

Lecture 10

Instrumental variables

Lectures in SDPE: *Econometrics I* on February 26, 2024

Markus Jäntti
Swedish Institute For Social Research
Stockholm University

Introduction

- We have endogeneity if

$$Y = X'\beta + e, \text{ and } E[Xe] \neq 0. \quad (1)$$

- This is meaningful only relative to more structure on the problem than that provided by just the linear projection interpretation.
- That is, β needs to have a *structural* interpretation.
- Substantively, we should distinguish between different types of endogeneity, such as:
 - measurement errors in regressors
 - two-way causation
 - simultaneity
 - omitted variables
- Much of this chapter will devoted to the “classical” IV setup in economics involving simultaneous equations.
- Much of the *estimation* and *asymptotic distribution* for IV will be covered under generalized method of moments.

Measurement error in regressors

- Let (Y, \mathbf{X}^*) be the random variables of interest and the object of study the conditional expectation $E[Y|\mathbf{X}^*] = \mathbf{X}^{*'}\beta$.
- We observe an error-ridden version $\mathbf{X} = \mathbf{X}^* + \mathbf{u}$ with the $k \times 1$ vector of measurement errors \mathbf{u} , independent of Y and \mathbf{X}^* .
- Our regression is

$$Y = \mathbf{X}^{*'}\beta + e = (\mathbf{X} - \mathbf{u})'\beta + e = \mathbf{X}'\beta + (e - \mathbf{u}'\beta) = \mathbf{X}'\beta + v \quad (2)$$

- Our moment condition is

$$E[Xv] = E[(\mathbf{X}^* + \mathbf{u})(e - \mathbf{u}'\beta)] = -E[\mathbf{u}\mathbf{u}']\beta \neq \mathbf{0} \quad (3)$$

as long as $\beta \neq 0$ and $E[\mathbf{u}\mathbf{u}'] \neq \mathbf{0}$.

Measurement error in regressors

- We then have

$$\beta^* = (E[XX'])^{-1}E[XY] = \beta - (E[XX'])^{-1}E[\mathbf{u}\mathbf{u}']\beta \neq \beta. \quad (4)$$

- Our LS estimator

$$\hat{\beta} = \left(n^{-1} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left(n^{-1} \sum_{i=1}^n \mathbf{X}_i Y_i \right) \quad (5)$$

converges to β^* , not β . This is known as *measurement error bias*.

Supply and demand

- Suppose we are studying supply and demand using

$$\begin{aligned}Q &= -\beta_1 P + e_1, & (\text{demand}) \\Q &= \beta_2 P + e_2, & (\text{supply})\end{aligned}\tag{6}$$

- Let $(e_1, e_2) = e$ be iid and $E[e] = \mathbf{0}$, $E[ee'] = \mathbf{I}_2$ for simplicity and $\beta_1 + \beta_2 = 1$. What happens if we regress Q on P ?
- First, express the equation as

$$\begin{bmatrix} 1 & \beta_1 \\ 1 & -\beta_2 \end{bmatrix} \begin{pmatrix} Q \\ P \end{pmatrix} = \begin{pmatrix} e_1 \\ e_2 \end{pmatrix}\tag{7}$$

Supply and demand

- Solve for (Q, P)

$$\begin{aligned}\begin{pmatrix} Q \\ P \end{pmatrix} &= \begin{bmatrix} 1 & \beta_1 \\ 1 & -\beta_2 \end{bmatrix}^{-1} \begin{pmatrix} e_1 \\ e_2 \end{pmatrix} = \begin{bmatrix} \beta_2 & \beta_1 \\ 1 & -1 \end{bmatrix} \begin{pmatrix} e_1 \\ e_2 \end{pmatrix} \\ &= \begin{pmatrix} \beta_2 e_1 + \beta_1 e_2 \\ e_1 - e_2 \end{pmatrix}\end{aligned}\tag{8}$$

- Projecting Q onto P yields

$$Q = \beta^* P + u, \quad E[Pu] = 0\tag{9}$$

with

$$\beta^* = \frac{E[PQ]}{E[P^2]} = \frac{\beta_2 - \beta_1}{2}\tag{10}$$

A LS estimate of this projection thus converges to $\frac{\beta_2 - \beta_1}{2}$, not β_2 or β_1 .
This is known as *simultaneous equations bias*.

Instrumental Variables

- Consider the linear regression, now a *structural equation*

$$Y = \mathbf{X}'\boldsymbol{\beta} + e = \mathbf{X}'_1\boldsymbol{\beta}_1 + \mathbf{X}'_2\boldsymbol{\beta}_2 + e, \text{E}[\mathbf{X}e] \neq \mathbf{0} \quad (11)$$

where we have endogeneity.

- Now, $\mathbf{X}_1, \boldsymbol{\beta}_1$ are $k_1 \times 1$ and $\mathbf{X}_2, \boldsymbol{\beta}_2$ are $k_2 \times 1$, and are partitioned such that

$$\text{E}[\mathbf{X}_1 e] = \mathbf{0} \quad \text{but} \quad \text{E}[\mathbf{X}_2 e] \neq \mathbf{0} \quad (12)$$

I.e., \mathbf{X}_2 has non-zero covariance with e and is *endogenous* while \mathbf{X}_1 is *exogenous*.

- To get an unbiased/consistent estimate, we need to solve $\boldsymbol{\beta}$ from the sample analog of

$$\text{E}[\mathbf{Z}e] = \text{E}\left[\begin{pmatrix} \mathbf{X}_1 \\ \mathbf{Z}_2 \end{pmatrix} e\right] = \mathbf{0} \quad (13)$$

where $\mathbf{X}_1 = \mathbf{Z}_1$ are the $k_1 \times 1$ *included exogenous regressors* and \mathbf{Z}_2 are the $l_2 \times 1$ *excluded exogenous regressors*.

Instrumental Variables

- If $l = k_1 + l_2 = k$ the model is *just-identified* (or exactly identified), if $l > k$ it is over-identified
- That you *need* $l_2 \geq k_2$ instruments to consistently estimate β does *not* mean they *exist*!

Solving measurement error bias

- Recall Slide 3, equation 3.
- Suppose there exist a set of l instruments ($l \geq k$), \mathbf{Z} , such that
 - it is uncorrelated with the regression error $E[\mathbf{Z}e] = 0$
 - it is uncorrelated with the measurement error $E[\mathbf{Z}\mathbf{u}] = 0$
 - it is correlated with the true regressors $E[\mathbf{Z}\mathbf{X}^{*'}] = \mathbf{Q}$ (where $\mathbf{Q} > 0$)
- An unbiased estimator of β can be constructed using the sample analog of the condition that

$$E[\mathbf{Z}(Y - \mathbf{X}'\beta)] = 0. \quad (14)$$

- The key issue is whether or not such a \mathbf{Z} exists.

Reduced Form

- In contrast to the structural form, we approach estimation of the parameters in the *reduced form* equations using linear projection.
- Consider the two sets of linear projection coefficients:

$$\begin{aligned}\Gamma &= E[\mathbf{Z}\mathbf{Z}']^{-1}E[\mathbf{Z}\mathbf{X}'] \\ \lambda &= E[\mathbf{Z}\mathbf{Z}']^{-1}E[\mathbf{Z}Y]\end{aligned}\tag{15}$$

(Γ is $l \times k$, λ is $l \times 1$.)

- From the first, we have the $k \times 1$ projection error:

$$\mathbf{u} = \mathbf{X} - \Gamma'\mathbf{Z} \Rightarrow \mathbf{X} = \Gamma'\mathbf{Z} + \mathbf{u}, E[\mathbf{Z}\mathbf{u}'] = \mathbf{0}.\tag{16}$$

- Substitute 16 into 11 to get the *reduced form* for Y :

$$Y = \beta'(\Gamma'\mathbf{Z} + \mathbf{u}) + e \Rightarrow Y = \lambda'\mathbf{Z} + v.\tag{17}$$

- The reduced form and structural parameters are related:

$$\lambda = \Gamma\beta,\tag{18}$$

as are the errors:

$$v = \beta'\mathbf{u} + e.\tag{19}$$

Reduced Form

- By construction, we have

$$E[\mathbf{Z}v] = E[\mathbf{Z}\mathbf{u}']\beta + E[\mathbf{Z}e] = 0. \quad (20)$$

- Using sample data, both sets of linear projection coefficients can be estimated by LS:

$$\begin{aligned}\hat{\Gamma} &= (\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{X} \\ \hat{\lambda} &= (\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{Y}\end{aligned} \quad (21)$$

Identification

- *Identification* is about this: can the structural parameter β be recovered from the reduced form coefficients?
- Solve β from eq. 18:

$$\beta = \Gamma^{-1}\lambda, l = k; \quad \beta = (\Gamma'W\Gamma)^{-1}\Gamma'W\lambda, l > k, W > 0. \quad (22)$$

For this to be possible, the *rank condition* must hold:

$$\text{rank}(\Gamma) = k \quad (23)$$

- Note that $\mathbf{Z}' = (\mathbf{X}'_1, \mathbf{Z}'_2)$ and $\mathbf{X}' = (\mathbf{X}'_1, \mathbf{X}'_2)$ so

$$\Gamma = \begin{bmatrix} \Gamma_{11} & \Gamma_{12} \\ \Gamma_{21} & \Gamma_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \Gamma_{12} \\ \mathbf{0} & \Gamma_{22} \end{bmatrix} \quad (24)$$

and the regression of \mathbf{X} on \mathbf{Z} can be re-written as

$$\begin{aligned} \mathbf{X}_1 &= \mathbf{Z}_1 \quad (\text{Note: } \mathbf{Z}_1 \equiv \mathbf{X}_1) \\ \mathbf{X}_2 &= \Gamma'_{12}\mathbf{Z}_1 + \Gamma'_{22}\mathbf{Z}_2 + \mathbf{u} \end{aligned} \quad (25)$$

- The key to identification is the rank of Γ_{22} which must be k_2 to ensure $\text{rank}(\Gamma) = k$.

Estimation

- We treat IV estimation and the distribution of the estimators in greater detail under GMM (chapter 12; lecture 13) but now cover the main estimators.

IV estimator

- With a just-identified model $l = k$, we have the moment condition

$$E[\mathbf{Z}e] = \mathbf{0}. \quad (26)$$

- With $e = Y - \mathbf{X}'\beta$, we have

$$\begin{aligned} E[\mathbf{Z}(Y - \mathbf{X}'\beta)] &= \mathbf{0} \Leftrightarrow \\ E[\mathbf{Z}Y] - E[\mathbf{Z}\mathbf{X}']\beta &= \mathbf{0}, \end{aligned} \quad (27)$$

so

$$\beta = (E[\mathbf{Z}\mathbf{X}'])^{-1}E[\mathbf{Z}Y]. \quad (28)$$

- The Instrumental Variables (IV) estimator $\hat{\beta}_{IV}$ is the “plug-in” estimator

$$\begin{aligned} \hat{\beta}_{IV} &= \left(\frac{1}{n} \sum_{i=1}^n \mathbf{Z}_i \mathbf{X}_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{Z}_i Y_i \right) \\ &= (\mathbf{Z}'\mathbf{X})^{-1}(\mathbf{Z}'\mathbf{Y}). \end{aligned} \quad (29)$$

- To generalize a bit, for any instruments \mathbf{W} , an IV estimator is

$$\hat{\beta}_{IV} = (\mathbf{W}'\mathbf{X})^{-1}(\mathbf{W}'\mathbf{Y}). \quad (30)$$

Two-stage least squares

- The IV estimator is simple to solve from the moment conditions as $l = k$ (just identified).
- With $l \geq k$, this will not work.
- The reduced form can be written as

$$Y = \mathbf{Z}'\mathbf{\Gamma}\beta + v, \quad \mathbf{E}[\mathbf{Z}v] = \mathbf{0}. \quad (31)$$

- Now let $\mathbf{w} = \mathbf{\Gamma}'\mathbf{Z}$ so

$$Y = \mathbf{w}'\beta + v, \quad \mathbf{E}[\mathbf{w}v] = \mathbf{0} \quad (32)$$

Two-stage least squares

- If Γ was known, the “natural” estimator would be

$$\begin{aligned}\widehat{\beta} &= (W'W)^{-1}(W'Y) \\ &= (\Gamma'Z'Z\Gamma)^{-1}(\Gamma'Z'Y).\end{aligned}\tag{33}$$

- Since this is not feasible, we must use $\widehat{\Gamma}$ in its place:

$$\begin{aligned}\widehat{\beta}_{2SLS} &= (\widehat{\Gamma}'Z'Z\widehat{\Gamma})^{-1}(\widehat{\Gamma}'Z'Y) \\ &= (X'Z(Z'Z)^{-1}Z'Z(Z'Z)^{-1}Z'X)^{-1}X'Z(Z'Z)^{-1}Z'Y \\ &= (X'Z(Z'Z)^{-1}Z'X)^{-1}X'Z(Z'Z)^{-1}Z'Y.\end{aligned}\tag{34}$$

- The parameters are estimated by least squares in two stages, hence the name.
- The properties of the 2SLS estimator will be examined in the next lecture.

Endogeneity tests

- The 2SLS estimator is needed as the regressors X_2 in the structural equation are thought to be endogenous ($E[X_2e] \neq 0$).
- This is a testable restriction, with hypotheses

$$\mathbb{H}_0 : E[X_2e] = \mathbf{0} \quad \text{against} \quad \mathbb{H}_0 : E[X_2e] \neq \mathbf{0} \quad (35)$$

- This can be approached using *control functions*.

Control functions

- We can write the structural and first-stage regressions as

$$\begin{aligned} Y &= X_1' \beta_1 + X_2' \beta_2 + e \\ X_2 &= \Gamma_{12}' Z_1 + \Gamma_{22}' Z_2 + \mathbf{u}_2. \end{aligned} \tag{36}$$

- The endogeneity of X_2 means that \mathbf{u}_2 and e are correlated
- Consider then the linear projection of e onto \mathbf{u}_2 :

$$e = \mathbf{u}_2' \alpha + \epsilon; \quad \alpha = (E[\mathbf{u}_2 \mathbf{u}_2'])^{-1} E[\mathbf{u}_2 e]; \quad E[\mathbf{u}_2 \epsilon] = \mathbf{0}. \tag{37}$$

- Substituting this back into the structural equation we have

$$\begin{aligned} Y &= X_1' \beta_1 + X_2' \beta_2 + \mathbf{u}_2' \alpha + \epsilon \\ E[X_1 \epsilon] &= 0; \quad E[X_2 \epsilon] = 0; \quad E[\mathbf{u}_2 \epsilon] = 0 \end{aligned} \tag{38}$$

- Note that X_2 is uncorrelated with ϵ . It's correlation with e is through \mathbf{u}_2 and ϵ is the error after e has been projected orthogonally onto \mathbf{u}_2 .

Control functions

- While we do not observe \mathbf{u}_2 , it can be estimated by the first-stage residual

$$\hat{\mathbf{u}}_2 = \mathbf{X}_2 - \widehat{\mathbf{\Gamma}}'_{12}\mathbf{Z}_1 + \widehat{\mathbf{\Gamma}}'_{22}\mathbf{Z}_2 \quad (39)$$

- The coefficients $(\beta_1, \beta_2, \alpha)$ can be estimated by LS on

$$Y = \mathbf{X}'\beta + \hat{\mathbf{u}}_2\alpha + \hat{\epsilon}. \quad (40)$$

Endogeneity tests

- Since $E[X_2e] = 0$ if and only if $E[\mathbf{u}_2e] = 0$, the above \mathbb{H}_0 can be restated as

$$\mathbb{H}_0 : \alpha = 0 \quad \text{against} \quad \mathbb{H}_1 : \alpha \neq 0.$$

- This can be tested using a Wald statistic; the problem is slightly more involved as the testing relies on the “generated regressor” $\widehat{\mathbf{u}}_2$ rather than the true values (see Hansen 2021, chs 12.26–12.27)
- An alternative approach is to use a so-called Hausman test (see Hansen 2021, ch 9.15)

Identification failure

- Recall from eq. 25

$$X_2 = \Gamma'_{12}Z_1 + \Gamma'_{22}Z_2 + u \quad (41)$$

that if Γ_{22} is not of rank k_2 , identification fails. (The closely related problem of weak instruments – $\Gamma_{22} \neq 0$ but is very small – will be studied in Econometrics II.)

- Suppose we study three scalar variables (Y, X, Z) :

$$\begin{aligned} Y &= X\beta + e \\ X &= Z\gamma + u \end{aligned} \quad (42)$$

- Suppose $\gamma = 0$, so $E[ZX] = 0$. What happens to the IV and LS estimators?
- Assume homoscedastic, unit variance and correlated errors (so X is endogenous),

$$\text{Var}\left(\begin{pmatrix} e \\ u \end{pmatrix} \middle| Z\right) = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}, \quad E[Z^2] = 1. \quad (43)$$

Identification failure

- Study $(Ze, Zu)'$, for which we know by the CLT that

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \begin{pmatrix} Z_i e_i \\ Z_i u_i \end{pmatrix} \xrightarrow{d} \begin{pmatrix} \xi_1 \\ \xi_2 \end{pmatrix} \sim N\left(\mathbf{0}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}\right). \quad (44)$$

- Define $\xi_0 = \xi_1 - \rho\xi_2$, which is normal and $\perp \xi_2$.
- For the (O)LS estimator of β

$$\widehat{\beta} - \beta = \frac{n^{-1} \sum_{i=1}^n u_i e_i}{n^{-1} \sum_{i=1}^n u_i^2} \xrightarrow{p} \rho \neq 0. \quad (45)$$

(I.e., the endogeneity of X renders LS inconsistent.)

- With identification failure, $\gamma = 0$, the asymptotic distribution of the IV estimator is

$$\widehat{\beta}_{IV} - \beta = \frac{\frac{1}{\sqrt{n}} \sum_{i=1}^n Z_i e_i}{\frac{1}{\sqrt{n}} \sum_{i=1}^n Z_i X_i} \xrightarrow{d} \frac{\xi_1}{\xi_2} = \rho + \frac{\xi_0}{\xi_2}. \quad (46)$$

Identification failure

- $\widehat{\beta}_{IV}$ is inconsistent: it converges to a random variable, not the correct constant (or indeed any constant)
- The ratio ξ_0/ξ_2 is symmetrically distributed around zero so $\widehat{\beta}_{IV}$ has median at $\beta + \rho$
- ξ_0/ξ_2 is a ratio of two independent normal variables so follows a Cauchy distribution; it therefore does not have a finite expectation

Identification failure

- For the t-statistic, we need

$$\begin{aligned}\hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n \left(Y_i - X_i \hat{\beta}_{IV} \right)^2 \\ &= \frac{1}{n} \sum_{i=1}^n e_i^2 - \frac{2}{n} \sum_{i=1}^n e_i X_i \left(\hat{\beta}_{IV} - \beta \right) + \frac{1}{n} \sum_{i=1}^n X_i^2 \left(\hat{\beta}_{IV} - \beta \right)^2 \\ &\xrightarrow{d} 1 - 2\rho \frac{\xi_1}{\xi_2} + \left(\frac{\xi_1}{\xi_2} \right)^2.\end{aligned}\tag{47}$$

- Then we have the asymptotic distribution of the t-statistic,

$$T = \frac{\hat{\beta}_{IV} - \beta}{\sqrt{\hat{\sigma}^2 \sum_{i=1}^n Z_i^2 / \sum_{i=1}^n |X_i Z_i|}} \xrightarrow{d} \frac{\xi_1}{\sqrt{1 - 2\rho \frac{\xi_1}{\xi_2} + \left(\frac{\xi_1}{\xi_2} \right)^2}}\tag{48}$$

- This is non-normal. E.g., if $\rho \rightarrow 1$, $\xi_1/\xi_2 \xrightarrow{p} 1$ and $\hat{\sigma}^2 \xrightarrow{p} 0$. As a consequence, the standard error of $\hat{\beta}_{IV}$ converges to zero and the t-statistic converges to ∞ ! so $T \rightarrow \infty$ as $\xi_1/\xi_2 \rightarrow 1$.

References



Hansen, Bruce E (2021). *Econometrics*. Madison, WI: University of Wisconsin.