

$$Y = \beta_0 + \beta_1 X + u \quad \left\{ \begin{array}{l} \text{Endogeneity: there exists some omitted} \\ \text{variable (s) } w \\ \text{i) in the error term } u \\ \text{(so } w \text{ is a determinant of } Y) \\ \text{ii) } w \text{ is correlated with } X \\ \therefore \text{ OVB in estimation} \end{array} \right.$$

• X is **endogenous** : $E[u|X] \neq 0$

\Rightarrow Violation of OLS assumption of β_1

\Rightarrow OVB for an endogenous variable

| Variation in X | |
|-------------------|----------------|
| Endogenous | Exogenous |
| $E[u X_n] \neq 0$ | $E[u X_x] = 0$ |

$\Rightarrow Z$ valid instrument

① Relevance : $\text{Cov}(X, Z) \neq 0$
 $\text{Cor}(X, Z) \neq 0$

② Exogenous : $E[u|Z] = 0$

① Two-stage approach

i) Regress $X = \pi_0 + \pi_1 Z + e$

$$\Rightarrow \hat{X} = \hat{\pi}_0 + \hat{\pi}_1 Z$$

$= f(Z)$, a linear function of Z , exogenous

$\therefore \hat{X}$ must be exogenous also

ii) Regress $Y = \beta_0 + \beta_1 \hat{X} + e$

$$= \beta_0 + \beta_1 (\hat{\pi}_0 + \hat{\pi}_1 Z) + e$$

$$= \underbrace{(\beta_0 + \beta_1 \hat{\pi}_0)}_{\text{constant}} + \beta_1 \underbrace{(\hat{\pi}_1 Z)}_{\text{exogenous}} + e$$

$\therefore \beta_1$ estimated without OVB

② Control Variable Approach

$$Y = \beta_0 + \beta_1 X + u$$

- X is endogenous : $E[u | X] \neq 0$
- In particular, there is a (set of) variable(s) W s.t.
 - W is contained in u : it is a determinant of Y
 - $\text{Cov}(X, W) \neq 0 \iff \text{Corr}(X, W) \neq 0 \therefore \hat{\beta}_{1,OLS}$ estimated with OVB

Goal: Isolate the effect of X on Y independent of W
 supposing we have data on W β_1

The OLS assumption for valid inference in multiple regression:

$$E[u | X, W] = E[u | W] \quad \text{"conditional mean assumption"} \quad \begin{array}{l} \text{holding } W \\ \text{constant,} \\ X \text{ is uncorrelated} \\ \text{with the error} \\ \text{and } \beta_1 \text{ can} \\ \text{be estimated} \\ \text{without OVB} \end{array}$$

Variation in X

| Endogenous | Exogenous |
|------------|-----------|
|------------|-----------|

Covariation with Z

all variation in X NOT explained by Z
 must be a comprehensive source of endogenous variation

$$\begin{aligned} E[u | \hat{x}] &= E[u | X - \hat{x}] \\ &= E[u | X - f(Z)] \\ &= E[u | X] \\ E[u | \hat{x}, X] &= E[u | X] \end{aligned}$$

Conditional mean assumption is satisfied

\Rightarrow Estimation of β_1 is unbiased with inclusion of W as a control variable

In fact, because we're using the same information (X, Z) to identify β_1 , the estimates will be

exactly the same as the first approach

$$\Rightarrow \hat{\beta}_{1,2S} = \hat{\beta}_{1,control} \sim \text{unbiased } (\beta_1)$$

① Regress $X = \pi_0 + \pi_1 Z + e$

$$\Rightarrow \hat{X} = \hat{\pi}_0 + \hat{\pi}_1 Z = f(Z), \text{ must be exogenous}$$

$$\Rightarrow (X - \hat{X}) \text{ will capture all endogenous variation in } X$$

$$\hat{u} \text{ residual} \Rightarrow \text{candidate for } W$$

② Regress $Y = \beta_0 + \beta_1 X + \beta_2 W + u$
 $= \beta_0 + \beta_1 X + \beta_2 (X - \hat{X}) + u$