

Lecture 12

①

Specification Analysis

The consequences of

- 1) Omitting a Relevant Variable
- 2) Including an Irrelevant Variable

1) Assume the correct (true) model is:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + u_i \quad (1)$$

but instead you estimate the model

$$y_i = \beta_0 + \hat{\beta}_1 x_{i1} + u_i. \quad (2)$$

Now

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1) y_i}{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2}$$

(2)

Therefore

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1)(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + u_i)}{SST_1}$$

$$\text{where } SST_1 = \sum_{i=1}^n (x_{i1} - \bar{x}_1)^2$$

Then

$$\begin{aligned}\hat{\beta}_1 &= \beta_0 \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1)}{SST_1} + \beta_1 \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1) x_{i1}}{SST_1} \\ &\quad + \beta_2 \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1) x_{i2}}{SST_1} + \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1) u_i}{SST_1} \\ &= \beta_0 \cdot \frac{0}{SST_1} + \beta_1 \frac{SST_1}{SST_1} + \beta_2 \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1) x_{i2}}{SST_1} \\ &\quad + \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1) u_i}{SST_1}\end{aligned}$$

Now taking the expectation we have

(3)

$$E(\hat{\beta}_1) = \beta_1 + \beta_2 \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1) x_{i2}}{SST_1}$$

$$+ \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1) E u_i}{SST_1}$$

$$= \beta_1 + \beta_2 \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1) x_{i2}}{SST_1}$$

$$= \beta_1 + \beta_2 \tilde{\delta}_1$$

where $\tilde{\delta}_1$ is the slope coefficient estimate obtained by applying ordinary least squares to the auxiliary equation

$$x_2 = \delta_0 + \delta_1 x_1 + v.$$

The term $\beta_2 \tilde{\delta}_1$ is called the omitted variable bias.

(4)

Consider the Following table which tells us the bias of using ordinary least squares to estimate β_1 in the misspecified equation (2) when in fact X_2 should have been included (see eq. (1)).

Table 3.2 in Wooldridge

	$\text{cov}(X_1, X_2) > 0$	$\text{cov}(X_1, X_2) < 0$
$\beta_2 > 0$	positive bias	negative bias
$\beta_2 < 0$	negative bias	positive bias

Now consider the case where the true model is

$$y_i = \beta_0 + \beta_1 x_{i1} + u_i . \quad (3)$$

However, suppose that you instead estimate the following model that includes the irrelevant variable x_2

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + u_i . \quad (4)$$

Let ~~$\hat{\beta}_1$~~ be the ordinary least squares estimator of β_1 in equation (3) and $\tilde{\beta}_1$ and $\tilde{\beta}_2$ be, respectively, the ordinary least squares estimators of β_1 and β_2 in equation (4).

Then, it can be shown that

(6)

$$E(\tilde{\beta}_2) = 0 \quad \therefore \text{unbiased}$$

$$E(\tilde{\beta}_1) = \beta_1 \quad \therefore \text{unbiased}$$

but

$$\text{Var}(\tilde{\beta}_1) \geq \text{Var}(\hat{\beta}_1)$$

and the ordinary least squares estimator of β_1 obtained from the overspecified equation is inefficient. The inefficiency

of $\tilde{\beta}_1$ is given by the equation

$$\frac{\text{Var}(\tilde{\beta}_1)}{\text{Var}(\hat{\beta}_1)} = \frac{1}{1 - r_{12}^2}, \quad 0 \leq r_{12}^2 \leq 1 \\ -1 \leq r_{12} \leq 1$$

where r_{12} is the sample correlation between X_1 and X_2 , namely

(7)

$$r_{12} = \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2)}{\sqrt{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2 \sum_{i=1}^n (x_{i2} - \bar{x}_2)^2}}.$$

Note that $\text{Var}(\tilde{\beta}_1) = \text{Var}(\hat{\beta}_1)$ only

when $r_{12} = 0$ and X_1 and X_2 are orthogonal to each other.