

RESEARCHER CHOICES AND BIAS

Data Analysis for Journalism and Political Communication
(Fall 2025)

Prof. Bell

On the Supreme Court's Emergency Docket, Sharp Partisan Divides

This is apparent in the overall numbers, with the Trump administration prevailing much more often than its predecessor had — 84 percent of the time, compared with 53 percent for the Biden administration. That is perhaps unsurprising, given that the court is dominated by six Republican appointees.

Same Data, Different Conclusions

Twenty-nine research teams were given the same set of soccer data and asked to determine if referees are more likely to give red cards to dark-skinned players. Each team used a different statistical method, and each found a different relationship between skin color and red cards.

Referees are
three times as
likely to give red
cards to
dark-skinned
players

Statistically
significant results
showing referees are
more likely to give red
cards to dark-skinned
players

Twice as likely

Equally likely

ONE RESEARCH TEAM

95% CONFIDENCE INTERVAL

Non-significant
results

RESEARCHER CHOICES

- ➊ What is my hypothesis?

RESEARCHER CHOICES

- ① What is my hypothesis?
 - ② What do I measure?

RESEARCHER CHOICES

- ① What is my hypothesis?
 - ② What do I measure?
 - ③ How do I collect my data?

RESEARCHER CHOICES

- ① What is my hypothesis?
- ② What do I measure?
- ③ How do I collect my data?
- ④ How much data do I collect?

RESEARCHER CHOICES

- ① What is my hypothesis?
- ② What do I measure?
- ③ How do I collect my data?
- ④ How much data do I collect?
- ⑤ What statistical analyses do I use?

RESEARCHER CHOICES

- ① What is my hypothesis?
- ② What do I measure?
- ③ How do I collect my data?
- ④ How much data do I collect?
- ⑤ What statistical analyses do I use?
- ⑥ How do I handle outliers, missing data, and other peculiarities?

WHAT IS MY HYPOTHESIS?

Choosing a hypothesis is all about avoiding error:

- **Type I error:** False positives
 - ▶ Sending the innocent to jail - this is bad!

WHAT IS MY HYPOTHESIS?

Choosing a hypothesis is all about avoiding error:

- **Type I error:** False positives
 - ▶ Sending the innocent to jail - this is bad!
- **Type II error:** False negatives
 - ▶ Letting the guilty go free - we can accept this

WHAT IS MY HYPOTHESIS?

Choosing a hypothesis is all about avoiding error:

- **Type I error:** False positives
 - ▶ Sending the innocent to jail - this is bad!
- **Type II error:** False negatives
 - ▶ Letting the guilty go free - we can accept this

Our goal is to reduce Type I error. Assume that the data is innocent (that the hypothesis is false) until it is proven guilty.

MEASURING TYPE I ERROR

P-value

The chance that we would get a particular result from our test if the true answer is false

MEASURING TYPE I ERROR

P-value

The chance that we would get a particular result from our test if the true answer is false

- The p-value is our chance of committing a Type I error - sending the innocent to jail

MEASURING TYPE I ERROR

P-value

The chance that we would get a particular result from our test if the true answer is false

- The p-value is our chance of committing a Type I error - sending the innocent to jail
- Common p-value cut-offs in scientific research: .01, .05, and .1 indicate **statistical significance**

MEASURING TYPE I ERROR

Hypothesis: A student cheated on an exam.

P-value (chance that we conclude that student cheated, but they did not): .50

MEASURING TYPE I ERROR

Hypothesis: A student cheated on an exam.

- The student performed much better on this exam than on previous exams

P-value (chance that we conclude that student cheated, but they did not): .25

MEASURING TYPE I ERROR

Hypothesis: A student cheated on an exam.

- The student performed much better on this exam than on previous exams
- The student finished their exam more quickly than other students

P-value (chance that we conclude that student cheated, but they did not): .15

MEASURING TYPE I ERROR

Hypothesis: A student cheated on an exam.

- The student performed much better on this exam than on previous exams
- The student finished their exam more quickly than other students
- The student's roommate saw them up all night studying before the exam

P-value (chance that we conclude that student cheated, but they did not): .40

MEASURING TYPE I ERROR

Hypothesis: A student cheated on an exam.

- The student performed much better on this exam than on previous exams
- The student finished their exam more quickly than other students
- The student's roommate saw them up all night studying before the exam
- The student missed the same questions as other students

P-value (chance that we conclude that student cheated, but they did not): .75

THE SCIENTIFIC METHOD

- P-values are a product of the data we use and our choices about what we include and exclude from the analysis.

THE SCIENTIFIC METHOD

- P-values are a product of the data we use and our choices about what we include and exclude from the analysis.
- We follow the scientific method: Theory ⇒ Hypothesis ⇒ Test ⇒ Analyze ⇒ Report

THE SCIENTIFIC METHOD

- P-values are a product of the data we use and our choices about what we include and exclude from the analysis.
- We follow the scientific method: Theory ⇒ Hypothesis ⇒ Test ⇒ Analyze ⇒ Report
- But in practice, no analysis plan survives contact with the data

ARE DEMOCRATS OR REPUBLICANS GOOD FOR THE ECONOMY?

Use (a recreation of) FiveThirtyEight's online tool to test what you think is the best approach to answering the question. There are no right or wrong answers - just select the model you think is best, and report your results in the form:

<https://bit.ly/smpa2152>



(The link to the tool is on the form.)

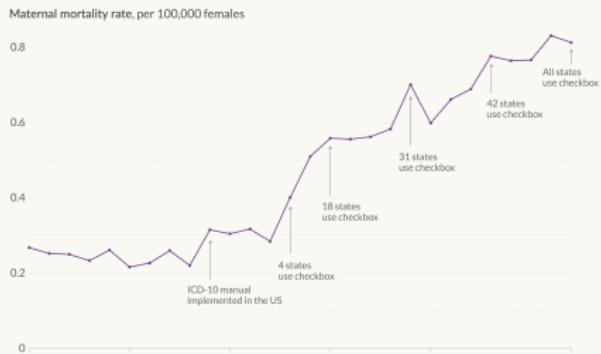
WHAT DO I MEASURE?

Operationalization

The process of defining a measurable version of a concept.

The US maternal mortality rate rose as more states adopted the “pregnancy checkbox”

As more states in the US adopted the “pregnancy checkbox” on death certificates — which asked if the deceased had been pregnant or recently pregnant — the reported maternal mortality rate rose.

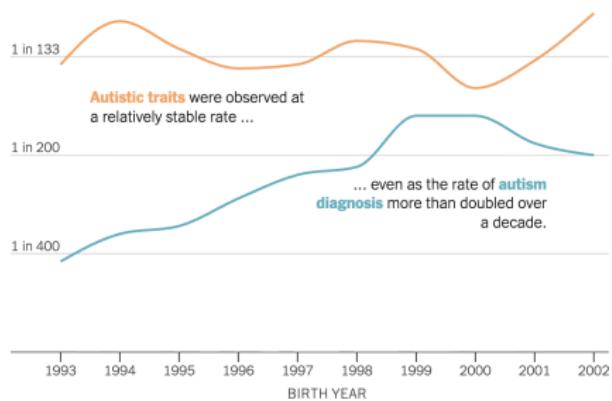


Source: WHO Mortality Database (2022). Adapted from ICES Joseph et al. (2022) Maternal mortality in the United States. Data includes “late maternal deaths”, which occur up to 1 year after the end of pregnancy.

OurWorldInData.org – Research and data to make progress against the world's largest problems.

Licensed under CC-BY by the author Saloni Dattani

A Swedish study illuminated the difference between autism as a diagnosis and autism as a condition



Source: Lundstrom et al., The BMJ (2015)

PRINCIPLES OF GOOD OPERATIONALIZATION

① Unambiguous



PRINCIPLES OF GOOD OPERATIONALIZATION

- ① Unambiguous
- ② Parsimonious

PRINCIPLES OF GOOD OPERATIONALIZATION

- 1 Unambiguous
- 2 Parsimonious
- 3 Accurate

→ NOTE: Please answer BOTH Question 5 about Hispanic origin and Question 6 about race. For this census, Hispanic origins are not races.
5. Is this person of Hispanic, Latino, or Spanish origin?

- No, not of Hispanic, Latino, or Spanish origin
- Yes, Mexican, Mexican Am., Chicano
- Yes, Puerto Rican
- Yes, Cuban
- Yes, another Hispanic, Latino, or Spanish origin — Print origin, for example, Argentinean, Colombian, Dominican, Nicaraguan, Salvadoran, Spaniard, and so on. ↗

6. What is this person's race? Mark one or more boxes.

- White
- Black, African Am., or Negro
- American Indian or Alaska Native — Print name of enrolled or principal tribe. ↗

- Asian Indian Japanese Native Hawaiian
- Chinese Korean Guamanian or Chamorro
- Filipino Vietnamese Samoa
- Other Asian — Print race, for example, Hmong, Laotian, Thai, Pakistani, Cambodian, and so on. ↗ Other Pacific Islander — Print race, for example, Fijian, Tongan, and so on. ↗

- Some other race — Print race. ↗

What is your race or ethnicity?
Select all that apply AND enter additional details in the spaces below.
Note, you may report more than one group.

□ WHITE — Provide details below.

- German Irish English
- Italian Polish French

Enter, for example, Scottish, Norwegian, Dutch, etc.

□ HISPANIC OR LATINO — Provide details below.

- Mexican American Puerto Rican Cuban
- Salvadoran Dominican Colombian

Enter, for example, Guatemalan, Honduran, Ecuadorian, etc.

□ BLACK OR AFRICAN AMERICAN — Provide details below.

- African American Jamaican Haitian
- Nigerian Ethiopian Somali

Enter, for example, Ghanaian, South African, Barbadian, etc.

□ ASIAN — Provide details below.

- Chinese Filipino Asian Indian
- Vietnamese Korean Japanese

Enter, for example, Pakistani, Cambodian, Hmong, etc.

□ AMERICAN INDIAN OR ALASKA NATIVE — Enter, for example, Navajo Nation, Blackfeet Tribe, Mayon, Aztec, Native Village of Barrow Inupiat Tribal Government, Pima, etc.

□ MIDDLE EASTERN OR NORTH AFRICAN — Provide details below.

- Lebanese Iranian Egyptian
- Syrian Moroccan Israeli

Enter, for example, Algerian, Iraqi, Kurdish, etc.

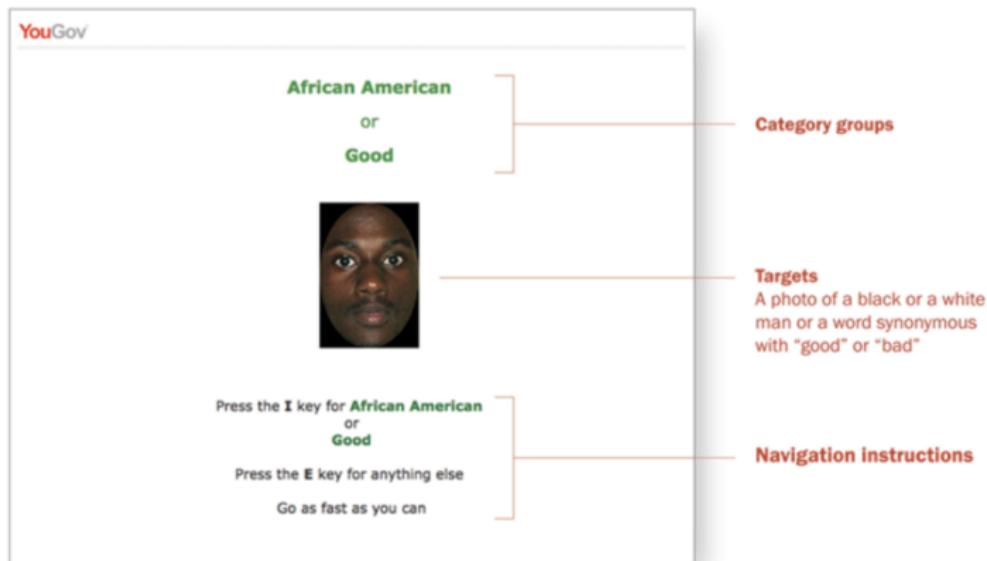
□ NATIVE HAWAIIAN OR PACIFIC ISLANDER — Provide details below.

- Native Hawaiian Samoan Chamorro
- Tongan Fijian Marshallese

Enter, for example, Palauan, Tahitian, Chuukese, etc.

PRINCIPLES OF GOOD OPERATIONALIZATION

- 1 Unambiguous
- 2 Parsimonious
- 3 Accurate
- 4 Reliable



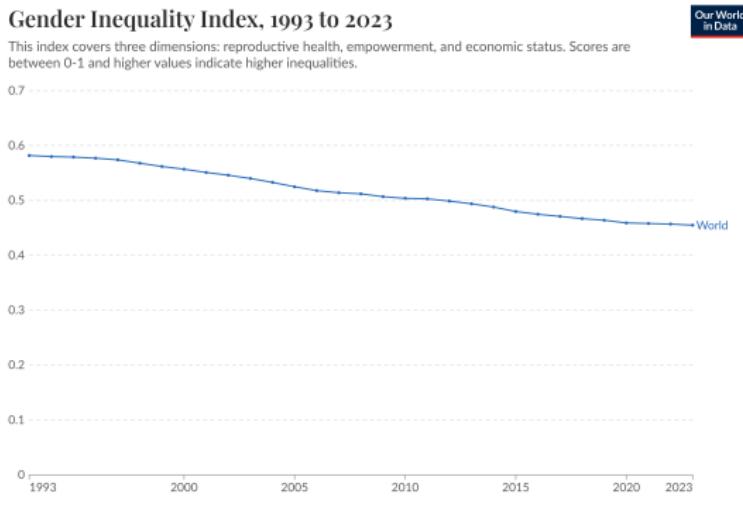
PRINCIPLES OF GOOD OPERATIONALIZATION

- ① Unambiguous
 - ② Parsimonious
 - ③ Accurate
 - ④ Reliable
 - ⑤ Feasible

PRINCIPLES OF GOOD OPERATIONALIZATION

Index (or Scale) Variable

An index (or scale) variable is a type of proxy measure in which the researcher combines multiple sources of data into a single representation of the concept under investigation using a pre-defined mathematical operation.



EXERCISE: OPERATIONALIZATION

- ① You want to measure how happy people are
 - ② You want to measure people's driving ability
 - ③ You want to measure the political ideology of a member of Congress

HOW DO I COLLECT MY DATA?

Data Generating Process

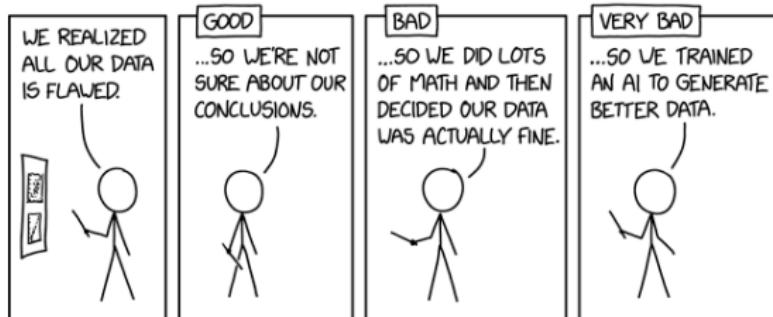
The rules and procedures that produce the data one is interested in

HOW DO I COLLECT MY DATA?

Data Generating Process

The rules and procedures that produce the data one is interested in

- No amount of statistical wizardry can compensate for bad data



HOW DO I COLLECT MY DATA?

Data Generating Process

The rules and procedures that produce the data one is interested in

- No amount of statistical wizardry can compensate for bad data
- The gold standard of data generating processes is the **random sample**

SAMPLING

- The group we are interested in studying is known as the **population**

SAMPLING

- The group we are interested in studying is known as the **population**
- Often, we are not able to count every unit in the population, so we take a **sample**

SAMPLING

- The group we are interested in studying is known as the **population**
- Often, we are not able to count every unit in the population, so we take a **sample**
- The key to any data analysis project is a quality sample, which is determined by two elements:

SAMPLING

- The group we are interested in studying is known as the **population**
- Often, we are not able to count every unit in the population, so we take a **sample**
- The key to any data analysis project is a quality sample, which is determined by two elements:
 - ① A **random sample** of the population

SAMPLING

- The group we are interested in studying is known as the **population**
- Often, we are not able to count every unit in the population, so we take a **sample**
- The key to any data analysis project is a quality sample, which is determined by two elements:
 - 1 A **random sample** of the population
 - 2 The **sample size** is sufficiently large

RANDOM SAMPLE

Definition

The probability of any given unit being drawn from the population is uniform (the same)

RANDOM SAMPLE

Definition

The probability of any given unit being drawn from the population is uniform (the same)

- A failure of each unit to have a uniform probability of being drawn from the population is known as **selection bias**

Definition

The probability of any given unit being drawn from the population is uniform (the same)

- A failure of each unit to have a uniform probability of being drawn from the population is known as **selection bias**
- Units are “selecting” into our data because they are more observable than other units

RANDOM SAMPLE

Definition

The probability of any given unit being drawn from the population is uniform (the same)

- A failure of each unit to have a uniform probability of being drawn from the population is known as **selection bias**
- Units are “selecting” into our data because they are more observable than other units
- Selection bias reduces our **generalizability** to the population because the data is not representative of the population

RANDOM SAMPLE

Can I randomly sample 10 students from this class to generalize to:

- the population of GW students?

RANDOM SAMPLE

Can I randomly sample 10 students from this class to generalize to:

- the population of GW students?
- the population of SMPA students?

RANDOM SAMPLE

Can I randomly sample 10 students from this class to generalize to:

- the population of GW students?
- the population of SMPA students?
- the population of Data Analysis students?

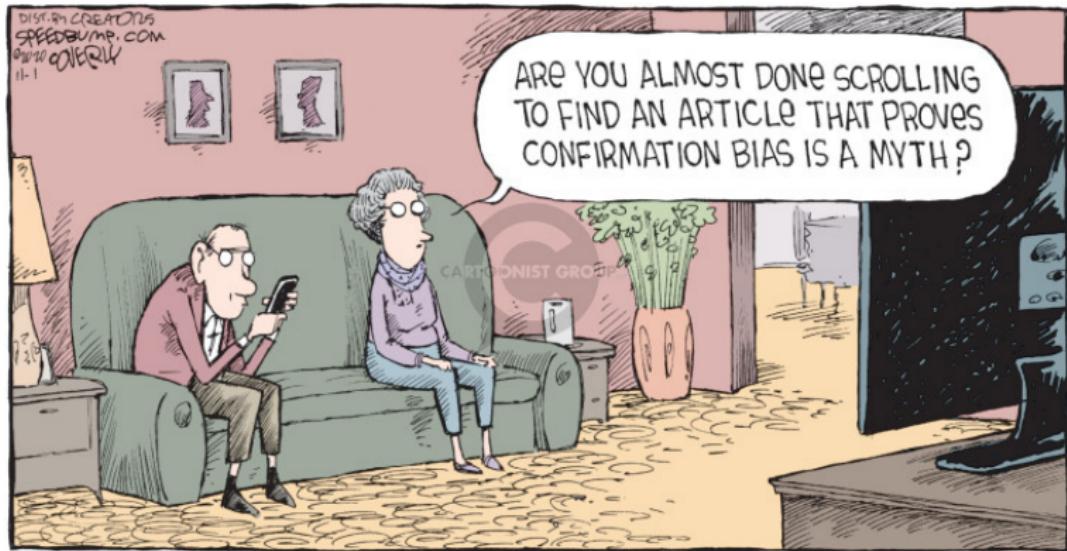
EXERCISE: SELECTION BIAS

BIASES IN RESEARCH

- ① Confirmation bias
- ② Desirability bias
- ③ Authority bias
- ④ Availability bias
- ⑤ Certainty bias

CONFIRMATION BIAS

We privilege evidence that supports our existing beliefs and discount evidence that challenges those beliefs.



©Dave Coverly. All rights reserved.

DESIRABILITY BIAS

We prefer evidence that supports a conclusion we want to be true and discount evidence that undermines that conclusion.

Pandemic in Retreat

And what else you need to know today.



By David Leonhardt
Feb. 11, 2021

THE MORNING NEWSLETTER

Covid, in Retreat

New cases in the U.S. have fallen by more than a third in the past month.



A medic cleans a stretcher in New York after transferring a Covid-19 patient last month. John Minchillo/Getty Images

By David Leonhardt
Oct. 4, 2021

David Leonhardt @DLeonhardt

Covid may now be in permanent retreat in the U.S.

It is not over, but after more than a year of death, sickness, grieving and isolation, you're allowed to feel joyful about the progress.
nytms.com/2021/05/21/bri...

Change in Daily U.S. Covid-19 Cases and Deaths Since Jan. 1

50%

THE MORNING NEWSLETTER

Omicron Is in Retreat

What's next?



Covid patients at a Brooklyn hospital last week. Wong J. Duke for The New York Times

By David Leonhardt
Jan. 19, 2022



AUTHORITY BIAS

We give greater weight to evidence offered by people in positions of authority.



AVAILABILITY BIAS

We give greater weight to evidence that is most memorable.



"THEY MUST HAVE A DEATH WISH TO SWIM IN THAT WATER."

CERTAINTY BIAS

We over- and under-state probabilistic evidence.



We're forecasting the election with three models

- Polls-plus forecast
What polls, the economy and historical data tell us about Nov. 8
- Polls-only forecast
What polls alone tell us about Nov. 8
- Now-cast
Who would win the election if it were held today

 National overview

President
Updated Nov. 8, 2016

Senate
Updated Nov. 8

Who will win the presidency?



Chance of winning



RETRO REPORT (2021) - WHAT'S IN A NUMBER?

What's in a Number? Some Research Shows That a Lower B.M.I. Isn't Always Better. | Retro Report

Katherine Flegal
Former Senior Scientist, CDC

MORE VIDEOS

1:10 / 10:42 • Obesity

RETRORREPORT SUBSCRIBE

CC YouTube

The video thumbnail shows a woman with long blonde hair, Katherine Flegal, sitting and gesturing with her hands while speaking. She is wearing a dark cardigan over a blue shirt. The background is a room with a window. The video title and her name are overlaid on the left side. The YouTube interface at the bottom includes a play button, volume icon, timestamp, and other video controls.

