

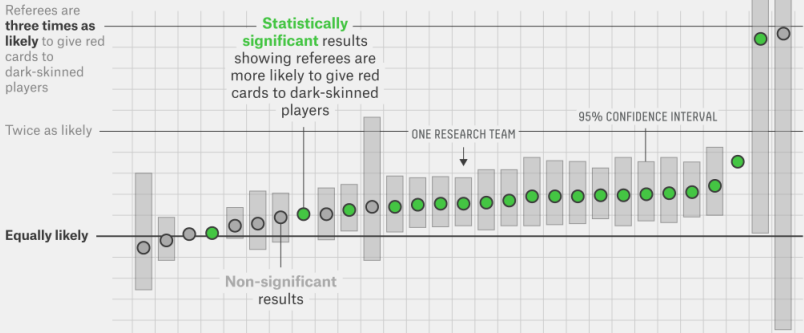
RESEARCHER CHOICES

Data Analysis for Journalism and Political Communication
(Spring 2026)

Prof. Bell

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.



FIVETHIRTYEIGHT

SOURCE: BRIAN NOSEK ET AL.

RESEARCHER CHOICES

- 1 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 \Rightarrow Hypothesis \Rightarrow Test \Rightarrow Analyze \Rightarrow 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 \Rightarrow Hypothesis \Rightarrow Test \Rightarrow Analyze \Rightarrow 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.)

ARE DEMOCRATS OR REPUBLICANS GOOD FOR THE ECONOMY?

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”

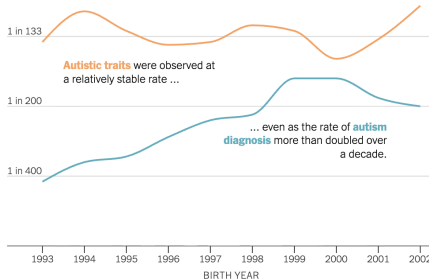
As more states in the US adopted the “pregnancy checkbox” on death certificates — which asked if the deceased had been pregnant or recently pregnant — the reported maternal mortality rate rose.

Maternal mortality rate, per 100,000 females



Source: WHO Mortality Database (2022). Adapted from KS Joseph et al. (2021) 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 Sakori 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. ↗

- | | | |
|--|-------------------------------------|--|
| <input type="checkbox"/> Asian Indian | <input type="checkbox"/> Japanese | <input type="checkbox"/> Native Hawaiian |
| <input type="checkbox"/> Chinese | <input type="checkbox"/> Korean | <input type="checkbox"/> Guamanian or Chamorro |
| <input type="checkbox"/> Filipino | <input type="checkbox"/> Vietnamese | <input type="checkbox"/> Samoan |
| <input type="checkbox"/> Other Asian — Print race, for example, Hmong, Laotian, Thai, Pakistani, Cambodian, and so on. ↗ | | |
| <input type="checkbox"/> 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.

<input type="checkbox"/> German	<input type="checkbox"/> Irish	<input type="checkbox"/> English
<input type="checkbox"/> Italian	<input type="checkbox"/> Polish	<input type="checkbox"/> French

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

☐ **HISPANIC OR LATINO** — Provide details below.

<input type="checkbox"/> Mexican or Mexican American	<input type="checkbox"/> Puerto Rican	<input type="checkbox"/> Cuban
<input type="checkbox"/> Salvadoran	<input type="checkbox"/> Dominican	<input type="checkbox"/> Colombian

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

☐ **BLACK OR AFRICAN AMERICAN** — Provide details below.

<input type="checkbox"/> African American	<input type="checkbox"/> Jamaican	<input type="checkbox"/> Haitian
<input type="checkbox"/> Nigerian	<input type="checkbox"/> Ethiopian	<input type="checkbox"/> Somali

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

☐ **ASIAN** — Provide details below.

<input type="checkbox"/> Chinese	<input type="checkbox"/> Filipino	<input type="checkbox"/> Asian Indian
<input type="checkbox"/> Vietnamese	<input type="checkbox"/> Korean	<input type="checkbox"/> Japanese

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

☐ **AMERICAN INDIAN OR ALASKA NATIVE** — Enter, for example, Navajo Nation, Blackfeet Tribe, Mayan, Aztec, Native Village of Barrow Inupiat Tribal Government, Iñupiat, etc.

☐ **MIDDLE EASTERN OR NORTH AFRICAN** — Provide details below.

<input type="checkbox"/> Lebanese	<input type="checkbox"/> Iranian	<input type="checkbox"/> Egyptian
<input type="checkbox"/> Syrian	<input type="checkbox"/> Moroccan	<input type="checkbox"/> Israeli

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

☐ **NATIVE HAWAIIAN OR PACIFIC ISLANDER** — Provide details below.

<input type="checkbox"/> Native Hawaiian	<input type="checkbox"/> Samoan	<input type="checkbox"/> Chamorro
<input type="checkbox"/> Tongan	<input type="checkbox"/> Fijian	<input type="checkbox"/> Marshallese

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

PRINCIPLES OF GOOD OPERATIONALIZATION

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

The screenshot shows a YouGov survey interface. At the top, the YouGov logo is in the upper left. The main content area displays the text "African American" in green, followed by "or" in black, and "Good" in green. Below this text is a square photograph of a Black man. At the bottom of the interface, there are three lines of text: "Press the I key for African American or Good", "Press the E key for anything else", and "Go as fast as you can".

Annotations with red lines pointing to the interface elements:

- Category groups**: Points to the text "African American or Good".
- Targets**: Points to the photograph of the Black man. Below the label, it says: "A photo of a black or a white man or a word synonymous with 'good' or 'bad'".
- Navigation Instructions**: Points to the three lines of text at the bottom: "Press the I key for African American or Good", "Press the E key for anything else", and "Go as fast as you can".

PRINCIPLES OF GOOD OPERATIONALIZATION

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

PRINCIPLES OF GOOD OPERATIONALIZATION

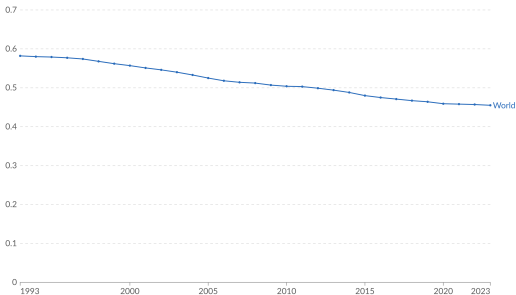
Index (or Scale) Variable

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.

Gender Inequality Index, 1993 to 2023

Our World
in Data

This index covers three dimensions: reproductive health, empowerment, and economic status. Scores are between 0-1 and higher values indicate higher inequalities.



Data source: UNDP, Human Development Report (2025)

OurWorldinData.org/economic-inequality-by-gender | CC BY

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

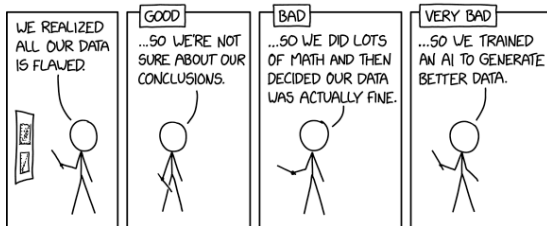
The rules and procedures that produce the data one is interested in

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**

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

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