

Specification Bias

$$E(UU') = \sigma^2 I$$

(homoscedasticity)

Digression \rightarrow Matrix approach to Linear Regression Model:

$$Y_1 = \beta_0 + \beta_1 X_{11} + \beta_2 X_{21} + \beta_3 X_{31} + \dots + \beta_k X_{k1} + U_1$$

$$Y_2 = \beta_0 + \beta_1 X_{12} + \beta_2 X_{22} + \beta_3 X_{32} + \dots + \beta_k X_{k2} + U_2$$

$$Y_3 = \dots$$

$$Y_n = \beta_0 + \beta_1 X_{1n} + \beta_2 X_{2n} + \beta_3 X_{3n} + \dots + \beta_k X_{kn} + U_n$$

Putting this equation in matrix form,

$$Y = X\beta + U$$

where

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} \quad X = \begin{bmatrix} 1 & X_{11} & \dots & X_{k1} \\ 1 & X_{12} & \dots & X_{k2} \\ \vdots & \vdots & \dots & \vdots \\ 1 & X_{1n} & \dots & X_{kn} \end{bmatrix}$$

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} \quad \text{and} \quad U = \begin{bmatrix} U_1 \\ U_2 \\ \vdots \\ U_n \end{bmatrix}$$

Assumptions:

① $E(U) = 0$

② $U = \begin{bmatrix} U_1 \\ U_2 \\ \vdots \\ U_n \end{bmatrix}$

$$U' = [U_1 \ U_2 \ \dots \ U_n]$$

$$\dots \dots \dots U_1 U_2 \dots U_1 U_n$$

$\downarrow U_n \downarrow$

$$UU' = \begin{bmatrix} \textcircled{U_1^2} & U_1 U_2 & U_1 U_3 & \dots & U_1 U_n \\ U_2 U_1 & \textcircled{U_2^2} & \dots & \dots & U_2 U_n \\ \vdots & \vdots & \textcircled{0} & \dots & \vdots \\ U_n U_1 & U_n U_2 & \dots & \dots & \textcircled{U_n^2} \end{bmatrix}$$

$$E(UU') = \begin{bmatrix} E(U_1^2) & E(U_1 U_2) & E(U_1 U_3) & \dots & E(U_1 U_n) \\ E(U_2 U_1) & E(U_2^2) & E(U_2 U_3) & \dots & E(U_2 U_n) \\ \vdots & \vdots & \vdots & \dots & \vdots \\ E(U_n U_1) & \dots & \dots & \dots & E(U_n^2) \end{bmatrix}$$

$$= \begin{bmatrix} \textcircled{\sigma^2} & \text{Cov}(U_1, U_2) & \text{Cov}(U_1, U_3) & \dots & \text{Cov}(U_1, U_n) \\ \text{Cov} & \textcircled{\sigma^2} & \text{Cov} & \dots & \text{Cov} \\ \vdots & \vdots & \textcircled{\sigma^2} & \dots & \vdots \\ \text{Cov} & \text{Cov} & \text{Cov} & \dots & \textcircled{\sigma^2} \end{bmatrix}$$

CRM: $E(U^2) = \sigma^2$
 $E(U_i U_j) = 0$

$$= \begin{bmatrix} \sigma^2 & 0 & 0 & 0 & \dots & 0 \\ 0 & \sigma^2 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \dots & \vdots \\ 0 & 0 & 0 & 0 & \dots & \sigma^2 \end{bmatrix}$$

$$= \sigma^2 \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix}$$

$E(UU') = \sigma^2 I_n$

* ——— *

Specification error: Case of Explanatory Variable

The use of ordinary least squares when some variables are left out may introduce bias into the estimates. Bias ... and specification bias.

The use of ~~omitted~~ variables left out may introduce bias into the estimates. This bias that originates in this way is called specification bias.

For instance, the true function explaining variation in y is given as:

$$y = \beta_1 x_1 + \beta_2 x_2 + u$$

However either due to ignorance of true relation or because of non-availability of data on x_2 following regression equation is estimated:

$$y = \beta_1^* x_1 + v$$

It can be shown that β_1^* is different from β_1 .

On applying OLS to $y = \beta_1^* x_1 + v$

we obtain $\beta_1^* = \frac{\sum x_1 y}{\sum x_1^2}$.

On the other hand, the normal equations of the true function: $y = \beta_1 x_1 + \beta_2 x_2 + u$ are

$$\begin{aligned} \sum x_1 y &= \beta_1 \sum x_1^2 + \beta_2 \sum x_1 x_2 && \text{--- divide by } \sum x_1^2 \\ \sum x_2 y &= \beta_1 \sum x_1 x_2 + \beta_2 \sum x_2^2 \end{aligned}$$

$$\frac{\sum x_1 y}{\sum x_1^2} = \beta_1 + \beta_2 \left(\frac{\sum x_1 x_2}{\sum x_1^2} \right)$$

Since $\beta_1^* = \frac{\sum x_1 y}{\sum x_1^2} \Rightarrow \beta_1^* = \beta_1 + \beta_2 \left(\frac{\sum x_1 x_2}{\sum x_1^2} \right)$

$\beta_1^* = \beta_1$ only if $\sum x_1 x_2 = 0$

$$\beta_1^* = \beta_1 \text{ only if } \frac{\sum x_1 x_2}{\sum x_1^2} = 0$$

$$\text{if } \sum x_1 x_2 = 0$$

$$\therefore \text{Specification error} = \underbrace{(\beta_1^* - \beta_1)}_{\downarrow} = \beta_2 \left(\frac{\sum x_1 x_2}{\sum x_1^2} \right)$$

It is hence proved that β_1^* from incorrect specification (through OLS procedure) is a biased estimate of parameter β_1 .

$$\text{The bias, which is equal to } \left[\beta_2 \left\{ \frac{\sum x_1 x_2}{\sum x_1^2} \right\} \right]$$

depends on two terms, namely, the regression coefficient of omitted variable in the true relation (β_2) and the covariance of the omitted variable with the included variable $\frac{\sum x_1 x_2}{\sum x_1^2}$.

Variance of β_1^*

$$\text{for } y = \beta_1^* x_1 + v$$

$$\text{var}(\beta_1^*) = \frac{\sigma_v^2}{\sum x_1^2} = \frac{\sum e_i^2 / (n-2)}{\sum x_1^2}$$

$$\text{var}(\beta_1^*) = \frac{\sum (y - \beta_1^* x_1)^2}{\sum x_1^2}$$

$$\text{var}(\beta_1^*) = \frac{\sum (y - \beta_1^* x_1)^2}{(n-2) \sum x_1^2}$$

$$y = \beta_1 x_1 + \beta_2 x_2 + u \quad (\text{true equation}) \Rightarrow \text{var}(\beta_1^*) = \frac{\sum (\beta_1 x_1 + \beta_2 x_2 - \beta_1^* x_1)^2}{(n-2) \sum x_1^2}$$

$$= \frac{\sum [-(\beta_1^* - \beta_1)x_1 + \beta_2 x_2]^2}{(n-2) \sum x_1^2}$$

$$= \frac{\sum \{(\beta_1^* - \beta_1)x_1\}^2}{(n-2) \sum x_1^2} + \frac{\sum \beta_2^2 x_2^2}{(n-2) \sum x_1^2}$$

$$= \text{var}(\beta_1^*) + \frac{\sum \beta_2^2 x_2^2}{n-2 \sum x_1^2}$$

This implies that estimator of β_1^* is positively biased.

Therefore, the usual tests of significance concerning β_1 shall be invalid in present circumstances.

The estimator of the constant intercept of the incorrectly specified function turns out to be biased.

$$\begin{aligned}
 E(\beta_0^*) &= E(\bar{y} - \beta_1 \bar{x}_1) \\
 &= E(\beta_0 + \beta_1 \bar{x}_1 + \beta_2 \bar{x}_2 - \beta_1 \bar{x}_1) \\
 &= E(\beta_0 + \beta_2 \bar{x}_2) \\
 E(\beta_0^*) &= \beta_0 + \beta_2 \bar{x}_2
 \end{aligned}$$

β_0^* will be unbiased if $\bar{x}_2 = 0$. ↳ biasedness

- The above discussion can now be summarised as follows:
- (i) If the omitted or left out explanatory variable is correlated with the included explanatory variable, the OLS estimator of β_1 will be biased and inconsistent.
 - (ii) If the omitted variable is not correlated with the included variable, the estimator of constant term will still be biased and inconsistent, but estimator of β_1 will be unbiased.
 - (iii) The variance of β_1 will contain an upward bias. Therefore, the test on significance of the estimator would not lead to correct conclusions.