

Unit 1: Introduction to Data

1. Data: Where it comes from,
and why that matters

9/28/2016

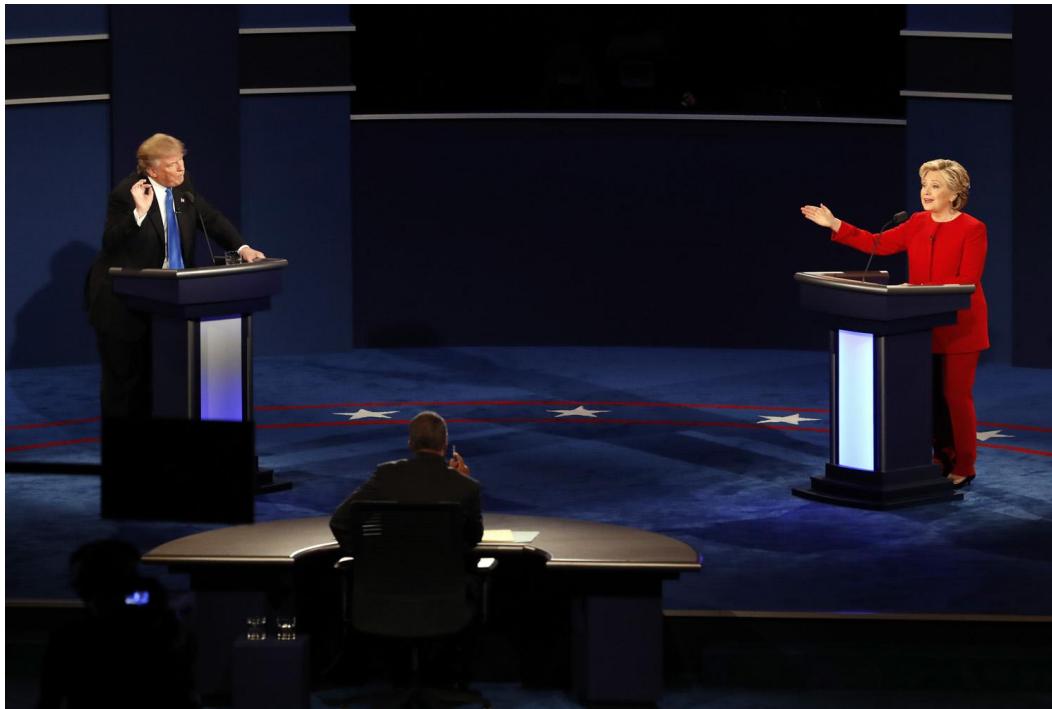
Sampling: The bridge between data and analysis



Key ideas

1. Using samples to make inferences about populations
2. The way you sample your data can change your inferences about the population
3. Experiments use random assignment to treatment groups,
observational studies do not
4. Random samples help with **generalizability**,
random assignment helps with **causality**

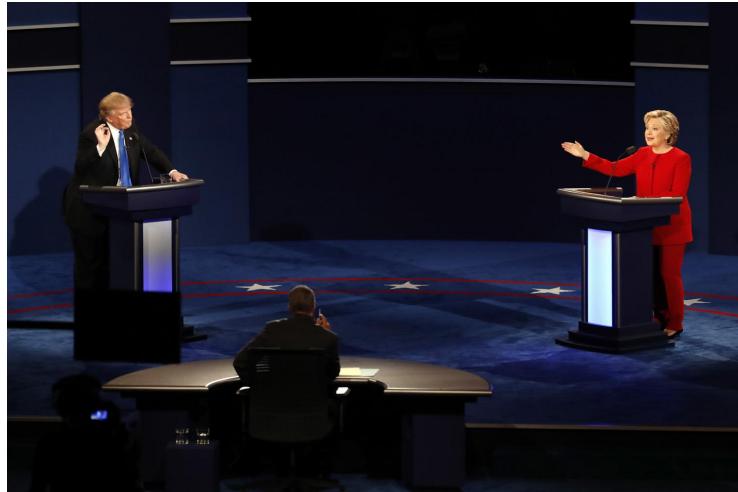
Who won the first presidential debate?



Who won the first presidential debate?

Why did I have you close your eyes?

I wanted to get **independent** samples



Are there any other sources of **measurement error**?

You might want to give an answer that you think I will like.

This is a **Demand characteristic**. E.g. the Bradley effect

Who do Americans think won the debate?



Hillary Clinton: 81%
Donald Trump 3%
42 votes cast

Each of these polls is a **sample**
But I want to make an inference
to the **population**



Hillary Clinton: 62%
Donald Trump 27%
521 votes cast

When I draw a conclusion about the population from a sample, I make an **inference**.

The way I collect my sample can lead me to different inferences.



Hillary Clinton: 25%
Donald Trump 67%
15905 votes cast

Which of these samples is the best?

Larger samples are better samples

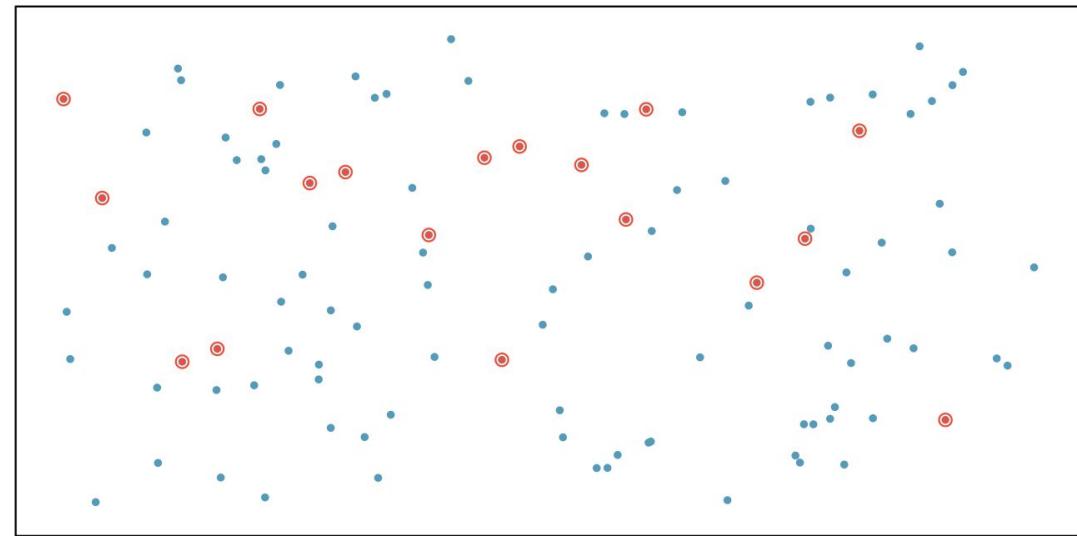
Why is bigger better?

Small samples are more **variable**.

There are 100 dots here, and 18 of them are red.

If I draw 3 dots, **more than half** the time 0 will be red.

If I draw 50 dots,
less than **1 out of 100 billion**
times 0 will be red



For random samples, larger samples are more **representative**

Are these random samples?



Hillary Clinton: 81%
Donald Trump 3%
42 votes cast



Hillary Clinton: 62%
Donald Trump 27%
521 votes cast

No! They are **convenience** samples



Hillary Clinton: 25%
Donald Trump 67%
15905 votes cast

How representative is the CNN/ORC sample?



Interviews with 521 registered voters who watched the presidential debate conducted by telephone (landline and cell) by ORC International on September 26, 2016. The margin of sampling error for results based on the total sample is plus or minus 4.5 percentage points.

Survey respondents were first interviewed as part of a random national sample conducted September 23-25, 2016. In those interviews, respondents indicated they planned to watch tonight's debate and were willing to be re-interviewed after the debate.

26% of the respondents who participated in tonight's survey identified themselves as Republicans, 41% identified themselves as Democrats, and 33% identified themselves as Independents.

How representative is the Washington Times sample?



Hillary Clinton: 25%

Donald Trump 67%

15905 votes cast

People could vote multiple times.
Why is this bad?

People voted if they read the Washington Times.
Why is this bad?

People from outside the US could vote.
Why is this bad?

Sampling bias in the polls: Landon vs. FDR



Alf Landon

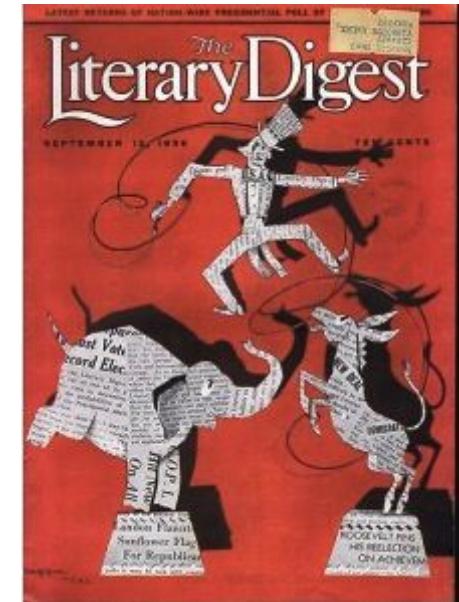


Franklin Delano Roosevelt

In 1936, Landon sought the Republican presidential nomination opposing the re-election of FDR.

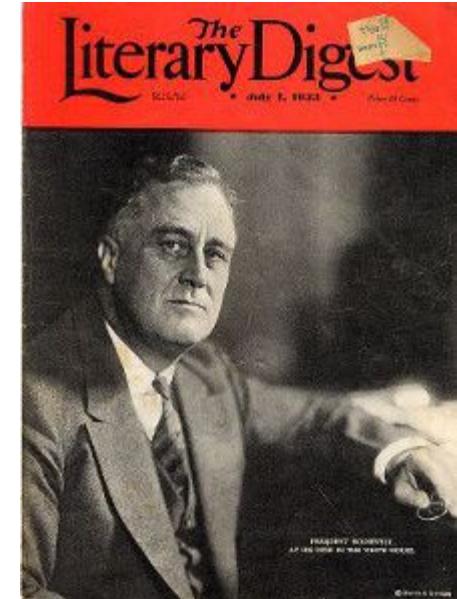
The Literary Digest poll

- The Literary Digest polled about 10 million Americans, and got responses from about 2.4 million.
- The poll showed that Landon would likely be the overwhelming winner and FDR would get only 43% of the votes.
- Election result: FDR won, with 62% of the votes.
- The magazine was completely discredited because of the poll, and was soon discontinued.



What went wrong?

- The magazine had surveyed
 - its own readers,
 - registered automobile owners,
 - registered telephone users, and
 - country club members
- These groups had incomes well above the national average--it was the Great Depression!
 - The sample was **not representative**
- This sample was huge--2.4 million people. But it was biased, and thus inaccurate.



A sampling metaphor



When you taste a spoonful of soup and decide the spoonful you tasted isn't salty enough, that's **exploratory analysis**

If you generalize and conclude that your entire soup needs salt, that's an **inference**

For your inference to be valid, the spoonful you tasted (the **sample**) needs to be **representative** of the entire pot (the **population**)

If the soup is not well stirred, it doesn't matter how large a spoon you have, it will still not taste right. If the soup is well stirred, a small spoon will suffice to test the soup.

Practice Question 1

A school district is considering whether it will no longer allow high school students to park at school after two recent accidents where students were severely injured. As a first step, they survey parents by mail, asking them whether or not the parents would object to this policy change. Of 6,000 surveys that go out, 1,200 are returned. Of these 1,200 surveys that were completed, 960 agreed with the policy change and 240 disagreed.

Which of the following statements are true?

1. Some of the mailings may have never reached the parents.
 2. The district has strong support from parents to move forward with the policy
 3. It is possible that majority of the parents disagree with the policy change.
 4. The survey results are unlikely to be biased because all parents were mailed a survey.
-
- (a) Only 1 (b) 1 and 2 (c) 1 and 3 (d) 3 and 4 (e) Only 4

Practice Question 1

A school district is considering whether it will no longer allow high school students to park at school after two recent accidents where students were severely injured. As a first step, they survey parents by mail, asking them whether or not the parents would object to this policy change. Of 6,000 surveys that go out, 1,200 are returned. Of these 1,200 surveys that were completed, 960 agreed with the policy change and 240 disagreed.

Which of the following statements are true?

1. Some of the mailings may have never reached the parents.
 2. The district has strong support from parents to move forward with the policy
 3. It is possible that majority of the parents disagree with the policy change.
 4. The survey results are unlikely to be biased because all parents were mailed a survey.
- (a) Only 1 (b) 1 and 2 (c) 1 and 3 (d) 3 and 4 (e) Only 4

Practice Question 2



Hillary Clinton: 81%
Donald Trump 3%
42 votes cast



Hillary Clinton: 62%
Donald Trump 27%
521 votes cast

I want to prediction who **UChicago Students as a whole** think won.

Which Sample should I use?

The Washington Times

Hillary Clinton: 25%
Donald Trump 67%
15905 votes cast

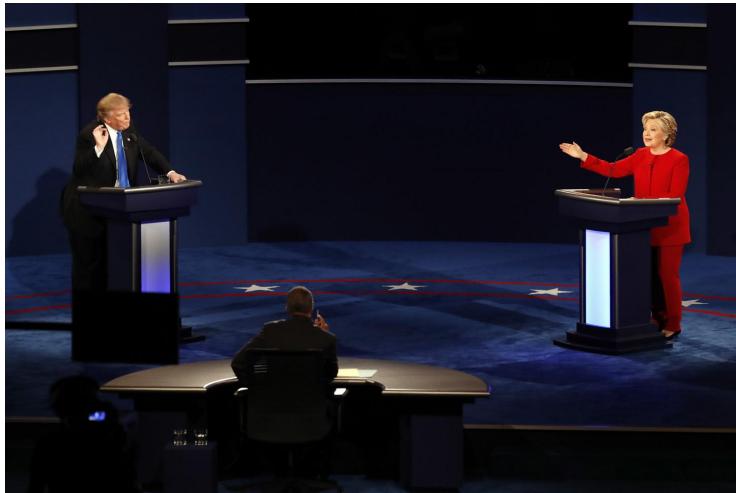
Do we know if watching the debate played a **causal** role?

What if we ask our non-watchers?

Hillary Clinton: XX%

Donald Trump XX%

xx votes cast



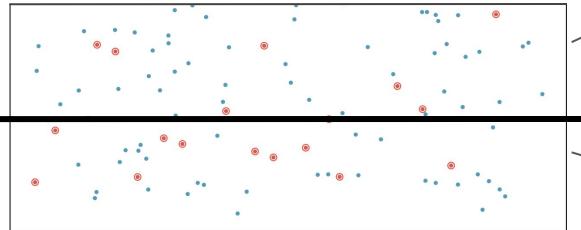
Can I just compare the watchers and non-watchers to each-other?

Can the non-watchers be a **control** group for the watchers?

They might actually be from a different **population**.

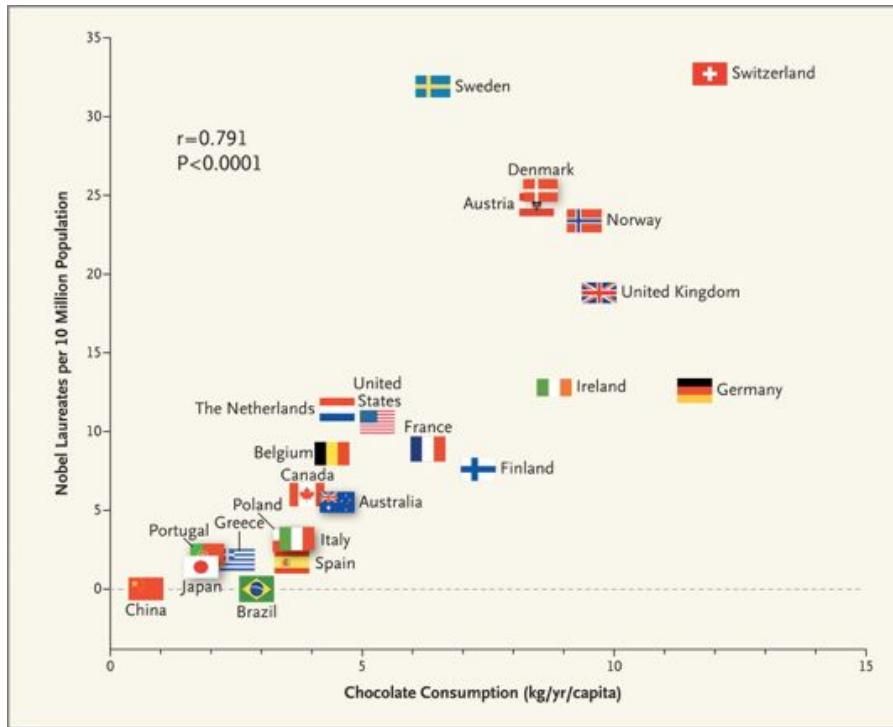
We did an **observational study**, not an **experiment**

To know that my **treatment** was causal is to use **random assignment**



Who won?

Consequences of non-random assignment



Chocolate makes you brilliant?

Brilliant people like chocolate?

What else could it be?

Practice Question 3

A study that surveyed a random sample of otherwise healthy adults found that people are more likely to get migraines when they're stressed. The study also noted that people drink more coffee and sleep less when they're stressed.

What type of study is this?

Observational

What is the conclusion of the study?

There is an association between increased stress & muscle cramps.

Can we conclude a **causal relationship between increased stress and migraines?**

Migraines might also be due to increased caffeine consumption or sleeping less – these are potential **confounding** variables.

Key ideas

1. Using samples to make inferences about populations
2. The way you sample your data can change your inferences about the population
3. Experiments use random assignment to treatment groups,
observational studies do not
4. Random samples help with **generalizability**,
random assignment helps with **causality**

Things to do:

Take the CAOS Test. Due Friday Night!

Start thinking about the homework.

Online Resources

Course Website:

<https://dyurovsky.github.io/psyc20100/>

- Find syllabus, slides, etc.

Google Classroom:

<https://classroom.google.com/>

- Submit assignments, post questions, etc.