

# What Does IV Identifies: An Overview

Yagan Hazard

February 16, 2026— Advanced IV Methods

# Credits

These slides were first constructed by Peter Hull.

# Outline

## **Preliminaries**

### **IV Mechanics**

Just-Identified IV

Overidentification

Weak vs. Many-Weak Bias

### **IV Interpretation**

LATE Fundamentals

Generalizations and Limitations

Characterizing Compliers

Diff-in-Diff and IV

# Preliminaries: Parameters, Estimands, and Estimators

Three distinct objects, not always clearly distinguished:

# Preliminaries: Parameters, Estimands, and Estimators

Three distinct objects, not always clearly distinguished:

- **Parameters** come from economic (or other) models of the world
  - E.g. a “structural” model of supply and demand, or a potential outcome model relating schooling to earnings
  - They set the target for empirical analyses: what we want to know

# Preliminaries: Parameters, Estimands, and Estimators

Three distinct objects, not always clearly distinguished:

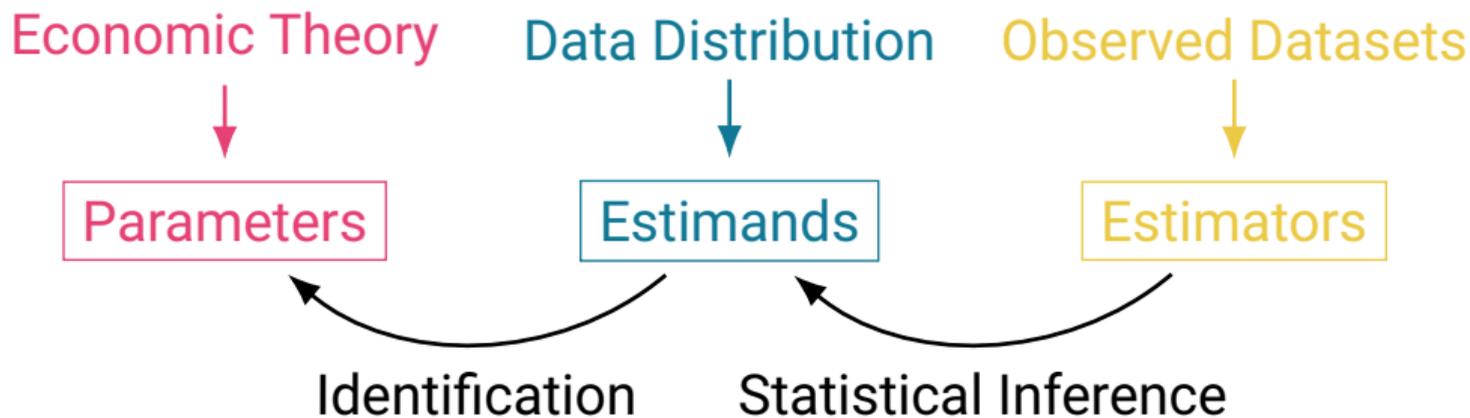
- **Parameters** come from economic (or other) models of the world
  - E.g. a “structural” model of supply and demand, or a potential outcome model relating schooling to earnings
  - They set the target for empirical analyses: what we want to know
- **Estimands** are functions of the population data distribution
  - E.g. a difference in means or ratio of population regression coef's
  - We make assumptions to link parameters & estimands (identification)

# Preliminaries: Parameters, Estimands, and Estimators

Three distinct objects, not always clearly distinguished:

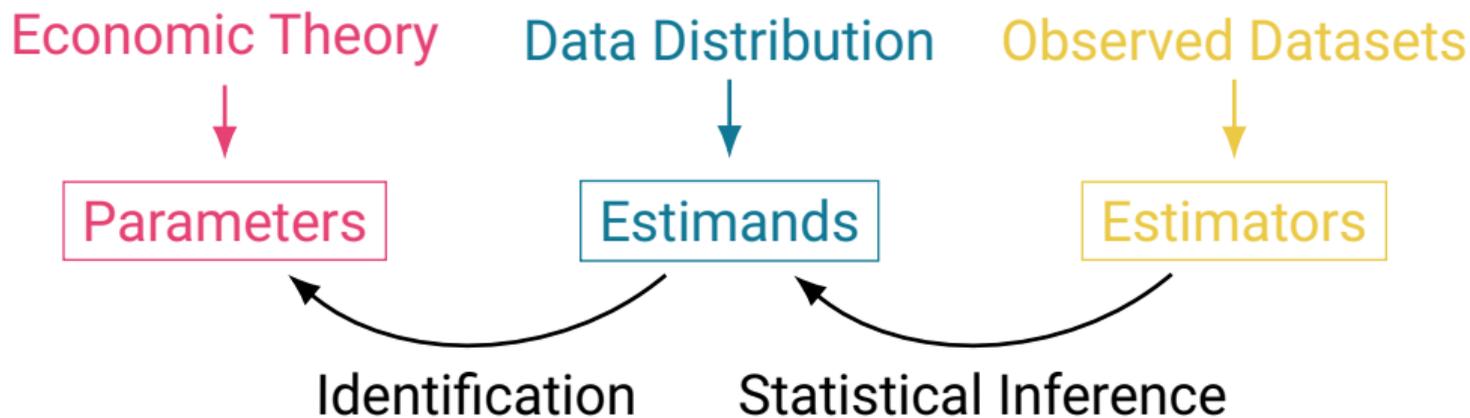
- **Parameters** come from economic (or other) models of the world
  - E.g. a “structural” model of supply and demand, or a potential outcome model relating schooling to earnings
  - They set the target for empirical analyses: what we want to know
- **Estimands** are functions of the population data distribution
  - E.g. a difference in means or ratio of population regression coef's
  - We make assumptions to link parameters & estimands (identification)
- **Estimators** are functions of observed data (i.e. the “sample”)
  - E.g. a difference in sample means or ratio of OLS coefficients
  - Since data are random, so are estimators. Each has a distribution
  - We use knowledge of estimator distributions to learn about estimands (inference) and thus identified parameters

# The Key Concepts in Statistical Inference



Separating out the very  $\neq$  tasks of identification and inference helps tremendously!

# The Key Concepts in Statistical Inference



Separating out the very  $\neq$  tasks of identification and inference helps tremendously!

Structure for today:

- recap of how IV estimands are structured [+ brief mention of estimation],
- then focus on what they *identify*.

## An Example

Human capital theory (e.g. Becker, 1957) tells us that taking econometrics classes will likely boost later-life productivity

# An Example

Human capital theory (e.g. Becker, 1957) tells us that taking econometrics classes will likely boost later-life productivity

- Parameter: returns to enrolling in this class  $\beta$ , measured in some outcome  $Y_i$  (e.g. lifetime earnings or # of top publications)
- Simple causal model:  $Y_i = \beta D_i + \varepsilon_i$ , where  $D_i \in \{0, 1\}$  indicates class enrollment and  $\varepsilon_i$  is  $i$ 's outcome without the class

# An Example

Human capital theory (e.g. Becker, 1957) tells us that taking econometrics classes will likely boost later-life productivity

- Parameter: returns to enrolling in this class  $\beta$ , measured in some outcome  $Y_i$  (e.g. lifetime earnings or # of top publications)
- Simple causal model:  $Y_i = \beta D_i + \varepsilon_i$ , where  $D_i \in \{0, 1\}$  indicates class enrollment and  $\varepsilon_i$  is  $i$ 's outcome without the class

We observe  $Y_i$  and  $D_i$  for some sample of individuals  $i = 1, \dots, N$

- We fire up Stata and `reg Y D, r`.

# An Example

Human capital theory (e.g. Becker, 1957) tells us that taking econometrics classes will likely boost later-life productivity

- Parameter: returns to enrolling in this class  $\beta$ , measured in some outcome  $Y_i$  (e.g. lifetime earnings or # of top publications)
- Simple causal model:  $Y_i = \beta D_i + \varepsilon_i$ , where  $D_i \in \{0, 1\}$  indicates class enrollment and  $\varepsilon_i$  is  $i$ 's outcome without the class

We observe  $Y_i$  and  $D_i$  for some sample of individuals  $i = 1, \dots, N$

- We fire up Stata and `reg Y D, r`. How do we interpret the results?

# Population Regression and Endogeneity

In large samples ( $N \rightarrow \infty$ ), the OLS estimator  $\hat{\beta}^{OLS}$  gets arbitrarily close to [i.e., consistently estimates] the regression estimand  $\beta^{OLS}$

# Population Regression and Endogeneity

In large samples ( $N \rightarrow \infty$ ), the OLS estimator  $\hat{\beta}^{OLS}$  gets arbitrarily close to [i.e., consistently estimates] the regression estimand  $\beta^{OLS}$

- The *identification* question is thus whether or not  $\beta^{OLS} = \beta$

# Population Regression and Endogeneity

In large samples ( $N \rightarrow \infty$ ), the OLS estimator  $\hat{\beta}^{OLS}$  gets arbitrarily close to [i.e., consistently estimates] the regression estimand  $\beta^{OLS}$

- The *identification* question is thus whether or not  $\beta^{OLS} = \beta$

By definition,  $\beta^{OLS} = \frac{Cov(Y_i, D_i)}{Var(D_i)}$ .

# Population Regression and Endogeneity

In large samples ( $N \rightarrow \infty$ ), the OLS estimator  $\hat{\beta}^{OLS}$  gets arbitrarily close to [i.e., consistently estimates] the regression estimand  $\beta^{OLS}$

- The *identification* question is thus whether or not  $\beta^{OLS} = \beta$

By definition,  $\beta^{OLS} = \frac{Cov(Y_i, D_i)}{Var(D_i)}$ . Plugging in our causal model:

$$\begin{aligned}\beta^{OLS} &= \frac{Cov(\beta D_i + \varepsilon_i, D_i)}{Var(D_i)} \\ &= \beta + \frac{Cov(\varepsilon_i, D_i)}{Var(D_i)}\end{aligned}$$

Thus, we have (regression) identification if and only if  $Cov(\varepsilon_i, D_i) = 0$ .

Otherwise, selection bias: people with certain potential outcomes  $\varepsilon_i$  are more/less likely to take this class, such that  $Cov(\varepsilon_i, D_i) \neq 0$

# Population Regression and Endogeneity

In large samples ( $N \rightarrow \infty$ ), the OLS estimator  $\hat{\beta}^{OLS}$  gets arbitrarily close to [i.e., consistently estimates] the regression estimand  $\beta^{OLS}$

- The *identification* question is thus whether or not  $\beta^{OLS} = \beta$

By definition,  $\beta^{OLS} = \frac{Cov(Y_i, D_i)}{Var(D_i)}$ . Plugging in our causal model:

$$\begin{aligned}\beta^{OLS} &= \frac{Cov(\beta D_i + \varepsilon_i, D_i)}{Var(D_i)} \\ &= \beta + \frac{Cov(\varepsilon_i, D_i)}{Var(D_i)}\end{aligned}$$

# Population Regression and Endogeneity

In large samples ( $N \rightarrow \infty$ ), the OLS estimator  $\hat{\beta}^{OLS}$  gets arbitrarily close to [i.e., consistently estimates] the regression estimand  $\beta^{OLS}$

- The *identification* question is thus whether or not  $\beta^{OLS} = \beta$

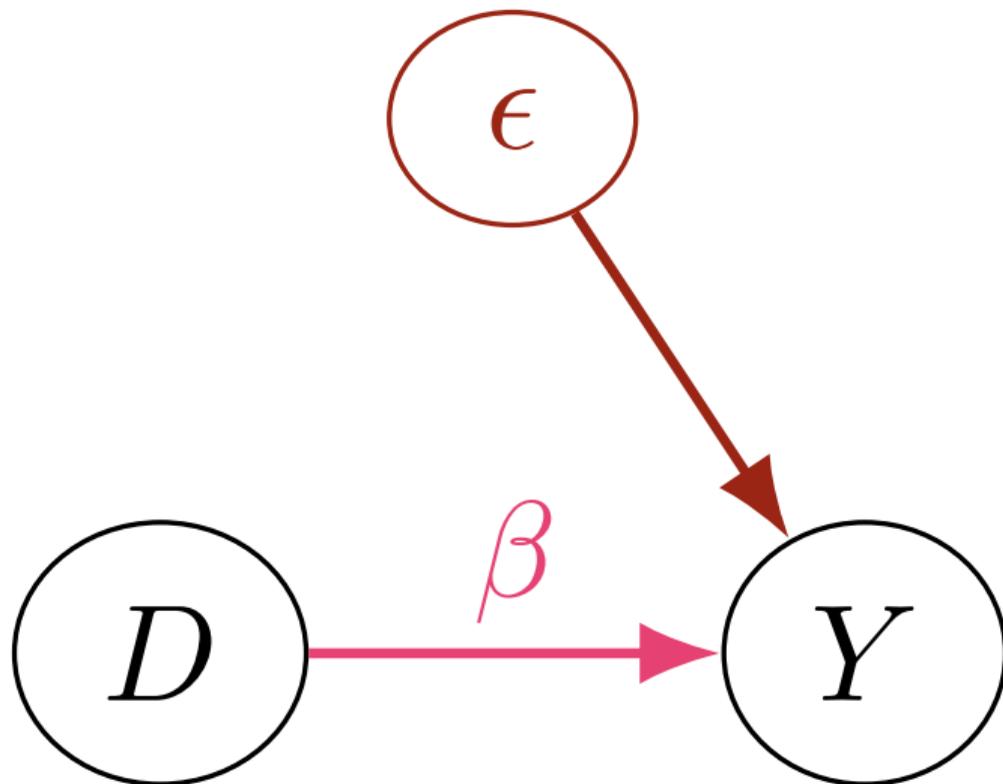
By definition,  $\beta^{OLS} = \frac{Cov(Y_i, D_i)}{Var(D_i)}$ . Plugging in our causal model:

$$\begin{aligned}\beta^{OLS} &= \frac{Cov(\beta D_i + \varepsilon_i, D_i)}{Var(D_i)} \\ &= \beta + \frac{Cov(\varepsilon_i, D_i)}{Var(D_i)}\end{aligned}$$

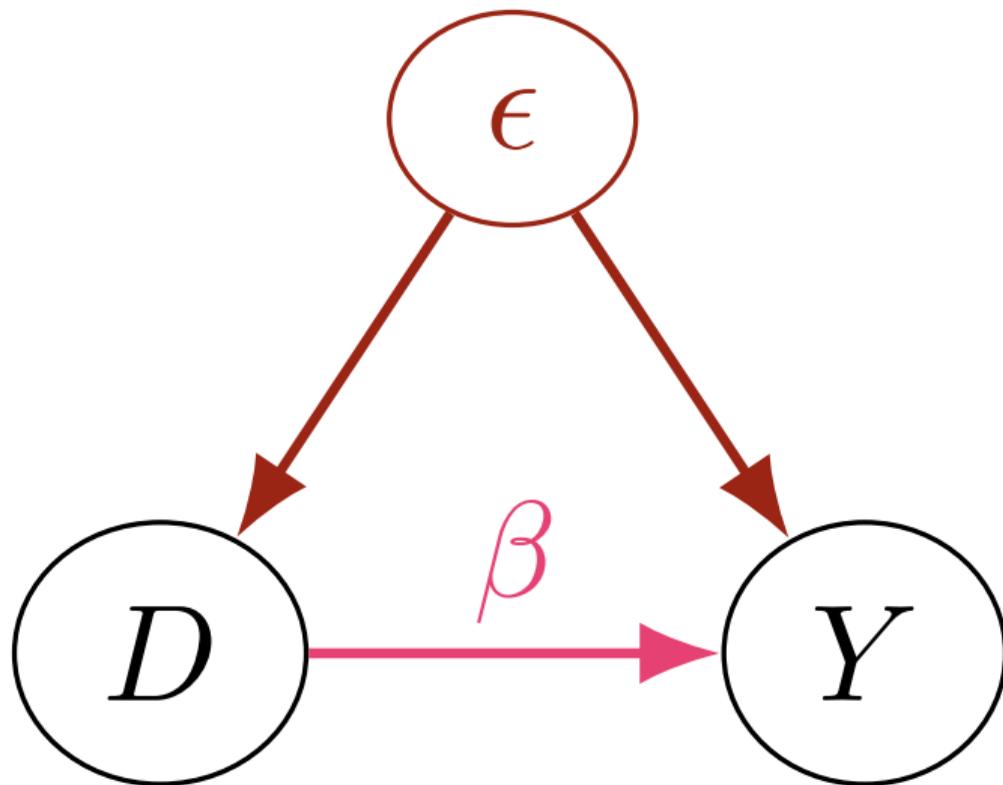
Thus, we have (regression) identification if and only if  $Cov(\varepsilon_i, D_i) = 0$ .

Otherwise, selection bias: people with certain potential outcomes  $\varepsilon_i$  are more/less likely to take this class, such that  $Cov(\varepsilon_i, D_i) \neq 0$

# Regression "Exogeneity"



# Regression "Endogeneity"



# Outline

## Preliminaries

## IV Mechanics

Just-Identified IV

Overidentification

Weak vs. Many-Weak Bias

## IV Interpretation

LATE Fundamentals

Generalizations and Limitations

Characterizing Compliers

Diff-in-Diff and IV

# The IV Solution

Suppose this class was “oversubscribed” (i.e., more people wanted to attend than we had space for), so seats were determined by lottery.

- $Z_i \in \{0, 1\}$  indicates randomized admission to the class

## The IV Solution

Suppose this class was “oversubscribed” (i.e., more people wanted to attend than we had space for), so seats were determined by lottery.

- $Z_i \in \{0, 1\}$  indicates randomized admission to the class
- Randomness + no direct effects of  $Z_i$  on  $Y_i$  implies  $Cov(Z_i, \varepsilon_i) = 0$ .

# The IV Solution

Suppose this class was “oversubscribed” (i.e., more people wanted to attend than we had space for), so seats were determined by lottery.

- $Z_i \in \{0, 1\}$  indicates randomized admission to the class
- Randomness + no direct effects of  $Z_i$  on  $Y_i$  implies  $Cov(Z_i, \varepsilon_i) = 0$ .

Plugging in the model for  $\varepsilon_i = Y_i - \beta D_i$ , we now have IV identification:

$$Cov(Z_i, Y_i - \beta D_i) = 0 \implies \beta = \frac{Cov(Z_i, Y_i)}{Cov(Z_i, D_i)} \equiv \beta^{IV},$$

so long as  $Cov(Z_i, D_i) \neq 0$ .

# The IV Solution

Suppose this class was “oversubscribed” (i.e., more people wanted to attend than we had space for), so seats were determined by lottery.

- $Z_i \in \{0, 1\}$  indicates randomized admission to the class
- Randomness + no direct effects of  $Z_i$  on  $Y_i$  implies  $Cov(Z_i, \varepsilon_i) = 0$ .

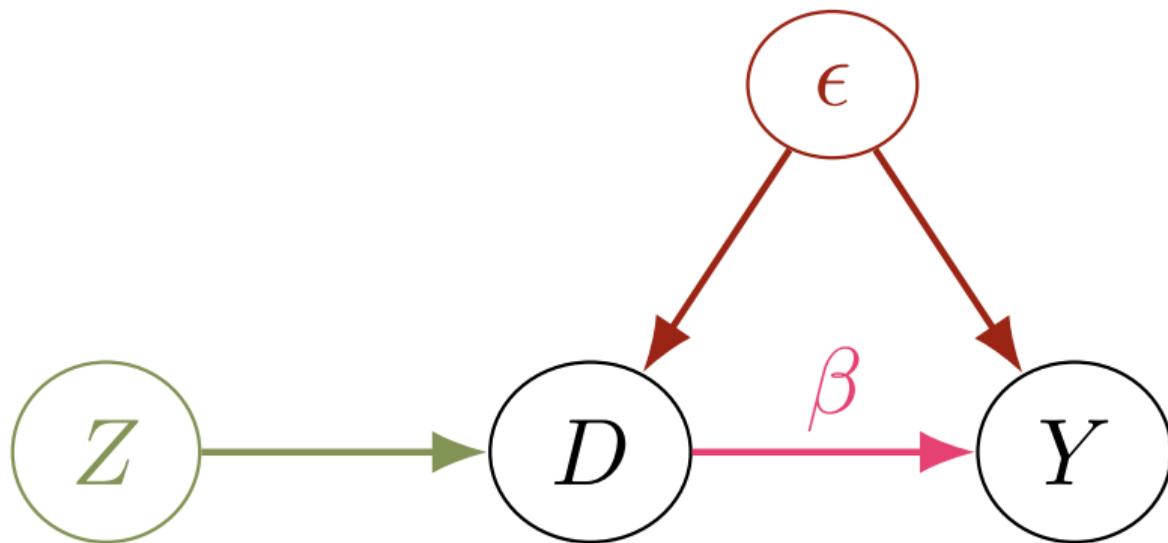
Plugging in the model for  $\varepsilon_i = Y_i - \beta D_i$ , we now have IV identification:

$$Cov(Z_i, Y_i - \beta D_i) = 0 \implies \beta = \frac{Cov(Z_i, Y_i)}{Cov(Z_i, D_i)} \equiv \beta^{IV},$$

so long as  $Cov(Z_i, D_i) \neq 0$ . We can estimate this by  $\hat{\beta}^{IV} = \frac{\widehat{Cov}(Z_i, Y_i)}{\widehat{Cov}(Z_i, D_i)} = (\mathbf{Z}'\mathbf{D})^{-1}\mathbf{Z}'\mathbf{Y}$

[Or, in Stata, `ivreg2 Y (D=Z), r]`

# The IV Solution



Note: no arrow connecting  $\epsilon$  and  $Z$  (“as-good-as-random assignment”), and no arrow from  $Z$  to  $Y$  directly (“exclusion”). We’ll come back to both.

## Reduced Form and First Stage

We're usually pretty comfortable w/regression; how does it connect to IV?

$$\beta^{IV} = \frac{Cov(Z_i, Y_i)}{Cov(Z_i, D_i)} = \frac{Cov(Z_i, Y_i)/Var(Z_i)}{Cov(Z_i, D_i)/Var(Z_i)} \equiv \frac{\rho}{\pi}$$

where  $\rho$  and  $\pi$  come from two population regressions:

$$Y_i = \kappa + \rho Z_i + \nu_i \quad \text{The "reduced form"}$$

$$D_i = \mu + \pi Z_i + \eta_i \quad \text{The "first stage"}$$

## Angrist (1990): Draft Lottery IV

Angrist (1990) [Vietnam draft lottery as IV to estimate earnings effects of military service]

- $Z_i \in \{0, 1\}$  an indicator for draft eligibility,
- $D_i \in \{0, 1\}$  an indicator for military service,
- $Y_i$  measures later-life earnings.

## Angrist (1990): Draft Lottery IV

Angrist (1990) [Vietnam draft lottery as IV to estimate earnings effects of military service]

- $Z_i \in \{0, 1\}$  an indicator for draft eligibility,
- $D_i \in \{0, 1\}$  an indicator for military service,
- $Y_i$  measures later-life earnings.

$$\beta^{IV} = \frac{Cov(Z_i, Y_i)/Var(Z_i)}{Cov(Z_i, D_i)/Var(Z_i)} = \frac{E[Y_i | Z_i = 1] - E[Y_i | Z_i = 0]}{E[D_i | Z_i = 1] - E[D_i | Z_i = 0]} \text{ as } Z_i \text{ is binary}$$

→ the famous **Wald estimand**

Regression on binary var.

## Angrist (1990): Draft Lottery IV

Angrist (1990) [Vietnam draft lottery as IV to estimate earnings effects of military service]

- $Z_i \in \{0, 1\}$  an indicator for draft eligibility,
- $D_i \in \{0, 1\}$  an indicator for military service,
- $Y_i$  measures later-life earnings.

$$\beta^{IV} = \frac{Cov(Z_i, Y_i)/Var(Z_i)}{Cov(Z_i, D_i)/Var(Z_i)} = \frac{E[Y_i | Z_i = 1] - E[Y_i | Z_i = 0]}{E[D_i | Z_i = 1] - E[D_i | Z_i = 0]} \text{ as } Z_i \text{ is binary}$$

→ the famous **Wald estimand**

Regression on binary var.

- $E[D_i | Z_i = 1] - E[D_i | Z_i = 0]$ : effect of eligibility on the *probability* of military service (because  $D_i$  is binary)
- $E[Y_i | Z_i = 1] - E[Y_i | Z_i = 0]$ : effect of eligibility on adult earnings [in 1971, 1981...]

IV interprets the latter causal effect [reduced form] in terms of the former [first stage].

#### IV Estimates of the Effects of Military Service on the Earnings of White Men born in 1950

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)	.267	.159 (.040)	-2,741 (1,324)
1971	3,338	-325.9 (46.6)			-2050 (293)
1969	2,299	-2.0 (34.5)			

Note: Adapted from Table 5 in Angrist and Krueger (1999) and author tabulations. 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.

## Adding Controls

We might only think our  $Z_i$  is exogenous controlling (linearly) for some vector  $W_i$

- Just add controls to the reduced form and first stage!  $\beta^{IV} = \frac{\rho}{\pi}$  for

$$Y_i = \kappa + \rho Z_i + W_i' \phi + \nu_i$$

$$D_i = \mu + \pi Z_i + W_i' \psi + \eta_i$$

with the estimator  $\hat{\beta}^{IV}$  defined analogously

## Adding Controls

We might only think our  $Z_i$  is exogenous controlling (linearly) for some vector  $W_i$

- Just add controls to the reduced form and first stage!  $\beta^{IV} = \frac{\rho}{\pi}$  for

$$Y_i = \kappa + \rho Z_i + W_i' \phi + \nu_i$$

$$D_i = \mu + \pi Z_i + W_i' \psi + \eta_i$$

with the estimator  $\hat{\beta}^{IV}$  defined analogously

- The Frisch-Waugh-Lovell theorem tells us the effective instrument is now  $\tilde{Z}_i$ : the residuals from regressing  $Z_i$  on  $W_i$  in the population

# Adding Controls

We might only think our  $Z_i$  is exogenous controlling (linearly) for some vector  $W_i$

- Just add controls to the reduced form and first stage!  $\beta^{IV} = \frac{\rho}{\pi}$  for

$$Y_i = \kappa + \rho Z_i + W_i' \phi + \nu_i$$

$$D_i = \mu + \pi Z_i + W_i' \psi + \eta_i$$

with the estimator  $\hat{\beta}^{IV}$  defined analogously

- The Frisch-Waugh-Lovell theorem tells us the effective instrument is now  $\tilde{Z}_i$ : the residuals from regressing  $Z_i$  on  $W_i$  in the population
- E.g. if  $W_i$  is a vector of dummies for randomization strata in an RCT, then  $\tilde{Z}_i$  captures the within-strata variation in  $Z_i$

# Multiple Treatments

We might be interested in a multi-dimensional model:  $Y_i = X_i' \beta + \varepsilon_i$

- Instrument vector  $Z_i$ , with  $L = \dim(Z_i) = \dim(X_i) = J$  (“just-identified”)
- Population residuals  $\tilde{Z}_i$ , given some vector of controls  $W_i$

## Multiple Treatments

We might be interested in a multi-dimensional model:  $Y_i = X_i' \beta + \varepsilon_i$

- Instrument vector  $Z_i$ , with  $L = \dim(Z_i) = \dim(X_i) = J$  (“just-identified”)
- Population residuals  $\tilde{Z}_i$ , given some vector of controls  $W_i$

Suppose  $Cov(\tilde{Z}_i, \varepsilon_i) = 0$ . Then, just as before, we have identification:

$$Cov(\tilde{Z}_i, Y_i - X_i' \beta) = 0 \implies \beta = Cov(\tilde{Z}_i, X_i)^{-1} Cov(\tilde{Z}_i, Y_i) \equiv \beta^{IV},$$

so long as  $Cov(\tilde{Z}_i, X_i)$  [a  $L \times L$  matrix] is full-rank.

# Multiple Treatments

We might be interested in a multi-dimensional model:  $Y_i = X_i' \beta + \varepsilon_i$

- Instrument vector  $Z_i$ , with  $L = \dim(Z_i) = \dim(X_i) = J$  (“just-identified”)
- Population residuals  $\tilde{Z}_i$ , given some vector of controls  $W_i$

Suppose  $Cov(\tilde{Z}_i, \varepsilon_i) = 0$ . Then, just as before, we have identification:

$$Cov(\tilde{Z}_i, Y_i - X_i' \beta) = 0 \implies \beta = Cov(\tilde{Z}_i, X_i)^{-1} Cov(\tilde{Z}_i, Y_i) \equiv \beta^{IV},$$

so long as  $Cov(\tilde{Z}_i, X_i)$  [a  $L \times L$  matrix] is full-rank. Equivalently,  $\beta^{IV} = \pi^{-1} \rho$  where:

$$Y_i = Z_i' \rho + W_i' \phi + \nu_i$$

$$X_i = \pi Z_i + W_i' \psi_i + \eta_i,$$

[Estimation as in simpler case: take residuals from  $J$  1<sup>st</sup>-stages and plug them into 2<sup>nd</sup> stage.]

## Multiple Instruments

What happens when  $\dim(Z_i) = L > J = \dim(X_i)$ ? **Overidentification:**

$$\text{Cov}(\tilde{Z}_i, Y_i - X_i' \beta) = 0 \implies \underbrace{\text{Cov}(\tilde{Z}_i, Y_i)}_{L \times 1} = \underbrace{\text{Cov}(\tilde{Z}_i, X_i)}_{L \times J} \underbrace{\beta}_{J \times 1}$$

so we can drop any  $L - J$  instruments and still identify  $\beta$ .

More generally, we can take any full-row rank linear combination  $\tilde{Z}_i^* = M \tilde{Z}_i$  [with  $M$  a  $J \times L$  full-row rank matrix] such that  $\text{Cov}(\tilde{Z}_i^*, X_i)$  is invertible

$$\begin{aligned} \underbrace{M \cdot \text{Cov}(\tilde{Z}_i, Y_i)}_{\text{Cov}(\tilde{Z}_i^*, Y_i)} &= \underbrace{M \cdot \text{Cov}(\tilde{Z}_i, X_i)}_{\text{Cov}(\tilde{Z}_i^*, X_i)} \beta \\ \implies \beta &= \left( M \cdot \text{Cov}(\tilde{Z}_i, X_i) \right)^{-1} M \cdot \text{Cov}(\tilde{Z}_i, Y_i) \end{aligned}$$

This defines a *class* of IV estimands/estimators, indexed by  $M$

## Multiple Instruments

$$\beta = \left( M \cdot Cov(\tilde{Z}_i, X_i) \right)^{-1} M \cdot Cov(\tilde{Z}_i, Y_i)$$

This defines a *class* of IV estimands/estimators, indexed by the  $J \times L$  matrix  $M$ .

Note that in the just-identified [i.e.,  $L = J$ ] case, the choice of  $M$  is irrelevant.

I.e., there is only one IV estimand,  $\beta^{IV} = Cov(\tilde{Z}_i, X_i)^{-1} Cov(\tilde{Z}_i, Y_i)$ .

Why?

## Multiple Instruments

$$\beta = \left( M \cdot Cov(\tilde{Z}_i, X_i) \right)^{-1} M \cdot Cov(\tilde{Z}_i, Y_i)$$

This defines a *class* of IV estimands/estimators, indexed by the  $J \times L$  matrix  $M$ .

Note that in the just-identified [i.e.,  $L = J$ ] case, the choice of  $M$  is irrelevant.

I.e., there is only one IV estimand,  $\beta^{IV} = Cov(\tilde{Z}_i, X_i)^{-1} Cov(\tilde{Z}_i, Y_i)$ .

**Why?** Because in this case,  $M$  is a full-row rank  $L \times L$  matrix  $\Rightarrow M$  is invertible, and therefore for any invertible  $M$ ,

$$\begin{aligned} \left( M \cdot Cov(\tilde{Z}_i, X_i) \right)^{-1} M \cdot Cov(\tilde{Z}_i, Y_i) &= Cov(\tilde{Z}_i, X_i)^{-1} \cdot M^{-1} \cdot M \cdot Cov(\tilde{Z}_i, Y_i) \\ &= Cov(\tilde{Z}_i, X_i)^{-1} \cdot Cov(\tilde{Z}_i, Y_i). \end{aligned}$$

## Two-Stage Least Squares (2SLS)

2SLS sets  $M = Cov(\tilde{Z}_i, X_i)' Var(\tilde{Z}_i)^{-1} = \pi$ : the  $J \times L$  matrix of first-stage coefficients

- this makes  $\tilde{Z}_i^* = \pi \tilde{Z}_i$  the (residualized) first-stage fitted values
- intuitively: it combines IVs according to their predictiveness of  $X_i$ .

[We should expect this to decrease the IV standard error by increasing fitted values variation.]

$$\beta^{2SLS} = \underbrace{Cov(\tilde{Z}_i, X_i)' Var(\tilde{Z}_i)^{-1} Cov(\tilde{Z}_i, X_i)}_M^{-1} \underbrace{Cov(\tilde{Z}_i, X_i)' Var(\tilde{Z}_i)^{-1} Cov(\tilde{Z}_i, Y_i)}_M$$

Since the first-stage from regressing  $X_i$  on this  $\tilde{Z}_i^*$  is one (by construction),  $\beta^{2SLS}$  can be obtained in two stages:

1. Regress  $X_i$  on  $Z_i$  and  $W_i$  (first stage)
2. Regress  $Y_i$  on first-stage fitted values and  $W_i$  (second stage)

Are we talking about an estimand or an estimator so far?

## Two-Stage Least Squares (2SLS)

2SLS sets  $M = Cov(\tilde{Z}_i, X_i)' Var(\tilde{Z}_i)^{-1} = \pi$ : the  $J \times L$  matrix of first-stage coefficients

- this makes  $\tilde{Z}_i^* = \pi \tilde{Z}_i$  the (residualized) first-stage fitted values
- intuitively: it combines IVs according to their predictiveness of  $X_i$ .

[We should expect this to decrease the IV standard error by increasing fitted values variation.]

$$\beta^{2SLS} = \underbrace{Cov(\tilde{Z}_i, X_i)' Var(\tilde{Z}_i)^{-1} Cov(\tilde{Z}_i, X_i)}_M^{-1} \underbrace{Cov(\tilde{Z}_i, X_i)' Var(\tilde{Z}_i)^{-1} Cov(\tilde{Z}_i, Y_i)}_M$$

Since the first-stage from regressing  $X_i$  on this  $\tilde{Z}_i^*$  is one (by construction),  $\beta^{2SLS}$  can be obtained in two stages:

1. Regress  $X_i$  on  $Z_i$  and  $W_i$  (first stage)
2. Regress  $Y_i$  on first-stage fitted values and  $W_i$  (second stage)

Are we talking about an estimand or an estimator so far? → the 2SLS estimand!

# Two-Stage Least Squares (2SLS)

Are we talking about an estimand or an estimator so far?

# Two-Stage Least Squares (2SLS)

Are we talking about an estimand or an estimator so far? → the 2SLS estimand!

Yet the exact same logic holds for the 2SLS estimator [i.e., two stages of OLS] as it simply takes the logic to the sample. But do not 2SLS by hand (let softwares do it)! Why? Manual 2SLS Mistakes

$$\begin{aligned}\hat{\beta}^{2SLS} &= [\widehat{\text{Cov}}(\tilde{Z}_i, X_i)' \widehat{\text{Var}}(\tilde{Z}_i)^{-1} \widehat{\text{Cov}}(\tilde{Z}_i, X_i)]^{-1} \widehat{\text{Cov}}(\tilde{Z}_i, X_i)' \widehat{\text{Var}}(\tilde{Z}_i)^{-1} \widehat{\text{Cov}}(\tilde{Z}_i, Y_i) \\ &= (X' P_{\tilde{Z}} X)^{-1} X' P_{\tilde{Z}} Y \quad \text{where } X \text{ and } Y \text{ stack obs. of } X_i' \text{ and } Y_i\end{aligned}$$

$P_{\tilde{Z}}$  is the sample projection matrix on (sample) residualized  $\tilde{Z}_i$ . Specifically,

$P_{\tilde{Z}} = \hat{\tilde{Z}} \left( \hat{\tilde{Z}}' \hat{\tilde{Z}} \right)^{-1} \hat{\tilde{Z}}'$  where  $\hat{\tilde{Z}}$  stacks residuals from sample proj. of  $Z_i$  on controls.

Since  $P_{\tilde{Z}}$  is an idempotent and symmetric matrix, we can rewrite

$$\begin{aligned}\hat{\beta}^{2SLS} &= (X' P_{\tilde{Z}} X)^{-1} X' P_{\tilde{Z}} Y = (X' P_{\tilde{Z}} P_{\tilde{Z}} X)^{-1} X' P_{\tilde{Z}} Y = (X' P_{\tilde{Z}}' P_{\tilde{Z}} X)^{-1} X' P_{\tilde{Z}}' Y \\ &= ((P_{\tilde{Z}} X)' P_{\tilde{Z}} X)^{-1} (P_{\tilde{Z}} X)' Y\end{aligned}$$

i.e., formula for OLS reg. of  $Y_i$  on  $\hat{X}_i = 1^{\text{st}}$ -stage fitted values (partialling out controls): 2.S.L.S.!

## 2SLS Is a Many-Splendored Thing

Another really useful way to understand 2SLS with multiple instruments [and possibly multiple treatments] is as a **weighted average of just-identified IVs**:

$\beta^{2SLS}$

$$\begin{aligned} &= \underbrace{\left( \text{Cov}(\tilde{Z}_i, X_i)' \text{Var}(\tilde{Z}_i)^{-1} \text{Cov}(\tilde{Z}_i, X_i) \right)^{-1}}_{\pi} \underbrace{\left( \text{Cov}(\tilde{Z}_i, X_i)' \text{Var}(\tilde{Z}_i)^{-1} \text{Cov}(\tilde{Z}_i, Y_i) \right)}_{\pi} \\ &= \left( \pi \text{Var}(\tilde{Z}_i) \underbrace{\text{Var}(\tilde{Z}_i)^{-1} \text{Cov}(\tilde{Z}_i, X_i)'}_{\pi'} \right)^{-1} \pi \text{Var}(\tilde{Z}_i) \underbrace{\text{Var}(\tilde{Z}_i)^{-1} \text{Cov}(\tilde{Z}_i, Y_i)}_{\rho} \text{ as } \text{Var}(\tilde{Z}_i) \text{Var}(\tilde{Z}_i)^{-1} = I \\ &= (\pi \text{Var}(\tilde{Z}_i) \pi')^{-1} \pi \text{Var}(\tilde{Z}_i) \rho, \end{aligned}$$

i.e., formula for a  $\text{Var}(\tilde{Z}_i)$ -weighted reg. of reduced-form coefficients  $\rho$  [ $L \times 1$  vector] on the (transposed) matrix of 1<sup>st</sup>-stage coefficients  $\pi'$  [ $\pi$  is  $J \times L$ ,  $\pi'$  is  $L \times J$ .] [w/o constant].

## 2SLS Is a Many-Splendored Thing

$$\underbrace{\beta^{2SLS}}_{J \times 1} = \left( \underbrace{\pi}_{J \times L} \underbrace{\text{Var}(\tilde{Z}_i)}_{L \times L} \pi' \right)^{-1} \pi \text{Var}(\tilde{Z}_i) \underbrace{\rho}_{L \times 1}$$

When  $J = 1$  [one treatment] and still  $L > 1$  [multiple instruments], this becomes

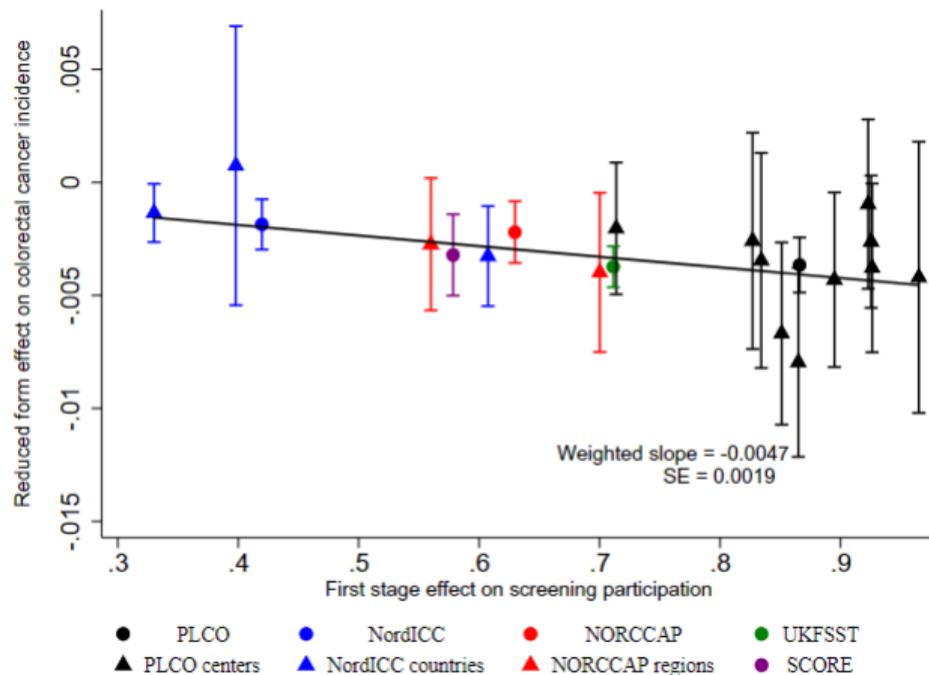
$$\beta^{2SLS} = \sum_{\ell} \omega_{\ell} \beta_{\ell}^{IV}$$

where  $\omega_{\ell} = (\pi \text{Var}(\tilde{Z}_i) \pi')^{-1} \underbrace{\pi}_{1 \times L} \underbrace{\text{Var}(\tilde{Z}_i)_{\cdot, \ell}}_{L \times 1} \pi_{\ell}$  and  $\beta_{\ell}^{IV} = \rho_{\ell} / \pi_{\ell}$ .

[ $\text{Var}(\tilde{Z}_i)_{\cdot, \ell}$  denotes the  $\ell^{\text{th}}$  column of  $\text{Var}(\tilde{Z}_i)$ .]

So 2SLS with one endogenous treatment and multiple instruments combines multiple “one-at-a-time” just-identified IVs  $\beta_{\ell}^{IV}$ .

# Angrist-Hull '23: "Visual IV" for Cancer Screening Trials



Each dot = a  $(\rho_\ell, \pi_\ell)$  for a trial  $\ell$  where randomized screening offers  $Z_i$  instrument for screening participation  $D_i$ . Slope of weighted line-of-best fit through 0 = 2SLS estimate.

## Overidentification Tests

Under the constant-effects causal model of  $Y_i = X_i'\beta + \varepsilon_i$ , overidentification gives a way to test instrument validity

- All just-identified IVs should coincide: i.e.  $\beta_\ell^{IV} = \beta$  for all  $\ell$
- Graphically: the  $R^2$  from visual IV plots should = 1 in large samples

## Overidentification Tests

Under the constant-effects causal model of  $Y_i = X_i'\beta + \varepsilon_i$ , overidentification gives a way to test instrument validity

- All just-identified IVs should coincide: i.e.  $\beta_\ell^{IV} = \beta$  for all  $\ell$
- Graphically: the  $R^2$  from visual IV plots should = 1 in large samples

Softwares like Stata automatically compute the  $p$ -value for this test when  $L > J$

- If  $p > 0.05$ , it means the  $\hat{\beta}_\ell^{IV}$ 's are all pretty similar to each other

# Overidentification Tests

Under the constant-effects causal model of  $Y_i = X_i'\beta + \varepsilon_i$ , overidentification gives a way to test instrument validity

- All just-identified IVs should coincide: i.e.  $\beta_\ell^{IV} = \beta$  for all  $\ell$
- Graphically: the  $R^2$  from visual IV plots should = 1 in large samples

Softwares like Stata automatically compute the  $p$ -value for this test when  $L > J$

- If  $p > 0.05$ , it means the  $\hat{\beta}_\ell^{IV}$ 's are all pretty similar to each other

Don't place too much stock in overidentification tests, however:

- They tend to have low power (b/c individual  $\hat{\beta}_\ell^{IV}$  tend to be noisy)
- Rejection doesn't tell us which IV is invalid (they all might be!)
- If they reject, it need not mean the instruments are invalid

[because of treatment effect heterogeneity → we're getting there!]

## Weak Instruments skip

When running just-identified IV, people worry about instrument “strength”

- Specifically the first stage F-statistic, which tests if  $\pi = 0$

## Weak Instruments skip

When running just-identified IV, people worry about instrument “strength”

- Specifically the first stage F-statistic, which tests if  $\pi = 0$

If  $\pi$  is small relative to its standard error, we say the IV is “weak”

- Typically use the rule-of-thumb of  $F < 10$  (Staiger and Stock 1997)
- In this case the second-stage SEs will be large and the 2SLS estimate will tend to be biased towards the corresponding OLS

## Weak Instruments skip

When running just-identified IV, people worry about instrument “strength”

- Specifically the first stage F-statistic, which tests if  $\pi = 0$

If  $\pi$  is small relative to its standard error, we say the IV is “weak”

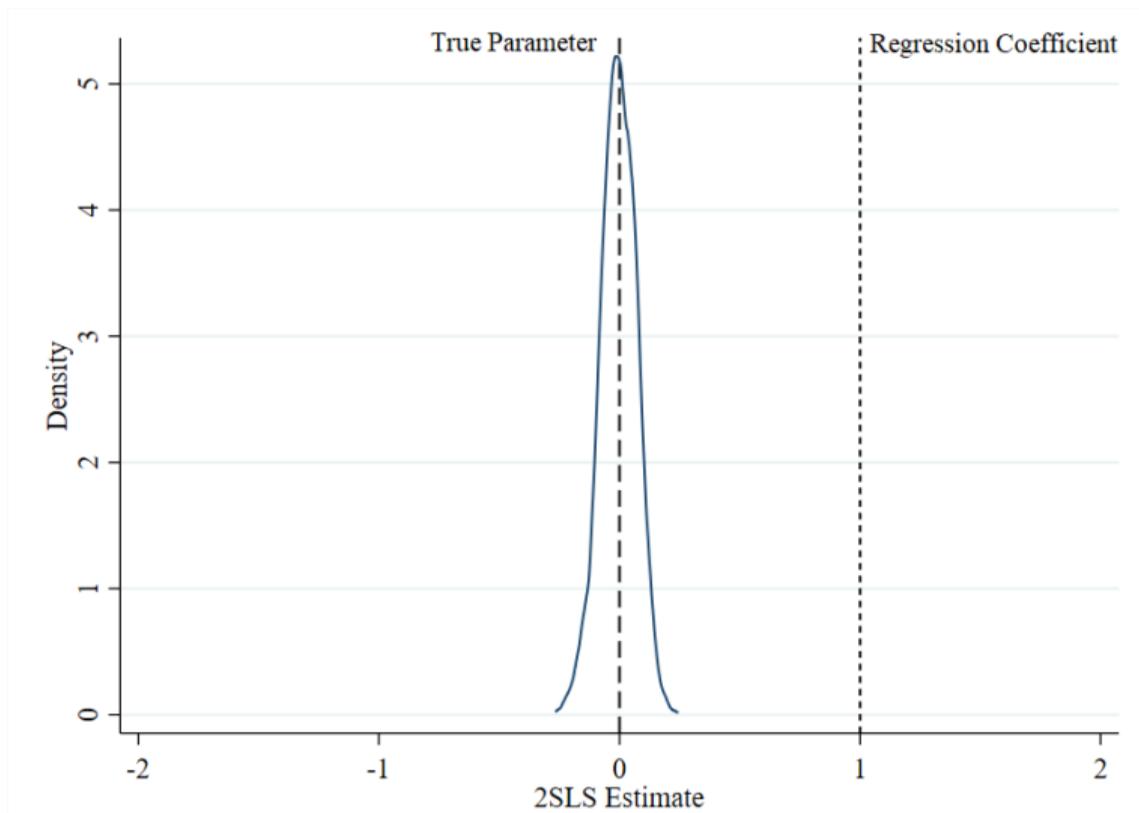
- Typically use the rule-of-thumb of  $F < 10$  (Staiger and Stock 1997)
- In this case the second-stage SEs will be large and the 2SLS estimate will tend to be biased towards the corresponding OLS

Much has been made of this over the years, but Angrist and Kolesár (2022) show that we shouldn't worry too much

- The SE increase tends to be large enough to “cover up” the bias, so you're unlikely to reject the null of  $\beta = 0$  spuriously

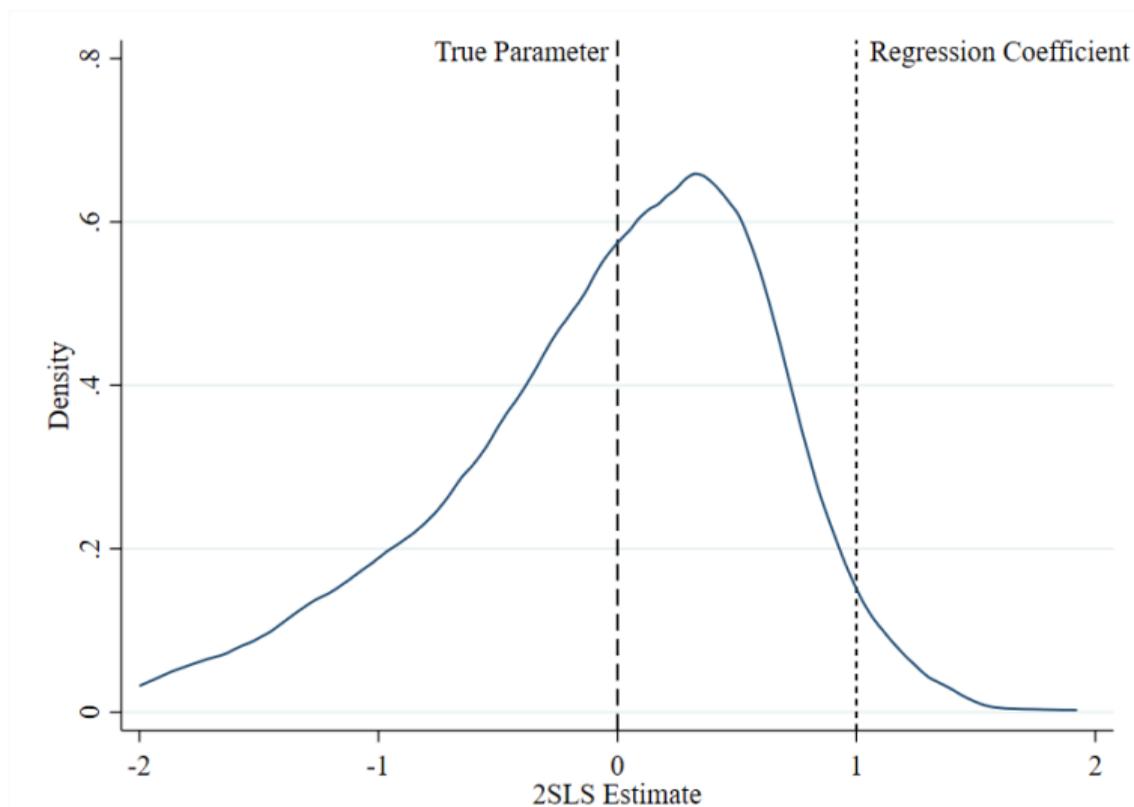
# Weak Instruments: Visualized

Monte Carlo:  $Y_i = 0 \cdot D_i + \varepsilon_i$ ,  $D_i = \pi Z_i + \eta_i$ :  $\pi = \text{Var}(\varepsilon_i) = \text{Var}(\eta_i) = 1$



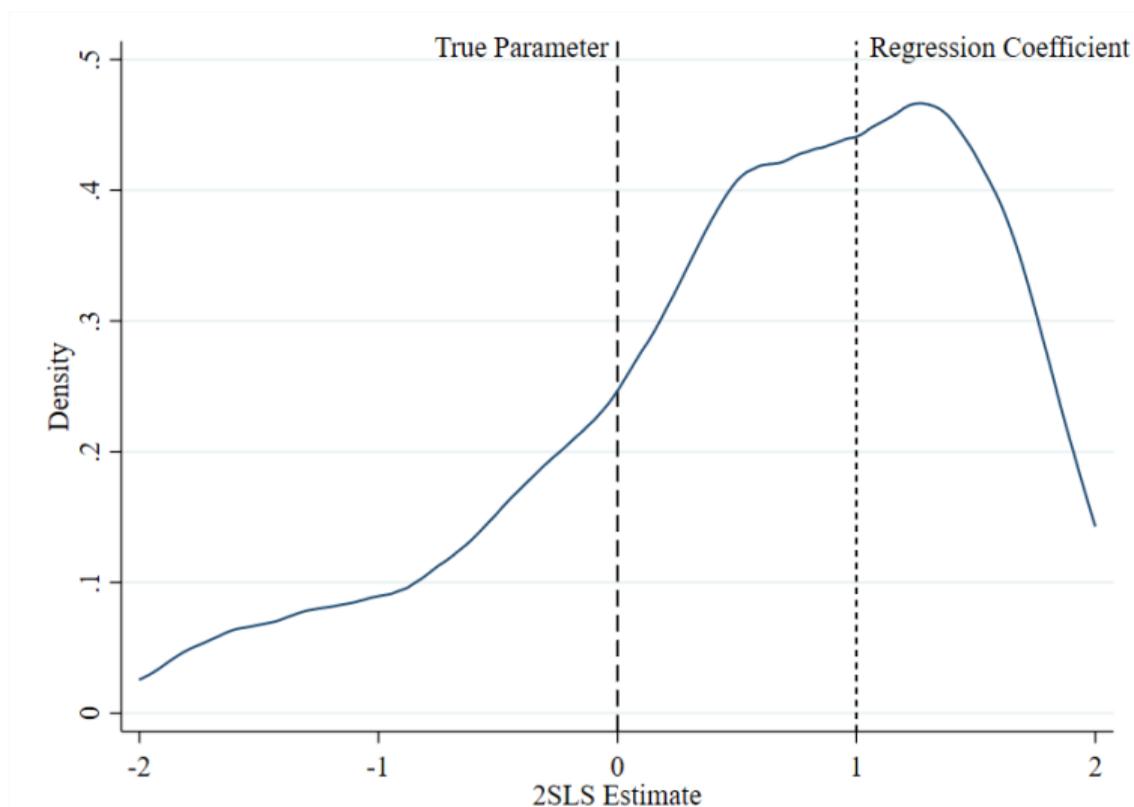
# Weak Instruments: Visualized

Monte Carlo:  $Y_i = 0 \cdot D_i + \varepsilon_i$ ,  $D_i = \pi Z_i + \eta_i$ :  $\pi = 0.1$  (Weaker)



# Weak Instruments: Visualized

Monte Carlo:  $Y_i = 0 \cdot D_i + \varepsilon_i$ ,  $D_i = \pi Z_i + \eta_i$ :  $\pi = 0.01$  (Very Weak)



# Many IVs

A thornier problem is many-weak bias, when overidentified

- This also tends to manifest in low first-stage  $F$ 's, and also causes 2SLS to be biased towards OLS

Unlike when just-id., however, with many weak IVs the SE's go *down*

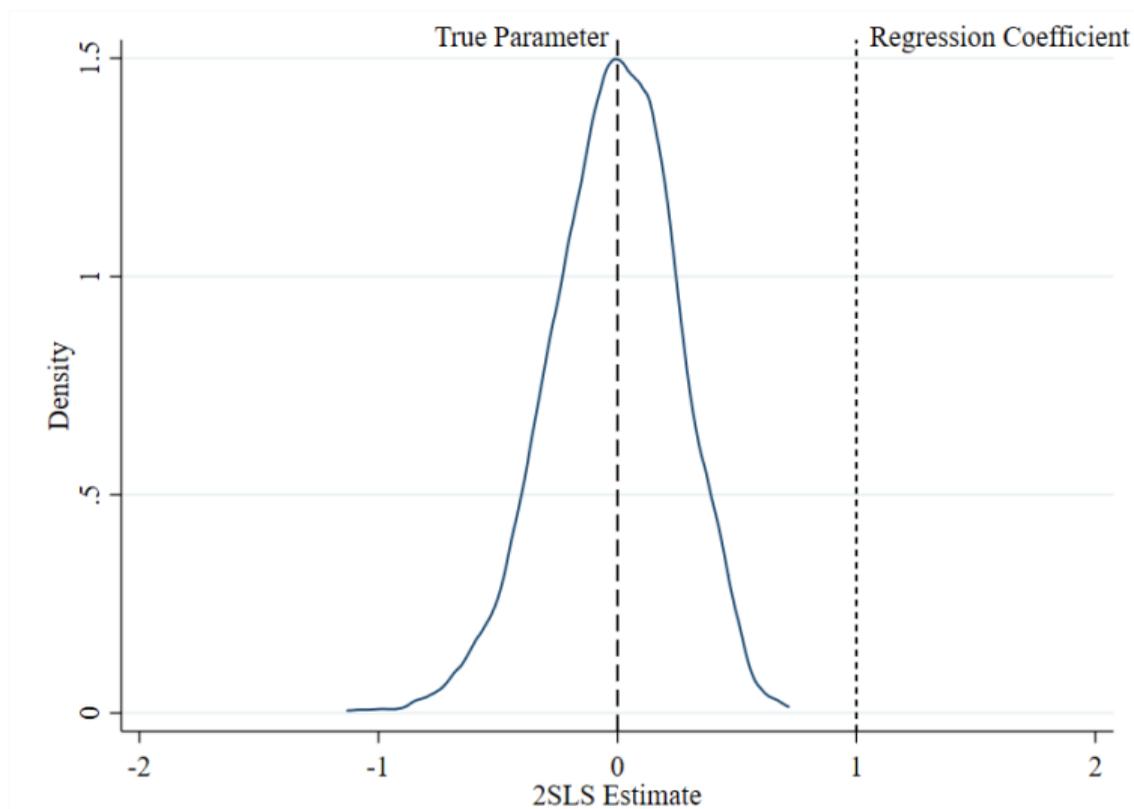
- Intuitively, a more flexible FS tends to fit  $D_i$  better  $\rightarrow$  more power
- But we can have *overfitting* with lots of instruments, which essentially recreates the (endogenous) variation in  $D_i$

This became a high-profile problem with Angrist-Krueger '91, where the QOB instrument was interacted with many state/year FEs

- These days folks don't make this mistake ... but many-IV bias can be lurking in other settings with constructed instruments (e.g. judge IV)

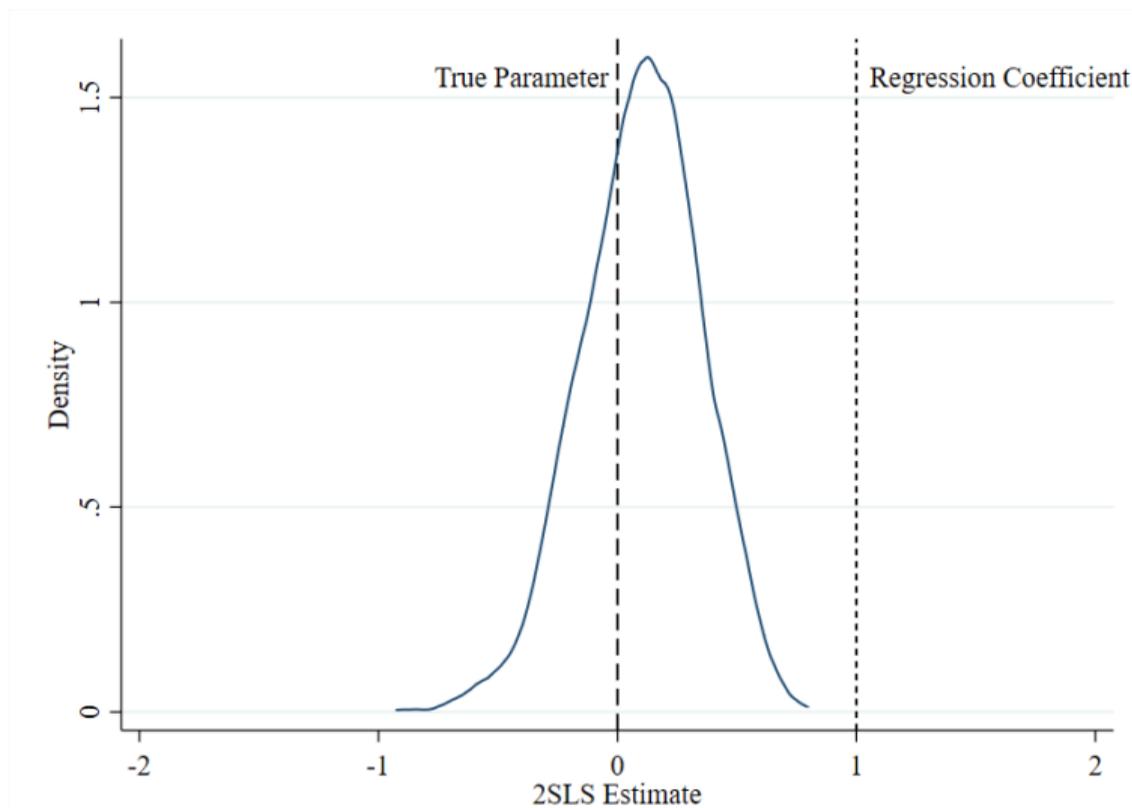
# Many Instruments: Visualized

$$Y_i = \beta \cdot D_i + \varepsilon_i, D_i = \pi Z_{i1} + \sum_{\ell > 1} \pi_{\ell} \cdot Z_{i\ell} + \eta_i: \text{IV w/ } Z_{i1} \text{ only}$$



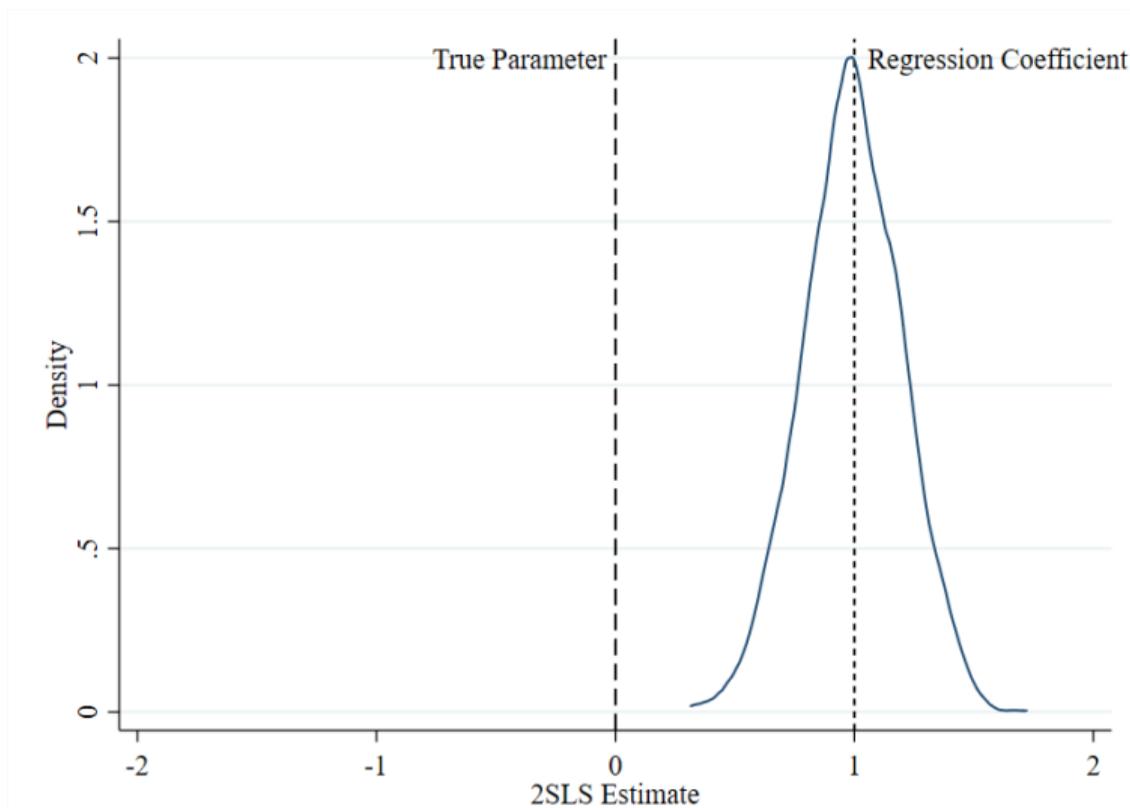
# Many Instruments: Visualized

$$Y_i = \beta \cdot D_i + \varepsilon_i, D_i = \pi Z_{i1} + \sum_{l>1} \pi_l \cdot Z_{il} + \eta_i: \text{IV w/ } Z_{i1}, \dots, Z_{i10}$$



# Many Instruments: Visualized

$$Y_i = \beta \cdot D_i + \varepsilon_i, D_i = \pi Z_{i1} + \sum_{l>1} \theta_l \cdot Z_{il} + \eta_i: \text{IV w/ } Z_{i1}, \dots, Z_{i100}$$



# What to Do?

Aim for few instruments, and check your  $F$ 's after every *ivreg*

- State of the art: Montiel Olea and Pflueger '15; `weakivtest` in Stata
- Staiger-Stock rule-of-thumb ( $F > 10$ ) still seems widely held
- Cf. Lee et al. (2020), Keane and Neal (2022) for discussions of additional subtleties

If your  $F$  is small, some things to consider:

- Is there a better functional form for your instrument?
- Do interactions with covariates help? (note: beware many-weak!)
- Does changing the covariate set help? (note: beware invalidity!)
- Check results w/a more robust approach (e.g. Anderson-Rubin, JIVE)

# Outline

## **Preliminaries**

## **IV Mechanics**

Just-Identified IV

Overidentification

Weak vs. Many-Weak Bias

## **IV Interpretation**

LATE Fundamentals

Generalizations and Limitations

Characterizing Compliers

Diff-in-Diff and IV

## What Does IV Identify, Really?

IV was invented for structural economic models (SEMs), typically w/ a single parameter  $\beta$  linearly relating  $Y_i$  to  $X_i$ . Yet modern view of  $Y_i = \beta X_i + \varepsilon_i$  is that it describes a **causal relationship** and imposes a (strong) **linear-and-constant-effects restriction**.

The Imbens-Angrist 1994 LATE result revolutionized our understanding of IV estimands, and clarified some subtle points around IV identification:

- $\beta^{IV}$  often identifies a convex average of heterogeneous effects under first-stage **monotonicity**:  $Z_i$  only affects  $X_i$  in one direction...  
though there are notable exceptions! [E.g., multiple treatments]
- IV “exogeneity” [ $Cov(Z_i, \varepsilon_i) = 0$ ] can arise from two conceptually different assumptions of instrument **independence** and **exclusion**.

## Basic (Binary Treatment, Binary Instrument) Setup

$Y_i(0), Y_i(1)$ : denote indiv.  $i$ 's potential outcomes given a binary treatment  $D_i \in \{0, 1\}$ .

→ Observed outcomes:  $Y_i = Y_i(0)(1 - D_i) + Y_i(1)D_i = (Y_i(1) - Y_i(0)) D_i + Y_i(0)$

$D_i(0), D_i(1)$ : denote indiv.  $i$ 's potential treatments given a binary instrument  $Z_i \in \{0, 1\}$ .

→ Observed treatment:  $D_i = (D_i(1) - D_i(0)) Z_i + D_i(0)$

Under what assumptions can we causally interpret the regression of  $Y$  on  $D$  instrumented by  $Z$ ?

# Imbens and Angrist (1994) Assumptions

1. *As-good-as-random assignment*:  $Z_i \perp (Y_i(0), Y_i(1), D_i(0), D_i(1))$

→ Consider the Angrist draft lottery, or Angrist-Krueger's QoB IV

2. *Exclusion*:  $Z_i$  only affects  $Y_i$  through its effect on  $D_i$

→ Implicit in the potential outcome notation:  $Y_i(d)$  is not indexed by  $Z_i$

3. *Relevance*:  $Z_i$  is correlated with  $D_i$

→ Equivalently, given Assumption 1,  $E[D_i(1) - D_i(0)] \neq 0$

4. *Monotonicity*:  $D_i(1) \geq D_i(0)$  for all  $i$  (i.e., almost-surely)

→ The instrument can only shift the treatment in one direction

I.e., there are only "Always-takers"  $D_i(1) = D_i(0) = 1$ , "Never-takers"  $D_i(1) = D_i(0) = 0$ , and "Compliers"  $1 = D_i(1) > D_i(0) = 0$  but no "Defiers"  $0 = D_i(1) < D_i(0) = 1$ .

# Local Average Treatment Effect (LATE) Identification

Imbens and Angrist (1994) showed that under these assumptions:

$$\beta^{IV} = E[Y_i(1) - Y_i(0) \mid D_i(1) > D_i(0)]$$

**Proof:** By independence,  $E[D_i \mid Z_i = 1] - E[D_i \mid Z_i = 0] = E[D_i(1) - D_i(0)]$

Similarly,

$$\begin{aligned} E[Y_i \mid Z_i = 1] - E[Y_i \mid Z_i = 0] &= E[Y_i(0) + (Y_i(1) - Y_i(0)) D_i(1) \mid Z_i = 1] \\ &\quad - E[Y_i(0) + (Y_i(1) - Y_i(0)) D_i(0) \mid Z_i = 0] \\ &= E[(Y_i(1) - Y_i(0)) (D_i(1) - D_i(0))] \end{aligned}$$

By monotonicity,  $D_i(1) - D_i(0) \in \{0, 1\}$ . Thus:

$$\frac{E[(Y_i(1) - Y_i(0)) (D_i(1) - D_i(0))]}{E[D_i(1) - D_i(0)]} = E[Y_i(1) - Y_i(0) \mid D_i(1) > D_i(0)]$$

# Implications

- Potential outcomes notation clearly distinguishes independence vs. exclusion  
→ A lottery'd  $Z_i$  can ensure independence, but exclusion can still fail.
- First-stage monotonicity becomes important under heterogeneous effects  
→ Otherwise,  $\beta^{IV}$  identifies a non-convex average: proof

$$\frac{E[(Y_i(1) - Y_i(0))(D_i(1) - D_i(0))]}{E[D_i(1) - D_i(0)]} = \frac{c}{c-d} E[Y_i(1) - Y_i(0) \mid D_i(1) > D_i(0)] + \frac{-d}{c-d} \underbrace{E[Y_i(1) - Y_i(0) \mid D_i(1) < D_i(0)]}_{\text{Avg. effect for defiers}}$$

- Monotonicity: clearly sensible in some settings [e.g., draft lottery], but can be questionable in others [e.g., judge IVs].
- The LATE result formalizes a key limitation of overidentification tests  
→ Two valid IVs can identify different LATEs [*internal vs. external validity*].

## Multivalued Ordered Treatments ("Variable Treatment Intensity")

Angrist and Imbens (1995) show that when  $D_i \in \{0, \dots, \bar{d}\}$  and a generalization of the LATE assumptions hold, IV identifies an "Average Causal Response" (ACR)

$$\frac{E[Y_i | Z_i = 1] - E[Y_i | Z_i = 0]}{E[D_i | Z_i = 1] - E[D_i | Z_i = 0]} = \sum_{d=1}^{\bar{d}} \omega_d E[Y_i(d) - Y_i(d-1) | D_i(1) \geq d > D_i(0)]$$

with  $\omega_d = \frac{\Pr(D_i(1) \geq d > D_i(0))}{\sum_{d'=1}^{\bar{d}} \Pr(D_i(1) \geq d' > D_i(0))}$  convex weights [under monotonicity].

Averages unit causal responses  $E[Y_i(d) - Y_i(d-1) | D_i(1) \geq d > D_i(0)]$

→ More weight on margins  $d$  with more first-stage "action"/"bite"

→ Note: "Compliers" with  $D_i(1) > D_i(0)$  can be double-counted at different margins

→ Weights are identified by difference in treatment CDF when  $Z_i = 1$  vs. when  $Z_i = 0$ .

# Multivalued Ordered Treatments

## Proof

Assumptions:

- Independence and Exclusion:  $(Y_i(0), \dots, Y_i(\bar{d}), D_i(0), D_i(1)) \perp\!\!\!\perp Z_i$
- First-Stage:  $E(D_i(1) - D_i(0)) \neq 0$
- Monotonicity:  $D_i(1) - D_i(0) \geq 0, \forall i$  (sign wlog).

Let  $I(A)$  be the indicator function for the event  $A$ . Define the following indicators:

$\lambda_{Zd} = I(D(Z) \geq d)$  for  $Z = 0, 1$  and  $d = 0, 1, 2, \dots, \bar{d} + 1$ . Note that  $\lambda_{Z0} = 1$  and  $\lambda_{Z\bar{d}+1} = 0$  for all  $Z$ . In terms of the  $\lambda_{Zd}$ ,  $Y$  can be written as [subscript  $i$  dropped to reduce notational burden]

$$\begin{aligned} Y &= Z \cdot Y(D(1)) + (1 - Z) \cdot Y(D(0)) \\ &= \left\{ Z \cdot \sum_{d=0}^{\bar{d}} Y(d) (\lambda_{1d} - \lambda_{1d+1}) \right\} + \left\{ (1 - Z) \cdot \sum_{d=0}^{\bar{d}} Y(d) (\lambda_{0d} - \lambda_{0d+1}) \right\}. \end{aligned}$$

as  $\lambda_{Zd} - \lambda_{Zd+1} = I(D(Z) \geq d) - I(D(Z) \geq d + 1) = I(d \leq D(Z) < d + 1)$

and  $Y(D(Z)) = \sum_{d=0}^{\bar{d}} Y(d) I(d \leq D(Z) < d + 1)$ .

# Multivalued Ordered Treatments

## Proof

Under independence,  $E[Y | Z = 1] - E[Y | Z = 0] = E[Y(D(1)) - Y(D(0))]$  which equals

$$\begin{aligned} E \left\{ \sum_{d=0}^{\bar{d}} Y_j \cdot [\lambda_{1d} - \lambda_{1d+1} - \lambda_{0d} + \lambda_{0d+1}] \right\} &= E \left\{ \sum_{d=1}^{\bar{d}} [(Y_d - Y_{d-1}) \cdot (\lambda_{1d} - \lambda_{0d})] + Y_0 \cdot \underbrace{(\lambda_{10} - \lambda_{00})}_{=0} \right\} \\ &= E \left\{ \sum_{d=1}^{\bar{d}} (Y_d - Y_{d-1}) \cdot (\lambda_{1d} - \lambda_{0d}) \right\} \end{aligned}$$

because  $\lambda_{Z0} = 1$  for  $Z = 0, 1$ . Note that  $\lambda_{1d} \geq \lambda_{0d}$  under monotonicity and that  $\lambda_{1d}$  and  $\lambda_{0d}$  equal 0 or 1. Therefore,  $\lambda_{1d} - \lambda_{0d}$  equals 0 or 1, and we can write the previous expression as

$$\begin{aligned} \sum_{d=1}^{\bar{d}} E [Y_d - Y_{d-1} | \lambda_{1d} - \lambda_{0d} = 1] \cdot \Pr(\lambda_{1d} - \lambda_{0d} = 1) \\ = \sum_{d=1}^{\bar{d}} E [Y_d - Y_{d-1} | D(1) \geq d > D(0)] \cdot \Pr(D(1) \geq d > D(0)) \end{aligned}$$

# Multivalued Ordered Treatments

## *Proof*

Similarly, for the denominator,  $D = Z \cdot D(1) + (1 - Z) \cdot D(0)$  and, because  $d$  plays the role played by  $Y(d)$  in previous derivations for the numerator,

$$\begin{aligned} & E[D \mid Z = 1] - E[D \mid Z = 0] \\ &= E \left\{ \sum_{d=0}^{\bar{d}} d \cdot (\lambda_{1d} - \lambda_{1d+1} - \lambda_{0d} + \lambda_{0d+1}) \right\} \\ &= E \left\{ \sum_{d=1}^{\bar{d}} \underbrace{(d - (d - 1))}_{=1} \cdot (\lambda_{1d} - \lambda_{0d}) \right\} = \sum_{d=1}^{\bar{d}} \Pr(D(1) \geq d > D(0)) \end{aligned}$$

**QED.**

# Multivalued Ordered Treatments

## Weights identification

Notice that each weight  $\omega_d = \frac{\Pr(D_i(1) \geq d > D_i(0))}{\sum_{d'=1}^{\bar{d}} \Pr(D_i(1) \geq d' > D_i(0))}$  are identified too.

The denominator is identified (as derived in the previous slide) by  $E[D | Z = 1] - E[D | Z = 0]$ .

As for the numerator of  $\omega_d$ , assuming that  $D_i$  is continuously distributed with no mass point,

$$\begin{aligned} P[D_i(1) \geq d > D_i(0)] &= P[D_i(1) \geq d] - P[D_i(0) \geq d] \\ &= P[D_i \geq d | Z_i = 1] - P[D_i \geq d | Z_i = 0] \quad \text{under independence} \\ &= E[I(D_i \geq d) | Z_i = 1] - E[I(D_i \geq d) | Z_i = 0] \end{aligned}$$

which corresponds to the coefficient on  $Z_i$  in a regression of  $I(D_i \geq d)$  on  $(1, Z_i)$ .

Notice this also equals the gap between the CDF of  $D_i$  when  $Z_i = 0$  vs.  $Z_i = 1$  [guaranteed to be positive under monotonicity]. Indeed,

$$P[D_i(1) \geq d] - P[D_i(0) \geq d] = \underbrace{P[D_i(0) < d]}_{F_{D_i(0)}(d)} - \underbrace{P[D_i(1) < d]}_{F_{D_i(1)}(d)}$$

# Continuous Treatments

Angrist, Graddy and Imbens (2000) extend this to a continuous treatment.

- Potential outcome function  $q_i(p)$ : demand for fish in market  $i$  at hypothetical price  $p$
- Treatment variable  $p_i$ : price of fish in market  $i$
- Observed outcome is  $q_i(p_i)$
- $\frac{\partial}{\partial p} q_i(p) = q'_i(p)$ : slope of the demand curve [price elasticity of demand if  $q_i$  and  $p_i$  are in logs], a.k.a. a  $i$ -specific (marginal) treatment effect of  $p$  on  $q_i$  [as  $q'_i(p) = \frac{q_i(p) - q_i(p-h)}{h}$ ]
- Instrument: bad weather indicator  $stormy_i \in \{0, 1\}$  [= exogenous supply shock].

Assumptions:

- Independence and Exclusion:  $(q_i(p), p_i(z)) \perp\!\!\!\perp stormy_i, \forall (p, z) \in Supp(p_i) \times \{0, 1\}$
- First-Stage:  $E(p_i(1) - p_i(0)) \neq 0$
- Monotonicity:  $p_i(1) - p_i(0) \geq 0, \forall i$  (sign wlog).

# Continuous Treatments

Then the 2SLS/Wald estimand identifies

$$\frac{E[q_i \mid stormy_i = 1] - E[q_i \mid stormy_i = 0]}{E[p_i \mid stormy_i = 1] - E[p_i \mid stormy_i = 0]} = \frac{\int E[q'_i(t) \mid p_i(1) \geq t > p_i(0)] \cdot P[p_i(1) \geq t > p_i(0)] dt}{\int P[p_i(1) \geq t > p_i(0)] dt}$$

a.k.a. the continuous average causal response formula.

# Continuous Treatments

## Proof

Let us use  $Z_i = \text{stormy}_i$  to lighten the notation. Let us also assume that  $p_i$  has a lower bound, and (wlog) that this lower bound is 0. Then by the fundamental theorem of calculus,

$$q_i = Z_i q_i(p_i(1)) + (1 - Z_i) q_i(p_i(0)) = Z_i \left[ q_i(0) + \int_0^{p_i(1)} q'_i(t) dt \right] + (1 - Z_i) \left[ q_i(0) + \int_0^{p_i(0)} q'_i(t) dt \right].$$

By independence + monotonicity [implying  $p_i(1) > p_i(0)$ ],  $E[q_i \mid Z_i = 1] - E[q_i \mid Z_i = 0]$  equals

$$E[q_i(p_i(1)) - q_i(p_i(0))] = E \left[ q_i(0) + \int_0^{p_i(1)} q'_i(t) dt \right] - E \left[ q_i(0) + \int_0^{p_i(0)} q'_i(t) dt \right].$$

Under monotonicity,  $p_i(1) > p_i(0)$  hence [assuming all necessary regularity conditions allowing the use of Fubini theorem to interchange the integration and expectation operators]

$$\begin{aligned} E[q_i \mid Z_i = 1] - E[q_i \mid Z_i = 0] &= E \left[ \int_{p_i(0)}^{p_i(1)} q'_i(t) dt \right] = E \left[ \int_0^\infty q'_i(t) I(p_i(1) \geq t > p_i(0)) dt \right] \\ &= \int_0^\infty E[q'_i(t) \mid p_i(1) \geq t > p_i(0)] P[p_i(1) \geq t > p_i(0)] dt \end{aligned}$$

# Continuous Treatments

## *Proof*

Similarly, for the denominator,  $p_i = Z_i p_i(1) + (1 - Z_i) p_i(0)$ , and because  $p$  plays the role of  $q_i(p)$  in previous derivations for the numerator [and  $\frac{\partial}{\partial p} = 1$ ]

$$\begin{aligned} E[p_i | Z_i = 1] - E[p_i | Z_i = 0] &= E \left[ \int_{p_i(0)}^{p_i(1)} \underbrace{p'_i(t)}_{=1} dt \right] = E \left[ \int_0^\infty I(p_i(1) \geq t > p_i(0)) dt \right] \\ &= \int_0^\infty P[p_i(1) \geq t > p_i(0)] dt \end{aligned}$$

**QED.**

# Continuous Treatments

## *Building intuition*

Two special cases help building intuition about the continuous ACR identification result.

1. **Linear causal response function.**  $q_i(p) = \alpha_{0i} + \alpha_{1i}p$  for random coefficients  $(\alpha_{0i}, \alpha_{1i})$ .

$$\begin{aligned} \frac{E[q_i | Z_i = 1] - E[q_i | Z_i = 0]}{E[p_i | Z_i = 1] - E[p_i | Z_i = 0]} &= \frac{E[q_i(p_i(1)) - q_i(p_i(0))]}{E[p_i(1) - p_i(0)]} \text{ by independence} \\ &= \frac{E[\alpha_{1i} (p_i(1) - p_i(0))]}{E[p_i(1) - p_i(0)]}, \end{aligned}$$

a weighted average (across markets) of the effect of price on demand with weights proportional to the price change induced by the weather instrument in market  $i$ .

Did we use monotonicity at this point?

# Continuous Treatments

## Building intuition

Two special cases help building intuition about the continuous ACR identification result.

1. **Linear causal response function.**  $q_i(p) = \alpha_{0i} + \alpha_{1i}p$  for random coefficients  $(\alpha_{0i}, \alpha_{1i})$ .

$$\begin{aligned} \frac{E[q_i | Z_i = 1] - E[q_i | Z_i = 0]}{E[p_i | Z_i = 1] - E[p_i | Z_i = 0]} &= \frac{E[q_i(p_i(1)) - q_i(p_i(0))]}{E[p_i(1) - p_i(0)]} \text{ by independence} \\ &= \frac{E[\alpha_{1i} (p_i(1) - p_i(0))]}{E[p_i(1) - p_i(0)]}, \end{aligned}$$

a weighted average (across markets) of the effect of price on demand with weights proportional to the price change induced by the weather instrument in market  $i$ .

**Did we use monotonicity at this point?** No. But if we want to ensure that the weights

$\frac{p_i(1) - p_i(0)}{E[p_i(1) - p_i(0)]}$  are non-negative, we need the monotonicity assumption.

# Continuous Treatments

## Building intuition

2. **Non-linear but homogeneous demand function.**  $q_i(p) = Q(p) + \eta_i$  where  $Q(\cdot)$  is a non-stochastic function and  $\eta_i$  an additive random error [independent from everything].

In this case,  $q'_i(p) = Q'(p)$  every day or in every market, and

$$\frac{E[q_i | Z_i = 1] - E[q_i | Z_i = 0]}{E[p_i | Z_i = 1] - E[p_i | Z_i = 0]} = \frac{\int (Q'(t) + \overbrace{E[\eta_i | p_i(1) \geq t > p_i(0)]}^{=0}) P[p_i(1) \geq t > p_i(0)] dt}{\int P[p_i(1) \geq t > p_i(0)] dt}$$
$$= \int Q'(t) \omega(t) dt, \text{ where } \omega(t) \equiv \frac{P[p_i(1) \geq t > p_i(0)]}{\int P[p_i(1) \geq r > p_i(0)] dr}$$

a weighted average *along* the length of the causal response function [assumed here common to all markets], placing more weight on derivatives at prices where the

instrument shifts CDF of prices most sharply. [Indeed, if  $(p_i(1), p_i(0))$  are continuously

distributed w/o point mass,  $P[p_i(1) \geq t > p_i(0)] = F_{p_i(0)}(t) - F_{p_i(1)}(t) \geq 0$  under monotonicity.]

## Continuous Treatment *and* Instrument

Borusyak and Hull (2024, AEA P&P) gives an even more general version of LATE, with a continuous treatment  $X_i$  and a continuous instrument  $Z_i$ :

$$\frac{\text{Cov}(Z_i, Y_i)}{\text{Cov}(Z_i, X_i)} = E \left[ \int_x Y_i'(x) \omega_i(x) dx \right]$$

where  $\omega_i(x) \geq 0$  under the appropriate first-stage monotonicity condition.

The result follows from the weighted-average-effect interpretation of reduced form and first stage, and chain rule:  $\frac{d}{dz} Y_i(X_i(z)) = Y_i'(x) X_i'(z)$ .

## Multiple Instruments

Things can get become tricky here. Recall: 2SLS identifies a weighted average of one-at-a-time IV estimands... Key question then: **when are these weights convex?**

Imbens and Angrist (1994) show convexity under **generalized monotonicity**:

- For all  $z, z' \in \text{Supp}(Z_i)$ , either  $D_i(z) \geq D_i(z')$  almost-surely or  $D_i(z) \leq D_i(z')$  almost-surely
- Implies instruments are "nested" (e.g. immediate + waitlist offer IVs)

Mogstad, Torgovitsky, and Walters (2021) show convexity may (or may not!) fail under a weaker **"partial monotonicity"** condition

- If  $z' \geq z$  component-wise, then  $D_i(z') \geq D_i(z)$  almost-surely
- E.g. a "price" and "distance" instrument for college (non-nested)

# Multiple Instruments

## *The Angrist and Imbens (1994) result*

Recall  $Z_i$  can be a vector now. Identifying assumptions:

1. Independence:  $Y_i(0), Y_i(1), D_i(z) \perp\!\!\!\perp Z_i, \forall z \in \text{Supp}(Z_i)$
2. Relevance:  $P(z) = E[D_i | Z_i = z]$  is a non-trivial fn. of  $z$  [i.e.,  $P(z)$  varies with  $z$ , testable]
3. Monotonicity: for all  $(z, w) \in \text{Supp}(Z_i)^2$ , either  $D_i(z) \geq D_i(w), \forall i$  or  $D_i(z) \leq D_i(w), \forall i$

Under 1.-2., and assuming wlog 3. holds for  $D_i(z) \geq D_i(w), \forall i$ , we get identification of

$$\alpha_{z,w} = E[Y_i(1) - Y_i(0) | D_i(z) \geq D_i(w)] = \frac{E[Y_i | Z_i = z] - E[Y_i | Z_i = w]}{P(z) - P(w)}.$$

Now, how could we exploit our vector-valued instrument  $Z_i$ ?

# Multiple Instruments

## *The Angrist and Imbens (1994) result*

Recall  $Z_i$  can be a vector now. Identifying assumptions:

1. Independence:  $Y_i(0), Y_i(1), D_i(z) \perp\!\!\!\perp Z_i, \forall z \in \text{Supp}(Z_i)$
2. Relevance:  $P(z) = E[D_i | Z_i = z]$  is a non-trivial fn. of  $z$  [i.e.,  $P(z)$  varies with  $z$ , testable]
3. Monotonicity: for all  $(z, w) \in \text{Supp}(Z_i)^2$ , either  $D_i(z) \geq D_i(w), \forall i$  or  $D_i(z) \leq D_i(w), \forall i$

Under 1.-2., and assuming wlog 3. holds for  $D_i(z) \geq D_i(w), \forall i$ , we get identification of

$$\alpha_{z,w} = E[Y_i(1) - Y_i(0) | D_i(z) \geq D_i(w)] = \frac{E[Y_i | Z_i = z] - E[Y_i | Z_i = w]}{P(z) - P(w)}.$$

Now, **how could we exploit our vector-valued instrument  $Z_i$ ?** One way is to estimate the ratio

$\text{Cov}(Y, g(Z)) / \text{Cov}(D, g(Z))$  for some scalar fn  $g(Z) : \text{Supp}(Z_i) \rightarrow \mathbb{R}$ .

→ If  $Z$  is a scalar r.v., then  $g(z) = z$  leads to the standard IV estimand

→ If  $Z$  is a vector,  $g(z)$  is often an estimate of  $P(z)$ .

# Multiple Instruments

*The Angrist and Imbens (1994) result*

What if instead we saturate the first stage w/ indicators  $I(Z_i = z)$  for all values of  $z \in \text{Supp}(Z_i)$ ?

[E.g., judge IVs where you include all judge indicators in the first-stage.]

# Multiple Instruments

## *The Angrist and Imbens (1994) result*

What if instead we saturate the first stage w/ indicators  $I(Z_i = z)$  for all values of  $z \in \text{Supp}(Z_i)$ ?  
[E.g., judge IVs where you include all judge indicators in the first-stage.]

Then since the first stage is saturated, it will capture  $E[D_i | Z_i] = P(Z_i)$ , and the 2SLS estimand is  $\frac{\text{Cov}(Y_i, P(Z_i))}{\text{Var}(P(Z_i))}$ . Yet  $\text{Cov}(D_i, P_i(Z_i)) = \text{Cov}(D_i - P_i(Z_i), P_i(Z_i)) + \text{Var}(P(Z_i))$ , this is also equivalent to the IV estimand instrumenting by  $P(Z_i)$  directly!

**Careful about inference though!** There are some subtleties there, as you want to take into account inferential consequences of first-stage estimation involving many instruments.

→ reason behind “leave-one-out” (JIVE) estimation is precisely the threat of many-IV bias

→ with additional controls, the UJIVE procedure developed by Kolesar (2013) is state-of-the-art, cf. [R package jive](#).

# Multiple Instruments

## *The Angrist and Imbens (1994) result*

An additional assumption on  $g(Z)$  is necessary to ensure convex weights on LATEs:

- 4.i For all  $(z, w) \in \text{Supp}(Z_i)^2$ ,  $\Pr(z) \leq P(w) \Rightarrow g(z) \leq g(w)$  or  $P(z) \leq P(w) \Rightarrow g(z) \geq g(w)$   
4.ii  $\text{Cov}(D, g(Z)) \neq 0$ .

Notice this assumption is trivially satisfied (i) if  $Z$  is binary, (ii) if  $Z$  is a scalar r.v. and both  $g(z)$  and  $P(z)$  are monotone in  $z$ , or (iii)  $g(z) = E[D | Z = z] = P(z)$  [then monotonicity implies 4.i].

Assume 1. to 4. Let  $Z$  be a discrete r.v. with support  $\{z_0, z_1, \dots, z_K\}$ , ordered in such a way that if  $l < m$  then  $P(z_l) \leq P(z_m)$ . Then, the IV estimand identifies

$$\alpha_g^{IV} = \text{Cov}(Y, g(Z)) / \text{Cov}(D, g(Z)) = \sum_{k=1}^K \lambda_k \cdot \alpha_{z_k, z_{k-1}},$$

with weights  $\lambda_k = \frac{(P(z_k) - P(z_{k-1})) \cdot \sum_{l=k}^K \pi_l \cdot (g(z_l) - E[g(Z)])}{\sum_{m=1}^K (P(z_m) - P(z_{m-1})) \cdot \sum_{l=m}^K \pi_l \cdot (g(z_l) - E[g(Z)])}$ ,

where  $\pi_k = \Pr(Z = z_k)$  and  $\alpha_{z_k, z_{k-1}} = E[Y_i(1) - Y_i(0) | D_i(z_k) = 1, D_i(z_{k-1}) = 0]$ .

Weights  $\lambda_k$  are nonnegative and add up to one.

# Multiple Instruments

## *The Angrist and Imbens (1994) result*

**Proof:** First, assume wlog that for all  $(z, w) \in \text{Supp}(Z_i)^2$ ,  $\Pr(z) \leq P(w) \Rightarrow g(z) \leq g(w)$ .

Given that the points of support are ordered, this implies that if  $l < m$ , then  $g(z_l) \leq g(z_m)$  and

that  $\text{Cov}(D, g(Z)) = \underbrace{\text{Cov}(D - E[D | Z], g(Z))}_{=0} + \text{Cov}(E[D | Z], g(Z)) > 0$  as

$E[D | Z = z] = P(z)$  and  $g(z)$  are both increasing in  $z$ . covariance property

Second, given the  $K + 1$  points in  $\text{Supp}(Z)$ , we can define  $\binom{K+1}{2} = K \times (K + 1)/2$  LATEs  $\alpha_{z_l, z_m}$ , one for each (unordered) pair of points of support  $(z_l, z_m)$ . These  $K \times (K + 1)/2$  LATEs are related in the following way (using the definition of  $\alpha_{z_m, z_l}$ ):

$$\alpha_{z_m, z_k} = \frac{P(z_l) - P(z_k)}{P(z_m) - P(z_k)} \cdot \alpha_{z_l, z_k} + \frac{P(z_m) - P(z_l)}{P(z_m) - P(z_k)} \cdot \alpha_{z_m, z_l} \quad \text{for all } k \neq l, k \neq m, \text{ and } l \neq m.$$

Third, the conditional expectation of  $Y$  given  $Z = z_k$  for  $k \geq 1$  can be written as  $E[Y | Z = z_k]$

$$= E[Y | Z = z_0] + \alpha_{z_k, z_0} \cdot (P(z_k) - P(z_0)) = E[Y | Z = z_0] + \sum_{l=1}^k \alpha_{z_l, z_{l-1}} \cdot (P(z_l) - P(z_{l-1})).$$

# Multiple Instruments

## *The Angrist and Imbens (1994) result*

The IV procedure estimates  $\alpha_g^{IV} = \frac{\text{Cov}(Y, g(Z))}{\text{Cov}(D, g(Z))} = \frac{E[Y \cdot (g(Z) - E[g(Z)])]}{E[D \cdot (g(Z) - E[g(Z)])]}$

First we analyze the numerator of this expression:

$$\begin{aligned} E[Y \cdot (g(Z) - E[g(Z)])] &= \sum_{l=0}^K \pi_l \cdot E[Y | Z = z_l] \cdot (g(z_l) - E[g(Z)]) \quad \text{by LIE} \\ &= \overbrace{\sum_{l=0}^K \pi_l \cdot E[Y | Z = z_0] \cdot (g(z_l) - E[g(Z)])}^{=0} \\ &\quad + \sum_{l=1}^K \pi_l \sum_{k=1}^l \alpha_{z_k, z_{k-1}} \cdot (P(z_k) - P(z_{k-1})) \cdot (g(z_l) - E[g(Z)]) \\ &= \sum_{k=1}^K \alpha_{z_k, z_{k-1}} \cdot (P(z_k) - P(z_{k-1})) \sum_{l=k}^K \pi_l \cdot (g(z_l) - E[g(Z)]) \end{aligned}$$

# Multiple Instruments

## *The Angrist and Imbens (1994) result*

Similar calculations  $\alpha_g^{IV}$  denominator shows it equals the denominator of the weights  $(\lambda_k)_{k=1}^K$ . [It is sufficient to observe that when replacing  $Y$  by  $D$  in the previous derivations, the terms analogous to  $\alpha_{z_k, z_m}$  all equal 1 for any  $(z_k, z_m)$ .]

The weights  $(\lambda_k)_{k=1}^K$  clearly add up to 1 by construction. They are nonnegative because  $P(z_k) \geq P(z_{k-1})$  and  $g(z_k) \geq g(z_{k-1})$  which follows from the ordering of the points of support, and this in turn implies that  $\sum_{l=k}^K \pi_l \cdot (g(z_l) - E[g(Z)]) \geq 0$  for all  $k$ . Indeed, for any  $k \geq 1$ ,

$$\begin{aligned} \sum_{l=k}^K \pi_l \cdot (g(z_l) - E[g(Z)]) &= \sum_{l=k}^K \pi_l \cdot g(z_l) - E[g(Z)]P[Z_i \geq k] \\ &= E[g(Z) \mid Z_i \geq k]P[Z_i \geq k] - E[g(Z)]P[Z_i \geq k] \\ &\geq 0 \text{ as } E[g(Z) \mid Z_i \geq k] \geq E[g(Z)] \end{aligned}$$

where we used  $\sum_{l=k}^K \pi_l = \sum_{l=k}^K P[Z_i = l] = P[Z_i \geq k]$ .

# Multiple Instruments

*The Mogstad, Torgovitsky and Walters (2021) critique*

MTW21 make two key observations.

## 1. Generalized monotonicity may be highly/too restrictive on choice behavior.

Suppose we have two binary instruments,  $Z_1$  and  $Z_2$ . For instance,

- $Z_1 = 1$  if you are offered a grant to attend college, 0 otherwise
- $Z_2 = 1$  if you live close to a university campus, 0 otherwise.

Let us assume  $(Z_1, Z_2)$  are valid to instrument college education  $D_i \in \{0, 1\}$  in a regression on later-life earnings  $Y_i$ .

Generalized monotonicity in IA94: for any  $(z_1, z_2), (z'_1, z'_2) \in \{0, 1\}^2 \times \{0, 1\}^2$ , we have either  $D_i(z_1, z_2) \geq D_i(z'_1, z'_2), \forall i$  or  $D_i(z_1, z_2) \leq D_i(z'_1, z'_2), \forall i$ .

→ Probably credible that all  $i$  have  $D_i(1, 1)$  greater than  $D_i(0, 0), D_i(1, 0), D_i(0, 1)$ .

→ Similarly, credible that all  $i$  have  $D_i(1, 0)$  or  $D_i(0, 1)$  greater than  $D(0, 0)$ .

# Multiple Instruments

*The Mogstad, Torgovitsky and Walters (2021) critique*

But what about, e.g.,  $D_i(1, 0)$  vs.  $D_i(0, 1)$ ?

# Multiple Instruments

*The Mogstad, Torgovitsky and Walters (2021) critique*

But what about, e.g.,  $D_i(1, 0)$  vs.  $D_i(0, 1)$ ? Generalized monotonicity assumes there can't be an individual whose decision to attend college is affected by the grant but not by distance while the reverse is true for another individual.

→ Essentially, this imposes some homogeneity in preferences.

Bottom line: when we think there are values of the instruments leading to “non-nested” choices of individuals [i.e, the set of instrument values leading individuals  $i$  and  $j$  to take up treatment are not necessarily included/“nested” one into the other] ...

... then generalized monotonicity fails [Cf. proposition 1 and 2 of MTW21] .

# Multiple Instruments

*The Mogstad, Torgovitsky and Walters (2021) critique*

2. What does 2SLS identifies under “partial monotonicity”? MTW21 show 2SLS may or may not be a convex weighted avr. of LATEs under a weaker “partial monotonicity”.

Partial monotonicity: if  $z' \geq z$  component-wise, then  $D_i(z') \geq D_i(z)$  almost-surely.

Then 2SLS saturating the first stage in the instruments may not identify a convex-weighted average of LATEs

- A negative weight could be placed on compliers that respond to one of the instrument only
- No problem if  $Cov(Z_1, Z_2) = 0$  or  $Cov(Z_1, Z_2) > 0$  though
- Whether or not weights are convex can be tested  $\rightarrow$  cf. `mivcausal` Stata module [for binary treatment and two binary instruments]
- Instruments can always be used separately, as long as the analysis is conditional on the other instrument [otherwise independence is likely to fail]

# Multiple Treatments

Things can get tricky when there are multiple treatments involved.

1. **Many treatments and one instrument.** [Behaghel et al. (2013), Kline and Walters (2016)].

Suppose you have a treatment indicator  $D_i \in \{n, c, h\}$  and a binary instrument  $Z_i$  increasing the likelihood to get into treatment  $h$ .

I.e.,  $Z_i$  is a valid instrument (satisfying in particular monotonicity) for treatment recoded as  $I(D_i = h)$ . So we can identify

$$LATE_h = E[Y_i(h) - Y_i(D_i(0)) | D_i(1) = h, D_i(0) \neq h]$$

But we don't know if the counterfactual is  $n$  or  $c$ !

# Multiple Treatments

Indeed,  $LATE_h$  equals

$$\begin{aligned} & E[Y_i(h) - Y_i(D_i(0)) | D_i(1) = h, D_i(0) \neq h] \\ = & E[Y_i(h) - Y_i(D_i(0)) | D_i(1) = h, D_i(0) \neq h, D_i(0) = c] P(D_i(0) = c | D_i(1) = h, D_i(0) \neq h) + \\ & E[Y_i(h) - Y_i(D_i(0)) | D_i(1) = h, D_i(0) \neq h, D_i(0) = n] P(D_i(0) = n | D_i(1) = h, D_i(0) \neq h) \\ = & E[Y_i(h) - Y_i(D_i(0)) | D_i(1) = h, D_i(0) = c] P(D_i(0) = c | D_i(1) = h, D_i(0) \neq h) + \\ & E[Y_i(h) - Y_i(D_i(0)) | D_i(1) = h, D_i(0) = n] P(D_i(0) = n | D_i(1) = h, D_i(0) \neq h) \\ = & LATE_{ch} P(D_i(0) = c | D_i(1) = h, D_i(0) \neq h) + \\ & LATE_{nh} P(D_i(0) = n | D_i(1) = h, D_i(0) \neq h) \end{aligned}$$

and using Bayes theorem,  $P(D_i(0) = c | D_i(1) = h, D_i(0) \neq h)$  equals

$$\frac{P(D_i(0) = c, D_i(1) = h, D_i(0) \neq h)}{P(D_i(1) = h, D_i(0) \neq h)} = \frac{P(D_i(0) = c, D_i(1) = h)}{P(D_i(1) = h, D_i(0) \neq h)} \equiv S_c$$

so all in all,  $LATE_h = S_c LATE_{ch} + (1 - S_c) LATE_{nh}$

# Multiple Treatments

**Is this an issue?** Well, it depends. Here, think of  $n$  as absence of treatment, and  $c$  as an alternative treatment. The quantity identified can only recover a weighted average of

- a (local) effect of going from no treatment to the treatment encouraged by the instrument
- a (local) effect of going from the alternative treatment to the one encouraged by the instrument.

In Kline and Walters (2016), this quantity turns out to be policy relevant.

But there may be cases where mixing the no-treatment/encouraged-treatment contrast with the alternative-treatment/encouraged-treatment contrast may be problematic.

**What if we have several instrument?**

# Multiple Treatments

**What if we have several instrument?** There are still challenges in this case. Behaghel et al (2013) and Bhuller and Sigstad (2024) show that in “just-identified” models [as many binary instruments as binary treatments] a sufficient and necessary condition for the TSLS estimand to identify proper LATEs for each treatment is that **each instrument affects only one treatment**.

**What compliance patterns do we rule out?** With  $D \in \{0, 1, 2\}$  and the associated instrument  $Z \in \{0, 1, 2\}$ , you can't have  $I(D(2) = 1) > I(D(2) = 0)$  [i.e., some individuals being induced to take up treatment 1 when encouraged to take treatment 2] .

Alternatives: MTE-based solutions (Mountjoy, 2022, Humphries et al, 2025). [Sadly, I doubt I will have time to cover them in this iteration of the class]

# Necessary Controls

There are settings where we may only be willing to assume that the instrument  $Z_i$  is as-good-as-random only *conditional* on controls  $W$ . Formally, this can be translated into assumptions of the form

$$\{Y_i(d)\}_{d \in \text{Supp}(D_i)}, \{D_i(z)\}_{z \in \text{Supp}(Z_i)} \perp\!\!\!\perp Z_i \mid W_i$$

or the weaker mean-independence condition

$$E[Z_i \mid W_i, \{Y_i(d)\}_{d \in \text{Supp}(D_i)}, \{D_i(z)\}_{z \in \text{Supp}(Z_i)}] = E[Z_i \mid W_i].$$

[Examples: encouragement-based RCTs with strata-specific propensity score, quasi-experiments where implicit propensity score is varying, judge designs controlling for courtrooms etc.]

**Can we still interpret TSLS with controls causally?** There are a few (recent and less recent) results on this topic.

# Necessary Controls

Can we still interpret TSLS with controls causally? There are a few (recent and less recent) results on this topic.

1. Capturing the “propensity score”  $E[Z_i | W_i]$ . In general, Blandhol et al. (2024) the key necessary and sufficient condition to ensure that TSLS captures a convex weighted average of conditional LATEs is that the first-stage controls sufficiently flexibly for  $W_i$  that it captures  $E[Z_i | W_i]$ , i.e., it implicitly “demean”  $Z_i$  by  $E[Z_i | W_i]$ .

[E.g., if one is willing to assume that  $E[Z_i | W_i]$  is linear in  $W_i$ , then controlling linearly for  $W_i$  in the first and second stage will do the job.]

Can you think of settings in which  $E[Z_i | W_i]$  will be linear by construction?

# Necessary Controls

2. **Discrete controls.** Angrist (1998) showed that if the necessary control  $W_i$  is discrete, in which case  $E[Z_i | W_i]$  is linear in  $W_i$  by construction, saturating in  $W_i$  only the first and second stages is sufficient to recover a convex-weighted average of conditional LATEs.

What's more, if  $E[Z_i | W_i]$  is actually constant, then the estimand identified by the regression described above will be the unconditional LATE.

## Necessary Controls

**Proof.** By FWL theorem, the first-stage reg. of  $D_i$  on  $Z_i$  controlling for  $W_i$  is equivalent to a reg. of  $D_i$  on the residual of a reg. of  $Z_i$  on  $W_i$ . Since  $E[Z_i | W_i]$  is assumed to be linear, this last reg. generates a residual  $Z_i - E[Z_i | W_i]$ . Then the TSLS estimand is given by

$$\beta^{TSLS} = \frac{\text{Cov}(Y_i, Z_i - E[Z_i | W_i])}{\text{Cov}(D_i, Z_i - E[Z_i | W_i])}.$$

[We can ignore the residualization of  $Y_i$  and  $D_i$  on  $W_i$  as the residual  $Z_i - E[Z_i | W_i]$  is mean-independent of any function of  $W_i$ .]

$$\begin{aligned} \frac{\text{Cov}(Y_i, Z_i - P(W_i))}{\text{Var}(Z_i - P(W_i))} &= \frac{E[\text{Cov}(Y_i, Z_i | W_i)]}{E[\text{Var}(Z_i | W_i)]} \quad \text{and denoting } \text{Var}(Z_i | W_i) \equiv \sigma_Z(W_i) \\ &= E \left[ \frac{\sigma_Z^2(W_i)}{E[\sigma_Z^2(W_i)]} (E[Y_i | Z_i = 1, W_i] - E[Y_i | Z_i = 0, W_i]) \right] \end{aligned}$$

$$\text{since } \frac{\text{Cov}(Y_i, Z_i | W_i)}{\text{Var}(Z_i | W_i)} = E[Y_i | Z_i = 1, W_i] - E[Y_i | Z_i = 0, W_i].$$

## Necessary Controls

A similar result is obtained for the first-stage, replacing  $Y_i$  by  $D_i$  in the derivations. Thus,

$$\begin{aligned}\beta^{TSLS} &= \frac{E[\sigma_Z^2(W_i)(E[Y_i | Z_i = 1, W_i] - E[Y_i | Z_i = 0, W_i])]}{E[\sigma_Z^2(W_i)(E[D_i | Z_i = 1, W_i] - E[D_i | Z_i = 0, W_i])]} \\ &= E\left[\frac{\sigma_Z^2(W_i)\pi(W_i)}{E[\sigma_Z^2(W_i)\pi(W_i)]}\frac{E[Y_i | Z_i = 1, W_i] - E[Y_i | Z_i = 0, W_i]}{E[D_i | Z_i = 1, W_i] - E[D_i | Z_i = 0, W_i]}\right] \\ &= E\left[\frac{\sigma_Z^2(W_i)\pi(W_i)}{E[\sigma_Z^2(W_i)\pi(W_i)]}LATE(W_i)\right],\end{aligned}$$

where  $\pi(W_i) \equiv E[D_i | Z_i = 1, W_i] - E[D_i | Z_i = 0, W_i]$ , and

$LATE(W_i) \equiv E[Y_i(1) - Y_i(0) | D_i(1) > D_i(0), W_i]$ . The last line follows provided that the LATE identifying assumptions, and in particular the monotonicity assumption, holds for all values of  $W_i$  (and in the same direction no matter the value of  $W_i$ ).

# Necessary Controls

3. Specifications fully saturating in  $W_i$ . Then identification of an estimand that weights disproportionately cells with higher share of compliers... (saturate and weight result).  
But more robust to monotonicity violations, as it can tolerate conditional monotonicity  
cf. Sloczynski 2024

## Necessary Controls

**Proof.** Theorem 3. Let  $g[W]$  be a design matrix constructed from indicator variables for each value of  $W$ . Consider the TSLS estimate computed using  $g(W)$  and a full set of interactions between  $g(W)$  and  $Z$  as instruments for a regression of  $Y$  on rows of  $g[W]$  and  $S$ . The resulting estimate is

$$\begin{aligned}\beta^{TSLS} &= \frac{E\{Y \cdot (E[D | W, Z] - E[D | W])\}}{E\{D \cdot (E[D | W, Z] - E[D | W])\}} \\ &= \frac{E\{\beta(W)\Theta(W)\}}{E\{\Theta(W)\}}.\end{aligned}$$

where  $\Theta(W) = E\{E[D | W, Z] \cdot (E[D | W, Z] - E[D | W]) | W\} = \text{Var}(E[D | Z, W] | W)$   
and  $\beta(X) = \frac{E\{Y \cdot (E[D | W, Z] - E[D | W]) | W\}}{E\{D \cdot (E[D | W, Z] - E[D | W]) | W\}}$ .

For the first equality, we used the fact that the TSLS regression described above is equivalent to a regression of  $Y$  on  $E[D | Z, W]$  controlling for  $g[W]$ , as the saturated first-stage captures by construction  $E[D | Z, W]$ . By FWL, this regression is itself equivalent to a reg. of  $Y$  on the residual of a reg. of  $E[D | Z, W]$  on  $g[W]$ . This residual of a reg. saturated in  $W$  will be  $E[D | Z, W] - E[E[D | Z, W] | W] = E[D | Z, W] - E[D | W]$ .

## Necessary Controls

Note that the coef. of a reg. of  $Y$  on  $E[D | Z, W] - E[D | W]$  is given by

$$\frac{\text{Cov}\{Y, E[D | W, Z] - E[D | W]\}}{\text{Var}\{E[D | W, Z] - E[D | W]\}} = \frac{E\{Y \cdot (E[D | W, Z] - E[D | W])\}}{\text{Var}\{(E[D | W, Z] - E[D | W])\}}$$

Yet notice that  $\text{Var}\{E[D | W, Z] - E[D | W]\} = E\{(E[D | W, Z] - E[D | W])^2\}$

since  $E\{(E[D | W, Z] - E[D | W])\} = 0$ . Indeed,  $E[D | W, Z] - E[D | W]$  is a regression residual of a saturated regression, so it has to be mean 0. It can be checked by direct computation using the LIE twice:  $E[D | W, Z] - E[D | W] = E\{E[D | W, Z] - E[D | W] | W\} = E\{E[D | W, Z] | W\} - E[D | W] = E[D | W] - E[D | W] = 0$ .

Now, we only need to observe that the numerator of  $\beta^{TSLs}$  actually equals

$$\begin{aligned} E\{D \cdot (E[D | W, Z] - E[D | W])\} &= E[E\{D \cdot (E[D | W, Z] - E[D | W]) | W, Z\}] \\ &= E\{E[D | W, Z] \cdot (E[D | W, Z] - E[D | W])\} \\ &= E\{(E[D | W, Z] - E[D | W])^2\} = \text{V}\{E[D | W, Z] - E[D | W]\} \end{aligned}$$

where in the third equality we used  $E\{E[D | W] \cdot (E[D | W, Z] - E[D | W])\} = 0$ .

## Necessary Controls

$$\begin{aligned}\beta^{TSLS} &= \frac{E\{Y \cdot (E[D | W, Z] - E[D | W])\}}{E\{D \cdot (E[D | W, Z] - E[D | W])\}} \\ &= \frac{E\{\beta(W)\Theta(W)\}}{E[\Theta(W)]}.\end{aligned}$$

where  $\Theta(W) = E\{E[D | W, Z] \cdot (E[D | W, Z] - E[D | W]) | W\} = \text{Var}(E[D | Z, W] | W)$

and  $\beta(W) = \frac{E\{Y \cdot (E[D | W, Z] - E[D | W]) | W\}}{E\{D \cdot (E[D | W, Z] - E[D | W]) | W\}}.$

For the second line, LIE and multiplying-dividing by  $E\{D \cdot (E[D | W, Z] - E[D | W]) | W\},$

$E\{Y \cdot (E[D | W, Z] - E[D | W])\} =$

$$E\left\{E\{D \cdot (E[D | W, Z] - E[D | W]) | W\} \frac{E\{Y \cdot (E[D | W, Z] - E[D | W]) | W\}}{E\{D \cdot (E[D | W, Z] - E[D | W]) | W\}}\right\}.$$

By the LIE again,

$E\{D \cdot (E[D | W, Z] - E[D | W]) | W\} = E\{E[D | W, Z] \cdot (E[D | W, Z] - E[D | W]) | W\} = \Theta(W),$

and therefore using the LIE at the denominator and using the above observation in both the numerator and the denominator, we get the result.

## Necessary Controls

This first part of the result was proven in Angrist and Imbens (1995). Sloczynski (2024) further notices that under exclusion restriction, independence of the instrument, and a weaker form of monotonicity, this  $\beta^{TSLS}$  estimand identifies a convex weighted average of LATEs.

**Assumption WM (Weak monotonicity).** There exists a subset of the support of  $X$  such that  $P[D(1) \geq D(0) | X] = 1$  on it and  $P[D(1) \leq D(0) | X] = 1$  on its complement.

Under this set of assumptions,  $\beta(w)$ , which is the TSLS using  $Z$  as instrument in the subsample of obs. with  $W = w$ , identifies the treatment effect among compliers (or defiers, depending on the way  $Z$  affects  $D$  for individuals with  $W = w$ ). Meanwhile, noticing that  $E[D | Z, W]$  equals  $E[D | Z = 0, W] + Z \cdot \{E[D | Z = 1, W] - E[D | Z = 0, W]\}$ , this gives  $\Theta(W)$  equal to  $V(E[D | Z, W] | W) = (E[D | Z = 1, W] - E[D | Z = 0, W])^2 V(Z | W) = \pi^2(W)V(Z | W)$  where  $\pi(w)$  denotes the share of compliers/defiers in cell  $W = w$ . Hence

$$\beta^{TSLS} = \frac{E[\pi^2(W)V(Z | W)LATE(W)]}{E[\pi^2(W)V(Z | W)]}$$

## Regression coefficient on a constant and binary variable

Recall that in a univariate regression of  $Y$  on  $(1, Z)$ , the coef. on  $Z$  is given by  $\frac{Cov(Y, Z)}{Var(Z)}$  [while the coef. on the constant is given by  $E(Y) - \frac{Cov(Y, Z)}{Var(Z)}E(Z)$ ].

Then one can observe that

$$\begin{aligned}Cov(Y, Z) &= E(YZ) - E(Y)E(Z) \\&= E(Y | Z = 1)P(Z = 1) \\&\quad - [E(Y | Z = 1)P(Z = 1) + E(Y | Z = 0)(1 - P(Z = 1))] P(Z = 1) \\&= \underbrace{P(Z = 1)(1 - P(Z = 1))}_{=Var(Z)} [E(Y | Z = 1) - E(Y | Z = 0)]\end{aligned}$$

## Avoid Manual 2SLS

Although easy, you should never literally run 2SLS in two stages. Cf. Mostly Harmless Econometrics (MHE) 4.6 for details, but here is summary/refresher:

1. Point estimates will be right, but s.e. generally won't be
2. Risk of omitting some controls from second stage in first stage (or vice versa)
3. Risk of “Forbidden regressions”: e.g. regressing  $Y_i$  on probit/logit fits for  $X_i$
4. Risk of regressing on  $\hat{X}_i$  and  $\hat{X}_i^2$ , instead of instrumenting  $X_i^2$  directly [e.g., by adding  $Z_i^2$  as extra instrument]

back

## Proof: LATE Identification Without Monotonicity

$$\begin{aligned} & E [(Y_i(1) - Y_i(0)) (D_i(1) - D_i(0))] \\ &= E [(Y_i(1) - Y_i(0)) | D_i(1) > D_i(0)] \cdot P[D_i(1) > D_i(0)] \\ &\quad - E [(Y_i(1) - Y_i(0)) | D_i(1) < D_i(0)] \cdot P[D_i(1) < D_i(0)] \end{aligned}$$

and  $E [(D_i(1) - D_i(0))] = 1 \cdot P(D_i(1) > D_i(0)) - 1 \cdot P(D_i(1) < D_i(0))$ .

Using the notation  $P(D_i(1) > D_i(0)) \equiv c$  and  $P(D_i(1) < D_i(0)) \equiv d$  and taking the ratio of these two quantities yields the result.

[back](#)

# Covariance Property

**Claim:** Let  $X$  be a random variable and  $f, g : \mathbb{R} \rightarrow \mathbb{R}$  be increasing functions. Then  $\text{Cov}(f(X), g(X)) \geq 0$ .

**Proof:** Since  $f$  and  $g$  are increasing, then  $(f(x) - f(y))(g(x) - g(y)) \geq 0$  for all  $x, y \in \mathbb{R}$ .

Assume  $X, Y$  are independent and identically distributed. By the above observation and the monotonicity of expectations, we get  $E[(f(X) - f(Y))(g(X) - g(Y))] \geq 0$ .

Expanding, we get  $E[f(X)g(X)] - E[f(X)g(Y)] - E[f(Y)g(X)] + E[f(Y)g(Y)] \geq 0$ .

Now due to independence, the LHS becomes

$$E[f(X)g(X)] - E[f(X)]E[g(Y)] - E[f(Y)]E[g(X)] + E[f(Y)g(Y)]$$

Then due to identically distributed, the LHS further becomes  $2E[f(X)g(X)] - 2E[f(X)]E[g(X)]$ .

So together we have  $2 \text{Cov}(f(X), g(X)) \geq 0$ .

**QED.**

back