# Introduction to Data Analysis

Data Analysis for Journalism and Political Communication
(Fall 2024)

Prof. Bell

# 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
- Humans do not do well with probability

# CAN YOU SEE THE UNCERTAINTY?

- We rarely get a complete count of everything
- 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
- Humans do not do well with probability

# CAN YOU SEE THE UNCERTAINTY?

- We rarely get a complete count of everything
- Humans do not do well with probability

# CAN YOU SEE THE UNCERTAINTY?

- We rarely get a complete count of everything
- Humans do not do well with probability



- We will learn about how to measure and communicate about uncertainty

# CAN WE LOOK BEYOND THE AVERAGES?

- There is always a trade-off between accuracy and simplicity when working with data
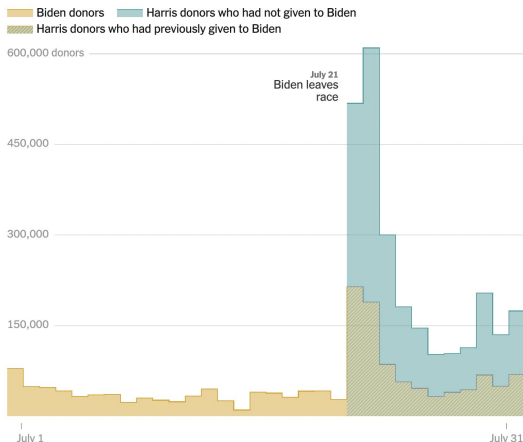
# CAN WE LOOK BEYOND THE AVERAGES?

- 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?

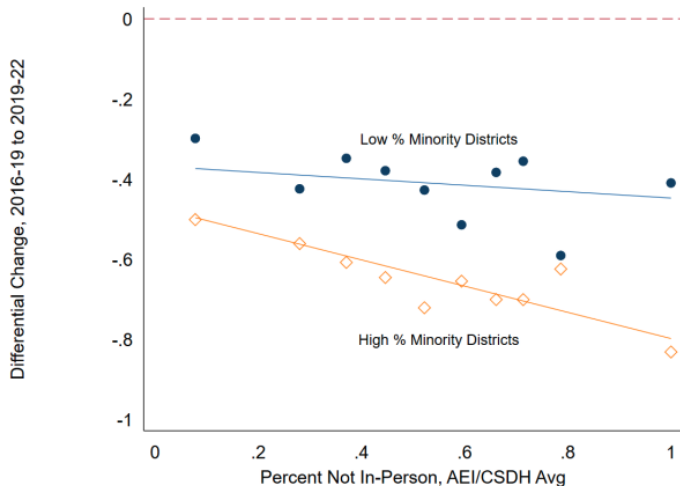**Donors both old and new gave to the newly renamed Harris campaign**

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.

Biden donors    Harris donors who had not given to Biden
Harris donors who had previously given to Biden

600,000 donors

July 21
Biden leaves
race

450,000

300,000

150,000

July 1

July 31

Source: Federal Election Commission  ·  The New York Times

# Can we look beyond the averages?



Figure 5. Math achievement losses vs percent not-in-person, by percent minority

Source: Fahle, et al. (2023). "School District and Community Factors Associated With Learning Loss During the COVID-19 Pandemic."

## 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 ($p = 0.008$)

*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 ◆

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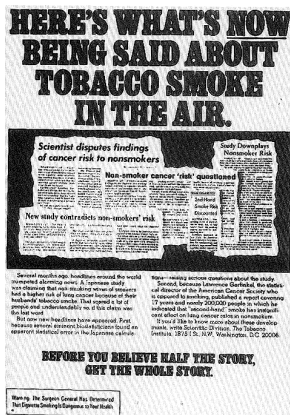
association between storks and the concept of women as bringers of life, and also in the bird's feeding habits, which were once regarded as a search for embryonic life in water (Cooper 1992). The legend lives on to this day, with neonate-bearing storks being a regular feature of greetings cards celebrating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

- Data is not objective – it is generated by humans

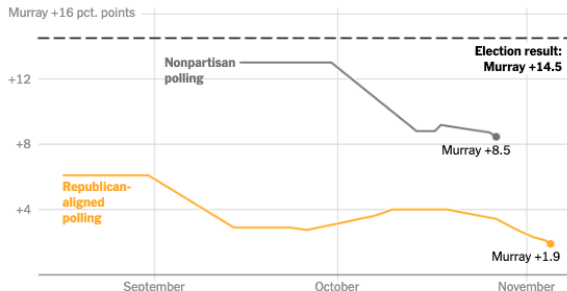- Data is not objective – it is generated by humans
- Some data is produced by unscrupulous actors

# How was the data collected?

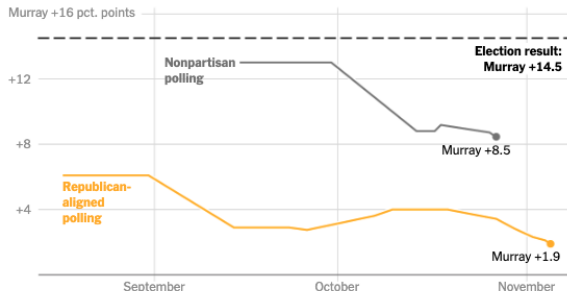- Data is not objective – it is generated by humans
- Some data is produced by unscrupulous actors



Murray +16 pct. points

Nonpartisan polling

Election result: Murray +14.5

+12

+8

Murray +8.5

Republican-aligned polling

+4

Murray +1.9

September          October          November

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
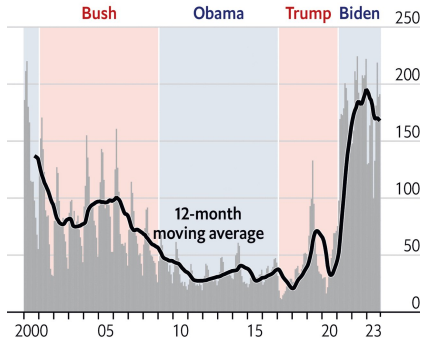
# HOW WAS THE DATA COLLECTED?

- 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?



**Illegal immigration** 1
Monthly encounters at the south-west
land border*, '000

Bush  Obama  Trump  Biden

12-month
moving average

*Only encounters between ports of entry. Since March
2020 monthly totals include apprehensions & expulsions.
Prior totals include apprehensions only
Source: US Customs and Border Protection
The Economist

- 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
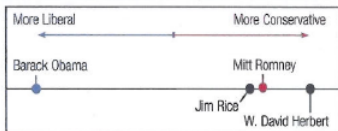
# How was the data collected?

# 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