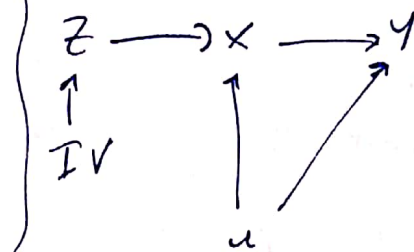
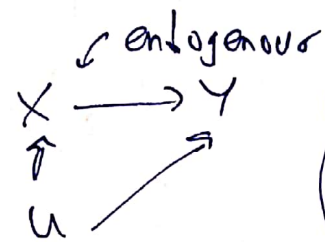
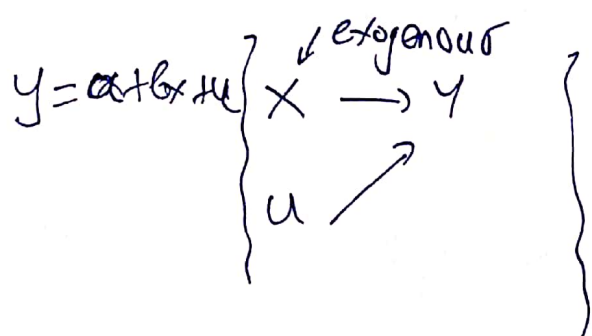


Instrumental Variables



Causes of endogeneity:

- measurement error
- omitted variable bias (OVB)

e.g. OVB

True model $y_i = a + b_1 x_i + b_2 z_i + u_i$ (1)

OVB model $y_i = a + b_1 x_i + \varepsilon_i$ (missing z_i)

In our model $\text{Cov}(x_i, \varepsilon_i) = \text{Cov}(x_i, b_2 z_i + u_i)$
 $= \text{Cov}(x_i, b_2 z_i) + \text{Cov}(x_i, u_i)$

where $\text{Cov}(x_i, u_i) = 0$ (x_i is exogenous on 1)

but $\text{Cov}(x_i, z_i) \neq 0$ if x_i and z_i are correlated
 so x_i is endogenous.

Side note

"Convergence in probability"

$$\lim P(|Y_n - \mu| > \epsilon) = 0$$

Let $\hat{\theta}$ be an estimator for θ . Then

If $\hat{\theta}$ converges in probability to θ
we call $\hat{\theta}$ a consistent estimator,

If $E[\hat{\theta}] = \theta$ we call $\hat{\theta}$ an unbiased estimator.

If all the OLS assumptions hold
(including $\text{Cov}(x_i, u_i) = 0 \forall x_i$) Then
OLS is both consistent and unbiased.

An instrumental variable \check{Z}_i should satisfy:

1) $\text{Cov}(Z_i, x_i) \neq 0$ (the higher $\text{Cov}(x_i, Z_i)$ is
the "stronger" the instrument).

2) $\text{Cov}(Z_i, u_i) = 0$

2 stage linear regression (2SLS)

Stage 1.

Regress endogenous variables on instrument.

$$x_i = z_i \delta + v_i$$

$$\text{get } \hat{\delta} = (Z^T Z)^{-1} Z^T X \xrightarrow{\text{get estimator } X} \hat{x} = Z \hat{\delta} = \underbrace{Z(Z^T Z)^{-1} Z^T}_{P_Z} X = P_Z X$$

Stage 2.

Regress y on estimated \hat{x}

$$y = \hat{x} \beta + \varepsilon_i = P_Z X \beta + \varepsilon_i$$

$$\hat{\beta} = ((P_Z X)^T P_Z X)^{-1} (P_Z X)^T y$$

$$\textcircled{*} = (X^T P_Z X)^{-1} X^T P_Z y$$

↑ substitute P_Z and compute...

⊛ Side note
you can show

$$P_Z^2 = P_Z$$

$$P_Z^T = P_Z$$

$$\text{So now } \hat{\beta} = (Z^T X)^{-1} Z^T y = \underbrace{(Z^T X)^{-1} Z^T}_{\text{substitute } y = P_Z X \beta + \varepsilon_i} (X \beta + \varepsilon) =$$

$$= (Z^T X)^{-1} Z^T X \beta + (Z^T X)^{-1} Z^T \varepsilon = \beta + \left(\frac{Z^T X}{n} \right)^{-1} \left(\frac{Z^T \varepsilon}{n} \right)$$

last term $\frac{1}{n} Z^T \varepsilon = \frac{1}{n} \sum z_i \varepsilon_i \xrightarrow{n \rightarrow \infty} E(z_i \varepsilon_i) = E(z_i) E(\varepsilon_i) = 0$

Since z_i and ε_i are uncorrelated.

from
central
limit theorem

Summary

IV estimator is consistent, which intuitively means that as we get more data points our estimator converge in probability to the ~~true~~ true ones.

Disclaimer: I used naively some theorems for convergence etc. This is not a formal proof but gives us an idea.