

Instrumental Variables Design

Prof. Tzu-Ting Yang
楊子霆

Institute of Economics, Academia Sinica
中央研究院經濟研究所

May 13, 2026

Main Idea

Main Idea of Instrumental Variables

- The Instrumental Variable (IV) is an exogenous sources of variation that drives the treatment D_i but unrelated to other confounding factors that affect outcome Y_i
- Intuitively, IV breaks variation of the treatment D_i into two parts
 - 1 A part that might be correlated with other confounding factors
 - ★ This part causes selection bias
 - 2 A part that is not
 - ★ This part could represent the clean causal effect
- Use the variation in D_i that is not correlated with other confounding factors to estimate causal effect of the treatment

Unobservable Omitted Variable

- Suppose the true model is:

$$Y_i = \delta + \alpha D_i + \beta_1 X_i + \beta_2 U_i + \epsilon_i$$

- X_i is the **observed characteristics**
 - ▶ We can directly control for it
- But U_i is the **unobserved characteristics**
 - ▶ e.g. ability, preference, health
- So we cannot include it into our regression and estimate the following model:

$$Y_i = \delta + \alpha D_i + \beta_1 X_i + \zeta_i$$

- where $\zeta_i = \beta_2 U_i + \epsilon_i$

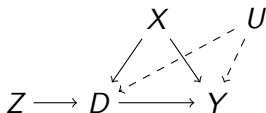
Unobservable Omitted Variable

- As mentioned before, failure to include key covariates will lead to omitted variable bias

$$\begin{aligned}\hat{\alpha} &\xrightarrow{p} \alpha + \frac{\text{Cov}(\zeta_i, D_i)}{V(D_i)} \\ &= \alpha + \beta_2 \frac{\text{Cov}(U_i, D_i)}{V(D_i)}\end{aligned}$$

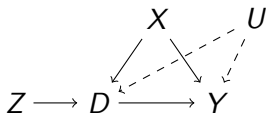
- Remember there is NO omitted variable bias (OVB) if:
 - U_i is unrelated to Y_i : $\beta_2 = 0$
 - U_i is unrelated to D_i : $\text{Cov}(U_i, D_i) = 0$
- To obtain eliminate OVB, we need a variation in D_i that is unrelated to the unobserved confounding factor U_i

Main Idea of Instrumental Variables



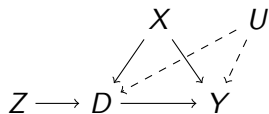
- Y is an outcome (e.g. earnings),
- Z is the instrument
- D is the treatment (e.g. college degree)
- X is the **observed** confounding factor (e.g. family income)
- U is the **unobserved** confounding factor (e.g. ability)

Main Idea of Instrumental Variables



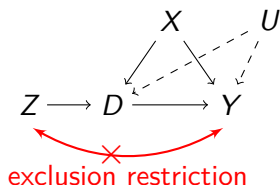
- **Unobserved ability** U might confound with the effect of college degree D
 - ▶ Since ability affects people to get college degree D and their earnings Y
- We need to find an IV that generate a variation in getting college degree D that is unrelated to ability U

Main Idea of Instrumental Variables



- IV initiates a causal chain: the instrument Z affects D , which in turn affects Y
- A valid IV needs to satisfy the following conditions:
 - 1 First-stage relationship (Instrument relevance): Z affects D

Main Idea of Instrumental Variables

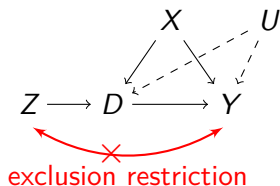


- A valid IV needs to satisfy the following conditions:

2 Exclusion restriction (Instrument exogeneity):

- ★ **No direct or indirect effect** of the instrument Z on the outcome Y NOT through the treatment variable D
- ★ The instrument Z affects the outcome Y **only through the treatment variable** D

Main Idea of Instrumental Variables



- We can test whether the **instrument relevance** is satisfied
- But the **instrument exogeneity** cannot be tested
 - ▶ You have to convince your audience that it is satisfied

Identification

Example

Effect of Military Service on Lifetime Income

Joshua D. Angrist (1990) “**Lifetime Earnings and the Vietnam Era Draft Lottery: Evidence from Social Security Administrative Records**” AER

- He wanted to examine the effect of military service on lifetime income.
- We will use Angrist's paper on the effects of military service (D_i) on earnings (Y_i) as an example to go through key concept of IV design

Example

Effect of Military Service on Lifetime Income

- Joining military service is a personal choice
- Is there any selection bias due to **unobservable confounding factors** in this example ?
 - ▶ **Time preference:**
 - ★ Less patient people may voluntarily join military service early
 - ★ This myopic thinking may have negative impact on their earnings (e.g. less human capital investment)
 - ▶ **Health condition:**
 - ★ Better health people can join military service
 - ★ Better health condition also have positive impact on their earnings
- We need a IV for the treatment variable of joining military service

Example

Effect of Military Service on Lifetime Income

- Angrist (1990) uses the **Vietnam draft lottery** (Z_i) as in IV for military service
 - ▶ In the 1960s and early 1970s, young American men were draft for military service to serve in Vietnam
 - ▶ Concerns about the fairness of the conscription policy lead to the introduction of a **draft lottery** in 1970

Example

Effect of Military Service on Lifetime Income

- From 1970 to 1972 **random sequence numbers** were assigned to each birth date in cohorts of 19-year-olds
 - ▶ Men with lottery numbers below a cutoff were eligible for military service
 - ▶ While men with numbers above the cutoff were ineligible
 - ★ $Z = \mathbf{I}[L < c]$
 - ★ L is lottery number and c is cutoff
- The eligibility did NOT perfectly determinate military service:
 - ▶ Many draft-eligible men were exempted for health and other reasons
 - ▶ Draft-ineligible men volunteered for service
- Next, we briefly discuss whether draft eligibility induced by lottery is a good IV or not

Example

Effect of Military Service on Lifetime Income

- The lottery used by the Selective Service to determine who would be drafted for Vietnam first



Source: Historic Photographs

Example

Effect of Military Service on Lifetime Income

- First-stage relationship (Instrument relevance): Z_i affects D_i
 - ▶ Vietnam veteran status (joining military service) was not completely determined by randomized draft eligibility
 - ▶ But draft eligibility is highly correlated with Vietnam veteran status
- Exclusion restriction (Instrument exogeneity):
 - ▶ The draft eligibility is determined by random numbers
 - ▶ These numbers should not affect one's earnings directly

Potential Outcomes Framework

Treatment Assignment

- Treatment Assignment

$$Z_i = \begin{cases} 1 & \text{if an individual } i \text{ is eligible for a treatment} \\ 0 & \text{if an individual } i \text{ is not eligible for a treatment} \end{cases}$$

- ▶ $Z_i = 1$: those who get draft eligibility
- ▶ $Z_i = 0$: those who do not get draft eligibility
 - ★ Due to lottery results

Potential Outcomes Framework

Potential Treatments

- Potential Treatments

- ▶ D_i^z : Potential treatment status given the value of Z

- ★ D_i^1 : Potential treatment status if eligible for a treatment

- ★ D_i^0 : Potential treatment status if not eligible for a treatment

- Observed Treatment

$$D_i = \begin{cases} D_i^1 & \text{if } Z_i = 1 \\ D_i^0 & \text{if } Z_i = 0 \end{cases}$$

- or, in a more compact notation: $D_i = Z_i D_i^1 + (1 - Z_i) D_i^0$

Potential Outcomes Framework

Potential Outcomes

- Potential Outcomes

- ▶ Y_i^1 : earnings if individual i serves in the military
- ▶ Y_i^0 : earnings if individual i does not serve in the military

- Observed Outcome

$$Y_i = D_i Y_i^1 + (1 - D_i) Y_i^0$$

- This notation already imposes the exclusion idea: once treatment status is fixed, Z_i has no separate effect on Y_i

Identification Roadmap

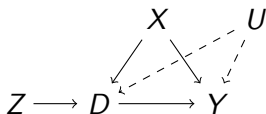
- We observe how outcomes and treatment rates differ by draft eligibility:

$$\alpha_{IV} = \frac{E[Y_i|Z_i = 1] - E[Y_i|Z_i = 0]}{E[D_i|Z_i = 1] - E[D_i|Z_i = 0]}$$

- To interpret this ratio causally, we need to answer three questions:
 - ▶ Why does comparing $Z_i = 1$ and $Z_i = 0$ recover causal variation?
 - ▶ Through which channel can Z_i affect earnings?
 - ▶ Whose treatment status is actually changed by Z_i ?

Identification Results for IV

Observed Wald Ratio



- IV uses the causal chain from draft eligibility Z_i to military service D_i to earnings Y_i
- Intuitively:

Effect of Z_i on Y_i

$$= (\text{Effect of } Z_i \text{ on } D_i) \times (\text{Effect of } D_i \text{ on } Y_i)$$

Identification Results for IV

Observed Wald Ratio

- Rearranging gives the Wald estimand:

$$\alpha_{IV} = \frac{E[Y_i|Z_i = 1] - E[Y_i|Z_i = 0]}{E[D_i|Z_i = 1] - E[D_i|Z_i = 0]}$$

- This ratio has two observable pieces:
 - ▶ Reduced-form effect: the numerator captures how draft eligibility changes average earnings
 - ▶ First-stage effect: the denominator captures how draft eligibility changes the probability of military service
- The remaining question: what causal effect does this ratio identify?

Identification Assumptions for IV

First-Stage Relationship

- **First-Stage Relationship:** Z_i affects treatment status D_i

$$E[D_i|Z_i = 1] - E[D_i|Z_i = 0] \neq 0$$

- ▶ Draft eligibility affects the probability of joining military

Identification Assumptions for IV

Independence

- **Independence Assumption:** Z_i is independent of potential outcomes and potential treatments

$$(Y_i^1, Y_i^0, D_i^1, D_i^0) \perp\!\!\!\perp Z_i$$

- ▶ Draft eligibility is as good as randomly assigned
- ▶ Therefore, potential earnings and potential treatment choices are balanced across $Z_i = 1$ and $Z_i = 0$
- ▶ This lets us replace conditional expectations given Z_i with unconditional expectations in the proof

Identification Assumptions for IV

Exclusion Restriction

- **Exclusion Restriction:** Z_i affects outcome Y_i only through changing treatment status D_i
 - ▶ The instrument has no direct effect on the outcome once treatment status is fixed
 - ▶ In the draft lottery example, draft eligibility should affect earnings only by changing military service

$$Y_i = D_i Y_i^1 + (1 - D_i) Y_i^0$$

Identification Assumptions for IV

Monotonicity Assumption

- **Monotonicity Assumption:** $D_i^1 \geq D_i^0$
 - ▶ Monotonicity says that the presence of the instrument never dissuades someone from taking the treatment
 - ★ This is sometimes called **no defiers**
 - ▶ In the draft lottery example, eligibility should weakly encourage military service

IV and Compliance Types

- Because draft eligibility does not perfectly determine military service, people can react differently to Z_i
- Each individual has two potential treatment statuses:
 - ▶ D_i^1 : treatment status if draft eligible
 - ▶ D_i^0 : treatment status if not draft eligible
- The pair (D_i^1, D_i^0) tells us whose behavior the instrument can change

IV and Compliers

- We can define four types of individuals based on whether they follow the draft eligibility results:
 - ▶ **Compliers:** $D_i^1 > D_i^0$ ($D_i^0 = 0$ and $D_i^1 = 1$)
 - ★ David got draft eligibility and joined military service
 - ★ Tim did not get draft eligibility and did not join military service
 - ▶ **Always-takers:** $D_i^1 = D_i^0 = 1$
 - ★ John always joined military service no matter the lottery results (whether he got draft eligibility)
 - ▶ **Never-takers:** $D_i^1 = D_i^0 = 0$
 - ★ Trump never joined military service no matter the lottery results (whether he got draft eligibility)
 - ▶ **Defiers:** $D_i^1 < D_i^0$ ($D_i^0 = 1$ and $D_i^1 = 0$)
 - ★ Jimmy got draft eligibility but did NOT join military service
 - ★ Jonson did NOT get draft eligibility but joined military service

Identification Results for IV

What the Wald Ratio Identifies

IV Identify LATE

$$\alpha_{IV} = \frac{E[Y_i|Z_i = 1] - E[Y_i|Z_i = 0]}{E[D_i|Z_i = 1] - E[D_i|Z_i = 0]} = E[Y_i^1 - Y_i^0 | D_i^1 > D_i^0]$$

- IV identifies the average causal effect for individuals whose treatment status is changed by the instrument
- In the draft lottery example, this is the effect of military service for men induced to serve because they became draft eligible

Identification Results for IV

Proof: Convert Observed Differences

Proof:

$$\begin{aligned}\alpha_{IV} &= \frac{E[Y_i|Z_i = 1] - E[Y_i|Z_i = 0]}{E[D_i|Z_i = 1] - E[D_i|Z_i = 0]} \\ &= \frac{E[Y_i^1 D_i^1 + Y_i^0(1 - D_i^1)|Z_i = 1] - E[Y_i^1 D_i^0 + Y_i^0(1 - D_i^0)|Z_i = 0]}{E[D_i^1|Z_i = 1] - E[D_i^0|Z_i = 0]} \\ &= \frac{E[Y_i^1 D_i^1 + Y_i^0(1 - D_i^1)] - E[Y_i^1 D_i^0 + Y_i^0(1 - D_i^0)]}{E[D_i^1] - E[D_i^0]} \\ &= \frac{E[(Y_i^1 - Y_i^0)(D_i^1 - D_i^0)]}{E[D_i^1] - E[D_i^0]}\end{aligned}$$

Identification Results for IV

Proof: Use Compliance Types

- Since D_i^z is a dummy, $D_i^1 - D_i^0$ can only be 1, 0, or -1
 - ▶ **Always-takers** and **never-takers**: $D_i^1 - D_i^0 = 0$
 - ▶ **Compliers**: $D_i^1 - D_i^0 = 1$
 - ▶ **Defiers**: $D_i^1 - D_i^0 = -1$
- Under monotonicity, defiers are ruled out:

$$\begin{aligned} & E[(Y_i^1 - Y_i^0)(D_i^1 - D_i^0)] \\ &= E[(Y_i^1 - Y_i^0) \times (1) | D_i^1 - D_i^0 = 1] \Pr(D_i^1 - D_i^0 = 1) \\ &= E[Y_i^1 - Y_i^0 | D_i^1 > D_i^0] \Pr(D_i^1 - D_i^0 = 1) \end{aligned}$$

Identification Results for IV

Proof: Interpret the First Stage

- The denominator is the first-stage effect:

$$\begin{aligned}E[D_i^1] - E[D_i^0] &= \Pr(D_i^1 = 1) - \Pr(D_i^0 = 1) \\&= \Pr(D_i^1 = 1, D_i^0 = 0) + \Pr(D_i^1 = 1, D_i^0 = 1) \\&\quad - \Pr(D_i^0 = 1, D_i^1 = 0) - \Pr(D_i^0 = 1, D_i^1 = 1) \\&= \Pr(D_i^1 = 1, D_i^0 = 0) - \Pr(D_i^0 = 1, D_i^1 = 0)\end{aligned}$$

- Under monotonicity, $\Pr(D_i^0 = 1, D_i^1 = 0) = 0$
- Therefore, the first stage is the share of compliers:

$$E[D_i^1] - E[D_i^0] = \Pr(D_i^1 = 1, D_i^0 = 0) = \Pr(D_i^1 - D_i^0 = 1)$$

Identification Results for IV

Proof: LATE

- Combine the numerator and denominator results

Proof:

$$\begin{aligned}\alpha_{IV} &= \frac{E[Y_i|Z_i = 1] - E[Y_i|Z_i = 0]}{E[D_i|Z_i = 1] - E[D_i|Z_i = 0]} \\ &= \frac{E[Y_i^1 D_i^1 + Y_i^0(1 - D_i^1)|Z_i = 1] - E[Y_i^1 D_i^0 + Y_i^0(1 - D_i^0)|Z_i = 0]}{E[D_i^1|Z_i = 1] - E[D_i^0|Z_i = 0]} \\ &= \frac{E[Y_i^1 D_i^1 + Y_i^0(1 - D_i^1)] - E[Y_i^1 D_i^0 + Y_i^0(1 - D_i^0)]}{E[D_i^1] - E[D_i^0]} \\ &= \frac{E[(Y_i^1 - Y_i^0)(D_i^1 - D_i^0)]}{E[D_i^1] - E[D_i^0]} \\ &= \frac{E[(Y_i^1 - Y_i^0)(1)|D_i^1 - D_i^0 = 1] \Pr(D_i^1 - D_i^0 = 1)}{\Pr(D_i^1 - D_i^0 = 1)} \\ &= E[Y_i^1 - Y_i^0 | D_i^1 > D_i^0] = \alpha_{\text{LATE}}\end{aligned}$$

Identification Results for IV

Why Local?

- **Never-takers** and **always-takers** do not change treatment status when the instrument changes
 - ▶ They do not contribute to the first stage
 - ▶ Their treatment effects are not revealed by this instrument
- Monotonicity rules out **defiers**
- Therefore, IV identifies the average treatment effect for **compliers**

Identification Results for IV

LATE

LATE

$\alpha_{\text{LATE}} = E[Y_i^1 - Y_i^0 | D_i^1 > D_i^0]$, the average treatment effect for compliers, is called the **Local Average Treatment Effect (LATE)**.

- It is local because it applies to the subpopulation whose treatment status is changed by the instrument
 - ▶ In this example: men induced into military service by draft eligibility
- Different instruments may move different groups of compliers
 - ▶ Whether LATE is interesting depends on whether those compliers are substantively important

Identification Results for IV

LATE and ATE

- Without further assumptions (e.g. constant causal effects), LATE is not informative about effects on never-takers or always-takers
 - ▶ Because the instrument does not affect their treatment status.
- In most applications we would be mostly interested in estimating the average treatment effect on the whole population (ATE).

$$E[Y_i^1 - Y_i^0]$$

- This is usually not possible with IV.

Estimation

Review: IV Estimation

A intuitive way

- Causal relationship of interest: the effect of military service on earnings

$$Y_i = \delta + \alpha D_i + u_i$$

- Remember we just derive:

$$\begin{aligned}\alpha_{IV} &= \text{Effect of treatment on outcome} \\ &= \frac{\text{Effect of instrument on outcome}}{\text{Effect of instrument on treatment}}\end{aligned}$$

Review: IV Estimation

A intuitive way

- We can estimate α_{IV} by running the following two regressions:
- Reduced form regression: the effect of lottery draft on earnings

$$Y_i = \mu + \alpha_{RF}Z_i + X'\delta + \varepsilon_i$$

$$\alpha_{RF} = \frac{\text{Cov}(Y_i, Z_i)}{V(Z_i)}$$

- First-Stage regression: the effect of lottery draft on military service

$$D_i = \kappa + \alpha_{FS}Z_i + X'\beta + \zeta_i$$

$$\alpha_{FS} = \frac{\text{Cov}(D_i, Z_i)}{V(Z_i)}$$

- The IV estimator is:

$$\hat{\alpha}_{IV} = \frac{\hat{\alpha}_{RF}}{\hat{\alpha}_{FS}} = \frac{\hat{\text{Cov}}(Y_i, Z_i)}{\hat{\text{Cov}}(D_i, Z_i)}$$

Review: IV Estimation

Two Stage Least Squares (TSLS)

- In practice, IV is often estimated using Two Stage Least Squares (TSLS).
- If identification assumptions hold only after conditioning on X , covariates are introduced in TSLS regression.
- TSLS involves two steps:
 1. First-stage regression:

$$D_i = \kappa + \alpha_{FS}Z_i + X_i'\beta + \nu_i$$

Estimate to obtain fitted values \hat{D}_i .

2. Second-stage regression:

$$Y_i = \delta + \alpha_{TSLS}\hat{D}_i + X_i'\gamma + u_i^*$$

Review: IV Estimation

Two Stage Least Squares (TSLS)

- However, following the two-step procedure naively would yield incorrect standard errors.
- Why? Because the error term u_i^* reflects the use of estimated regressors:

- ▶ Standard errors based on u_i^* are incorrect:

$$u_i^* = Y_i - \delta - X_i' \gamma - \alpha_{TSLS} \hat{D}_i$$

- ▶ What we need is based on the true error:

$$u_i = Y_i - \delta - X_i' \gamma - \alpha_{TSLS} D_i$$

- Standard software packages (e.g., Stata, R) automatically compute correct standard errors, adjusting for the first-stage estimation error.

Review: Inference in TSLS

- Under standard assumptions, the TSLS estimator is consistent and asymptotically normal in large samples.
- Inference (hypothesis testing, confidence intervals) proceeds as in OLS, using the asymptotic normality of the estimator.

Summary of Hypothesis Testing for TSLS Regression

- We estimate the following regressions:

$$Y_i = \delta + \alpha_{TSLS}D_i + X_i'\gamma + u_i$$

$$D_i = \kappa + \alpha_{FS}Z_i + X_i'\beta + \zeta_i$$

1. Check IV relevance (first stage):

- ▶ Test whether α_{FS} is statistically significantly different from zero.
- ▶ More formally, check whether the first-stage F-statistic exceeds 10 to avoid weak instrument problems.

2. Choose a null hypothesis to test:

- ▶ $H_0 : \alpha_{TSLS} = 0$ or $H_0 : \alpha_{TSLS} = \mu$
- ▶ A claim we would like to reject.

Summary of Hypothesis Testing for Regression

3. Choose a test statistic:

$$\blacktriangleright t = \frac{(\hat{\alpha}_{TSLs} - \alpha_{TSLs})}{\hat{SE}(\hat{\alpha}_{TSLs})}$$

4. Compute the standard error of $\hat{\alpha}_{TSLs}$:

$$\hat{SE}(\hat{\alpha}_{TSLs}) = \sqrt{\frac{\hat{\sigma}_u^2}{\sum_{i=1}^n (\hat{D}_i - \bar{\hat{D}})^2}}$$

where:

- $\hat{\sigma}_u^2$ is the estimated variance of the second-stage residuals.
- \hat{D}_i are the fitted values from the first-stage regression.
- $\bar{\hat{D}}$ is the mean of \hat{D}_i .
- n is the sample size.

Factors Reducing the Standard Error of $\hat{\alpha}_{TSLS}$

- **Lower Residual Variance ($\hat{\sigma}^2$):** Better model fit leads to less variability around the regression line.
- **Greater Variability in Fitted Values (\hat{D}_i):** More dispersion in \hat{D}_i improves the precision of $\hat{\alpha}_{TSLS}$.
- **Larger Sample Size (n):** Reduces sampling variability and improves estimate precision.

Summary of Hypothesis Testing for Regression

5. Determine the distribution of the test statistic under the null hypothesis
 - ▶ If sample size is sufficient large, using CLT, t-statistic will have standard normal distribution
6. Calculate the probability of wrongly reject null hypothesis given null hypothesis is true (p-value)
 - ▶ We reject the null hypothesis $H_0 : \alpha_{TSLs} = 0$ against the alternative $H_1 : \alpha_{TSLs} \neq 0$ at the 5% significance level if $|t| > 1.96$

Summary of Findings on Vietnam Draft Lottery

1. First Stage Results:

- ▶ Having a low lottery number (being eligible for the draft) increases the probability of veteran status by about 16 percentage points.
- ▶ Note that the mean veteran status is about 27%.
- ▶ This confirms the relevance condition for the instrumental variable.

2. Second Stage Results:

- ▶ Serving in the army lowers annual earnings by between \$2,050 and \$2,741.
- ▶ This estimate captures the LATE for compliers.

3. Placebo Test:

- ▶ There is no evidence of an association between draft eligibility (low lottery number) and earnings in 1969.
- ▶ 1969 earnings were realized before the 1970 draft lottery, supporting the validity of the exclusion restriction.

Summary of Findings on Vietnam Draft Lottery

Earnings year	Earnings		Veteran Status		Wald Estimate of Veteran Effect
	Mean	Eligibility Effect	Mean	Eligibility Effect	
	(1)	(2)	(3)	(4)	(5)
1981	16,461	-435.8 (210.5)	0.267	0.159 (0.040)	-2,741 (1,324)
1971	3,338	-325.9 (46.6)			-2050 (293)
1969	2,299	-2.0 (34.5)			

Notes: Adapted from Angrist (1990), Tables 2 and 3. Standard errors are shown in parentheses. Earnings data are from Social Security administrative records. Figures are in nominal dollars. Veteran status data are from the Survey of Program Participation. There are about 13,500 individuals in the sample.

STATA Example

Acemoglu, Johnson, and Robinson (2001) “**The Colonial Origins of Comparative Development: An Empirical Investigation**” AER

- They want to examine the effect of institutions on economic development
 - ▶ Do countries with better institutions achieve a greater level of income?
- **Good institutions** include:
 - ▶ Strong protection of property rights
 - ▶ Less distortionary policies (e.g., high tax rate)
 - ▶ Checks and balances on government power
 - ▶ Investment-friendly policies (e.g., rule of law, contract enforcement)

Acemoglu, Johnson, and Robinson (2001)

Motivation



Overview

Program and Data

- See IV.do
- Use AJR_table4.dta
- If you want to know the complete programs and data for this paper, please use the following:
 - ▶ AJR-IV.do
 - ▶ Use AJR_table2.dta to AJR_table7.dta

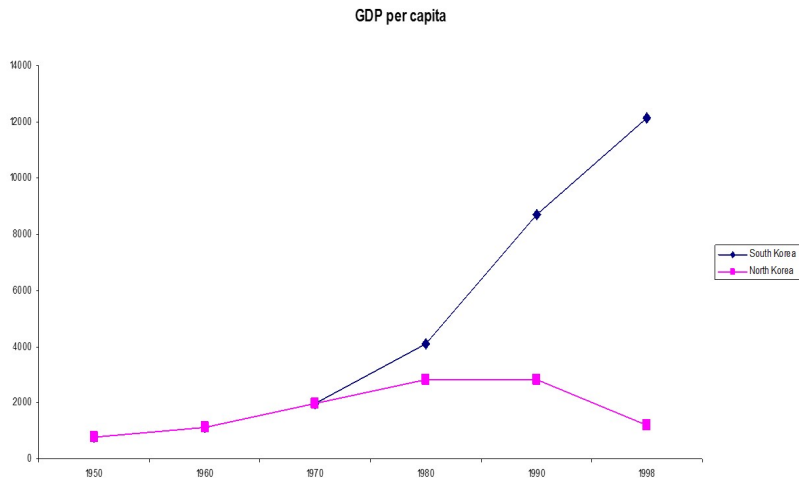
Acemoglu, Johnson, and Robinson (2001)

Motivation

- At some level it is obvious that institutions matter
- Witness, for example, the divergent paths of North and South Korea, or East and West German
 - ▶ central planning and collective ownership V.S. private property and market economy

Acemoglu, Johnson, and Robinson (2001)

Motivation



Acemoglu, Johnson, and Robinson (2001)

Motivation

- Nevertheless, we lack reliable estimates of the effect of institutions on economic performance
 - ▶ **Selection bias 1:** It is quite likely that rich economies choose or can afford better institutions
 - ▶ **Selection bias 2:** Economies that are different for a variety of reasons will differ both in their institutions and in their income per capita
- Need to eliminate selection bias

Acemoglu, Johnson, and Robinson (2001)

Identification Strategy

- They propose an IV to generate an exogenous variation in institution based on theory plus history
- They look only among former European colonies

Acemoglu, Johnson, and Robinson (2001)

Identification Strategy

- Their theory is based on the following facts:
 - 1 In some colonies, Europeans had good survival, in others not
 - 2 Where Europeans could survive, they put down roots, established good institutions
 - ★ Replicate European institutions, with strong emphasis on private property and checks against government power

- 3 Where they were dying like flies, they set up extractive institutions
- 4 Extractive institutions designed to get resources quickly
 - ▶ Extractive institutions did not introduce much protection for private property, nor did they provide checks and balances against government expropriation
- These institutions persisted after decolonialization

Acemoglu, Johnson, and Robinson (2001)

Identification Strategy

- They use **mortality rates expected by the first European settlers** in the colonies as an IV for current institutions in these countries
- Malaria and yellow fever were the major sources of European mortality in the colonies
- Their theory is:

Health environment \Rightarrow Settler mortality \Rightarrow European settlement \Rightarrow
Early institutions \Rightarrow Current institutions \Rightarrow Output today

Acemoglu, Johnson, and Robinson (2001)

Identification Strategy

- **Key assumption – exclusion restriction:** settler mortality can NOT affect output today by any other channel
- Possible threats to identification:
 - ▶ Health environment might affect output today directly
 - ▶ Having European settlers affects output through some channel other than institutions (language, human capital)

Acemoglu, Johnson, and Robinson (2001)

Identification Strategy

- Yellow fever and malaria had much less effect on native inhabitants, who had acquired and genetic immunity
 - ▶ The prevalence of these diseases depend on the microclimate of an area (e.g. temperature and humidity)
- This suggests that mortality rates faced by Europeans are unlikely to be a proxy for some simple geographic or climatic feature of the country

Acemoglu, Johnson, and Robinson (2001)

TSLS estimation

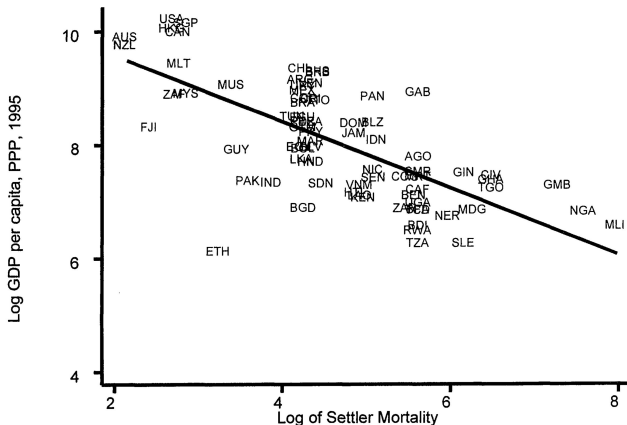
- They conduct a TSLS estimation
 - ▶ **First stage:** current economic institutions = $g(\text{settler mortality})$
 - ▶ **Second stage:** log income per capita = $f(\text{current economic institutions})$

- Mortality rates of soldiers stationed in the colonies in the early 19th century
 - ▶ They got it from historian Philip Curtin
- **Current economic institutions** proxied by **protection against expropriation risk**
 - ▶ Average **protection against expropriation risk** is measured on a scale from 0 to 10
 - ▶ A higher score means more protection against expropriation, averaged over 1985 to 1995, from Political Risk Services

Acemoglu, Johnson, and Robinson (2001)

Reduced-form relationship

- Income and Settler Mortality



Acemoglu, Johnson, and Robinson (2001)

First-stage relationship

- Protection for Expropriation Risk and Settler Mortality

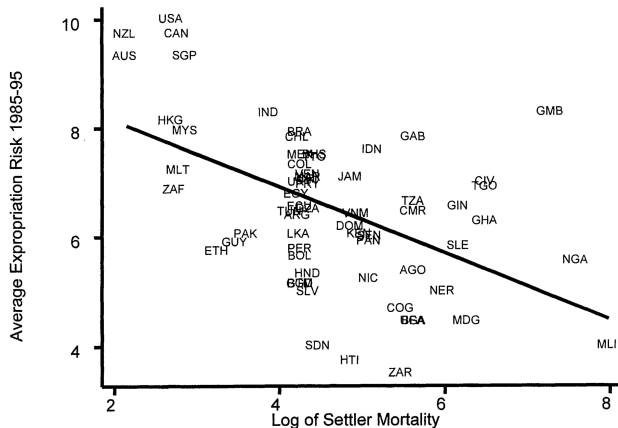


FIGURE 3. FIRST-STAGE RELATIONSHIP BETWEEN SETTLER MORTALITY AND EXPROPRIATION RISK

Acemoglu, Johnson, and Robinson (2001)

Second-stage relationship

- Protection for Expropriation Risk and Income

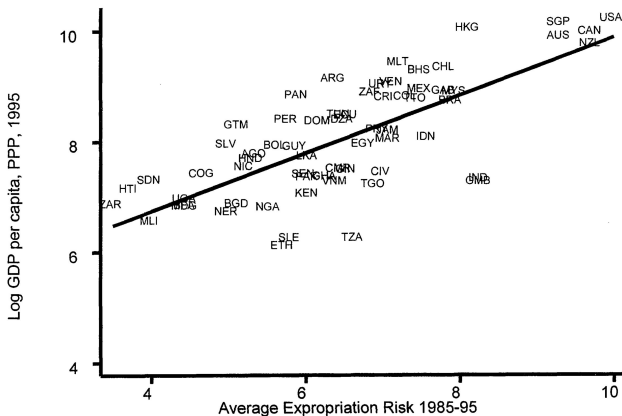


FIGURE 2. OLS RELATIONSHIP BETWEEN EXPROPRIATION RISK AND INCOME

Acemoglu, Johnson, and Robinson (2001)

Descriptive Statistics

TABLE 1—DESCRIPTIVE STATISTICS

	Whole world	Base sample	By quartiles of mortality			
			(1)	(2)	(3)	(4)
Log GDP per capita (PPP) in 1995	8.3 (1.1)	8.05 (1.1)	8.9	8.4	7.73	7.2
Log output per worker in 1988 (with level of United States normalized to 1)	-1.70 (1.1)	-1.93 (1.0)	-1.03	-1.46	-2.20	-3.03
Average protection against expropriation risk, 1985–1995	7 (1.8)	6.5 (1.5)	7.9	6.5	6	5.9
Constraint on executive in 1990	3.6 (2.3)	4 (2.3)	5.3	5.1	3.3	2.3
Constraint on executive in 1900	1.9 (1.8)	2.3 (2.1)	3.7	3.4	1.1	1
Constraint on executive in first year of independence	3.6 (2.4)	3.3 (2.4)	4.8	2.4	3.1	3.4
Democracy in 1900	1.1 (2.6)	1.6 (3.0)	3.9	2.8	0.19	0
European settlements in 1900	0.31 (0.4)	0.16 (0.3)	0.32	0.26	0.08	0.005
Log European settler mortality	n.a.	4.7 (1.1)	3.0	4.3	4.9	6.3
Number of observations	163	64	14	18	17	15

Notes: Standard deviations are in parentheses. Mortality is potential settler mortality, measured in terms of deaths per annum per 1,000 “mean strength” (raw mortality numbers are adjusted to what they would be if a force of 1,000 living people were kept in place for a whole year, e.g., it is possible for this number to exceed 1,000 in episodes of extreme mortality as those who die are replaced with new arrivals). Sources and methods for mortality are described in Section III, subsection B, and in the unpublished Appendix (available from the authors; or see Acemoglu et al., 2000). Quartiles of mortality are for our base sample of 64 observations. These are: (1) less than 65.4; (2) greater than or equal to 65.4 and less than 78.1; (3) greater than or equal to 78.1 and less than 280; (4) greater than or equal to 280. The number of observations differs by variable; see Appendix Table A1 for details.

$$\log(Y_i) = \mu + \alpha R_i + X_i' \gamma + \varepsilon_i$$

- Y_i is income per capita in country i
- R_i is the protection against expropriation measure
- X_i' is a vector of other covariates
- α represents the **effect of institutions on income per capita**

OLS Results

STATA Implementation

```
. regress logpgp95 avexpr lat_abst africa asia other if baseco==1, robust
```

```
Linear regression                               Number of obs   =           64
                                                F(5, 58)       =           53.79
                                                Prob > F       =           0.0000
                                                R-squared     =           0.7139
                                                Root MSE     =           .58163
```

logpgp95	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]	
avexpr	.4012826	.0640653	6.26	0.000	.2730419	.5295233
lat_abst	.8752965	.6142898	1.42	0.160	-.3543381	2.104931
africa	-.8806805	.1555275	-5.66	0.000	-1.192003	-.5693585
asia	-.5767549	.2991853	-1.93	0.059	-1.175639	.0221296
other	.107205	.2229782	0.48	0.632	-.3391344	.5535445
_cons	5.736729	.3893229	14.74	0.000	4.957415	6.516044

OLS Results

TABLE 2—OLS REGRESSIONS

	Whole world (1)	Base sample (2)	Whole world (3)	Whole world (4)	Base sample (5)	Base sample (6)	Whole world (7)	Base sample (8)
	Dependent variable is log GDP per capita in 1995						Dependent variable is log output per worker in 1988	
Average protection against expropriation risk, 1985–1995	0.54 (0.04)	0.52 (0.06)	0.47 (0.06)	0.43 (0.05)	0.47 (0.06)	0.41 (0.06)	0.45 (0.04)	0.46 (0.06)
Latitude			0.89 (0.49)	0.37 (0.51)	1.60 (0.70)	0.92 (0.63)		
Asia dummy				-0.62 (0.19)		-0.60 (0.23)		
Africa dummy				-1.00 (0.15)		-0.90 (0.17)		
“Other” continent dummy				-0.25 (0.20)		-0.04 (0.32)		
R^2	0.62	0.54	0.63	0.73	0.56	0.69	0.55	0.49
Number of observations	110	64	110	110	64	64	108	61

Notes: Dependent variable: columns (1)–(6), log GDP per capita (PPP basis) in 1995, current prices (from the World Bank’s World Development Indicators 1999); columns (7)–(8), log output per worker in 1988 from Hall and Jones (1999). Average protection against expropriation risk is measured on a scale from 0 to 10, where a higher score means more protection against expropriation, averaged over 1985 to 1995, from Political Risk Services. Standard errors are in parentheses. In regressions with continent dummies, the dummy for America is omitted. See Appendix Table A1 for more detailed variable definitions and sources. Of the countries in our base sample, Hall and Jones do not report output per worker in the Bahamas, Ethiopia, and Vietnam.

Acemoglu, Johnson, and Robinson (2001)

IV Results

- First stage

$$R_i = \mu + \alpha \log(M_i) + X_i' \vartheta + \eta_i$$

- Second stage

$$\log(Y_i) = \mu + \alpha R_i + X_i' \gamma + \varepsilon_i$$

- M_i is mortality rates faced by settler at country i

STATA Command: ivregress

- Syntax:

```
1 ivregress estimator depvar [varlist1] (varlist2 =  
   varlistiv) [if] [in] [weight] [, options]
```

- Example:

```
1 ivregress 2sls logpgp95 lat_abst africa asia other_cont  
   (avexpr=logem4), first  
2 estat firststage
```

- **varlist1** is the list of exogenous variables.
- **varlist2** is the list of endogenous variables.
- **varlistiv** is the list of exogenous variables used with varlist1 as instruments for varlist2.
- **2sls**: two-stage least squares

STATA Command: ivregress

- options:
 - ▶ **estat firststage**: report first-stage F-statistic
 - ▶ **level(#)**: set confidence level; default is level(95)
 - ▶ **first**: requests that the first-stage regression results be displayed

IV Results

STATA Implementation

```
ivregress 2sls logppg95 (avexpr=logem4) f_brit f_french, first
```

First-stage regressions

	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
avexpr						
f_brit	.629348	.3664792	1.72	0.091	-.1037196	1.362416
f_french	.0474048	.4295458	0.11	0.912	-.8118147	.9066243
logem4	-.5343989	.139576	-3.83	0.000	-.8135924	-.2552053
_cons	8.746647	.6904157	12.67	0.000	7.36561	10.12768

```
Number of obs = 64
F( 3, 60) = 8.91
Prob > F = 0.0001
R-squared = 0.3081
Adj R-squared = 0.2736
Root MSE = 1.2518
```

Instrumental variables (2SLS) regression

```
Number of obs = 64
Wald chi2(3) = 32.21
Prob > chi2 = 0.0000
R-squared = 0.0483
Root MSE = 1.0099
```

	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
logppg95						
avexpr	1.07785	.210709	5.12	0.000	.6648682	1.490832
f_brit	-.7777037	.343026	-2.27	0.023	-1.450022	-.1053851
f_french	-.1169738	.3435071	-0.34	0.733	-.7902354	.5562878
_cons	1.372403	1.34394	1.02	0.307	-1.261672	4.006477

```
Instrumented: avexpr
Instruments: f_brit f_french logem4
```

```
. estat firststage
```

First-stage regression summary statistics

Variable	R-sq.	Adjusted R-sq.	Partial R-sq.	F(1,60)	Prob > F	

STATA Command: ivreghdfe

- Syntax:

```
1 ivreghdfe depvar [controls] (endogenous_var =  
    instruments) [if] [in] [weight], absorb(fixed  
    effects) [options]
```

- Example:

```
1 ivreghdfe logpgp95 africa asia other_cont (avexpr=  
    logem4), first
```

- Basic structure:

- ▶ **controls**: Exogenous control variables.
- ▶ **endogenous_var**: Endogenous variable(s) to be instrumented.
- ▶ **instruments**: Instrumental variable(s).

STATA Command: ivreghdfe

- Key options:

- ▶ `absorb(fixed effects)`: Specify fixed effects to control for (e.g., country, year).
- ▶ `cluster()`: Cluster standard errors at a specified level.
- ▶ `first`: Display first-stage regression results and diagnostics (optional).

- Advantages of `ivreghdfe`:

- ▶ Handles multiple high-dimensional fixed effects efficiently.
- ▶ Automatically clusters standard errors for valid inference.
- ▶ Useful for large datasets and widely used in applied microeconomic research.

Robustness Checks

- Control for latitude
- Maybe settler mortality just means bad disease environment
 - ▶ Can't control **life expectancy** because clearly this is endogenous (affected by treatment): **bad control**
 - ▶ Include measures of temperature and humidity meant to capture disease environment
 - ▶ Also put in measures of soil quality
- IV result survives all of these robustness checks

Test of Exclusion Restriction

- Theory:
 - ▶ Settler mortality (M) → Settlements (S) → Early institutions (C) → Current institutions (R) → Income ($\log y$).
- Exclusion restriction assumption:
 - ▶ Mortality (M) affects income ($\log y$) only through institutions (R).

Test of Exclusion Restriction

- Testing strategy:
 - ▶ Test whether M, S, or C have a direct effect on $\log y$ after controlling for R.
 - ▶ Add log settler mortality directly as a regressor in the second stage.
- Result:
 - ▶ Mortality is not significantly related to income after controlling for institutions.
 - ▶ Supports the validity of the exclusion restriction.

Test Exclusion Restriction

TABLE 8—OVERIDENTIFICATION TESTS

	Base sample (1)	Base sample (2)	Base sample (3)	Base sample (4)	Base sample (5)	Base sample (6)	Base sample (7)	Base sample (8)	Base sample (9)	Base sample (10)
Panel A: Two-Stage Least Squares										
Average protection against expropriation risk, 1985–1995	0.87 (0.14)	0.92 (0.20)	0.71 (0.15)	0.68 (0.20)	0.72 (0.14)	0.69 (0.19)	0.60 (0.14)	0.61 (0.17)	0.55 (0.12)	0.56 (0.14)
Latitude		-0.47 (1.20)		-0.34 (1.10)		0.31 (1.05)		-0.41 (0.92)		-0.16 (0.81)
Panel B: First Stage for Average Protection Against Expropriation Risk										
European settlements in 1900	3.20 (0.62)	2.90 (0.83)								
Constraint on executive in 1900			0.32 (0.08)	0.26 (0.09)						
Democracy in 1900					0.24 (0.06)	0.20 (0.07)				
Constraint on executive in first year of independence							0.25 (0.08)	0.22 (0.08)		
Democracy in first year of independence									0.19 (0.05)	0.17 (0.05)
R ²	0.30	0.30	0.20	0.24	0.24	0.26	0.19	0.25	0.26	0.30
Panel C: Results from Overidentification Test										
p-value (from chi-squared test)	[0.67]	[0.96]	[0.09]	[0.20]	[0.11]	[0.28]	[0.67]	[0.79]	[0.22]	[0.26]
Panel D: Second Stage with Log Mortality as Exogenous Variable										
Average protection against expropriation risk, 1985–1995	0.81 (0.23)	0.88 (0.30)	0.45 (0.25)	0.42 (0.30)	0.52 (0.23)	0.48 (0.28)	0.49 (0.23)	0.49 (0.25)	0.4 (0.18)	0.41 (0.19)
Log European settler mortality	-0.07 (0.17)	-0.05 (0.18)	-0.25 (0.16)	-0.26 (0.17)	-0.21 (0.15)	-0.22 (0.16)	-0.14 (0.16)	-0.14 (0.15)	-0.19 (0.13)	-0.19 (0.12)
Latitude		-0.52 (1.15)		0.38 (0.89)		0.28 (0.86)		-0.38 (0.84)		-0.17 (0.73)

R Example

Overview

Program and Data

- See IV.R
- Use AJR_table4.dta

Install Required Packages

- Install and load required packages:

```
1 library(AER)           # for IV regression
2 library(haven)        # for reading Stata files
3 library(texreg)       # for regression tables
4 library(dplyr)        # for data manipulation
5 library(tibble)       # for rownames handling
6 library(xml2)         # for HTML processing
7 library(rvest)        # for HTML table extraction
8 library(writexl)      # for Excel output
```

- **AER**: Implements instrumental variables regression
- **texreg**: Creates publication-quality regression tables
- **haven**, **dplyr**, **tibble**: Data import and manipulation
- **xml2**, **rvest**, **writexl**: Convert and export results

R Commands: ivreg

```
1 # Run IV regression
2 ivreg(logpgp95 ~ avexpr + lat_abst | logem4 + lat_abst, data
      = data)
```

- **formula:** outcome endogenous + exogenous | instruments + exogenous
- **data:** specify the data frame
- Equivalent to STATA's **ivregress 2sls**

Practical Issues

Practical Tips for IV Design

Checking for Weak Instruments

1. Check IV relevance:

- ▶ Is the instrument theoretically plausible?
- ▶ Are the first-stage coefficients of the expected sign and reasonable magnitude?
- ▶ Report the first-stage F-statistic for instrument strength.

● Rule of thumb:

- ▶ If $F\text{-statistic} > 10$: Instruments are considered strong; TSLS is reliable.
- ▶ If $F\text{-statistic} < 10$: Instruments are weak; TSLS may be biased.

● Consequences of weak instruments:

- ▶ TSLS estimator becomes biased toward OLS.
- ▶ Standard t-statistics and confidence intervals are invalid (non-normal distribution).

Testing for Weak Instruments

- **First-stage F-statistic** (homoskedastic case):
 - ▶ Standard F-test from first-stage regression
 - ▶ Compare to Stock-Yogo critical values
- **Kleibergen-Paap rk statistic** (heteroskedastic case):
 - ▶ Robust version of Cragg-Donald F-statistic
 - ▶ Still compare to Stock-Yogo critical values (though conservative)
- **Interpretation:** If test statistic $<$ critical value, instruments are likely weak, requiring robust inference methods.

Statistical Inference Robust to Weak Instruments

- When instruments are weak (as identified by the previous tests), standard TSLS inference is invalid:
- Alternative inference methods that remain valid with weak instruments:
 - ▶ **Anderson-Rubin test:**
 - ★ Valid regardless of instrument strength
 - ▶ **Weak IV robust confidence intervals:**
 - ★ Construct confidence sets using methods that maintain correct coverage regardless of instrument strength (e.g., AR-based CI)
 - ★ May be wider than conventional CI, reflecting the uncertainty from weak instruments

- Practical implementation using Stata:
 - ▶ `ivreghdfe` supports weak IV diagnostics.
 - ▶ Use the `first` option to obtain Cragg-Donald F-statistic and Anderson-Rubin tests

Statistical Inference Robust to Weak Instruments

Summary results for first-stage regressions

Variable				(Underid)		(Weak id)	
	F(1, 62)	P-val	SW	Chi-sq(1)	P-val	SW	F(1, 62)
avexpr	22.95	0.0000		23.69	0.0000		22.95

Stock-Yogo weak ID F test critical values for single endogenous regressor:

10% maximal IV size	16.38
15% maximal IV size	8.96
20% maximal IV size	6.66
25% maximal IV size	5.53

Source: Stock-Yogo (2005). Reproduced by permission.

NB: Critical values are for Sanderson-Windmeijer F statistic.

Underidentification test

H₀: matrix of reduced form coefficients has rank=K1-1 (underidentified)

H_a: matrix has rank=K1 (identified)

Anderson canon. corr. LM statistic Chi-sq(1)=17.29 P-val=0.0000

Weak identification test

H₀: equation is weakly identified

Cragg-Donald Wald F statistic 22.95

Stock-Yogo weak ID test critical values for K1=1 and L1=1:

10% maximal IV size	16.38
15% maximal IV size	8.96
20% maximal IV size	6.66
25% maximal IV size	5.53

Source: Stock-Yogo (2005). Reproduced by permission.

Weak-instrument-robust inference

Tests of joint significance of endogenous regressors B1 in main equation

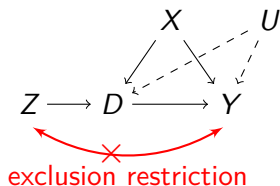
H₀: B1=0 and orthogonality conditions are valid

Anderson-Rubin Wald test F(1,62)= 56.60 P-val=0.0000

Anderson-Rubin Wald test Chi-sq(1)= 58.43 P-val=0.0000

Stock-Wright LM S statistic Chi-sq(1)= 30.54 P-val=0.0000

Practical Tips For IV Papers



2. Check exclusion restriction

- The exclusion restriction cannot be tested directly, but it can be falsified
- **Placebo test**
 - ▶ Test the reduced form effect of Z_i on Y_i in situations where it is impossible or extremely unlikely that Z_i could affect D_i
 - ▶ Because Z_i can't affect D_i , then the exclusion restriction implies that this placebo test should have zero effect.

Practical Tips For IV Papers

3. If you have many IVs pick your best instrument and report the just identified model
4. Look at the reduced form
 - ▶ Directly estimate the effect of instrument Z on outcome Y
 - ▶ If you can't see the causal relationship of interest in the reduced form it is probably not there

Suggested Readings

- Chapter 3, Mastering Metrics: The Path from Cause to Effect
- Chapter 4, Mostly Harmless Econometrics
- Chapter 7, Causal Inference: The Mixtape