

# SAMPLING

Data Analysis for Journalism and Political Communication  
(Fall 2024)

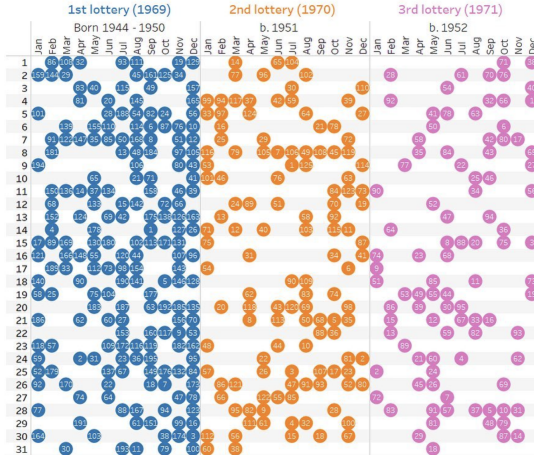
Prof. Bell

# 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: @visyuvai

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**

# SAMPLING

- There are many ways to derive an estimate from a sample, but recall that “garbage in = garbage out”: no amount of statistical wizardry can compensate for bad data

# SAMPLING

- There are many ways to derive an estimate from a sample, but recall that “garbage in = garbage out”: no amount of statistical wizardry can compensate for bad data
- The key to any data analysis project is a quality sample, which is determined by two elements:



# SAMPLING

- There are many ways to derive an estimate from a sample, but recall that “garbage in = garbage out”: no amount of statistical wizardry can compensate for bad data
- 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

- There are many ways to derive an estimate from a sample, but recall that “garbage in = garbage out”: no amount of statistical wizardry can compensate for bad data
- The key to any data analysis project is a quality sample, which is determined by two elements:
  - 1 A **random sample** of the population

## Definition

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

# SAMPLING

- There are many ways to derive an estimate from a sample, but recall that “garbage in = garbage out”: no amount of statistical wizardry can compensate for bad data
- 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

# 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

- There is a mathematical formula to estimate the confidence interval for a continuous mean, but more frequently we estimate a proportion

# MARGIN OF ERROR

- There is a mathematical formula to estimate the confidence interval for a continuous mean, but more frequently we estimate a proportion
- The confidence interval for a proportion is also called the **margin of error (MOE)**

# MARGIN OF ERROR

- There is a mathematical formula to estimate the confidence interval for a continuous mean, but more frequently we estimate a proportion
- The confidence interval for a proportion 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

# SAMPLE SIZE AND THE MARGIN OF ERROR

- E.g. Washington Post-University of Maryland poll, December 14-18, 2023

# SAMPLE SIZE AND THE MARGIN OF ERROR

- E.g. Washington Post-University of Maryland poll, December 14-18, 2023

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

# SAMPLE SIZE AND THE MARGIN OF ERROR

- E.g. Washington Post-University of Maryland poll, December 14-18, 2023

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

$$1.96 * \sqrt{.33 * (1 - .33) / 1024} = .029$$

# SAMPLE SIZE AND THE MARGIN OF ERROR

- E.g. Washington Post-University of Maryland poll, December 14-18, 2023

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

$$1.96 * \sqrt{.33 * (1 - .33) / 1024} = .029$$

- We report the estimate with the MOE, i.e., 33 +/- 2.9%.



# SAMPLE SIZE AND THE MARGIN OF ERROR

- E.g. Washington Post-University of Maryland poll, December 14-18, 2023

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

$$1.96 * \sqrt{.33 * (1 - .33) / 1024} = .029$$

- We report the estimate with the MOE, i.e., 33 +/- 2.9%.
- This means that there is a 5% chance that the true population proportion is outside of about 30-36%.

# SAMPLE SIZE AND THE MARGIN OF ERROR

- E.g. Washington Post-University of Maryland poll, December 14-18, 2023

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

$$1.96 * \sqrt{.33 * (1 - .33)/1024} = .029$$

- We report the estimate with the MOE, i.e., 33 +/- 2.9%.
- This means that there is a 5% chance that the true population proportion is outside of about 30-36%.  
*(Note: this does not mean that every value in the MOE is equally likely!)*

# SAMPLE SIZE AND THE MARGIN OF ERROR

- E.g. Washington Post-University of Maryland poll, December 14-18, 2023

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

$$1.96 * \sqrt{.33 * (1 - .33)/1024} = .029$$

- We report the estimate with the MOE, i.e., 33 +/- 2.9%.
- This means that there is a 5% chance that the true population proportion is outside of about 30-36%.  
*(Note: this does not mean that every value in the MOE is equally likely!)*
- Typically, pollsters will use a proportion ( $p$ ) of .5 to calculate an MOE for the entire poll, rather than individual questions

# SAMPLE SIZE AND THE MARGIN OF ERROR

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

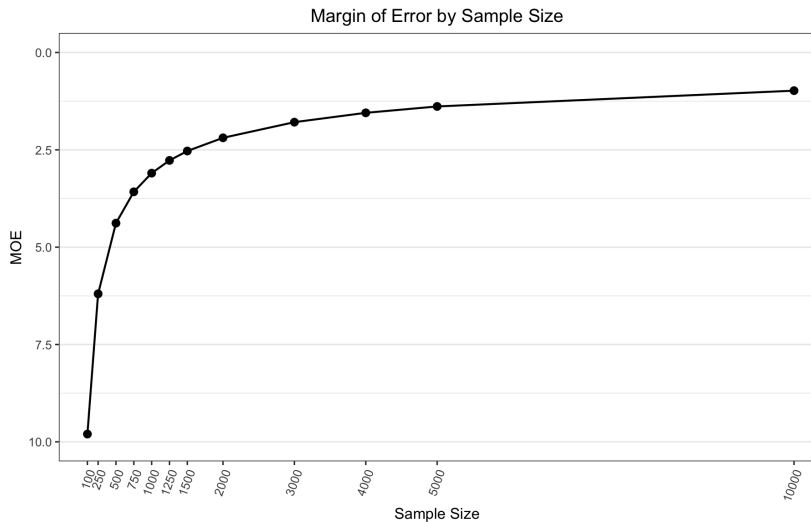
# SAMPLE SIZE AND THE MARGIN OF ERROR

- 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

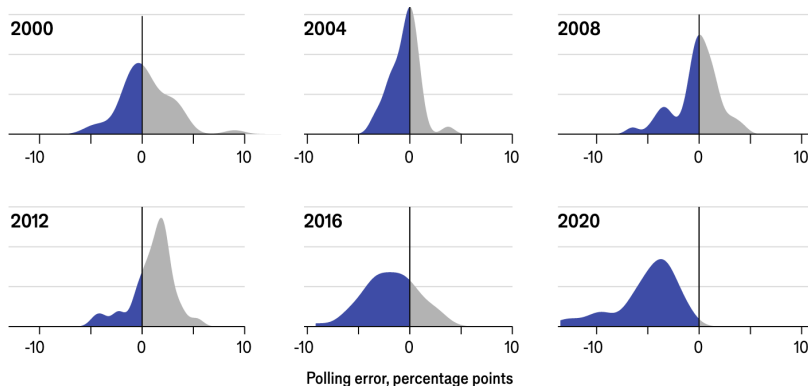


# SAMPLE SIZE AND THE MARGIN OF ERROR

## Distribution of polling errors

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

Overestimated Democrats Underestimated Democrats



Source: The Economist

# GETTING A RANDOM SAMPLE

How would you generate a random sample of American voters?

# GETTING A RANDOM SAMPLE

- **Systematic Sampling:** Select every  $n^{th}$  unit of a population (e.g. exit polls)

# GETTING A RANDOM SAMPLE

- **Systematic Sampling:** Select every  $n^{th}$  unit of a population (e.g. exit polls)
- **List-based Sampling:** Randomly select units from an existing list of the population (e.g. registered voter lists)

# GETTING A RANDOM SAMPLE

- **Systematic Sampling:** Select every  $n^{\text{th}}$  unit of a population (e.g. exit polls)
- **List-based Sampling:** Randomly select units from an existing list of the population (e.g. registered voter lists)
- **Address-based Sampling (ABS):** Randomly select households from a list of addresses provided by the U.S. Postal Service

# GETTING A RANDOM SAMPLE

- **Systematic Sampling:** Select every  $n^{\text{th}}$  unit of a population (e.g. exit polls)
- **List-based Sampling:** Randomly select units from an existing list of the population (e.g. registered voter lists)
- **Address-based Sampling (ABS):** Randomly select households from a list of addresses provided by the U.S. Postal Service
- **Random-digit Dialing (RDD):** Randomly select area codes, and then random digits are added to the end to create 10-digit phone numbers

# GETTING A RANDOM SAMPLE

- **Systematic Sampling:** Select every  $n^{\text{th}}$  unit of a population (e.g. exit polls)
- **List-based Sampling:** Randomly select units from an existing list of the population (e.g. registered voter lists)
- **Address-based Sampling (ABS):** Randomly select households from a list of addresses provided by the U.S. Postal Service
- **Random-digit Dialing (RDD):** Randomly select area codes, and then random digits are added to the end to create 10-digit phone numbers
- **Non-probability/Quota Sampling:** Pseudo-randomly selecting, from an opt-in pool of respondents, a sample that approximates the make-up of the general population