

# Lecture notes VIII

1

## Assessing Studies Based on Multiple Regression

- SW ch. 9

Definition: Internal validity means that the (statistical) inferences are valid for the population actually studied.

External validity means that the inferences can be generalized beyond the population actually studied, i.e. to other populations and studies

Examples: - Time series analysis and prediction  
- Country specific studies

What is actually the underlying population in non-experimental data?

Focus on internal validity :

Internal validity specifically has to do with the following issues :

1. ~~Are~~ the estimators unbiased and (more importantly) consistent ?

- can we interpret the <sup>interest-</sup>parameter as representing a partial effect (i.e. causal effect of a partial change in the variable of interest)

2. Are the standard errors correct ?

- robustness standard errors

3. Do confidence intervals have the correct coverage and do hypothesis tests have the correct level of significance ?

5 main reasons for violation of internal validity :

- 1. Omitted variable bias
- 2. Misspecification of functional form
- 3. Errors in variables
- 4. Sample selection
- 5. Simultaneous causality

1. Omitted variable bias :

- one variable determining  $\tau$  and which is correlated with an included variable is excluded

example: "public sector" dummy in estimates of returns to schooling.

Remedy: Test whether additional variables change estimates of interest parameter. Are they significant determinants of  $\tau$ ?  $\bar{R}^2$ -adj may be helpful.

Try out different specifications and disclose the results!

## 2. Misspecification of functional form

- Example 1: The marginal returns to schooling is not constant  $\Rightarrow$  log-log model is misspecified.
- Example 2: Relation between R&D subsidies obtained under SkatteFUNN and number of employees in the firm not a linear-log model!

## 3. Errors-in-variables

The regressor is measured with error:

$$X^{obs} = X + \varepsilon, \quad E(\varepsilon | X) = 0$$

where  $X$  true value of  $X$  and  $\varepsilon$  measurement error.

Regression:

$$Y = \beta_0 + \beta_1 X + u, \quad E(u | X) = 0$$

$$= \beta_0 + \beta_1 (X^{obs} - \varepsilon) + u$$

$$= \beta_0 + \beta_1 X^{obs} + \tilde{u}, \quad \tilde{u} = u - \beta_1 \varepsilon$$

Then

$$\text{bias } \hat{\beta}_1 = \frac{\text{Cov}(X^{obs}, u)}{\text{Var}(X^{obs})} = \frac{\text{Cov}(X + \epsilon, u - \beta_1 \epsilon)}{\text{Var}(X + \epsilon)}$$

$$= \frac{-\beta_1 \text{Var}(\epsilon)}{\text{Var}(X) + \text{Var}(\epsilon)} = -\beta_1 \frac{\sigma_\epsilon^2}{\sigma_X^2 + \sigma_\epsilon^2}$$

- bias does not vanish when  $n \rightarrow \infty$

### 4. Sample selection

Whether or not a potential observation unit is included in the sample depends on the error term in the regression equation.

- but no problem if this depends on observed variables

~~Example: Truncation: Unit i is observed only if  $\gamma_i \geq a$ , where a is some threshold.~~

Example: Truncation: Unit i is observed only if  $\gamma_i \geq a$ , where a is some threshold.

i.e. units with low values of  $\gamma_i$  are never observed

Assume:

$$Y_i = \beta_0 + \beta_1 X_i + u_i, \quad E(u_i | X_i) = 0$$

If unit  $i$  is observed, then  $Y_i \geq a$

$$\Leftrightarrow \beta_0 + \beta_1 X_i + u_i \geq a \quad (\text{depends on } u_i)$$

$$\Rightarrow E(u_i | X_i, "i \text{ is observed}") =$$

$$E(u_i | X_i, Y_i \geq a) = E(u_i | X_i, \beta_0 + \beta_1 X_i + u_i \geq a)$$

$$= E(u_i | u_i \geq a - \beta_0 - \beta_1 X_i) > 0$$

### 5. Simultaneous causality

-  $X$  causes  $Y$ , but also  $Y$  causes  $X$

Example: Supply and Demand

$$\textcircled{1} \quad X^S = \beta_0 + \beta_1 P + \varepsilon, \quad X^S \equiv \text{supply}$$

$P$  is price,  $\varepsilon$  error term

$$\textcircled{2} \quad X^D = \lambda_0 + \lambda_1 P + u, \quad X^D \equiv \text{demand}, \quad u \text{ is error}$$

$$\text{Cov}(\varepsilon, u) = 0$$

In equilibrium:  $X^S = X^D$

7

$$\Leftrightarrow \beta_0 + \beta_1 P + \varepsilon = \lambda_0 + \lambda_1 P + u$$

$$\Rightarrow (\beta_1 - \lambda_1)P = \lambda_0 - \beta_0 + u - \varepsilon$$

$$\Rightarrow P = \frac{\lambda_0 - \beta_0}{\beta_1 - \lambda_1} + \frac{u - \varepsilon}{\beta_1 - \lambda_1}$$

If we estimate ① by OLS, we get biased estimates, since

$$\text{Cov}(\varepsilon, P) = \text{Cov}\left(\varepsilon, \frac{u - \varepsilon}{\beta_1 - \lambda_1}\right) = -\frac{\text{Var}(\varepsilon)}{\beta_1 - \lambda_1}$$

even if  $\text{Cov}(\varepsilon, u) = 0$

Similarly, if we estimate ② by OLS, the regressor and error is also correlated:

$$\text{Cov}(u, P) = \text{Cov}\left(u, \frac{u - \varepsilon}{\beta_1 - \lambda_1}\right) = \frac{\text{Var}(u)}{\beta_1 - \lambda_1}$$