# QA: Different strokes for different systems.

Michael Hilton

institute for SOFTWARE RESEARCH

**Carnegie Mellon University**
School of Computer Science

# Learning goals

- Understand challenges for QA of ML systems
- Be able to test assumptions about the data
- Understand and choose different fairness approaches
- Detect changes in the model over time

# Administrivia

- Homework 4D due today
- Homework 5 released soon
- Review retrospective answers
- Visitor on Thursday

institute for SOFTWARE RESEARCH

**Carnegie Mellon University**
School of Computer Science

# What should we start doing?

- Issues with slack (QA via google docs?
- Unclear Homework questions
- Recitation Recordings
- Bigger breakout rooms
- Direct instruction on decision making process
- Docker and Kubernetes via lecture
- Piazza
- WSL ☹
- More TA's
- More technical work
- More bingo-like activities
- More Technical office hours

# What should we stop doing?

- Talking about candy we would have had without the pandemic ☹

- Google Slides DDOS, creating a document is stressful
- Less break out rooms
- More break out rooms
- Some people feel the grading is too harsh
- Homework 3 was too short (not enough time)

# What should we keep doing?

- Having in class questions. Makes me have a reason to pay extra attention in the lecture.
- Show more of the cat. :-)
- Guest speakers
- Slack
- Informal polls
- Breakout rooms
- Support of non/US East students
- Banter

institute for SOFTWARE RESEARCH

**Carnegie Mellon University**
School of Computer Science

# QA for ML

# What does it mean to do QA for a ML System?



Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small black box in the middle. The required surrounding infrastructure is vast and complex.

# Broad considerations when testing ML

- Data debugging, validation, and testing
- Model  debugging, validation, and testing
- Service debugging, validation, and testing
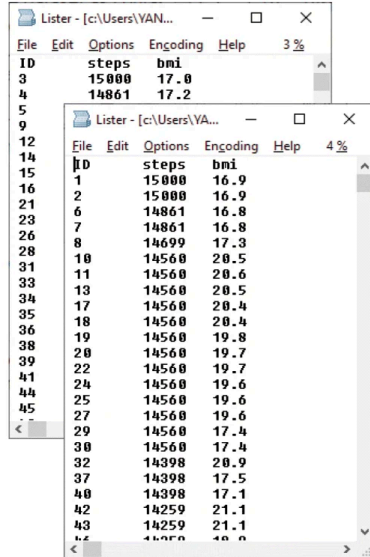  - Traditionally testing, Design docs, already covered

# Data Debugging

- Validate Input Data Using a Data Schema
  - For your feature data, understand the range and distribution. For categorical features, understand the set of possible values.
  - Encode your understanding into rules defined in the schema.
  - Test your data against the data schema.
- Test Engineered Data: For example:
  - All numeric features are scaled, for example, between 0 and 1.
  - One-hot encoded vectors only contain a single 1 and N-1 zeroes.
  - Missing data is replaced by mean or default values.
  - Data distributions after transformation conform to expectations.
  - Outliers are handled, such as by scaling or clipping.

institute for SOFTWARE RESEARCH

**Carnegie Mellon University**
School of Computer Science
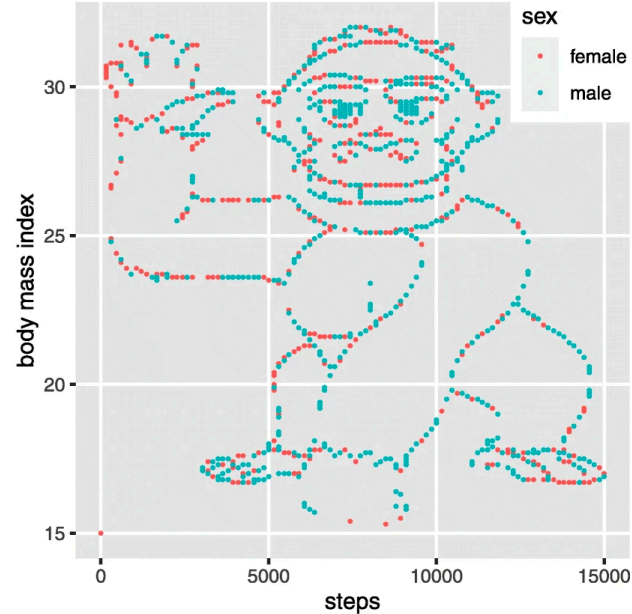
# Data Debugging Cont...

- Is your data sampled in a way that represents your users (e.g., will be used for all ages, but you only have training data from senior citizens) and the real-world setting (e.g., will be used year-round, but you only have training data from the summer

- Training-serving skew—the difference between performance during training and performance during serving. During training, try to identify potential skews and work to address them. During evaluation, continue to try to get evaluation data that is as representative as possible of the deployed setting.

- Are any features in your model redundant or unnecessary? Use the simplest model possible.

- Data bias is another important consideration

institute for SOFTWARE RESEARCH

**Carnegie Mellon University**
School of Computer Science

https://ai.google/responsibilities/responsible-ai-practices/

# Examine your data…

a



b



c

|  | Gorilla <u>not</u> discovered | Gorilla discovered |
|---|---|---|
| Hypothesis-focused | 14 | 5 |
| Hypothesis-free | 5 | 9 |

# Model Debugging

- Check that the data can predict the labels.
  - Use 10 examples from your dataset that the model can easily learn from.  Alternatively, use synthetic data.
- Establish a baseline
  - Use a linear model trained solely on most predictive feature
  - In classification, always predict the most common label
  - In regression, always predict the mean value

institute for SOFTWARE RESEARCH

**Carnegie Mellon University**
School of Computer Science

https://ai.google/responsibilities/responsible-ai-practices/

# Recommended practices

Use a human-centered design approach ⌄

Identify multiple metrics to assess training and monitoring ⌄

When possible, directly examine your raw data ⌄

Understand the limitations of your dataset and model ⌄

Test, Test, Test ⌄

Continue to monitor and update the system after deployment ⌄

# Test, Test, Test

Learn from software engineering best test practices and quality engineering to make sure the AI system is working as intended and can be trusted.

- Conduct rigorous **unit tests** to test each component of the system in isolation.

- Conduct **integration tests** to understand how individual ML components interact with other parts of the overall system.

- Proactively detect **input drift** by testing the statistics of the inputs to the AI system to make sure they are not changing in unexpected ways.

# Test, Test, Test - Continued

- Use a gold standard dataset to test the system and ensure that it **continues to behave as expected**. Update this test set regularly in line with changing users and use cases, and to reduce the likelihood of training on the test set.

- Conduct iterative user testing to incorporate a diverse set of users' needs in the development cycles.

- Apply the quality engineering principle of [poka-yoke](): build quality checks into a system, so that unintended failures either cannot happen or **trigger an immediate response** (e.g., if an important feature is unexpectedly missing, the AI system won't output a prediction).

institute for SOFTWARE RESEARCH

**Carnegie Mellon University**
School of Computer Science

# Beyond Correctness….



**Prediction probabilities**

atheism 0.59
christian 0.41

atheism          christian

Posting 0.16
Host 0.13
NNTP 0.10
edu 0.05
have 0.01
There 0.01

**Text with highlighted words**
From: johnchad@triton.unm.edu (jchadwic)
Subject: Another request for Darwin Fish
Organization: University of New Mexico, Albuquerque
Lines: 11
NNTP-Posting-Host: triton.unm.edu

Hello Gang,

There have been some notes recently asking where to obtain the DARWIN fish.
This is the same question I have and I have not seen an answer on the
net. If anyone has a contact please post on the net or email me.

Thanks,

john chadwick

**Carnegie Mellon University**
School of Computer Science

# https://github.com/drivendataorg/deon

# Where ML has gone wrong

- Data Collection
  - [Facebook uses phone numbers provided for two-factor authentication to target users with ads.](#)
  - [Personal information on taxi drivers can be accessed in poorly anonymized taxi trips dataset released by New York City.](#)
- Data Storage
  - [FedEx exposes private information of thousands of customers after a legacy s3 server was left open without a password.](#)
- Analysis
  - [A widely used commercial algorithm in the healthcare industry underestimates the care needs of black patients, assigning them lower risk scores compared to equivalently sick white patients.](#)
  - [Strava heatmap of exercise routes reveals sensitive information on military bases and spy outposts.](#)
  - [Excel error in well-known economics paper undermines justification of austerity measures.](#)
- Modeling
  - [Google Photos tags two African-Americans as gorillas.](#)
  - [Google searches involving black-sounding names are more likely to serve up ads suggestive of a criminal record than white-sounding names.](#)

https://deon.drivendata.org/examples/

# ML Fairness

- Getting answers is the easy part… Asking the right questions is the hard part.

institute for SOFTWARE RESEARCH | **Carnegie Mellon University** School of Computer Science
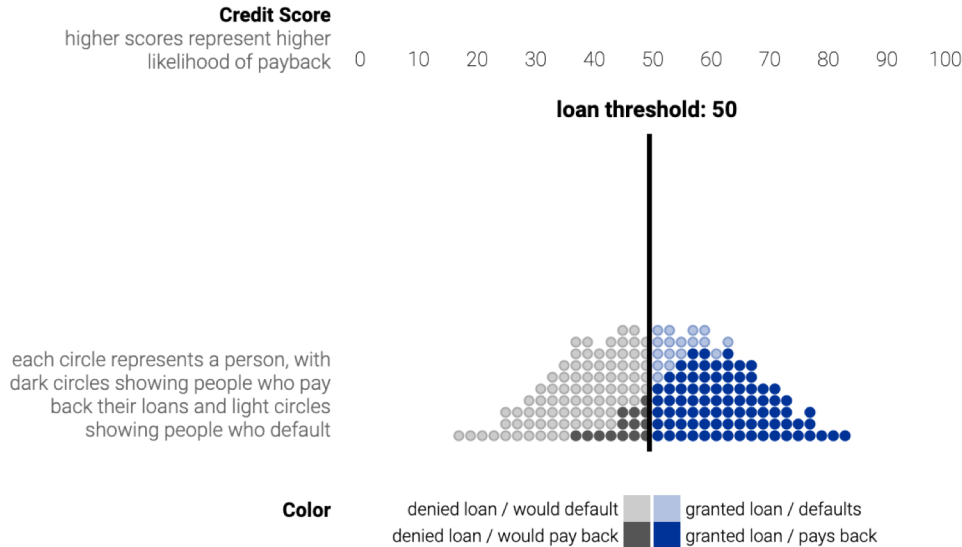
# Varieties of fairness

- Group unaware
  - Ignore group data (one group could get excluded)
- Group thresholds
  - Different rules per group (rules differ by group)
- Demographic parity
  - Same percentage in pool as outcomes (might result in random selection)
- Equal opportunity
  - Equal chance out positive outcomes regardless of groups (focus on individual, rules differ per group)
- Equal accuracy
  - Equal chance of both outcomes per group (focus on group, rules differ per group)

# Explainability



Simulating loan thresholds

Drag the black threshold bars left or right to change the cut-offs for loans.

Threshold Decision / Outcome

# Activity

Consider the different approaches to fairness. Can you come up with different scenarios where each fairness approach might be appropriate?

Remember the fairness approaches are:

- Group unaware
- Group thresholds
- Demographic parity
- Equal opportunity
- Equal accuracy