

Data Science 100

Principles & Techniques of Data Science

Slides by:
Deborah Nolan deborah.nolan@berkeley.edu
 Spring 2018 updates
Fernando Pérez fernando.perez@berkeley.edu

1/21/18 Data 100 - Lecture 2 1


Announcements for Today

- Midterm is scheduled for **March 8 during class**
- We will try using Google forms today, but may need to resort to clickers/cell phone
- Slides and notes from lecture available online at <http://ds100.org/sp18>
- HW 1 will be released next Tuesday (Jan 23).


1/21/18 Data 100 - Lecture 2 2

A quick bit about me


- New faculty in statistics department (Fall 2017)
- Physics, applied math background
- Open source tools for science, reproducible research, interactive data science



Fernando Pérez
 419 Evans
 OH: Tues 2-4pm



Over 1,400 contributors!



1/21/18 Data 100 - Lecture 2

Let's finish name/age model from Tues.

1/21/18 Data 100 - Lecture 2 4

Three themes in DS 100

- *Problem formulation*
- Data Provenance
- Implications of scale

Each class will address one or more of these themes

1/21/18

Data 100 - Lecture 2

6

Data Provenance

We take a narrow approach to this topic

1/21/18

Data 100 - Lecture 2

6

How were these data generated?

- The mechanisms by which the data arose impacts whether we can answer the question of interest
- Can we generalize beyond what we observe in our data?

1/21/18

Data 100 - Lecture 2

7

How were these data generated?

- Do they form a census – information on all subjects?
- Are they a subset or sample?
 - Was a chance mechanism used to select the sample?
 - Or, is this a self-selected, convenience, or judgment sample?
 - Or, an administrative database?
- Were the data generated from
 - a collection in nature
 - an industrial process
 - a social science study


1/21/18

Data 100 - Lecture 2

8

Baby Names Data

- > Do they form a census?
- > Can we use our findings to make generalizations about the students in this class?
- > SSN system created 1937

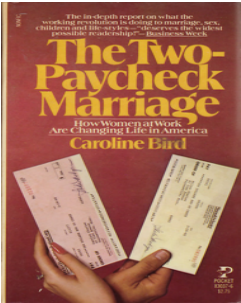



1/21/18 Data 100 - Lecture 2 9

Statistical Science – A first encounter (D. Nolan)

Data collected from questionnaire published in Women's Day magazine

Thousands of women responded to questions about their lifestyle

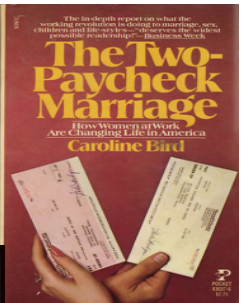




1/21/18 Data 100 - Lecture 2 10

Statistical Science – A first encounter

Do they form a census – information on all subjects? **NOPE**

Were there enough responses to overcome any issues with a nonprobability sample? **NOPE**

1/21/18 Data 100 - Lecture 2 11

Computational Science: powerful tools



1/21/18 Data 100 - Lecture 2 19

The purpose of computing is insight, not numbers.

R. Hamming
Numerical Methods for Scientists and Engineers (1962)

1/21/18

Data 100 - Lecture 2

14

Lessons Learned

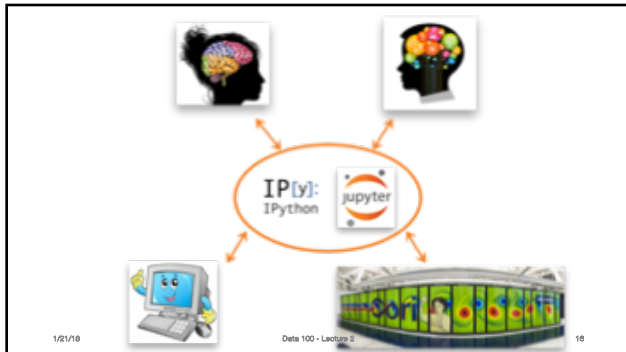
Good Data Analysis \neq
Simple Application of a Statistics Recipe

Good Data Analysis \neq
Simple Application of Statistical Software

1/21/18

Data 100 - Lecture 2

18



To Answer this question we need an estimate of fertility rates of Harvard mothers

Are first-borns more likely to attend Harvard?

Between 75% and 80% of students at Harvard are first-borns. Do first-born children work harder academically, and so end up overrepresented at top universities? So claims noted philosopher Michael Sandel. But **Antony Millner** and **Raphael Calel** find a simple fault in the statistical reasoning and give a more plausible explanation.

1/21/18

Data 100 - Lecture 2

19

DETOUR:

1. What is a simple random sample?
2. Why is it so desirable?
3. Can we make up for not having a SRS with big data?

1/21/18

Data 100 - Lecture 2

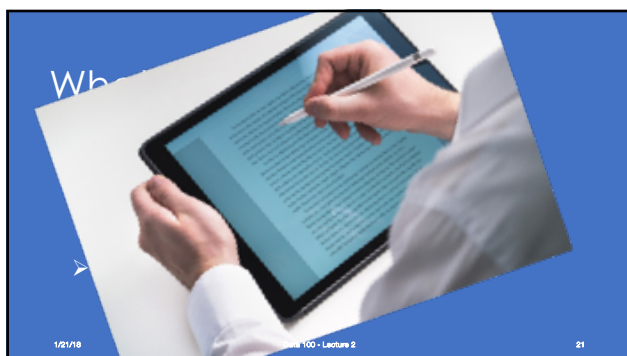
19

What is a Simple Random Sample?

1/21/18

Data 100 - Lecture 2

20



1/21/18

Data 100 - Lecture 2

21

The Simple Random Sample

- Suppose we have a population with N subjects
- We want to sample n of them
- The SRS is a random sample where every possible unique subset of n subjects has the same chance of appearing in the sample
- This means each person is equally likely to be in the sample

1/21/18

Data 100 - Lecture 2

22

Why is the SRS so Useful?

1/21/18

Data 100 - Lecture 2

24

The Advantages of a SRS

- *Representative: The sample tends to look like the population*
- *Statistics based on the sample tend to be close to statistics based on the population*
- *We can provide typical deviations of sample statistics from population values.*
- *AND MORE...*

1/21/18

Data 100 - Lecture 2

25

Fertility of mothers aged 40-44

1/21/18

Data 100 - Lecture 2

26

Fertility of mothers aged 40-44

- *8.7 million mothers aged 40-44 in the US (CPS)*
- *The have 1, 2, 3, or 4+ children (we treat 4+ as 4)*

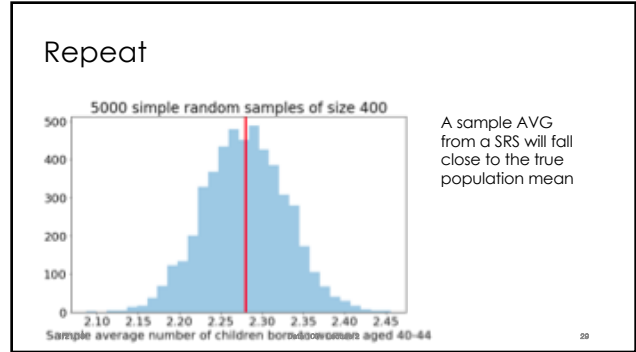
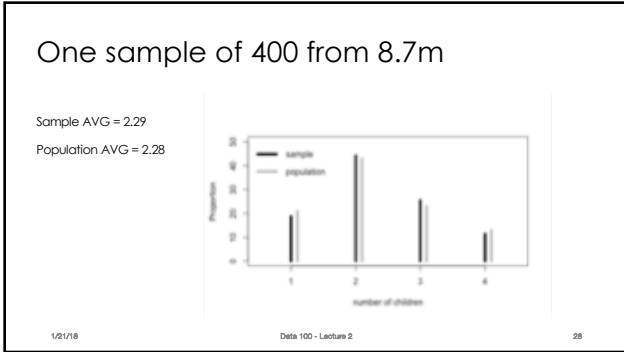
	Number of Children			
	1	2	3	4+
Percent mothers	21%	43%	23%	13%

- *Take a Simple Random Sample of 400 (simulate)*

1/21/18

Data 100 - Lecture 2

27



Can we make up for no SRS with Big Data?

1/21/18 Data 100 - Lecture 2 30

Harvard Professor Meng asked a room full of statisticians at an international conference something like this:

Suppose we want to estimate a proportion of people in the US that have a particular characteristic. We can either take a SRS of 100 people, or examine an administrative dataset where those who have the characteristic are slightly more likely to be included in the dataset. How large does the AD need to be to provide a more accurate estimate than the SRS?

Survey: <http://bit.ly/d100-sp18-afd>

A. 5% (16m) B. 20% (64m)
C. 50% (160m) D. 90% (206m)

1/21/18 More: Meng 2013. A trio of inference problems that could win you a Nobel Prize in statistics...

Please also answer the following questions:

- Are you your birth mother's first born?
- How many children does your birth mother have? (including you)

1/21/18

Data 100 - Lecture 2

82

Three possibilities

- SRS of 400 mothers
- SRS of 400 of the children
- An administrative dataset with 80,000 mothers, where those with more children are slightly more likely to be included in the database

Which of these three approaches better?

1/21/18

Data 100 - Lecture 2

83

SRS of 400 children

- Compare the chance a child from a 2 children family is chosen to the chance an only child is selected
- What does this imply about the estimate for the average number of children?
- Can we fix this problem?

This is size-biased sampling
Why?
We can fix the problem if we know how many siblings each child has

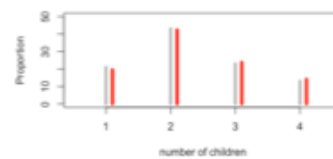
1/21/18

Data 100 - Lecture 2

84

Administrative Data of 80000 mothers

- Compare the proportions in the true population to our administrative dataset



This can't possibly make a difference, right?

1/21/18

Data 100 - Lecture 2

85

Administrative Data of 80,000 mothers

- Compare the accuracy based on 1000 simulations of Administrative Correlated Data

SRS of 400 vs Administrative Sample of 80,000

sample average number of children born to women aged 40-44

1/21/18 Data 100 - Lecture 2 86

Are first-borns more likely to attend Harvard?

And What about Berkeley?

1/21/18 Data 100 - Lecture 2 87

First-borns predominate at Harvard

- Professor Sandel found 75% of the students in his class are first born
- We know from our study of mothers aged 40-44

	Number of Children			
	1	2	3	4+
Percent mothers	21%	43%	23%	13%

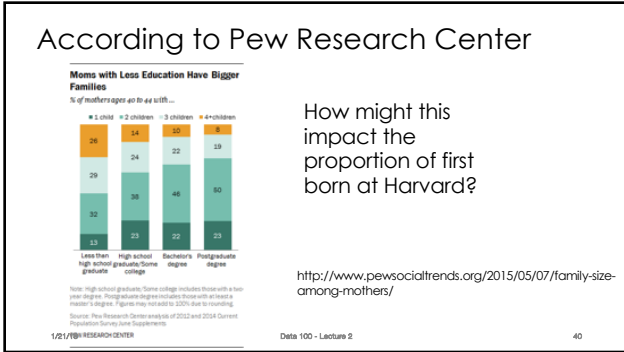
- The proportion of first born children in this group of children: $100 / (21 + 86 + 69 + 52) = 1/2.3 = 0.44$ or 44%

1/21/18 Data 100 - Lecture 2 88

First-borns predominate at Harvard

- BUT that's the wrong comparison
- What about Data Provenance?
- How large are Harvard families?
- Do we have any reason to believe that mothers of Harvard students might not have the same fertility rate as all US mothers aged 40-44?

1/21/18 Data 100 - Lecture 2 89



Are first-borns more likely to attend Harvard?

1/21/18 Data 100 - Lecture 2 42

Children of mothers with at least one child at Harvard

1/21/18 Data 100 - Lecture 2 43

Children of mothers with at least one child at Harvard

1/21/18 Data 100 - Lecture 2 44

Children of mothers with at least one child at Harvard

Want $P(H|F)$

1/21/18 Data 100 - Lecture 2 46

Find the chance first born is at Harvard

We want to find the chance you are a Harvard student (H) given you are first born (F) and

$$P(H|F) = \frac{P(H \cap F)}{P(F)}$$

Compare it to the chance you are a Harvard student given you are not first born

$$P(H|F^c) = \frac{P(H \cap F^c)}{P(F^c)}$$

1/21/18 Data 100 - Lecture 2 46

Find the chance first born is at Harvard

Recall $P(H \cap F) = P(H)P(F|H)$

Which implies Bayes Rule $P(H|F) = \frac{P(H)P(F|H)}{P(F)}$

So the ratio becomes: $\frac{P(H|F)}{P(H|F^c)} = \frac{P(F|H)}{P(F^c|H)} \times \frac{1-P(F)}{P(F)}$

Professor observed $\frac{3}{4}$ for $P(F|H)$. But we need the fertility rate for Harvard moms

1/21/18 Data 100 - Lecture 2 47

Find the chance first born is at Harvard

So the ratio becomes:

$$P(F) = \frac{1}{\text{fertility}}$$

$$r = \frac{P(H|F)}{P(H|F^c)} = \frac{3/4}{1/4} \times \frac{1-P(F)}{P(F)} = 3(\text{fertility}-1)$$

We need the fertility rate for Harvard moms...

If we plug in the US value for 1994 we get: $r = 3.9$

If we plug in the US highest educated value from Pew we get: $r = 3.2$

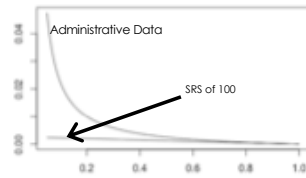
1/21/18 Data 100 - Lecture 2 48

Can Big Data compensate for no SRS?

Big Data is **not** enough

With a slight correlation between the presence of the characteristic and the appearance in the dataset, the accuracy of a SRS of 100 outperforms a AD of 90%

I Got More Data, My Model Is More Refined..., Meng & Xie (2014)
Econometrics Review



1/21/18

Data 100 - Lecture 2

40