_i$$

Stoch. Regressors TGM assumptions :

$$\oplus \quad E(X_i u_i) = 0$$

↳ X_i - exogenous

otherwise

X_i - endogenous

Problem (Dourgerty textbook 1st edition). A social scientist thinks that the level of activity in the shadow economy, Y_t , depends either positively on the level of the tax burden, X_t , or negatively on the level of government expenditure to discourage shadow economy activity, Z_t . The value Y_t may also depend on the X_t and Z_t simultaneously. There are annual time series data for 20 years, where the value of Y_t , X_t and Z_t are measured in the same units. Sociologist builds regression dependence (1): Y_t only on the value of X_t , (2): Y_t only on the value of Z_t and (3): Y_t from both variables X_t and Z_t , in relation to each city, with the following results (in parentheses are standard errors). Having carried out the appropriate statistical tests, write a short report advising the social scientist how to interpret these results.

	Constant	Estimated coefficients		R^2
		X_t	Z_t	
City A				
X	315.7 (18.5)	1.54 n.s. (0.97)	-	0.12
2	128.6 (50.9)	-	-0.96 * (0.06)	0.94
3	218.0 (76.6)	2.85 * (0.25)	-1.21 * (0.03)	0.99
City B				
1	197.6 (16.8)	2.86 * (0.25)	-	0.88
X	512.2 (202.6)	-	-0.05 n.s. (0.08)	0.02
3	230.8 (82.5)	2.94 * (0.27)	-0.01 (n.s.) (0.03)	0.88

City A : model (3)

X_t - relevant

$\Rightarrow \hat{\beta}_2$ - unbiased

City B : model (1), since Z_t - irrelevant

Question 4 (UoL Exam). Explain the RESET test as a general test for functional form misspecification and

discuss the drawbacks and advantages of this test. In your answer consider the following multiple linear regression model:

$$y_i = \beta_1 + \beta_2 x_{2i} + \beta_3 x_{3i} + u_i, \quad i = 1, \dots, n,$$

where x_{2i} and x_{3i} are exogenous variables known to affect $E(y_i)$.

$$1) \quad y_i | x_2, x_3 \Rightarrow \hat{y}_i \quad \quad \quad x_2^2 = x_4$$

$$2) \quad y_i = \beta_1 + \beta_2 x_{2i} + \beta_3 x_{3i} + \beta_4 \hat{y}_i^2 + (\beta_5 \hat{y}_i^3 \dots) + u_i$$

H_0 : correctly spec. $H_0: \beta_4 = 0$

Disadvantage: does not indicate correct form

Advantage: do not lose d.o.f

no problem of multicollinearity

Question 5 (ICEF Exam).

A researcher has data on output per worker, y , and capital per worker, k , both measured in thousands of pounds, for 50 firms in the manufacturing sector of the U.K. for 2016. She hypothesizes that output per worker depends on capital per worker and perhaps also the technological sophistication of the firm, $tech$:

$$y = \beta_1 + \beta_2 k + \beta_3 tech + u$$

where u is a disturbance term. She is unable to measure $tech$ and decides to use expenditure per worker on research and development in 1998, exp , as a proxy for it.

(a) ☐ What do you mean by good or poor proxy?

Proxy: $cov(exp, tech) \neq 0$

$$cov(exp, u) = 0$$

☐ Explain the consequences of using exp as a proxy for $tech$ if it is a good proxy.

☐ Explain the consequences of using exp as a proxy for $tech$ if it is a poor proxy.

Goal: minimize bias for endogenous var.

$$tech = \lambda + \mu \cdot exp + v$$

$$y = \beta_1 + \beta_2 \cdot k + \beta_3 (\lambda + \mu \cdot exp + v) + u$$

$$= (\beta_1 + \beta_3 \lambda) + \underbrace{\beta_2}_{\hat{\beta}_2} \cdot k + \underbrace{\beta_3 \cdot \mu}_{\text{bias}} \cdot exp + \underbrace{\beta_3 \cdot v + u}_{u^*}$$

The researcher fits the following regressions (standard errors in parentheses):

$$(1) \quad \hat{y} = 1.02 + 0.32k \quad R^2 = 0.749$$

(0.45) (0.04) ↓

$$(2) \quad \hat{y} = 0.34 + 0.29k + 0.05 exp \quad R^2 = 0.750$$

(0.61) (0.22) (0.15)

bias > 0

The correlation coefficient for k and exp was 0.92.

(b) ☐ Discuss these regression results assuming that y does depend on both k and $tech$.

exp - poor proxy

☐ Discuss these regression results assuming that y depends only on k .

$$y = \alpha + \beta_1 k + \beta_2 \cdot tech + u$$

$$E(\hat{y}_1) = \beta_1 + \beta_2 \cdot \frac{cov(k, tech)}{var(k)}$$

" bias

$$\text{bias: } \beta_2 > 0 \quad cov(k, tech) > 0 \quad \text{bias} > 0$$

