

## Biases due to Omitted Variables (Specification Bias Analysis)

Let the true equation be

$$y_t = a_0 + a_1 x_{1t} + a_2 x_{2t} + u_t$$

where  $u_t$  satisfies all assumptions. Suppose we estimate instead:

$$y_t = a'_0 + a'_1 x_{1t} + v_t$$

Clearly,  $v_t = (a_0 - a'_0) + (a_1 - a'_1)x_{1t} + a_2 x_{2t} + u_t$ . Define the OLS estimator of  $a'_1$ :

$$\hat{a}'_1 = \frac{\sum (x_{1t} - \bar{x}_1) y_t}{\sum (x_{1t} - \bar{x}_1)^2} = \frac{\sum (x_{1t} - \bar{x}_1) [a_0 + a_1 x_{1t} + a_2 x_{2t} + u_t]}{\sum (x_{1t} - \bar{x}_1)^2}$$

$$= \frac{a_1 \sum (x_{1t} - \bar{x}_1) x_{1t} + a_2 \sum (x_{1t} - \bar{x}_1) x_{2t} + \sum (x_{1t} - \bar{x}_1) u_t}{\sum (x_{1t} - \bar{x}_1)^2}$$

$$= a_1 + a_2 \frac{\sum (x_{1t} - \bar{x}_1) x_{2t}}{\sum (x_{1t} - \bar{x}_1)^2} + \frac{\sum (x_{1t} - \bar{x}_1) u_t}{\sum (x_{1t} - \bar{x}_1)^2}$$

Note that  $\frac{\sum (x_{1t} - \bar{x}_1) x_{2t}}{\sum (x_{1t} - \bar{x}_1)^2}$  is the covariance (see p. 75)

This is so because:

$$\hat{\alpha}_1 = \frac{\sum (y_t - \bar{y})(x_{1t} - \bar{x}_1)}{\sum (x_{1t} - \bar{x}_1)^2} = \frac{\sum (x_{1t} - \bar{x}_1)y_t}{\sum (x_{1t} - \bar{x}_1)^2} + \frac{\sum (x_{1t} - \bar{x}_1)\bar{y}}{\sum (x_{1t} - \bar{x}_1)^2}$$

The second term on the r.h.s. is zero as:

$$\frac{\sum (x_{1t} - \bar{x}_1)\bar{y}}{\sum (x_{1t} - \bar{x}_1)^2} = \bar{y} \frac{\sum (x_{1t} - \bar{x}_1)}{\sum (x_{1t} - \bar{x}_1)^2} = 0 \text{ since } \sum (x_{1t} - \bar{x}_1) = 0$$

Substitute the true model:

$$\hat{\alpha}_1 = \frac{\sum (x_{1t} - \bar{x}_1)y_t}{\sum (x_{1t} - \bar{x}_1)^2} = \frac{\sum (x_{1t} - \bar{x}_1)(\alpha_0 + \alpha_1 x_{1t} + \alpha_2 x_{2t} + u_t)}{\sum (x_{1t} - \bar{x}_1)^2}$$

$$= \alpha_1 + \alpha_2 \frac{\sum (x_{1t} - \bar{x}_1)x_{2t}}{\sum (x_{1t} - \bar{x}_1)^2} + \frac{\sum (x_{1t} - \bar{x}_1)u_t}{\sum (x_{1t} - \bar{x}_1)^2}$$

$$\text{as: } (*) \frac{\alpha_0 \sum (x_{1t} - \bar{x}_1)}{\sum (x_{1t} - \bar{x}_1)^2} = 0 \text{ since } \sum (x_{1t} - \bar{x}_1) = 0$$

$$(*) \alpha_1 \frac{\sum (x_{1t} - \bar{x}_1)x_{1t}}{\sum (x_{1t} - \bar{x}_1)^2} = \alpha_1 \frac{\sum (x_{1t} - \bar{x}_1)(x_{1t} - \bar{x}_1)}{\sum (x_{1t} - \bar{x}_1)^2} = \alpha_1$$

between  $x_1$  and  $x_2$  divided by the variance of  $x_2$

$$E(\hat{a}_1) = a_1 + a_2 \frac{\sigma_{x_1 x_2}}{\sigma_{x_1}^2}$$

since  $E(x_2 v) = 0$ .

Consider the auxiliary regression of the omitted variable against included variables:

$$x_{2t} = p_0 + p_1 x_{1t} + \epsilon_t$$

The OLS estimator is:

$$\hat{p}_1 = \frac{\sigma_{x_1 x_2}}{\sigma_{x_1}^2}$$

We write:

$$E(\hat{a}_1) = a_1 + a_2 p_1$$

Hence: The bias  $E(\hat{a}_1) - a_1 = a_2 p_1$  is the product of the true coefficient on excluded variable and the auxiliary coefficient.

Note: The estimator  $\hat{\alpha}'$  is unbiased if either

i)  $\alpha_2 = 0$

ii)  $\rho_{y, x_2} + \sigma_{x_1, x_2}^2 / \sigma_{x_1}^2 = 0$

### Applications to Bivariate Case

#### 1. Production function estimation

Suppose true p.f. is:

$$q_t = \alpha l_t + \beta c_t + u_t$$

where small letters denote logs. Suppose true measure of labor  $L$  is the measured labor  $L^o$  times a quality index  $\bar{z}$ :

$$L_t = \bar{z}_t L_t^o$$

Then  $l_t = \bar{z}_t l_t^o$  Suppose we estimate:

$$\hat{q}_t = \alpha \hat{l}_t + \beta c_t + v_t$$

The true p.f. in terms of observables is

$$y_t = \alpha x_t^0 + \alpha z_t + \beta c_t + u_t$$

and we have effectively omitted  $z_t$ . Then

$$E(\hat{\alpha}) = \alpha + \alpha P_{z^0}$$

$$E(\hat{\beta}) = \beta + \alpha P_{z^0 c}$$

where the auxiliary regression is

$$z_t = P_{z^0 c} c_t + P_{z^0} x_t^0 + \varepsilon_t$$

Can we say something about probable signs of  $P_{z^0 c}$  and  $P_{z^0}$ ? Seems plausible that  $P_{z^0 c} > 0$  and  $P_{z^0} > 0$ . So we expect that both  $\alpha$  and  $\beta$  to be overestimated here.

Note: Both  $\hat{\alpha}$  and  $\hat{\beta}$  are generally biased.

## Error in Measurement and Specification Bias

Let true model be:

$$y_t = a_0 + a_1 x_t + u_t \quad ; \quad u_t \text{ Gauss-Markov.}$$

1. Suppose we observe the independent variable with random error

$$x_t^o = x_t + \epsilon_t \quad \text{where } \epsilon_t \sim N(0, \sigma_\epsilon^2)$$

and  $E(\epsilon_t x_t) = 0$ .

What happens if we estimate using  $x_t^o$  instead of  $x_t$ ?  
First write the true model in terms of observed variable,

$$y_t = a_0 + a_1 (x_t^o - \epsilon_t) + u_t$$

or 
$$y_t = a_0 + a_1 x_t^o + u_t - a_1 \epsilon_t$$

In a sense, we estimate this model but omit the (random) variable  $\epsilon_t$ , which is not observable.

By specification bias analysis we obtain:

$$E(\hat{a}_1) = a_1 - a_1 \rho_1$$

where auxiliary regression is  $\epsilon_t = \rho_0 + \rho_1 x_t^o + \eta_t$ .

But here we can say more about  $\rho_1$ . Because:

$$\rho_1 = \frac{\sigma_{\epsilon x^0}}{\sigma_{x^0}^2} = \frac{\sigma_{\epsilon}^2}{\sigma_{x^0}^2}$$

where second equality uses fact that  $x_t^0 = x_t + \epsilon_t$   
and  $E(\epsilon x) = 0$ .

Then we have:

$$E(\hat{a}_1) = a_1 - a_1 \frac{\sigma_{\epsilon}^2}{\sigma_{x^0}^2} = a_1 \left( 1 - \frac{\sigma_{\epsilon}^2}{\sigma_{x^0}^2} \right)$$

Since  $\sigma_{x^0}^2 = \sigma_x^2 + \sigma_{\epsilon}^2$ ,  $0 \leq \sigma_{\epsilon}^2 / \sigma_{x^0}^2 \leq 1$ . Therefore,  
we conclude that a random error in measurement  
in the independent variable biases the  
coefficient toward zero.

Note that measurement error can be analyzed as a special case of specification bias analysis, where we can say something definite about the auxiliary regression coefficient.

## 2. Error in Measurement : Expectations Case

Suppose wage rates,  $w$ , are a function of the current productivity level,  $y$ , and the expected price level,  $x$ :

$$w_t = \alpha y_t + \beta x_t + u_t$$

We do not observe the expected price level,  $x_t$ , but only the actual level,  $x_t^o$ . Suppose:

$$x_t^o = x_t + v_t \quad v_t \sim N(0, \sigma_v^2)$$

$$E(v_t x_t) = 0$$

i.e. that expectations are on average correct but there is a random noise component,  $v_t$ , which determines the actual level of  $x_t$ ,  $x_t^o$ . In terms of observables, the true model is:

$$w_t = \alpha y_t + \beta (x_t^o - v_t) + u_t$$

i.e.

$$w_t = \alpha y_t + \beta x_t^o - \beta v_t + u_t$$



Hence, the error in measurement is like an omitted variable,  $v_t$ . We analyze as before:

$$E(\hat{\alpha}) = \alpha - \beta P_{vy}$$

$$E(\hat{\beta}) = \beta - \beta P_{vx^0}$$

where

$$y_t^0 = P_{vy} y_t + P_{vx^0} x_t^0 + \xi_t$$

However, in this case we can get a more exact statement of  $P_{vx^0}$ . From

$$x_t^0 = x_t + v_t$$

it follows:

$$P_{vx^0} = \frac{\sigma_{vx^0}}{\sigma_{x^0}^2} = \frac{\sigma_v^2}{\sigma_x^2 + \sigma_v^2}$$

Note  $0 \leq P_{vx^0} < 1$

Hence

$$E(\hat{\beta}) = \beta (1 - P_{vx^0})$$

which we also write as:

$$E(\hat{\beta}) = \beta \left( 1 - \frac{\sigma_v^2}{\sigma_x^2 + \sigma_v^2} \right)$$

or 
$$E(\hat{\beta}) = \beta \left( \frac{\sigma_x^2}{\sigma_x^2 + \sigma_v^2} \right) < \beta$$

The ratio  $\sigma_x^2 / (\sigma_x^2 + \sigma_v^2)$  is called the noise ratio.

We conclude that an error in measurement biases the coefficient of the mismeasured variable toward zero.

A macroeconomic example: consumption function

Permanent Income Hypothesis

$$C_t^P = \alpha Y_t^P + u_t$$

However, we observe not permanent consumption and income but the sum of permanent and transitory movements. Define:

$$C_t^O = C_t^P + C_t^T$$

$$Y_t^O = Y_t^P + Y_t^T$$

and assume:

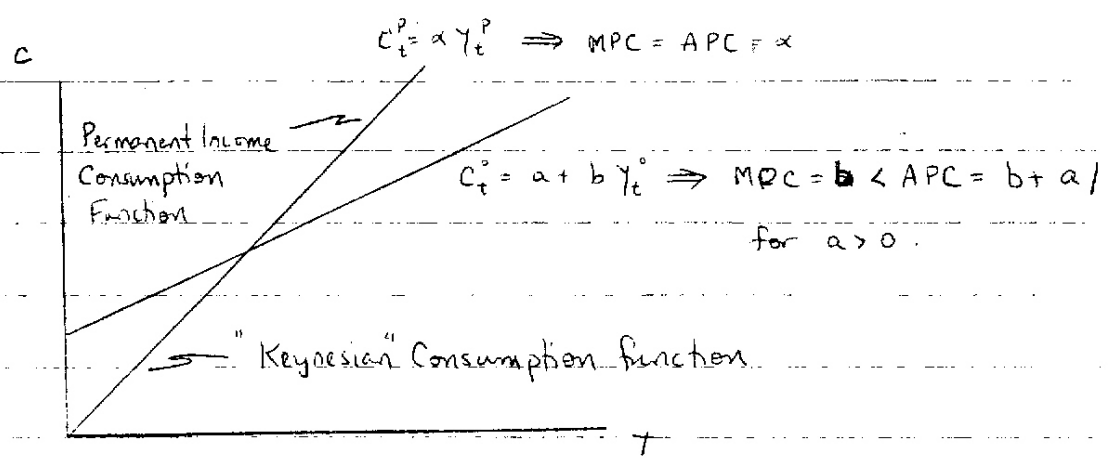
$$E(Y_t^T, Y_s^T) = 0 \quad t \neq s$$

$$E(C_t^T, C_s^T) = 0 \quad t \neq s$$

$$E(Y_t^T, Y_t^P) = E(C_t^T, C_t^P) = E(Y_t^T, C_t^T) = 0$$

$$E(Y_t^P, u_t) = 0$$

$$E(C_t^T) = E(Y_t^T) = 0$$



In terms of observables, the model is

$$C_t^o - C_t^T = \alpha (Y_t^o - Y_t^T) + U_t$$

or 
$$C_t^o = \alpha Y_t^o + (U_t - \alpha Y_t^T + C_t^T)$$

The standard approach is to estimate:

$$C_t^o = a + b Y_t^o + \varepsilon_t$$

The OLS estimator of  $b$  is:

$$\hat{b} = \frac{\sum (C_t^o - \bar{C}^o)(Y_t^o - \bar{Y}^o)}{\sum (Y_t^o - \bar{Y}^o)^2}$$

Substituting for  $C_t^o$  from the true model:

$$\hat{b} = \frac{\sum [\alpha (Y_t^o - \bar{Y}^o) + (U_t - \alpha Y_t^T + C_t^T)] (Y_t^o - \bar{Y}^o)}{\sum (Y_t^o - \bar{Y}^o)^2}$$

$$= \alpha + \frac{\sum (C_t^T - \alpha Y_t^T) (Y_t^o - \bar{Y}^o)}{\sum (Y_t^o - \bar{Y}^o)^2}$$

But note

$$y_t^o - \bar{y}^o = y_t^T$$

Hence

$$\sum c_t^T (y_t^o - \bar{y}^o) = 0$$

and

$$\alpha \sum y_t^T (y_t^o - \bar{y}^o) = \alpha \sigma_{y^T}^2$$

whence

$$\hat{b} = \alpha \left( 1 - \frac{\sigma_{y^T}^2}{\sigma_{y^o}^2} \right) < \alpha$$

Since  $\hat{b} < \alpha$ , we would also conclude that

$$\hat{a} > 0$$

Conclusion: Failure to consider transitory components leads to a downward bias in the estimate of the MPC and a spurious finding of a positive intercept.

## Instrumental Variable (IV) Estimation

Consider an equation:

$$y_t = b x_t^* + u_t$$

$u_t$  Gauss Markov.

$$x_t = x_t^* + \epsilon_t$$

$$E(\epsilon_t x_t^*) = E(\epsilon_t) = 0$$

$$E(\epsilon_t^2) = \sigma_\epsilon^2$$

Then:

$$y_t = b x_t + (u_t - b \epsilon_t) = b x_t + v_t$$

Note!

$$\begin{aligned} E(x_t v_t) &= E[x_t (u_t - b \epsilon_t)] = E x_t u_t - b E x_t \epsilon_t \\ &= -b E(x_t^* + \epsilon_t) \epsilon_t = -b \sigma_\epsilon^2 \end{aligned}$$

Then OLS estimation yields:

$$\hat{b}_{OLS} = \frac{\sum x_t y_t}{\sum x_t^2} = b + \frac{\sum x_t u_t}{\sum x_t^2} - b \frac{\sum x_t \epsilon_t}{\sum x_t^2}$$

$$\text{plim } \hat{b}_{OLS} = b - b \frac{\sigma_\epsilon^2}{\sigma_x^2} = b \left( 1 - \frac{\sigma_\epsilon^2}{\sigma_x^2} \right) < b$$

How to derive a consistent estimator of  $b$ ?

Define: Instrumental variable  $z_t$

- i)  $E(z_t x_t^*) \neq 0$  ..
- ii)  $E(z_t \epsilon_t) = 0$ .

Consider IV estimator

$$\hat{b}_{iv} = \frac{\sum y_t z_t}{\sum x_t z_t} = \frac{b \sum x_t z_t + \sum u_t z_t - b \sum \epsilon_t z_t}{\sum x_t z_t}$$

$$= b + \frac{\sum u_t z_t}{\sum x_t z_t} - \frac{b \sum \epsilon_t z_t}{\sum x_t z_t}$$

$$\text{plim } \hat{b}_{iv} = b.$$

IV Estimator: Covariance between  $y_t$  and  $z_t$   
divided by covariance between  
 $x_t$  and  $z_t$ .



Another interpretation: Two Stage Least Squares (2SLS)

Consider two step procedure:

- i) Regress  $x_t$  against  $z_t$  and take fitted value  $\hat{x}_t$ .
- ii) Regress  $y_t$  against  $\hat{x}_t$ .

Proposition: This two stage estimator is equivalent to IV estimator.

Proof.

$$x_t = A z_t + \epsilon_t$$

$$\hat{A} = \frac{\sigma_{xz}}{\sigma_z^2}$$

$$\hat{x}_t = \hat{A} z_t = \frac{\sigma_{xz}}{\sigma_z^2} z_t$$

Then take

$$y_t = b \hat{x}_t + v_t$$

$$\hat{b} = \frac{\sigma_{y\hat{x}}}{\sigma_{\hat{x}}^2} = \frac{\sigma(\gamma, \sigma_{xz}/\sigma_z^2 \cdot z)}{\sigma^2(\sigma_{xz}/\sigma_z^2 \cdot z)}$$

$$= \frac{(\sigma_{xz}/\sigma_z^2) \sigma_{yz}}{(\sigma_{xz}^2/\sigma_z^4) \sigma_z^2} = \frac{\sigma_{yz}}{\sigma_{xz}} = \hat{b}_{yx}$$

qed.

## Inclusion of Irrelevant Variables

Consider the ~~true~~ model

$$i) \quad y_t = a_1 x_{1t} + a_2 x_{2t} + U_t$$

and the estimated version

$$ii) \quad y_t = a_1 x_{1t} + a_2 x_{2t} + a_3 x_{3t} + U_t$$

Clearly, provided  $E(U_t x_{3t}) = 0$ , all estimators will be unbiased. But are they BLUE?

Consider the variances of say  $\hat{a}_1$ :

$$\text{From (i)} \quad \sigma_{\hat{a}_1}^2 = \sigma_u^2 / \sum \hat{v}_t^2$$

where  $\hat{v}_t$  is defined by

$$x_{1t} = \hat{\alpha}_0 + \hat{\alpha}_1 x_{2t} + \hat{v}_t$$

$$\text{From (ii)} \quad \sigma_{\hat{a}_1}^2 = \sigma_u^2 / \sum \hat{w}_t^2$$

where  $\hat{w}_t$  is defined by

$$x_{1t} = \hat{\alpha}_0 + \hat{\alpha}_1 x_{2t} + \hat{\alpha}_2 x_{3t} + \hat{w}_t$$

Clearly, provided  $\hat{\beta}_2 \neq 0$ ,  $\sum \hat{v}_{1t}^2 > \sum \hat{w}_{1t}^2$ . Hence

$$\sigma_{\hat{\beta}_1}^{2(i)} = \frac{\sigma_u^2}{\sum \hat{v}_{1t}^2} < \sigma_{\hat{\beta}_1}^{2(ii)} = \frac{\sigma_u^2}{\sum \hat{w}_{1t}^2}$$

Conclusion: Including irrelevant variables in general leaves unbiasedness property intact. However, unless the included irrelevant variable is uncorrelated with true exogenous variables, the variances of estimated coefficients are increased.