# Instrumental Variables

## EC 425/525, Set 8

Edward Rubin
08 May 2019

# Prologue

# Schedule

## Last time

Mathching and propensity-score methods

- Conditional independence
- Overlap

## Today

Instrumental variables (and two-stage least squares)

## Upcoming

- Admin: Assignment/project proposal this weekend
- Admin: Midterm very soon

# Research designs

# Research designs

## Selection on observables and/or unobservables

We've been focusing on **selection-on-observables designs**, *i.e.*,

$$(\mathbf{Y}_{0i},\ \mathbf{Y}_{1i}) \perp\!\!\!\perp \mathbf{D}_i | \mathbf{X}_i$$

for **observable** variables $\mathbf{X}_i$.

**Selection-on-unobservable designs** replace this assumption with two new (but related) assumptions

1. $(\mathbf{Y}_{0i},\ \mathbf{Y}_{1i}) \perp \mathbf{Z}_i$

2. $\mathrm{Cov}(\mathbf{Z}_i,\ \mathbf{D}_i) \neq 0$

# Research designs

## Selection on observables and/or unobservables

Our main goal in causal-inference minded (applied) econometrics boils down to isolating **"good" variation** in $D_i$ (exogenous/as-good-as-random) from **"bad" variation** (the part of $D_i$ correlated with $Y_{0i}$ and $Y_{1i}$).

(We want to avoid selection bias.)

- **Selection-on-observables designs** assume that we can control for all *bad variation* (selection) in $D_i$ through a known (observed) $X_i$.

- **Selection-on-unobservables designs** assume that we can extract part of the *good variation* in $D_i$ (generally using some $Z_i$) and then use this *good* part of $D_i$ to estimate the effect of $D_i$ on $Y_i$. We throw away the *bad variation* in $D_i$ (it's bad).

# Research designs

## Which route?

So set of research designs is more palatable?

1. There are plenty of bad applications of both sets.
   Violated assumptions, bad controls, *etc.*

1. **Selection on observables** assumes we know *everything* about selection into treatment—we can identify *all* of the good (or bad) variation in $\mathbf{D}_i$.
   Tough in non-experimental settings. Difficult to validate in practice.

1. **Selection on unobservables** assumes we can isolate *some* good/clean variation in $\mathbf{D}_i$, which we then use to estimate the effect of $\mathbf{D}_i$ on $\mathbf{Y}_i$.
   Seems more plausible. Possible to validate. May be underpowered.

# Instrumental variables

## Introduction

**Instrumental variables** (IV)[†] is the canonical selection-on-unobservables design—isolating *good variation* in $\mathbf{D}_i$ via some magical instrument $\mathbf{Z}_i$.

Consider some model (structural equation)

$$\mathbf{Y}_i = \beta_0 + \beta_1 \mathbf{D}_i + \varepsilon_i \tag{1}$$

To guarantee consistent OLS estimates for $\beta_1$, want $\mathbf{Cov}(\mathbf{D}_i, \varepsilon_i) = 0$.
In general, this is a heroic assumption.

*Alternative:* Estimate $\beta_1$ via instrumental variables.

[†] For the moment, we're lumping together IV and two-stage least squares (2SLS) together—as many people do—even though they are technically different.

# Instrumental variables

## Definition

For our model

$$\mathbf{Y}_i = \beta_0 + \beta_1 \mathbf{D}_i + \varepsilon_i \tag{1}$$

A valid **instrument** is a variable $\mathbf{Z}_i$ such that

1. $\mathrm{Cov}(\mathbf{Z}_i, \mathbf{D}_i) \neq 0$
   our instrument correlates with treatment (so we can keep part of $\mathbf{D}_i$)

1. $\mathrm{Cov}(\mathbf{Z}_i, \varepsilon_i) = 0$
   our instrument is uncorrelated with other (non-$\mathbf{D}_i$) determinants of $\mathbf{Y}_i$,
   *i.e.*, $\mathbf{Z}_i$ is excludable from equation (1). (exclusion restriction)

# Instrumental variables

## Example

Back to the returns to a college degree,

$$\text{Income}_i = \beta_0 + \beta_1 \text{Grad}_i + \varepsilon_i$$

OLS is likely biased.

What if that state conducts a (random) **lottery** for scholarships?

Let $\text{Lottery}_i$ denote an indicator for whether $i$ won a lottery scholarship.[†]

1. $\text{Cov}(\text{Lottery}_i, \text{Grad}_i) \neq 0 \ (> 0)$ if scholarships increase grad. rates.

1. $\text{Cov}(\text{Lottery}_i, \varepsilon_i) = 0$ since the lottery is randomized.

[†] We'll have to focus on families who were eligible/who applied.

# Instrument variables

## The IV estimator

The IV estimator for our model

$$\mathbf{Y}_i = \beta_0 + \beta_1 \mathbf{D}_i + \varepsilon_i \tag{1}$$

with (valid) instrument $\mathbf{Z}_i$ is

$$\hat{\beta}_{\mathrm{IV}} = \left(\mathbf{Z}'\mathbf{D}\right)^{-1} \left(\mathbf{Z}'\mathbf{Y}\right)$$

If you have no covariates, then

$$\hat{\beta}_{\mathrm{IV}} = \frac{\mathrm{Cov}(\mathbf{Z}_i,\ \mathbf{Y}_i)}{\mathrm{Cov}(\mathbf{Z}_i,\ \mathbf{D}_i)}$$

# Instrument variables

## The IV estimator

The IV estimator for our model

$$\mathbf{Y}_i = \beta_0 + \beta_1 \mathbf{D}_i + \varepsilon_i \tag{1}$$

with (valid) instrument $\mathbf{Z}_i$ is

$$\hat{\beta}_{\mathrm{IV}} = \left(\mathbf{Z}'\mathbf{D}\right)^{-1}\left(\mathbf{Z}'\mathbf{Y}\right)$$

If you have additional (exogenous) covariates $\mathbf{X}_i$, then

$$\mathbf{Z} = \begin{bmatrix} \mathbf{Z}_i & \mathbf{X}_i \end{bmatrix}$$

$$\mathbf{D} = \begin{bmatrix} \mathbf{D}_i & \mathbf{X}_i \end{bmatrix}$$

# Instrumental variables

## Proof: Consistency

With a valid instrument $\mathbf{Z}_i$, $\hat{\beta}_{\mathrm{IV}}$ is a consistent estiamtor for $\beta_1$ in

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \tag{1}$$

$$\mathrm{plim}\left(\hat{\beta}_{IV}\right)$$

$$= \mathrm{plim}\left(\left(\mathbf{Z}'\mathbf{D}\right)^{-1}\left(\mathbf{Z}'\mathbf{Y}\right)\right)$$

$$= \mathrm{plim}\left(\left(\mathbf{Z}'\mathbf{D}\right)^{-1}\left(\mathbf{Z}'\mathbf{D}\beta + \mathbf{Z}'\varepsilon\right)\right)$$

$$= \mathrm{plim}\left(\left(\mathbf{Z}'\mathbf{D}\right)^{-1}\left(\mathbf{Z}'\mathbf{D}\right)\beta\right) + \mathrm{plim}\left(\frac{1}{N}\mathbf{Z}'\mathbf{D}\right)^{-1}\mathrm{plim}\left(\frac{1}{N}\mathbf{Z}'\varepsilon\right)$$

$$= \beta \quad \textcolor{magenta}{\checkmark}$$

# Two-stage least squares

# Two-stage least squares

## Setup

You'll commonly see IV implemented as a two-stage process known as **two-stage least squares** (2SLS).

**First stage** Estimate the effect of the instrument $\mathbf{Z}_i$ on our endogenous variable $\mathbf{D}_i$ and (predetermined) covariates $\mathbf{X}_i$. Save $\widehat{\mathbf{D}}_i$.

$$\mathbf{D}_i = \gamma_1 \mathbf{Z}_i + \gamma_2 \mathbf{X}_i + u_i$$

**Second stage** Estimate model we wanted—but only using the variation in $\mathbf{D}_i$ that correlates with $\mathbf{Z}_i$, *i.e.*, $\widehat{\mathbf{D}}_i$.

$$\mathbf{Y}_i = \beta_1 \widehat{\mathbf{D}}_i + \beta_2 \mathbf{X}_i + \varepsilon_i$$

*Note* The controls $\mathbf{X}_i$ must match in the first and second stages.

# Two-stage least squares

## IV estimation

This two-step procedure, with a valid instrument, produces an estimator $\hat{\beta}_1$ that is consistent for $\beta_1$.

$$\hat{\beta}_{\text{2SLS}} = \left(\mathbf{D}'\mathbf{P_Z}\mathbf{D}\right)^{-1}\left(\mathbf{D}'\mathbf{P_Z}\mathbf{Y}\right)$$

$$\mathbf{P_Z} = \mathbf{Z}\left(\mathbf{Z}'\mathbf{Z}\right)^{-1}\mathbf{Z}'$$

where $\mathbf{D}$ is a matrix of our treatment and predetermined covariates $(\mathbf{X}_i)$ and $Z$ is a matrix of our instrument and our predetermined covariates.

# Two-stage least squares

## IV estimation

Important notes

- The controls ($\mathbf{X}_i$) must match in the first and second stages.

- If you have exactly **one instrument** and exactly **one endogenous variable**, then 2SLS and IV are identical.

- Your second-stage standard errors are not correct.

# Table of contents

## Admin

## Instrumental variables

## Two-stage least squares