

• When do we violate  $E(u_i | x_i) = 0$  assumption?

④ Measurement Error among regressor.

$$\underline{y_i = \alpha + \beta x_i^* + u_i} \quad (x_i^* : \text{Not observed})$$

→ Until now, when we think of data, those are considered as "observed".

Thus, we can refer the data as "fixed ones".

If  $X$  can be controlled,  $X$ : fixed variable.

If  $X$  cannot be controlled,  $X$ : stochastic variable.

→ Now, let's think of a case:  $x_i^*$  cannot be observed. (Assume  $E(u_i | x_i^*) = 0$ )

Instead, we can see  $x_i = x_i^* + v_i$  under  $E(u_i | x_i^*) = 0$

↓ implies  
 $E(x_i^* u_i) = 0$

•  $x_i^*, v_i$  are independent. ( $x_i^* \perp v_i$  and  $E(v_i) = 0$ )

⇒  $E(x_i^* v_i) = 0$

•  $v_i, u_i$  are also independent. ( $v_i \perp u_i$  and  $E(u_i) = 0$ )

⇒  $E(u_i v_i) = 0$

Classical Measurement error model. (That is, assume

- |   |                    |
|---|--------------------|
| ① | $E(x_i^* u_i) = 0$ |
| ② | $E(x_i^* v_i) = 0$ |
| ③ | $E(u_i v_i) = 0$   |

• used model:  $y_i = \alpha + \beta x_i + \varepsilon_i$  (Not true model)

$$y_i = \alpha + \beta x_i^* + u_i$$

$$= \alpha + \beta (x_i - v_i) + u_i = \alpha + \beta x_i + (u_i - \beta v_i)$$

check  $E(x_i \varepsilon_i) = 0$  :  $E(x_i \varepsilon_i) = E[(x_i^* + v_i)(u_i - \beta v_i)]$

$$= E[(x_i^* u_i + v_i u_i - \beta x_i^* v_i - \beta v_i^2)] = 0$$

↓  $E(x_i^* u_i) = 0$  and  $E(x_i^* v_i) = 0$

$$= E(v_i u_i) - \beta E(v_i^2) = \underset{\uparrow}{-} \beta E(v_i^2) = -\beta \sigma_v^2 \neq 0$$

By assumption,  $v_i, u_i$  are independent.

Note Several assumption. (Not Lecture Note)

Strong  
assump.

	$E(\hat{\beta}) = \beta$ (unbiased)	$\hat{\beta} \xrightarrow{P} \beta$ (consistent)
$X_i^*$ : constant (Always $X_i^*, u_i$ independent)	○	○
$X_i^*$ : stochastic & $X_i^*, u_i$ are independent ( $X_i \perp u_i$ and $E(u_i) = 0$ )	○	○
$X_i^*$ : stochastic & $E(u_i   X_i) = 0$	○	○
$X_i^*$ : stochastic & $\text{COV}(X_i, u_i) = 0$ , $E(u_i) = 0$	Not always	○

①  $X_i$ : stochastic &  $X_i, u_i$  are independent.

$$\begin{aligned}
 E(\hat{\beta}) - \beta &= E\left[\left(\sum X_i^2\right)^{-1} \sum X_i u_i\right] = E\left(\frac{\sum X_i}{\sum X_i^2} u_i\right) = E[\sum W_i u_i] \text{ where } W_i = \frac{X_i}{\sum X_i^2} \\
 &= \sum E(W_i u_i) = \sum E(W_i) E(u_i) = 0. \quad (\because X_i, u_i \text{ independent} \Rightarrow E(u_i) = 0)
 \end{aligned}$$

Therefore,  $E(\hat{\beta}) = \beta$  (unbiased)

$$\hat{\beta} = \left(\sum X_i^2\right)^{-1} \sum X_i y_i = \beta + \frac{\sum X_i u_i}{\sum X_i^2} = \beta + \frac{\frac{1}{N} \sum X_i u_i}{\frac{1}{N} \sum X_i^2} \xrightarrow{P} \frac{E(X_i u_i)}{E(X_i^2)} \text{ by LLN}$$

By  $X_i^*, u_i$  independent,  $E(X_i u_i) = E(X_i) E(u_i) = E(X_i) \cdot 0 = 0$

Therefore,  $\hat{\beta} \xrightarrow{P} \beta$  (consistent).

②  $X_i$ : stochastic &  $E(u_i | X_i) = 0$

$$\begin{aligned}
 E(\hat{\beta}) &= \beta + E\left(\frac{\sum X_i u_i}{\sum X_i^2}\right) = \beta + E(\sum W_i u_i) \quad (\text{where } W_i = \frac{X_i}{\sum X_i^2}) \\
 &= \beta + E(E(\sum W_i u_i | X_i)) \quad \text{by law of Iterative Expectation} \\
 &= \beta + E(\sum W_i \underbrace{E(u_i | X_i)}_{=0}) = \beta
 \end{aligned}$$

$$\hat{\beta} = \beta + \frac{\frac{1}{N} \sum X_i u_i}{\frac{1}{N} \sum X_i^2} \xrightarrow{P} \frac{E(X_i u_i)}{E(X_i^2)} \text{ by LLN. and then}$$

$$E(X_i u_i) = E(E(X_i u_i | X_i)) = E(X_i - E(u_i | X_i)) = 0 \quad \text{Thus, } \hat{\beta} \xrightarrow{P} \beta$$

③  $x_i$  : stochastic &  $\text{cov}(x_i, u_i) = 0$ ,  $E(u_i) = 0$

$$\hat{\beta} = \beta + \frac{\frac{1}{N} \sum x_i u_i}{\frac{1}{N} \sum x_i^2} \xrightarrow{P} \frac{E(x_i u_i)}{E(x_i^2)} \quad \text{by LLN.}$$

We need  $E(x_i u_i)$  for consistency.

$$\because \text{cov}(x_i, u_i) = 0$$

$$\text{cov}(x_i, u_i) = E(x_i - E(x_i))(u_i - E(u_i)) = E[x_i u_i - E(x_i)u_i] = E(x_i u_i) \stackrel{\downarrow}{=} 0$$

Therefore,  $E(x_i u_i) = 0$ . Thus,  $\hat{\beta} \xrightarrow{P} \beta$

\* We usually use assumption for " $x_i$  is stochastic" :  $E(u_i | x_i) = 0$ .

(For unbiasedness & consistency)

Under classical measurement error assumption for regressor,  $E(x_i^* v_i) = 0$

$$y_i = \beta \cdot x_i^* + u_i, \quad E(u_i | x_i^*) = 0$$

↑ we cannot see this

$$E(x_i^* u_i) = 0$$

$$E(u_i v_i) = 0$$

$$= \beta \cdot (x_i - v_i) + u_i = \beta \cdot x_i + \underbrace{(u_i - \beta v_i)}_{\substack{\uparrow \\ \text{we can see this.}}} = \varepsilon_i \quad (\text{Assume } \sigma_x^2 = E(x_i^{*2}), \sigma_v^2 = E(v_i^2))$$

$$\hat{\beta} = \beta + \frac{\frac{1}{N} \sum x_i \varepsilon_i}{\frac{1}{N} \sum x_i^2}$$

$$\textcircled{1} \frac{1}{N} \sum x_i^2 = \frac{1}{N} \sum (x_i^* + v_i)^2 = \frac{1}{N} \sum x_i^{*2} + \frac{2}{N} \sum x_i^* v_i + \frac{1}{N} \sum v_i^2 \xrightarrow{P} \sigma_x^2 + \sigma_v^2$$

$\nearrow \begin{matrix} E(x_i^* v_i) \text{ by LLN} \\ = 0 \end{matrix}$

$$\textcircled{2} \frac{1}{N} \sum x_i \varepsilon_i = \frac{1}{N} \sum (x_i^* + v_i)(u_i - \beta v_i) = \frac{1}{N} \sum x_i^* u_i + \frac{1}{N} \sum u_i v_i + \frac{1}{N} \sum x_i^* (-\beta v_i) - \beta \frac{1}{N} \sum v_i^2$$

$\begin{matrix} \downarrow & & \downarrow & & \downarrow & & \downarrow \\ E(x_i^* u_i) = 0 & & E(u_i v_i) = 0 & & -\beta E(x_i^* v_i) = 0 & & -\beta E(v_i^2) \end{matrix}$

$$\xrightarrow{P} -\beta \cdot \sigma_v^2$$

$$\hat{\beta} = \beta + \frac{\frac{1}{N} \sum x_i \varepsilon_i}{\frac{1}{N} \sum x_i^2} \xrightarrow{P} \beta + \frac{-\beta \sigma_v^2}{\sigma_x^2 + \sigma_v^2} = \left(1 - \frac{\sigma_v^2}{\sigma_x^2 + \sigma_v^2}\right) \beta < \beta$$

$\Rightarrow$  If  $\exists$  measurement error,  $\hat{\beta}$  converge to less number than  $\beta$ .

(Not consistent)



## c.f) Measurement Error on the dependent variable

→ If the measurement error is on the dependent variable, the same argument for the previous case does not hold.

$$\underline{y_i^* = \alpha + \beta x_i + u_i} \quad (y_i^*: \text{Not observed}) \quad (\text{Assume } E(u_i | x_i) = 0)$$



$$y_i = y_i^* + v_i : y_i \text{ observed.}$$

$$y_i = y_i^* + v_i = \alpha + \beta x_i + \boxed{u_i + v_i}$$

Logically/Usually  
 $x_i, v_i$  are not correlated.  
 maybe.  
 ↓

$$E(x_i \cdot (u_i + v_i)) = E(x_i u_i) + E(x_i v_i) = E(x_i v_i) = 0.$$

$$\hookrightarrow E(x_i u_i) = E(E(x_i u_i | x_i)) = E(x_i E(u_i | x_i)) = 0.$$

∴ Therefore, we could assume that measurement error on  $y_i^*$  and  $x_i$  are Not Correlated without falling into logical inconsistency.

## c.f 2) Dummy variable case.

→ If  $x_i^*$  is discrete (so that  $x_i$  is also discrete) then the classical measurement error does not hold.

$$\text{ex) } x_i^* = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \quad x_i = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

$$x_i = x_i^* + v_i \Rightarrow \begin{array}{ll} x_i^* = 1 & v_i = 0 \text{ or } -1 \\ x_i^* = 0 & v_i = 0 \text{ or } 1 \end{array} \quad \Bigg) \text{ Not independent.}$$

Negatively correlated. : This shows dummy variable case  
 (Maybe,  $E(x_i^* v_i) \leq 0$ ). break classical assumption for measurement error.



This case is a little different from the previous dummy.

: Comment on proxy

$$y_i = \alpha + \beta x_i + u_i, \quad E(u_i | x_i) = 0$$

When we use proxy, we can make  $Z_i$  as follows;

$$Z_i = \begin{cases} 1 & \text{if } f(x_i) \geq 0 \\ 0 & \text{if } f(x_i) < 0 \end{cases}$$

$$\begin{aligned} y_i &= \underbrace{\theta_0 + \theta_1 Z_i}_{E(y_i | Z_i)} + \underbrace{\varepsilon_i}_{y_i - E(y_i | Z_i)} \\ &= E(y_i | Z_i) + y_i - E(y_i | Z_i) \end{aligned}$$

→ In this case, we don't want to get  $\beta$  so that there is No measurement error for the proxy.

• When do we violate  $E(u_i|x_i)=0$  assumption?

⑤ Lagged Dependent Variable among regressors  
& serial correlation in the residual term.

$$\left( \begin{array}{l} \underline{y_t = \alpha + \beta \cdot y_{t-1} + u_t} \\ \text{and } \underline{u_t = \rho u_{t-1} + v_t} \end{array} \right. \quad \begin{array}{l} \text{When we see } y_{t-1}, \\ y_{t-1} = \alpha + \beta y_{t-2} + \underline{u_{t-1}} \end{array}$$

Thus,  $y_t$  and  $y_{t-1}$  are generally correlated.

\* Solution for each case of the violation of  $E(u_i|x_i)=0$  assumption.

① Measurement Error

② Simultaneity

③ Lagged Dependent Variable & Serial Correlation.

$\Rightarrow$  Instrumental Variable

• Panel Data analysis.

④ Sample selection  $\Rightarrow$  MLE, Control function Approach, Semi-parametric Analysis.

⑤ Functional Misspecification  $\Rightarrow$  Non-parametric or semi-parametric Analysis.

$\Rightarrow$  As I showed, we should be careful to take a regression when  $E(u_i|x_i) \neq 0$ , which implies it is possible to violate consistency.

- What happens to the OLS Estimator when  $E(u_i|x_i) \neq 0$ ?

$$\hat{\beta} = \beta + (X'X)^{-1}X'u = \beta + \left(\frac{1}{N}X'X\right)^{-1} \frac{1}{N}X'u \xrightarrow{P} \beta + E(X_iX_i')^{-1} E(X_i u_i)$$

$$\begin{array}{ccc} \searrow P & & \searrow P \\ E(X_iX_i') & & E(X_i u_i) = E(E(X_i u_i | X_i)) = E(X_i E(u_i | X_i)) \\ & & = E(X_i) \cdot E(u_i | X_i) \end{array}$$

→ Thus, when  $E(u_i|x_i) \neq 0$ , inconsistency is possible.

- Effect of function misspecification

① Including a regressor which is not needed among the regressors.

$$E(y|x) = \beta_1 + \beta_2 x_2 + \beta_3 x_3 \quad \& \quad \beta_3 = 0.$$

→ Even though  $\beta_3 = 0$  is true, we cannot know this.

Thus, we usually put  $x_2, x_3$  for check the real structure for data.

→ In this case, even if we eliminate  $\beta_3 x_3$ , the above model still holds so that conditional unbiasedness and consistency for OLS estimator hold.

$$\begin{aligned} \text{ex)} \quad u_i &= y_i - E(y_i | x_i) \\ &= y_i - [\beta_1 + \beta_2 x_{2i} + \beta_3 x_{3i}] \end{aligned}$$

$$y_i = \beta_1 + \beta_2 x_{2i} + \beta_3 x_{3i} + u_i, \quad E(u_i | x_i) = 0 \Rightarrow \text{Still classic.}$$

but the conditional variance is bigger.

$$\text{ex2)} \quad \text{Var}(\hat{\beta}_2 | X) = \frac{\sigma^2}{\sum_{i=1}^N \hat{v}_{2i}^2}, \quad \hat{v}_{2i}: \text{the residual of auxiliary regression.}$$



② Omitting a regressor which should be included?

$$E(y|x) = \beta_1 + \beta_2 x_2 + \beta_3 x_3, \quad \beta_3 \neq 0. \quad (\text{True model})$$

But regress  $y$  on  $1, x_2$  :  $E(y|x) = \beta_1 + \beta_2 x_2$  [short regression]

↓

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

$$y_i = \beta_1 + \beta_2 x_{2i} + \boxed{\beta_3 x_{3i} + u_i}$$

This reg. problem is, in this case,  
 $u_i$  and  $x_{2i}$  correlated.

If  $x_{2i}, x_{3i}$  are correlated,

then  $x_{2i}, \varepsilon_i$  are correlated, which implies  $\hat{\beta}_2$  is inconsistent.

("omitted variable bias" causes Direction of inconsistency problem.)

• How do we assess the direction of the omitted variable problem?

⇒ use auxiliary regression analysis.

$$y_i = \beta_1 + \beta_2 x_{2i} + \beta_3 x_{3i} + u_i \quad (\text{True model}).$$

$$y_i = \hat{\beta}_1 + \hat{\beta}_2 x_{2i} + \hat{\beta}_3 x_{3i} + \hat{u}_i \quad (\text{Ideal estimation})$$

Auxiliary regression  $x_{3i}$  on  $1$  and  $x_{2i}$

$$x_{3i} = \pi_0 + \pi_1 x_{2i} + \hat{v}_i$$

$$y_i = \hat{\beta}_1 + \hat{\beta}_2 x_{2i} + \hat{\beta}_3 x_{3i} + \hat{u}_i$$

$$= \hat{\beta}_1 + \hat{\beta}_2 x_{2i} + \hat{\beta}_3 [\hat{\pi}_0 + \hat{\pi}_1 x_{2i} + \hat{v}_i] + \hat{u}_i$$

$$= \hat{\beta}_1 + \hat{\beta}_3 \hat{\pi}_0 + \underline{(\hat{\beta}_2 + \hat{\beta}_3 \hat{\pi}_1)} x_{2i} + \hat{u}_i + \hat{\beta}_3 \hat{v}_i$$

↑  $\beta_2$  : inconsistent if  $\hat{\beta}_3 > 0, \hat{\pi}_1 > 0$ .

(Continuing...)