

Empirical Exercise 2: Medicare Advantage Data and Instrumental Variables

Ian McCarthy | Emory University
Econ 771

Table of contents

1. Overview
2. Instrumental Variables
3. Data
4. Analysis
5. Extras

Overview

Goals of this assignment

1. Work with data on insurance markets (Medicare Advantage)
2. Employ 2SLS with real data
3. Implement some sensitivity analyses relevant for 2SLS in practice

Plus the sub-goal for all assignments...practice your Git/GitHub workflow, version control, and replicability

Specific "research question"

Does more private insurance lead to lower costs for Medicare?

- Economics: Spillover effects from MA policy onto how Medicare patients are treated
- Important for quantifying full effects of changes to MA policy

Instrumental Variables

What is instrumental variables

Instrumental Variables (IV) is a way to identify causal effects using variation in treatment participation that is due to an *exogenous* variable that is only related to the outcome through treatment.

Why bother with IV?

Two reasons to consider IV:

1. Selection on unobservables
2. Reverse causation

Either problem is sometimes loosely referred to as *endogeneity*

Simple example

- $y = \beta x + \varepsilon(x)$,
where $\varepsilon(x)$ reflects the dependence between our observed variable and the error term.
- Simple OLS will yield
$$\frac{dy}{dx} = \beta + \frac{d\varepsilon}{dx} \neq \beta$$

What does IV do?

- The regression we want to do:

$$y_i = \alpha + \delta W_i + \gamma A_i + \epsilon_i,$$

where W_i is treatment (think of schooling for now) and A_i is something like ability.

- A_i is unobserved, so instead we run:

$$y_i = \alpha + \beta W_i + \epsilon_i$$

- From this "short" regression, we don't actually estimate δ . Instead, we get an estimate of

$$\beta = \delta + \lambda_{ws}\gamma \neq \delta,$$

where λ_{ws} is the coefficient of a regression of A_i on W_i .

Intuition

IV will recover the "long" regression without observing underlying ability

IF our IV satisfies all of the necessary assumptions.

More formally

- We want to estimate

$$E[Y_i|W_i = 1] - E[Y_i|W_i = 0]$$

- With instrument Z_i that satisfies relevant assumptions, we can estimate this as

$$E[Y_i|W_i = 1] - E[Y_i|W_i = 0] = \frac{E[Y_i|Z_i=1] - E[Y_i|Z_i=0]}{E[W_i|Z_i=1] - E[W_i|Z_i=0]}$$

- In words, this is effect of the instrument on the outcome ("reduced form") divided by the effect of the instrument on treatment ("first stage")

IVs in practice

Easy to think of in terms of randomized controlled trial...

Measure	Offered Seat	Not Offered Seat	Difference
Score	-0.003	-0.358	0.355
% Enrolled	0.787	0.046	0.741
Effect			0.48

Angrist *et al.*, 2012. "Who Benefits from KIPP?" *Journal of Policy Analysis and Management*.

What is IV *really* doing

Think of IV as two-steps:

1. Isolate variation due to the instrument only (not due to endogenous stuff)
2. Estimate effect on outcome using only this source of variation

In regression terms

Interested in estimating δ from $y_i = \alpha + \beta x_i + \delta W_i + \varepsilon_i$, but W_i is endogenous (no pure "selection on observables").

Step 1: With instrument Z_i , we can regress W_i on Z_i and x_i ,
$$W_i = \lambda + \theta Z_i + \kappa x_i + \nu,$$
and form prediction \hat{W}_i .

Step 2: Regress y_i on x_i and \hat{W}_i ,
$$y_i = \alpha + \beta x_i + \delta \hat{W}_i + \xi_i$$

In regression terms

But in practice, *DON'T* do this in two steps. Why?

Because standard errors are wrong...not accounting for noise in prediction, \hat{W}_i .
The appropriate fix is built into most modern stats programs.

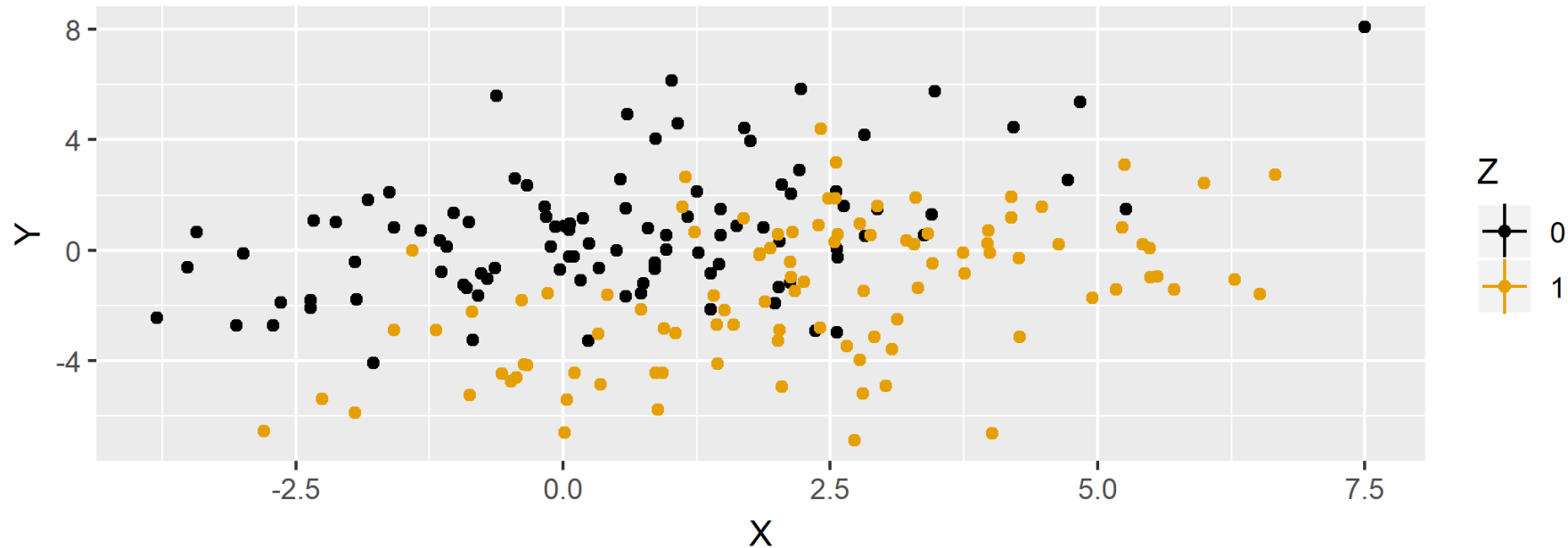
Key IV assumptions

1. *First-stage*: Instrument is correlated with the endogenous variable
2. *Exogeneity*: Instrument is uncorrelated with the error term
3. *Exclusion*: Instruments do not directly affect your outcome
4. *Monotonicity*: Treatment more (less) likely for those with higher (lower) values of the instrument

Assumptions 1 and 3 sometimes grouped into an *only through* condition.

Animation for IV

The Relationship between Y and X, With Binary Z as an Instrumental Variable
1. Start with raw data. Correlation between X and Y: 0.251



Simulated data

```
n ← 5000
b.true ← 5.25
iv.dat ← tibble(
  z = rnorm(n,0,2),
  eps = rnorm(n,0,1),
  w = (z + 1.5*eps>0.15),
  y = 2.5 + b.true*w + eps + rnorm(n,0,0.5)
)
```

- endogenous `eps`: affects treatment and outcome
- `z` is an instrument: affects treatment but no direct effect on outcome

Results with simulated data

Recall that the *true* treatment effect is 5.25

```
##  
## Call:  
## lm(formula = y ~ w, data = iv.dat)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -3.5482 -0.6674 -0.0037  0.6840  3.3803   
##  
## Coefficients:  
##              Estimate Std. Error t value
```

```
##  
## Call:  
## ivreg(formula = y ~ w | z, data = iv.dat)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -4.009846 -0.744281 -0.005625  0.757291   
##  
## Coefficients:  
##              Estimate Std. Error t value
```

Checking instrument

- Check the 'first stage'

```
##  
## Call:  
## lm(formula = w ~ z, data = iv.dat)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -1.02567 -0.30855 -0.00491  0.31129  1.6
```

- Check the 'reduced form'

```
##  
## Call:  
## lm(formula = y ~ z, data = iv.dat)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -8.4742 -2.1306 -0.0198  2.1180  8.7368
```

Two-stage equivalence

```
step1 ← lm(w ~ z, data=iv.dat)
w.hat ← predict(step1)
step2 ← lm(y ~ w.hat, data=iv.dat)
summary(step2)
```

```
##
```

```
## Call:
```

```
## lm(formula = y ~ w.hat, data = iv.dat)
```

```
##
```

Data

Data sources

1. Medicare Advantage, available [here](#)
2. Area Health Resources Files, available [here](#)

What is Medicare Advantage

- Private provision of health insurance for Medicare beneficiaries
- Medicare "replacement" plans
- It's just private insurance for Medicare folks

Medicare Advantage History

- Existed since 1980s, formalized in the 1990s, expanded in 2000s
- Medicare+Choice as part of Balanced Budget Act in 1997
- Largest expansion: Medicare Modernization Act in 2003 (also brought Medicare Part D)

Medicare Advantage Details

In its current form...

- Insurers submit plan details and a price needed to cover traditional Medicare ("bid")
- If approved, Medicare pays risk-adjusted bid *or* benchmark
- Bid < benchmark, insurer gets a rebate
- Bid > benchmark, insurer charges premium
- Seperate bidding for Part D

Medicare Advantage in Real Life

Let's take a look at the Medicare Advantage plan options...

Medicare Plan Finder

Datasets

My code files to read in the data are available here...

1. MA benchmark data, [1_benchmark.R](#)
2. MA market share data, [2_penetration.R](#)
3. AHRF data, [3_ahrf.R](#)

MA Benchmarks

- Reflects a payment that CMS sets for each enrollee of a MA plan
- Varies by county
- CMS doesn't actually pay this rate to each plan...final payment depends on the plan's "bid"
- Technically varies by star rating beginning 2012, but for our purposes, I take the benchmark rate for a 3-star plan

MA penetration data

- Enrollments published at the plan/county/year/month level
- Not much change across months in the same year
- For this, I take the penetration data for December of each year see [MA repo](#) for code extracting all months and averaging

AHRF data

- Easy to download but can be tricky to read into `R`
- Use the `SAScii` package to read in the full ascii (.asc) file using the .sas input file provided in the documentation
- Variable of interest, *F15299*

Analysis

Basic IV specification

- Lots of packages, including `ivreg`
- Since we have panel data, I'll use `felm` with the instrument specification,

```
felm(medicare_ffs_exp ~ 0 | year + ssa | (penetration ~ bench_pay) | ssa,  
     data=full.data)
```

First stage and reduced form

These terms are frustrating, but they are common in the literature

- First stage: Regression of your endogenous variable on all exog. covariates and instrument
- Reduced form: Regression of outcome on all exog. variables and the instrument

IV in practice?

- Ideally, we would spend more time talking about the assumptions, testing those assumptions, and what to do when those assumptions fail
- For now, curious minds can take a look at my [Navigating Empirical Microeconomics](#) project. It is very raw but the IV section is somewhat filled out.

What do we estimate?

- *Not* an average treatment effect
- Estimate Local Average Treatment Effect (LATE) if monotonicity assumption holds
- Otherwise, not clear...need to do more work

Extensions

Lots of new methods

- Sensitivity to outliers, `riv`
- Violations of the exclusion restriction, "plausibly exogenous" work from Conley et al. (2012). `plausexog` in Stata.
- Restore unbiasedness in your IV estimate with assumptions on sign of the first stage, Andrews and Armstrong (2017)