

CS 188

Scalable Internet Services

Andrew Mutz

November 21, 2019



Announcements

Today: **Intelligent Systems**

November 26, 2019 **Microservices, Containers, Kubernetes**

November 28, 2019 **Thanksgiving day, No class**

December 3, 2019 **CDNs and Course conclusion**

December 5, 2019 All papers due. Final Presentations

December 6, 2019 Final Presentations



Announcements

Papers (Dec 5) and presentations (Dec 5 and 6) are due 2 weeks from today

There will be lab tomorrow, but no lab next week (due to Thanksgiving).

If you are having trouble, make sure to reach out early and often on Piazza!



Disclaimer

- I have two years experience building production ML systems, but I am not an expert



The Machine Learning Surprise

- Why are we talking about AI & ML in a course on Scalable Internet Services?
- "The Surprising Truth About What it Takes to Build a Machine Learning Product" by Josh Cogan*

*<https://medium.com/thelaunchpad/the-ml-surprise-f54706361a6c>



The Machine Learning Surprise

Effort Allocation

Expectation

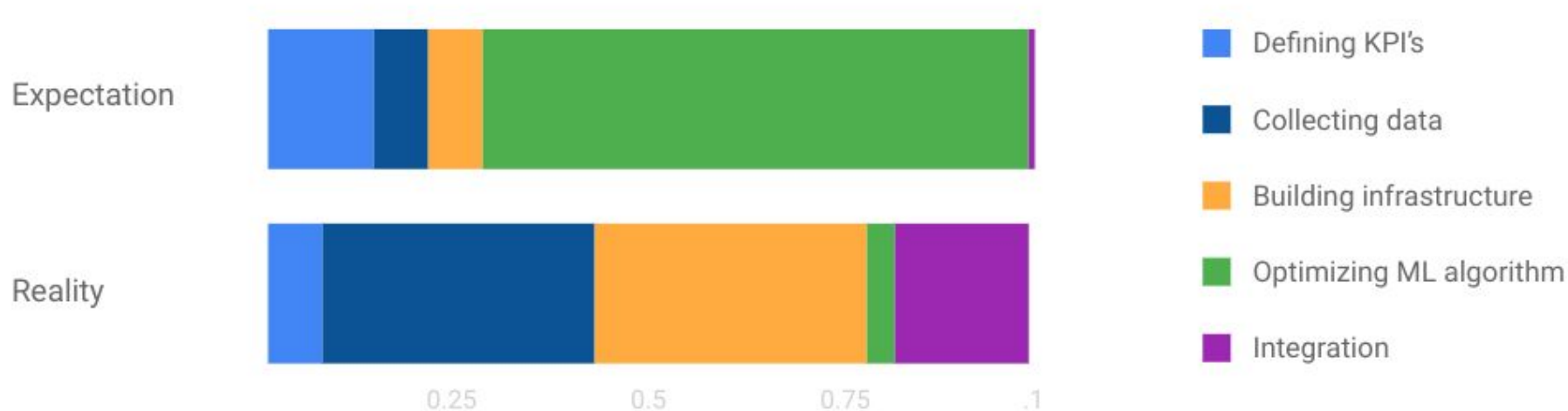


- Defining KPI's
- Collecting data
- Building infrastructure
- Optimizing ML algorithm
- Integration

*Informally based on my many conversations with new ML practitioners

The Machine Learning Surprise

Effort Allocation



Why is ML such a big deal?

- Let's take a step back before we dive in to reflect on why this is so important
- Why is ML such a widely discussed area these days?



Why is ML such a big deal?

- ML systems are delivering breakthroughs in the fundamental capabilities of software systems:
 - The ability to understand the content images and video
 - The ability to understand human natural language (spoken and written)
 - The ability to play and win a game from only a set of rules
 - The ability to generate plausible images and language
- A naive reflection on these breakthroughs can see where they can be applied:
 - Organize your photos
 - Talk to Alexa
 - Better Dota adversaries
- But this is so much more...



Why is ML such a big deal?

- Today's products are the result of many constraints: what machines can and can't do
- These new capabilities **significantly** change the constraints that technologists work within
- If the form and function of today's products are the result of working within constraints, and many of these constraints fall away, the form and function of these products may change radically
 - Every one of these changes is an opportunity for the people sitting in this room!



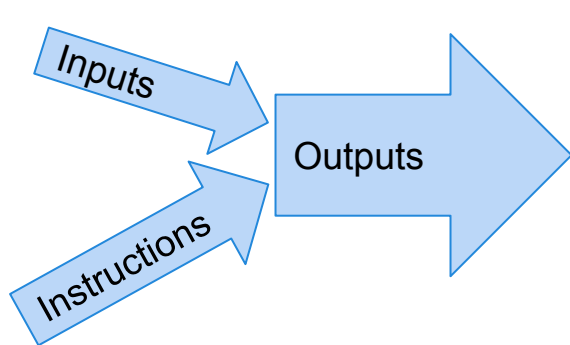
Why is ML such a big deal?

- Example: driving a car
 - Before AI: GPS gives you turn by turn directions
 - After AI: The car drives itself
- Why such a large change? What constraints changed?
 - Before this transition the machines could not see the road
 - Before this transition the machines could not learn to drive from playing a simulation
- Before AI: A human sees the road and instructs the machine how to turn. A human practices driving in order to learn how the controls affect the car
- After AI: Because these crucial steps can be automated by technology, the product radically improves

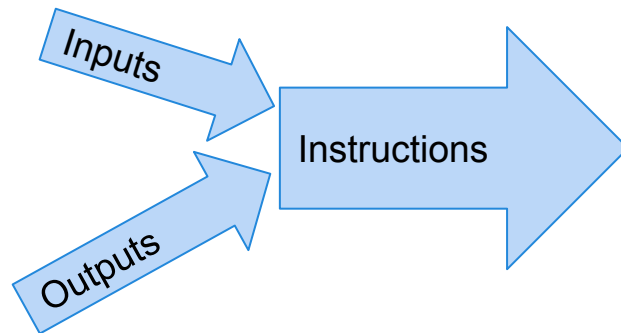


Terminology

- How does building software with machine learning compare to traditional software?
- Traditional software: Programmers told computers how to solve problems
- Machine Learning: Programmers give millions of examples of correct behavior, and the computer figures out the best rules



vs.



Terminology

- Training vs evaluation
 - Training is when you present a large number of examples to the machine, and it determines a set of weights that best map the inputs presented to the outputs (or "labels")
 - The output of training is referred to as a "model"
 - Evaluation is when you run a new set of inputs on a previously-generated set of weights in order to produce a new output.
 - The output of evaluation is referred to as a "prediction"

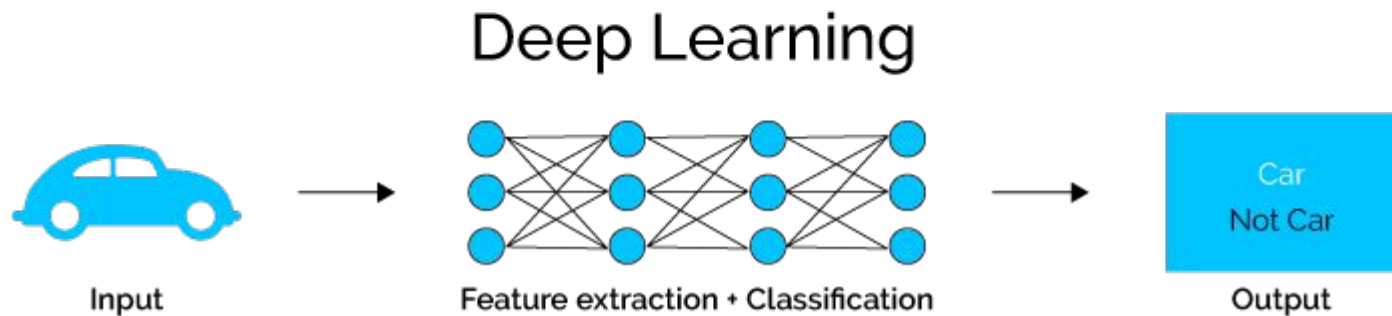
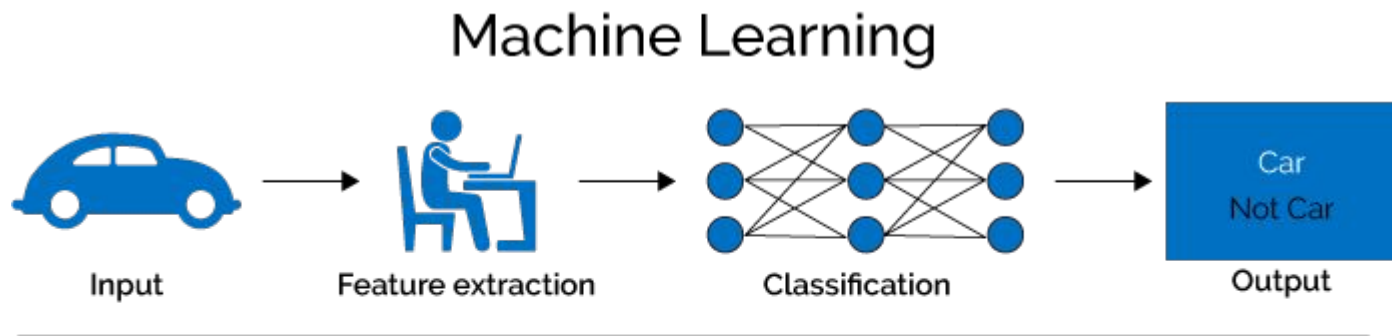


Terminology

- Deep Neural Networks
 - The source of the most impressive breakthroughs in recent years
 - Millions of connected artificial neurons, inspired by the human brain
 - Generally less emphasis on feature extraction
 - Accuracy scales well with size of the data set
 - Examples: CNNs, RNNs, LSTMs, GANs
- Traditional Machine Learning
 - Feature extraction and engineering done by the engineer
 - Generally easier to reason about and debug
 - Better performance on smaller data sets
 - Examples: Random Forest, Naive Bayes, Support Vector Machine



Terminology



Where does ML fit in?

- So we have...
 - Training data examples (MB or GB)
 - Python code that, when given training data, generates a model (MB or GB) over minutes to days
 - Python code that, when given an input and a model, generates a prediction (usually small: bytes or kb), over milliseconds to a few seconds
- How should we fit this in to our tech stack?



Where does ML fit in?

- In order to figure out where and how the ML fit in, we first need to answer some questions:
 - How is ML being used in the product?
 - Where does training data come from?
 - How hard is it to train the model?
 - How often do we retrain?
- Let's look at a made-up example...



Where does ML fit in?

- Cat v dog detector
 - User provides a photo
 - Web page says "This is a cat!" (or dog)
 - User interface allows the user to confirm or reject



How is catvdog using ML?

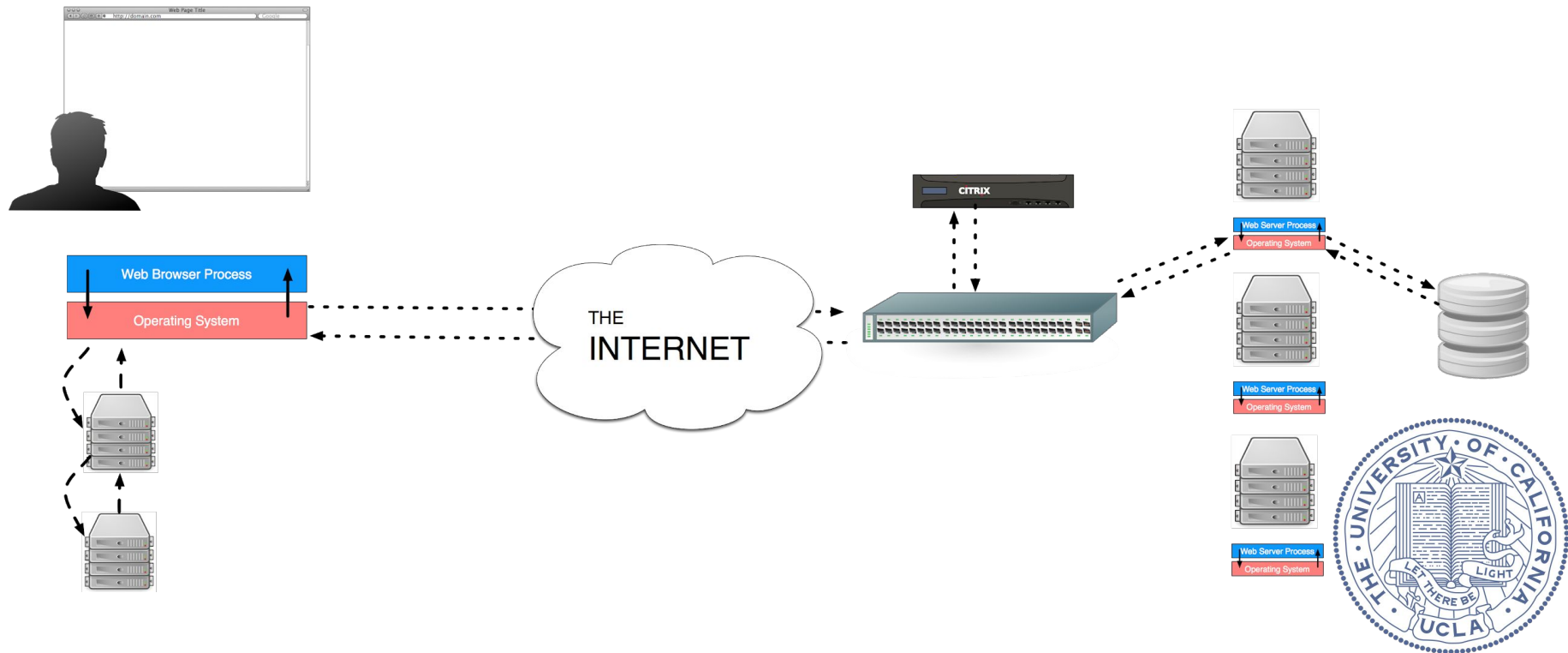
Let's talk about each of these questions as they relate to the catvdog app

How is ML being used?

- There are two convolutional neural networks that classify each photo:
 - Does it have a dog?
 - Does it have a cat?
- Where should we run the model?
 - What are our options?



Where does ML fit in?



How is catvdog using ML?

Let's talk about each of these questions as they relate to the catvdog app

How is ML being used?

- There are two convolutional neural networks that classify each photo:
 - Does it have a dog?
 - Does it have a cat?
- Where should we run the model?
 - What are our options?
 - We can run on the client
 - We can run inside our app server process
 - We can run in a separate process (service oriented)
 - What are the advantages and disadvantages of each?



How is catvdog using ML?

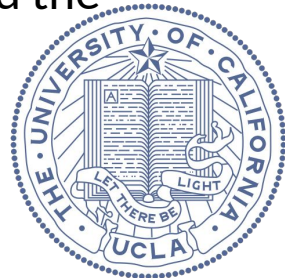
- We can run on the client
 - Pros: Some user interfaces require it
 - What type of user interface would require this for catvdog?
 - Cons: Restricts your technical choices (RAM, CPU, etc)
- We can run inside our app server process
 - Pros: We don't have to deal with services (monoliths are easy)
 - Cons: Your tech stack for the ML and the Internet Service need to match
- We can run in a separate process (service oriented)
 - More on the next slide...



How is catvdog using ML?

We can run in a separate process (service oriented)

- This is the most common approach
- Provide a RESTful interface to model evaluation (and any feature extraction)
- Cons:
 - Extra complexity for deploying additional service
- Pros:
 - It decouples your technology choices between the internet and the ML
 - It can run wherever you want it to
 - It can scale independently



How is catvdog using ML?

It decouples your technology choices between the internet and the ML

- Platform choices are hard, and their consequences last many years
- In 2019 what is the best internet platform to start your project on?
- In 2019 what is the best ML language and framework for your project?
- What if these aren't the same?



How is catvdog using ML?

It can run wherever you want it to

- The best hardware to run your web app may not be the best hardware to run your models
- Models vary and your model may...
 - Require a lot of memory
 - Require a lot of CPU
 - Require a GPU (or TPU)
- You may want to outsource your model evaluation completely
 - Google Cloud AI Platform, AWS Sagemaker are both great tools for evaluating and monitoring your models
- Models are semi-portable
 - Pickle
 - PMML (predictive model markup language)



How is catvdog using ML?

It can scale independently

- This is true of any service, but why is this particularly true of a service that evaluates ML models?



How is catvdog using ML?

It can scale independently

- This is true of any service, but why is this particularly true of a service that evaluates ML models?
- Model evaluation is stateless!
 - In this class we have demonstrated how easy it is to scale out stateless application servers and how hard it is to scale out stateful databases
 - Put a load balancer in front of as many application servers as you want, and scale easily



Where does training data come from?

You can't build a model without training data

- If you were going to build this, where would you get photos?
- And where would you get labels?



Where does training data come from?

You can't build a model without training data

- If you were going to build this, where would you get photos?
- And where would you get labels?

You can do a lot of things

- Take photos yourself
- Find images online by hand
- Label images yourself
- Buy data labeling services online

These are great, but some issues arise:

- None of these deliver massive data sets
- Are these photos and labels representative of what we will see in production?



Where does training data come from?

The ideal is when you can get your training data directly from your users

- There is no real limit to the amount of data you can gather
- The data is by definition representative of what your users will upload

But we have a chicken and the egg problem:

- We need a model to get users
- We need user data to get a model

How do we solve this?



Where does training data come from?

How do we solve this?

- Start with a poor model built with non-ideal data
- Get enough usage to improve the data set
- Improve your model with the new data
- A better model will mean a better product
- A better product will mean a more usage

The Virtuous Cycle of AI

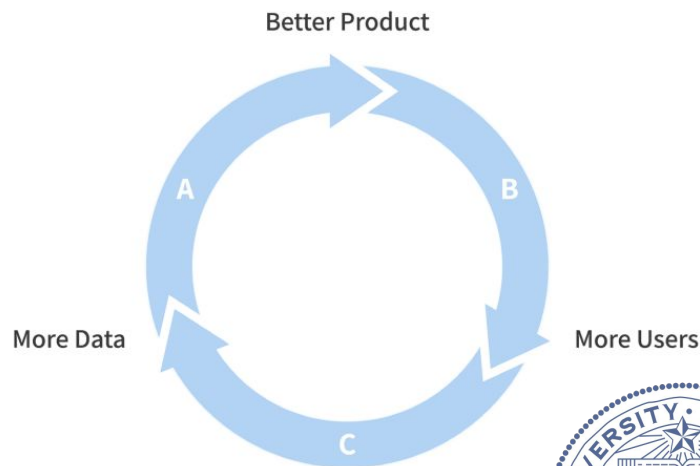


Image source: Landing AI's transformation playbook

Where does training data come from?

So how would we do this for catvdog?

- Take lots of photos of cats and dogs
- Label them yourself
- Build a so-so model on this small data set
- Get some people to try using your product
- When you have enough user data, switch to building the model with user data



Where does training data come from?

- Technical specifics of getting training data
 - Developer does a SQL query on the live DB
 - ETL process to a data warehouse
 - Unstructured data lake



How hard is it to train the model?

Some questions to ask in order to understand the training requirements of your system

- Can your data preprocessing be completed by a single machine within a reasonable amount of time?
- Can your training be completed by a single machine within a reasonable amount of time?
- Can your training fit within memory on a single machine at all?



How hard is it to train the model?

- Can your data preprocessing be completed by a single machine within a reasonable amount of time?
 - If not, you need to use a cluster computing framework
 - Hadoop: Open source map reduce implementation
 - Spark: JVM based cluster computing based around distributed data
 - Apache Beam (dataflow): Data processing pipelines
 - Unified batch and streaming architecture
- For catvdog, we likely could preprocess data until we had significant scale



How hard is it to train the model?

Can your training be completed by a single machine within a reasonable amount of time?

- If so, you can use your workstation
- NNs can be significantly sped up by using GPUs
 - Anywhere from 5 to 100x speedup
- For catvdog, we likely be using CNNs over images and as a result a GPU would be effective



How hard is it to train the model?

- Can your training fit within memory on a single machine at all?
 - SparkML
 - tensorflow distributed
- For catvdog, we likely would eventually need a scalable way to train our neural network, so would probably use TF distributed



How often do we need to retrain?

How often do we retrain?

- The model is learning a function, and that function can degrade over time. Would you expect these to degrade?
 - Learning to recognize a face in a photo?
 - Learning to identify a face in a photo?
- Do you think catvdog would need to be retrained often? What might change that answer?



How often do we need to retrain?

Why is frequency of retraining important?



How often do we need to retrain?

Why is frequency of retraining important?

- It affects how developers interact with training
- If you train a model infrequently, a model can be thought of as the output of development work
 - The developer is creating a model
- If you are frequently training models, they start to resemble databases
 - The developer is creating code to create and manage models



Where does ML fit in?

- Lets see some examples of how ML fits in to real products:
 - Smart Bill Entry
 - Lisa AI



- User uploads invoices via PDF
- Machine interprets the invoice
- User is presented with the invoice, and the (hopefully) correct accounting data to be entered

Where does ML fit in to SBE?

How is ML being used in the product?

- Uploaded invoice fed to multiple models in order to extract information

Where does training data come from?

- Users reviewing the invoices

How hard is it to train the model?

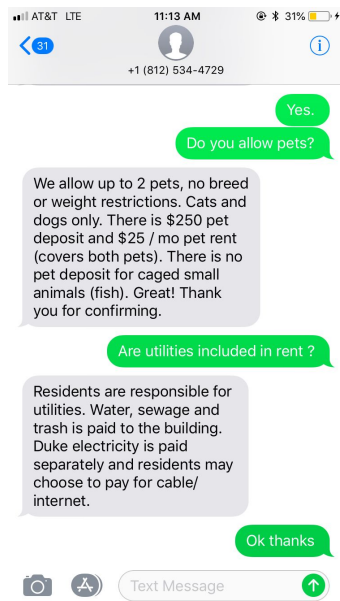
- Any particular model is small, but we tune models per-customer and have >10K customers

How often do we retrain?

- Models degrade quickly, so we retrain weekly



Where does ML fit in?



- Prospect has a text-based conversation with an AI system
- Prospects ask questions about apartments
- Conversationally schedule a time to come see the apartment in person
- Light pre-qualification
- Follow up after the showing



Where does ML fit in to Lisa?

How is ML being used in the product?

- All human communication runs through multiple models to understand it

Where does training data come from?

- Behind the scenes operators

How hard is it to train the model?

- Few models, take 2-4 hours each to train

How often do we retrain?

- Models degrade slowly, so every few months



For Next Time...

Lab tomorrow and no lab next week!

Stay active on Piazza. Don't get stuck. Don't wait until the last minute.

