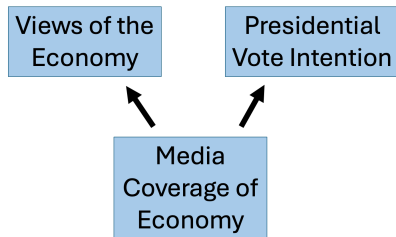


# ETHICS AND SAMPLING

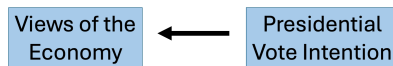
Data Analysis for Journalism and Political Communication  
(Fall 2025)

Prof. Bell

# CORRELATION VS. CAUSATION REVIEW



Confounding



Reverse Causation

# CORRELATION VS. CAUSATION REVIEW

The amount of money that a political candidate raises for their campaign increases the likelihood that they will win the election.

# CORRELATION VS. CAUSATION REVIEW

Living near electric power lines increases the likelihood of developing leukemia.

# CORRELATION VS. CAUSATION REVIEW

Jurisdictions that adopt ranked choice voting tend to elect more moderate candidates.

# CORRELATION VS. CAUSATION REVIEW

Taking a course for pass/fail credit typically decreases the amount of learning that a student obtains from that course.

# ETHICS REVIEW

## Respect for Persons

Individuals should be treated as autonomous agents, and persons with diminished autonomy are entitled to protection.

## Beneficence

(1) Do not harm and (2) maximize possible benefits and minimize possible harms.

## Justice

Groups who bear the burden of research should also be the beneficiaries of that research.

# ETHICS CASE STUDIES

- 1 Home DNA Testing
- 2 Crisis Text Line
- 3 Diversity in Faces (DiF) dataset



# ETHICS CASE STUDIES

- 1 What are the relevant ethical principles and practices in this case?
- 2 In what ways/why are there concerns about a violation of ethical principles in this case?
- 3 What are some ways that data could have been used more ethically in this case?

# ETHICS CASE STUDIES

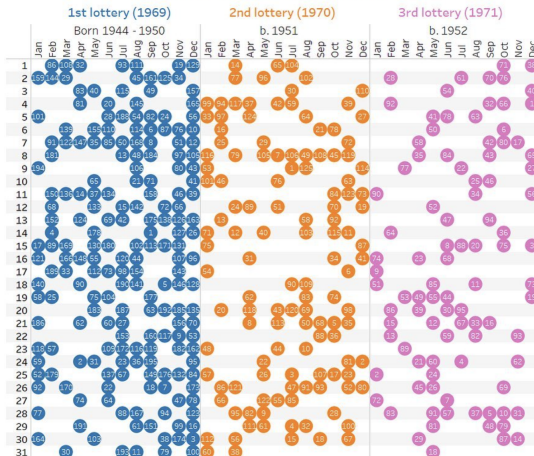
- 1 Home DNA Testing
- 2 Crisis Text Line
- 3 Diversity in Faces (DiF) dataset

# 1970 VIETNAM WAR DRAFT



# 1970 VIETNAM WAR DRAFT

Birthdates of US servicemen drafted into the Vietnam War as a result of birthdate lotteries held in 1969, 1970 and 1971



Source: [@visyual](#)

Note: The numbers denote the order that the birthdates were drawn, as this determined the order of call. The highest lottery number called for duty in the 1st, 2nd and 3rd lotteries was 195, 125 and 95, respectively.

# DEFINITIONS

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

# DEFINITIONS

- 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**

# DEFINITIONS

- 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**
- Our best guess about the population based on our sample is the **estimate**

# DEFINITIONS

- 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**
- Our best guess about the population based on our sample is the **estimate**
- The key to a good estimate is a quality sample, which is determined by two elements:
  - ① A **random sample** of the population
  - ② The **sample size** is sufficiently large



# In-class exercise

# SAMPLE SIZE

- How many units should you sample from the population?

# SAMPLE SIZE

- How many units should you sample from the population?  
**It depends on your desired level of certainty.**

# SAMPLE SIZE

- How many units should you sample from the population?  
**It depends on your desired level of certainty.**
- The most common level of certainty is 95% (the inverse of a **p-value** of .05, meaning that there is a 5% chance we are committing Type I error)

# SAMPLE SIZE

- How many units should you sample from the population?  
**It depends on your desired level of certainty.**
- The most common level of certainty is 95% (the inverse of a **p-value** of .05, meaning that there is a 5% chance we are committing Type I error)
- In other words, there is a 5% chance that the true population value is outside of the **confidence interval**

# SAMPLE SIZE

- How many units should you sample from the population?  
**It depends on your desired level of certainty.**
- The most common level of certainty is 95% (the inverse of a **p-value** of .05, meaning that there is a 5% chance we are committing Type I error)
- In other words, there is a 5% chance that the true population value is outside of the **confidence interval**
- If we re-sampled the population 100 times, 95 of our estimates would fall within the confidence interval (let's see this in action!)

# MARGIN OF ERROR

- The confidence interval for a proportion<sup>1</sup> is also called the **margin of error (MOE)**

---

<sup>1</sup>There is a different formula for continuous variables.

# MARGIN OF ERROR

- The confidence interval for a proportion<sup>1</sup> is also called the **margin of error (MOE)**
- The 95% MOE is calculated as:

$$1.96 * \sqrt{p * (1 - p) / n}$$

where  $p$  is the proportion and  $n$  is the sample size

---

<sup>1</sup>There is a different formula for continuous variables.



# MARGIN OF ERROR

- The confidence interval for a proportion<sup>1</sup> is also called the **margin of error (MOE)**
- The 95% MOE is calculated as:

$$1.96 * \sqrt{p * (1 - p) / n}$$

where  $p$  is the proportion and  $n$  is the sample size

- Typically, pollsters will use a proportion ( $p$ ) of .5 to calculate an MOE for the entire poll:

$$1.96 * \sqrt{.5 * (1 - .5) / 1000} = .031$$

---

<sup>1</sup>There is a different formula for continuous variables.

# MARGIN OF ERROR

- The confidence interval for a proportion<sup>1</sup> is also called the **margin of error (MOE)**
- The 95% MOE is calculated as:

$$1.96 * \sqrt{p * (1 - p) / n}$$

where  $p$  is the proportion and  $n$  is the sample size

- Typically, pollsters will use a proportion ( $p$ ) of .5 to calculate an MOE for the entire poll:

$$1.96 * \sqrt{.5 * (1 - .5) / 1000} = .031$$

- We report the estimate with the MOE, e.g., 45 +/- 3.1%.

---

<sup>1</sup>There is a different formula for continuous variables.

# SAMPLE SIZE AND THE MARGIN OF ERROR

- Mathematically, the MOE gets smaller as the sample size increases

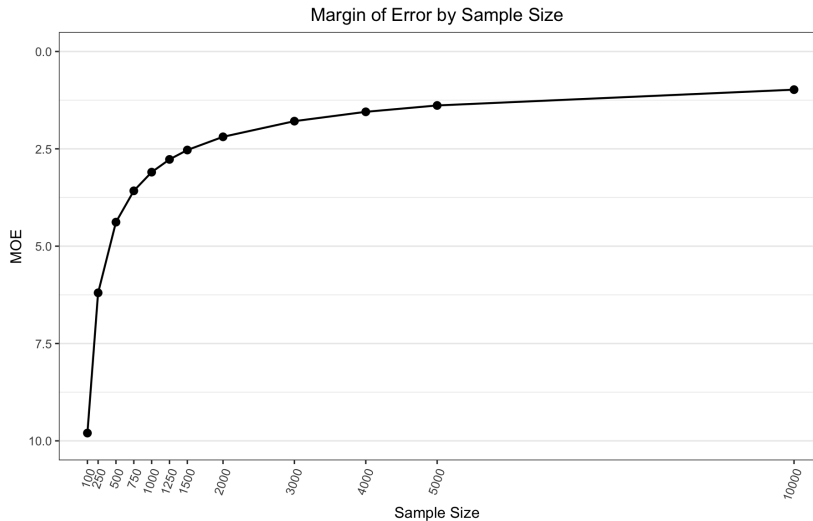
# SAMPLE SIZE AND THE MARGIN OF ERROR

- Mathematically, the MOE gets smaller as the sample size increases
- Intuitively, the larger our sample size, the less likely we are to draw an “outlier” sample and the more likely we are to get a representative sample

# SAMPLE SIZE AND THE MARGIN OF ERROR

- Mathematically, the MOE gets smaller as the sample size increases
- Intuitively, the larger our sample size, the less likely we are to draw an “outlier” sample and the more likely we are to get a representative sample
- But the marginal improvement in the MOE from adding units to the sample decreases as the sample size grows

# SAMPLE SIZE AND THE MARGIN OF ERROR



Note: 95% confidence level

# SAMPLE SIZE AND THE MARGIN OF ERROR

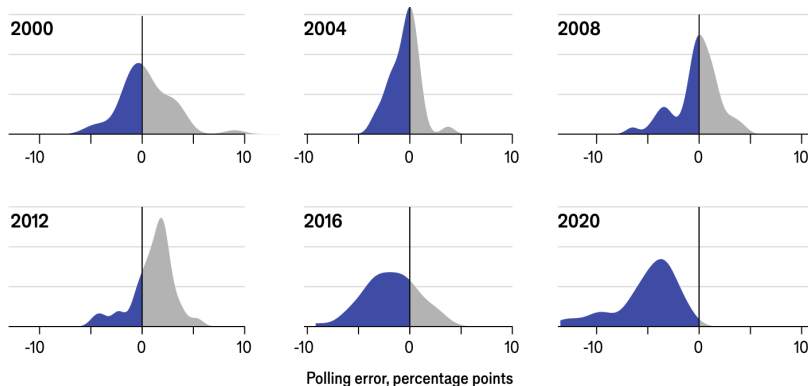
- Mathematically, the MOE gets smaller as the sample size increases
- Intuitively, the larger our sample size, the less likely we are to draw an “outlier” sample and the more likely we are to get a representative sample
- But the marginal improvement in the MOE from adding units to the sample decreases as the sample size grows
- Remember that the MOE only takes into account the sample size, not the potential for selection bias

# SAMPLE SIZE AND THE MARGIN OF ERROR

## Distribution of polling errors

Democratic share of the two-party vote in each state minus predicted share

Overestimated Democrats Underestimated Democrats



Source: The Economist