

ORIE 4741: Learning with Big Messy Data

Introduction

Professor Udell

Operations Research and Information Engineering
Cornell

October 16, 2021

Outline

Logistics

Stories

Definitions

Kinds of learning

Syllabus

ORIE 4741/5741: Learning with Big Messy Data

want to take this class?

- ▶ **ASAP:**

- ▶ enroll (or drop) (or get on wait list)
- ▶ fill out course survey
- ▶ sign up for discussion forum
- ▶ sign up for iClicker

- ▶ **Thursday 9/2/2021:** homework 0

links on course website:

<https://people.orie.cornell.edu/mru8/orie4741/>

Course staff

- ▶ Prof. Madeleine Udell
- ▶ TA: Richard Phillips (CS PhD)
- ▶ TA: Tao Jiang (ORIE PhD)
- ▶ TA: Connor Lawless (ORIE PhD)
- ▶ TA: Yuanping Du (ORIE MEng)
- ▶ TA: Max de Ledebur (ORIE MEng)
- ▶ TA: Jody Zhu (ORIE+CS Undergraduate)
- ▶ TA: Tara Khanna (ORIE Undergraduate)
- ▶ TA: Kevin Jiang (CS Undergraduate)

Tech stack

- ▶ In person or Zoom for lectures, section, and office hours
- ▶ Course website for course materials (syllabus, schedule, homework, project, etc)
- ▶ iClicker for polls
- ▶ Zulip for Q&A
- ▶ Gradescope for quizzes, submitting homework, grades, solutions
- ▶ Github for code (demos, projects, and hw starter code)

Course requirements and grading

course website:

(grading, course requirements, lectures, homework, etc)

<https://people.orie.cornell.edu/mru8/orie4741/>

- ▶ (15%) Participation: for every lecture (after this one), use
 - ▶ iClicker for sync lectures
 - ▶ participation form for async lectures
- ▶ (30%) Homework
 - ▶ due every two weeks or so
 - ▶ first one due next Thursday
- ▶ (15%) Quizzes
 - ▶ 30 min quiz every week or so
- ▶ (40%) Project

Course requirements and grading

course website:

(grading, course requirements, lectures, homework, etc)

<https://people.orie.cornell.edu/mru8/orie4741/>

- ▶ (15%) Participation: for every lecture (after this one), use
 - ▶ iClicker for sync lectures
 - ▶ participation form for async lectures
- ▶ (30%) Homework
 - ▶ due every two weeks or so
 - ▶ first one due next Thursday
- ▶ (15%) Quizzes
 - ▶ 30 min quiz every week or so
- ▶ (40%) Project

FAQ:

- ▶ yes, you can take the class online (even async)
- ▶ yes, you can take section online (even async), or not take the section

ORIE 5741 vs 4741

- ▶ 5741 has same course material as 4741

ORIE 5741 vs 4741

- ▶ 5741 has same course material as 4741
- ▶ expectations and rubric for course project differ
 - ▶ more business-oriented project
 - ▶ more detailed problem formulation
 - ▶ project presentation required (in addition to report)

ORIE 5741 vs 4741

- ▶ 5741 has same course material as 4741
- ▶ expectations and rubric for course project differ
 - ▶ more business-oriented project
 - ▶ more detailed problem formulation
 - ▶ project presentation required (in addition to report)
- ▶ rubric will reflect MEng learning outcomes:
 1. Mastery and Application of Core Disciplinary Knowledge
 2. Problem Formulation and Organization and Planning of the Solution Process
 3. Collaborative Problem Solving and Issue Resolution
 4. Communication of Knowledge, Ideas, and Decision Justification
 5. Self-Directed Learning and Professional Development

ORIE 5741 vs 4741

- ▶ 5741 has same course material as 4741
- ▶ expectations and rubric for course project differ
 - ▶ more business-oriented project
 - ▶ more detailed problem formulation
 - ▶ project presentation required (in addition to report)
- ▶ rubric will reflect MEng learning outcomes:
 1. Mastery and Application of Core Disciplinary Knowledge
 2. Problem Formulation and Organization and Planning of the Solution Process
 3. Collaborative Problem Solving and Issue Resolution
 4. Communication of Knowledge, Ideas, and Decision Justification
 5. Self-Directed Learning and Professional Development
- ▶ Everyone in a project group with 5741 students will be graded according to 5741 rubric.

Questions

during lecture:

- ▶ ask out loud
- ▶ zoom chat (to everyone, or to a TA)

outside of lecture:

- ▶ ask at office hours
- ▶ ask on discussion forum
- ▶ don't send email

Outline

Logistics

Stories

Definitions

Kinds of learning

Syllabus

Oh, you work with big messy data? Maybe you could help us out...?

My career in big data

academic

- ▶ B.S. in Mathematics and Physics at Yale
- ▶ Ph.D. in Computational and Mathematical Engineering at Stanford
- ▶ postdoctoral fellow at the Center for the Mathematics of Information at Caltech
- ▶ professor in ORIE at Cornell

My career in big data

academic

- ▶ B.S. in Mathematics and Physics at Yale
- ▶ Ph.D. in Computational and Mathematical Engineering at Stanford
- ▶ postdoctoral fellow at the Center for the Mathematics of Information at Caltech
- ▶ professor in ORIE at Cornell

applied work

- ▶ finance: Goldman Sachs, BlackRock, Capital One, Schonfeld, Two Sigma, . . .
- ▶ tech: Google, Retina.ai, Marketing Attribution
- ▶ cybersecurity: DARPA, Expanse (formerly Qadium)
- ▶ healthcare: Apixio, Ontario
- ▶ clean energy: Aurora
- ▶ politics: Obama 2012

Data table: politics

age	gender	state	income	education	voted?	support	...
29	F	CT	\$53,000	college	yes	Biden	...
57	?	NY	\$19,000	high school	yes	?	...
?	M	CA	\$102,000	masters	no	Trump	...
41	F	NV	\$23,000	?	yes	Trump	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

Data table: politics

age	gender	state	income	education	voted?	support	...
29	F	CT	\$53,000	college	yes	Biden	...
57	?	NY	\$19,000	high school	yes	?	...
?	M	CA	\$102,000	masters	no	Trump	...
41	F	NV	\$23,000	?	yes	Trump	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

goals:

- ▶ detect demographic groups?
- ▶ find typical responses?
- ▶ identify related features?
- ▶ impute missing entries?

Medicine

Noble, James - Armstrong - SRS EHR v6.3.13.0 - Workstation AMANDA

File Drawers View Mail Reports Tools Help

Appointments

6/20/2012 | Armstrong, Max: Multiple Loca

Time	Name	Appointment Type
09:00 AM	Noble, James	Consult
09:30 AM	Carpen, Carolyn	Exam - New Patient
10:00 AM	Harris, Richard	Follow Up
10:30 AM	Meil, Leslie	Exam - Est. Patient
10:30 AM	Moss, Pete	Consult
10:45 AM	Copeland, Elizabeth	Consult
11:00 AM	Deak, Marlin	Exam - New Patient
01:00 PM	Palhow, May	Consult
01:15 PM	Gold, Alan	Exam - Est. Patient
01:30 PM	Dipiero, Drew	Consult
01:30 PM	Glass, Tyler	Exam - Est. Patient
02:15 PM	Thompson, Brian	Exam - Est. Patient
02:30 PM	Carle, James	Consult
03:30 PM	Albani, Daniela	Follow Up
04:00 PM	Brown, Kevin	Exam - Est. Patient
04:15 PM	Newsome, Gina	Exam - Est. Patient
04:30 PM	Newsome, Jenna	Exam - Est. Patient

Mail Status

Normal	11
Rx	2
Transcription	1

Clinical Summary

Noble, James



Demographics

Address: 4516 West Huron Street
Chicago, IL 60607

Email: noblej@comcast.com
DOB: 04.05.1962
Age: 50y
Patient ID: 88501

Primary Ins. Aetna U.S. Healthcare - Master
Pharmacy: D/SigPharmacy #86871
510 College Mall Pk
Bloomington, IN 47401 P(812) 336-7306
F(812) 335-9347

User Defined Fields

Name	Value
IBNY	Shoulder pain
Employer Name	WBI
Employer Fax	201-867-5309
Insurance Carrier	QBCO
Adjuster Fax Num...	201-555-8735
Diagnosis	727.3
Position	Truck Driver
CareTracker ID #	36817221

Diagnosis

Diagnosis	Status	Date	ICD-9	Notes
ACL Tear		06.20.2012	844.2	
Lumbago		11.17.2011	724.2	

Vitals

Smoking Status

Status	ID	Date
never smoker	4	06.18.2012

Procedures

73000	73090	99211	99212	99213	99214	99215
Procedure	Date	CPT	Type	Status		
Knee Arthroscopy/Surg	06.20.2012	29581	Internal	Performed		

Rx History

Status	Date	Drug	Strength	Instructions
✓	06.18.2012	Keflex	250 MG	1 tab po BID
✓	11.04.2010	CeleBREX	100 MG	1 tablet by mouth in the AM
✗	04.13.2006	OxyCODONE/HCl	20 MG	1 tab po q-4 hr pain

Non-Drug Allergies

Description	Reaction	Notes
Latex	Severe rash	
Shellfish	Hives	

Family History

Relationship	Deceased	Notes
Sister	No	JFA
Maternal Grandmother	Yes	Osteoarthritis

Surgeries

Description	Date	Surgeon	Notes
Arthroscopy	12.16.2011	Armstrong	N/A

Appointments

Date	Time	Doctor	Reason	Type	Location	Notes
06.20.2012	09:00 AM	Max Armstrong	Shoulder		QBCO Main	

Current encounter: 6/20/2012 9:00:00 AM Transfer Encounter

Noble, James Desktop Message Ctr Scan Place Forms

Chat Notes	Radiology Reports	History	Digital X-Rays	Correspondence Received	Consult Letters Sent	Orders	Injections
Op-Reports	Discharge Summaries	PT Assessments	Charge Capture	SRS PACS	Surgery Documents	Surgery	DASH

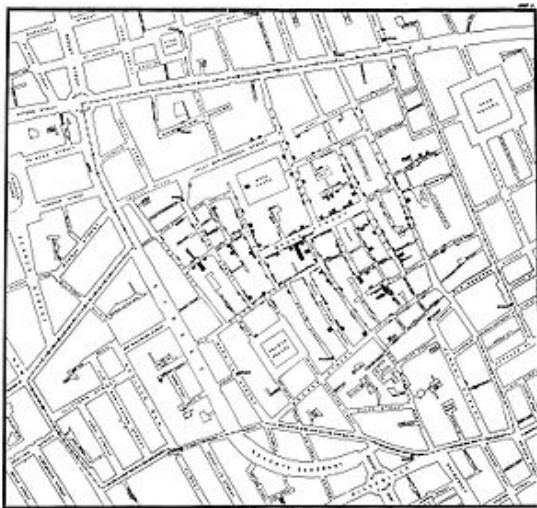
13 / 38

Data table: medicine

age	gender	heart disease	statins?	...
29	F	yes	no	...
57	?	no	no	...
?	M	no	no	...
41	F	yes	yes	...
⋮	⋮	⋮	⋮	

- ▶ find similar patients?
- ▶ understand systemic healthcare needs?
- ▶ use symptoms to detect which patients have COVID-19?
- ▶ detect patients who had series of mini-strokes?

Pollution



[Snow, 1854]

Pollution

location	time	CO2	O2	O3	...
1	1	.7	.9	?	...
1	2	.5	.7	?	...
1	3	.4	.5	1.4	...
⋮	⋮	⋮	⋮		

Marketing

Today's Recommendations For You

Here's a daily sample of items recommended for you. Click here to see all recommendations.

The Little Boy Who Swam to France

Amazon.com has these recommendations for you based on items you purchased or sold on our site.

Today's Recommendations For You

Here's a daily sample of items recommended for you. Click here to see all recommendations.

The Little Boy Who Swam to France

Amazon.com has these recommendations for you based on items you purchased or sold on our site.

Recommended For You

Amazon.com has these recommendations for you based on items you purchased or sold on our site.

The Little Boy Who Swam to France

Amazon.com has these recommendations for you based on items you purchased or sold on our site.

Recommended For You

Amazon.com has these recommendations for you based on items you purchased or sold on our site.

The Little Boy Who Swam to France

Amazon.com has these recommendations for you based on items you purchased or sold on our site.

Recommended For You

Amazon.com has these recommendations for you based on items you purchased or sold on our site.

The Little Boy Who Swam to France

Amazon.com has these recommendations for you based on items you purchased or sold on our site.

Recommended For You

Amazon.com has these recommendations for you based on items you purchased or sold on our site.

The Little Boy Who Swam to France

Amazon.com has these recommendations for you based on items you purchased or sold on our site.

Recommended For You

Amazon.com has these recommendations for you based on items you purchased or sold on our site.

The Little Boy Who Swam to France

Amazon.com has these recommendations for you based on items you purchased or sold on our site.

Customers Who Bought This Item Also Bought

The Little Boy Who Swam to France

Amazon.com has these recommendations for you based on items you purchased or sold on our site.

Today's Recommendations For You

Here's a daily sample of items recommended for you. Click here to see all recommendations.

The Little Boy Who Swam to France

Amazon.com has these recommendations for you based on items you purchased or sold on our site.

Today's Recommendations For You

Here's a daily sample of items recommended for you. Click here to see all recommendations.

The Little Boy Who Swam to France

Amazon.com has these recommendations for you based on items you purchased or sold on our site.

Personalized Features in Amazon

The Little Boy Who Swam to France

Amazon.com has these recommendations for you based on items you purchased or sold on our site.

Customers Who Bought This Item Also Bought

The Little Boy Who Swam to France

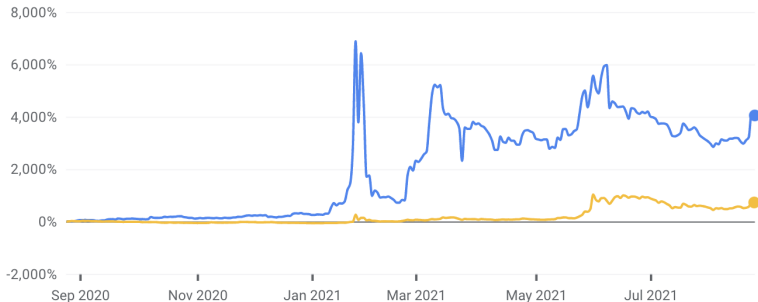
Amazon.com has these recommendations for you based on items you purchased or sold on our site.

Marketing

customer	product 1	product 2	product 3	...
1	yes	?	yes	...
2	yes	yes	?	...
3	?	?	yes	...
⋮	⋮	⋮	⋮	⋮

Finance

1D 5D 1M 6M YTD 1Y 5Y MAX



GameStop Corp.

\$206.59

+\$201.61

↑ 4,048.39%

AMC Entertainmen...

\$45.71

+\$40.17

↑ 725.09%



Finance

ticker	t_1	t_2	\dots
AAPL	.05	-.21	\dots
GOOG	-.11	.24	\dots
FB	.07	-.18	\dots
\vdots	\vdots	\vdots	\ddots

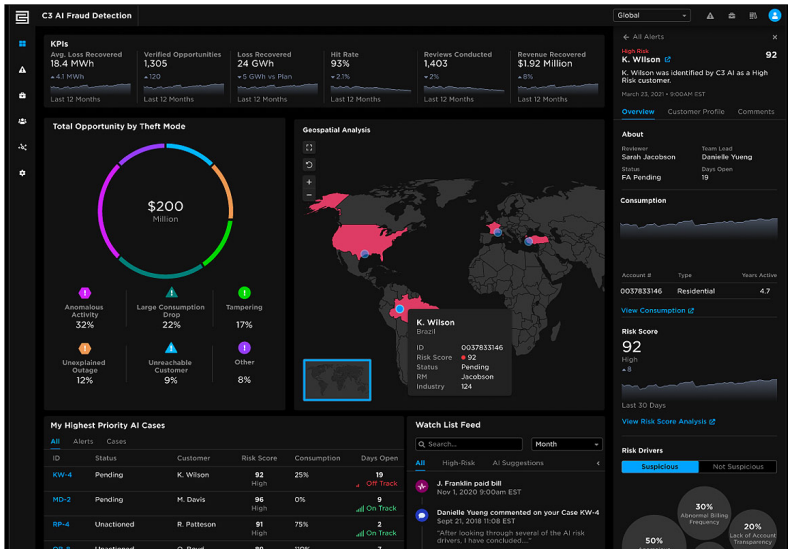
Environmental, social and governance (ESG) data

- ▶ one row for each asset at each time
- ▶ one column per key performance indicator (KPI)
 - ▶ carbon emissions
 - ▶ e-waste management
 - ▶ climate change risk
 - ▶ worker safety
 - ▶ ...
- ▶ values: numerical ratings $1, \dots, 10$ or boolean $\{0, 1\}$
- ▶ triangular missing pattern: KPI/asset coverage increases with time

goals for ESG analysis:

- ▶ impute missing items?
- ▶ audit/improve on vendor data quality?
- ▶ predict long term returns?

Fraud detection



Autocomplete

Your AI pair programmer

With GitHub Copilot, get suggestions for whole lines or entire functions right inside your editor.

Sign up >

```
sentiment.ts  write_sql.go  parse_expenses.py  addresses.rb

1 #!/usr/bin/env ts-node
2
3 import { fetch } from "fetch-h2";
4
5 // Determine whether the sentiment of text is positive
6 // Use a web service
7 async function isPositive(text: string): Promise<boolean> {
8   const response = await fetch(`http://text-processing.com/api/sentiment/`, {
9     method: "POST",
10    body: `text=${text}`,
11    headers: {
12      "Content-Type": "application/x-www-form-urlencoded",
13    },
14  });
15  const json = await response.json();
16  return json.label === "pos";
17 }
```

Application areas

- ▶ health
- ▶ politics
- ▶ governance
- ▶ advertising
- ▶ retail
- ▶ ecommerce
- ▶ finance
- ▶ ...

Outline

Logistics

Stories

Definitions

Kinds of learning

Syllabus

Big

- ▶ NASA, 1997: “taxing the capacities of main memory, local disk, and even remote disk”

¹image courtesy of Kim Minor © IBM

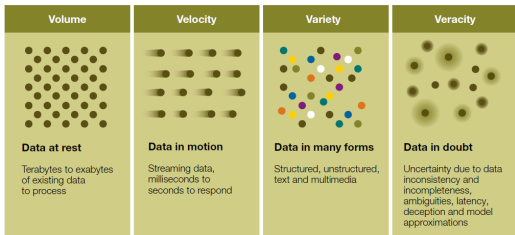
Big

- ▶ NASA, 1997: “taxing the capacities of main memory, local disk, and even remote disk”
- ▶ OED, 2015: “data of a very large size, typically to the extent that its manipulation and management present significant logistical challenges”

¹image courtesy of Kim Minor @ IBM

Big

- ▶ NASA, 1997: “taxing the capacities of main memory, local disk, and even remote disk”
- ▶ OED, 2015: “data of a very large size, typically to the extent that its manipulation and management present significant logistical challenges”
- ▶ 4 Vs:

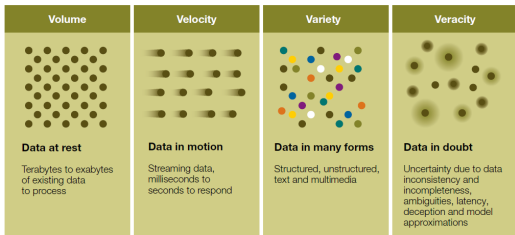


1

¹image courtesy of Kim Minor @ IBM

Big

- ▶ NASA, 1997: “taxing the capacities of main memory, local disk, and even remote disk”
- ▶ OED, 2015: “data of a very large size, typically to the extent that its manipulation and management present significant logistical challenges”
- ▶ 4 Vs:



1

- ▶ 5th V: value

¹image courtesy of Kim Minor @ IBM

Big: our definition

Definition

An algorithm for **big data** is one with computational and memory requirements that scale **linearly** (or nearly linearly) in the size of the data.

Big: our definition

Definition

An algorithm for **big data** is one with computational and memory requirements that scale **linearly** (or nearly linearly) in the size of the data.

why this definition? independent of

- ▶ hardware
- ▶ business

Big: our definition

Definition

An algorithm for **big data** is one with computational and memory requirements that scale **linearly** (or nearly linearly) in the size of the data.

why this definition? independent of

- ▶ hardware
- ▶ business

if you use only algorithms for **big data**, then you're working with **big data**

Messy

- ▶ noisy: some (or all) values suffer errors, inaccuracies, or malicious corruption

Messy

- ▶ noisy: some (or all) values suffer errors, inaccuracies, or malicious corruption
- ▶ missing: some values are missing, inconsistent, not recorded, or lost

Messy

- ▶ noisy: some (or all) values suffer errors, inaccuracies, or malicious corruption
- ▶ missing: some values are missing, inconsistent, not recorded, or lost
- ▶ heterogeneous: values of many different types
 - ▶ continuous values (e.g., 4.2, π)
 - ▶ discrete values (e.g., 0, 4, 994)
 - ▶ nominal values (e.g., apple, banana, pear)
 - ▶ ordinal values (e.g., rarely, sometimes, often)
 - ▶ graphs or networks (e.g., person 1 is friends with person 2)
 - ▶ text (e.g., doctor's note describing symptoms)
 - ▶ sets (e.g., items purchased)

Learning

Learning

- ▶ machine learning?

Learning

- ▶ machine learning?
- ▶ human learning?

Learning

- ▶ machine learning?
- ▶ human learning?
- ▶ when data is **big** and **messy**,
machine help is essential for human learning!

Outline

Logistics

Stories

Definitions

Kinds of learning

Syllabus

Supervised learning

- ▶ identify one column of data that we want to predict

$$\begin{bmatrix} A \end{bmatrix} = \begin{bmatrix} x_{11} & \cdots & x_{1d-1} & y_1 \\ \vdots & \ddots & \vdots & \vdots \\ x_{n1} & \cdots & x_{nd-1} & y_n \end{bmatrix} = \begin{bmatrix} X & y \end{bmatrix}$$

- ▶ $x_i \in \mathcal{X}$ for $i = 1, \dots, n$ are rows of X
- ▶ $y_i \in \mathcal{Y}$ for $i = 1, \dots, n$ are entries of y

Supervised learning

- ▶ identify one column of data that we want to predict

$$\begin{bmatrix} A \end{bmatrix} = \begin{bmatrix} x_{11} & \cdots & x_{1d-1} & y_1 \\ \vdots & \ddots & \vdots & \vdots \\ x_{n1} & \cdots & x_{nd-1} & y_n \end{bmatrix} = \begin{bmatrix} X & y \end{bmatrix}$$

- ▶ $x_i \in \mathcal{X}$ for $i = 1, \dots, n$ are rows of X
- ▶ $y_i \in \mathcal{Y}$ for $i = 1, \dots, n$ are entries of y
- ▶ we believe there is a mapping $f : \mathcal{X} \rightarrow \mathcal{Y}$

$$y_i \approx f(x_i)$$

- ▶ our goal is to learn f

Example: supervised learning for credit decisioning

- ▶ goal: decide which credit card applicants should be approved
- ▶ input space: entries of $\mathcal{X} \in \mathbf{R}^d$ correspond to fields in credit application
 - ▶ e.g., salary, years in residence, outstanding debt, number of credit lines, ...
- ▶ output space: $\mathcal{Y} = \{+1, -1\}$
 - ▶ +1 means approve
 - ▶ -1 means reject
- ▶ data: $\mathcal{D} = (x_1, y_1), \dots, (x_n, y_n)$
applications of previous customers, and credit approval decisions made by humans

Example: supervised learning for credit decisioning

- ▶ goal: decide which credit card applicants should be approved
- ▶ input space: entries of $\mathcal{X} \in \mathbf{R}^d$ correspond to fields in credit application
 - ▶ e.g., salary, years in residence, outstanding debt, number of credit lines, ...
- ▶ output space: $\mathcal{Y} = \{+1, -1\}$
 - ▶ +1 means approve
 - ▶ -1 means reject
- ▶ data: $\mathcal{D} = (x_1, y_1), \dots, (x_n, y_n)$
applications of previous customers, and credit approval decisions made by humans

Q: what are potential problems with using a model built with this data?

Example: supervised learning for credit decisioning

- ▶ goal: decide which credit card applicants should be approved
- ▶ input space: entries of $\mathcal{X} \in \mathbf{R}^d$ correspond to fields in credit application
 - ▶ e.g., salary, years in residence, outstanding debt, number of credit lines, ...
- ▶ output space: $\mathcal{Y} = \{+1, -1\}$
 - ▶ +1 means approve
 - ▶ -1 means reject
- ▶ data: $\mathcal{D} = (x_1, y_1), \dots, (x_n, y_n)$
applications of previous customers, and credit approval decisions made by humans

Q: what are potential problems with using a model built with this data?

A: wrong objective: human decision may not be correct decision;
covariate shift: future data may look unlike past data; ...

Exercise: formalizing real problems

- ▶ identify a prediction goal
- ▶ identify the input space \mathcal{X}
- ▶ identify the output space \mathcal{Y}
- ▶ identify the data $\mathcal{D} = (x_1, y_1), \dots, (x_n, y_n)$ you'd like to use
- ▶ what kinds of noise do you expect in the data?

Outline

Logistics

Stories

Definitions

Kinds of learning

Syllabus

Course objectives (I)

- ▶ plot
- ▶ predict
- ▶ cluster
- ▶ impute
- ▶ denoise
- ▶ recommend
- ▶ understand

Course objectives (II)

this course is about

- ▶ algorithms for big messy data
- ▶ learning to ask the right questions

at the end of the course, you should have learned

- ▶ at least one method to solve any problem
- ▶ machine learning is not magic; it's math
- ▶ when **not** to trust your solution

Course objectives (II)

this course is about

- ▶ algorithms for big messy data
- ▶ learning to ask the right questions

at the end of the course, you should have learned

- ▶ at least one method to solve any problem
- ▶ machine learning is not magic; it's math
- ▶ when **not** to trust your solution

the rest you can learn online. . .

ORIE 4741: Learning with Big Messy Data

want to take this class?

- ▶ **ASAP:**

- ▶ enroll (or drop) (or get on wait list)
- ▶ fill out course survey
- ▶ sign up for discussion forum
- ▶ sign up for iClicker

- ▶ **Thursday 9/2/2021:** homework 0

links on course website:

<https://people.orie.cornell.edu/mru8/orie4741/>