

## Lecture 8: Instrumental Variables Part I

Wages are a function of schooling:

$$Y_{si} \equiv f(s_i)$$

$$Y_{si} = \alpha + \rho s_i + \eta_i \quad \text{model we estimate}$$

$$\eta_i = \underbrace{A_i}' \gamma + v_i$$

omitted variable: ability.

OLS:

(1)  $\varepsilon_i \sim iid(0, \sigma^2) \Rightarrow E(\varepsilon) = 0, \text{Var}(\varepsilon) = \sigma^2 \cdot I_{n \times n}$

~~(2)~~  $E(\varepsilon_i | X_{ik}) = 0$  for all  $k$

(3)  $X$  has full column rank  $K$

not satisfied.  $\text{Corr}(\eta_i, s_i) \neq 0$  because  
 $\text{Corr}(A_i, s_i)$  likely  $> 0$ .

We could determine causal effect of  $s_i$  on  $Y_i$  using OLS only if we had all the variables that determined  $Y_i$  and could be correlated with  $s_i$ .

### Propensity Score Matching

- Valid when  $X_i$  is binary  $\{0,1\}$ .  
Cannot use with a continuous variable like years of schooling

### Regression Discontinuity:

- Valid when  $X_i$  is binary  $\{0,1\}$  and there's some discontinuous rule:

$$X_j \geq \bar{X}_j \Rightarrow X_i = 1$$

$$X_j < \bar{X}_j \Rightarrow X_i = 0.$$

- Cannot use with a continuous variable and when there's no known rule that affects continuous variable.

### What's next?

INSTRUMENTAL  
VARIABLES

## Instrumental Variables:

$$y_i = X_i' \beta + \varepsilon_i$$

What is an I.V.?

A variable (or variables)  $z$  that:

"strong"

{ (1) are strongly correlated with  
endogenous  $X$

"valid"

{ (2) do not otherwise affect  $Y$   
except through  $X$ .

$$\hookrightarrow \text{Cov}(\varepsilon_i, z_i) = 0$$

What is our structural relation of interest:

$$(1) \quad Y_i = \alpha + \rho s_i + \eta_i \quad (\text{structural})$$

→ we cannot estimate this by OLS,  
but we have an instrument  $z_i$   
that affects  $s_i$  but does not otherwise  
affect  $Y_i$

$$(2) \quad s_i = \pi_0 + \pi_1 z_i + u_i \quad (\text{first stage})$$

Substitute (2) into (1):

$$(3) \quad Y_i = \alpha + \rho(\pi_0 + \pi_1 z_i + u_i) + \eta_i$$
$$= (\alpha + \rho \pi_0) + \rho \pi_1 z_i + (\rho u_i + \eta_i)$$

$$Y_i = \pi_{20} + \pi_{21} z_i + \varepsilon_i \quad (\text{reduced form})$$

Note that:

$$\frac{\text{red. form}}{\text{first stage}} = \frac{\pi_{21}}{\pi_{11}} = \frac{\frac{\text{Cov}(Y_i, z_i)}{\text{Var}(z_i)}}{\frac{\text{Cov}(S_i, z_i)}{\text{Var}(z_i)}} = \frac{\text{Cov}(Y_i, z_i)}{\text{Cov}(S_i, z_i)} = \frac{\rho \pi_{11}}{\pi_{11}} = \boxed{\rho}$$

we obtained our causal estimate of  $\rho$ !

[This is known as "direct least squares: DLS"]

What happens when more than one instrument?

$$(2a) \quad S_i = \pi_{i0} + \pi_{i1} z_{i1} + \pi_{i2} z_{i2} + u_i$$

$$(3a) \quad Y_i = \alpha + \rho(\pi_{i0} + \pi_{i1} z_{i1} + \pi_{i2} z_{i2} + u_i) + \eta_i$$

$$\begin{aligned} &= \underbrace{(\alpha + \rho\pi_{i0})}_{\pi_{20}} + \underbrace{\rho(\pi_{i1} z_{i1} + \pi_{i2} z_{i2})}_{\rho \left[ \begin{smallmatrix} \text{a weighted} \\ \text{avg. of} \\ \text{instruments.} \end{smallmatrix} \right]} + \underbrace{(\rho u_i + \eta_i)}_{\xi_i} \\ &= \pi_{20} + \rho \left[ \begin{smallmatrix} \text{a weighted} \\ \text{avg. of} \\ \text{instruments.} \end{smallmatrix} \right] \xi_i \end{aligned}$$

can no longer do OLS:  $\rho = \frac{\pi_{21}}{\pi_{11}} = \frac{\rho \pi_{11}}{\pi_{11}}$

This is where 2SLS comes in. From (2a), calculate  $\hat{S}_i$ .

Estimate:  $Y_i = \tilde{\alpha} + \tilde{\rho} \hat{S}_i + \tilde{\eta}_i$

Output from 2SLS:

F-statistic.

We need to make sure that  $\text{corr}(x_i, \tilde{z}_i) \neq 0$ .  
Why?

$$\rho = \frac{\pi_{21}}{\pi_{11}} = \frac{\cancel{\rho} \pi_{11}}{\pi_{11}} \quad \text{If } \pi_{11} \approx 0, \text{ then we have } \frac{0}{0}.$$

undefined.

What is an F-statistic?

$$s_i = \pi_0 + \pi_{11} z_{i1} + \pi_{21} z_{i2} + u_i$$

Joint test that  $\pi_{11}$  or (not and)  $\pi_{21}$   
not equal to zero.

Rule of thumb: F-statistic  $\geq 10$ .

The larger, the better.

## Underidentification Test

Suppose  $K$  endog. variables:  $X_1, \dots, X_K$ .

Need at least  $K$  instruments, one for each endog. variable.

F-stat will tell you that  $\pi_{11}, \dots, \pi_{1K}$  jointly different from zero, but we also want to know that each  $\pi_{11}, \dots, \pi_{1K}$  are different from zero and that one of  $z_k$  is not a linear combination of the others.

Hansen J-Statistic. (Sargan when errors not robust).

For over-identified case. (more instruments than endogenous regressors).

*Implication: if instruments are valid, then the OLS estimates should be unbiased and efficient.*

- (1) Estimate equation using IV.
- (2) If I.V. valid, should obtain "good" parameter estimates and "good" estimates of original errors.
- (3) Instruments should be uncorrelated with errors. Regress:

$$\hat{\varepsilon}_i = \alpha + \beta_1 z_1 + \beta_2 z_2 + \underbrace{X_i}_{\text{exogenous}}$$

Perform F-test to see if  $\beta_1$  and  $\beta_2$  are jointly different from zero.

$$NR^2 \sim \text{chi sq}(\underbrace{z's - \text{endog. } x's}_{\text{degrees of freedom}})$$