

ECON 480-3
LECTURE 4: ENDOGENEITY

Ivan A. Canay
Northwestern University



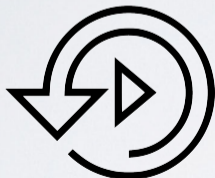
PAST & FUTURE

SO FAR

- ▶ Three Interpretations of β
- ▶ Solving and estimating sub-vectors of β
- ▶ Properties of LS
- ▶ Estimating \mathbb{V}
- ▶ Classical Problems that lead to $E[XU] \neq 0$

TODAY

- ▶ Instrumental Variables
- ▶ The IV Estimator
- ▶ The 2SLS Estimator
- ▶ Properties of 2SLS
- ▶ Estimating \mathbb{V}



INSTRUMENTAL VARIABLES

- ▶ Let (Y, X, U) be a random vector where Y and U take values in \mathbf{R} and X takes values in \mathbf{R}^{k+1} . Assume further that the first component of X is constant and equal to one, i.e., $X = (X_0, X_1, \dots, X_k)'$ with $X_0 = 1$. Let $\beta = (\beta_0, \beta_1, \dots, \beta_k)' \in \mathbf{R}^{k+1}$ be such that

$$Y = X'\beta + U.$$

- ▶ We do not assume $E[XU] = 0$. Any X_j such that $E[X_jU] = 0$ is said to be *exogenous*; any X_j such that $E[X_jU] \neq 0$ is said to be *endogenous*. Normalizing β_0 if necessary, we view X_0 as exogenous.
- ▶ **Instrument**: to overcome the difficulty associated with $E[XU] \neq 0$, we assume that there is an additional random vector Z taking values in $\mathbf{R}^{\ell+1}$ with $\ell + 1 \geq k + 1$ such that $E[ZU] = 0$.
- ▶ Any exogenous component of X is contained in Z (the so-called *included instruments*). In particular, we assume the first component of Z is constant equal to one, i.e., $Z = (Z_0, Z_1, \dots, Z_\ell)'$ with $Z_0 = 1$.
- ▶ We also assume that $E[ZX'] < \infty$, $E[ZZ'] < \infty$ and that there is no perfect collinearity in Z .

INSTRUMENTAL VARIABLES

- ▶ We assume (1) $E[ZU] = 0$, (2) $E[ZX'] < \infty$, (3) $E[ZZ'] < \infty$, and (4) there is no perfect collinearity in Z .
- ▶ The requirement that $E[ZU] = 0$ is termed *instrument exogeneity*.
- ▶ We further assume (5) the rank of $E[ZX']$ is $k + 1$. This is termed *instrument relevance* or *rank condition*.
- ▶ A necessary condition for (5) to be true is $\ell \geq k$. This is referred to as the *order condition*.
- ▶ Using that $U = Y - X'\beta$ and $E[ZU] = 0$, we see that β solves the system of equations

$$E[ZY] = E[ZX']\beta .$$

- ▶ Since $\ell + 1 \geq k + 1$, this may be an *over-determined* system of equations.

A USEFUL LEMMA

LEMMA

Suppose there is no perfect collinearity in Z and let Π be such that $BLP(X|Z) = \Pi'Z$. $E[ZX']$ has rank $k + 1$ if and only if Π has rank $k + 1$. Moreover, the matrix $\Pi'E[ZX']$ is invertible.

SOLVING FOR β

β solves: $E[Z'Y] = E[Z'X]\beta$ or $\Pi'E[Z'Y] = \Pi'E[Z'X]\beta$

Using the previous lemma and $\Pi = E[Z'Z]^{-1}E[Z'X]$, we can derive **three formulae** for β

INTERPRETING THE RANK CONDITION

Interpretation: Consider the case where $k = \ell$ and only X_k is endogenous. Let $Z_j = X_j$ for all $0 \leq j \leq k - 1$. In this case,

$$\Pi' = \begin{pmatrix} 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \\ \pi_0 & \pi_1 & \cdots & \pi_{\ell-1} & \pi_\ell \end{pmatrix}.$$

The rank condition therefore requires $\pi_\ell \neq 0$: the instrument Z_ℓ must be “correlated with X_k after controlling for X_0, X_1, \dots, X_{k-1} .”

PARTITION OF β : ENDOGENOUS COMPONENTS

- ▶ Partition X into X_1 and X_2 , where X_2 is *exogenous*. Partition Z into Z_1 and Z_2 and β into β_1 and β_2 analogously.

- ▶ Note that $Z_2 = X_2$ are *included* instruments and Z_1 are *excluded* instruments. Then,

$$Y = X_1' \beta_1 + X_2' \beta_2 + U .$$

- ▶ We can conveniently re-write this by projecting (BLP) on $Z_2 = X_2$. Consider the case $k = \ell$

$$\text{BLP}(Y|Z_2) = \text{BLP}(X_1|Z_2)' \beta_1 + X_2' \beta_2 .$$

- ▶ Define $Y^* = Y - \text{BLP}(Y|Z_2)$ and $X_1^* = X_1 - \text{BLP}(X_1|Z_2)$ so that

$$E[Z_1 Y^*] = E[Z_1 X_1^{*'}] \beta_1 + E[Z_1 U]$$

- ▶ It follows that

$$\beta_1 = E[Z_1 X_1^{*'}]^{-1} E[Z_1 Y^*] .$$

QUESTIONS?



ESTIMATING β : THE IV ESTIMATOR

- ▶ **Just identified case:** $k = \ell$. Denote by P the marginal distribution of (Y, X, Z) .
- ▶ Let $(Y_1, X_1, Z_1), \dots, (Y_n, X_n, Z_n)$ be an i.i.d. sequence of random variables with distribution P .
- ▶ By analogy with $\beta = E[ZX']^{-1}E[Z Y]$, the natural estimator of β is simply

$$\hat{\beta}_n = \left(\frac{1}{n} \sum_{1 \leq i \leq n} Z_i X_i' \right)^{-1} \left(\frac{1}{n} \sum_{1 \leq i \leq n} Z_i Y_i \right).$$

- ▶ This estimator is called the *instrumental variables (IV)* estimator of β . Note that $\hat{\beta}_n$ satisfies

$$\frac{1}{n} \sum_{1 \leq i \leq n} Z_i (Y_i - X_i' \hat{\beta}_n) = 0.$$

In particular, $\hat{U}_i = Y_i - X_i' \hat{\beta}_n$ satisfies

$$\frac{1}{n} \sum_{1 \leq i \leq n} Z_i \hat{U}_i = 0.$$

THE IV ESTIMATOR

Insight on the IV estimator: assume $X_0 = 1$ and $X_1 \in \mathbf{R}$. An interesting interpretation of the IV estimator of β_1 is obtained by multiplying and dividing by $\frac{1}{n} \sum_{i=1}^n (Z_{1,i} - \bar{Z}_{1,n})^2$, i.e.,

$$\hat{\beta}_{1,n} = \frac{\frac{1}{n} \sum_{i=1}^n (Z_{1,i} - \bar{Z}_{1,n}) Y_i / \frac{1}{n} \sum_{i=1}^n (Z_{1,i} - \bar{Z}_{1,n})^2}{\frac{1}{n} \sum_{i=1}^n (Z_{1,i} - \bar{Z}_{1,n}) X_{1,i} / \frac{1}{n} \sum_{i=1}^n (Z_{1,i} - \bar{Z}_{1,n})^2}$$

THE IV ESTIMATOR: MATRIX NOTATION

This estimator may be expressed more compactly using matrix notation. Define

$$\mathbf{Z} = (Z_1, \dots, Z_n)'$$

$$\mathbf{X} = (X_1, \dots, X_n)'$$

$$\mathbf{Y} = (Y_1, \dots, Y_n)' .$$

In this notation, we have

$$\hat{\beta}_n = (\mathbf{Z}'\mathbf{X})^{-1}(\mathbf{Z}'\mathbf{Y}) .$$

THE TWO-STAGE LEAST SQUARES (TSLS) ESTIMATOR

- ▶ **Over-identified case:** $\ell > k$

- ▶ The expressions we derived for β in this case, like

$$\beta = E[\Pi' E[ZX']]^{-1} \Pi' E[Z Y],$$

all involved the matrix Π , where

$$\text{BLP}(X|Z) = \Pi' Z.$$

- ▶ An estimate of Π can be obtained by OLS.
- ▶ Since $\Pi = E[ZZ']^{-1} E[ZX']$, a natural estimator of Π is

$$\hat{\Pi}_n = \left(\frac{1}{n} \sum_{1 \leq i \leq n} Z_i Z_i' \right)^{-1} \left(\frac{1}{n} \sum_{1 \leq i \leq n} Z_i X_i' \right).$$

THE TWO-STAGE LEAST SQUARES (TSLS) ESTIMATOR

$$\text{Let } X_i = \hat{\Pi}'_n Z_i + \hat{V}_i \quad \text{where} \quad \hat{\Pi}_n = \left(\frac{1}{n} \sum_{1 \leq i \leq n} Z_i Z_i' \right)^{-1} \left(\frac{1}{n} \sum_{1 \leq i \leq n} Z_i X_i' \right).$$

With this estimator of Π , a natural estimator of β is simply

THE TSLS ESTIMATOR

- ▶ Note that $\hat{\beta}_n$ satisfies

$$\frac{1}{n} \sum_{1 \leq i \leq n} \hat{\Pi}'_n Z_i (Y_i - X_i' \hat{\beta}_n) = 0.$$

- ▶ In particular, $\hat{U}_i = Y_i - X_i' \hat{\beta}_n$ satisfies

$$\frac{1}{n} \sum_{1 \leq i \leq n} \hat{\Pi}'_n Z_i \hat{U}_i = 0.$$

- ▶ This implies that \hat{U}_i is orthogonal to all of the instruments equal to an exogenous regressors, but may not be orthogonal to the other regressors.

- ▶ It is termed the TSLS estimator because it may be obtained in the following way:

① regress (each component of) X_i on Z_i to obtain $\hat{X}_i = \hat{\Pi}'_n Z_i$;

② regress Y_i on \hat{X}_i to obtain $\hat{\beta}_n$. However, in order to obtain proper standard errors, it is recommended to compute the estimator in one step (see the following section).

THE TSLS ESTIMATOR: MATRIX NOTATION

This estimator may be expressed more compactly using matrix notation. Define

$$\begin{aligned}\mathbf{Z} &= (Z_1, \dots, Z_n)' \\ \mathbf{X} &= (X_1, \dots, X_n)' \\ \mathbf{Y} &= (Y_1, \dots, Y_n)' \\ \hat{\mathbf{X}} &= (\hat{X}_1, \dots, \hat{X}_n)' \\ &= \mathbf{P}_Z \mathbf{X},\end{aligned}$$

where

$$\mathbf{P}_Z = \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'$$

is the projection matrix onto the column space of \mathbf{Z} . In this notation, we have

$$\begin{aligned}\hat{\beta}_n &= (\hat{\mathbf{X}}'\mathbf{X})^{-1}(\hat{\mathbf{X}}'\mathbf{Y}) \\ &= (\hat{\mathbf{X}}'\hat{\mathbf{X}})^{-1}(\hat{\mathbf{X}}'\mathbf{Y}) \\ &= (\mathbf{X}'\mathbf{P}_Z\mathbf{X})^{-1}(\mathbf{X}'\mathbf{P}_Z\mathbf{Y}).\end{aligned}$$

QUESTIONS?



PROPERTIES OF TWO-STAGE LEAST SQUARES

- ▶ Let (Y, X, U) be a random vector where Y and U take values in \mathbf{R} and X takes values in \mathbf{R}^{k+1} . Assume further that the first component of X is constant and equal to one, i.e., $X = (X_0, X_1, \dots, X_k)'$ with $X_0 = 1$. Let $\beta = (\beta_0, \beta_1, \dots, \beta_k)' \in \mathbf{R}^{k+1}$ be such that

$$Y = X'\beta + U.$$

- ▶ We assume ① $E[ZU] = 0$, ② $E[ZX'] < \infty$, ③ $E[ZZ'] < \infty$, and ④ there is no perfect collinearity in Z , and ⑤ the rank of $E[ZX']$ is $k + 1$
- ▶ Let $(Y_1, X_1, Z_1), \dots, (Y_n, X_n, Z_n)$ be an i.i.d. sequence of random variables with distribution P .
- ▶ Under these assumptions the TSLS estimator is **consistent** for β , and under the additional requirement that $\text{Var}[ZU] < \infty$, it is **asymptotically normal** with limiting variance

$$\mathbb{V} = E[\Pi'ZZ'\Pi]^{-1}\Pi'\text{Var}[ZU]\Pi E[\Pi'ZZ'\Pi]^{-1}.$$

CONSISTENCY OF TSLS

$$\hat{\beta}_n = \left(\hat{\Pi}'_n \left(\frac{1}{n} \sum_{1 \leq i \leq n} Z_i X_i' \right) \right)^{-1} \hat{\Pi}'_n \left(\frac{1}{n} \sum_{1 \leq i \leq n} Z_i Y_i \right) \xrightarrow{P} \beta \text{ as } n \rightarrow \infty .$$

ASYMPTOTIC NORMALITY OF TSLS

Assume that $\text{Var}[ZU] = E[ZZ'U^2] < \infty$. Then, as $n \rightarrow \infty$,

$$\sqrt{n}(\hat{\beta}_n - \beta) \xrightarrow{d} N(0, \mathbb{V}) .$$

ESTIMATION OF \mathbb{V}

A natural estimator of \mathbb{V} is given by

$$\hat{\mathbb{V}}_n = \left(\hat{\Pi}'_n \left(\frac{1}{n} \sum_{1 \leq i \leq n} Z_i Z_i' \right) \hat{\Pi}_n \right)^{-1} \times \hat{\Pi}'_n \left(\frac{1}{n} \sum_{1 \leq i \leq n} Z_i Z_i' \hat{U}_i^2 \right) \hat{\Pi}_n \times \left(\hat{\Pi}'_n \left(\frac{1}{n} \sum_{1 \leq i \leq n} Z_i Z_i' \right) \hat{\Pi}_n \right)^{-1},$$

where $\hat{U}_i = Y_i - X_i' \hat{\beta}_n$.

- ▶ **Primary difficulty** in establishing the consistency of this estimator lies in showing that

$$\frac{1}{n} \sum_{1 \leq i \leq n} Z_i Z_i' \hat{U}_i^2 \xrightarrow{P} \text{Var}[ZU]$$

as $n \rightarrow \infty$. The complication lies in the fact that we do not observe U_i and therefore have to use \hat{U}_i .

- ▶ However, the desired result can be shown by arguing exactly as in the second part of this class.
- ▶ **Note:** $\hat{U}_i = Y_i - X_i' \hat{\beta}_n \neq Y_i - \hat{X}_i' \hat{\beta}_n$, so the standard errors from two repeated applications of OLS will be incorrect.

QUESTIONS?

