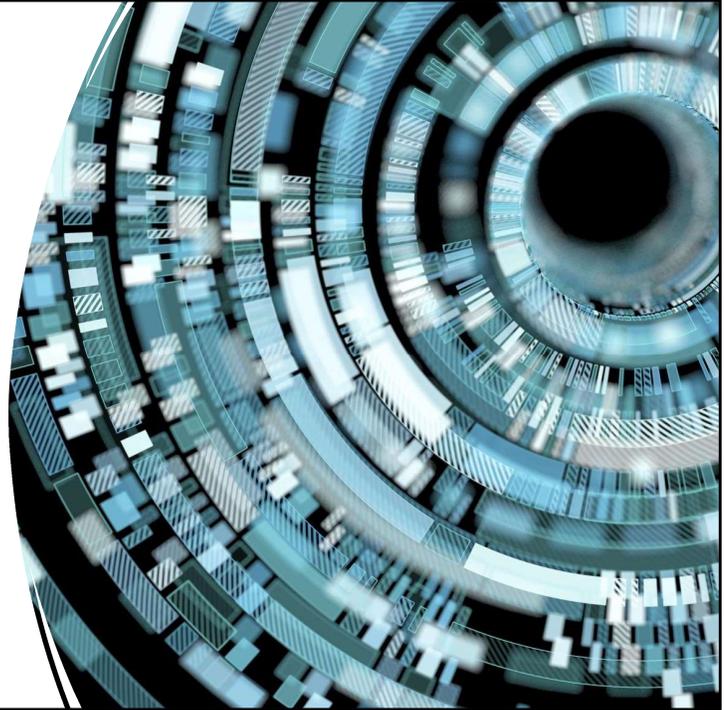


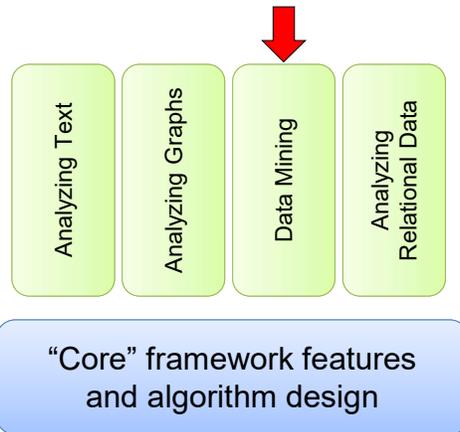
Data-Intensive
Distributed
Computing
CS431/451/631/651

Module 6 – Data Mining /
Machine Learning



Woah, trippy

Structure of the Course



What the, we skipped relational data?

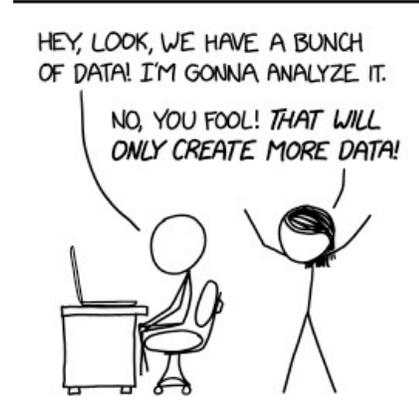
(Yes, the assignments flow more easily this way...maybe I should change the graphic...)

Data Mining

- [v] - Analyzing a large dataset in order to gain "actionable insights"
 - Whatever THAT means

Or "We know what we're looking for but it's more complicated than writing a regex!"

Or "We don't know what we're looking for, find me something good!"



The first "or" – we use supervised machine learning to find what we're looking for by example

The second "or" – We use unsupervised machine learning to hopefully extract interesting things.

How to Train Your Function

(with a ~~small~~ set of example input
and expected output)

AKA

Supervised Machine Learning



Classification



Classification Functions

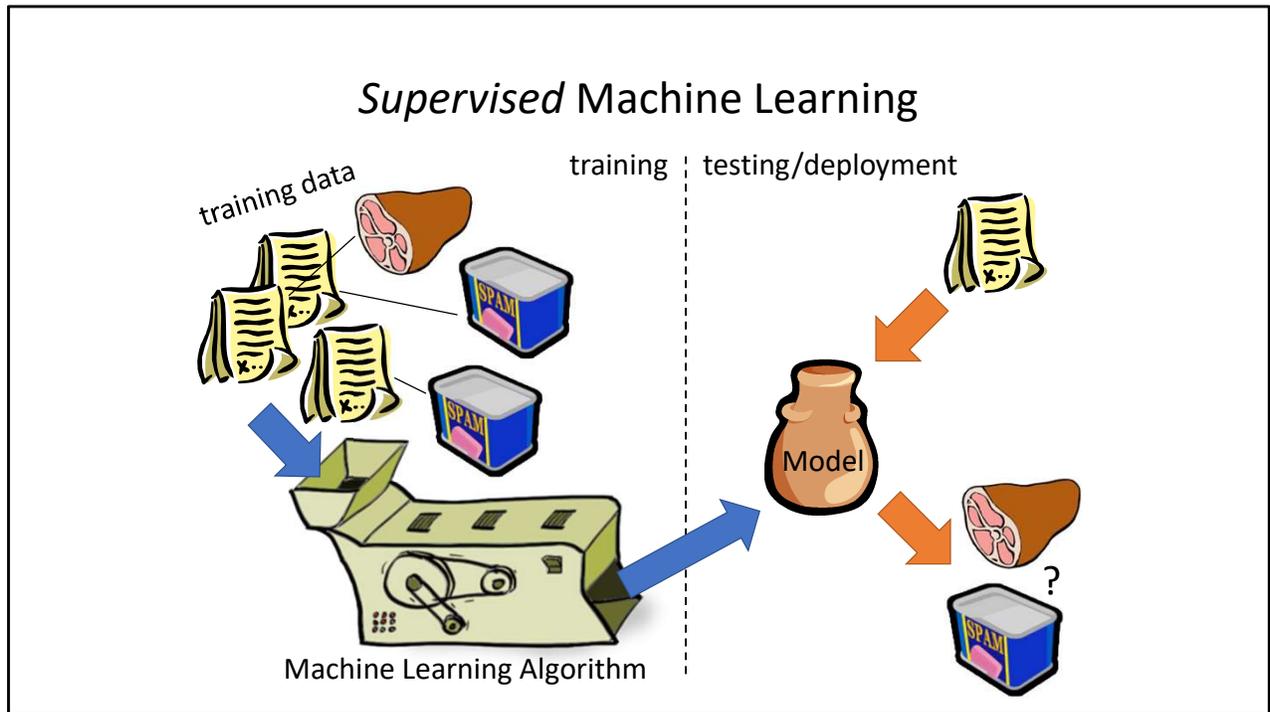
Input: Object

Output: Label(s)

Applications

- Spam Detection
- Sentiment Analysis
- Topic classification
- Link Prediction
- Document Ranking
- Object Recognition
- Fraud Detection
- NSFW Filter

There's so many more, but I'm out of space before it shrinks the font on me



Supervised: You give is a training set, where emails have been classified as “SPAM” or “HAM” (meaning not spam...it’s a bit of a joke you see)

Then, some sort of MATH happens, and out comes a model that can classify things.

Object Representation

Q: When is an email not an email?

A: When it's a "feature vector"

Objects are represented as a vector of features:

- Dense Features: sender IP, timestamp, # of recipients, etc.
- Sparse Features: message contains "Viagra", subject contains "URGENT", etc.

In terms of numbers (and everything's a number, really)

Dense Feature: zeros are rare or non-existent.

Sparse Feature: non-zeros are rare

There tend to be a LOT of potential sparse features. How can the computer know about all of them??? There are a lot of words...

What Features are Important For:

- Spam Detection
- Sentiment Analysis
- Content Classification
- IR Rankings
- Object Recognition
- Fraud Detection

**For the same
object, each
application will
have its own set
of Features!**

Spam: certain phrases, certain words, are very likely to be spam associated.

Sentiment Analysis: words that aren't important to spam/not spam are very important to inferring the writer's tone

Content classification: (Might overlap a bit with spam)

On Feature Selection

Horse or Crocodile?
How to tell the difference



	Horse	Crocodile
Eyes:	2	2
Ears:	Quite pointy	Not particularly pointy
Teeth:	Yes	Yes
Weight:	<250,000kg	<250,000kg
Location:	Earth	Earth
Attire:	None	None
Likelihood of eating a sugar cube if offered:	High	High
Culpable for the death of Princess Diana?	No	No
Any involvement in the overthrowing of the Russian government in 1917?	No	No

Conclusion

As we can see, there are very few differences between horses and crocodiles so the key thing to look for is the pointiness of the ears.

It's important that you select good features. This is a whole topic all on its own

Generally you want features that are independent. If two features strongly correlate you don't get much out of having both.

Embeddings – ML Feature Extraction

- An embedding is obtained by machine learning
- Example: word2vec – convert a word to a 384 dimension vector
- Goal: King – Man + Woman \approx Queen
- LLMs (e.g. ChatGPT) train the embedding (encoder/decoder) as well as the language model concurrently or alternating
 - Diffusers like StableDiffusion do as well

An embedding is a specific kind of dense feature vector – one where each dimension is NOT a distinct feature (per se) but one where the vector was learned by some kind of Deep Learning model.

Of course, as mentioned, different applications require different embeddings. A Queen is not drawn as a gender swapped King, but with different attire entirely (usually) – the Stable Diffusion embeddings don't follow this approximate equality

ML Solution Looks Like:

1. Acquire Data
 - Clean Data
2. Label data (by hand!)
3. Determine features
 - Maybe by hand
4. Pick an Architecture
5. Train the Model
6. Repeat as needed



So Many Architectures

- Naïve Bayes
- Logistic regression
- SVM
- Random Forest
- Neural Networks
 - Perceptrons / FFNs
 - RNN
 - CNN
 - Transformer
 - Diffuser



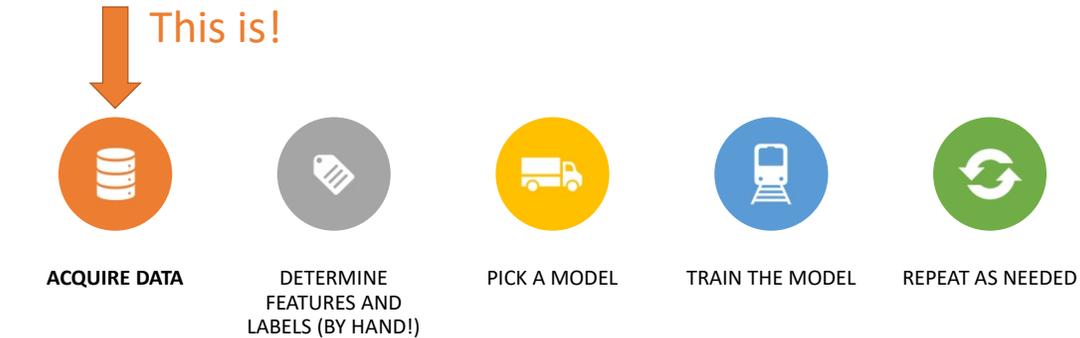
Pictured: A Random Forest (haha, get it, it's a joke! PowerPoint suggested the image. It's funnier than I am and it's not even trying)

Simple Spam Filter (aka Idiot's Spam Filter)

1. Put emails 1 per line in a text file, with SPAM/HAM label
2. Do a modified Word Count from A1
 - Keep two counts – across all emails and across spam
3. Do some math
 - $\Pr[\text{spam}] = N_s / N$
 - $\Pr[w] = C(w) / N$
 - $\Pr[w | \text{spam}] = C_s(w) / N_s$
4. Bayes' Theorem : $\Pr[\text{spam} | w] = \frac{\Pr[\text{spam}] \Pr[w | \text{spam}]}{\Pr[w]}$

Back in the 90s ISPs didn't have spam filters so you had to do this yourself. I downloaded a Bayes Filter plugin for Outlook Express. You had to train it yourself on your own emails. It's not really "training" though, in the same way that our N-Gram language model wasn't "training" - just counting.

What's Important?

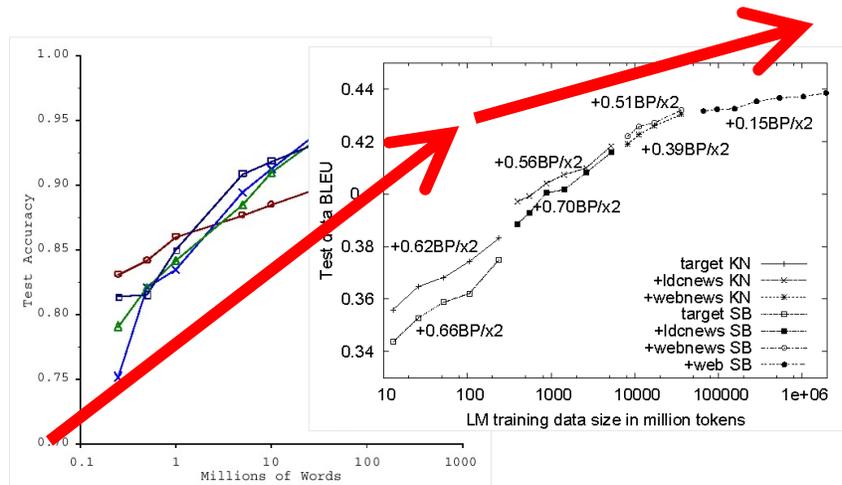


In a Big Data course, did you expect anything to be more important than Big Data?
With enough data, even a bad model will perform fairly well. That's the theory anyway.

It's true though. Actually, probably the hand labelling is more important...BUT, how can you label data you haven't acquired? Exactly.

(In fact "acquire high quality, clean, and labelled data" is the most important part, grouping #1 and #2 together with a secret third thing – cleaning up all the garbage).

No data like more data!



(Banko and Brill, ACL 2001)
(Brants et al., EMNLP 2007)

Left grammar checker.
Right machine translation

Take away – more data, more good. [that's another joke, because the diagram is about a grammar checker – though tragically this phrasing is common in tech circles somehow...]



AI Humor



A: Have you heard of this new deep learning model?

B: Yeah but I'm sticking to logistic regression

A: What? That's so old and simple!

B: If the performance is sufficient for your use case, simple is best.

- Mixtral AI, not understanding the concept of a "data science joke"

Wait, Hold Up, did
you say BY HAND?

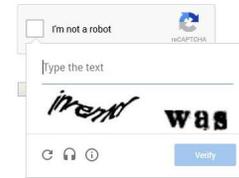
I did: "Determine ...
Labels **(by hand!)**"

Doesn't seem like it will
scale...



How to do things by hand at scale

Crowdsourcing



Bootstrapping (and other semi-supervised techniques)

Exploiting User Behaviour

(e.g. emojis are self-labelled sentiment analysis)

Crowdsourcing – Might be hiring people on Mechanical Turk, paying Google for surveys, etc.

Funny Captcha anecdote – BLIP2 (I think?) has a better accuracy than actual humans for those annoying “click the pictures with traffic lights” sort of captcha)

Bootstrapping – Label a reasonable amount by hand, and train a high-precision model (precision = $TP / (FP + TP)$)

Use model to generate pseudo-labels

Repeat on pseudo-labelled data, until pseudo-labels converge

OR – Use GPT4 to label your data [BUT, the ClosedAI TOS says you’re not allowed to use GPT4 for AI so you’ll be a wanted criminal maybe – still, it’s nice to be wanted]

New in Winter 2024: Yi 34B and Mixtral 7Bx8 are pretty close to GPT3.5 in terms of few-shot classification. Not SOTA but open weights you can run on a commodity GPU (or even on CPU if you don’t care about speed)

Binary Classification

Label is a single binary value. Yes or no. Spam or Not Spam.

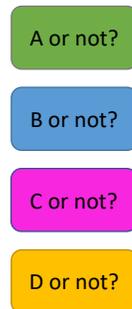
You can make more nuanced classifiers out of binary classifiers

Hotdog or not Hotdog

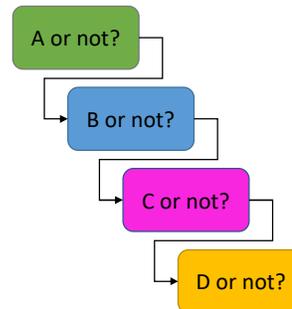
Binary Classifiers as Building Blocks

Example: four-way classification

One vs. rest classifiers



Classifier cascades



The above shows 4 way classification made from 4 “is or isn’t” classifiers.

“Spam or not spam”, “Notification or not notification”, “Personal or not personal”, etc.

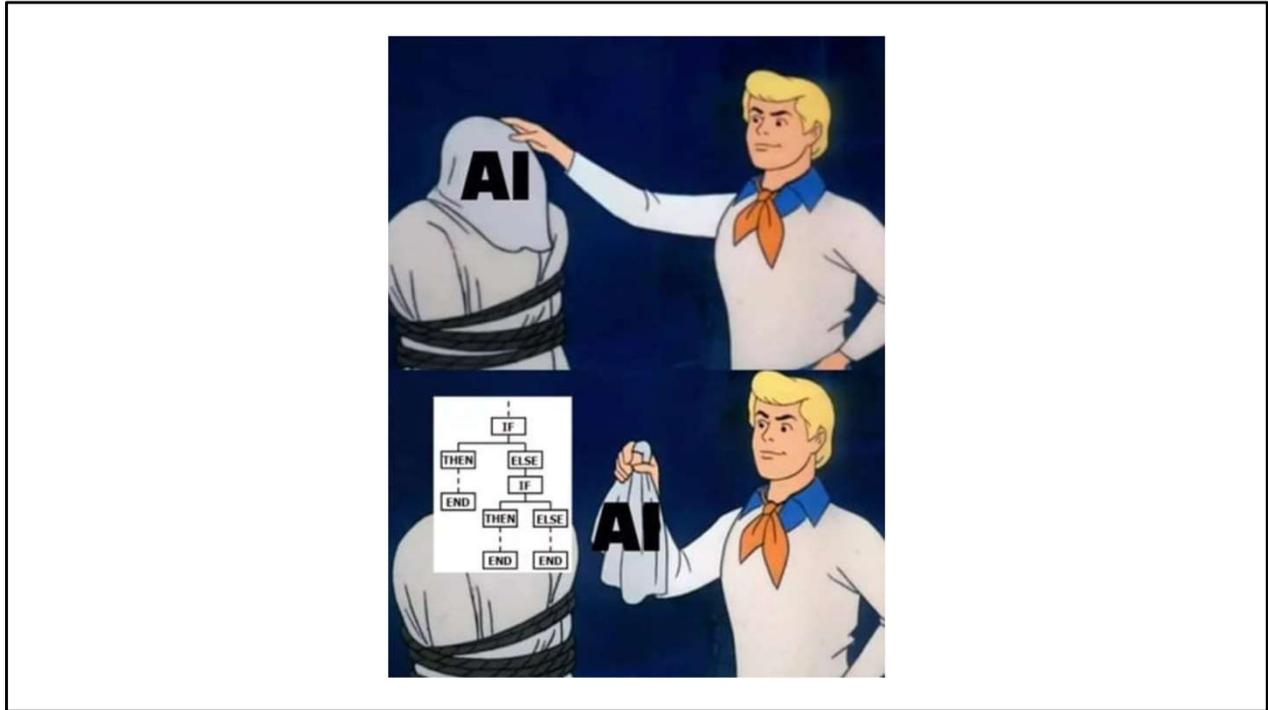
You can also use binary classifiers to create a classifier tree, if the labels are not mutually exclusive.

Botany books have decision trees like that:

“Leaf is serrated or not” – “flower is single or cluster” etc. Those might seem more like features than labels, and you’re right! But if you’re working from a PICTURE, you have to use ML to find the “features” of the image, and then again to classify based on those features! Neat.

- Fun word : Variegated. If a book asks you “is the foliage variegated?” it means “are the leaves stripy / splochy?”

Dan intends to draw something on the board. Let’s see if he does.



Binary Classifier Cascades aren't the only sort of AI.

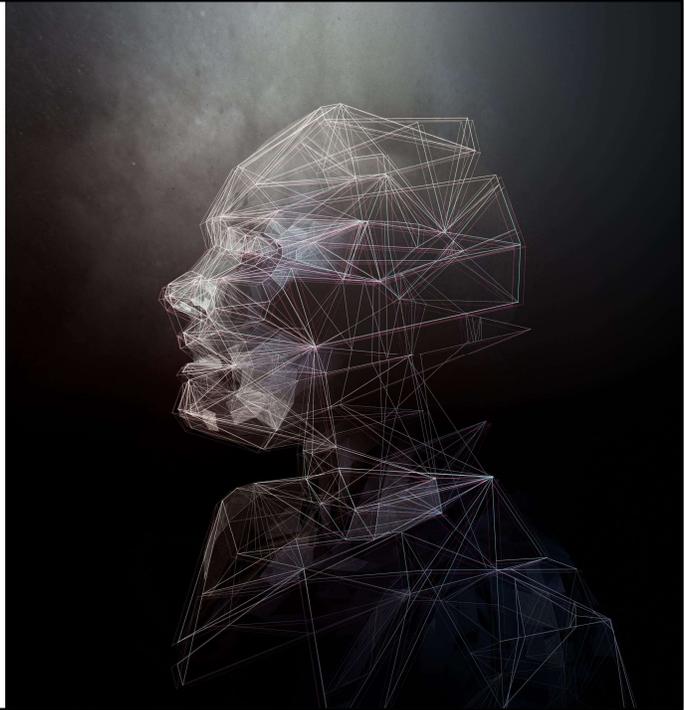
Generative Adversarial Networks, Deep Convolutional Networks, etc. don't really fit with the meme...so we'll ignore them.

(this meme is a joke on the LLaMA subreddits discussing 1-bit quants where maybe neural nets might be reduced to if/else trees???)

Also, how did you pick WHICH binary classifiers to train? You have to do that before you can train the decision tree.

—
Another AI Generated
Joke

Machine Learning:
The art of learning
without
understanding.



(it's not wrong, understanding requires consciousness and a bunch of tensors isn't conscious)

Training a Classifier

$D = \{(x_i, y_i)\}$

D: Training Data, x_i : feature vector, y_i : label

Want to find $f: X \rightarrow Y$

A function that maps feature vectors to labels

(But not any f , an f that minimizes “entropy” or “loss”)

$$\frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$$

Where ℓ is a “loss” function

A loss function for a Boolean classifier might be as simple as $\ell(a,b) = (a \text{ xor } b)$ i.e. 1 if they differ, 0 if they're equal. Of course you can apply different weights to false positives and false negatives if you want.



Computer Facts

@computerfact



concerned parent: if all your friends
jumped off a bridge would you
follow them?
machine learning algorithm: yes.

2:20 PM · Mar 15, 2018

Labels – what a person says

Machine Learning – copy what people are saying

Everyone: AI art will make designers obsolete

AI accepting the job:



Tried the meme myself...



Worked out a BIT better

Parameterization

There are LOTS of functions.

We'd rather have a parameterized function.

You sometimes see $f(x|\theta)$ when the function is returning probabilities

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=0}^n \ell(f(x_i; \theta), y_i)$$

CS folk often get confused. "Wait, is the vector of arguments, and theta is the set of parameters? But I thought argument was a value, and parameter was the name of that value?"

No, that's programming! This is math!

Less Writing

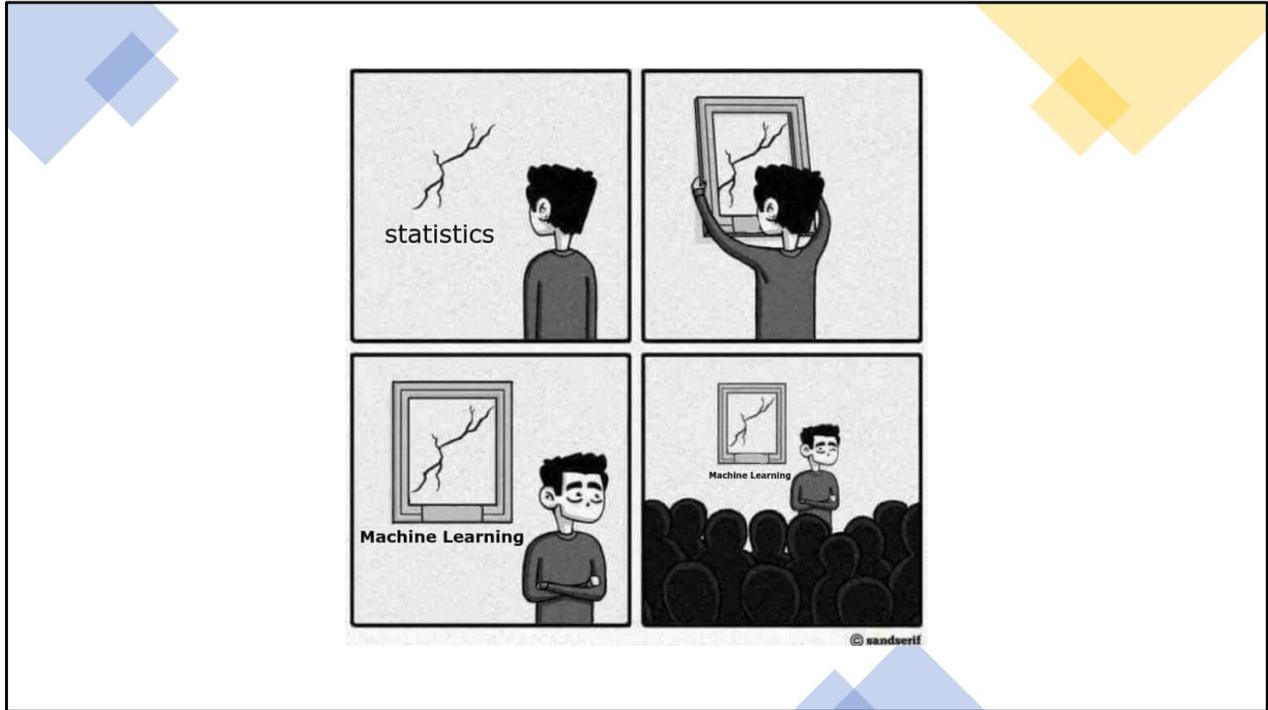
$$\arg \min_{\theta} \frac{1}{n} \sum_{i=0}^n \ell(f(x_i; \theta), y_i)$$

$$L(\theta) = \frac{1}{n} \sum_{i=0}^n \ell(f(x_i; \theta), y_i)$$

Trying to find argmin $L(\theta)$

Note from Past Dan: Dan has a bad habit of double clicking his slide remote and things get confusing if this slide is missed.

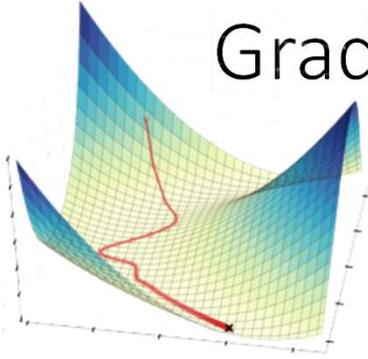
Note to self: Try REALLY hard to not accidentally skip this slide, otherwise nobody knows what capital L is



There are lots of jokes about “AI is just...” . Everything is “just” something else with a lot of messy details. Deal with it.
Besides, it’s mostly linear algebra, not statistics!

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=0}^n \ell(f(x_i; \theta), y_i)$$

There probably isn't a closed form solution to this so...

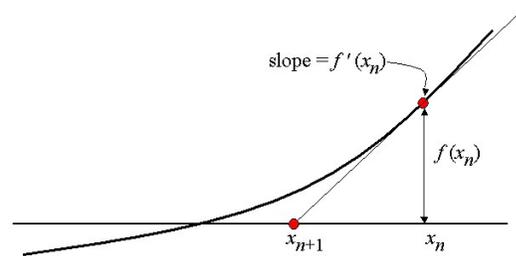


Gradient Descent

Intuition

Remember Newton's
Iteration? Sure ya do!

1. Get the slope at point x_n
2. Compute the intercept for the tangent line at point x_n , call it x_{n+1}
3. Repeat

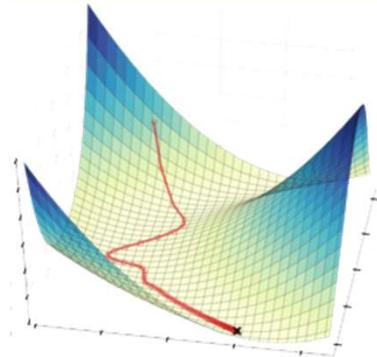


Intuition

Gradient Descent is the same thing

It just has many dimensions, not just x !

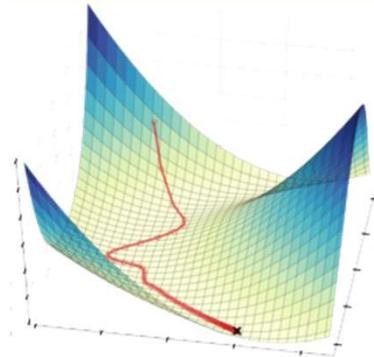
There's that "just" word again



Intuition

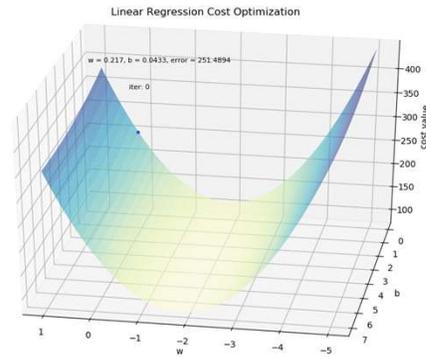
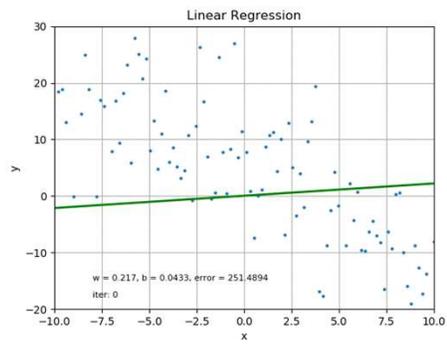
For each component x of θ :

1. Compute $\frac{\partial}{\partial x} L$
2. Compute “downhill” (like Newton’s Method with many dimensions)
3. Compute new θ



Ugh, differential equations. This is the one course I almost failed in undergrad...

Linear Regression + Gradient Descent



Gradient Descent

To find: $\operatorname{argmin}_{\theta} L(\theta)$

Compute Gradient: $\nabla L(\theta) = \left[\frac{\partial L}{\partial w_0}, \frac{\partial L}{\partial w_1}, \dots, \frac{\partial L}{\partial w_d} \right]$

So at any point θ ...

$$\theta' = \theta - \gamma \nabla L(\theta)$$

$$L(\theta) \geq L(\theta')$$

For “small” gamma.
“Assumption of local
linearity”

Wait, gamma?

Yes, it’s “step”, or how far we move along
the gradient

Also known as “Learning
Rate”

Replaced the complicated function with $L(x)$ – since nothing else changes, we only really need one parameter.

The inequality is “for sufficiently small gamma”.

Besides “gamma” and “step”, this is often called “LR” or “Learning Rate”

Missing Details



What step should we use?



What about local minima?



How fast does this converge?

A1: How small is 'sufficiently small'? – There are approaches to picking gamma. It'll depend on if the function is convex or concave too.

A2: it will get trapped in them ☹️

A3: not very quickly

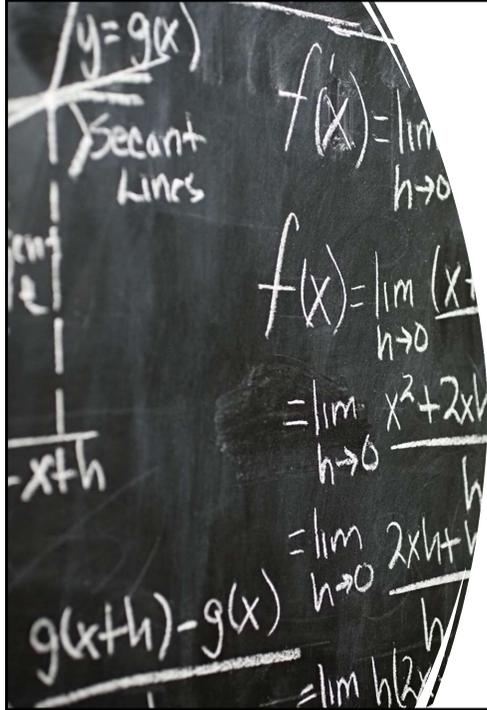
More Details

Gradient Descent is “first order”

- Local linear approximations
- Slow convergence

Fixes to avoid local minima

- Momentum
- Stochastic Gradient Descent
 - Randomly select a subset of training data
 - Gradient depends only on selected subset
- Can still get stuck



Other Approaches

Second order methods (Quasi-Newton)

- Quadratic approximation instead of linear
- Requires the Hessian (matrix of second order partial derivatives)
- Usually Impractical

I tried this for fitting database snippets into a protein with missing parts. It was tedious and then didn't work very well. In a bit we'll get to something somewhat similar to what I did! (Not really...but close enough for a detour)

Weren't we talking Boolean Classifiers?

Right, sorry.

$$f(X | W) : \mathbb{R}^d \rightarrow \{0, 1\}$$

Parameters to optimize:

- Weight vector W
 - θ is used as the abstract view of "The Model", vector W is the concrete realization

Loss?

$$f(X; W) : \mathbb{R}^d \rightarrow \{0, 1\}$$
$$\ell(y, t) = \frac{1}{2} (y - t)^2$$

y – predicted value
y = f(x), ya know?

t – actual value
t for “true”

Loss is 0 if prediction is right, 1 if prediction is wrong.

Problem: Not continuous, not differentiable.
Gradient Descent won't work.

Solution: **surrogate loss function**
Something that's not as accurate, but easier to optimize

My PhD work involved “squishing” protein snippets to make them fit into gaps, and the problem was something called “Multi-dimensional scaling”. Loss was called “stress” in the papers I read. The algorithm I used was called “SMACOF – Scaling by Majorizing a COMplicated Function”

“Majorizing” is a sort of surrogate function...where you find a way to create a simpler function that:

1. Is equal to the complicated one at a given point
2. Must be \leq the complicated function at all other points.

With SMACOF the majorizing function has a closed form minimum! No gradient descent needed, it goes straight to the minimum. Neat.

However, finding the minimum requires computing the pseudo-inverse of an $n \times n$ matrix, which is $O(n^3)$

Linear Function

$$f(X; W, b) : \mathbb{R}^d \rightarrow \mathbb{R} = W \cdot x + b$$
$$\ell(y, t) = \frac{1}{2}(y - t)^2$$

y: predicted value (any real number)

t: true value (0 or 1)

Problem: Cost is unbounded!

This seems like we've broken the classification function...but we haven't. Just set some cutoff point to map the values into "true" and "false"

Unbounded? If our classifier is "really confident" then it might return 100. If that's wrong, the loss will be immense!

We'll train the algorithm to never be confident. Might not be great.

Sigmoidal (Logistic) Function

$$\sigma(z) = \frac{e^z}{e^z + 1}$$

- Sigmoidal functions are nice

$$f(X; W) : \mathbb{R}^d \rightarrow [0, 1] = \sigma(w \cdot x)$$

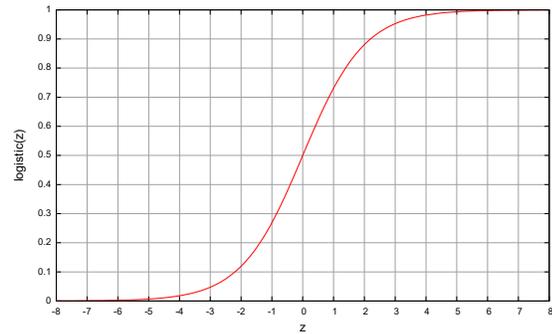
$$\ell(y, t) = \frac{1}{2}(y - t)^2$$

Prediction is between 0 and 1

Loss is between 0 and 1

Continuous and differentiable:

$$\frac{d}{dz} \sigma(z) = \sigma(z)(1 - \sigma(z))$$



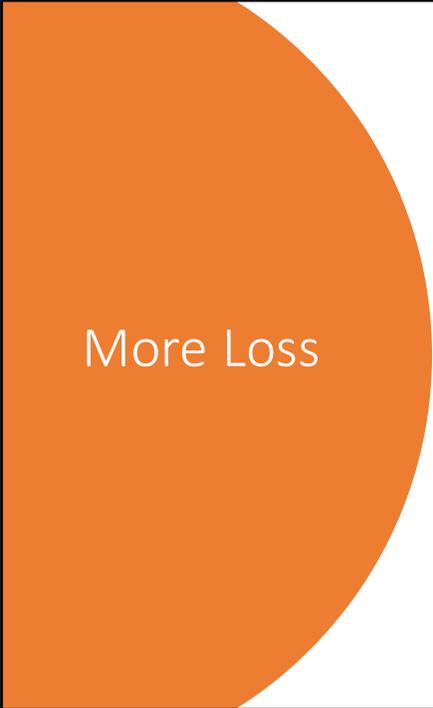
Not just differentiable, but easy! Gotta love exponentials.

More Maths

$$\frac{\partial L}{\partial z} = (y - t)y(1 - y)$$

$$\frac{\partial L}{\partial w_i} = (y - t)y(1 - y)x_i$$

Now we can do gradient descent



More Loss

A better loss function to use here is Binary Cross Entropy (BCE)

$$\ell(y, t) = -[t \log y + (1 - t) \log(1 - y)]$$

Giving the gradient

$$\nabla L(W) = (t - y)X$$

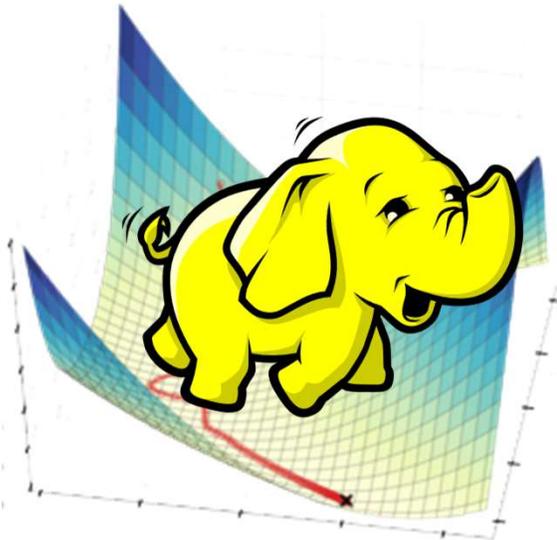
The logs and exponents cancel 😊



To Learn More...

Take an actual ML course. This is just a tribute

- Lots of other optimization techniques
- Lots of different loss functions
 - Logistic function has its own problems
 - 0.00001 vs 0.0000001, almost same loss
- Fixed step size is bad
 - Keyword: step scheduler (cosine, adaptive, etc.)



Gradient Descent on Hadoop

MapReduce is bad at iteration

- Startup Costs
- Retain State between Iterations
- Straggler Line – Iteration is bounded by the slowest worker

Gradient Descent: Must have only one reducer (bottleneck)

Spark to the Rescue? Again?

```
val points =  
spark.textFile(...).map(...).cache()  
var w = ... // initial vector  
for (i <- 1 to ITERATIONS) {  
  val gradient = points.map{ p => {  
    val score = p.x * w  
    val y = 1 / (1 + exp(-score))  
    p.x * (p.y - y)}}.reduce((a,b) => a+b)  
  w -= gradient * LR  
}
```

- Is that really better?

Yes

- But why?

1. Cache data
2. No Shuffles
3. Driver only does final reduce: 1 value per partition
 1. But it's a BIG VALUE.

Here the data point is split into x: feature vector, y: label

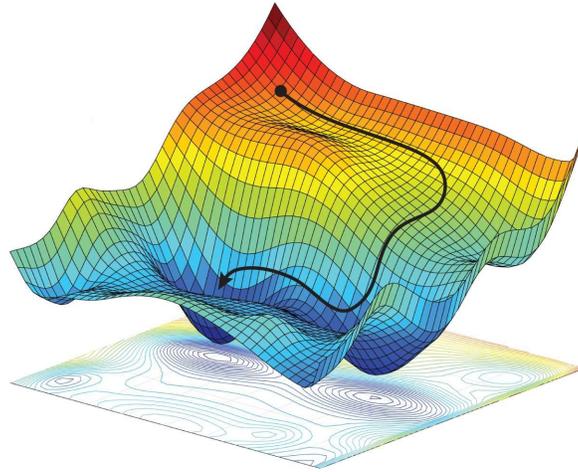
I'm assuming w and x are vector types that support dot products with each other, and scalar products with doubles.

Note that this is still unusably slow so it's never done.

Why: You need thousands of iterations and Spark iteration might be better than Hadoop MapReduce, but it's still not good.

Also why: If your model is 1 million parameters then shipping the gradient to the driver is ~4MB (or 2MB if using fp16 or bf16 for parameters, possibly 1MB if using Int8). That's a pretty small model though.

Stochastic Gradient Descent



Source: Wikipedia (Water Slide)

Stochastic Gradient Descent

- Pick a random subset, only compute gradient for that subset of values
- If the subset is small enough for one machine...why have the cluster?

If we have k mappers. Have each compute the gradient for its subset, then transform independently?

Ensemble Learning

(This was called “consensus modelling” when I was doing it for protein folding).



Ensemble Learning

1. Train multiple **INDEPENDENT** Models
2. Combine Predictions
 1. Voting
 2. Merging Models

E.g: Partition test set, train classifiers independently

Very parallel, perfect for a Hadoop cluster

You can't always merge models. For a simple linear / sigmoidal classifier function, averaging the weights is probably OK.

Other kinds of models are not so easily merged. Cannot merge models of different types. Surprisingly – You can merge neural networks sometimes – IF they're all finetunes of the same base model!

- This works very well with Stable Diffusion finetunes. With LLaMA it sometimes works but is more involved...

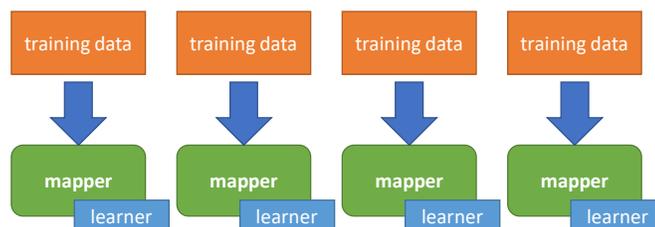


Ensemble Learning

Why does this work?

- Errors won't be correlated (hope)
 - Less likely that most models are all wrong (hope)
 - Reduced variance
-

Online Stochastic Gradient Descent as Ensemble



No iteration!

(Well, each mapper iterates over its partition, that's how MapReduce works!)

What's going on here: Technically, a single value is a set of 1!

We can select the elements in random order and compute the gradient ONE training value at a time. Combined with Inertia we might be on to something... (In fact you can compute the gradient for a constant number of elements at a time, it doesn't need to be only one.)

Why does this work:

If you have 10,000 data points, you can do 1 iteration through all 10,000 and sum the gradient up (10,000 gradient calcs, 1 multiplication (LR), 10,000 additions) or you can do one iteration per data point (10,000 gradient calcs, 10,000 multiplications (LR) 10,000 additions)

So for not a whole lot extra work you get 10,000x the iterations. They are very low quality iterations compared to the "all data" method, BUT, not 10,000x worse! In practice you should incorporate as many data points as will fit in memory at once.



In MapReduce

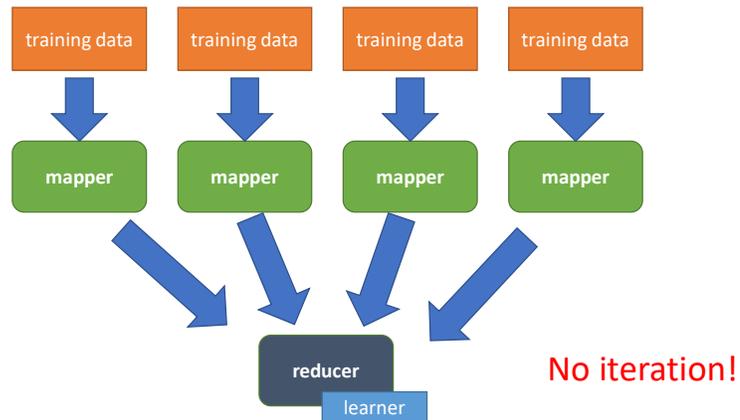
Each mapper holds the current parameter set in memory (create in setup)

Trains model by computing gradient for one element, updating parameters (in map)

Outputs Model (in cleanup)

What to do with all these models?
Ensemble learning!

Online Stochastic Gradient, One Model



If you only want one model, the mappers become “parsers” (parsing the strings and emitting feature – label pairs)
The learner is the run on the reducer.

“No iteration if each instance has its own key so there’s only one value” – On the assignment you do in fact need iteration to traverse the partition.



In MapReduce, One Model

Mappers are parsers only

Reducer becomes the learner:

- creates model (setup)
 - Trains model (loop in reduce)
 - Emits model
-

No more ensemble learning!

(Unless we set the number of reducers to something greater than 1)

In Spark

```
model = LogisticRegressionWithSGD(data,  
                                   iterations,  
                                   step)
```

That's right! The SparkML package has lots of models, estimators, loss functions. No need to write your own.

Except on the assignment, obviously

Sentiment Analysis Case Study

Binary polarity classification: {positive, negative} sentiment



Use the “emoticon trick” to gather data



Data

Test: 500k positive/500k negative tweets from 9/1/2011
Training: {1m, 10m, 100m} instances from before (50/50 split)

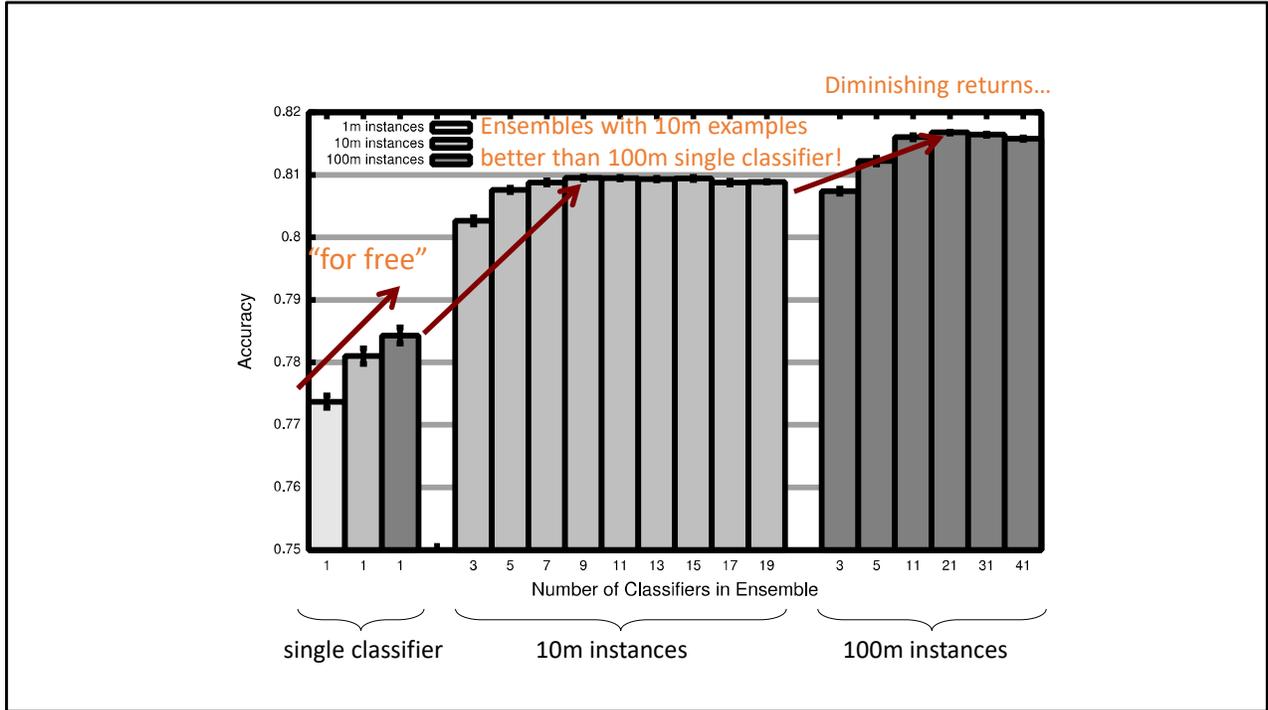
Features:

Sliding window byte-4grams

Models + Optimization:

Logistic regression with SGD (L2 regularization)
Ensembles of various sizes (simple weighted voting)

Source: Lin and Kolcz. (2012) Large-Scale Machine Learning at Twitter. SIGMOD.



Evaluation

How can we evaluate the model?

It's not enough to minimize loss.

Why?

Measure Accuracy?

Still not enough!

Why?

Loss is a measure of “how closely does it match the training data”. An overfitted model can have 0 loss, but still be useless.

Accuracy isn't as cut and dry. There's really FOUR different measures that are all interesting, and they're at cross purposes. (See 2 slides forward)

Overfitting

My model on training data



My model on test dataset



“Overfitting” is the key word. A model that’s PERFECT at predicting the training set might end up WORSE.

(Big Data helps this...the bigger the training set, the more “representative” it is.)



What to do about overfitting

Too much training = overfitted model = poor results

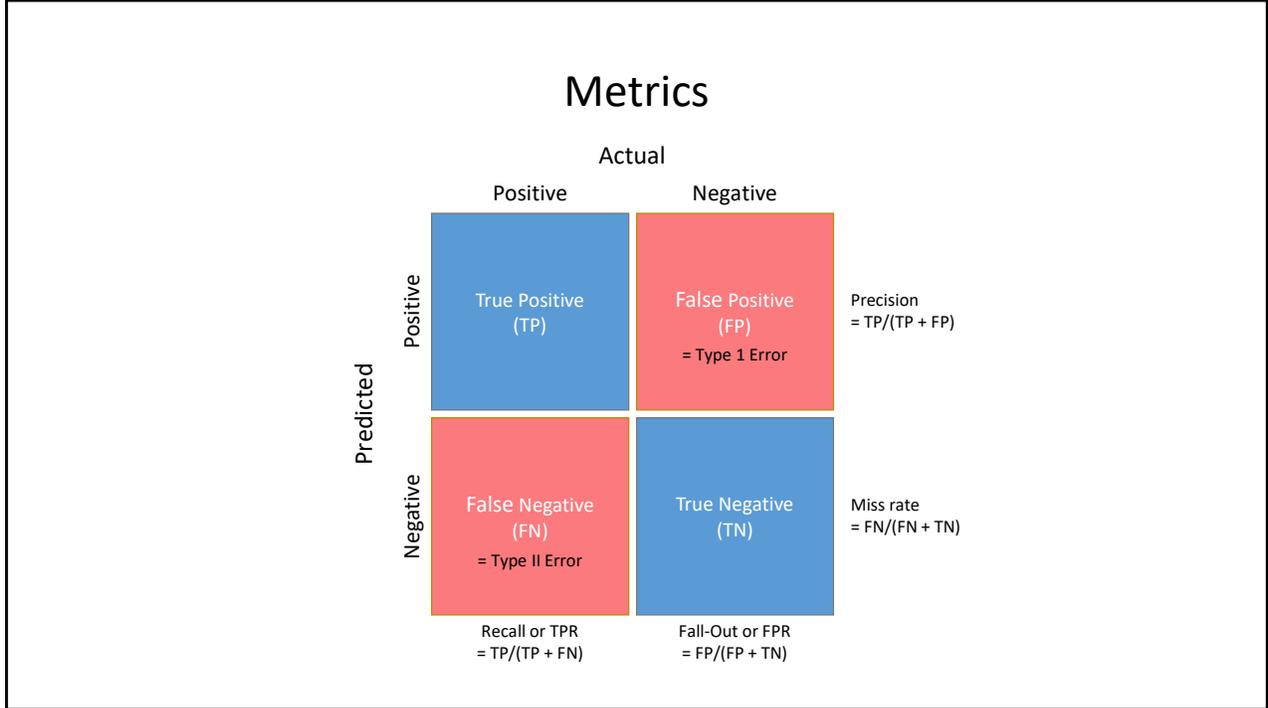
But: There's no way to know how much is too much except by comparing results!

Jargon:

- Epoch: 1 pass through the training data
- Checkpoint: Save the current model to disk

Usually you checkpoint every X epochs (the smaller X the better, but also the more HDD space you use)





Statistics again!

False Positive – Predictor says “yes”, but it’s a no

False Negative – Predictor says “no”, but it’s a yes

All measures shown (precision, recall, fall-out, miss rate) are interesting! But improving one can harm another



ROC and PR Curves

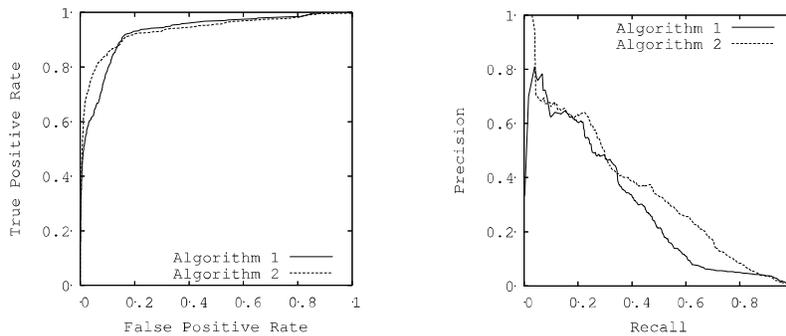
- ROC: Graph Recall (True Positive Rate) vs Fall-Out (False Positive Rate)
- PR: Precision vs Recall

In both cases, illustrates how the performance changes as you vary the threshold

A **receiver operating characteristic curve**, or **ROC curve**, is a graphical plot that illustrates the diagnostic ability of a binary classifier system as its discrimination threshold is varied.

Dan, draw them on the board.

ROC and PR Curves



Source: Davis and Goadrich, (2006) The Relationship Between Precision-Recall and ROC curves

ROC represents “as you vary the threshold, how many more false positives slip in, and how many more true positives are successfully found.”

What’s usually done is the area under the curve, or ROCA. A ROCA of 0.5 means it’s random chance. Lowering the threshold introduces as many true positives as it does false positives. ROCA of 1 means it’s a perfect predictor : regardless of threshold, false positive rate is 0%, true positive rate is 100%

So summing things up:

“Don’t use accuracy! There are many statistical metrics.”

Which is best?

“A ROC curve”

Can I put a number to it?

“Yes, Area-under-curve (AUC), or ROCA for short”

And that number represents...?

“OK so yes it represents accuracy. But like, threshold invariant accuracy. It’s different.”

Problem: Big Data Isn't Big Enough!

What's the testing set?

- Some data that wasn't in the training set

“Holdout Method”

The less data we train on, the worse our model!

The less data we test on, the less we trust our model!

K-Fold Cross Validation

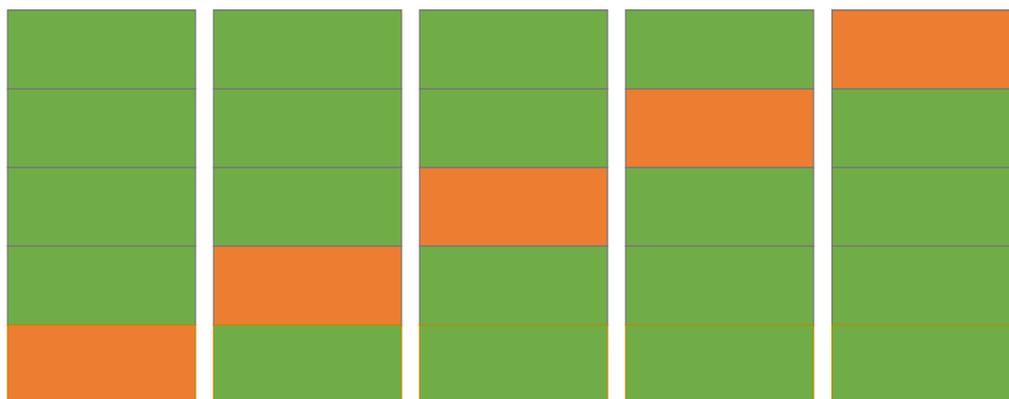
Divide Data into K subsets

Repeat holdout method k times. Each time, one subset is the test set, the rest are the training sets

Each set is used to validate ONCE

Each set is used to train K-1 times

5-Fold Cross Validation



Bonus? We have 5 models, can do 5 way voting? Or just merge the models. It depends on the kind of model.

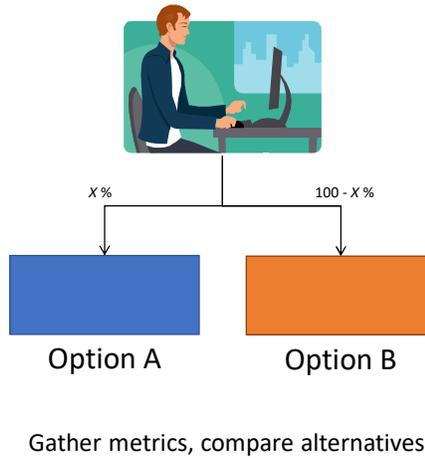


Why Cross-Validate?

This kind of ML is “offline” – When the technique is validated, then it can be deployed on live data in real-time.

Cross Validation gives you confidence that the technique will work on typical data sets.

A/B (Online) Testing



A/B testing is usually used for medical treatments. “Double Blind” – neither the patient nor the physician administering treatment knows which treatment it is. Usually one option (control) is a placebo or an already proven treatment.

Used in Marketing to check for a single variable.

A/B Testing for ML

Step 1: Offline training and evaluation (holdout, cross-validate, etc)

Step 2: A/B testing vs other methods

Return to Step 1 if needed



You can use this to compare different feature selection methods or compare to current best-practice models.

Applied ML in Academia

Download interesting dataset (comes with the problem)

Run baseline model

Train/Test

Build better model

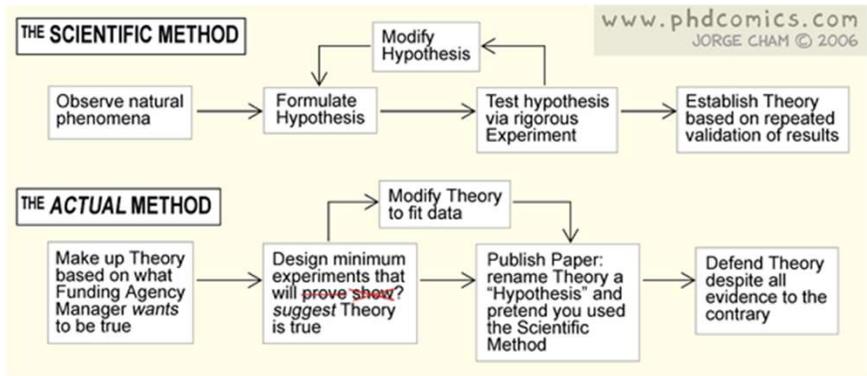
Train/Test

Does new model beat baseline?

Yes: publish a paper!

No: try again!

These few slides date back to Jimmy Lin. I'm not a ML person...but...ok, guilty...if you train a model and it's worse than the current standard, you can't publish that! So you keep trying until you can.



DATA

Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

FROM THE OCTOBER 2012 ISSUE

Fantasy

😎 Extract features

😎 Develop cool ML
technique

😎 #Profit

Reality

😬 What's the task?

😬 Where's the data?

😬 What's in this dataset?

😬 What's all the f#\$!* crap?

😬 Clean the data

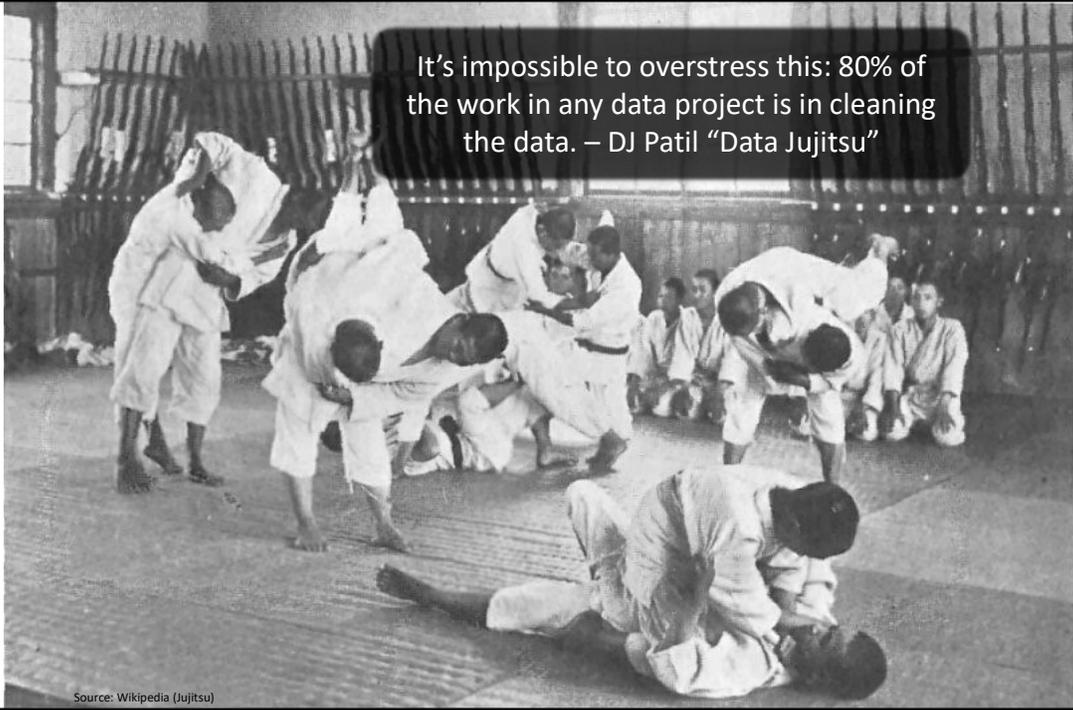
😬 Extract features

😬 "Do" machine learning

😬 Fail, iterate...

Dirty little secret. Data science is barely about ML at all!

It's impossible to overstate this: 80% of the work in any data project is in cleaning the data. – DJ Patil "Data Jujitsu"



Source: Wikipedia (Jujitsu)

On finding things...



On naming things...

- CamelCase
- smallCamelCase
- snake_case
- camel_Snake
- dunder__snake

uid UserId
userid userid
user_id user_id



Word embeddings actually help with this immensely!

OK, My
Model
(Finally)
Works on
the Training
Set

Good, you've made
a good first step!



I did A/B Testing,
and it has good
Precision & Recall!

OK, you're half-way there!





What's the other half?



Production



Source: Wikipedia (Oil refinery)

Production

What are your dependencies?

How / When are your job(s) scheduled?

Do you have enough resources?

How do you know if it's working?

What happens if it stops working?

All about infrastructure

TO BE CONTINUED...

(The file size is getting too big, going to make this two files)