

ORIE 5355: People, Data, & Systems

Lecture 5: Data collection module epilogue

Nikhil Garg

Announcements


- HW1 due Tuesday evening (submit via Gradescope)
 - Don't wait until the last minute!
 - Go to office hours
 - **Remember to “tag” your answers to each question**
- Quiz 1 released tomorrow, due Friday evening (via Canvas)
- HW 2 posted

Case study: Ratings and recommendations

Overview

- So far, we've talked about explicit opinion collection in polling
- The same challenges apply in other settings
- Some differences
 - Often we don't care about "absolute" opinion but "relative" opinions
 - We care a lot about "heterogeneous" opinions
 - We often have other "implicit" data on people's opinions
- Briefly discuss some of these challenges in context of ratings and recommendations

Rating systems



Detailed Seller Ratings (last 12 months)

Criteria	Average rating
Item as described	★★★★★
Communication	★★★★★
Shipping time	★★★★★
Shipping and handling charges	★★★★★



Customer Reviews


★★★★★ 4
4.6 out of 5 stars ▾

5 star	<div style="width: 75%;"><div style="width: 75%;"></div></div>	75%
4 star	<div style="width: 25%;"><div style="width: 25%;"></div></div>	25%
3 star	<div style="width: 0%;"><div style="width: 0%;"></div></div>	0%
2 star	<div style="width: 0%;"><div style="width: 0%;"></div></div>	0%
1 star	<div style="width: 0%;"><div style="width: 0%;"></div></div>	0%

[See all verified purchase reviews ▸](#)

Share your thoughts with other customers

Write a customer review



Private Feedback

This feedback will be kept anonymous and never shared directly with the freelancer. [Learn more](#)

Reason for ending contract:
Please select...

Would you hire this freelancer again, if you had a similar project?
 Definitely Not Probably Not Probably Yes Definitely Yes

Public Feedback

This feedback will be shared on your freelancer's profile only after they've left feedback for you. [Learn more](#)

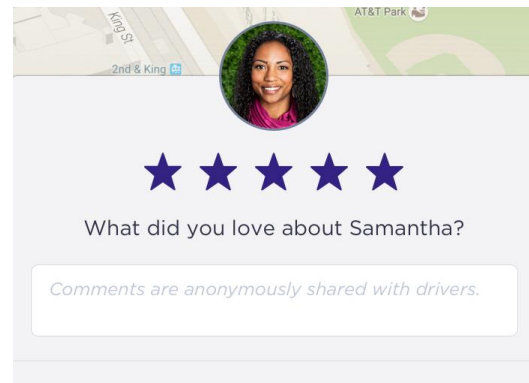

Feedback to Freelancer

- ★★★★★ Skills
- ★★★★★ Quality of Work
- ★★★★★ Availability
- ★★★★★ Adherence to Schedule
- ★★★★★ Communication
- ★★★★★ Cooperation

Total Score: **0.00**

Share your experience with this freelancer to the oDesk community:

[See an example of appropriate feedback](#)




Map showing location: 2nd & King, AT&T Park

★★★★★

What did you love about Samantha?

Comments are anonymously shared with drivers.



14 Reviews

★★★★★

Search reviews

Summary	Accuracy ★★★★★	Location ★★★★★
	Communication ★★★★★	Check In ★★★★★
	Cleanliness ★★★★★	Value ★★★★★

Translate reviews to English

Great location next to République stop. Nice communication from the

Measurement error: Ratings Inflation

4.68★
DRIVER RATING
 Unfortunately, your driver rating last week was **below average.**



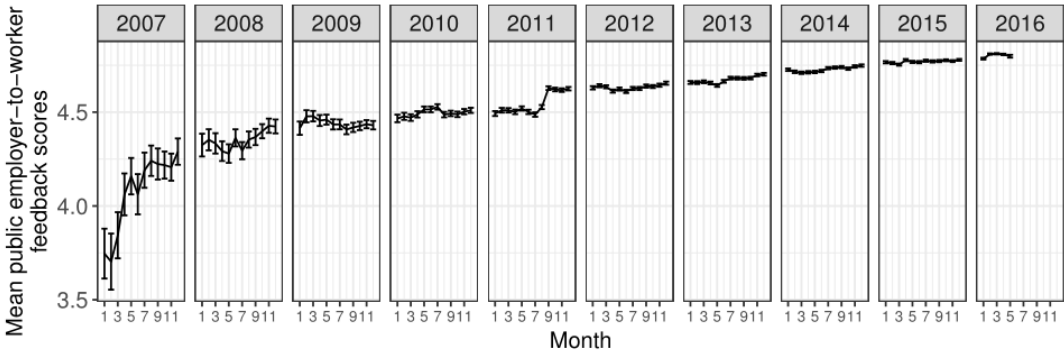
DON'T FORGET TO RATE 5 STARS
 FACT: WHEN A DRIVER'S RATING FALLS BELOW 4.7 THEY BECOME DEACTIVATED.
MORE DRIVERS MEANS LESS SURGES

OK, could be better

★★★★★

How can Sajid improve?

Figure 2: Monthly average public feedback scores assigned to workers by employers on completed projects.



[Filippas, Horton, Golden 2017]

When 4.3 Stars Is Average: The Internet's Grade-Inflation Problem

THE WALL STREET JOURNAL The Wall Street Journal April 5, 2017

UNDERSTANDING ONLINE STAR RATINGS:

- ★★★★★ [HAS ONLY ONE REVIEW]
- ★★★★★ EXCELLENT
- ★★★★☆ OK
- ★★★★☆
- ★★★★☆
- ★★★☆☆
- ★★★☆☆
- ★★★☆☆
- ★★★☆☆
- ★★★☆☆

CRAP

<https://xkcd.com/1098/>

Why ratings inflation & what to do about it?

- Many hypotheses for why ratings inflate
 - Explicit pressure from sellers – worry about retaliation
 - Implicit pressure – don't want to hurt people's livelihoods
 - Either misreport, or selection – less likely to report after bad experience
- Inflation is a type of measurement error:
 - The “quality” scale doesn't match well to the “rating” scale
 - Inflation over time – mapping from quality to rating changes over time
 - Why does it matter? We ask you this in the homework
- What to do about it:
 - Try to reduce some of the pressure
 - Weighting to tackle selection: paper in the homework: [Nosko & Tadelis]
 - Change the rating scale: [Garg and Johari]

Experiment Description

Status quo: Clients hire freelancers, rate them at contract end

Form includes a numeric rating from 0 to 10, with avg >8/10

Challenge: Can we induce different (non-inflated) ratings by changing the question we ask on the rating form?

Experiment design

- Add additional question to private portion of the form (6 treatments)
Randomization at the *client* level
- Observe ratings for 3 months (180k jobs, 60k clients, 80k freelancers)

Treatment groups

Treatment	Question Phrasing	Answer choices
Numeric	How would you rate this freelancer overall?	0 – 5

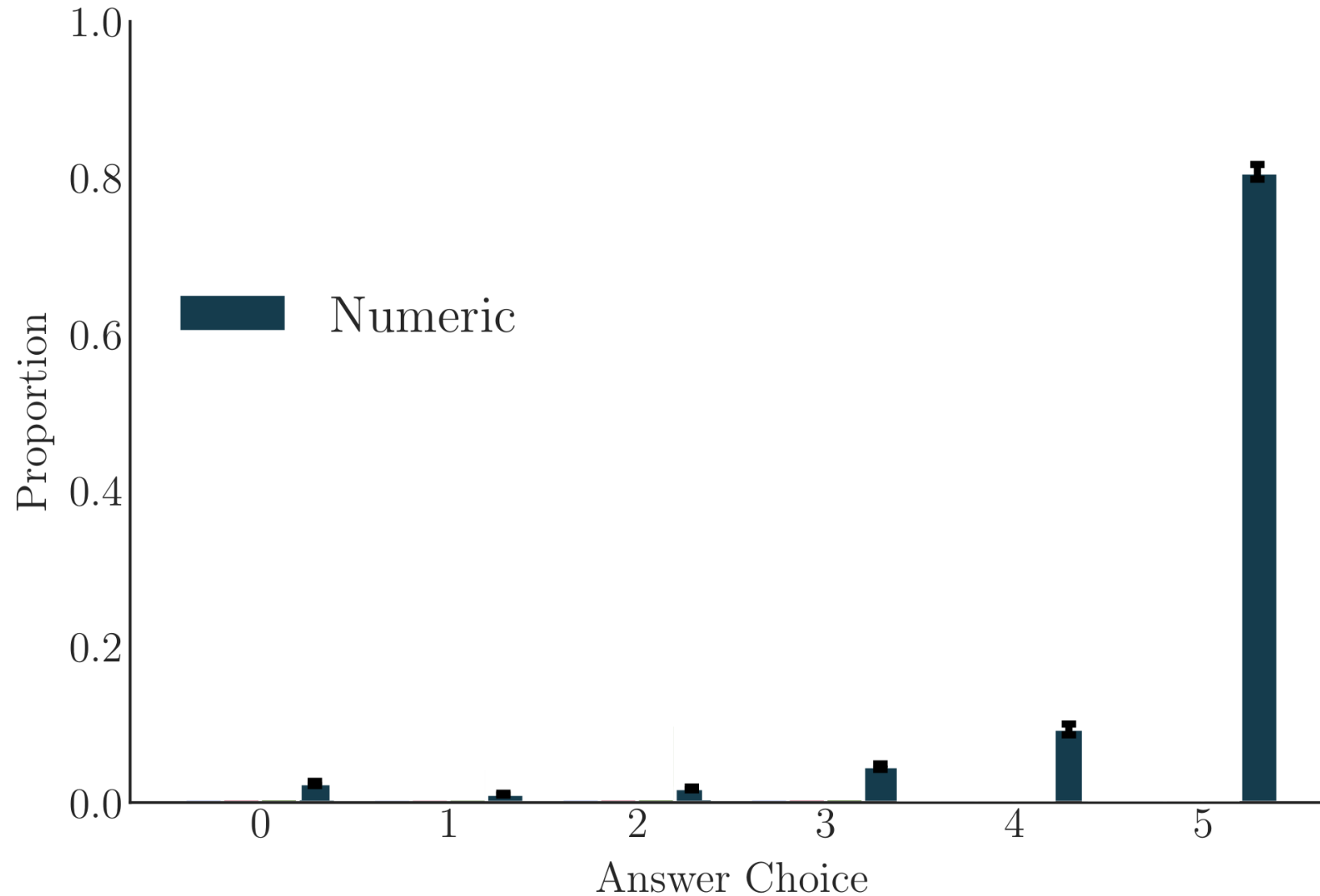
Treatment groups

Treatment	Question Phrasing	Answer choices
Numeric	How would you rate this freelancer overall?	0 – 5
Adjectives	How would you rate this freelancer overall?	Terrible Mediocre Good Great Phenomenal Best possible freelancer!

Treatment groups

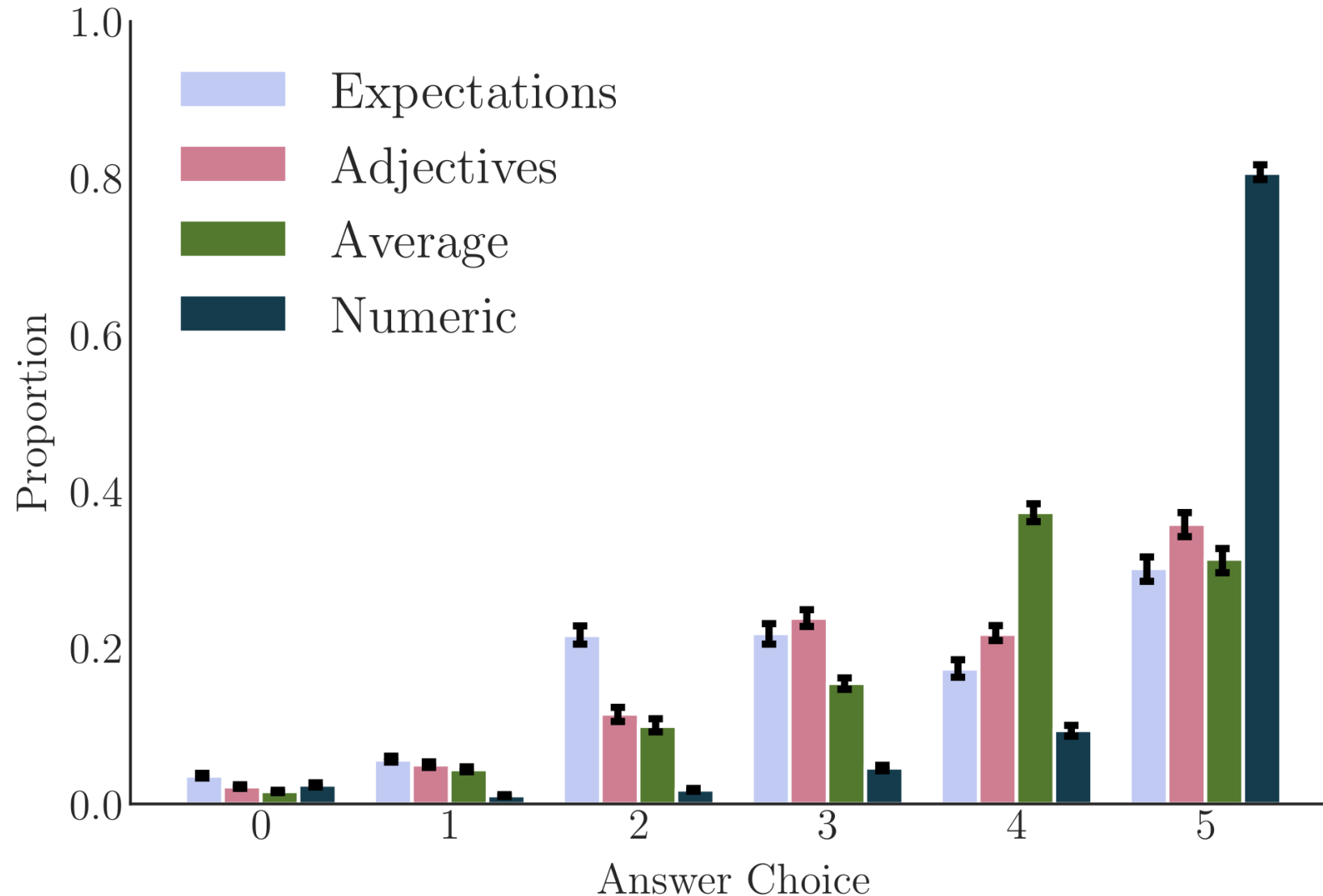
Treatment	Question Phrasing	Answer choices
Numeric	How would you rate this freelancer overall?	0 – 5
Adjectives	How would you rate this freelancer overall?	Terrible Mediocre Good Great Phenomenal Best possible freelancer!
Expectations	How did this freelancer compare to your expectations?	Much worse than I expected ... Beyond what I could have expected
Average	How does this freelancer compare to others you have hired?	Worst Freelancer I've Hired
Average, random order		Below Average Average Above Average
Average, not affect score		Well Above Average Best Freelancer I've hired

Result: marginal rating distributions



“Designing Informative Rating Systems: Evidence from an Online Labor Market”
Nikhil Garg and Ramesh Johari

Result: marginal rating distributions



“Designing Informative Rating Systems: Evidence from an Online Labor Market”
Nikhil Garg and Ramesh Johari

Ratings heterogeneity

- There is much ratings “heterogeneity”
 - Different people have different opinions on the same item
 - Different ‘categories’ of items might have different average ratings
- Why does this matter?
 - You want to give each person a personalized “rating” or recommendation
 - You want to compare items across categories
- What to do about it?
 - Personalized recommendations → starting next time
 - “Standardize” ratings across categories
 - Communicate to customers – e.g., “relative” ratings instead of “absolute” ones

Implicit data collection in recommendations

- You have many implicit signals about people's opinions
 - Do they finish watching the show, or start watching the next episode?
 - Do they keep coming back and buying other things
 - Did they browse other items instead of putting something in their cart?
 - Do they re-hire the same freelancer/work with the same client again?
- These give *different* information than do explicit ratings
 - From a different population of users
 - Often more numerous, but harder to analyze
 - “revealed preference” – might be more predictive of future behavior
- Using such data
 - Train models to predict different future behavior, using various signals
 - Might take away “user agency” – what if they want to change their behavior?

Case study 2: Crowdsourcing

Government service allocation

Local government manages many services

~8k miles of streets in NYC

~700k trees lining streets in NYC

Housing, sanitation, transportation, etc.

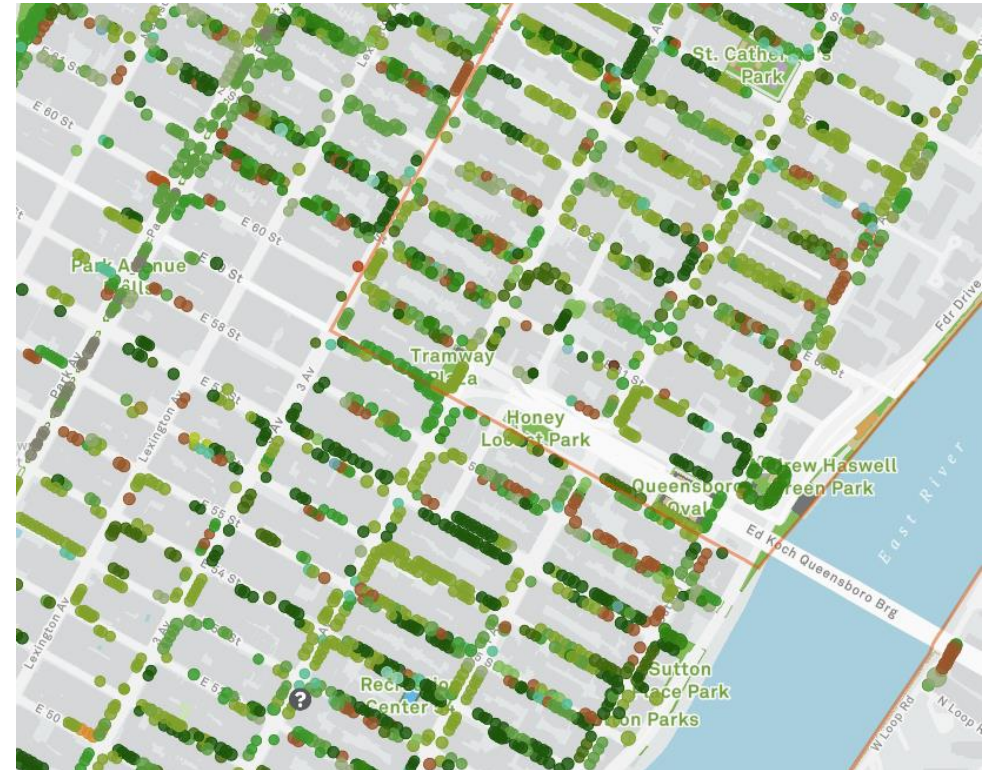
Operational tasks

[Learning] What problems are there?

[Allocation] Which ones to address?

[Auditing] Did we do a good job?

Desiderata: Efficiency & Equity



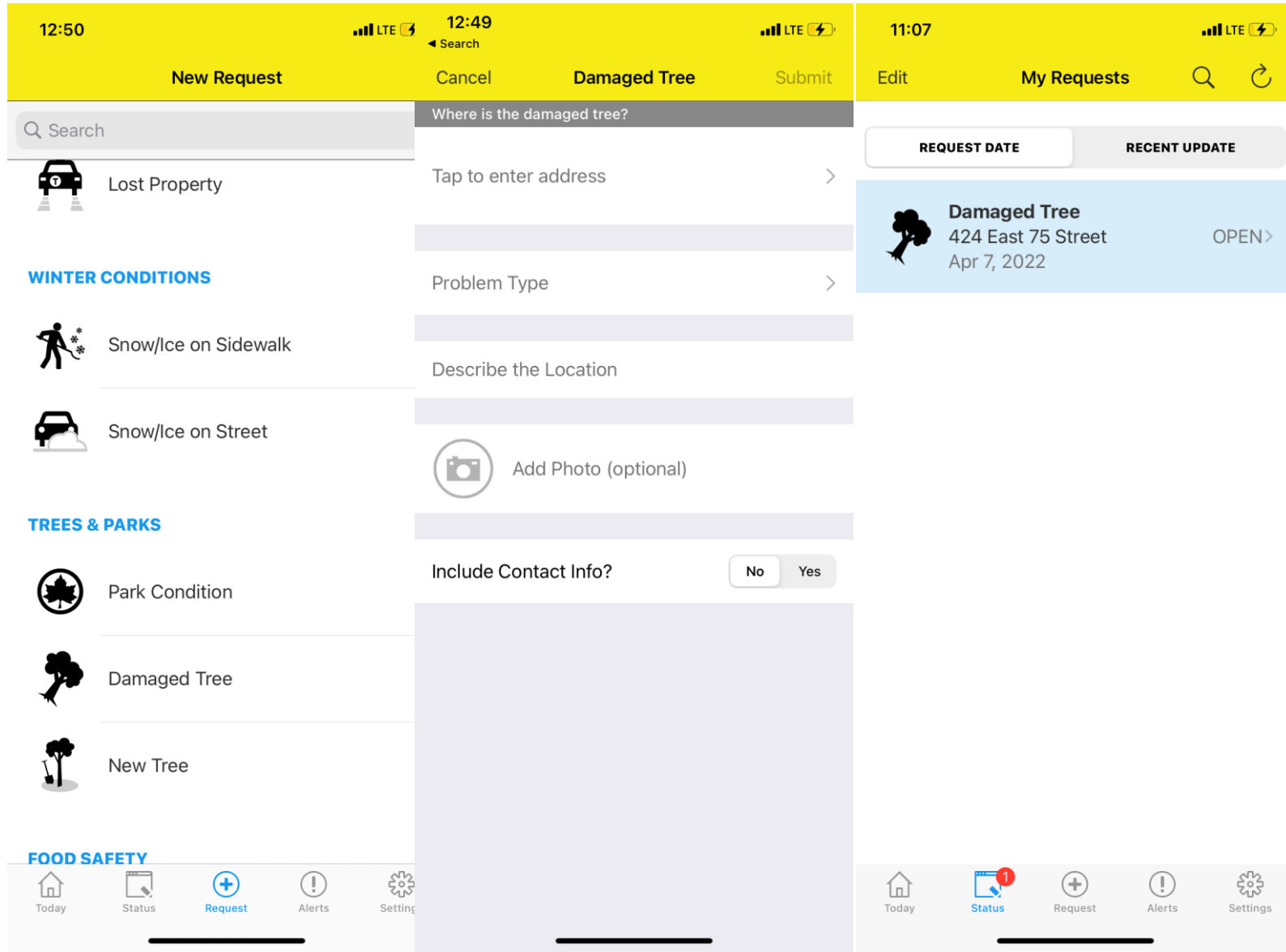
[Street trees on Upper East Side in NYC](#)

311 (crowdsourcing) systems

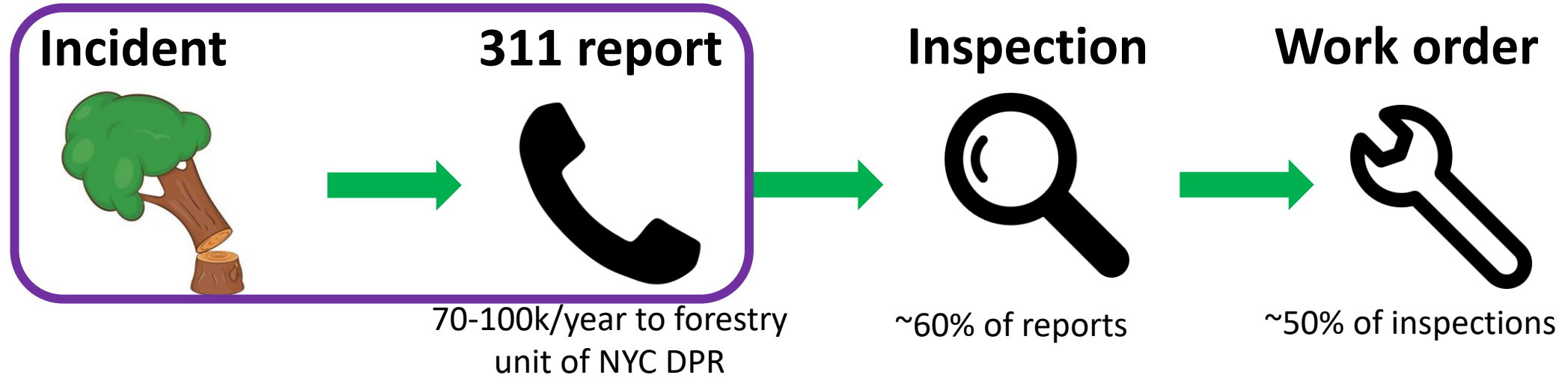
Cities have a phone number & app to complain to the local government

NYC's 311 system received about **2.7 million** requests 2021

These are the primary way the government learns about problems



Pipeline: from incident to work orders



Why is this hard? Uncertainty, heterogeneous + strategic behavior, distribution shifts over time, capacity constraints, pipelined decisions

Reporting behavior: If there are crowdsourcing differences (who reports what), then there will be downstream differences in decision-making

Possible data collection heterogeneity

Underreporting: The same number of problems exist in 2 neighborhoods, but one neighborhood reports more problems, faster.

Mis-reporting: Same types of problems in 2 neighborhoods, but people in one tend to exaggerate incident type/risk to get faster service.

In each case, we'll have disparities in what work gets done! (bad allocation of government services)!

Research agenda: How do we understand these reporting differences and then correct for them?

Miscellaneous topics in data and
data collection



(Differential) Privacy

- What if you're asking about a sensitive attribute?
For example, an insurance company wants to estimate the percentage of their policy holders who smoke
- Goal: collect data in a way such that you learn very little about any individual person, but you are accurate across population
- How? Add noise to each response
- Example: Tell each person, "roll a 6-sided dice. If it's 1 or 2, lie about whether you smoke. Otherwise, tell the truth." If fraction Y people tell you that they smoke, then we know that the truth X satisfies:

$$Y = \frac{4}{6}X + \frac{2}{6}(1 - X)$$

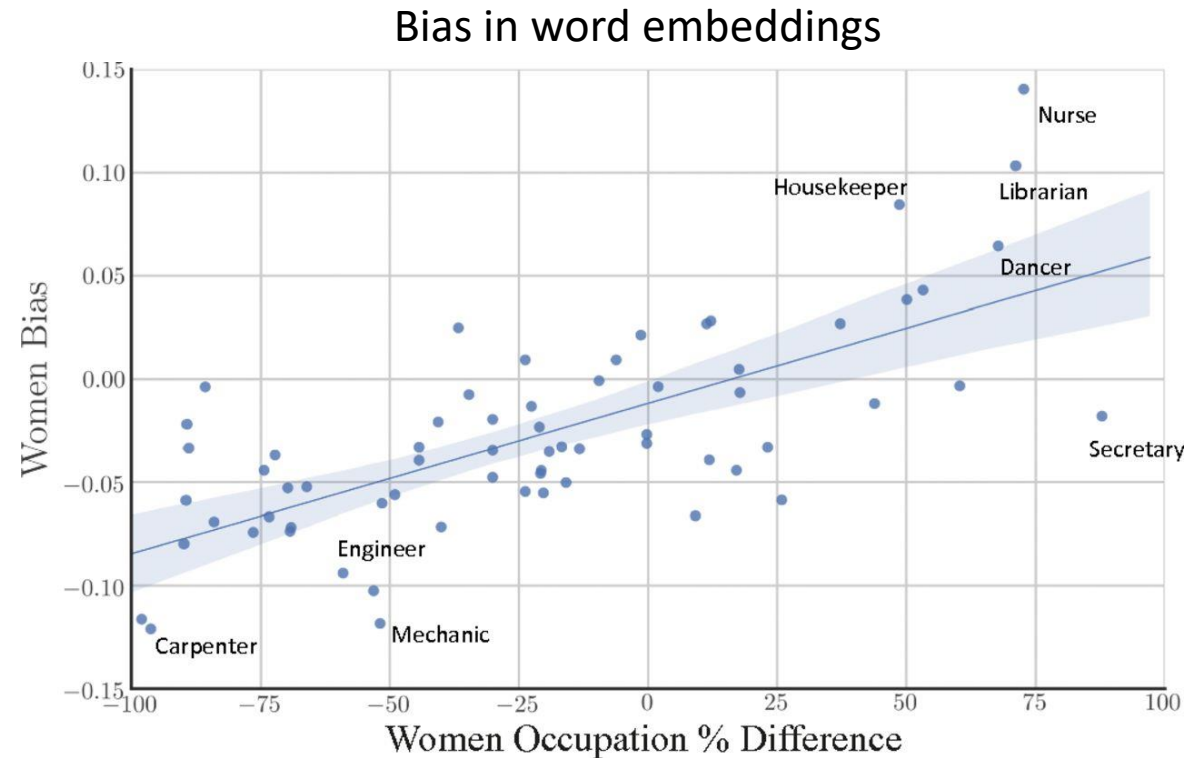
- Similar ideas used to collect and share data at Apple and the US Census

Eliciting complex opinions

- So far, we've talked about soliciting "low-dimensional" opinions
 - Binary opinions, or one of a small number of options
- What if we want to solicit opinions on complicated things?
 - How your town should spend \$2M budget across parks, sports teams, art festivals, etc.
 - When should we schedule these five events over 10 time slots?
- You can't ask people to rank every option
- Several standard techniques
 - Participatory budgeting
 - Pairwise comparisons
- More generally, many cool techniques in crowdsourcing

Using biased data

- The world is full of historic inequities
 - Some neighborhoods are over-policed compared to others → data will have more “crimes there”
 - Every possible opinion expressed on forums like Reddit
 - Who succeeded at a university
- Models trained using this data will *reflect* and *amplify* these biases
- Many techniques to audit and mitigate such biases in models



“Word Embeddings Quantify 100 Years of Gender and Ethnic Stereotypes” by Nikhil Garg, Londa Schiebinger, Dan Jurafsky, and James Zou

Module Summary

- Measurement error: The construct you care about is never perfectly captured by the data that you have
- Selection effects/differential non-response happens everywhere you're collecting opinions from people
- You can use stratification and weighting to mitigate selection effects *on known covariates*
- On unknown covariates, quantify uncertainty!

Never take opinion data at face value. Always ask:

- (1) What did I measure, versus what did I care to measure?
- (2) Who answered versus what's the population of interest
- (3) What am I going to *do* with the data, and how does that affect data collection?

Will show up in the rest of the course!

Questions?

(especially regarding homework)

Recommendation systems

Module overview

Part 1 – Prediction

How much will a given user like an item?

- Problem formulation and some algorithms
- Data challenges

Part 2 – From predictions to decisions

How to use the predictions to recommend items in practice?

- Capacity constraints
- Recommendations in *2 sided* markets
- Feedback loops in recommendations



Gregory Alan Isakov - If I Go, I'm Going (Acoustic Cover)
 Chase Eagleson 🎸
 713K views • 1 year ago



【理性讨论小组】2021 画空间【艺术跨学科对话】
 理性讨论小组
 7 views • 2 days ago



[SUB] Steamed Custard Buns :: Soft & fluffy :: Easy Recipe
 매일맛나 delicious day
 2M views • 4 months ago



Franz Schubert Octet in F Major, D 803
 Hochrhein Musikfestival
 1.4M views • 3 years ago

Recommendations for you



Your Orders



Grocery & Gourmet Food



Patio, Lawn & Garden



Electronics

Neighbory LIC Residents
 57 members • 3 posts a day

Join Group

JOIN OUR SOLO FEMALE TRAVELERS TRIPS
 All destinations <https://bit.ly/SFTTrips>

Solo Female Travelers (FIRST FB group for women who travel solo!)
 87K members • 60 posts a day

Join Group



WALKING FOR PLEASURE.
 15K members • 200 posts a day

Join Group

Types of Recommendations

Editorial and hand curated

- List of favorites
- Lists of “essential” items

Simple aggregates

Top 10, Most Popular, Recent Uploads

Tailored to individual users (Personalized recommendations)

Amazon, Netflix, ...

Personalized recommendations

- Motivation: filter the content to be more relevant for each individual
- Data Inferred from signals
 - Direct: ratings, feedbacks, etc
 - Indirect: purchase history, access patterns, etc
- Intermediate Goal: *predict* the relevance of each item for each user

Formal Model

- X = set of **Users**
- S = set of **Items**

Utility function $u: X \times S \rightarrow R$

R = Ratings that a user *would* give to an item if watched

R is a totally ordered set

e.g., **0-5** stars, real number in **[0,1]**

Ratings Matrix: suppose we have data \hat{R}

	Avatar	LOTR	Matrix	Pirates
Alice	1		0.2	
Bob		0.5		0.3
Carol	0.2		1	
David				0.4

In reality, the vast majority of entries are missing

Goal: fill in the missing entries!

Metric: mean squared error

Questions?