Help Section Week#4[1]

Why we want instrument variable?:

OVB Problem

- Consider regression:

$$\text{(long)} \ \ Y = \beta_0^{long} + \beta_1^{long} X_1 + \gamma^{long} X_2 + \epsilon$$
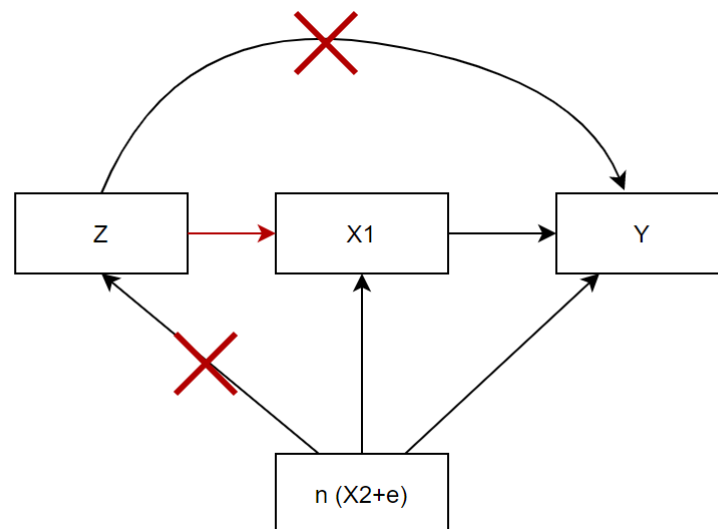$$\text{(short)} \ Y = \beta_0^{short} + \beta_1^{short} X_1 + \eta$$
$$\implies \hat{\beta_1}^{short} = \hat{\beta_1}^{long} + \hat{\gamma}^{long} \frac{\text{cov}(X_1, X_2)}{\text{var}(X_1)}$$

- In this case, we are not able to estimate the marginal effect of $X_1$ on $Y$:

  If $X_1$ in crease by one unit, $X_2$ will change by some unit, thus the change of $Y$ includes both the effect from changes in $X_1$ and $X_2$. In order to estimate the marginal effect of $X_1$ on $Y$, we want "something" that can change $X_1$ without affecting other unobserved/ignore variables (errors) or affecting $Y$ directly.

- Consider $X_1$'s instrument variable $Z$ satisfies:



  (i) Relevant: $\text{Cov}(Z, X_1) \neq 0$
  (ii) Exogeneity: $\text{Cov}(Z, \eta) = 0$ and Z should be related to Y only through X

  These assumptions ensure Z to be a valid instrument variable: by (i), we can change the value of $X_1$ through Z; by (ii), we are able to change $X_1$ without affecting other variables and the changes of Y is generated purely through the changes on $X_1$.

---

[1]Please watch for typos and errors

Variable with Measure Error

- Consider regression:

$$Y = \beta_0 + \beta_1 X^* + \eta$$

where $X^*$ meets all of the standard exogeneity assumptions, however, we do not observe true $X^*$, instead we observe X:

$$X = X^* + \mu_X$$

where $\mu_X$ is the measure error.

We are interested in the marginal effect of $X^*$, but we can only run regression on X:

$$\begin{aligned} Y &= \beta_0 + \beta_1 X^* + \eta \\ &= \beta_0 + \beta_1(X - \mu_X) + \eta \\ &= \beta_0 + \beta_1 X + (\eta - \beta_1 \mu_X) \end{aligned}$$
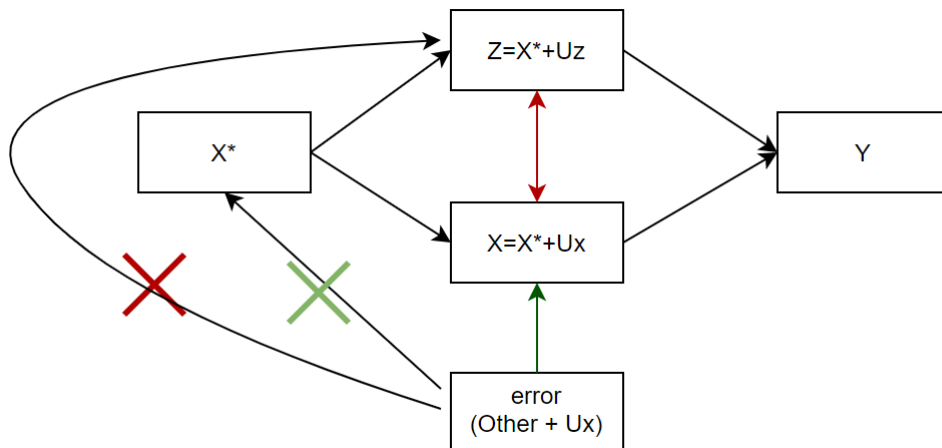
Since X is related to the composite error term $(\eta - \beta_1\mu_X)$, thus OLS will produce biased estimator of $\beta_1$.

- Suppose Z is another measure of $X^*$ with measure error:

$$Z = X^* + \mu_Z$$

where $\mu_Z$ is the measure error and we have $\mathrm{Cov}(\mu_Z, \eta) = 0$ and $\mathrm{Cov}(\mu_Z, \mu_X) = 0$. We can use Z as an instrument if:

(i) Relevant: $\mathrm{Cov}(Z, X) = \mathrm{Cov}(X^* + \mu_Z, X^* + \mu_X) = \mathrm{Var}(X^*) \neq 0$

(ii) Exogeneity: $\mathrm{Cov}(Z, \eta) = \mathrm{Cov}(X^* + \mu_Z, \eta) = \mathrm{Cov}(X^*, \eta) + \mathrm{Cov}(\mu_Z, \eta) = 0$

Implement:

- Two Stage Least Squares(2SLS)

$$(\text{First stage}) \quad X_1 = \pi_0 + \pi_1 Z + \mu$$
$$(\text{Second stage}) \quad Y = \beta_0 + \beta_1 \hat{X}_1 + \eta$$

The estimator of 2SLS:

$$\begin{aligned}
\text{Cov}(Y, Z) &= \text{Cov}(\beta_0^{\text{long}} + \beta_1^{\text{long}} X_1 + \gamma^{\text{long}} X_2 + \epsilon, Z) \\
&= \text{Cov}(\beta_0^{\text{long}}, Z) + \text{Cov}(\beta_1^{\text{long}} X_1, Z) + \text{Cov}(\gamma^{\text{long}} X_2, Z) + \text{Cov}(\epsilon, Z) \\
&= 0 + \beta_1^{\text{long}} \text{Cov}(X_1, Z) + \text{Cov}(\eta, Z) \\
\implies \beta_1^{\text{long,IV}} &= \frac{\text{Cov}(Y, Z)/\text{Var}(Z)}{\text{Cov}(X_1, Z)/\text{Var}(Z)} + \frac{\text{Cov}(\eta, Z)}{\text{Cov}(X_1, Z)} = \beta_1^{\text{long}} + \frac{\text{Cov}(\eta, Z)}{\text{Cov}(X_1, Z)} \\
&= \beta_1^{\text{long}} \quad (\text{if Z is exogenous, i.e. } \text{Cov}(\eta, Z) = 0)
\end{aligned}$$

- Compare IV and OLS:

  (i) Coefficient $\beta_1^{\text{IV}}$ with $\beta_1^{\text{OLS2}}$:

$$\beta_1^{\text{OLS}} = \beta_1^{\text{long}} + \frac{\text{Cov}(\eta, X_1)}{\text{Var}(X_1)} = \beta_1^{\text{long}} + \text{Corr}(X_1, \eta) \times \frac{\sigma_\eta}{\sigma_{X_1}}$$
$$\beta_1^{\text{IV}} = \beta_1^{\text{long}} + \frac{\text{Cov}(\eta, Z)}{\text{Cov}(X_1, Z)} = \beta_1^{\text{long}} + \frac{\text{Corr}(Z, \eta)}{\text{Corr}(Z, X_1)} \times \frac{\sigma_\eta}{\sigma_{X_1}}$$

  Suppose the Exogeneity assumption is violated (i.e. $\text{Cov}(Z, \eta) \neq 0$), if the instrument is "weak" (i.e. $\text{Corr}(Z, X_1)$ is small), the bias in the IV estimator will be large.

  (ii) Variance $\text{Var}(\beta_1^{\text{OLS}})$ and $\text{Var}(\beta_1^{\text{IV}})$[3]:

$$\text{Var}(\beta_1^{\text{OLS}}) = \frac{\sigma_\eta^2}{\sum_{i=1}^N (X_{1i} - \bar{X}_1)^2} = \frac{\sigma_\eta^2}{\text{SST}_{X_1}}$$
$$\text{Var}(\beta_1^{\text{IV}}) = \frac{\sigma_\eta^2}{\text{SST}_{X_1} * R_{X_1, Z}^2}$$

  $\text{SST}_{X_1}$ is the sum of residual $R_{X_1, Z}^2$ comes from an auxiliary regression of $X_1$ upon Z and captures the relationship instrument relevance (first stage). The variance of the IV estimator is (almost) always larger than that of OLS, suggesting that we should run OLS when possible.

---

[2]Estimate from : $Y = \beta_0^{\text{short}} + \beta_1^{\text{short}} X_1 + \eta$
[3]$\text{Var}(\beta_1^{\text{IV}}) = \sigma_\eta^2 (X_1' P_Z X_1)^{-1}$

<u>Testing:</u>

Test for Exogeneity

- IV might generate estimators with large standard errors, so we should not apply IV unless necessary. Consider regression:

$$\text{(short) } Y = \beta_0^{\text{short}} + \beta_1^{\text{short}} X_1 + \lambda^{\text{short}} \text{Exogenous} + \eta$$

  If we have $\text{Cov}(X_1, \eta) = 0 \implies$ no endogeneity problem (indicating there is no OVB) $\implies$ OLS can give us an unbiased estimator $\implies$ we should estimate the short via OLS.

- (i) A quick diagnostic is to estimate via OLS and via IV. If they are similar, then you may not have an endogeneity problem

  (ii) A formal rest – Hausman test

$$X_1 = \pi_0 + \pi_1 Z + \pi_2 \text{Exogenous} + v_1$$
$$\implies \text{Cov}(X_1, \eta) = \text{Cov}(\pi_0 + \pi_1 Z + \pi_2 \text{Exogenous} + v_1, \eta) = \text{Cov}(v_1, \eta)$$
$$\eta = \delta v_1 + v_2 \implies \text{Cov}(X_1, \eta) = 0 \text{ if } \delta = 0$$
$$Y = \beta_0^{\text{short}} + \beta_1^{\text{short}} X_1 + \lambda^{\text{short}} \text{Exogenous} + \delta \hat{v}_1 + \text{error}$$

  $H_0 : \delta = 0 (\text{Exogeneity})$ v.s. $H_0 : \delta \neq 0 (\text{Endogeneity})$.

Weak Instruments Test

- Consider the auxiliary regression:

$$X_1 = \pi_0 + \pi_1 Z + \pi_2 \text{Exogenous} + \text{error}$$

  We want $\text{Cov}(X_1, Z) \neq 0$. Since $\pi_1 \approx \frac{\text{Cov}(X_1, Z)}{\text{Var}(Z)}$ so Z is relevant if and only if $\pi_1 \neq 0$. $H_0 : \pi_1 = 0 (\text{Weak})$ v.s. $H_0 : \pi_1 \neq 0 (\text{Relevant})$.