# In-person session 6

**February 19, 2026**

PMAP 8521: Program evaluation
Andrew Young School of Policy Studies

# Plan for today

**Validity and *p*-hacking**

**DAGs**

***p*-values and confidence intervals**

# Validity and
# *p*-hacking

# What's going on with the assignments?

**andhs.co/nullworlds**

# What is *p*-hacking?

**andhs.co/hack**

# Are the results from p-hacking actually a threat to validity?

# Is a little exploratory p-hacking okay?

# Do people actually post their preregistrations?

# Yes!

## OSF

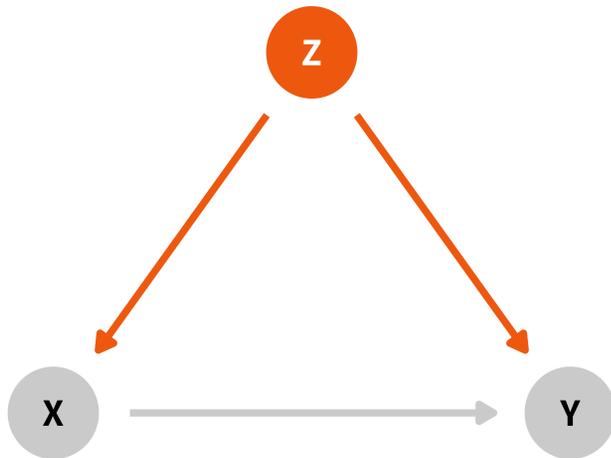See **this** and **this** for examples

## As Predicted

See **this**

# DAGs

# What about cycles?

Example time-based DAG

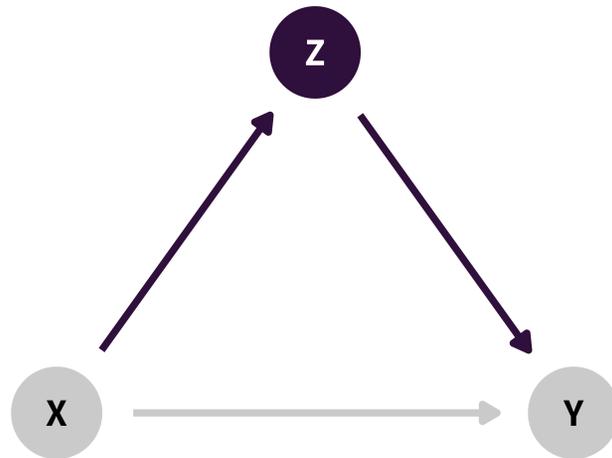# How do I know which of these is which?
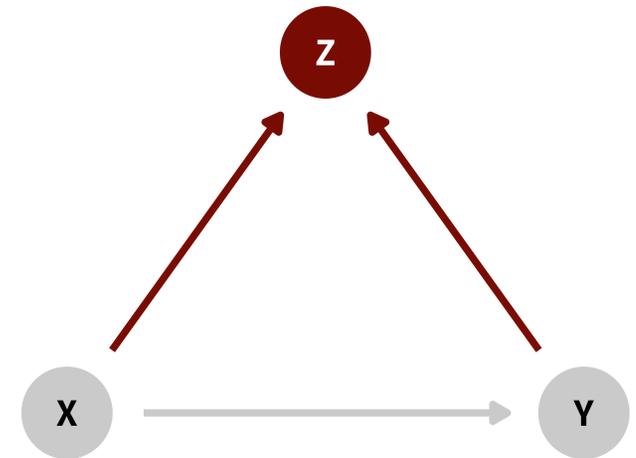
**Confounder**
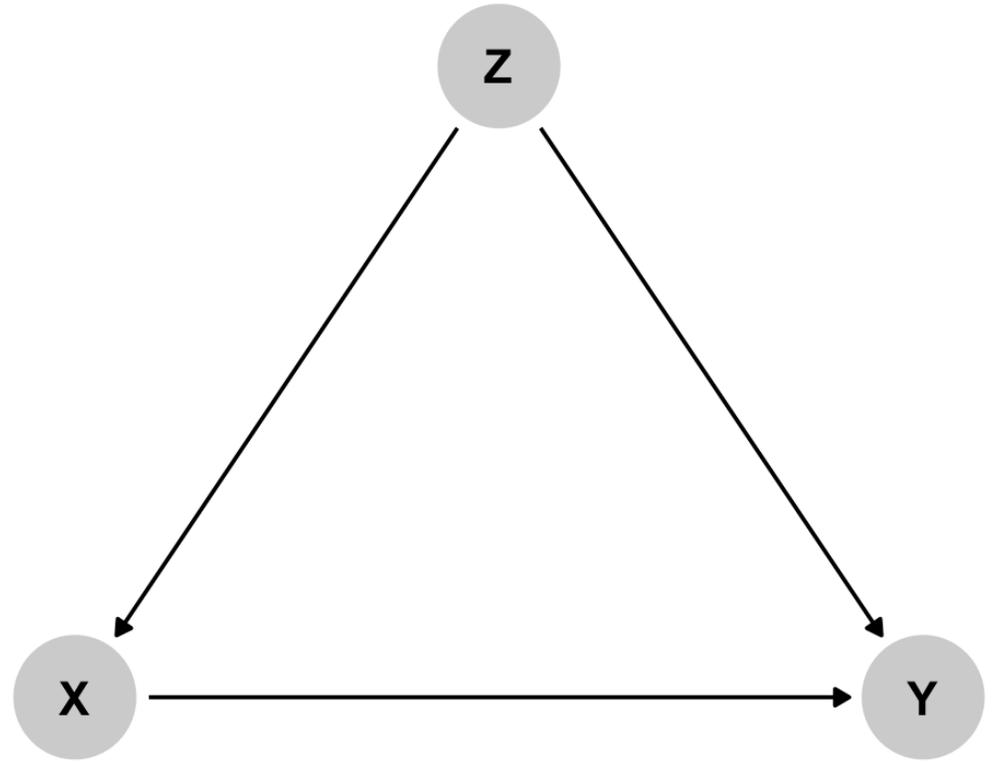(Fork)

**Mediator**
(Chain)

**Collider**
(Inverted fork)

**andhs.co/dags**

# d-separation

Except for the one arrow between X and Y,
no statistical association can flow between X and Y

This is **identification**—
all alternative stories are ruled out
and the relationship is isolated

# What's the difference between logic models and DAGs?

**Can't I just remake my logic model in Dagitty and be done?**

# DAGs vs. logic models

## DAGs are a *statistical* tool

Describe a data-generating process
and isolate/identify relationships

## Logic models are a *managerial* tool

Oversee the inner workings of a program and its theory

# DAGs and identification

DAGs are a statistical tool, but they don't tell you what statistical method to use

DAGs help you with the **identification strategy**

# Easiest identification

**Identification through research design**

**RCTs**

**When treatment is randomized, delete all arrows going into it**

# Most other identification

## Identification through do-calculus

**Rules for treating causal interventions as observations**

**Backdoor adjustment and frontdoor adjustment are special common patterns of do-calculus**

# Where can we learn more about *do*-calculus?

## Here!

# Adjusting for backdoor confounding

Causal effect
of $x$ on $y$

Conditional
mean of $y$,
given $x$ and $z$...

... weighted
by $z$

$$\mathbf{E}(y \mid \mathrm{do}(x)) = \sum_{z} \mathbf{E}(y \mid x, z) \times \mathbf{P}(z)$$

Sum across
all values of $z$

When things are identified, there are still arrows leading into Y.
What do we do with those?
How do you explain those relationships?

Outcomes have multiple causes.
How do you justify that your proposed cause is the most causal factor?

# How exactly do we close backdoors?

# What is front door adjustment??

# How exactly do colliders mess up your results?

## It looks like you can still get the effect of X on Y

# Facebook sent flawed data to misinformation researchers.



Mark Zuckerberg, chief executive of Facebook, testifying in Washington in 2018. Tom Brenner/The New York Times

# Does niceness improve appearance?

# Collider distorts the true effect!

# Colliders



It's ME hi I'm the collider it's ME

Lucy D'Agostino
McGowan
Wake Forest
University

**Taylor Swift fandom**

**Eras tour attendance**

**Income**

No relationship between income and Taylor Swift fandom

Attended Era's tour

Attended Era's tour

Taylor Swift fandom

Income

Eras tour attendance

# What we want:

$$\widehat{\text{fandom}} = \widehat{\beta}_0 + \widehat{\beta}_1 \text{income}$$

## What we want:

$$\widehat{\text{fandom}} = \widehat{\beta}_0 + \widehat{\beta}_1 \text{income}$$

## What we have:

$$\widehat{\text{fandom}} = \widehat{\beta}_0^* + \widehat{\beta}_1^* \text{income} + \widehat{\beta}_2^* \text{eras tour}$$

$$\widehat{\text{fandom}} = \widehat{\beta_0^*} + \widehat{\beta_1^*}\text{income} + \widehat{\beta_2^*}\text{eras tour}$$

$$\widehat{\mathrm{fandom}} = \widehat{\beta_0^*} + \widehat{\beta_1^*}\mathrm{income} + \widehat{\beta_2^*}\mathrm{eras\ tour}$$

$$\widehat{\text{fandom}} = \widehat{\beta_0^*} + \widehat{\beta_1^*}\text{income} + \widehat{\beta_2^*}\text{eras tour}$$

$$\widehat{\text{fandom}} = \widehat{\beta_0^*} + \widehat{\beta_1^*}\text{income} + \widehat{\beta_2^*}\text{eras tour}$$

# Effect of race on police use of force using administrative data

# Effect of race on police use of force using administrative data

## Administrative Records Mask Racially Biased Policing

DEAN KNOX   *Princeton University*
WILL LOWE   *Hertie School of Governance*
JONATHAN MUMMOLO   *Princeton University*

**R**esearchers often lack the necessary data to credibly estimate racial discrimination in policing. In particular, police administrative records lack information on civilians police observe but do not investigate. In this article, we show that if police racially discriminate when choosing whom to investigate, analyses using administrative records to estimate racial discrimination in police behavior are statistically biased, and many quantities of interest are unidentified—even among investigated individuals—absent strong and untestable assumptions. Using principal stratification in a causal mediation framework, we derive the exact form of the statistical bias that results from traditional estimation. We develop a bias-correction procedure and nonparametric sharp bounds for race effects, replicate published findings, and show the traditional estimator can severely underestimate levels of racially biased policing or mask discrimination entirely. We conclude by outlining a general and feasible design for future studies that is robust to this inferential snare.

**C**oncern over racial bias in policing, and the public availability of large administrative data sets documenting police–civilian interactions, have prompted a raft of studies attempting to quantify the effect of civilian race on law enforcement behavior. These studies consider a range of outcomes including ticketing, stop duration, searches, and the use of force (e.g., Antonovics and Knight 2009; Fryer 2019; Ridgeway 2006; Nix et al. 2017). Most research in this area attempts to adjust for omitted variables that may correlate with suspect race and the outcome of interest. In contrast, this study addresses a more fundamental problem that remains even if the vexing issue of omitted variable bias is solved: the inevitable statistical bias that results from studying racial discrimination using records that are themselves the product of racial discrimination (Angrist and Pischke 2008; Elwert and Winship 2014; Rosenbaum 1984). We show that when there is any biased absent additional data and/or strong and untestable assumptions.

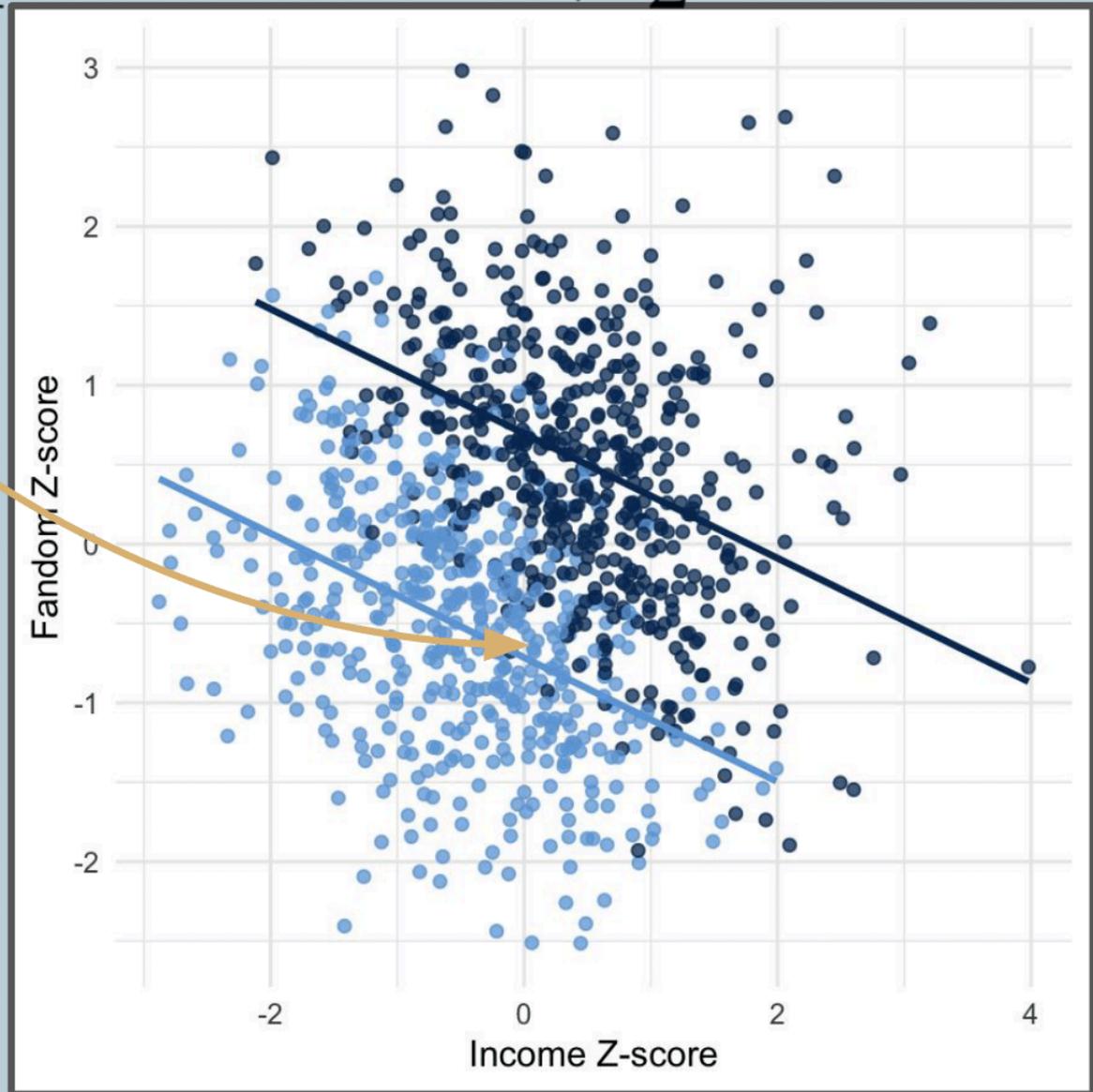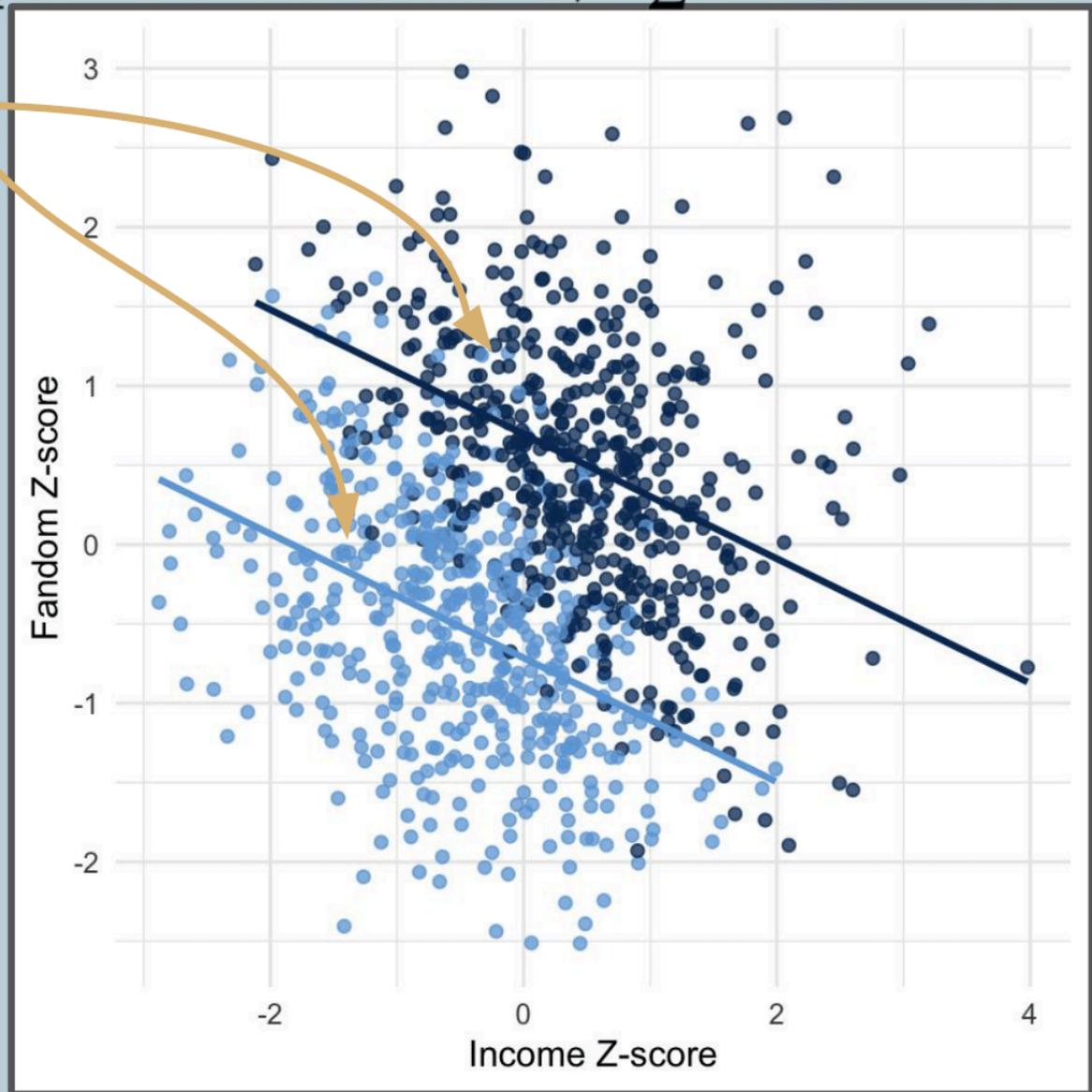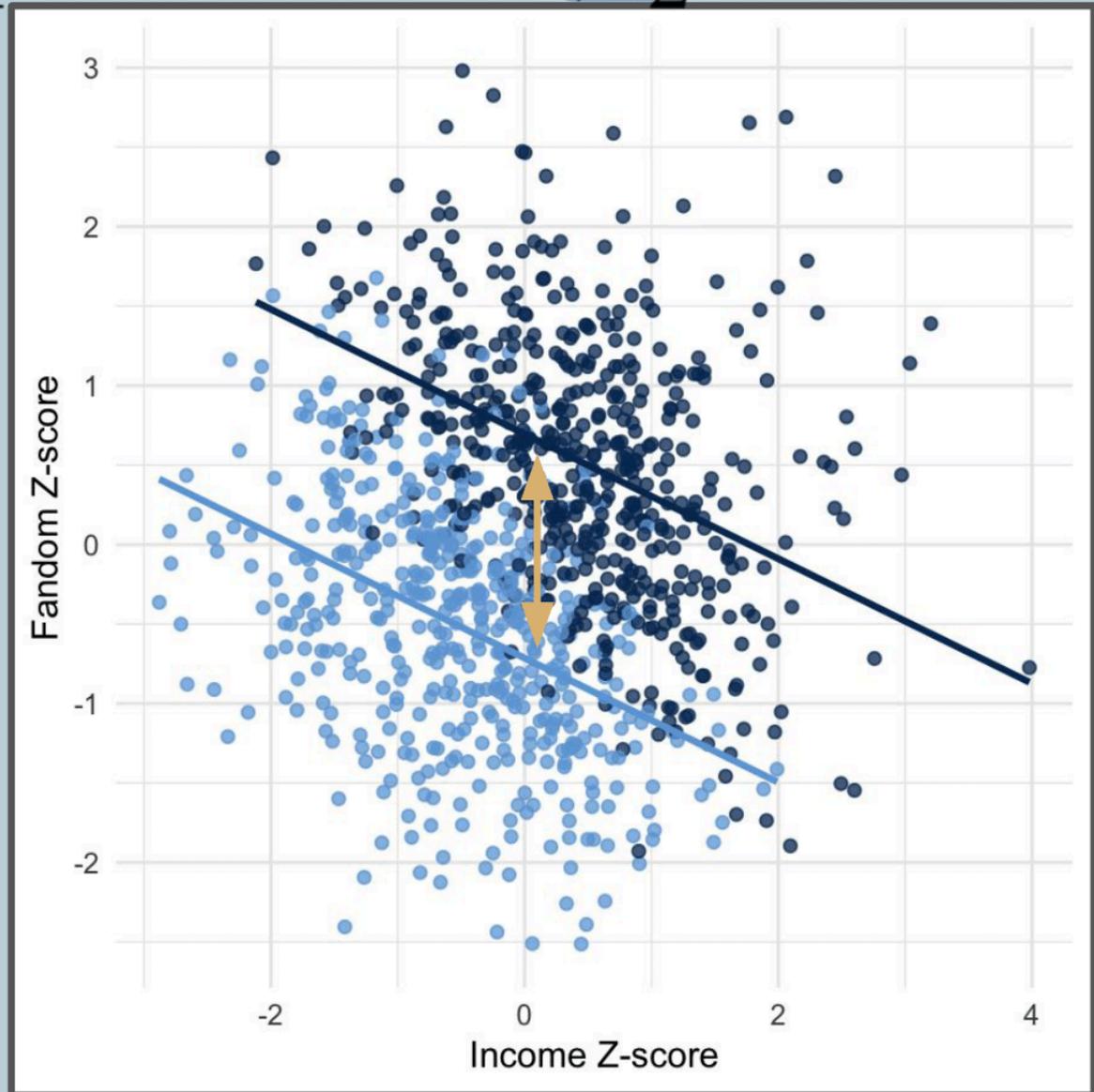This study makes several contributions. We clarify the causal estimands of interest in the study of racially discriminatory policing—quantities that many studies appear to be targeting, but are rarely made explicit—and show that the conventional approach fails to recover any known causal quantity in reasonable settings. Next, we highlight implicit and highly implausible assumptions in prior work and derive the statistical bias when they are violated. We proceed to develop informative nonparametric sharp bounds for the range of possible race effects, apply these in a reanalysis and extension of a prominent article on police use of force (Fryer 2019), and present bias-corrected results that suggest this and similar studies drastically underestimate the level of racial bias in police–civilian interactions. Finally, we outline strategies for future data collection and re-

# Can we make another DAG together?

# Potential outcomes vs. do() notation

# Expectations

$$\mathrm{E}(\cdot), \mathbf{E}(\cdot), \mathbb{E}(\cdot) \quad \text{vs.} \quad \mathrm{P}(\cdot)$$

**Basically a fancy way of saying "average"**

**andhs.co/potential-outcomes**

# Causal effects with potential outcomes

**Potential outcomes notation:**

$$\delta = \frac{1}{n} \sum_{i=1}^{n} Y_i(1) - Y_i(0)$$

or alternatively with $\mathbf{E}$

$$\delta = \mathbf{E}[Y_i(1) - Y_i(0)]$$

# Causal effects with do()

**Pearl notation:**

$$\delta = \mathbf{E}[Y_i \mid \text{do}(X = 1) - Y_i \mid \text{do}(X = 0)]$$

or more simply

$$\delta = \mathbf{E}[Y_i \mid \text{do}(X)]$$

$$\mathbf{E}[Y_i \mid \mathrm{do}(X)]$$

$$=$$

$$\mathbf{E}[Y_i(1) - Y_i(0)]$$

## We can't see this

$$\mathbf{E}[Y_i \mid \mathrm{do}(X)] \quad \text{or} \quad \mathbf{E}[Y_i(1) - Y_i(0)]$$

## So we find the average causal effect (ACE)

$$\hat{\delta} = \mathbf{E}[Y_i \mid X = 1] - \mathbf{E}[Y_i \mid X = 0]$$

The average population-level change in $y$ when *directly intervening* (or doing) $x$

$$\mathbf{E}(y \mid \mathrm{do}(x)) \qquad \neq \qquad \mathbf{E}(y \mid x)$$

The average population-level change in $y$ when accounting for *observed x*

Causation

Correlation

# *p*-values and confidence intervals

In the absence of *p*-values, I'm confused about how we report... significance?

# Imbens and *p*-values

**Nobody really cares about *p*-values**

**Decision makers want to know
a number or a range of numbers—
some sort of effect and uncertainty**

**Nobody cares how likely a number would be
in an imaginary null world!**

# Imbens's solution

**Report point estimates and some sort of range**

> "It would be preferable if reporting standards emphasized confidence intervals or standard errors, and, even better, Bayesian posterior intervals."

**Point estimate**

**The single number you calculate (mean, coefficient, etc.)**

**Uncertainty**

**A range of possible values**

# Population parameter

## Truth = Greek letter

**An single unknown number that is true for the entire population**

Proportion of left-handed students at GSU

Median rent of apartments in Atlanta

Proportion of red M&Ms produced in a factory

Treatment effect of your program

# Samples and estimates

We take a sample and make a guess

This single value is a *point estimate*

(This is the Greek letter with a hat)

# Variability

You have an estimate,
but how different might that
estimate be if you take another sample?

# Nets and confidence intervals

How confident are we that the sample picked up the population parameter?

Confidence interval is a net

We can be X% confident that our net is picking up that population parameter

If we took 100 samples, at least 95 of them would have the true population parameter in their 95% confidence intervals

A city manager wants to know the true average property value of single-owner homes in her city. She takes a random sample of 200 houses and builds a 95% confidence interval. The interval is ($180,000, $300,000).

**We're 95% confident that the interval ($180,000, $300,000) captured the true mean value**

# WARNING

It is way too tempting to say "We're 95% sure that the population parameter is X"

People do this all the time! People with PhDs!

YOU will do this too

# Nets

If you took lots of samples,
95% of their confidence intervals
would have the single true value in them

# Frequentism

This kind of statistics is called "frequentism"

The population parameter θ is fixed and singular while the data can vary

$$P(\mathrm{Data} \mid \theta)$$

You can do an experiment over and over again; take more and more samples and polls

# Frequentist confidence intervals

"We are 95% confident that this net captures the true population parameter"

~~"There's a 95% chance that the true value falls in this range"~~

# Bayesian statistics



Rev. Thomas Bayes

$$P(\theta \mid \mathrm{Data})$$

$$\color{orange}P(\mathrm{H} \mid \mathrm{E}) \color{black}= \frac{\color{red}P(\mathrm{H}) \color{black}\times \color{blue}P(\mathrm{E} \mid \mathrm{H})}{P(\mathrm{E})}$$

$$P(\text{H} \mid \text{E}) = \frac{P(\text{H}) \times P(\text{E} \mid \text{H})}{P(\text{E})}$$

$$P(\text{Hypothesis} \mid \text{Evidence}) =$$

$$\frac{P(\text{Hypothesis}) \times P(\text{Evidence} \mid \text{Hypothesis})}{P(\text{Evidence})}$$

# But the math is too hard!

**So we simulate!**

**(Monte Carlo Markov Chains, or MCMC)**

# Bayesianism and parameters

**In the world of frequentism,**
**there's a fixed population parameter**
**and the data can hypothetically vary**

$$P(\mathrm{Data} \mid \theta)$$

**In the world of Bayesianism,**
**the data is fixed** (you collected it just once!)
**and the population parameter can vary**

$$P(\theta \mid \mathrm{Data})$$

# Bayesian credible intervals

**(AKA posterior intervals)**

**"Given the data, there is a 95% probability that the true population parameter falls in the credible interval"**

# Intervals

## Frequentism

**There's a 95% probability that the range contains the true value**

**Probability of the range**

**Few people naturally think like this**

## Bayesianism

**There's a 95% probability that the true value falls in this range**

**Probability of the actual value**

**People *do* naturally think like this!**

# Thinking Bayesianly

We all think Bayesianly,
even if you've never heard of Bayesian stats

Every time you look at a confidence interval, you inherently think that the parameter is around that value, but that's wrong!
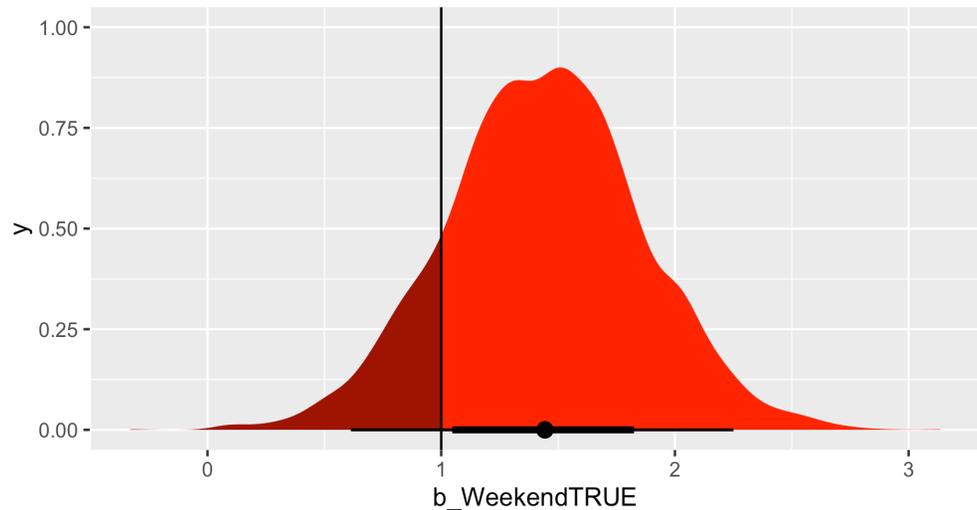
BUT Imbens cites research that
that's actually generally okay

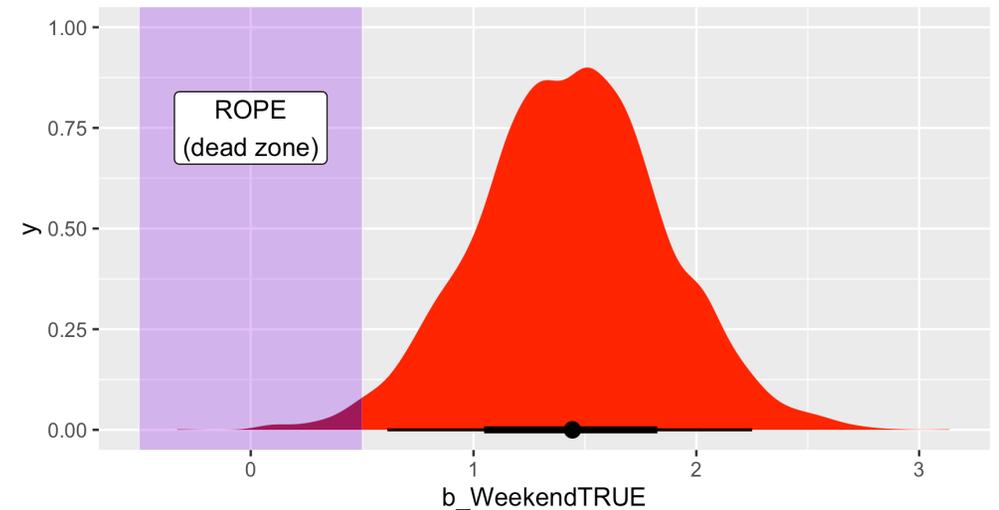Often credible intervals are super similar to confidence intervals

# Bayesian inference

## Inference without *p*-values!

### Probability of direction



Point shows median value;
thick black bar shows 66% credible interval;
thin black bar shows 95% credible interval

### Region of practical equivalence (ROPE)



ROPE (dead zone)

Point shows median value;
thick black bar shows 66% credible interval;
thin black bar shows 95% credible interval