# DS-UA 112
# Introduction to Data Science

Week 3: Lecture 1

Tables – Arranging Data in Rows and Columns

How can tables help us to summarize data?

# DS-UA 112
# Introduction to Data Science

Week 3: Lecture 1

Tables – Arranging Data in Rows and Columns

*Adapted from Nolan, Hug, and Salganik*

# Announcements

- ▶ Please check Week 3 agenda on NYU Classes
  - ▶ Homework 1
  - ▶ Lab 3
  - ▶ Grader Office Hours
- ▶ Remember to post to Piazza

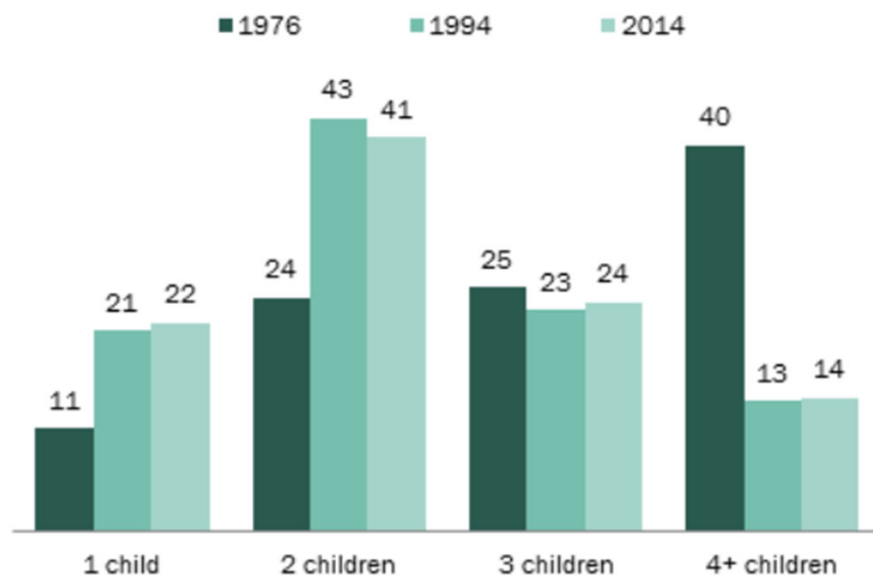Check the Calendar linked to NYU Classes for important dates

# Review

## Among Mothers, Family Size is Shrinking

*% of mothers ages 40 to 44 with...*

Legend: 1976, 1994, 2014

| | 1 child | 2 children | 3 children | 4+ children |
|---|---|---|---|---|
| 1976 | 11 | 24 | 25 | 40 |
| 1994 | 21 | 43 | 23 | 13 |
| 2014 | 22 | 41 | 24 | 14 |

## Moms with Less Education Have Bigger Families

*% of mothers ages 40 to 44 with ...*

Legend: 1 child, 2 children, 3 children, 4+children

| | Less than high school graduate | High school graduate/Some college | Bachelor's degree | Postgraduate degree |
|---|---|---|---|---|
| 1 child | 13 | 23 | 22 | 23 |
| 2 children | 32 | 38 | 46 | 50 |
| 3 children | 29 | 24 | 22 | 19 |
| 4+ children | 26 | 14 | 10 | 8 |

4

Base Rate Fallacy

All children

Firstborns
F

Students at
college A

Children of mothers with at least
one child at college A

Firstborns
F

college A
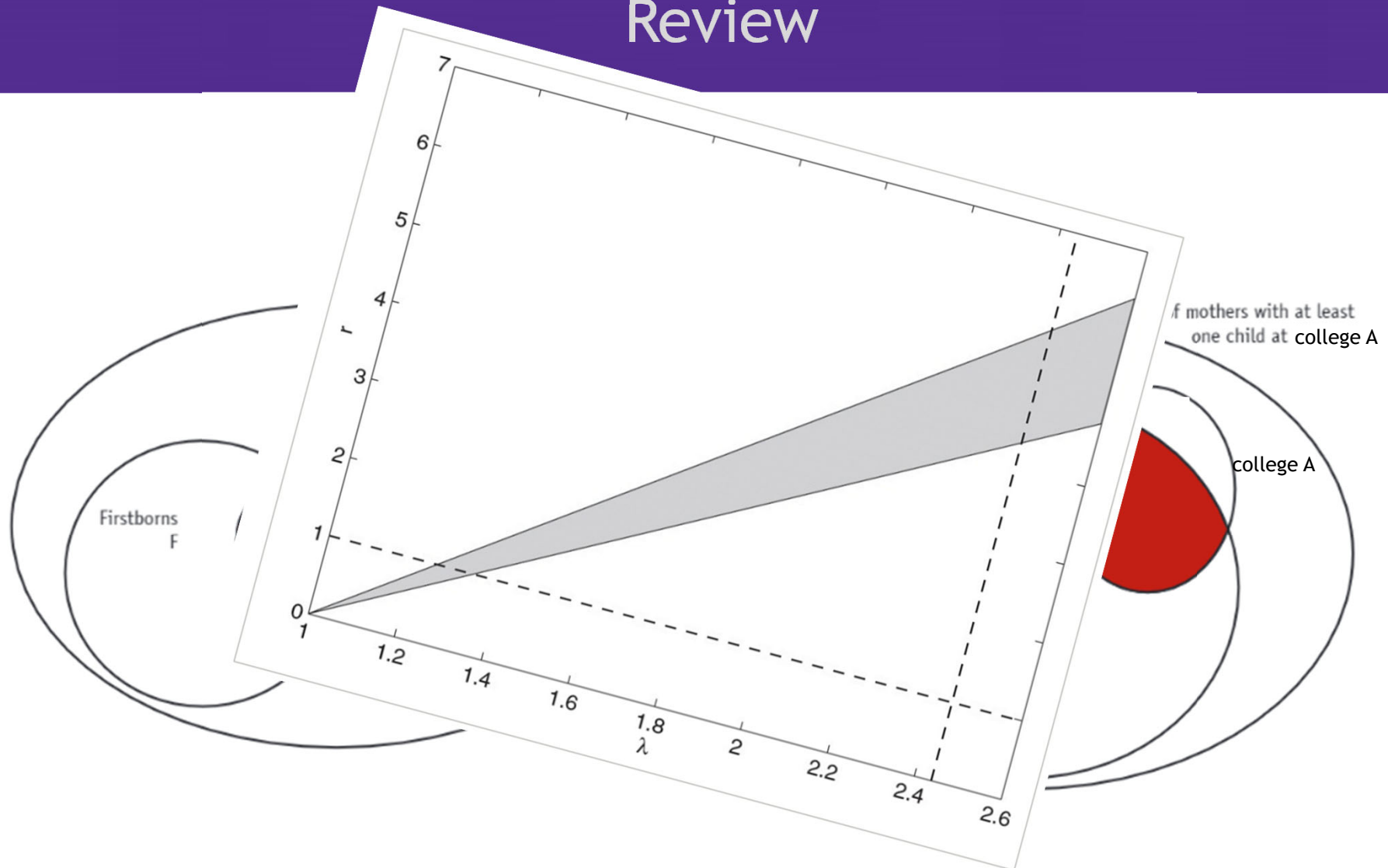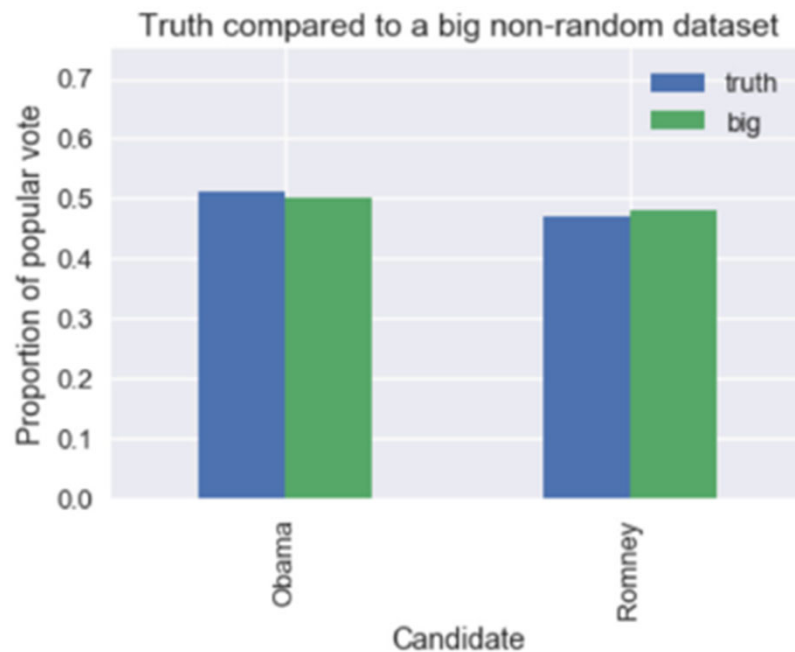
# Agenda

- ▶ Probability
  - ▶ Addition, Multiplication, Complement Rules
  - ▶ Summarize with average value
- ▶ Confounded Data
  - ▶ Adjusting for Bias
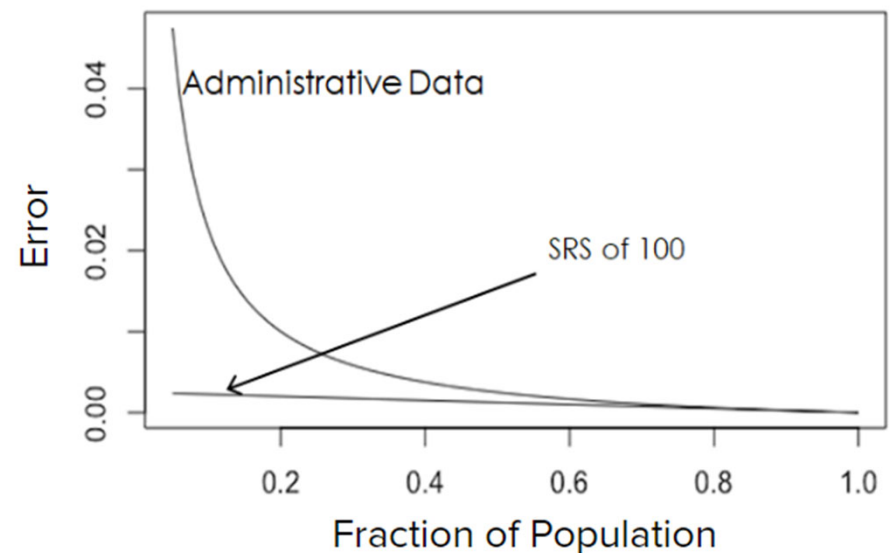- ▶ Messy Data
  - ▶ Arranging into Rows and Columns

**References**
- ▶ Nolan, Lau, Gonzalez (Chapter 2, 3.1)
  - ▶ https://cp71.github.io/textbook
- ▶ Salganik (Chapter 3)

## Truth compared to a big non-random dataset



If a dataset is not representative, then it may or may not be suitable for a study. Sometimes it causes bias in the analysis.

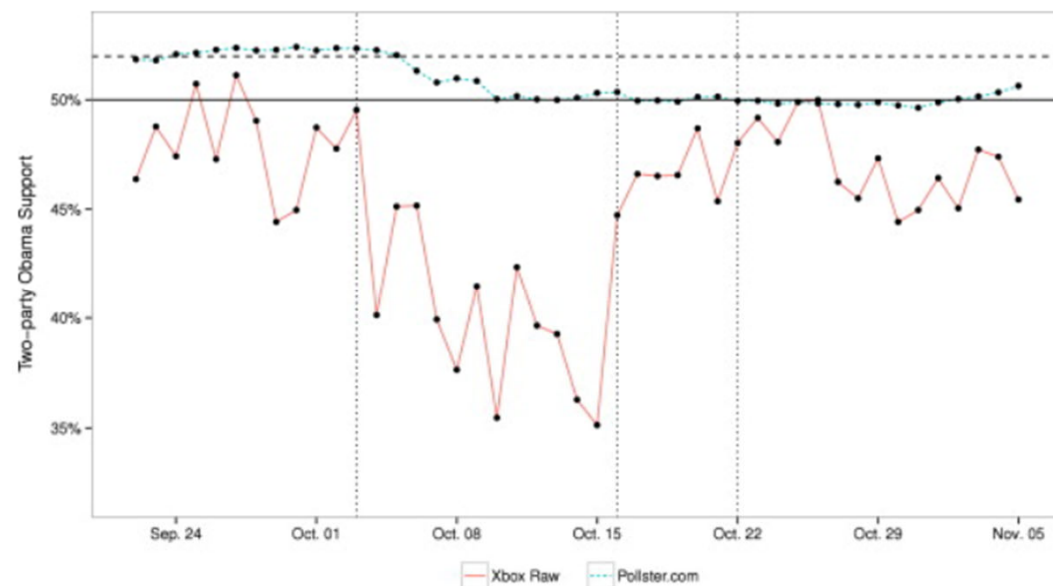Approach to data collection could indicate bias to us.



8

# Not Representative

Sampling Frame may not lie in the Population

## Forecasting elections with non-representative polls

Wei Wang [a] [⊠], David Rothschild [b] [⊠], Sharad Goel [b] [⊠], Andrew Gelman [a, c] [⊠]

## Abstract

Election forecasts have traditionally been based on representative polls, in which randomly sampled individuals are asked who they intend to vote for. While representative polling has historically proven to be quite effective, it comes at considerable costs of time and money. Moreover, as response rates have declined over the past several decades, the statistical benefits of representative sampling have diminished. In this paper, we show that, with proper statistical adjustment, non-representative polls can be used to generate accurate election forecasts, and that this can often be achieved faster and at a lesser expense than traditional survey methods. We demonstrate this approach by creating forecasts from a novel and highly non-representative survey dataset: a series of daily voter intention polls for the 2012 presidential election conducted on the Xbox gaming platform. After

# Not Representative

Adjustments made by stratifying following the data collection

# Forecasting elections with non-representative polls

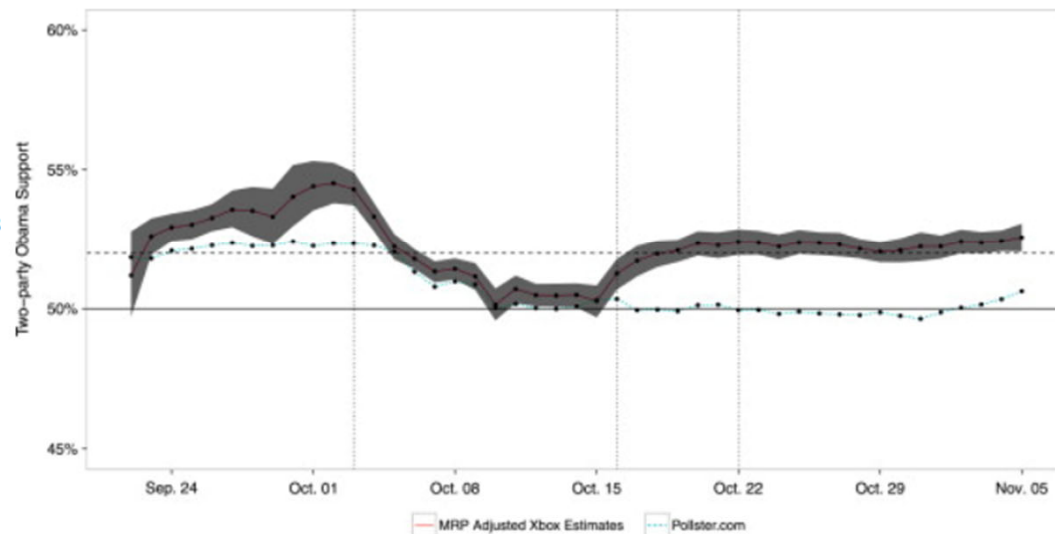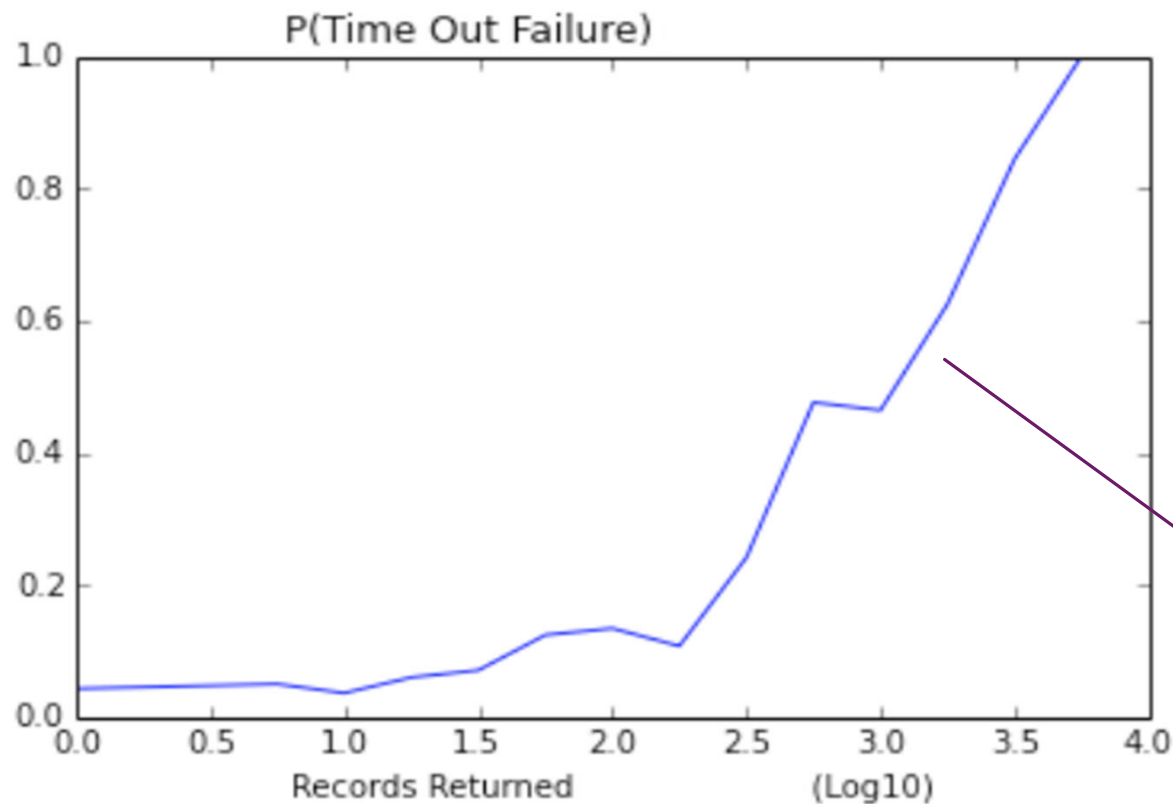Wei Wang [a], David Rothschild [b], Sharad Goel [b], Andrew Gelman [a, c]

## Abstract

Election forecasts have traditionally been based on representative polls, in which randomly sampled individuals are asked who they intend to vote for. While representative polling has historically proven to be quite effective, it comes at considerable costs of time and money. Moreover, as response rates have declined over the past several decades, the statistical benefits of representative sampling have diminished. In this paper, we show that, with proper statistical adjustment, non-representative polls can be used to generate accurate election forecasts, and that this can often be achieved faster and at a lesser expense than traditional survey methods. We demonstrate this approach by creating forecasts from a novel and highly non-representative survey dataset: a series of daily voter intention polls for the 2012 presidential election conducted on the Xbox gaming platform. After

# Confounded Data
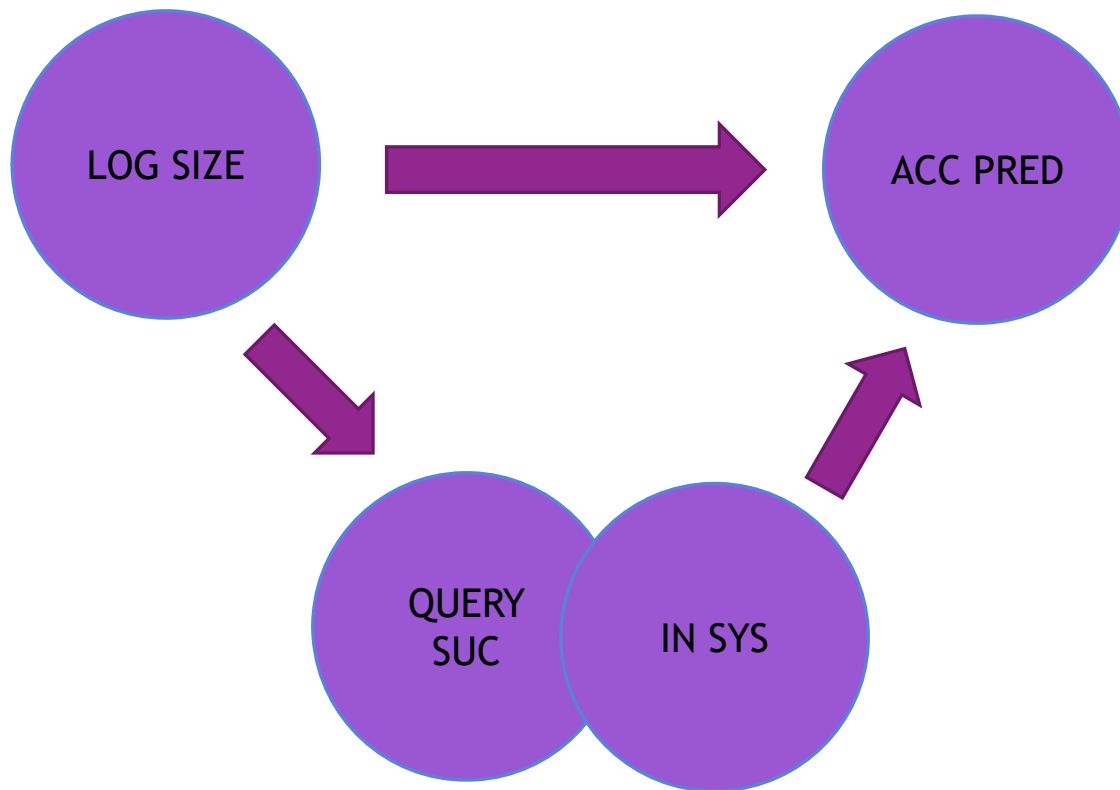
P(Time Out Failure)



Records Returned                    (Log10)

Here the probability of inclusion of the record in the sample changes depending on the size of the record

How could this lead to bias? In particular, why might predictions reflect habits of new customers

# Confounded Data



ACC PRED
Accurate Prediction of Least Popular
LOG SIZE
Size of Database Record
QUERY SUC
Whether Query Returned Successfully
IN SYS
Whether Record in the System

# Confounded Data

How to adjust the probabilities?

P(ACC PRED | LOG SIZE) =

P(ACC PRED | LOG SIZE, QUERY SUC) x P(QUERY SUC | LOG SIZE)

+

P(ACC PRED | LOG SIZE,NOT QUERY SUC) x P(NOT QUERY SUC | LOG SIZE) =

How to adjust the probabilities?

$$P(ACC\ PRED \mid LOG\ SIZE) =$$

$$P(ACC\ PRED \mid LOG\ SIZE,\ QUERY\ SUC) \times P(QUERY\ SUC \mid LOG\ SIZE)$$

$$+$$

$$P(ACC\ PRED \mid LOG\ SIZE,\ NOT\ IN\ SYS) \times P(NOT\ QUERY\ SUC \mid LOG\ SIZE) =$$

How to adjust the probabilities?

P(ACC PRED | LOG SIZE)

=

P(ACC PRED | LOG SIZE, QUERY SUC) x P(QUERY SUC | LOG SIZE)

+

P(ACC PRED | LOG SIZE, NOT IN SYS) x P(NOT QUERY SUC | LOG SIZE)

=

P(ACC PRED | LOG SIZE, QUERY SUC)P(QUERY SUC | LOG SIZE)

+

(0) P(NOT QUERY SUC | LOG SIZE)

How to adjust the probabilities?

P(ACC PRED | LOG SIZE)

=

P(ACC PRED | LOG SIZE, QUERY SUC) x P(QUERY SUC | LOG SIZE)

Without this quantity the two sides would not be equal.

How to adjust the probabilities?

$$P(\text{ACC PRED} \mid \text{LOG SIZE})$$

$$=$$

$$P(\text{ACC PRED} \mid \text{LOG SIZE}, \text{QUERY SUC}) \times P(\text{QUERY SUC} \mid \text{LOG SIZE})$$

$$\frac{P(\text{ACC PRED} \mid \text{LOG SIZE})}{P(\text{QUERY SUC} \mid \text{LOG SIZE})}$$

$$=$$

$$P(\text{ACC PRED} \mid \text{LOG SIZE}, \text{QUERY SUC})$$

Without this quantity the two sides would not be equal.

17

# Confounded Data

## Gender Achievement Gaps in U.S. School Districts

**Author/s:** Sean F. Reardon , Erin Fahle , Demetra Kalogrides , Anne Podolsky , Rosalía C. Zárate
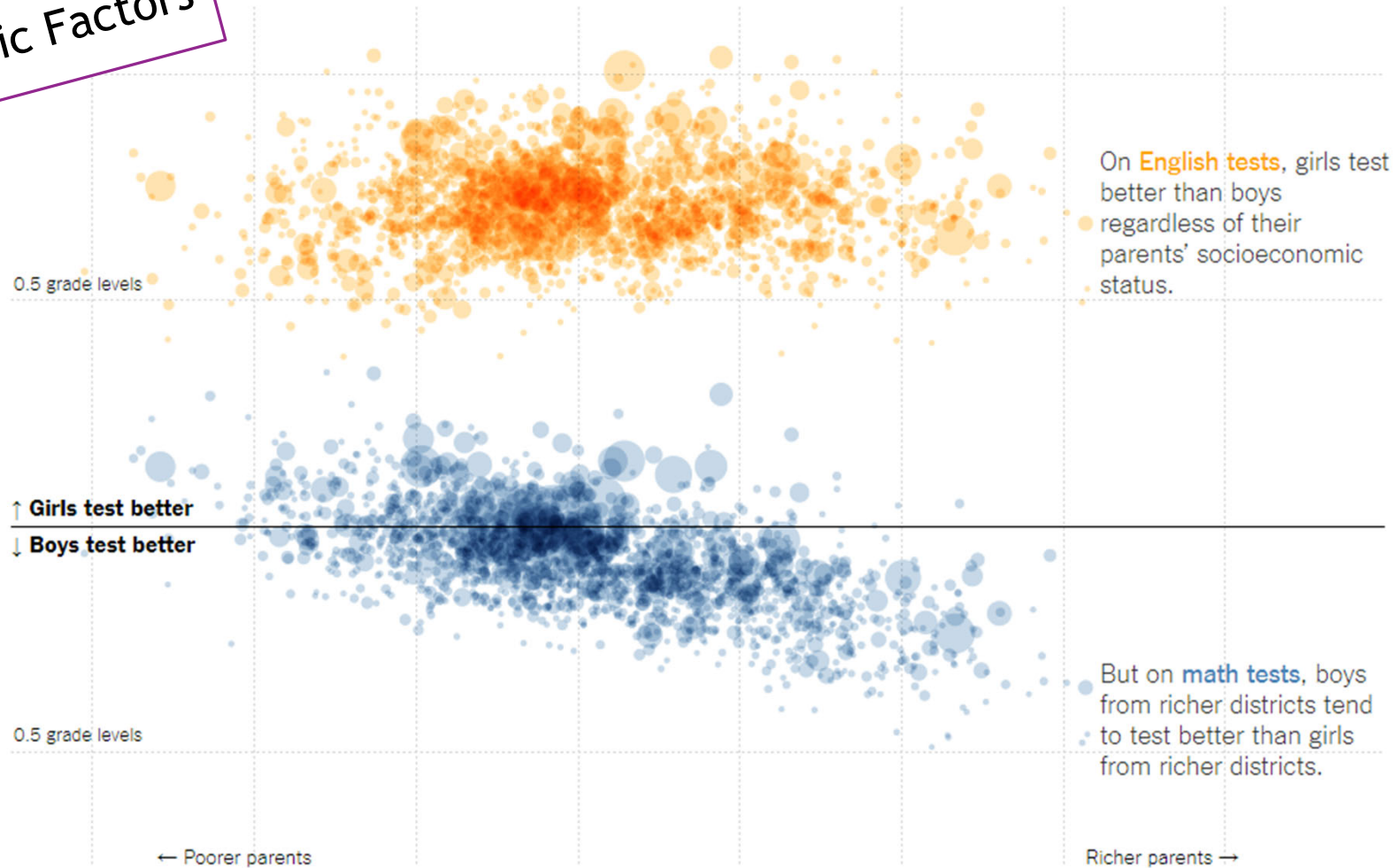
**Year of Publication:** 2018

In the first systematic study of gender achievement gaps in U.S. school districts, we estimate male-female test score gaps in math and English Language Arts (ELA) for nearly 10,000 school districts in the U.S. We use state

▶ What is the population?

▶ What is the question under study?

▶ What is the sampling frame?

▶ What could lead to confounded data?

18

Economic Factors



On **English tests**, girls test better than boys regardless of their parents' socioeconomic status.

0.5 grade levels

↑ **Girls test better**

↓ **Boys test better**

0.5 grade levels

But on **math tests**, boys from richer districts tend to test better than girls from richer districts.

← Poorer parents

Richer parents →
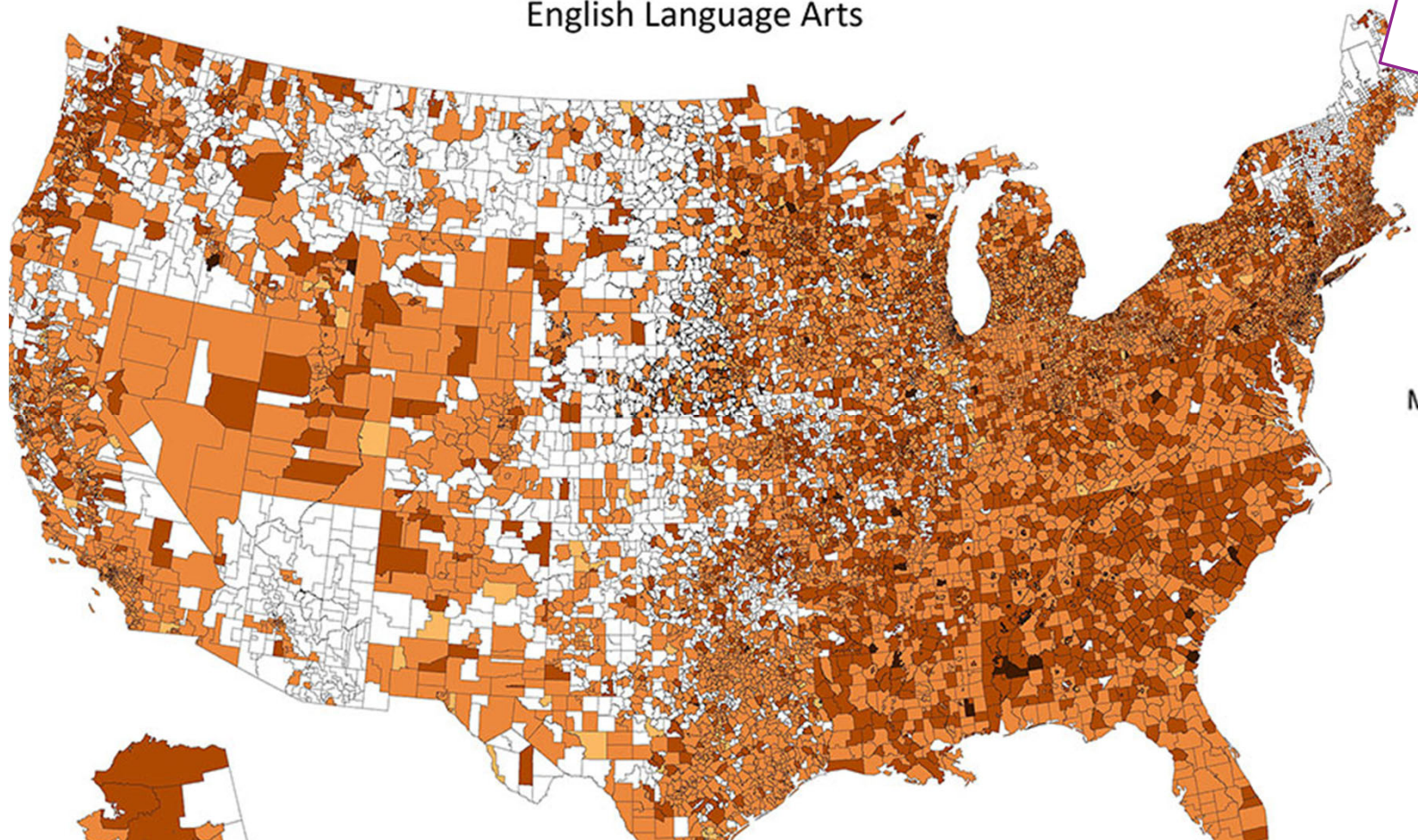
19

English Language Arts

Geographic Factors

Male-Female Gap, NS Scale
- More than 0.25 SDs
- 0.15 to 0.25 SDs
- 0.05 to 0.15 SDs
- -0.05 to 0.05 SDs
- -0.15 to -0.05 SDs
- -0.25 to -0.15 SDs
- -0.35 to -0.25 SDs
- Less than -0.35 SDs
- missing

20

# Questions

- ▶ Questions on Piazza?
- ▶ Question for You!

Should the word data be understood as singular or plural?

In Latin, data is the plural of datum and, historically and in specialized scientific fields , it is also treated as a plural in English, taking a plural verb, as in the data were collected and classified . In modern non-scientific use, however , despite the complaints of traditionalists, it is often not treated as a plural. Instead, it is treated as a mass noun, similar to a word like information, which cannot normally have a plural and which takes a singular verb. Sentences such as data was (as well as data were ) collected over a number of years are now widely accepted in standard English.

# Questions

▶ Questions on Piazza?

▶ Question for You!

Should the word data be understood as singular or plural?



In Lati... historically and in specialized
scienti... ...taking a plural
verb, ...scientific
use, ...
trea... ...a
lik... ...ed
si...
o...