



CHULALONGKORN
BUSINESS SCHOOL
FLAGSHIP FOR LIFE



MSF
Chula*

Financial Econometrics

Lecture VI:

Endogeneity, Instrumental Variable,
and 2SLS estimator

Narapong Srivisal, Ph.D.

Outline

causes of

Endogeneity

try to put a lot of controls to the model as regressors so that only noises remain in error term

- Omitted Variable Bias left some random variables that can explain y
↳ in the error term → endogeneity problem
- Measurement Errors
- Simultaneous Equations
- Reversed causality → easily correct

Resolutions for Endogeneity Problem

- Instrument
- Instrumental Variable Estimator
- 2SLS Estimator
- Over-identifying Test

Endogeneity

- Suppose that we want to study *ceteris paribus* effects on Y by using the linear regression model:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_k X_k + U$$

- We call the regressor X_j exogenous ^{variable} if $Cov(X_j, U) = 0$
- We call the regressor X_j endogenous if $Cov(X_j, U) \neq 0$
- We have an endogeneity problem ^{at least 1 random variable x is endogenous} if there is an endogenous regressor.
- This problem results in the OLS estimator being inconsistent and biased.

$$y = \beta_0 + \beta_1 x + (\underbrace{\gamma w + U}_{\text{error}})$$

Suppose we run reg y x

if w and x are correlated → endogeneity problem!

$$* \text{cov}(x, x) = \text{var}(x)$$

$$\hat{\beta}_{OLS} = \frac{\text{cov}(\hat{y}, x)}{\text{var}(\hat{x})} \xrightarrow{w \text{ and } U} \frac{\text{cov}(y, x)}{\text{var}(x)}$$

$$\text{cov}(\beta_0 + \beta_1 x + \gamma w + U, x) = \beta_1 \text{var}(x) + \gamma \text{cov}(w, x) + \text{cov}(U, x)$$

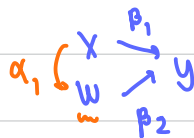
$$\hat{\beta}_{OLS} = \frac{\beta_1 \text{var}(x) + \gamma \text{cov}(w, x) + \text{cov}(U, x)}{\text{var}(x)}$$

$$\hat{\beta}_{OLS} \rightarrow \beta_1 + \frac{\gamma \text{cov}(w, x)}{\text{var}(x)}$$

can omit this iff $\text{cov}(w, x) = 0$

$$y = \beta_0 + \beta_1 x + \beta_2 w + U \quad ; \quad \text{cov}(x, U) = 0$$

$$\text{cov}(w, U) = 0$$



diff

$$w = \alpha_0 + \alpha_1 x + \varepsilon$$

When x moves 1 unit, w moves α_1

not get pure effect of x

$$y = \gamma_0 + \gamma_1 x + \eta$$

When x ↑ 1 unit, how much y? if not control for w?

$$\underbrace{\gamma_1 + \alpha_1 \beta_2}_{\text{contaminated effect from w}}$$

get only correlation indicator not pure effect

ceteris paribus effect of x on y

Omitted Variable Bias

- Suppose that we want to study the effect of X on Y , and the true underlying model is

$$Y = X'\beta + \gamma W + U$$

- If we omit the regressor W in our OLS estimation, the OLS estimator for β is then

$$\hat{\beta}_{omitted} = \left(\sum_{i=1}^n X_i X_i' \right)^{-1} \left(\sum_{i=1}^n X_i Y_i \right)$$

- By the WLLN and CMT, suppose $E[XU] = 0$, this converge in probability to

$$\begin{aligned} E[XX']^{-1}E[XY] &= E[XX']^{-1}E[X(X'\beta + \gamma W + U)] \\ &= \beta + \gamma E[XX']^{-1}E[XW] \end{aligned}$$

- So, our **OLS will be consistent only if $E[XW] = 0$** , that is

$$E[W] = E[X_1W] = E[X_2W] = \dots = E[X_kW] = 0$$

- Or only if the omitted variable W is uncorrelated with each of the other regressors
- Similarly for unbiasedness,

$$\begin{aligned} E[\hat{\beta} | X_1, \dots, X_n] &= E \left[\left(\frac{1}{n} \sum_{i=1}^n X_i X_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n X_i (X_i' \beta + \gamma W + U_i) \right) | X \right] \\ &= \beta + \left(\frac{1}{n} \sum_{i=1}^n X_i X_i' \right)^{-1} \frac{1}{n} \sum_{i=1}^n X_i (\gamma E[W_i | X] + E[U_i | X]) \end{aligned}$$

- So, we **need a stronger condition that $E[W|X] = 0$ in addition to $E[U|X] = 0$ to get unbiased OLS estimator** in this case.

Correlation vs Causation Revisit

- Suppose that we want to study the effect of education on wage, and the true underlying model is

$$\ln(wage) = \beta_0 + \beta_1 educ + \beta_2 exper + U$$

where, supposedly, the correlations between these two regressors and the error term are zero.

- Then, β_1 captures the ceteris paribus effect of education on wage.
- Suppose that education and experience are correlated, that is $\alpha_1 \neq 0$ in the regression model:

$$exper = \alpha_0 + \alpha_1 educ + \epsilon$$

- If we omit experience in the regression, **OLS will not consistently estimate the effect (β_1) but how $\ln(wage)$ is correlated with education**, which **include both the effect of education and the effect of experience**:

* Interpretation of correlation vs causation

* -----

reg lwage educ, nohead ^{without exper}

lwage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
γ_1 educ	.0827444	.0075667	10.94	0.000	.0678796	.0976092
_cons	.5837727	.0973358	6.00	0.000	.3925563	.7749891

reg lwage educ exper, nohead

lwage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
β_1 educ	.0979356 ^{effect of educ on wage}	.0076224	12.85	0.000	.0829613	.1129099
β_2 exper	.0103469 ^{effect of exper}	.0015551	6.65	0.000	.0072919	.013402
_cons	.2168544	.108595	2.00	0.046	.0035183	.4301904

reg exper educ, nohead

exper	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
$\hat{\alpha}_1$ educ	-1.468182	.2042881	-7.19	0.000	-1.869507	-1.066858
_cons	35.4615	2.627905	13.49	0.000	30.29898	40.62402

Remark: .0827444 = .0979356 + .0103469*(-1.468182)

$$\hat{\gamma}_1 = \hat{\beta}_1 + \hat{\beta}_2 \hat{\alpha}_1$$

total movement direct effect Contaminated effect from exper

Measurement Errors

Suppose we want to estimate the impact of X on Y according to the linear regression model

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

Measurement Error in a Regressor

- We **cannot measure X precisely but with an error**. That is we can only observe and use \tilde{X} as the regressor, where

$$\tilde{X} = X + \epsilon; \quad E[\epsilon|X] = 0; \quad E[\epsilon|Y] = 0; \quad E[\epsilon U] = 0$$

- So, in term of \tilde{X} , your true model is

$$Y = \beta_0 + \beta_1(\tilde{X} - \epsilon) + U = \beta_0 + \beta_1 \tilde{X} + \underbrace{(U - \beta_1 \epsilon)}_{\substack{\text{error term} \\ \text{if regressing } Y \text{ on} \\ \tilde{X} \text{ instead of } X}}$$

$$\text{wage} = \beta_0 + \beta_1 \text{edu} + \text{controls} + (A + \varepsilon)$$

left in error term bcs it cannot be captured
 cannot be separate from ε
 innate capabilities of people — cause endogeneity
 ↓
 solved by use other than OLS

$$y = \beta_0 + \beta_1 x + u$$

↑ brokerage account size
 ↑ income (cannot be observed precisely)
 only observe reported income = income + noise (measurement error)
 cannot really regress y on x
 can only reg y on \tilde{x}
 try to write in form of y & \tilde{x}

$$\tilde{x} = x + \varepsilon$$

observed not not
 ↑ not

$$y = \beta_0 + \beta_1 (\tilde{x} - \varepsilon) + u$$

$$y = \beta_0 + \beta_1 \tilde{x} + (u - \beta_1 \varepsilon)$$

cov(\tilde{x}, ε) ≠ 0!!! ⇒ endogeneity problem
 impact of reported income on y ⇒ cheating
 not income

Measurement error in x

reg y on \tilde{x}

$$\hat{\beta}_{OLS} = \frac{\text{cov}(\hat{y}, \tilde{x})}{\text{var}(\tilde{x})} \xrightarrow{\text{converge to}} \frac{\text{cov}(y, \tilde{x})}{\text{var}(\tilde{x})}$$

$$\frac{\beta_1 \text{var}(x) - \beta_1 \text{cov}(x, \varepsilon)}{\text{var}(x)}$$

$$\beta_1 - \beta_1 \left(\frac{\text{cov}(x, \varepsilon)}{\text{var}(x)} \right)$$

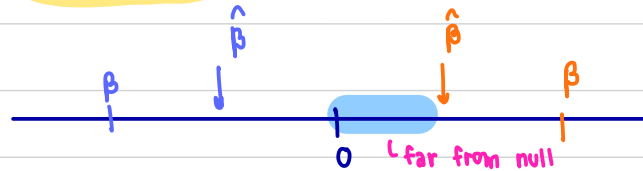
if $\text{cov}(x, \varepsilon) = 0$
 what remains:
 $\rightarrow \text{var}(\varepsilon)$
 $\rightarrow \text{var}(x) + \text{var}(\varepsilon)$

$$\beta_1 - \beta_1 \left(\frac{\text{var}(\varepsilon)}{\text{var}(x) + \text{var}(\varepsilon)} \right) \rightarrow \beta_1 \left(1 - \frac{\text{var}(\varepsilon)}{\text{var}(x) + \text{var}(\varepsilon)} \right) = \beta_1 \left(\frac{\text{var}(x) + \text{var}(\varepsilon) - \text{var}(\varepsilon)}{\text{var}(x) + \text{var}(\varepsilon)} \right)$$

$$= \beta_1 \left(\frac{\text{var}(x)}{\text{var}(x) + \text{var}(\varepsilon)} \right)$$

> 0
 < 0

converge to just a fraction of β



"Biased toward zero" \rightarrow estimate is closer to zero than what it is suppose to be

want to do inference

$$H_0: \beta = 0 \quad \text{vs} \quad H_1: \beta \neq 0$$

reject H_0 with $\hat{\beta}_{OLS}$ \rightarrow don't have to move to use diff estimator

(safely say that we reject this in population as well as we know the direction in bias)

safely make conclusion that $\beta \neq 0$

if not reject H_0 with $\hat{\beta}_{OLS}$ \rightarrow have to use new estimator
 not easy!
 cannot make any conclusion whether we can reject or not reject

- Since we use \tilde{X} instead of X as the regressor, our OLS is

$$\begin{aligned}\hat{\beta}_1 &= \frac{\hat{\sigma}_{\tilde{X}Y}}{\hat{\sigma}_{\tilde{X}}^2} \xrightarrow{P} \frac{\text{cov}(\tilde{X}, Y)}{\text{Var}(\tilde{X})} = \frac{\text{cov}(\tilde{X}, \beta_0 + \beta_1 \tilde{X} + (U - \beta_1 \epsilon))}{\text{Var}(\tilde{X})} \\ &= \beta_1 + \frac{\text{cov}(X + \epsilon, U - \beta_1 \epsilon)}{\text{Var}(X + \epsilon)} \\ &= \beta_1 - \beta_1 \left(\frac{\text{Var}(\epsilon)}{\text{Var}(X) + \text{Var}(\epsilon)} \right) \\ &= \beta_1 \left(\frac{\text{Var}(X)}{\text{Var}(X) + \text{Var}(\epsilon)} \right)\end{aligned}$$

- Notice that the parenthesis is always greater than zero and less than one. So, our OLS estimator is inconsistent and **biased toward zero**.
- If the error term has a larger variance relative to the correct measure of regressor X , the estimate tend to get closer to zero

→ not have much bias → get closer to β_1
 → get whole thing close to zero
 if large
 if this is close to zero
 if data is noisy
 closer to zero
 cannot detect any effect of x on y

Measurement Error in the Regressand

- We cannot measure Y precisely but observe \tilde{Y} for use as the regressand, where

$$\tilde{Y} = Y + \eta; \quad E[\eta|X] = 0; \quad E[\eta|Y] = 0; \quad E[\eta U] = 0$$

- So, in term of \tilde{Y} , our true model is

$$\tilde{Y} = \beta_0 + \beta_1 X + \underbrace{(U - \eta)}_{\substack{\text{error term} \\ \text{if regressing} \\ \tilde{Y} \text{ instead of } Y \text{ on } X}}$$

- Then, our OLS estimator is

$$\hat{\beta}_1 = \frac{\hat{\sigma}_{X\tilde{Y}}}{\hat{\sigma}_X^2} \xrightarrow{P} \frac{\text{cov}(X, \tilde{Y})}{\text{Var}(X)} = \frac{\text{cov}(X, \beta_0 + \beta_1 X + (U - \eta))}{\text{Var}(X)}$$

- So, in this case, our **OLS is still consistent under the regular assumption that $E[XU] = 0$.**

$$y = \beta_0 + \beta_1 x + u$$

↑
port size
↑
income

if this is noisy problem ← cannot be observed

$$\tilde{y} = y + \eta$$

(reported)
← measurement error in y

Need to run reg \tilde{y} x

$$\tilde{y} - \eta = \beta_0 + \beta_1 x + u$$

error term

$$\tilde{y} = \beta_0 + \beta_1 x + u + \eta$$

not correlated
typically bcs of diff r.v.

if η is not related to x

$$\text{cov}(x, \eta) = 0 \Rightarrow \text{no endogeneity problem}$$

Reversed causality

$$y = \beta_0 + \beta_1 x + u$$

reg x y \Rightarrow coefficient will not converge to $\frac{1}{\beta_1}$ \Rightarrow cannot make conclusion that $\beta = \frac{10}{7}$

$$0.7 = \frac{7}{10}$$

rewrite
 x on LHS
 y on RHS

$$\beta_1 x = -\beta_0 + y - u$$

u and y are correlated \rightarrow when $u \uparrow$, $y \uparrow$

$$x = \frac{-\beta_0}{\beta_1} + \left(\frac{1}{\beta_1}\right)y - \frac{u}{\beta_1}$$

endogeneity problem

\downarrow just use y as regressand

may happen with

Endogeneity

- Omitted Variable Bias

control as another regressor in model

- Measurement Errors

in x : biased toward zero

- Simultaneous Equations

in y : not a problem

happen at the same time ; x & y satisfy various equations at the same time

Reversed causality problem \Rightarrow switch back



require instrument to solve this

Simultaneous Equations

- Suppose we have two equations:

$$Y = \beta_0 + \beta_1 X + U \text{ and } Y = \gamma_0 + \gamma_1 X + \eta$$

- If **both equations hold simultaneously**, we can solve for X from the equation: $\beta_0 + \beta_1 X + U = Y = \gamma_0 + \gamma_1 X + \eta$ and get

$$X = -\frac{\beta_0 - \gamma_0}{\beta_1 - \gamma_1} + \frac{\eta}{\beta_1 - \gamma_1} - \frac{U}{\beta_1 - \gamma_1}$$

- Then,

$$\text{cov}(X, \eta) = \text{cov}\left(\frac{\eta}{\beta_1 - \gamma_1}, \eta\right) = \frac{\text{Var}(\eta)}{\beta_1 - \gamma_1} \neq 0$$

$$\text{cov}(X, U) = \text{cov}\left(-\frac{U}{\beta_1 - \gamma_1}, U\right) = -\frac{\text{Var}(U)}{\beta_1 - \gamma_1} \neq 0$$

- Therefore, if we only observe (X, Y) when both equations hold simultaneously, OLS will neither be a consistent estimator for the regression $Y = \beta_0 + \beta_1 X + U$ nor $Y = \gamma_0 + \gamma_1 X + \eta$.

- This problem is often present when estimating demand or supply curves, *if you use transaction data $\Rightarrow Q^{\text{supply}} = Q^{\text{demand}} \Rightarrow Q \text{ \& p satisfy both equations at the same time}$*

$$Q^{\text{Demand}} = \beta_0 + \beta_1 P + U$$

$$Q^{\text{Supply}} = \gamma_0 + \gamma_1 P + \eta$$

- We would expect $\beta_1 < 0$ and $\gamma_1 > 0$.
- However, available data are typically transactions occurred at market equilibria where price equates quantity demanded and supplied:

$$Q^{\text{Demand}} = Q^{\text{equilibrium}} = Q^{\text{Supply}}$$

$$p^{\text{Demand}} = p^{\text{equilibrium}} = p^{\text{Supply}}$$

- With equilibrium transaction data only, we would see the estimated coefficient of price between β_1 and γ_1 .

$$q^{\text{demand}} = \beta_0 + \beta_1 x + u \quad \xrightarrow{\text{correlated}} \quad = \gamma_0 + \gamma_1 x + \eta \quad \xrightarrow{\text{correlated}} \quad = q^{\text{supply}}$$

$$(\beta_1 - \gamma_1)x = -\beta_0 + \gamma_0 + \eta - u$$

$$x = - \left(\frac{\beta_0 - \gamma_0}{\beta_1 - \gamma_1} \right) + \frac{1}{\beta_1 - \gamma_1} (\eta - u) \quad \text{cannot use OLS estimator!}$$

When people use more credit card → people will use less cash

1st equation E.g. Cash holding in economy

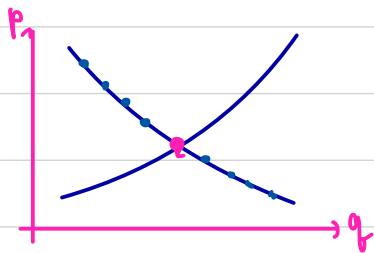
$$= \beta_0 + \beta_1 \text{card} + \text{controls} + u$$

⊖ usage

can also interpret as causation as well !!!
no cash → force to use credit card

2nd equation :

$$\text{card} = \gamma_0 + \gamma_1 \text{cash} + \text{controls} + \eta$$



1st solution: move from transaction data, to survey data e.g. find willingness to pay for a given price
very carefully decide this,
not straight forward

Resolution for Endogeneity

- Recall that we use the moment condition $E[U] = E[XU] = 0$ to solve for β coefficients.
- and use the sample-counterpart conditions from FOCs $\sum_{i=1}^n \hat{U}_i = \sum_{i=1}^n (X_i \hat{U}_i) = 0$ to solve for the OLS estimator.
- E.g. for bivariate case: $\beta_1 = \frac{\text{cov}(X,Y)}{\text{var}(X)} \Rightarrow \hat{\beta}_1^{OLS} = \frac{\hat{\sigma}_{XY}}{\hat{\sigma}_X^2}$
- With the endogeneity problem, we do not have $E[XU] = 0$.
- A resolution for this is to **find another moment condition to solve for the β and construct a consistent estimator from a sample counterpart of the new solution.**

Instrument need to find $z \rightarrow$ not recommend

- A valid **instrument** or **instrumental variable**, Z , is a random variable satisfying:
 - **Instrumental exogeneity**: it is uncorrelated with the error term

$$\underline{cov(Z, U) = 0 \text{ or } E[ZU] = 0}$$

- **Instrumental relevance**: relevant to X it is relevance by still correlated with the endogenous variable X

$$\underline{cov(Z, X) \neq 0}$$

- We can use a valid instrument to help estimate a model with endogeneity problem.

OLS estimator

start with assumption

- $y = \beta x + u$
- $E[u] = 0$ for all interpretation
- ~~$E[xu] = 0$~~ X if endogeneity

↑
 $E[u] = 0$

} solve for $\beta \Rightarrow$

~~$\bar{y} \quad \hat{\beta}_1 \quad \bar{x}$ use sample counterpart
converge to population~~
 ~~$\beta_0 = E[y] - \beta_1 E[x]$~~
 ~~$\beta_1 = \frac{\text{cov}(y, x)}{\text{var}(x)}$~~
 ~~\hat{v}^2 formula change~~

Instrumental Variable Estimator

A resolution to the endogeneity problem is to use the moment conditions related to instruments to solve for the coefficient parameter β 's and find a sample counterpart as the estimator.

Bivariate Model with One Endogenous Variable and One Instrument

- Use $E[U] = 0$ and $E[ZU] = 0$ instead of $E[XU] = 0$ to solve for β_0 and β_1 :

$$\star E[U] = 0 \Rightarrow E[Y - (\beta_0 + \beta_1 X)] = 0 \Rightarrow \beta_0 = E[Y] - \beta_1 E[X]$$

$$\star E[ZU] = 0 \Rightarrow 0 = E[Z(Y - \beta_0 - \beta_1 X)]$$

Use this to solve for $\beta_1 = E[Z(Y - E[Y] + \beta_1 E[X] - \beta_1 X)]$

$$\Rightarrow \beta_1 = \frac{E[Z(Y - E[Y])]}{E[Z(X - E[X])]} = \frac{\text{cov}(Z, Y)}{\text{cov}(Z, X)}$$

still need $\text{cov}(Z, X) \neq 0$ if not, we cannot find β

z : Instrumental variable of estimator

- Hence, we can get for the sample counterpart,

$$\hat{\beta}_1^{IV} = \frac{\hat{\sigma}_{ZY}}{\hat{\sigma}_{ZX}} \xrightarrow[\text{WLLN \& Continuous mapping theorem}]{\text{converge to}} \frac{\text{cov}(z, y)}{\text{cov}(z, x)}$$

$$\hat{\beta}_0^{IV} = \bar{Y}_n - \left(\frac{\hat{\sigma}_{ZY}}{\hat{\sigma}_{ZX}} \right) \bar{X}_n$$

- This is called the **instrumental variable (IV) estimator**.
- Under the regular conditions for WLLN, the IV estimator for β_1 and β_0 is consistent because

$$\hat{\beta}_1^{IV} = \frac{\hat{\sigma}_{ZY}}{\hat{\sigma}_{ZX}} \xrightarrow{P} \frac{\text{cov}(Z, Y)}{\text{cov}(Z, X)} = \beta_1$$

$$\hat{\beta}_0^{IV} = \bar{Y}_n - \left(\frac{\hat{\sigma}_{ZY}}{\hat{\sigma}_{ZX}} \right) \bar{X}_n \xrightarrow{P} E[Y] - \beta_1 E[X] = \beta_0$$

- However, the IV estimator is typically biased

Why it's hard to find z ?

e.g. $\ln(\text{wage}) = \beta_0 + \beta_1 \text{edu} + \text{controls} + (A + U)$

I need to find z that correlated to edu
 but not A

innate
 capability
 that correlate
 with edu

\downarrow Z
the distance btw — birthplace genius can be borned anywhere
 \ the near river good education places near river

GG!

* Recommend to review a lot of literature

- face endogeneity → what is the cause?
- replicate from other literature

eg. $\ln q_f^D = \beta_0 + \beta_1 \ln p + \text{budget controls} + U$

A hand-drawn diagram illustrating the relationship between various factors. At the top left, 'Z' is written above 'tax, fees' and 'cost of production', which are enclosed in a dashed pink oval. A blue arrow points from 'lifestyle' (written to the right) to the oval. Below the oval, a blue arrow points to the text 'this may not be instrument e.g. ...'. To the right of 'lifestyle', a green arrow points down to a list containing 'behaviors' and 'attitude'.

Z
 tax, fees
 cost of production
 lifestyle
 can correlate to controls
 ↓
 behaviors
 attitude
 this may not be instrument
 e.g. ...

1) Think what's in U?

2) try to think of x -related that avoid such factors in U .

we may have many instruments for one endogenous x

Bivariate Model with One Endogenous Variable **but Several Instruments**

- Suppose that we have m valid instruments Z_1, Z_2, \dots, Z_m

- Then, $cov(Z_1, U) = 0, cov(Z_2, U) = \dots = cov(Z_m, U) = 0,$

$$cov(Z_1, X) \neq 0, cov(Z_2, X) \neq 0, \dots, cov(Z_m, X) \neq 0$$

- We can use either each of Z_1, Z_2, \dots, Z_m or their linear combinations as an instrument in the IV estimator function.

- Caution: if using a linear combination $\alpha_0 + \alpha_1 Z_1 + \alpha_2 Z_2 + \dots + \alpha_m Z_m$, we need to check that

$$\alpha_1 cov(Z_1, X) + \alpha_2 cov(Z_2, X) + \dots + \alpha_m cov(Z_m, X) \neq 0$$

- The most efficient way** is to use the best linear ^{most correlate with x} projection of X given Z_1, Z_2, \dots, Z_m (the second interpretation) as the instrument:

$$X = \underbrace{\pi_0 + \pi_1 Z_1 + \pi_2 Z_2 + \dots + \pi_m Z_m}_{\text{BLP}[X|Z_1, \dots, Z_m]} + \eta$$

want to find best linear approximation of x

Linear combinations:

e.g. 2 valid instruments: $\text{cov}(z_1, x) = \frac{1}{2}$
 $\text{cov}(z_2, x) = \frac{1}{3}$ } if it correlate with $x \rightarrow$ can use it

$$\alpha_1 z_1 + \alpha_2 z_2 = \bar{z}$$

$$\text{cov}(\bar{z}, u) = \alpha_1 \text{cov}(z_1, u) + \alpha_2 \text{cov}(z_2, u)$$

$$\cdot \text{cov}(\bar{z}, x) = \alpha_1 \text{cov}(z_1, x) + \alpha_2 \text{cov}(z_2, x)$$

\downarrow $\frac{1}{2}$ \downarrow $\frac{1}{3}$
 -2 $+3$

$\hat{\beta}^{iv}$

\rightarrow consistent

BIASED \rightarrow may not equal to population

need large enough sample size
 (might be biased)

Weak instrument \rightarrow $\text{cov}(x, z)$ low \rightarrow biased & slow consistent
 \uparrow
 want this to be high

- However, since we don't see the whole population of X, Z_1, \dots, Z_m , we don't know $\pi_0, \pi_1, \dots, \pi_m$ but can use the OLS to consistently estimate them.
- In other words, use $\hat{X}^{OLS} = \hat{\pi}_0 + \hat{\pi}_1 Z_1 + \dots + \hat{\pi}_m Z_m$ as the instrument.
- Then, our IV estimator is

$$\hat{\beta}_0^{IV} = \bar{Y}_n - \hat{\beta}_1^{IV} \bar{X}_n; \quad \hat{\beta}_1^{IV} = \frac{cov(\hat{X}, Y)}{cov(\hat{X}, X)} = \frac{\hat{\sigma}_{\hat{X}Y}}{\hat{\sigma}_{\hat{X}X}} \left\{ \frac{\hat{\sigma}_{\hat{X}Y}}{Var(\hat{x})} \right.$$

- Remarks, by the properties of OLS,

$$\sum_{i=1}^n \hat{\eta}_i = \sum_{i=1}^n (Z_1 \hat{\eta}_i) = \dots = \sum_{i=1}^n (Z_m \hat{\eta}_i) = 0 \Rightarrow \sum_{i=1}^n (\hat{X}_i \hat{\eta}_i) = 0$$

- This implies that the sample covariance of \hat{X} and $\hat{\eta}$ is zero.
- Therefore,

$$\hat{\sigma}_{\hat{X}X} = cov(\hat{X}, X) = cov(\hat{X}, \hat{X} + \hat{\eta}) = Var(\hat{X}) + \underbrace{cov(\hat{X}, \hat{\eta})}_{=0} = \hat{\sigma}_{\hat{X}}^2$$

by property
of OLS estimator

if you \hat{x} is
instrument

2SLS or TSLS Estimator

no one do this → ask STATA to do this

- Because $\hat{\sigma}_{\hat{X}X} = \hat{\sigma}_{\hat{X}}^2$, the IV estimator is equal to

$$\hat{\beta}_1^{IV} = \frac{\hat{\sigma}_{\hat{X}Y}}{(\hat{\sigma}_{\hat{X}X})} = \frac{\hat{\sigma}_{\hat{X}Y}}{(\hat{\sigma}_{\hat{X}}^2)}$$

2 steps least square

prove this is the same

- which is the OLS estimator of β_1 of the regression $Y = \beta_0 + \beta_1 \hat{X} + U$.
- This means that we can also get the estimate by doing the 2 steps of the OLS estimation method:
 - Regress X on Z_1, \dots, Z_m and get the fitted value \hat{X}**
 - Regress Y on \hat{X} to get consistent estimates for β_0, β_1**
- Hence, people call this method as **Two-Step Least Squares (TSLS or 2SLS) estimator**, which gives exactly the same estimated values as the IV estimator with $BLS[X|Z_1, \dots, Z_m]$ as the instrument.

Multivariate Model

One Endogenous Variable, One Instrument

- Suppose that we have a multivariate regression with X_1 as the only endogenous regressor, and the other regressors are exogenous:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_k X_k + U$$

- Again, two conditions are required for Z to be a valid instrument:
 - Instrumental exogeneity:** $Cov(Z, U) = 0$
 - Instrumental relevance:** $Cov(Z, X_1 | X_2, \dots, X_k) \neq 0$
- The instrumental relevance says that Z is required to be correlated with X_1 after controlling for the other exogenous regressors.
- This means π_1 in the following best linear projection cannot be zero:

$$X_1 = \pi_0 + \pi_1 Z + \pi_2 X_2 + \cdots + \pi_k X_k + \eta$$

- So, for this model we have $k + 1$ moment conditions to solve for β :

$$E[U] = E[ZU] = E[X_2U] = \dots = E[X_kU] = 0$$

- Or in matrix notation $E[WU] = 0$, where $W = (1, Z, X_2, \dots, X_k)$.
- As we did in case of the OLS before, we can solve for β

$$E[WU] = 0 \Rightarrow E[W(Y - X'\beta)] = 0 \Rightarrow E[WY] - E[WX']\beta = 0$$

$$\beta = E[WX']^{-1}E[WY]$$

- And the instrumental variable estimator is the sample counterpart:

$$\hat{\beta}^{IV} = \left(\frac{1}{n} \sum_{i=1}^n W_i X_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n W_i Y_i \right)$$

- Similar to the bivariate case, if we use $\hat{X}_1 = \hat{\pi}_0 + \hat{\pi}_1 Z + \hat{\pi}_2 X_2 + \dots + \hat{\pi}_k X_k$ as the instrument instead, $W = (1, \hat{X}_1, X_2, \dots, X_k)$

$$\hat{\beta}^{IV} = \left(\frac{1}{n} \sum_{i=1}^n W_i X_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n W_i Y_i \right) = \left(\frac{1}{n} \sum_{i=1}^n W_i W_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n W_i Y_i \right) = \hat{\beta}^{2SLS}$$

Multivariate Model with Several Instruments

- First, suppose that we have $m > 1$ instruments for a single endogenous regressor X_1 . Again, we can construct an instrument \hat{X}_1 from the OLS fitted value of

$$X_1 = \pi_0 + \gamma_1 Z_1 + \cdots + \gamma_m Z_m + \pi_2 X_2 + \cdots + \pi_k X_k + \eta$$

- Next, if the model have several, say q , endogenous regressors:

$$Y = \beta_0 + \underbrace{\beta_1 X_1 + \cdots + \beta_q X_q}_{\text{endogenous terms}} + \underbrace{\beta_{q+1} X_{q+1} + \cdots + \beta_k X_k}_{\text{exogenous terms}} + U$$

- For this case, **we need to have $m \geq q$ instruments** and find q fitted values, for $1 \leq j \leq q$:

$$X_j = \pi_{j0} + \gamma_{j1} Z_1 + \cdots + \gamma_{jm} Z_m + \pi_{j(q+1)} X_{q+1} + \cdots + \pi_k X_k + \eta$$

easier to do this — no need to do 2 steps

. webuse hsg2, clear
(1980 Census housing data)

. ivregress 2sls y i .region pcturban popden (hsngval = faminc hsg)

rent	$\hat{\beta}$	Coef.	Std. Err.	$\frac{\hat{\beta}}{se(\hat{\beta})}$	P> z	[95% Conf. Interval]	
hsngval		.0041469	.0009073	4.57	0.000	.0023686	.0059252
region							
N Cntrl		-15.666	18.82128	-0.83	0.405	-52.55502	21.22303
South		-5.606589	17.2554	-0.32	0.745	-39.42656	28.21338
west		-70.46023	28.48546	-2.47	0.013	-126.2907	-14.62974
pcturban		-.1464167	.5531242	-0.26	0.791	-1.23052	.9376869
popden		-.0057151	.0038934	-1.47	0.142	-.013346	.0019158
_cons		76.19836	32.39852	2.35	0.019	12.69843	139.6983

Instrumented: hsngval

Instruments: 2.region 3.region 4.region pcturban popden faminc hsg

exogenous variable is instrument by itself

. quietly reg **hsngval** faminc **hsng** i. **region pcturban popden**, nohead

. **predict xhat, xb**

. reg **rent** **xhat** i. **region pcturban popden**, nohead

rent	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
xhat	.0041469	.00041	10.11	0.000	.00332	.0049739
region						
N Cntrl	-15.666	8.505803	-1.84	0.072	-32.81958	1.48759
South	-5.606587	7.798145	-0.72	0.476	-21.33305	10.11987
West	-70.46022	12.87329	-5.47	0.000	-96.42169	-44.49876
pcturban						
popden	-.1464166	.2499706	-0.59	0.561	-.6505303	.3576971
_cons	-.0057151	.0017595	-3.25	0.002	-.0092635	-.0021667
	76.19836	14.64169	5.20	0.000	46.67057	105.7261

↓
 $\hat{\beta}$

wrong

no. of
endogenous variable

q

instrumental variable

m

3 cases
~~~~~

$q$

<

$m$  : over identified

=

just identified

>

under identified

many people use this  
bcs it's hard to find  $z$

cannot be  
estimated

But there's problem:

don't observe  $\rightarrow$  cannot test, need story

$$\text{COV}(U, z) = 0$$

$$\text{COV}(X, z) \neq 0$$

test this because we observe  $\Rightarrow$  test: ?

overidentifying test

# Overidentifying Test

- Recall an instrumental variable,  $Z$ , must satisfy:
  - Instrumental exogeneity:  $cov(Z, U) = 0$  or  $E[ZU] = 0$
  - Instrumental relevance:  $cov(Z, X) \neq 0$
- Note also that we observe samples of  $X, Y, Z$  but we do not observe the error term  $U = Y - \beta X$  because we do not know the parameter  $\beta$ .
- Hence, we could check the instrumental relevance condition by estimating  $cov(Z, X)$  but **we could not estimate  $cov(Z, U)$**  to check the instrumental exogeneity condition.
- However, if the number of candidate instruments we have is larger than the number of endogenous variables (i.e.  $m > q$  or overidentifying case), we can conduct an **Over-identifying Test** to see if we have sufficient valid instruments out of the candidates.

. webuse hsn2, clear  
(1980 Census housing data)

*after run regression*

. Quietly ivregress 2sls rent i.region pcturban popden ///  
(hsngval = faminc hsn2)

. estat overid

Tests of overidentifying restrictions:

Sargan (score) chi2(1) = .036164 (p = 0.8492)

Basmann chi2(1) = .0304 (p = 0.8616)

*Hansen test also popular*

*test statistics*

*p-value of the test*

*fail to reject  $H_0$*

- Both Sargan and Basmann Tests for overidentifying have the **null hypothesis** that there are sufficient valid instruments. *→ null hypothesis*

- Therefore, rejecting the null hypothesis (p-value < 10%, 5%, or 1%) implies that we do not have enough valid instruments.
- Note: degree of freedom of both tests equals to the number of instruments minus the number of endogenous variables in the model

Review from previous class:

$$y = \beta x + u$$

↖ uncorrelated: exogenous  
↘ correlated: endogenous

↑ endogeneity problem  
because of

biased toward zero

← in  $y$   
← in  $x$

- omitted variable
- measurement error
- reversed causality
- simultaneous equations

1<sup>st</sup> put control

2<sup>nd</sup> can you make inference?

3<sup>rd</sup> try change  $x \rightleftharpoons y$

4<sup>th</sup> create new  $z$  (instrument) → cannot test this bcs we cannot observe error term  $u$

①  $\text{cov}(u, z) = 0$  exogeneity

②  $\text{cov}(x, z) \neq 0$  relevance

↓ There're 3 tests

• Overidentifying test #IV > #endo

↳ Hansen test  
↳ Sargan test  
↳ Basman

|| enough valid IV  
 $H_0, H_1$

↓  
we prefer not to reject  $H_0$

$$\text{Wage} = \beta_0 + \beta_1 \text{edu} + \text{control} + (\text{cap} + \overbrace{u}^{\text{error}})$$

↑  
only 1 endogenous variable

↑ emotional  
↑ cannot be captured but correlate to edu  
↑ not count this

IV estimator = 2SLS estimator