

# Introductory Econometrics

## Endogeneity and Instrument Variable Estimation

---

Yaohan Chen

School of Big Data and Statistics, Anhui University

Spring, 2025

# Endogeneity

---

- We say that there is *endogeneity* in the linear regression model

$$Y_i = \mathbf{X}_i' \boldsymbol{\beta} + u_i$$

if  $\boldsymbol{\beta}$  is the parameter of interest and

$$E(\mathbf{X}_i u_i) \neq 0.$$

- Sources of endogeneity.
  - Omitted Variable Bias
  - Simultaneous Equations Bias
  - Measurement error in regressors

# Endogeneity Issue: Omitted Variable Bias

---

- Suppose the true regression model is

$$Y_i = \mathbf{X}_{i1}'\boldsymbol{\beta}_1 + \mathbf{X}_{i2}'\boldsymbol{\beta}_2 + u_i$$

where  $\mathbf{X}_{i1}$  and  $\mathbf{X}_{i2}$  are  $k_1 \times 1$  and  $k_2 \times 1$  vectors respectively.

- If we omit  $\mathbf{X}_{i2}$  then OLS estimator for  $\boldsymbol{\beta}_1$  is biased and inconsistent.

# Endogeneity Issue: Omitted Variable Bias

- Suppose the true regression model is

$$Y_i = \mathbf{X}'_{i1}\boldsymbol{\beta}_1 + \mathbf{X}'_{i2}\boldsymbol{\beta}_2 + u_i$$

where  $\mathbf{X}_{i1}$  and  $\mathbf{X}_{i2}$  are  $k_1 \times 1$  and  $k_2 \times 1$  vectors respectively.

- If we omit  $\mathbf{X}_{i2}$  then OLS estimator for  $\boldsymbol{\beta}_1$  is biased and inconsistent.
- This because  $\mathbf{X} = [\mathbf{X}^{(1)}, \mathbf{X}^{(2)}]$

$$\begin{aligned}\tilde{\boldsymbol{\beta}}_1 &= \left(\mathbf{X}^{(1)'}\mathbf{X}^{(1)}\right)^{-1}\mathbf{X}^{(1)'}\mathbf{Y} \\ &= \boldsymbol{\beta}_1 + \left(\mathbf{X}^{(1)'}\mathbf{X}^{(1)}\right)^{-1}\mathbf{X}^{(1)'}\mathbf{X}^{(2)}\boldsymbol{\beta}_2 + \left(\mathbf{X}^{(1)'}\mathbf{X}^{(1)}\right)^{-1}\mathbf{X}^{(1)'}\mathbf{u}.\end{aligned}$$

# Endogeneity Issue: Measurement Error in Regressors

- Consider the following “true” regression model

$$Y_i = \widetilde{\mathbf{X}}_i' \boldsymbol{\beta} + u_i$$

where  $E(\widetilde{\mathbf{X}}_i u_i) = 0$ . But  $\widetilde{\mathbf{X}}_i$  is not measured accurately or cannot be directly measured.

- Instead, we observe  $\mathbf{X}_i = \widetilde{\mathbf{X}}_i + \mathbf{V}_i$ , where  $\mathbf{V}_i$  denotes a  $k \times 1$  vector collecting measurement errors. Therefore,

$$Y_i = \mathbf{X}_i' \boldsymbol{\beta} + \underbrace{(u_i - \mathbf{V}_i' \boldsymbol{\beta})}_{\varepsilon_i}$$

# Endogeneity Issue: Measurement Error in Regressors

- Let  $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)'$ ,  $\mathbf{U} = (u_1, \dots, u_n)'$ ,  $\mathbf{V} = (\mathbf{V}_1, \dots, \mathbf{V}_n)'$ ,  $\mathbf{Q} = \text{E}(\widetilde{\mathbf{X}}_i \widetilde{\mathbf{X}}_i')$ , and  $\mathbf{\Omega} = \text{E}(\mathbf{V}_i \mathbf{V}_i')$ , then

$$\begin{aligned}\hat{\beta} &= \beta + (n^{-1} \mathbf{X}' \mathbf{X})^{-1} [n^{-1} \mathbf{X}' (\mathbf{U} - \mathbf{V} \beta)] \\ &\xrightarrow{p} \beta - (\mathbf{Q} + \mathbf{\Omega})^{-1} \mathbf{\Omega} \beta \neq \beta,\end{aligned}$$

under the assumptions that  $\text{E}(\widetilde{\mathbf{X}}_i \mathbf{V}_i') = 0$ ,  $\text{E}(\widetilde{\mathbf{X}}_i u_i) = 0$ , and  $\text{E}(\mathbf{V}_i u_i) = 0$ .

# Endogeneity Issue: Simultaneous Equation Bias

- Consider following supply and demand equations

$$Q_t^d = \beta_0 + \beta_1 P_t + \varepsilon_t^d, \quad \beta_1 < 0$$

$$Q_t^s = \alpha_0 + \alpha_1 P_t + \varepsilon_t^s, \quad \alpha_1 > 0$$

$$Q_t^d = Q_t^s \quad (= Q_t)$$

and  $\varepsilon_t^d \stackrel{\text{i.i.d.}}{\sim} (0, \sigma_d^2)$ ,  $\varepsilon_t^s \stackrel{\text{i.i.d.}}{\sim} (0, \sigma_s^2)$ ,  $\text{Cov}(\varepsilon_t^d, \varepsilon_t^s) = 0$ .

- Suppose you observe  $\{P_t\}_{t=1}^T$  and  $\{Q_t\}_{t=1}^T$  and run following regression

$$Q_t = \delta_0 + \delta_1 P_t + u_t.$$

# Endogeneity Issue: Simultaneous Equation Bias

- We cannot either pin down  $\beta_1$  or  $\alpha_1$  using OLS estimator and  $\hat{\delta}_1$ .
- First, we solve  $Q_t$  and  $P_t$  by representing the simultaneous equations as follows

$$\begin{bmatrix} 1 & -\beta_1 \\ 1 & -\alpha_1 \end{bmatrix} \begin{bmatrix} Q_t \\ P_t \end{bmatrix} = \begin{bmatrix} \beta_0 \\ \alpha_0 \end{bmatrix} + \begin{bmatrix} \varepsilon_t^d \\ \varepsilon_t^s \end{bmatrix}$$

so that

$$\begin{bmatrix} Q_t \\ P_t \end{bmatrix} = \begin{bmatrix} 1 & -\beta_1 \\ 1 & -\alpha_1 \end{bmatrix}^{-1} \left[ \begin{pmatrix} \beta_0 \\ \alpha_0 \end{pmatrix} + \begin{pmatrix} \varepsilon_t^d \\ \varepsilon_t^s \end{pmatrix} \right].$$



# Endogeneity Issue: Simultaneous Equation Bias

- Next, note that

$$\begin{bmatrix} \hat{\delta}_0 \\ \hat{\delta}_1 \end{bmatrix} = \begin{bmatrix} T & \sum_{t=1}^T P_t \\ \sum_{t=1}^T P_t & \sum_{t=1}^T P_t^2 \end{bmatrix}^{-1} \begin{bmatrix} \sum_{t=1}^T P_t \\ \sum_{t=1}^T P_t Q_t \end{bmatrix}$$

hence,

$$\hat{\delta}_1 = \frac{-\left(\sum_{t=1}^T P_t\right)^2 + T \sum_{t=1}^T P_t Q_t}{T \sum_{t=1}^T P_t^2 - \left(\sum_{t=1}^T P_t\right)^2}$$

- We can show that (Exercise)

$$\hat{\delta}_1 \xrightarrow{p} \frac{(\alpha_0 - \beta_0)(\beta_1 \alpha_0 - \beta_0 \alpha_1)}{(\beta_1 - \alpha_1)^2 (\sigma_s^2 + \sigma_d^2)} + \frac{\sigma_s^2}{\sigma_s^2 + \sigma_d^2} \beta_1 + \frac{\sigma_d^2}{\sigma_s^2 + \sigma_d^2} \alpha_1.$$

# Instrument Variable Estimation

---

- One major problem arising from endogeneity is that coefficients of interested variables are not consistent.
- Suppose that  $X_j$  refers to the endogenous variable in linear regression model

$$Y_i = \mathbf{X}_i' \boldsymbol{\beta} + u_i,$$

an instrument variable  $Z$  is a random variable such that

- (1)  $\text{Cov}(Z, X_j) \neq 0$ .
  - (2)  $\text{Cov}(Z, u) = 0$ .
  - (3)  $Z$  is not highly correlated with other variables  $X_{\setminus j}$ .
- The key idea of instrument variable estimation originates from moment estimation methods.

# Instrument Variable Estimation

- For the simple univariate linear regression model,  $E(u_i) = 0$  and  $\text{Cov}(Z_i, u_i) = 0$ , sample moment condition is

$$\frac{1}{n} \sum (Y_i - \tilde{\beta}_0 - \tilde{\beta}_1 X_i) = 0 \quad \frac{1}{n} \sum Z_i (Y_i - \tilde{\beta}_0 - \tilde{\beta}_1 X_i) = 0.$$

Hence the normal system of equations can be established

$$\begin{aligned} \sum Y_i &= n\tilde{\beta}_0 + \tilde{\beta}_1 \sum X_i \\ \sum Z_i Y_i &= \tilde{\beta}_0 \sum Z_i + \tilde{\beta}_1 \sum Z_i X_i \end{aligned}$$

and finally solves

$$\tilde{\beta}_1 = \frac{\sum z_i y_i}{\sum z_i x_i}, \quad \tilde{\beta}_0 = \bar{Y} - \tilde{\beta}_1 \bar{X}.$$

# Instrument Variable Estimation

---

- In general, suppose that we can partition

$$\mathbf{X}_i = (1, X_{i1}, X_{i2}, \dots, X_{ik})'$$

into two parts such that

$$\mathbf{X}_i = \begin{pmatrix} \mathbf{X}_{i1} \\ \mathbf{X}_{i2} \end{pmatrix} \begin{matrix} k_1 \times 1 \\ k_2 \times 1 \end{matrix}$$

where  $\mathbf{X}_{i1}$  collects exogenous variables such that  $E(\mathbf{X}_{i1}u_i) = 0$  and  $\mathbf{X}_{i2}$  collects endogenous variables such that  $E(\mathbf{X}_{i2}u_i) \neq 0$ .

# Instrument Variable Estimation

- A random vector  $\mathbf{Z}_i$  ( $l \times 1$ ) is an instrument variable if
  - (1)  $E(\mathbf{Z}_i u_i) = 0$
  - (2)  $E(\mathbf{Z}_i \mathbf{X}_i') \neq 0$
- Given the definition of  $\mathbf{Z}_i$ ,  $\mathbf{X}_{i1}$  should be included in  $\mathbf{Z}_i$ . We can denote it as  $\mathbf{Z}_{i1} \equiv \mathbf{X}_{i1}$ , the **included exogenous variables**.
- We can define the remained part of  $\mathbf{Z}_i$ , with  $\mathbf{Z}_{i1}$  excluded, as  $\mathbf{Z}_{i2}$  so that for each  $i$  we have  $\mathbf{Z}_i = (\mathbf{Z}_{i1}', \mathbf{Z}_{i2}')'$ .  $\mathbf{Z}_{i2}$  is referred to as the **excluded exogenous variables**.
- We define  $\mathbf{Z} = (\mathbf{Z}_1, \dots, \mathbf{Z}_n)'$ , an  $n \times l$  matrix.
- In matrix notation:

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}^{(1)} & \mathbf{X}^{(2)} \end{bmatrix}, \quad \mathbf{Z} = \begin{bmatrix} \mathbf{Z}^{(1)} & \mathbf{Z}^{(2)} \end{bmatrix}$$

$n \times k$                        $n \times k_1$     $n \times k_2$                        $n \times l$                        $n \times l_1$     $n \times l_2$

# Instrument Variable Estimation

---

- Recall that

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u} \quad (\dagger)$$

and we can assume that

$$\mathbf{X} = \mathbf{Z}\boldsymbol{\Gamma} + \mathcal{E} \quad (\ddagger)$$

where  $\mathcal{E} = (\boldsymbol{\epsilon}_1, \dots, \boldsymbol{\epsilon}_n)'$  is a  $n \times k$  matrix collecting projection error term.

- $E(\mathbf{Z}_i \boldsymbol{\epsilon}_i') = 0$  given the definition of projection error.

# Instrument Relevance Condition (\*)

- From  $(\dagger)$ , for a specific  $i$  we have

$$\mathbf{X}_i = \mathbf{\Gamma}' \mathbf{Z}_i + \boldsymbol{\varepsilon}_i.$$

- Since  $E(\mathbf{Z}_i \boldsymbol{\varepsilon}_i') = 0$ , we can derive that

$$\mathbf{\Gamma} = E(\mathbf{Z}_i \mathbf{Z}_i')^{-1} E(\mathbf{Z}_i \mathbf{X}_i').$$

and hence  $E(\mathbf{Z}_i \mathbf{X}_i') \neq 0 \Rightarrow \mathbf{\Gamma} \neq 0$ .

- If we further assume that  $E(\boldsymbol{\varepsilon}_i) = 0$ , we can further derive that

$$\begin{aligned} \mathbf{\Gamma} &= \left\{ E[\mathbf{Z}_i - E(\mathbf{Z}_i)] [\mathbf{Z}_i - E(\mathbf{Z}_i)]' \right\}^{-1} \\ &\quad \times E[\mathbf{Z}_i - E(\mathbf{Z}_i)] [\mathbf{X}_i - E(\mathbf{X}_i)]'. \end{aligned}$$

and hence  $\text{Cov}(\mathbf{Z}_i, \mathbf{X}_i') \neq 0 \Rightarrow \mathbf{\Gamma} \neq 0$ .

- In this sense,  $\text{Cov}(\mathbf{Z}_i, \mathbf{X}_i') \neq 0$  is a more restrictive instrument relevance condition.

# Instrument Variable Estimation

- By substituting  $(\ddagger)$  into  $(\dagger)$ , we can establish the relationship between  $\mathbf{Y}$  and  $\mathbf{Z}$

$$\begin{aligned}\mathbf{Y} &= (\mathbf{Z}\mathbf{\Gamma} + \mathbf{\mathcal{E}})\boldsymbol{\beta} + \mathbf{u} \\ &= \mathbf{Z}\mathbf{\Gamma}\boldsymbol{\beta} + \mathbf{\mathcal{E}}\boldsymbol{\beta} + \mathbf{u} \\ &= \mathbf{Z}\boldsymbol{\lambda} + \mathbf{V}\end{aligned}\tag{*}$$

- OLS for  $(\ddagger)$  yields  $\hat{\mathbf{\Gamma}} = (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{X}$ .
- OLS for  $(*)$  yields  $\hat{\boldsymbol{\lambda}} = (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{Y}$ .
- It would be natural to ask whether we can solve  $\boldsymbol{\beta}$  from  $\mathbf{\Gamma}\boldsymbol{\beta} = \boldsymbol{\lambda}$  as an estimator for  $\boldsymbol{\beta}$ .



# Instrument Variable Estimation: Identification

---

- If one can solve or recover  $\Gamma\beta = \lambda$ , we say that  $\beta$  is **identified**.
- Assume that  $\text{rank}(\Gamma) = k$ , then
  - if  $l = k$ ,  $\beta = \Gamma^{-1}\lambda$ , which refers to the **just-identified** case.
  - if  $l > k$ , then for any positive definite matrix  $W$ ,

$$\beta = (\Gamma'W\Gamma)^{-1} \Gamma'W\lambda,$$

which refers to the **over-identified** case.

- A necessary condition (not sufficient) for  $\text{rank}(\Gamma) = k$  is  $l \geq k$ .

# Instrument Variable Estimation: Identification

---

- For the just-identified case, we have

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

- $\hat{\beta}_{IV}$  is consistent for  $\beta$  under the regular conditions.
- Finding a good IV(s) is more like **Art** rather than **Science**.  
For instance, Qian (2008, QJE) “*Missing Women and the Price of Tea in China: The Effect of Sex-Specific Earnings on Sex Imbalance*”

# Two Stage Least Squares (2SLS)

- **Stage 1:** Regress  $\mathbf{X}$  on  $\mathbf{Z}$  to obtain  $\hat{\Gamma} = (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{X}$  and save the predicted value  $\hat{\mathbf{X}} = \mathbf{Z}\hat{\Gamma} = P_Z\mathbf{X}$ .
- **Stage 2:** Regress  $\mathbf{Y}$  on  $\hat{\mathbf{X}}$  to obtain the 2SLS estimator  $\hat{\beta}_{2SLS} = (\hat{\mathbf{X}}'\hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}'\mathbf{Y}$ .
- In fact,

$$\begin{aligned}\hat{\beta}_{2SLS} &= (\mathbf{X}'P_Z\mathbf{X})^{-1} \mathbf{X}'P_Z\mathbf{Y} \\ &= \left[ \frac{\mathbf{X}'\mathbf{Z}}{n} \left( \frac{\mathbf{Z}'\mathbf{Z}}{n} \right)^{-1} \frac{\mathbf{Z}'\mathbf{X}}{n} \right]^{-1} \frac{\mathbf{X}'\mathbf{Z}}{n} \left( \frac{\mathbf{Z}'\mathbf{Z}}{n} \right)^{-1} \frac{\mathbf{Z}'\mathbf{Y}}{n}.\end{aligned}$$

- Under the regular conditions,  $\hat{\beta}_{2SLS}$  is consistent for  $\beta$ .

For the just-identified case, i.e.  $l = k$

- $\hat{\beta}_{2SLS} = \hat{\beta}_{IV}$ .

$$\begin{aligned}\hat{\beta}_{2SLS} &= \left[ \mathbf{X}'\mathbf{Z} (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{X} \right]^{-1} \mathbf{X}'\mathbf{Z} (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{Y} \\ &= (\mathbf{Z}'\mathbf{X})^{-1} (\mathbf{Z}'\mathbf{Z}) (\mathbf{X}'\mathbf{Z})^{-1} \mathbf{X}'\mathbf{Z} (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{Y} \\ &= (\mathbf{Z}'\mathbf{X})^{-1} \mathbf{Z}'\mathbf{Y}\end{aligned}$$

- We can solve IV estimator as  $\hat{\beta}_{IV} = \hat{\Gamma}^{-1} \hat{\lambda}$ .

$$\begin{aligned}\hat{\beta}_{ILS} &= \hat{\Gamma}^{-1} \hat{\lambda} \\ &= \left[ (\mathbf{Z}'\mathbf{Z})^{-1} (\mathbf{Z}'\mathbf{X}) \right]^{-1} (\mathbf{Z}'\mathbf{Z})^{-1} (\mathbf{Z}'\mathbf{Y}) \\ &= (\mathbf{Z}'\mathbf{X})^{-1} (\mathbf{Z}'\mathbf{Z}) (\mathbf{Z}'\mathbf{Z})^{-1} (\mathbf{Z}'\mathbf{Y}) \\ &= (\mathbf{Z}'\mathbf{X})^{-1} \mathbf{Z}'\mathbf{Y} = \hat{\beta}_{IV}\end{aligned}$$

Background:

- How economic conditions affect sex imbalance?
- **Identification problems:** areas with higher female income may have higher income precisely because women's status is higher for other reasons.
- How to find appropriate proxy variables?

- Two reforms after 1979: (1) increase in procurement prices of cash crops; (2) Household Production Responsibility System (HPRS).
- Women have a comparative advantage in producing tea, whereas men have a comparative advantage in producing orchard fruits. Therefore, areas suitable for tea cultivation experienced an increase in female-generated income, whereas areas suitable for orchard cultivation experienced an increase in male-generated income.

Key regression function:

$$\begin{aligned} sex_{ic} = & (tea_i \times post_c) \beta + (orchard_i \times post_c) \delta \\ & + (cashcrop_i \times post_c) \rho + Han_{ic} \zeta + \alpha + \psi_i + \gamma_c + \varepsilon_{ic} \end{aligned}$$

- $sex_{ic}$ : fraction of males in country  $i$  of cohort  $c$  (individuals born after 1979 or not,  $post_c$  as the dummy variable).
- $tea_i$ : the amount of tea planted.
- $orchard_i$ : the amount of orchard planted.
- $cashcrop_i$ : the amount of cash crops planted.
- $Han_i$ : the fraction of ethnically Han.

Two reasons for IV:

- 1997 agricultural data to proxy for agricultural conditions in earlier years introduces measurement error.
- Many other reasons, not necessarily only the increased value of female labor, that may lead to increase of amount of tea planted after 1979 reform. In other words, we potentially have omitted variable bias.

Inspiring IV in Qian (2008, QJE):

- Using the **average slope** of each county as the IV for tea planting.



# Qian (2008, QJE)

Darker shades correspond to more tea planted per household.



Darker shades correspond to steeper regions.

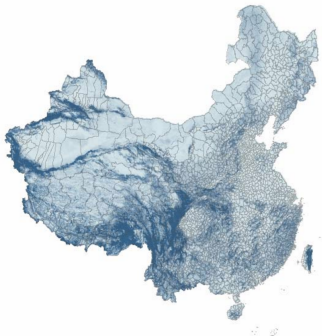


TABLE III  
OLS AND 2SLS ESTIMATES OF THE EFFECT OF PLANTING TEA AND ORCHARDS ON SEX  
RATIOS CONTROLLING FOR COUNTY LEVEL LINEAR COHORT TRENDS

	Dependent variables					
	Fraction of males			Tea $\times$ post	Fraction of males	
	(1) OLS	(2) OLS	(3) OLS	(4) 1st	(5) IV	(6) IV
Tea $\times$ post	-0.012 (0.007)	-0.013 (0.006)	-0.012 (0.005)		-0.072 (0.031)	-0.011 (0.007)
Orchard $\times$ post	0.005 (0.002)					
Slope $\times$ post	-0.002 (0.002)			0.26 (0.057)		
Linear trend	No	No	Yes	Yes	No	Yes
Observations	28,349	37,756	37,756	37,756	37,756	37,756

# Test of Endogeneity

Durbin-Wu-Hausman tests:

- Consider testing:

$$H_0 : E(\mathbf{X}_i u_i) = 0$$

$$H_1 : E(\mathbf{X}_i u_i) \neq 0, \quad E(\mathbf{Z}_i u_i) = 0$$

- Recall that

$$\underset{n \times k}{\mathbf{X}} = \begin{bmatrix} \underset{n \times k_1}{\mathbf{X}^{(1)}} & \underset{n \times k_2}{\mathbf{X}^{(2)}} \end{bmatrix}, \quad \underset{n \times l}{\mathbf{Z}} = \begin{bmatrix} \underset{n \times l_1}{\mathbf{Z}^{(1)}} & \underset{n \times l_2}{\mathbf{Z}^{(2)}} \end{bmatrix}$$

and run following **augmented** regression

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + P_{\mathbf{Z}}\mathbf{X}^{(2)}\boldsymbol{\delta} + \boldsymbol{\varepsilon}.$$

# Test of Endogeneity

---

- **Step 1.** Run the OLS regression of  $\mathbf{X}^{(2)}$  on the instrumental variables  $\mathbf{Z}$  and obtain the fitted value  $\widehat{\mathbf{X}}^{(2)} = P_{\mathbf{Z}}\mathbf{X}^{(2)}$ .
- **Step 2.** Run the OLS regression of  $\mathbf{Y}$  on  $\mathbf{X}$  and  $\widehat{\mathbf{X}}^{(2)}$  and test whether the coefficient vector  $\boldsymbol{\delta}$  of  $\widehat{\mathbf{X}}^{(2)}$  is 0.
- Under  $H_0$ , the  $F$ -test statistic  $F_n$  based on the **augmented** regression is asymptotically distributed as the Wald statistic:  $W_n \equiv k_2 F_n$  is asymptotically distributed as  $\chi^2(k_2)$ .