

ORCA 2500/4500

cjd 2119@columbia.edu

Hw 251.

Midterm 251.

Final Exam 251.

Final Project 251.

extra credit

class participation

course objectives

① teach basic python, (data stack)
SQL

② teach basic statistics and probability

③ regression

④ machine learning

What is data science?

→ multidisciplinary field

that combines the tools of compy,
statistics and machine learning to
solve problems!

→ approach-agnostic

→ what is statistics?

what is machine learning?

① descriptive statistics

↳ summarize and visualize

data

② inferential statistics

↳ making conclusions in

a mathematically rigorous
way about data

③ learning patterns in data and

applying it to making predictions

populations and sampling

Typically we are interested in understanding the difference between two populations, or learning about some parameters of a population.

Instead of looking at every single member of a population, we're going to look at a subset.

↳ how do we generate the subset

Causality

- Not enough to understand that two different populations are different → what to understand why?

Simpson's paradox

- sick patients given the option to buy a new drug
- 350 take it, 350 don't
(78%) (83%)

	Drug	No drug
M	81/87 (93%)	234/270 (87%)
F	192/263 (73%)	55/80 (69%)
Combined	78%	83%

gold standard for determining
causality →
randomized experiment

Observational study

→ a study in which scientists make conclusions based on data they have observed but had no hand in generating

confounding variable

variable that is correlated with both treatment and the outcome

Experiment

study in which scientists determine the assignment of the treatment