# Computational Thinking and Inferential Thinking

## *Foundations of Data Science*

Michael Jordan
michael_jordan@berkeley.edu

University of California, Berkeley
November 2020

# Our Vision

- All students should be able to make decisions based on data
- Data science education should be a conceptual blend of traditions in statistics and computer science, with a strong outward-looking component
- All students should be able to tackle significant problems, the complexity increasing with the students' background knowledge and interests

# The Vision, Realized

- Foundations of Data Science (data8.org)
  - Connector courses
- Principles and Techniques of Data Science (data100.org)
- Probability for Data Science (prob140.org)
- Data, Inference, and Decisions (data102.org)
- Human Context and Ethics (https://data.berkeley.edu/file/2062)

# Data, Inference, and Decisions

# The Students



**2019-2020**
- 2800 in Data 8
- 1600 in Data 100
- 500 in Prob 140
- 500 in Data 102

The first Data Science majors graduated in Spring 2019.

Over 300 DS majors graduated in Spring 2020, along with several hundred who got the DS minor.

# Data 8: *Foundations of Data Science*

- In the 2015-2016 academic year, we taught over 500 freshman-level students in our new Data Science course
  - a wide range of students, pre-disciplinary
  - it was a major success

- In the 2016-2017 academic year, we taught over 1000 students, and over 2500 in 2019-2020

- The course is a conceptual blend of traditions in statistics and computer science, with a strong outward-looking component ("connector courses")

# Some shapers of our perspective

- The increasingly prominent use of "data science" in job descriptions in the tech industry
  - and our undergraduates' awareness of the term

- The rise of the "CS for All" movement
  - and its partial failure at many universities

- Success in teaching inference at the graduate level via resampling and sampling ideas
  - surely this can work at the undergrad level too?

- The explosive growth of iPython (Jupyter)
  - and related browser-based computing environments

# Computational thinking

- Modern programming and problem-solving with a computer requires a new style of thinking
  - the core concepts: abstraction, modularity, scalability, robustness, etc.
  - concepts that are taught beautifully and successfully in freshman-level courses in computer science

- What's missing in computational thinking:
  - much of what we teach in our best statistics courses

# Inferential thinking

- Analyzing data requires more than just computation, but it also requires a style of thinking
  - consider the data-generating process behind the data
  - being clear on the sampling pattern, with its possible biases
  - understanding how to do causal inference
  - understanding how to capture uncertainty

- What's missing in inferential thinking:
  - much of what we teach in our best computer science courses

# Blending

- We believe that the blend of computational thinking and inferential thinking is:
  - powerful
  - engaging to students, by empowering them to solve real problems
  - what industry and government are really asking for when they refer to "data science" or "AI"
  - a kind of liberal arts for the 21st century

# Analyzing data: three classical steps

- Formulate the question, from some domain; reasonable assumptions about the data; choice of method

- Visualization and calculations

- Interpretation of the results in the language of the domain, ideally without statistical jargon

# Classical teaching of data analysis

- The question, from some domain; reasonable assumptions about the data; choice of method

- Visualization and formulas

- Interpretation of the results in the language of the domain, without statistical jargon

# The Data 8 way

- The question, from some domain; reasonable assumptions about the data; choice of method

- Visualization and computation

- Interpretation of the results in the language of the domain, without statistical jargon

# The approach

- A course for freshmen
- No prerequisites

- Teach data science as a way of thinking:
  - computational thinking
  - inferential thinking
  - these are blended in every lesson

# The Syllabus

- Cause and Effect
- Tables
- Data Types
- Census
- Histograms, Charts
- Functions
- Joins
- Iterations
- Chance
- Sampling
- Models
- Comparing Distributions
- Decisions and Uncertainty
- P-Values
- A/B Testing
- Causal Inference
- Confidence Intervals
- Designing Experiments
- Regression
- Classification

# Three 2-week projects

- **Water use in California:** mapping the water districts and overlaying IRS taxable income data by zip code

- **Murder rates and the death penalty:** nonparametric inference; the importance of visualization

- **Classification of song lyrics:** hip-hop or country?

# "Connector" Courses

- Data science in a particular domain
- Much smaller: 5-50 students
- Half the time commitment of a typical undergraduate course
- Can be taken concurrently with Data 8, or later

# Connector Courses, 2015-2016

- Data Science for Smart Cities
- Data Science and the Mind
- Data Science in Ecology and the Environment
- Exploring Geospatial Data
- How Does History Count?
- Data and Ethics
- Health, Human Behavior, and Data
- Race, Policing, and Data Science
- Literature and Data
- Matrices and Graphs in Data Science
- Computational Structures in Data Science
- Probability and Mathematical Statistics for Data Science

data.berkeley.edu

# Connections Fall 2016

# Data 8 Spring 2016 survey

- About 450 students, 418 responses
  - 42% women
  - 17% underrepresented ethnic or racial minority
  - 52 different intended majors

- Does the course fulfill a major requirement?
  - Yes:                  8.6%
  - Maybe:                10.4%
  - I don't know:         10.6%
  - No:                   70.4%

# Always DataBears



Of the roughly 1210 students in Foundations of Data Science in Fall 2018, about 70 worked on the Spring 2019 offering of the course and its connectors, as lab assistants, tutors, graders, and course developers.

# Environment and course materials

- Jupyter notebooks; Python 3

- JupyterHub
  - Multi-user server for Jupyter notebooks
  - Browser-based computation in the cloud

Course materials: data8.org

Source: https://github.com/data-8

# Data, Inference, and Decisions

# Data 102: *Data, Inference and Decisions*

- In the 2019-2020 academic year, we initiated a new senior-level course in Data Science course

- The course is again a conceptual blend of ideas from statistics and computer science, with a strong outward-looking component

- But it also brings in economics and control theory, offering an integrated perspective on decision-making under uncertainty

# The Syllabus

- False Discovery Rate Control
- Online False Discovery Rate Control
- Bayesian and Frequentist Decision-Making
- Probability Interpretation of Linear and Logistic Models
- Bayesian Hierarchical Models
- Importance Sampling and Gibbs Sampling
- Bootstrap
- Concentration Inequalities
- Causal Inference
- Experimental Design
- Bandits: Greedy and UCB Algorithms
- Thompson Sampling
- Matching Markets
- Game Theory
- Dynamic Programming and Q-Learning
- Introduction to LQR and Control Theory
- Privacy

# *Example*: Privacy

# Privacy and Data Analysis

● Individuals are not generally willing to allow their personal data to be used without control on how it will be used and now much privacy loss they will incur
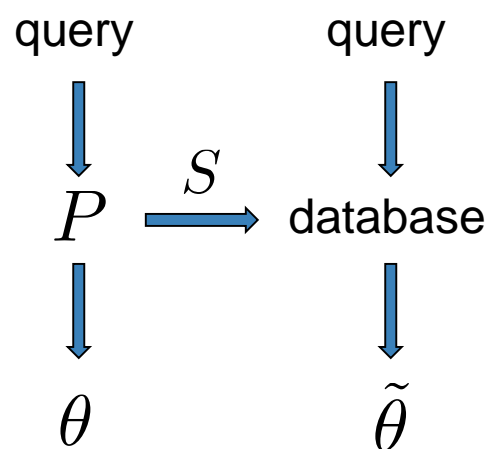
● "Privacy loss" can be quantified (e.g.) via differential privacy

● We want to trade privacy loss against the value we obtain from "data analysis"

● The question becomes that of quantifying such value and juxtaposing it with privacy loss

# Privacy

query                          query

↓                              ↓

database   $\xrightarrow{Q}$   privatized
                               database

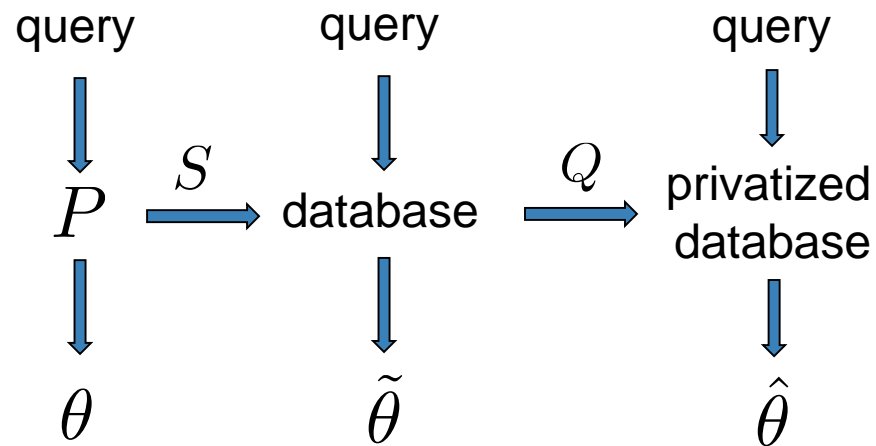↓                              ↓

$\tilde{\theta}$               $\hat{\theta}$

Classical problem in differential privacy:  show that $\hat{\theta}$ and $\tilde{\theta}$
are close under constraints on $Q$

# Inference

query          query

$$P \xrightarrow{\;S\;} \text{database}$$

$$\theta \qquad\qquad\qquad \tilde{\theta}$$

Classical problem in statistical theory: show that $\tilde{\theta}$ and $\theta$ are close under constraints on $S$

# Privacy and Inference



The privacy-meets-inference problem: show that $\theta$ and $\hat{\theta}$ are close under constraints on $Q$ and on $S$

# Conclusion

- Data Science involves blending computational thinking with inferential thinking

- Education in Data Science involves teaching this blend, in the context of solving real-world, substantive problems

- I believe that Data Science is a core part of a modern liberal education

# Computational Thinking and Inferential Thinking

## *Foundations of Data Science*

Michael Jordan
michael_jordan@berkeley.edu

University of California, Berkeley
November 2020