

Science Before Statistics

# CAUSAL INFERENCE

Richard McElreath MPI-EVA Leipzig





Horoscope of Prince Iskandar, grandson  
of Tamerlane, born 25 April 1384

# Horoscopes

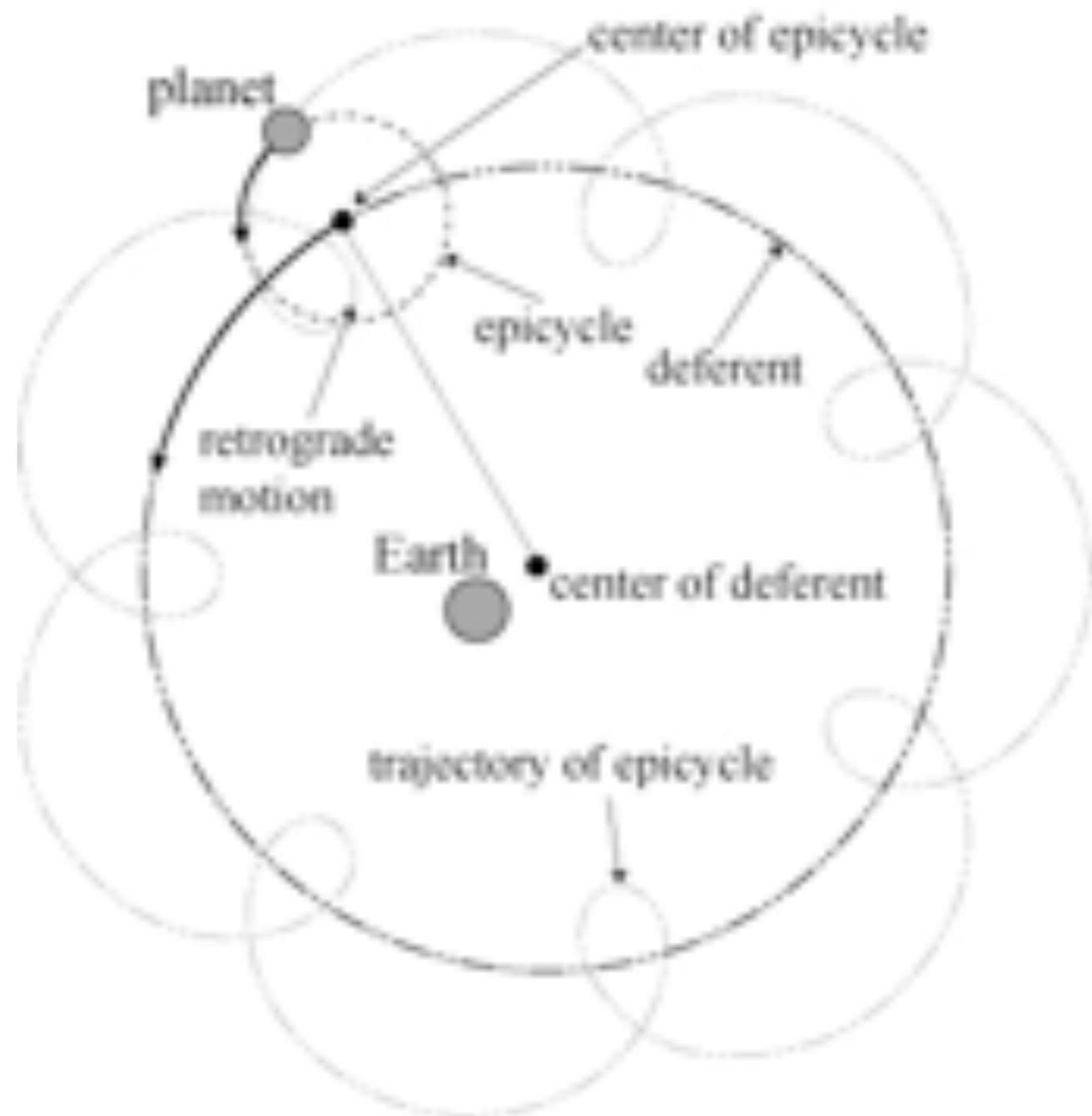
- Tell me your birthday, and I'll tell you how your life will unfold
- Tell me the data, and I'll tell you how to extract a significant result
- Both are superstition, false prophecy
- Prophecy must be vague, in order to sound correct



Mars and Saturn, wanderers

# Horoscopes

- Models can be accurate without being correct
- The data themselves do not contain information about causes
- It is not possible to reliably learn causes from data alone



A dramatic scene from the movie King Kong. The giant ape, King Kong, is shown in profile, facing right, and appears to be attacking or crashing through a city at night. The city is filled with smoke, fire, and debris. In the foreground, a small airplane is visible, having crashed into the ground. The sky is dark and filled with smoke and fire. The overall atmosphere is one of chaos and destruction.

**BAYES**

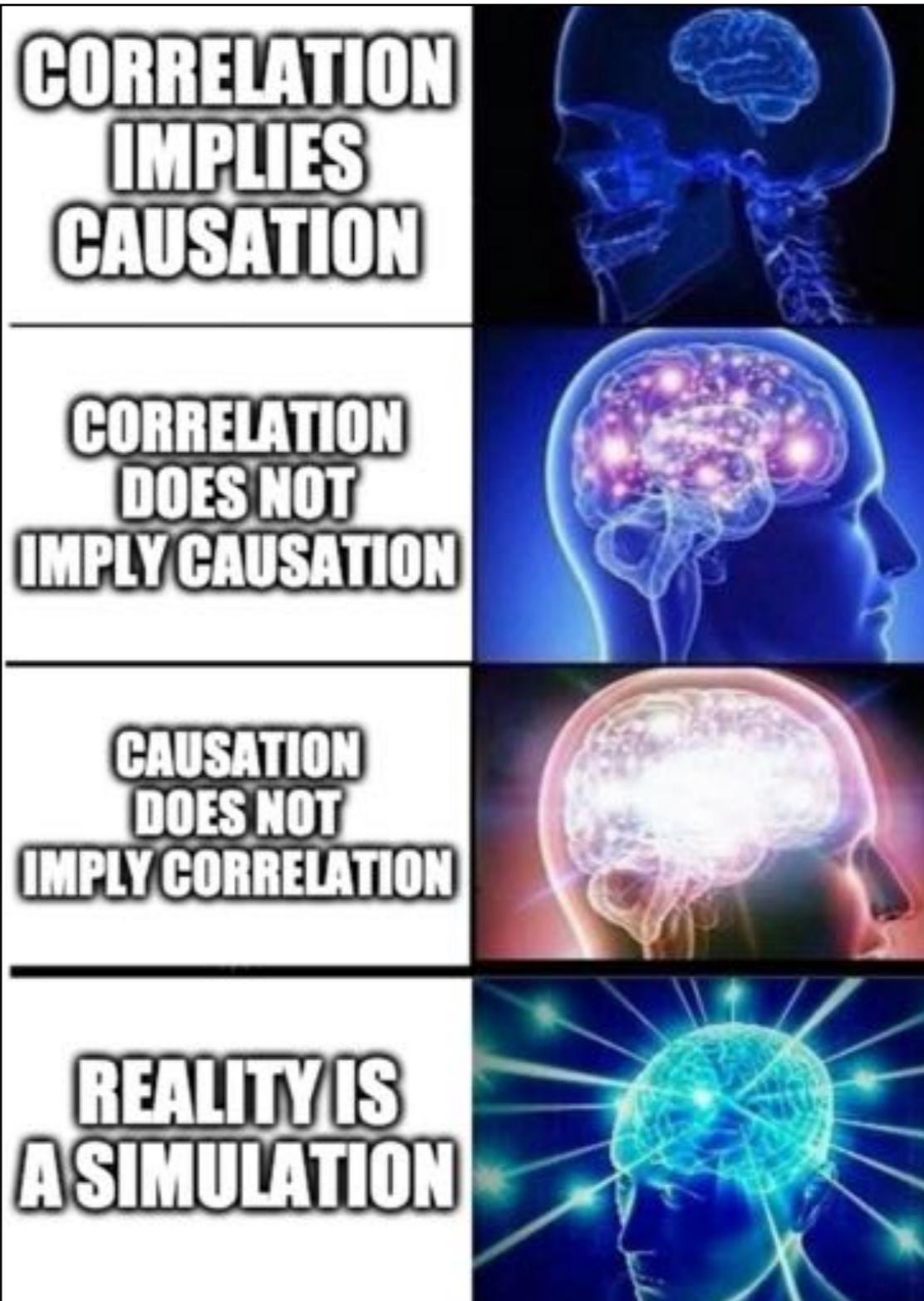
**FREQUENTISM**

A photograph of a coastal scene at sunset. The sky is filled with warm, orange and yellow hues, transitioning into cooler blues and purples. In the foreground, dark, silhouetted shapes of rocks and a small structure, possibly a lighthouse, are visible against the bright horizon. The ocean waves are visible in the distance.

# **CAUSAL INFERENCE**

# What is causal inference?

- Causal inference is more than learning an **association**.
- Two concepts:
  - Causal inference is **prediction of intervention**
  - Causal inference is **missing data imputation**



# Causal prediction



Knowing a **cause** means being able to predict  
the consequences of an **intervention**.  
*What if I do this?*

# Causal imputation



Knowing a **cause** means being able to construct unobserved **counterfactual outcomes**.  
*What if I had done something else?*

# Experiments are no refuge

- Why do experiments work? When do they work?
- Should you test for balance?
- What if treatment were imperfect?
- Should you control for anything? Everything?
- What is the causal effect in the target population?
- Answers depend upon causal assumptions.

# Description is no refuge

- Descriptive research requires causal inference
  - To do more than describe a **sample**, need causal assumptions
  - Sampling bias, stratified sampling, post-stratification, missing data, measurement error
  - We should not compare samples but **populations**

description



inference

# Today's Agenda

- Part 1: **Causal Salad**  
You must unlearn what you have learned
- Part 2: **Causal Design**  
Drawing our assumptions
- Part 3: **Full Luxury Bayesian Inference**  
Doing the computations





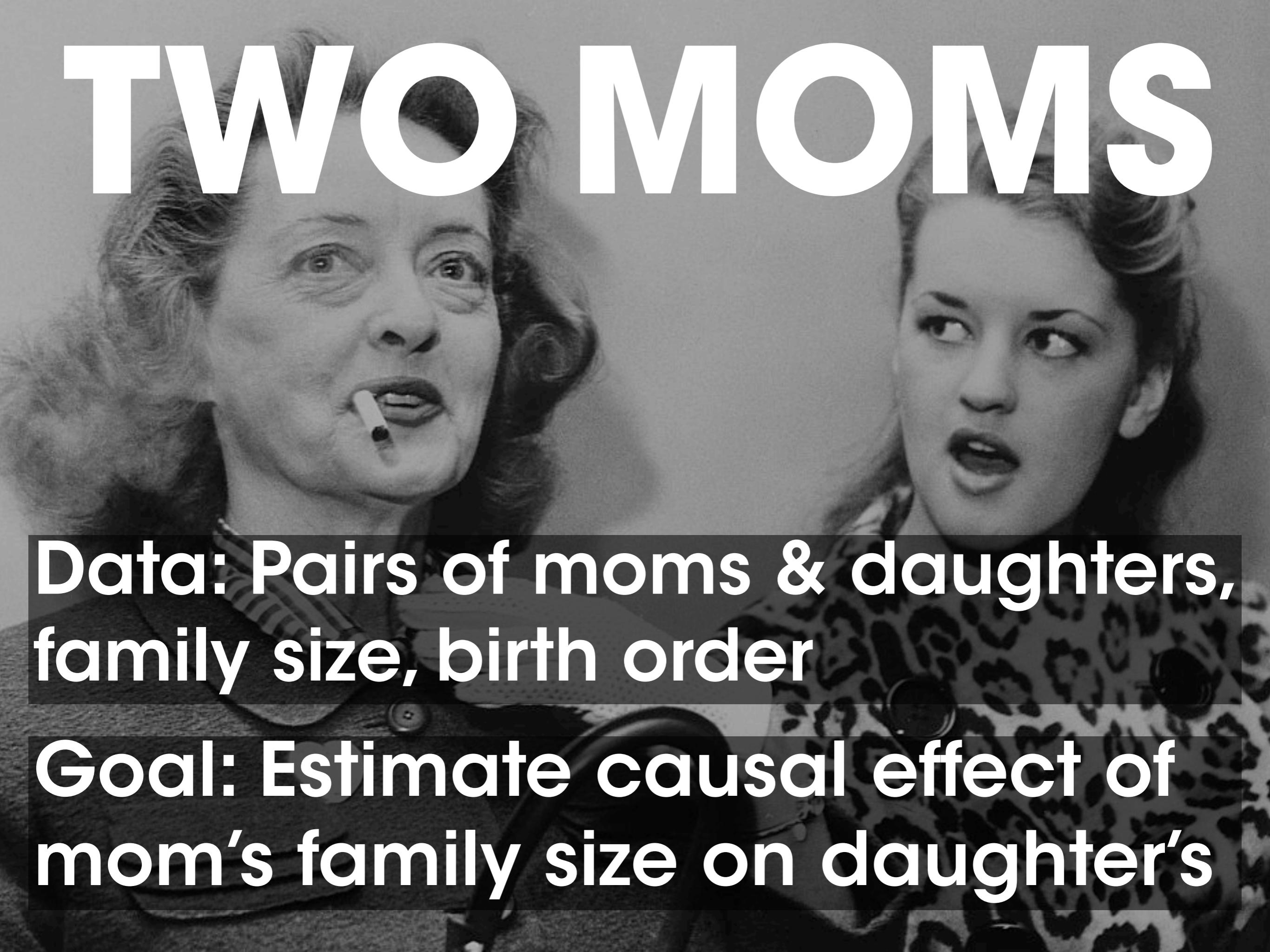
# TWO MOMS



# TWO MOMS

Data: Pairs of moms & daughters,  
family size, birth order

# TWO MOMS



**Data:** Pairs of moms & daughters,  
family size, birth order

**Goal:** Estimate causal effect of  
mom's family size on daughter's

# PEER BIAS



# PEER BIAS



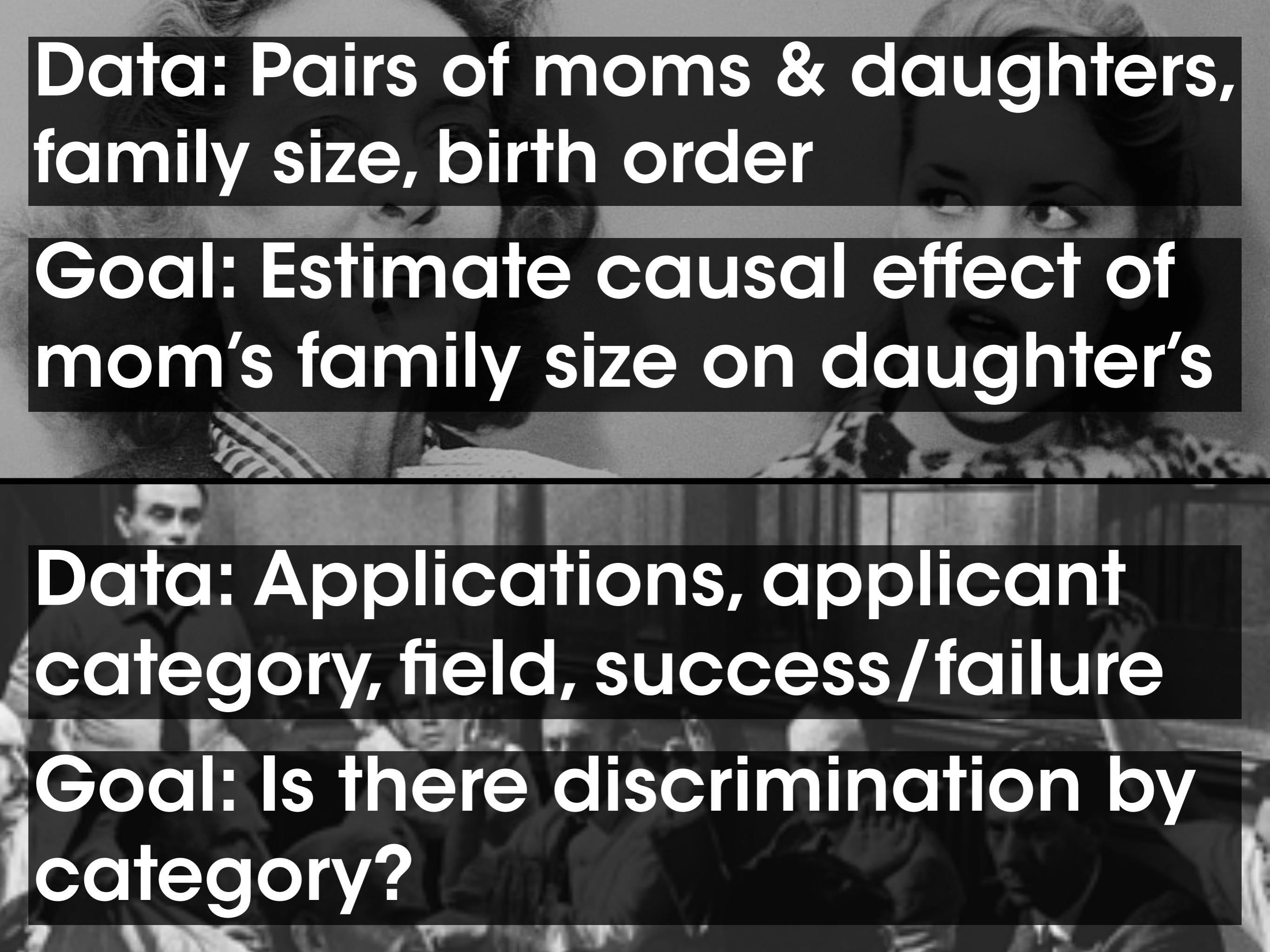
Data: Applications, applicant  
category, field, success/failure

# PEER BIAS



Data: Applications, applicant category, field, success/failure

Goal: Is there discrimination by category?



**Data: Pairs of moms & daughters,  
family size, birth order**

**Goal: Estimate causal effect of  
mom's family size on daughter's**

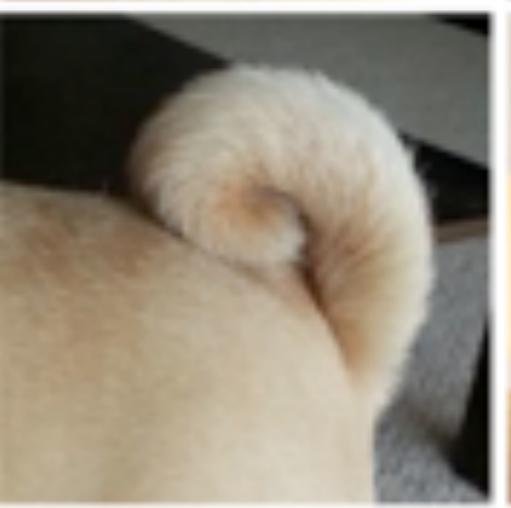


**Data: Applications, applicant  
category, field, success/failure**

**Goal: Is there discrimination by  
category?**

# CAUSAL SALAD





# Robots are cause-blind



20

"panda"  
57.7% confidence

# Robots are cause-blind



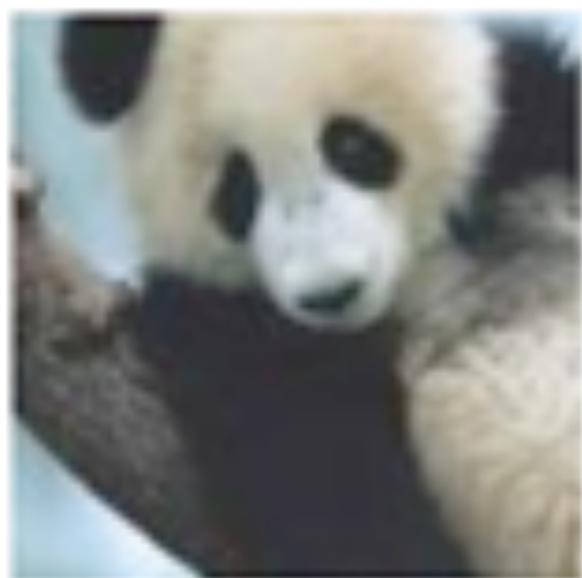
$+ .007 \times$



$x$   
"panda"  
57.7% confidence

$\text{sign}(\nabla_x J(\theta, x, y))$   
"nematode"  
8.2% confidence

# Robots are cause-blind



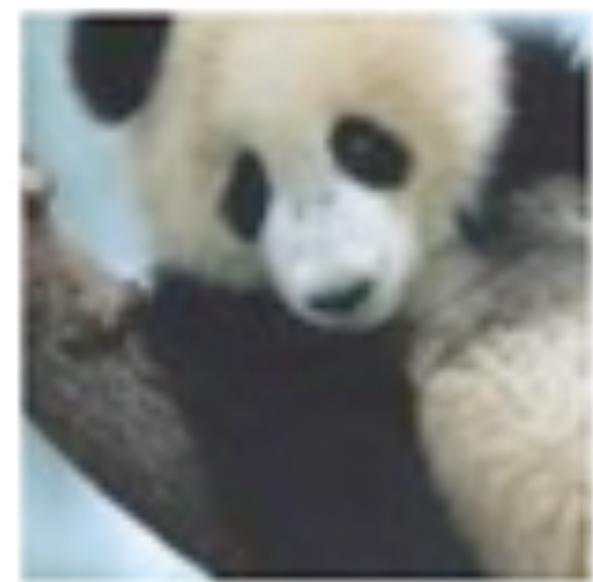
$x$   
“panda”  
57.7% confidence

+ .007 ×



$\text{sign}(\nabla_x J(\theta, x, y))$   
“nematode”  
8.2% confidence

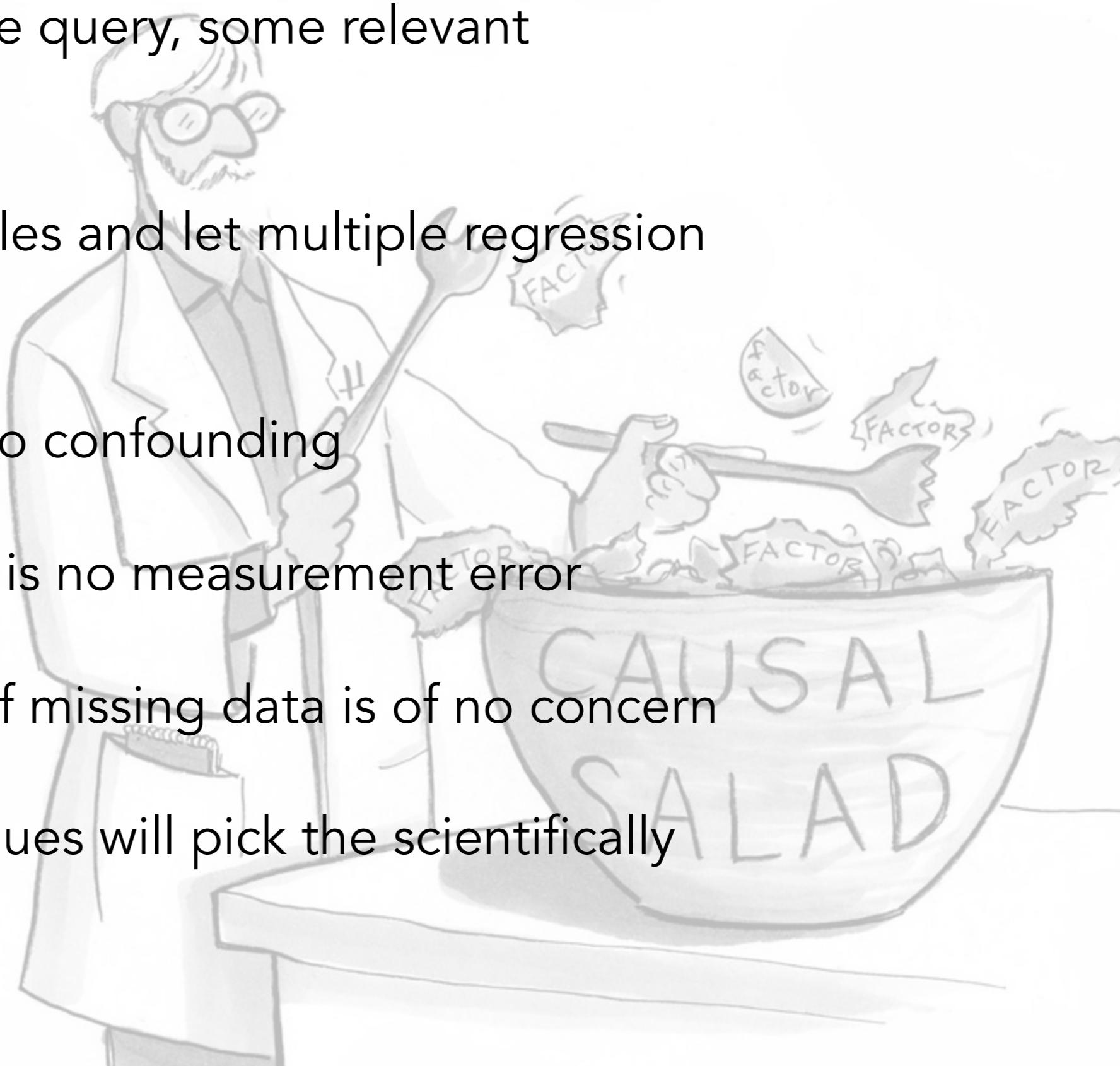
=



$x +$   
 $c \text{sign}(\nabla_x J(\theta, x, y))$   
“gibbon”  
99.3 % confidence

# Causal Salad

- Ingredients: Vague query, some relevant variables
- Add all the variables and let multiple regression sort it out
- Pretend there is no confounding
- Pretend there are is no measurement error
- Pretend pattern of missing data is of no concern
- Pretend AIC/p-values will pick the scientifically correct model



# Multiverse analysis

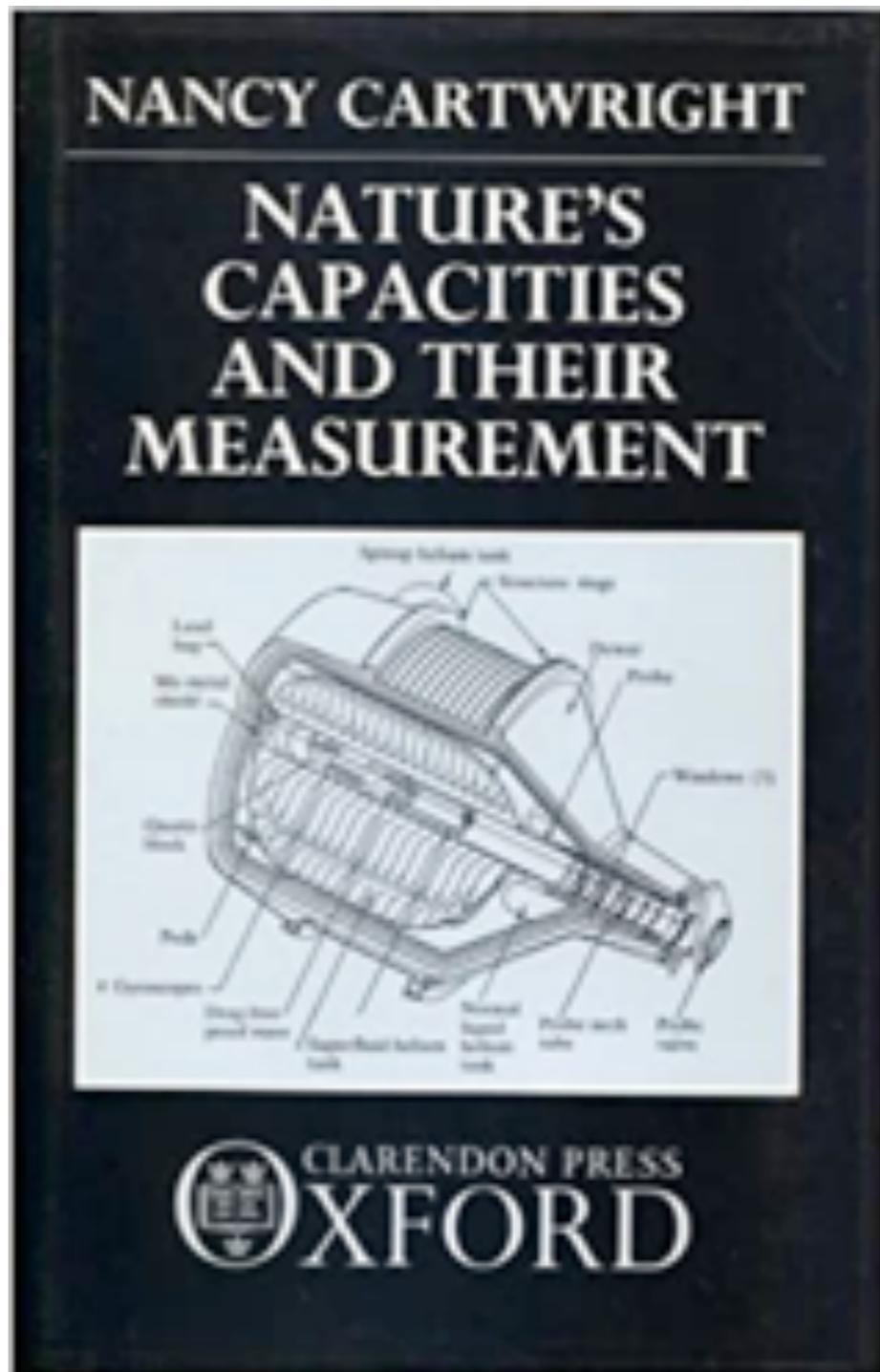
- Multiverse analysis: “performing all analyses across the whole set of alternatively processed data sets corresponding to a large set of **reasonable** scenarios”
- No surprise that results are sensitive to model structure and data processing
- Conditions that make a scenario “reasonable” are causal, unstated
- Does not prevent causal salad

**Table 1.** Processing choices

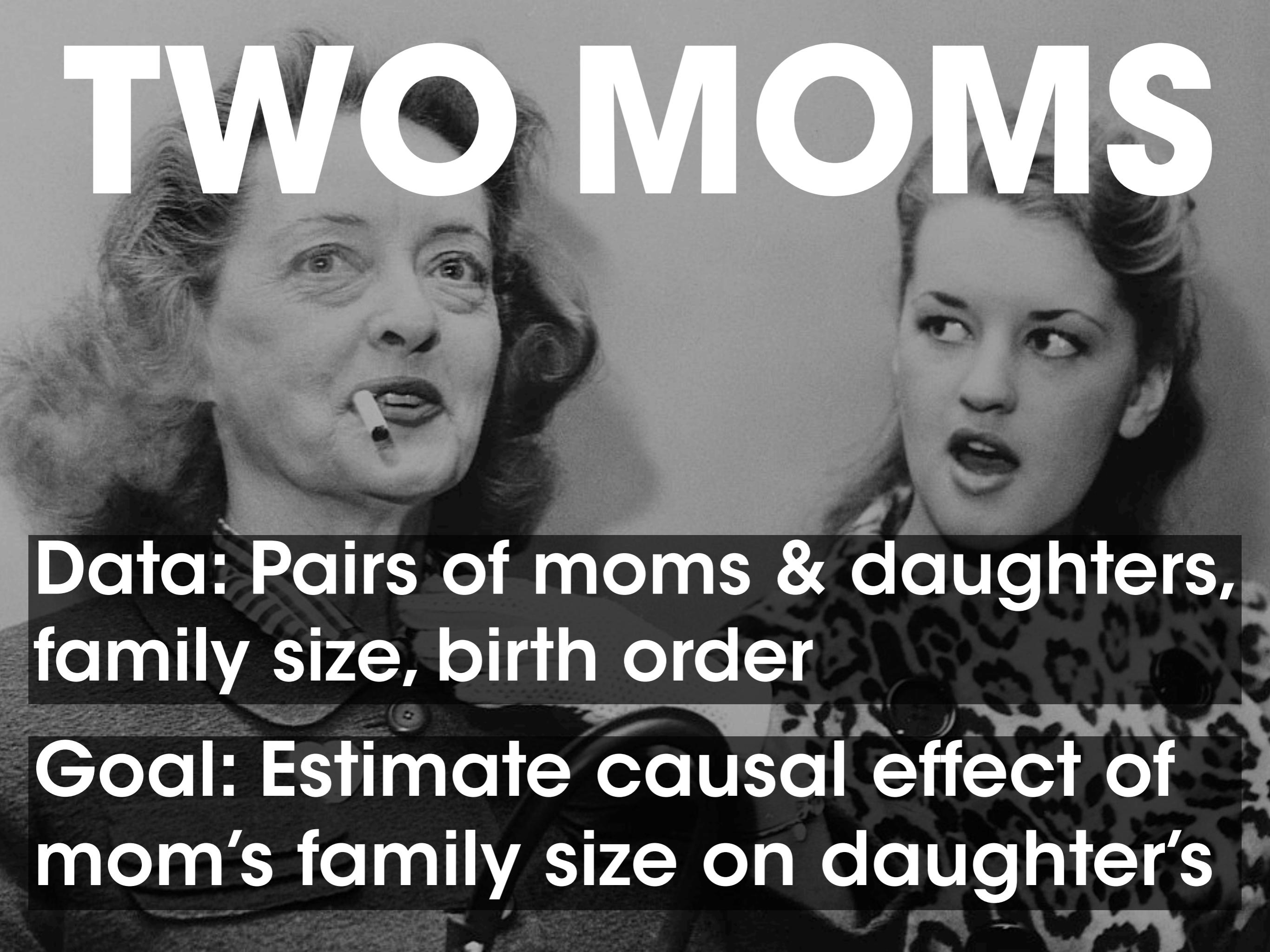
1. Assessment of fertility (F)—high vs low.
  - (a) F1: high = cycle days 7–14; low = cycle days 17–25
  - (b) F2: high = cycle days 6–14; low = cycle days 17–27
  - (c) F3: high = cycle days 9–17; low = cycle days 18–25
  - (d) F4: high = cycle days 8–14; low = cycle days 1–7 and 15–28
  - (e) F5: high = cycle days 9–17; low = cycle days 1–8 and 18–28
2. Next menstrual onset (NMO)
  - (a) NMO1: reported start date previous menstrual onset + computed cycle length
  - (b) NMO2: reported start date previous menstrual onset + reported cycle length
  - (c) NMO3: reported estimate of next menstrual onset
3. Assessment of relationship status (R) (single vs relationship)
  - (a) R1: single = response options 1 and 2; relationship = response options 3 and 4
  - (b) R2: single = response option 1; relationship = response options 2, 3, and 4
  - (c) R3: single = response option 1; relationship = response options 3 and 4
4. Exclusion of women based on cycle length (ECL)
  - (a) ECL1: no exclusion based on cycle length
  - (b) ECL2: exclusion of participants with computed cycle length greater than 25 or less than 35 days
  - (c) ECL3: exclusion of participants with reported cycle length greater than 25 or less than 35 days
5. Exclusion of women based on certainty ratings of start dates of two previous menstrual periods (EC)
  - (a) EC1: no exclusion based on certainty ratings
  - (b) EC2: exclusion of participants who are not certain about at least one start date (i.e., sure less than 6)

# No causes in; No causes out

- Statistical models alone insufficient;  
They do not contain causal information
- Multiple regression does not  
distinguish causes from confounds
- p-values are not causal statements
- AIC etc are purely predictive



# TWO MOMS



**Data:** Pairs of moms & daughters,  
family size, birth order

**Goal:** Estimate causal effect of  
mom's family size on daughter's

# Two Moms & a Regression

Variables:

Family sizes **M** and **D**

Birth orders **B1** (mom's)  
and **B2** (daughter's)

How would you construct a  
regression to estimate  
influence of **M** on **D**?

# Two Moms & a Regression

- Key question: Should you add **B1** and **B2** to the model?

$$M \sim D$$

$$M \sim D + B1$$

$$M \sim D + B2$$

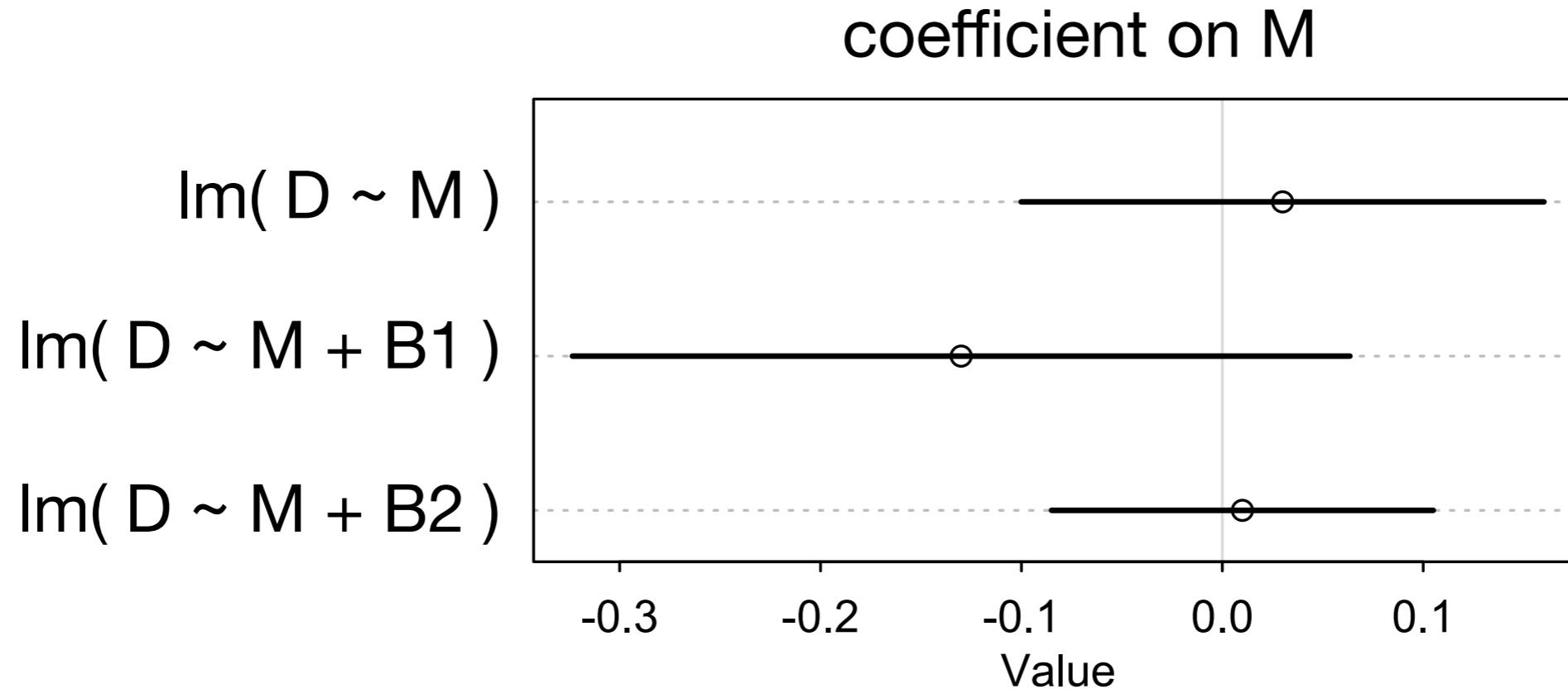
$$M \sim D + B1 + B2$$

$$M \sim D * B1 * B2$$

# Two Moms & a Regression

- To know the right thing to do, need a simulation
- Assume that mom's family size has ZERO influence on daughter's
- Assume first borns have higher fertility
- Simulate 200 mother-daughter pairs

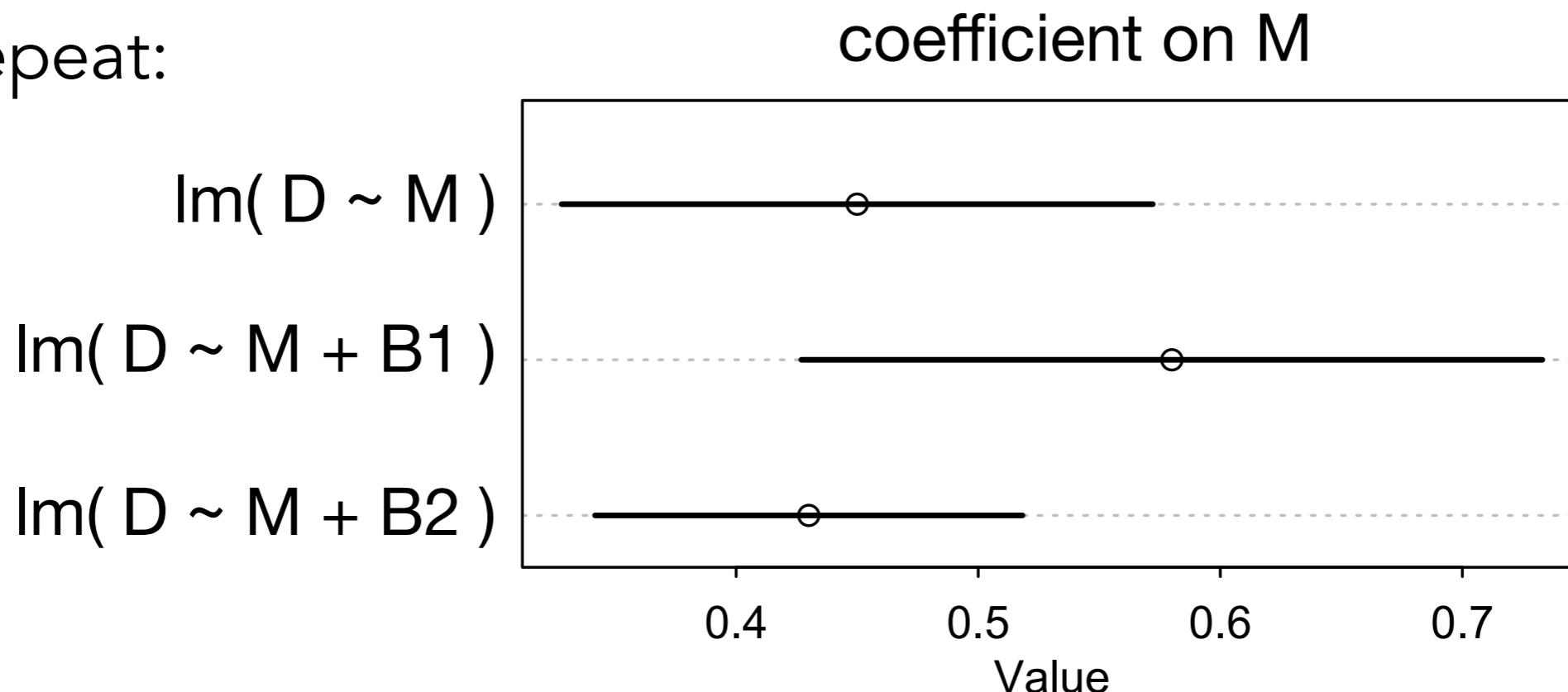
# Two Moms & a Regression



Just an example, but representative of average result.  
Adding **B1** hurts inference, adding **B2** helps. Why?

# Two Moms & a Confound

- Now suppose there is an unobserved confound:  
A common cause of **M** and **D** such as wealth or education
- Repeat:



When there is a confound, **B1** exaggerates the bias

# Two Moms & an Info Criterion

- Predictive criteria like AIC do not help

AIC(  $\text{Im}( D \sim M )$  ): 764

AIC(  $\text{Im}( D \sim M + B_1 )$  ): 757

AIC(  $\text{Im}( D \sim M + B_2 )$  ): 635

AIC(  $\text{Im}( D \sim M + B_1 + B_2 )$  ): 629

- Why? A model can make good predictions without correct causal structure. Association is not causation.
- Don't ask about  $p$ -values — they aren't even predictive criteria

# Two Moms Summary

- Assuming birth order influences family size:
  - Including mom's birth order (**B1**) hurts inference
  - Including daughter's birth order (**B2**) helps
  - AIC etc select models that include **B1**
  - When **M** and **D** are confounded, adding **B1** exaggerates the confound

# PEER BIAS



Data: Applications, applicant category, field, success/failure

Goal: Is there discrimination by category?

# Measuring Peer Bias

Variables:

Applicant category **X**

Field/department **E**

Outcome **Y**

How would you construct a regression to estimate influence of **X** on **Y**?

# Peer Bias

- Long history of attempts to address this question.
- Key finding: Including field/department in the model changes the result.
- But what does this change mean?



# An Empirical Analysis of Racial Differences in Police Use of Force

Roland G. Fryer Jr.

Harvard University and National Bureau of Economic Research

This paper explores racial differences in police use of force. On non-lethal uses of force, blacks and Hispanics are more than 50 percent more likely to experience some form of force in interactions with police. Adding controls that account for important context and civilian behavior reduces, but cannot fully explain, these disparities. On the most extreme use of force—officer-involved shootings—we find no racial differences either in the raw data or when contextual factors are taken into account. We argue that the patterns in the data are consistent with a model in which police officers are utility maximizers, a fraction of whom have a preference for discrimination, who incur relatively high expected costs of officer-involved shootings.

We can never be satisfied as long as the Negro is the victim of the unspeakable horrors of police brutality. (Martin Luther King Jr., August 28, 1963)

## I. Introduction

From “Bloody Sunday” on the Edmund Pettus Bridge to the public beatings of Rodney King, Bryant Allen, and Freddie Helms, the relationship

This work has benefited greatly from discussions and debate with Chief William Evans, Chief Charles McClelland, Chief Martha Montalvo, Sergeant Stephen Morrison, Jon Murad, Lynn Overmann, Chief Bud Riley, and Chief Scott Thomson. I am grateful to David Card, Kerwin Charles, Christian Dustmann, Michael Greenstone, James Heckman, Richard Holden, Lawrence Katz, Steven Levitt, Jens Ludwig, Glenn Loury, Kevin Murphy, Derek Neal, John Overdeck, Jesse Shapiro, Andrei Shleifer, Jorg Spenkuch, Max Stone, John Van Reenan, Christopher Winship, and seminar participants at Brown University, University of Chicago, London

## 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 uncontrollable 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 racial discrimination in the decision to detain civilians—a decision that determines which encounters appear in police administrative data at all—then estimates of the effect of civilian race on subsequent police behavior are

biased absent additional data and/or strong and uncontrollable 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 research design that can mitigate these threats to inference. These are discussed in the context of a detailed and feasible proposed study of racial bias in traffic stops.

As we show in this article, the difficulty of estimating racial bias using police records stems from a thorny combination of mediation (Hernán, Hernández-Díaz, and Robins 2004; Imai et al. 2011; Pearl 2001; Robins, Hernán, and Brumback 2000; VanderWeele 2009) and selection (Heckman 1979; Lee 2009): the effect of civilian race on the outcome of a police encounter is mediated by whether the civilian is stopped by police, but analysts only have data for one level of the mediator—that is, data on stopped individuals. Because of this, police records do not contain a representative sample of all individuals that police observe, but rather only those civilian encounters which escalated to the point of triggering a reporting requirement. If a civilian’s race affects whether officers choose to stop that civilian (Gelman, Fagan, and Kiss 2007; Glaser 2014), then analyzing administrative police records amounts to conditioning on a variable that is itself affected by suspect race, namely, whether a suspect appears in the data at all. This could occur if officers have a higher

---

Dean Knox , Assistant Professor of Politics, Princeton University, [dcknox@princeton.edu](mailto:dcknox@princeton.edu).

Will Lowe , Senior Research Scientist, Hertie School of Governance, [lowe@hertie-school.org](mailto:lowe@hertie-school.org).

Jonathan Mummolo , Assistant Professor of Politics and Public Affairs, Princeton University, [jmummolo@princeton.edu](mailto:jmummolo@princeton.edu).

We thank Matt Blackwell, Chuck Cameron, Tom Clark, Scott Cunningham, Lauren Davenport, Naoki Egami, Jeffrey Fagan, Avi Feller, Adam Glynn, Phillip Atiba Goff, Justin Grimmer, Andy Hall, Anna Harvey, Dan Hopkins, Matias Iaryczower, Kosuke Imai, Damon Jones, Dorothy Kronick, Shiro Kuriwaki, Neil Malhotra, Moritz Marbach, Nolan McCarty, Cyrus Samii, Maya Sen, Tara Slough, Rocio Titiunik, Tyler VanderWeele, Vesla Weaver, and Sean Westwood for helpful feedback. We thank Michael Pomirchy for research assistance. Replication files are available at the American Political Science Review Dataverse: <https://doi.org/10.7910/DVN/KFQOCV>.

Received: April 26, 2019; revised: October 10, 2019; accepted: January 8, 2020.

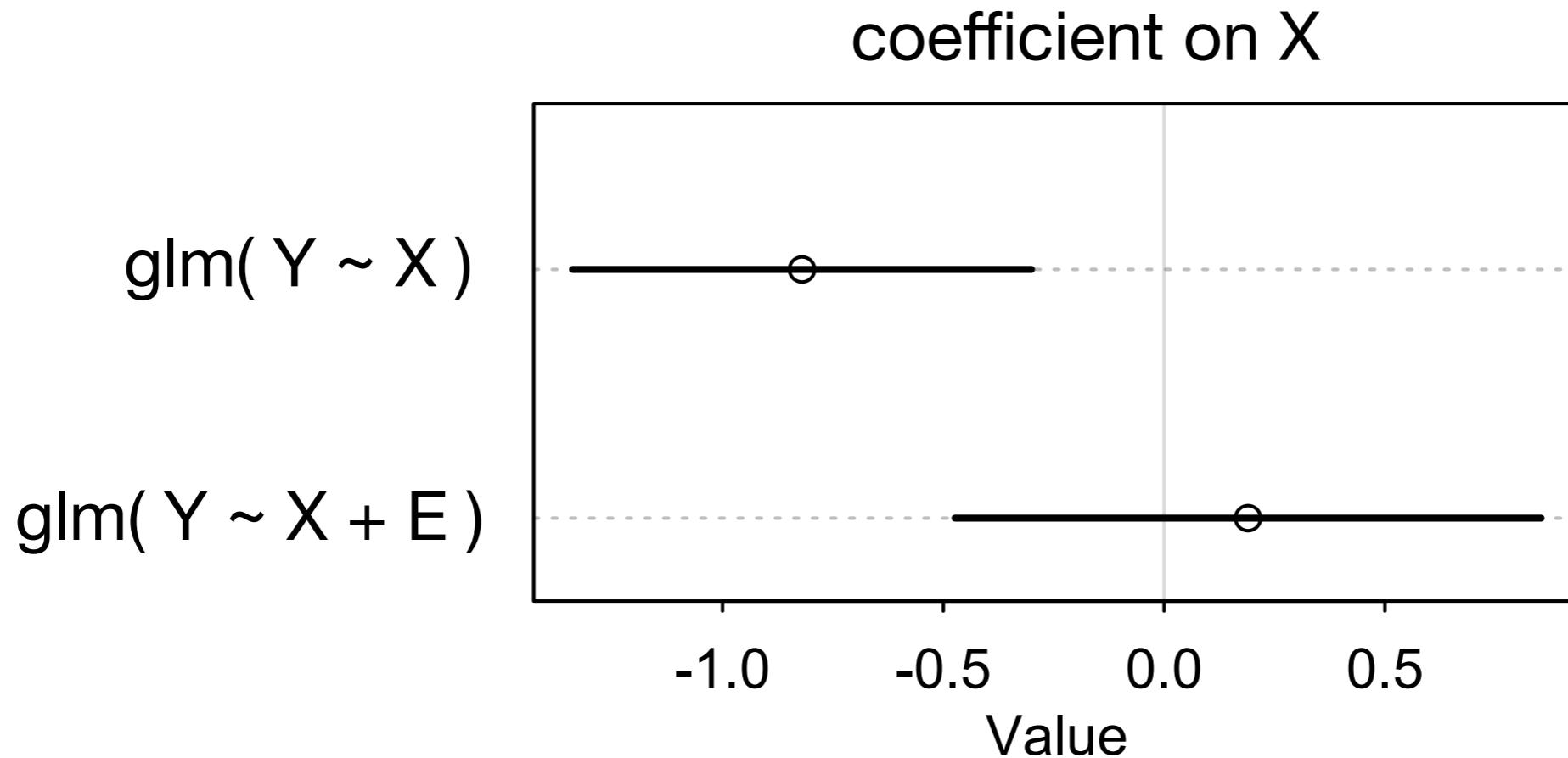
# Peer Bias

- Simulation example:
  - 500 applications in two subjects (**E**)
  - Subjects vary by average acceptance rate
  - Category **X** is NOT a target of discrimination
  - Outcome **Y** accepted/rejected
- Consider two regressions:

$$Y \sim X$$

$$Y \sim X + E$$

# Peer Bias



Why does  $E$  make false evidence of discrimination vanish?

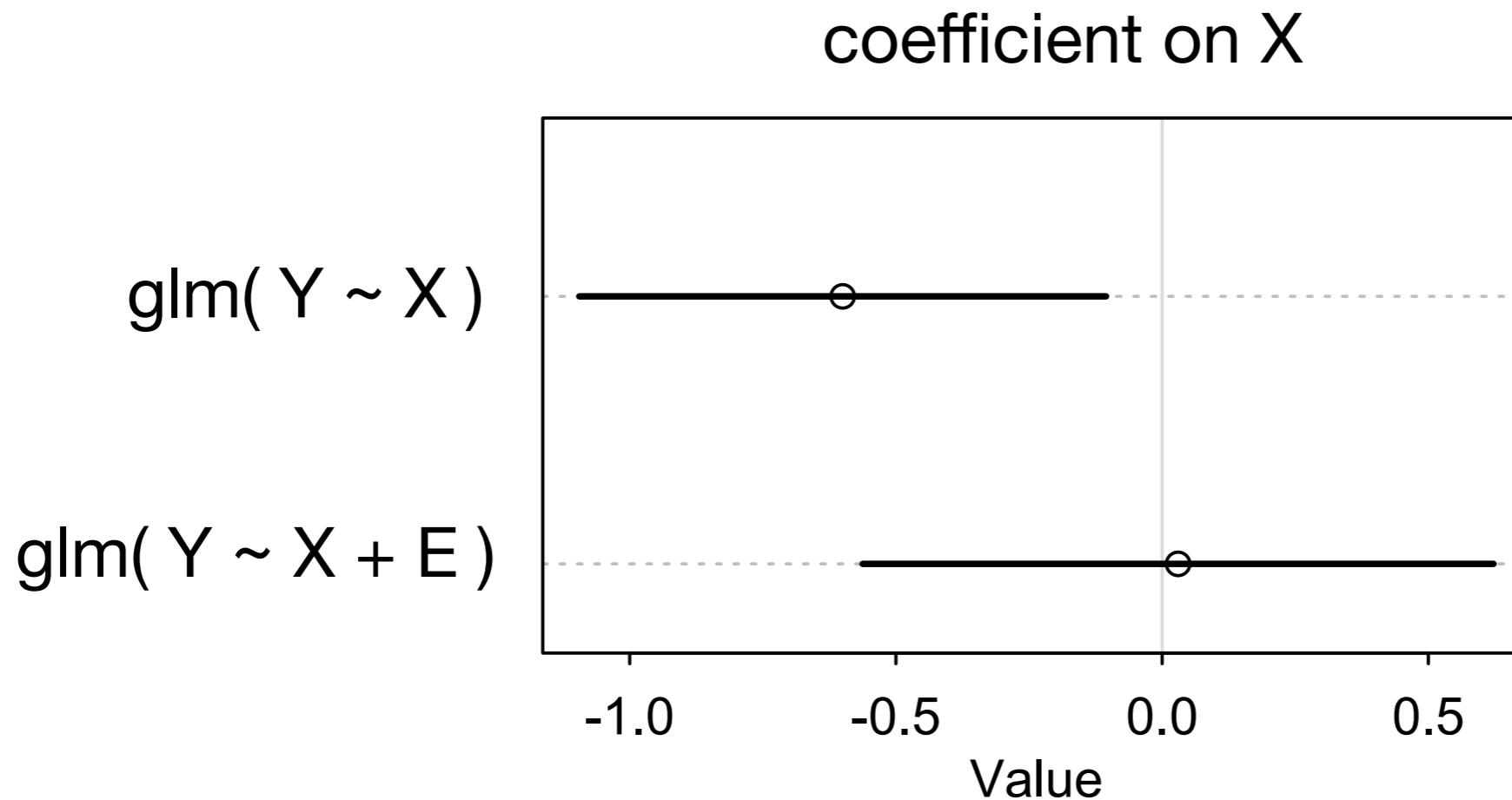
# Peer Bias Again

- ANOTHER Simulation example:
  - 500 applications in two subjects (**E**)
  - Subjects vary by average acceptance rate
  - Category **X** target of discrimination in both subjects
  - Outcome **Y** accepted/rejected
- Consider two regressions:

$$Y \sim X$$

$$Y \sim X + E$$

# Peer Bias Again



Why does  $E$  make discrimination vanish, even though there is discrimination in the simulation?

# Peer Bias Summary

- Assuming no discrimination against **X**
  - Model without field (**E**) shows discrimination
  - Model with **E** finds no evidence of it
- Assuming there is discrimination based on category **X**
  - Model without field (**E**) finds lower success for **X**
  - Model with **E** finds no evidence of lower success
- In real data, we don't know the truth
- Change of coefficient could be revelation or deception