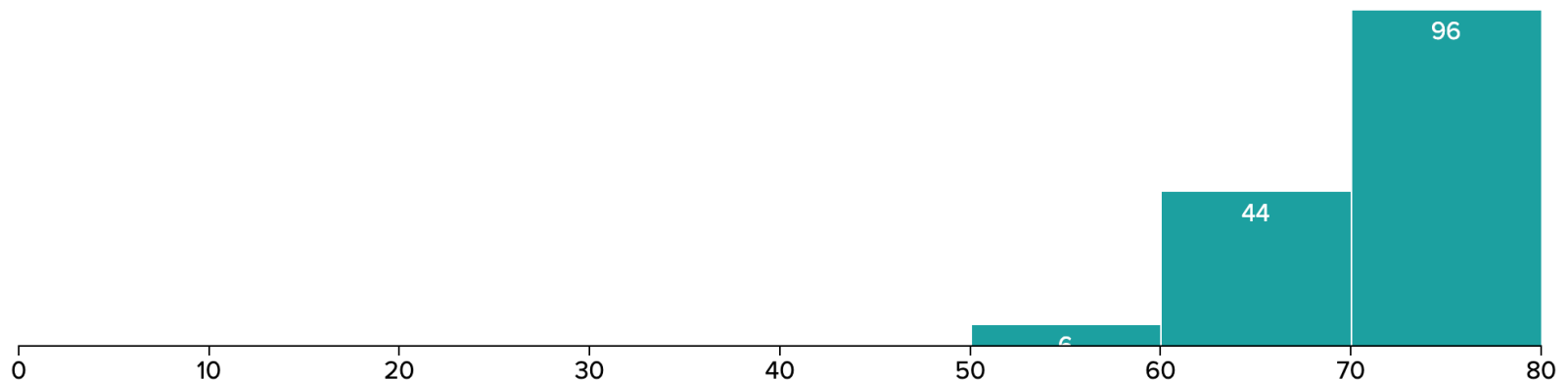# SE for ML

17-313 Fall 2022

# Administrivia

- ○ **HW4 Released**
  - ■ **3 checkpoints.  Note, for checkpoint 1, tests don't need to pass/run**
- ○ **Midterm is graded**
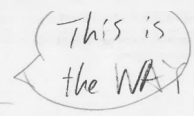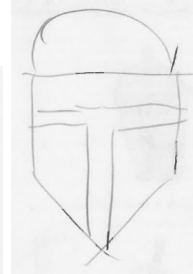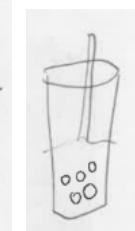


| MEDIAN | MAXIMUM | MEAN | STD DEV |
|--------|---------|------|---------|
| **72.0** | **80.0** | **71.01** | **5.33** |

I love
cool!! 0 0

(boba)

Moo!

uooo midterm finished!

17313
software
Engineering

→ only boba
boba drink

sun

(thx for the candy)

idk:/

This is
the WAY

MILK

3D

# Retrospectives

- "the purpose of the Sprint Retrospective is to plan ways to increase quality and effectiveness." –Scrum.org

- We often use three questions:
- What should we:
  - Start doing?
  - Stop doing?
  - Keep doing?

# Learning goals

- Identify differences between traditional software development and development of ML systems.
- Understand the stages that comprise the typical ML development pipeline.
- Identify challenges that must be faced within each stage of the typical ML development pipeline.

S3D

Quick poll:
**Have you taken a machine learning course before?**

S3D

# Machine Learning in One Slide

(Supervised)



Lots of labelled data
(Inputs, outputs)

Training

Model

Input

Output

"Bird"

Input

Output

"Bird"

# Traditional Software Development

"It is easy. You just chip away the stone that doesn't look like David." –(probably not) Michelangelo

# ML Development

- Observation
- Hypothesis
- Predict
- Test
- Reject or Refine Hypothesis



S3D

# Black-box View
# of Machine Learning



Image: https://xkcd.com/1838/

# Microsoft's view of Software Engineering for ML

# Three Fundamental Differences:

- Data discovery and management

- Customization and Reuse

- No modular development of model itself

S3D

# Case Study

- Case study developed by
- Christian Kästner
- https://ckaestne.github.io/seai/



**Machine Learning in Production / AI Engineering (17-445/17-645/17-745/11-695)**

*Formerly **Software Engineering for AI-Enabled Systems (SE4AI)**, CMU course that covers how to build, deploy, assure, and maintain applications with machine-learned models. Covers **responsible AI** (safety, security, fairness, explainability, …) and **MLOps**.*

**Fundamentals of Engineering AI-Enabled Systems**

Holistic system view: AI and non-AI components, pipelines, stakeholders, environment interactions, feedback loops

| **Requirements:** | **Architecture + design:** | **Quality assurance:** | **Operations:** |
|---|---|---|---|
| System and model goals | Modeling tradeoffs | Model testing | Continuous deployment |
| User requirements | Deployment architecture | Data quality | Contin. experimentation |
| Environment assumptions | Data science pipelines | QA automation | Configuration mgmt. |
| Quality beyond accuracy | Telemetry, monitoring | Testing in production | Monitoring |
| Measurement | Anticipating evolution | Infrastructure quality | Versioning |
| Risk analysis | Big data processing | Debugging | Big data |
| Planning for mistakes | Human-AI design | | DevOps, MLOps |

Teams and process: Data science vs software eng. workflows, interdisciplinary teams, collaboration points, technical debt

**Responsible AI Engineering**

| Provenance, versioning, reproducibility | Safety | Security and privacy | Fairness | Interpretability and explainability | Transparency and trust |
|---|---|---|---|---|---|

Ethics, governance, regulation, compliance, organizational culture

S3D

The Next Generation of Spectacles

# Qualities of Interest?



A



B



The Next Generation of Spectacles

C

# Typical ML Pipeline

- Static
    - Get labeled data (data collection, cleaning and, labeling)
    - Identify and extract features (feature engineering)
    - Split data into training and evaluation set
    - Learn model from training data (model training)
    - Evaluate model on evaluation data (model evaluation)
    - Repeat, revising features
- with production data
    - Evaluate model on production data; monitor (model monitoring)
    - Select production data for retraining (model training + evaluation)
    - Update model regularly (model deployment)

# Example Data

# Learning Data

似乎格式有問題



English output

**parallel corpus**

网站资讯分析网数据显示的主域名为全世界访问量最高的站点除此之外搜索在其他国家或地区域名下的多个站点等等及旗下的等

The corporation has been estim to run more than one million pag in data centers around the world to process over one billion searc requests and about twenty-four i of user-generated data each dat December 2012 Alexa listed as

**monolingual corpus**

started functioning in 1928 and established the tradition of large exhibitions and trade fairs held in Brno, and nowadays also ranks among the sights of the city. Brno is also known for hosting big motorbike and other races on the Masaryk Circuit, a tradition established in 1930 in which the Road Racing World Championship Grand Prix is one of the most prestigious races. Another notable cultural tradition is an international fireworks competition.

# Example Data

| UserId | PickupLocation | TargetLocation | OrderTime | PickupTime |
|--------|----------------|----------------|-----------|------------|
| 5      | ….             | …              | 18:23     | 18:31      |
| …      |                |                |           |            |
|        |                |                |           |            |
|        |                |                |           |            |
|        |                |                |           |            |
|        |                |                |           |            |
|        |                |                |           |            |
|        |                |                |           |            |

# Feature Engineering

- Identify parameters of interest that a model may learn on
- Convert data into a useful form
- Normalize data
- Include context
- Remove misleading things

**S3D**

**Features?**

# Feature Extraction

- In OCR/translation:
    - Bounding boxes for text of interest
    - Character boundaries
    - Line segments for each character
    - GPS location of phone (to determine likely source language)

S3D

Features?

# Feature Extraction

- ## In surge prediction:
  - Location and time of past surges
  - Events
  - Number of people traveling to an area
  - Typical demand curves in an area
  - Demand in other areas
  - Weather

S3D

## Data Cleaning

- Removing outliers
- Normalizing data
- Missing values

- ...

S3D

**Learning**

- Build a predictor that best describes an outcome for the observed features

S3D

# Evaluation

- Prediction accuracy on learned data vs
- Prediction accuracy on unseen data
  - Separate learning set, not used for training

- For binary predictors: false positives vs. false negatives, precision vs. recall
- For numeric predictors: average (relative) distance between real and predicted value
- For ranking predictors: top-K, etc.

S3D

Evaluation Data and Metrics?

Evaluation Data and Metrics?

# Learning and Evaluating in Production

- Beyond static data sets, **build telemetry**
- Design challenge: identify mistakes in practice

- Use sample of live data for evaluation
- Retrain models with sampled live data regularly
- Monitor performance and intervene

S3D

# ML Model Tradeoffs

- Accuracy
- Capabilities (e.g. classification, recommendation, clustering…)
- Amount of training data needed
- Inference latency
- Learning latency; incremental learning?
- Model size
- Explainable? Robust?
- …

S3D

# Where should the model live?

Glasses

Phone

Cloud

OCR Component

Translation Component

# Where should the model live?

Vehicle

Phone

Cloud

Surge
Prediction

S3D

# Considerations

- How much data is needed as input for the model?
- How much output data is produced by the model?
- How fast/energy consuming is model execution?
- What latency is needed for the application?
- How big is the model? How often does it need to be updated?
- Cost of operating the model? (distribution + execution)
- Opportunities for telemetry?
- What happens if users are offline?

S3D

# Typical Designs

- Static intelligence in the product
  - difficult to update
  - good execution latency
  - cheap operation
  - offline operation
  - no telemetry to evaluate and improve
- Client-side intelligence
  - updates costly/slow, out of sync problems
  - complexity in clients
  - offline operation, low execution latency

S3D

# Typical Designs

- Server-centric intelligence
  - latency in model execution (remote calls)
  - easy to update and experiment
  - operation cost
  - no offline operation
- Back-end cached intelligence
  - precomputed common results
  - fast execution, partial offline
  - saves bandwidth, complicated updates
- Hybrid models

S3D

# Other Considerations

- Coupling of ML pipeline parts
- Coupling with other parts of the system
- Ability for different developers and analysists to collaborate
- Support online experiments
- Ability to monitor

S3D

# Reactive Systems

- Responsive
    - consistent, high performance
- Resilient
    - maintain responsive in the face of failure, recovery, rollback
- Elastic
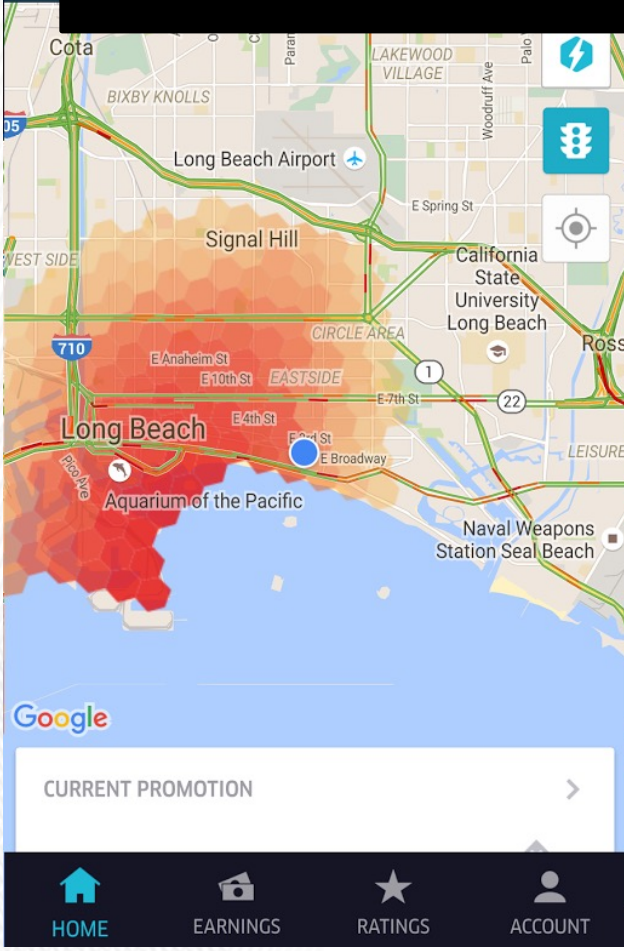    - scale with varying loads

S3D

# Updating Models

- Models are rarely static outside the lab
- Data drift, feedback loops, new features, new requirements
- When and how to update models?
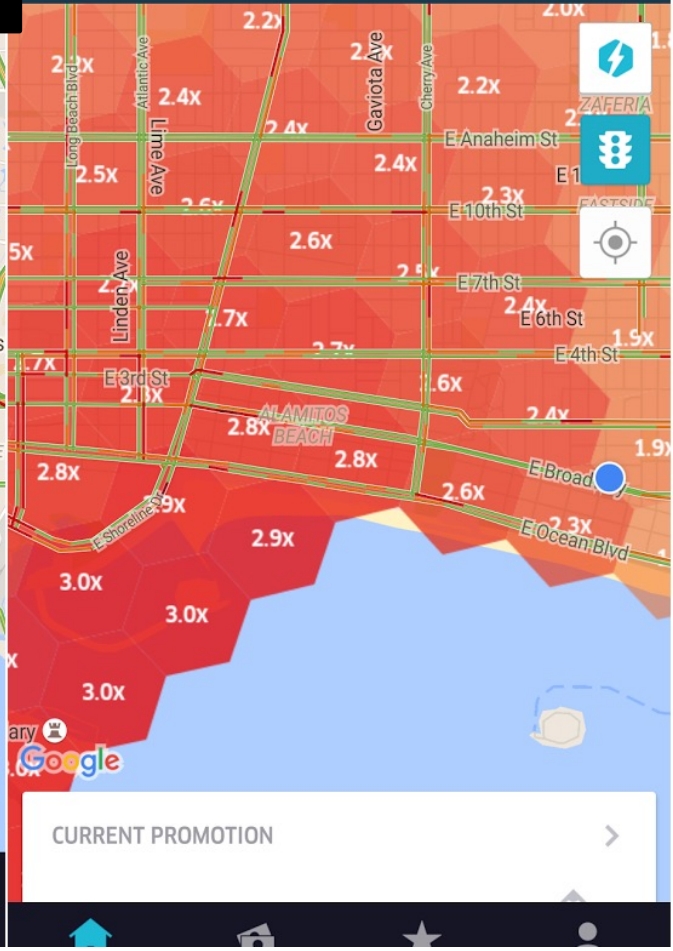- How to version? How to avoid mistakes?

S3D

**Update Strategy?**

Update Strategy?

# Mistakes will happen

- No specification
- ML components detect patterns from data (real and spurious)
- Predictions are often accurate, but mistakes always possible
- Mistakes are not predicable or explainable or similar to human mistakes
- Plan for mistakes
- Telemetry to learn about mistakes?

S3D

# How Models can Break

- System outage
- Model outage
  - model tested? deployment and updates reliable? file corrupt?
- Model errors
- Model degradation
  - data drift, feedback loops

S3D

# Hazard Analysis

- Worst thing that can happen?
- Backup strategy? Undoable? Nontechnical compensation?
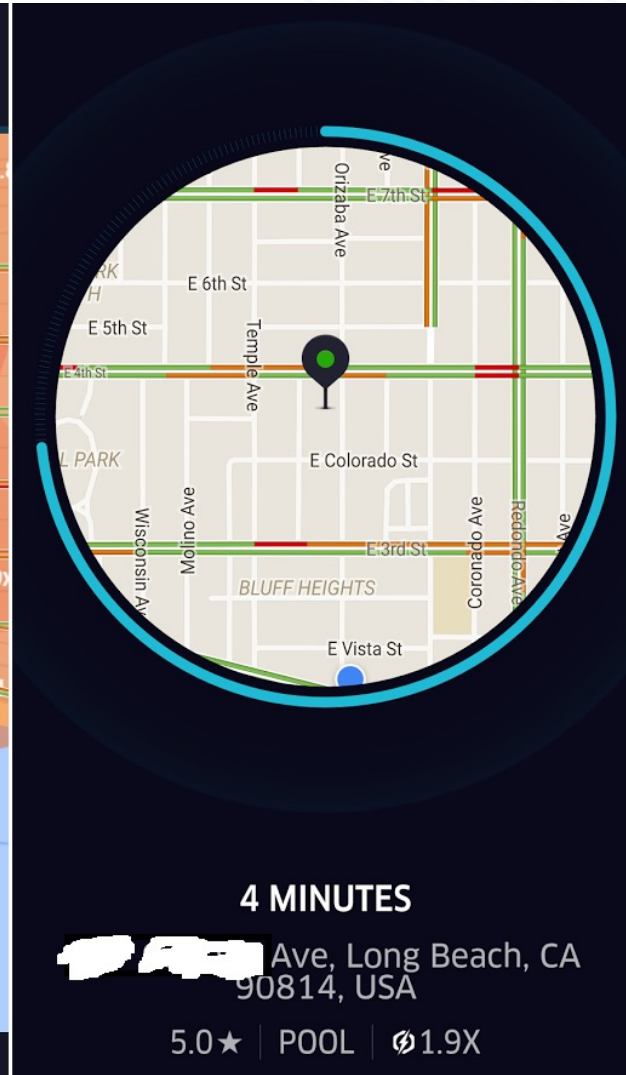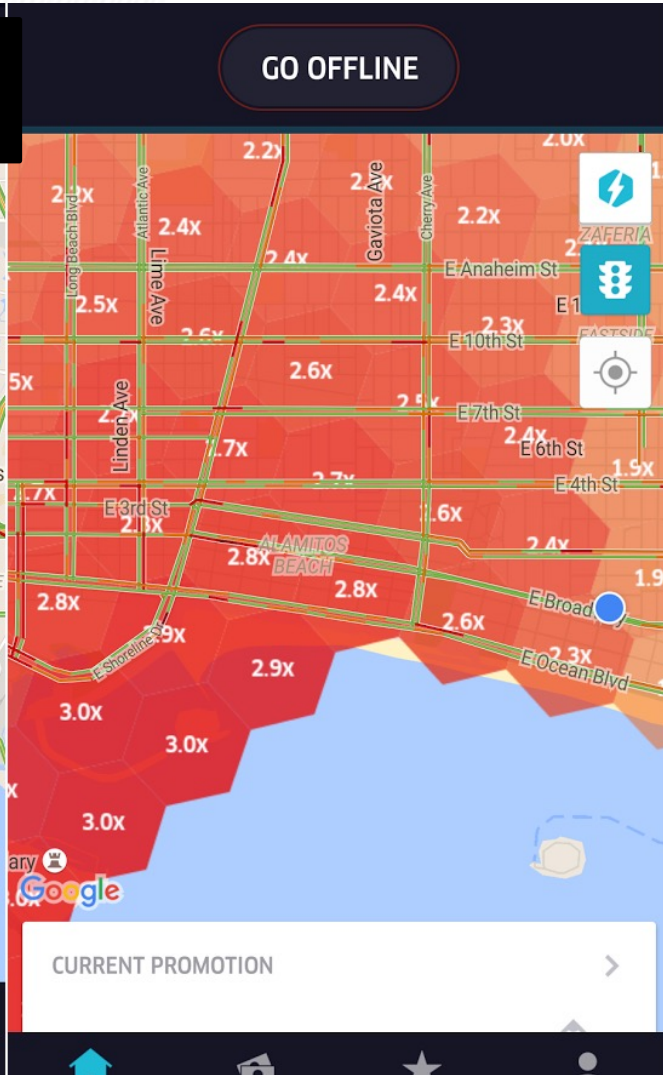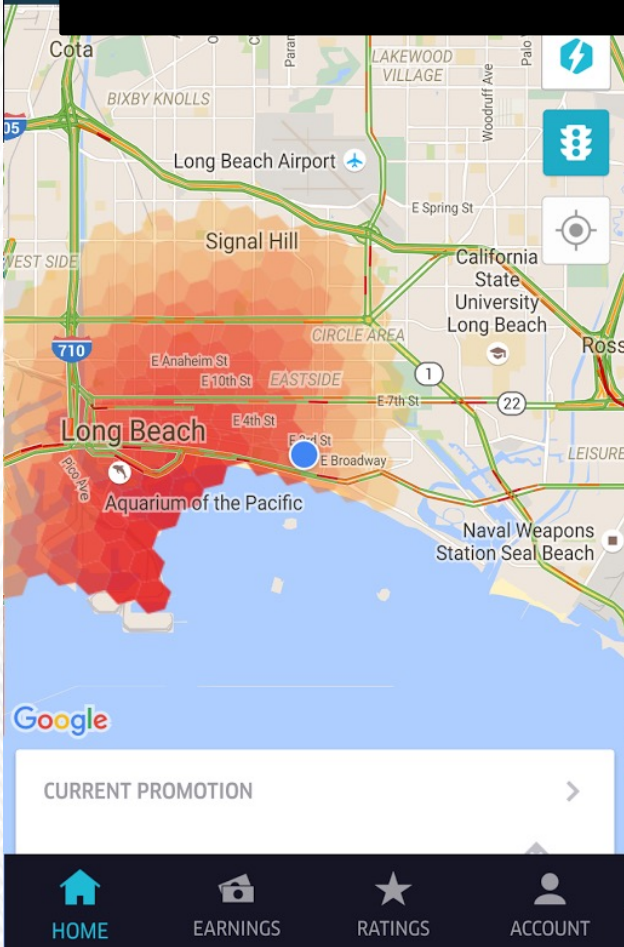
S3D

# Mitigating Mistakes

- Investigating in ML
  - e.g., more training data, better data, better features, better engineers
- Less forceful experience
  - e.g., prompt rather than automate decisions, turn off
- Adjust learning parameters
  - e.g., more frequent updates, manual adjustments
- Guardrails
  - e.g., heuristics and constraints on outputs
- Override errors
  - e.g., hardcode specific results

S3D

Mistakes?

Mistakes?

# Telemetry

- Purpose:
  - monitor operation
  - monitor success (accuracy)
  - improve models over time (e.g., detect new features)
- Challenges:
  - too much data – sample, summarization, adjustable
  - hard to measure – intended outcome not observable? proxies?
  - rare events – important but hard to capture
  - cost – significant investment must show benefit
  - privacy – abstracting data

S3D

# Requirements and estimation

- Talking to stakeholders

S3D

# Summary

- Machine learning in production systems is challenging
- Many tradeoffs in selecting ML components and in integrating them in larger system
- Plan for updates
- Manage mistakes, plan for telemetry