

# INTRODUCTION TO DATA ANALYSIS

Data Analysis for Journalism and Political Communication  
(Spring 2025)

Prof. Bell

# EXPECTATIONS

- **An A effort:** Starts homework early and asks questions/attends office hours; reviews the professor's feedback; reads assignment instructions and grading rubrics carefully; attends class; being perfect at coding is not an expectation of an A effort.
- **A B effort:** Completes all homework assignments; does not carefully read assignment instructions; uses outside resources to complete the homework; attends class; a B effort meets minimum expectations.
- **A C effort or lower:** Does not complete all homework assignments; does not carefully read assignment instructions; does not communicate with the professor; does not consistently attend class; I will reach out if you meet these criteria.

# WHEN YOU GET STUCK

- 1 Look back at the lecture materials. Everything required to complete the homework is in the lecture materials.
- 2 Google the error (this is different than Googling “how do I do X?”).
- 3 Email the professor. Be sure to include enough code in your email that I can recreate the problem.

# CHALABI: 3 WAYS TO SPOT A BAD STATISTIC



# CHALABI: 3 WAYS TO SPOT A BAD STATISTIC



- 1 Can you see uncertainty?
- 2 Can we look beyond the averages?
- 3 How was the data collected?

# CAN YOU SEE THE UNCERTAINTY?

- We rarely get a complete count of everything

# CAN YOU SEE THE UNCERTAINTY?

- We rarely get a complete count of everything
- How do we know that the **sample** we choose is a good representative of the whole **population**?

# CAN YOU SEE THE UNCERTAINTY?

- We rarely get a complete count of everything
- How do we know that the **sample** we choose is a good representative of the whole **population**?
- We will learn how to measure and communicate about uncertainty, which is both art and science



# CAN YOU SEE THE UNCERTAINTY?

- We rarely get a complete count of everything
- How do we know that the **sample** we choose is a good representative of the whole **population**?
- We will learn how to measure and communicate about uncertainty, which is both art and science
- Humans do not do well with probability

# CAN YOU SEE THE UNCERTAINTY?

- We rarely get a complete count of everything
- How do we know that the **sample** we choose is a good representative of the whole **population**?
- We will learn how to measure and communicate about uncertainty, which is both art and science
- Humans do not do well with probability

“There are only five probabilities the average human can handle: 99 percent, one percent, 100 percent, zero, and 50-50. That’s it.”

- Richard Thaler, Nobel Laureate in Economics

# CAN YOU SEE THE UNCERTAINTY?

- We rarely get a complete count of everything
- How do we know that the **sample** we choose is a good representative of the whole **population**?
- We will learn how to measure and communicate about uncertainty, which is both art and science
- Humans do not do well with probability



# CAN YOU SEE THE UNCERTAINTY?

- We rarely get a complete count of everything
- How do we know that the **sample** we choose is a good representative of the whole **population**?
- We will learn how to measure and communicate about uncertainty, which is both art and science
- Humans do not do well with probability

**FiveThirtyEight**  
2016 Election Forecast

President  
Updated Nov. 8, 2016

Senate  
Updated Nov. 8, 2016

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

## Who will win the presidency?

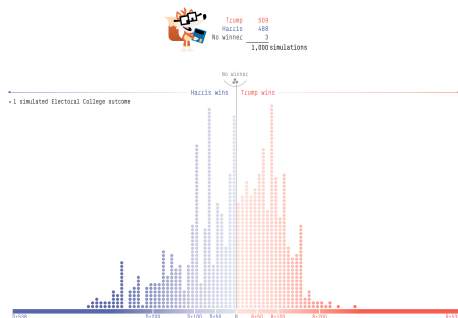


### Chance of winning



# CAN YOU SEE THE UNCERTAINTY?

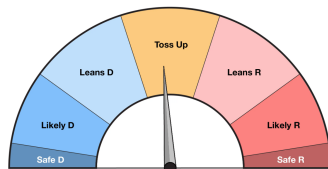
- We rarely get a complete count of everything
- How do we know that the **sample** we choose is a good representative of the whole **population**?
- We will learn how to measure and communicate about uncertainty, which is both art and science
- Humans do not do well with probability



# CAN YOU SEE THE UNCERTAINTY?

- We rarely get a complete count of everything
- How do we know that the **sample** we choose is a good representative of the whole **population**?
- We will learn how to measure and communicate about uncertainty, which is both art and science
- Humans do not do well with probability

## 2024 Presidential Forecast



# CAN WE LOOK BEYOND THE AVERAGES?

- There is always a trade-off between accuracy and simplicity when working with data

# CAN WE LOOK BEYOND THE AVERAGES?

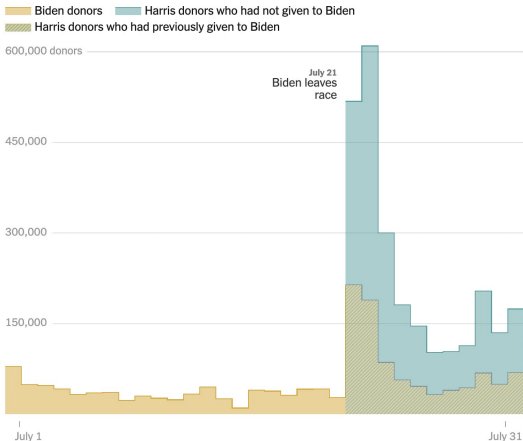
- There is always a trade-off between accuracy and simplicity when working with data
- We aggregate data to make it easier to comprehend, but we may also lose important context



# CAN WE LOOK BEYOND THE AVERAGES?

## Donors both old and new gave to the newly renamed Harris campaign

Both donors who had given to the Biden re-election campaign and new people who had not previously contributed rushed to donate to the Harris campaign.



Source: Federal Election Commission • The New York Times

# CAN WE LOOK BEYOND THE AVERAGES?

- There is always a trade-off between accuracy and simplicity when working with data
- We aggregate data to make it easier to comprehend, but we may also lose important context
- We will talk about using data visualization to communicate about data, as well as researcher choices and biases

# CAN WE LOOK BEYOND THE AVERAGES?

- There is always a trade-off between accuracy and simplicity when working with data
- We aggregate data to make it easier to comprehend, but we may also lose important context
- We will talk about using data visualization to communicate about data, as well as researcher choices and biases
- We will also talk about the importance of theory in understanding data, especially correlation vs. causation

## Storks Deliver Babies ( $p = 0.008$ )

### KEYWORDS:

Teaching;  
Correlation;  
Significance;  
 $p$ -values.

### Robert Matthews

Aston University, Birmingham, England.  
e-mail: rajm@compuserve.com

### Summary

This article shows that a highly statistically significant correlation exists between stork populations and human birth rates across Europe. While storks may not deliver babies, unthinking interpretation of correlation and  $p$ -values can certainly deliver unreliable conclusions.

### ◆ INTRODUCTION ◆

Introductory statistics textbooks routinely warn of the dangers of confusing correlation with causation, pointing out that while a high correlation coefficient is indicative of (linear) association, it cannot be taken as a measure of causation. Such

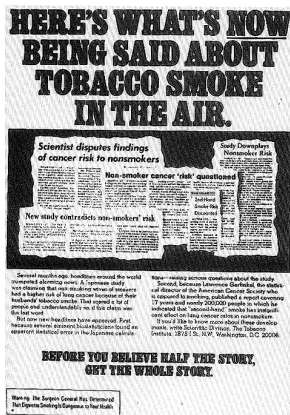
association between storks and the concept of women as bringers of life, and also in the bird's feeding habits, which were once regarded as a search for embryonic life in water (Cooper 1992). The legend lives on to this day, with neonate-bearing storks being a regular feature of greetings cards celebrating births.

# HOW WAS THE DATA COLLECTED?

- Data is not objective – it is generated by humans

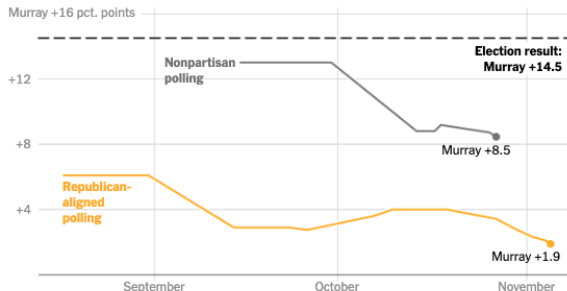
# HOW WAS THE DATA COLLECTED?

- Data is not objective – it is generated by humans
- Some data is produced by unscrupulous actors



# HOW WAS THE DATA COLLECTED?

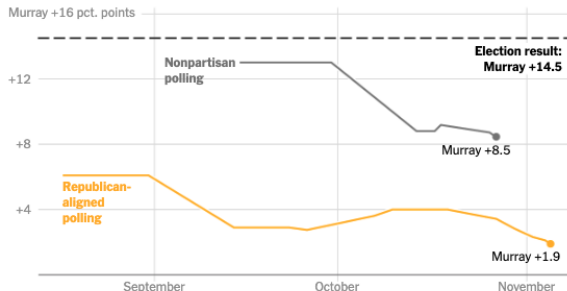
- Data is not objective – it is generated by humans
- Some data is produced by unscrupulous actors



Source: New York Times analysis of Washington Senate race polls aggregated by FiveThirtyEight • Notes: Trends are calculated with a 14-day average. Polling groups considered Republican-aligned include those identified by The New York Times and FiveThirtyEight. Polling groups considered nonpartisan are those not known to be aligned with or funded by a political party. • By Jason Kao

# HOW WAS THE DATA COLLECTED?

- Data is not objective – it is generated by humans
- Some data is produced by unscrupulous actors



Source: New York Times analysis of Washington Senate race polls aggregated by FiveThirtyEight. Notes: Trends are calculated with a 14-day average. Polling groups considered Republican-aligned include those identified by The New York Times and FiveThirtyEight. Polling groups considered nonpartisan are those not known to be aligned with or funded by a political party. By Jason Kao

- But most of the time, poor analysis is not nefarious – humans are imperfect

# HOW WAS THE DATA COLLECTED?

- Garbage in = garbage out: no amount of statistical wizardry can compensate for bad data



# HOW WAS THE DATA COLLECTED?

- Garbage in = garbage out: no amount of statistical wizardry can compensate for bad data
- We will spend a lot of time thinking about the **data generating process** and how it can bias our results

## HOW WAS THE DATA COLLECTED?

### Support for mass deportation varies depending on how the question is asked

A sample of different questions asked about deportation this year show significant variation in levels of support — sometimes, even within the same survey.

Pollster	Wording	Support deportations ▼	Do not support deportations
<b>CBS News/YouGov</b> <i>Registered voters, June 5-7</i>	Would you favor or oppose the U.S. government starting a new national program to deport all undocumented immigrants currently living in the U.S. illegally?	62%	38%
<b>Marquette Law School</b> <i>Registered voters, October 1-10</i>	Do you favor or oppose deporting immigrants who are living in the United States illegally back to their home countries?	58%	42%
<b>ABC/Ipsos</b> <i>US adults, October 4-8</i>	There are at least 11 million undocumented immigrants living in the United States. Would you support or oppose an effort by the federal government to deport all these undocumented immigrants and send them back to their home countries?	56%	43%
<b>CNN</b> <i>US adults, January 25-30</i>	If Donald Trump becomes president again, would you favor or oppose him trying to...detain and deport millions of undocumented immigrants?	48%	52%
<b>Gallup</b> <i>US adults, June 3-23</i>	Please tell me whether you strongly favor, favor, oppose, or strongly oppose each of the following proposals...Deporting all immigrants who are living in the United States illegally back to their home country	47%	51%
<b>Marquette Law School</b> <i>Registered voters, October 1-10</i>	Do you favor or oppose deporting immigrants who are living in the United States illegally back to their home countries even if they have lived here for a number of years, have jobs and no criminal record?	40%	60%
<b>Pew Research</b> <i>US adults, April 8-14</i>	Which comes closer to your view about how to handle undocumented immigrants who are now living in the U.S.? (They should not be allowed to stay in the country legally/There should be a way for them to stay in the country legally, if certain requirements are met)		
	If "not be allowed." "Do you think there should be a national law enforcement effort to deport all immigrants who are now living in the U.S. illegally?"	33%	67%
	("Support" percentage includes those who say there should be a national deportation effort, "do not support" includes all others)		

# HOW WAS THE DATA COLLECTED?

- Garbage in = garbage out: no amount of statistical wizardry can compensate for bad data
- We will spend a lot of time thinking about the **data generating process** and how it can bias our results
- We will also discuss our ethical responsibilities around data

# HOW WAS THE DATA COLLECTED?

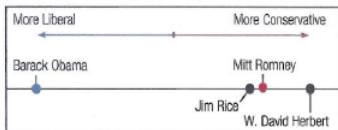


## 2014 Montana General Election Voter Information Guide

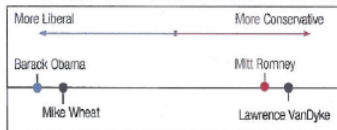
Election Date: November 4, 2014



### Nonpartisan Supreme Court Justice #1 Race



### Nonpartisan Supreme Court Justice #2 Race



For more information on how these figures were created, please see <http://data.stanford.edu/cime>. Please note that this guide is non-partisan and does not endorse any candidate or party. This guide was created as part of a joint research project at Stanford and Dartmouth.

Paid for by researchers at Stanford University and Dartmouth College, 616 Serra Street, Stanford, CA 94305

## Take this to the polls!

# GROUP DISCUSSION

Introduce yourself to your neighbor(s) and take a few minutes to review these additional graphs from Mona Chalabi. Do any of these stand out to you as being good (or bad) examples of our three questions for spotting a bad statistic?

- 1 Can you see uncertainty?
- 2 Can we look beyond the averages?
- 3 How was the data collected?

On your notecard, please write:

- 1 Preferred name
- 2 Preferred pronouns
- 3 Year in school and major
- 4 Your background in coding and/or statistics
- 5 One thing you hope to get out of this class