

Data-Intensive Distributed Computing

CS 431/631 451/651 (Fall 2021)

Part 7: Data Mining (2/4)

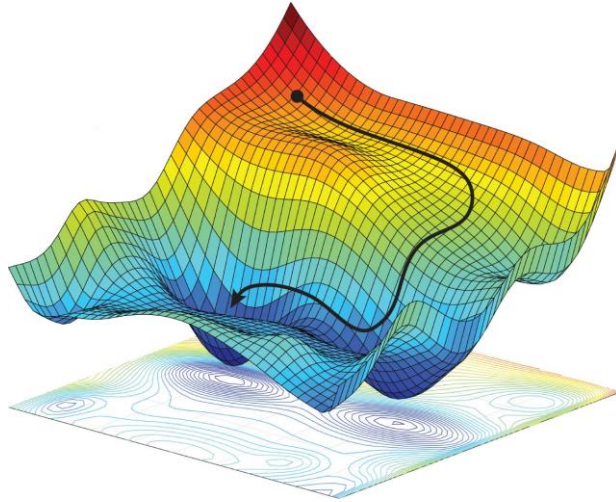
Ali Abedi

These slides are available at <https://www.student.cs.uwaterloo.ca/~cs451>



This work is licensed under a Creative Commons Attribution-NonCommercial-Share Alike 3.0 United States
See <http://creativecommons.org/licenses/by-nc-sa/3.0/us/> for details

Stochastic Gradient Descent



Stochastic Gradient Descent

Gradient Descent

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \frac{1}{n} \sum_{i=0}^n \nabla \ell(f(x_i; \theta^{(t)}), y_i)$$

Stochastic Gradient Descent (SGD)

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \nabla \ell(f(x; \theta^{(t)}), y)$$

Stochastic Gradient Descent

Gradient Descent

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \frac{1}{n} \sum_{i=0}^n \nabla \ell(f(x_i; \theta^{(t)}), y_i)$$

Considers all training instances in every iteration

Stochastic Gradient Descent (SGD)

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \nabla \ell(f(x; \theta^{(t)}), y)$$

Considers a random instance in every iteration

Stochastic Gradient Descent

Batch

Gradient Descent

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \frac{1}{n} \sum_{i=0}^n \nabla \ell(f(x_i; \theta^{(t)}), y_i)$$

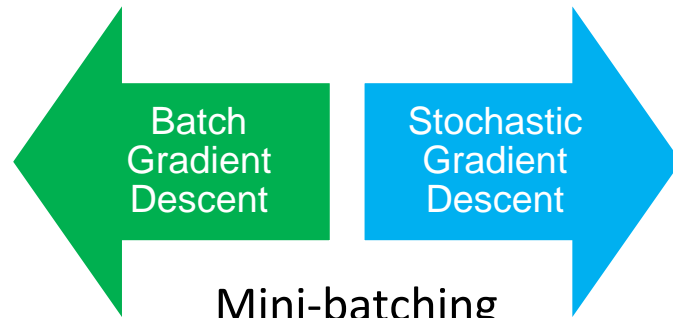
Considers all training instances in every iteration

Online

Stochastic Gradient Descent (SGD)

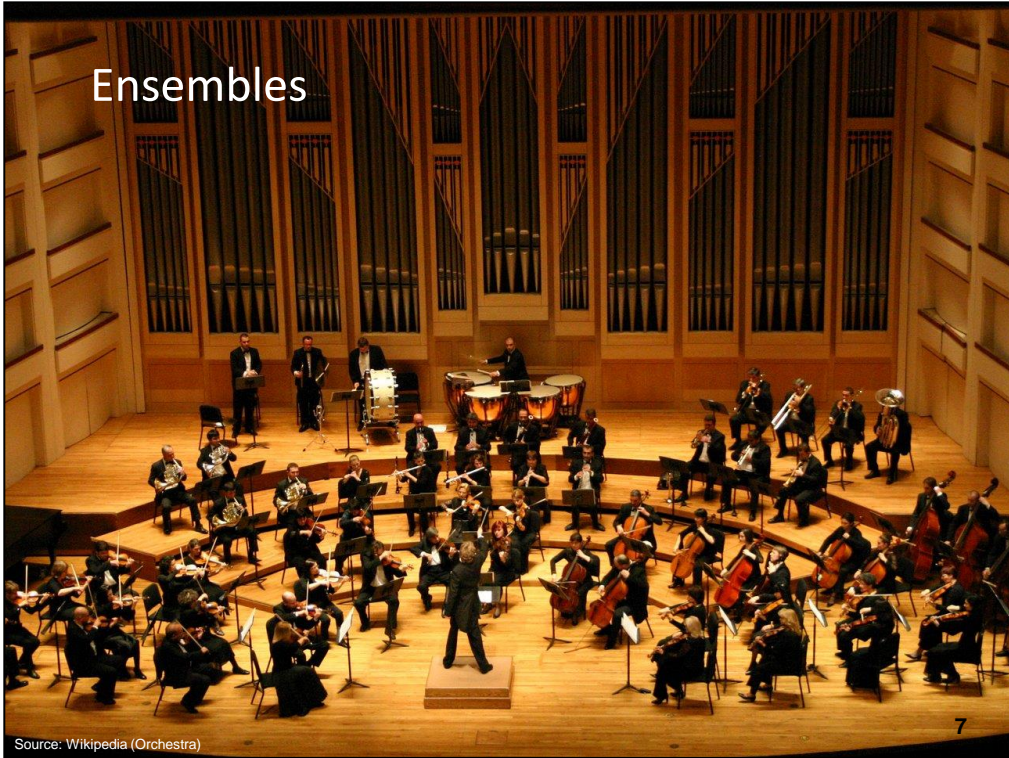
$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \nabla \ell(f(x; \theta^{(t)}), y)$$

Considers a random instance in every iteration



Considers a random subset of instances in every iteration

Ensembles



Source: Wikipedia (Orchestra)

Ensemble Learning

Learn *independent* multiple models, combine results from different models to make prediction

Common implementation:

Train classifiers on different input partitions of the data
Embarrassingly parallel!

Combining predictions:

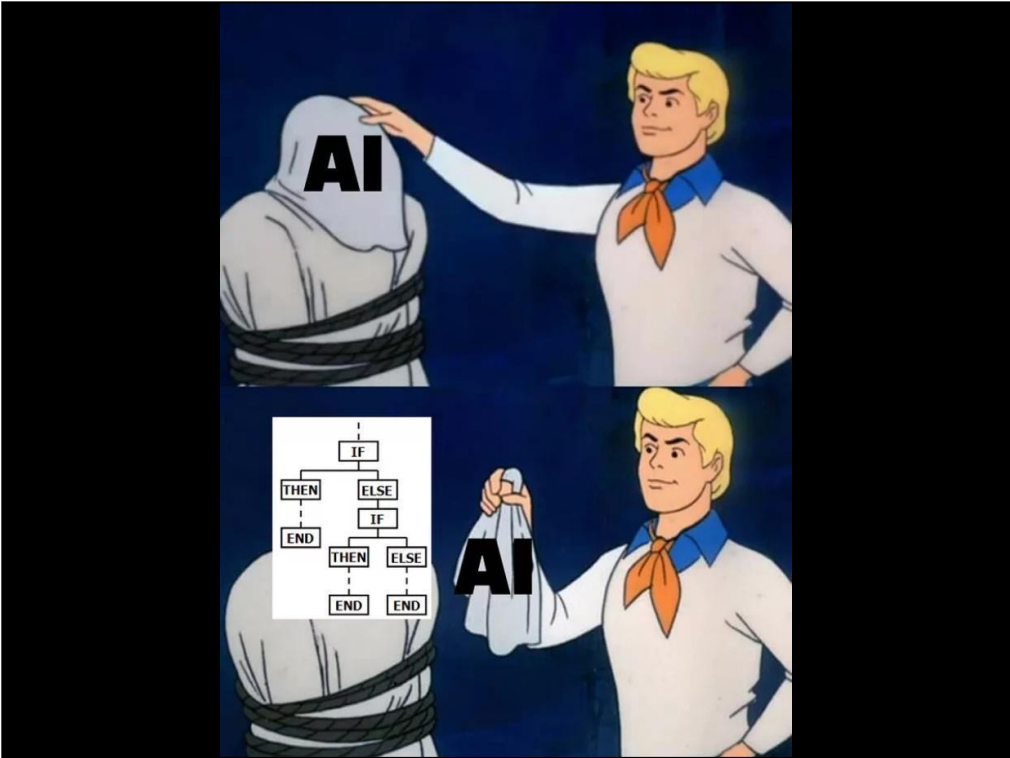
Majority voting
Model averaging

Ensemble Learning

Learn *independent* multiple models, combine results from different models to make prediction

Why does it work?

If errors uncorrelated, multiple classifiers being wrong is less likely
Reduces the variance component of error

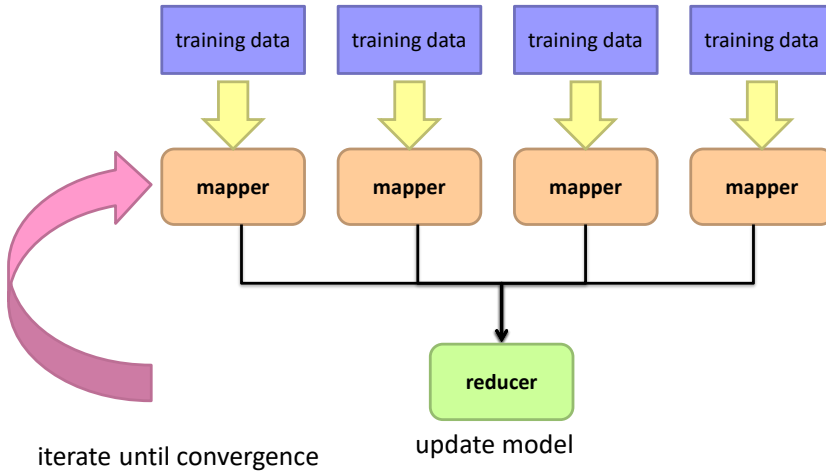


MapReduce Implementation

Reminder

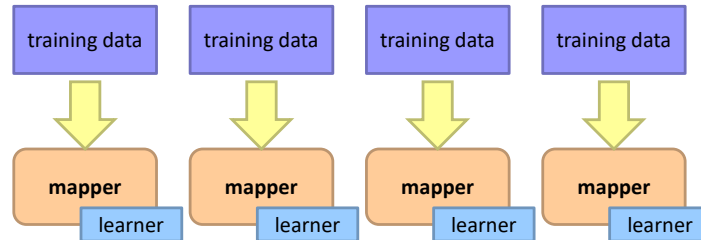
Gradient Descent

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \frac{1}{n} \sum_{i=0}^n \nabla \ell(f(\mathbf{x}_i; \theta^{(t)}), y_i)$$



Stochastic Gradient Descent

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \nabla \ell(f(\mathbf{x}; \theta^{(t)}), y)$$



No iteration!

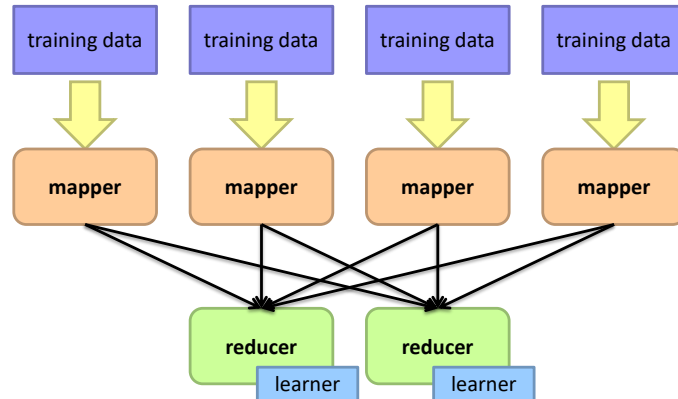
13

This is great because we no longer need iterations!

Mappers go through the record and apply the stochastic gradient descent rule on that record and update the model. This process continues for all records

Stochastic Gradient Descent

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \nabla \ell(f(\mathbf{x}; \theta^{(t)}), y)$$



No iteration!

MapReduce Implementation

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \nabla \ell(f(\mathbf{x}; \theta^{(t)}), y)$$

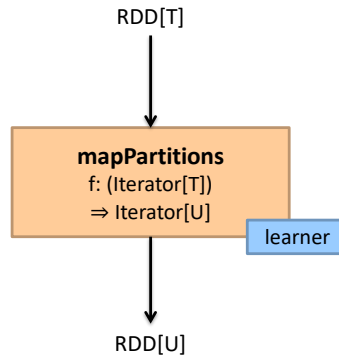
How do we output the model?

Option 1: write model out as “side data”

Option 2: emit model as intermediate output

What about Spark?

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \nabla \ell(f(\mathbf{x}; \theta^{(t)}), y)$$



In practice ...

Data scientists usually use provided transformations in Spark ML

```
val model = LinearRegressionWithSGD.train(parsedData, numIterations, stepSize)
```

```
val prediction = model.predict(point.features)
```

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)

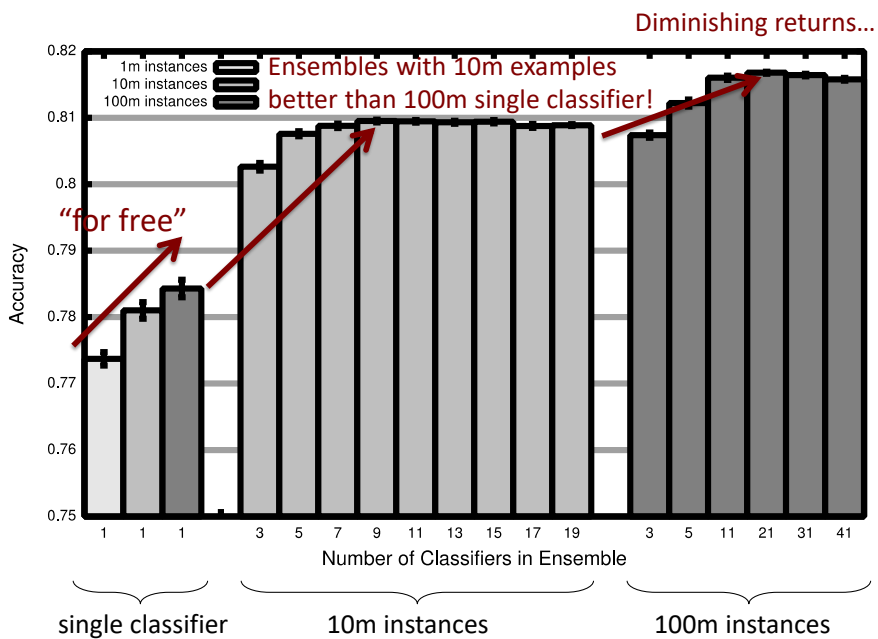
Features:

Sliding window byte-4grams

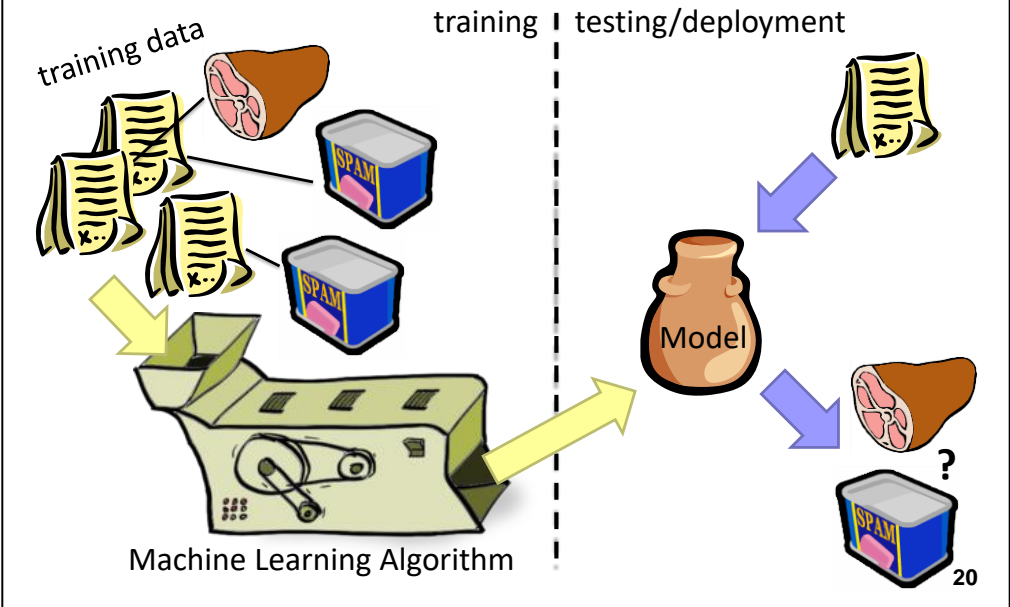
Models + Optimization:

Logistic regression with SGD (L2 regularization)

Ensembles of various sizes (simple weighted voting)



Supervised Machine Learning



Evaluation

How do we know how well we're doing?

Why isn't this enough?

Induce: $f : X \rightarrow Y$

Such that loss is minimized

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

We need end-to-end metrics!

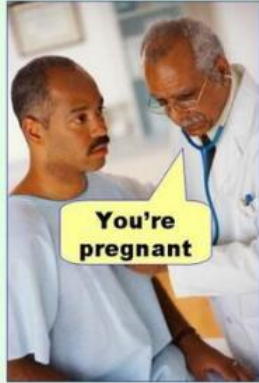
Obvious metric: accuracy

Why isn't this enough?

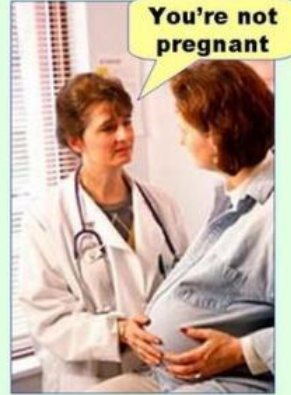
Metrics

		Actual		
		Positive	Negative	
Predicted	Positive	True Positive (TP)	False Positive (FP) = Type 1 Error	Precision = $TP / (TP + FP)$
	Negative	False Negative (FN) = Type II Error	True Negative (TN)	Miss rate = $FN / (FN + TN)$
		Recall or TPR = $TP / (TP + FN)$	Fall-Out or FPR = $FP / (FP + TN)$	

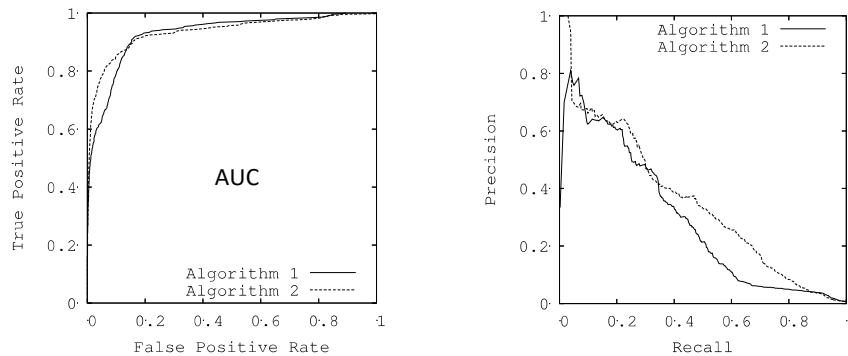
Type I error
(false positive)



Type II error
(false negative)



ROC and PR Curves

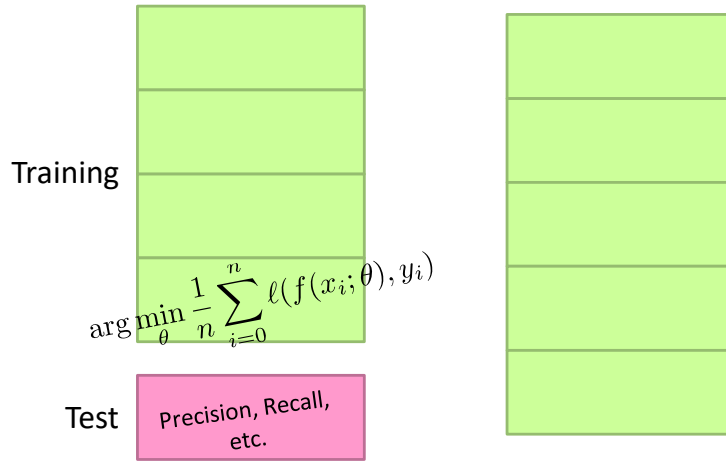


Source: Davis and Goadrich, (2006) The Relationship Between Precision-Recall and ROC curves

24

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.

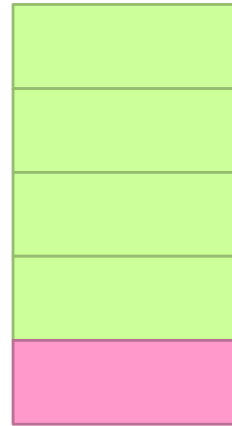
Training/Testing Splits



What happens if you need more?

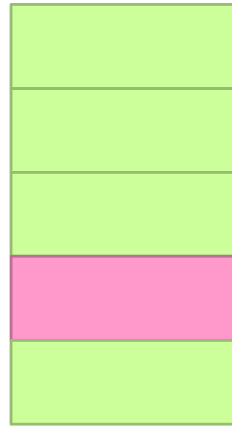
Cross-Validation

Training/Testing Splits



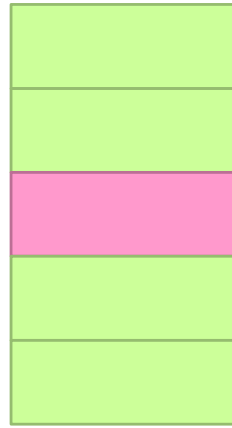
Cross-Validation

Training/Testing Splits



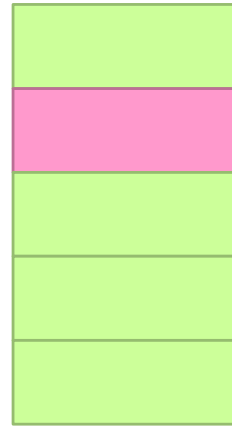
Cross-Validation

Training/Testing Splits



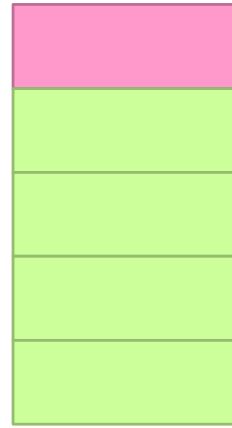
Cross-Validation

Training/Testing Splits



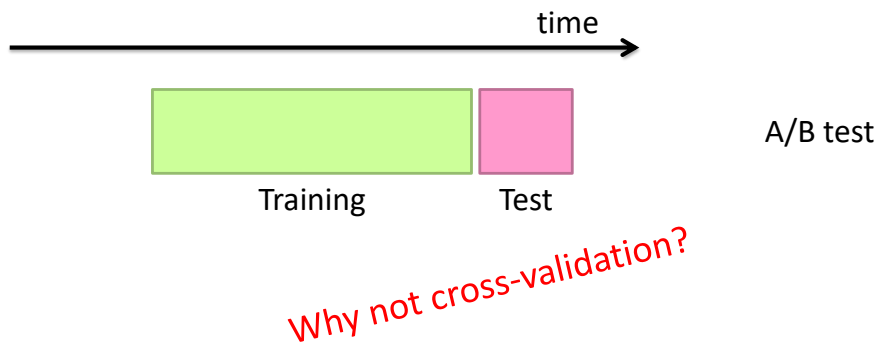
Cross-Validation

Training/Testing Splits

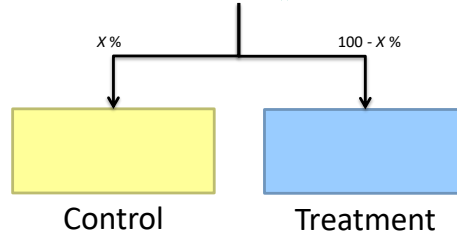


Cross-Validation

Typical Industry Setup



A/B Testing



Gather metrics, compare alternatives

A/B Testing: Complexities

Properly bucketing users

Novelty

Learning effects

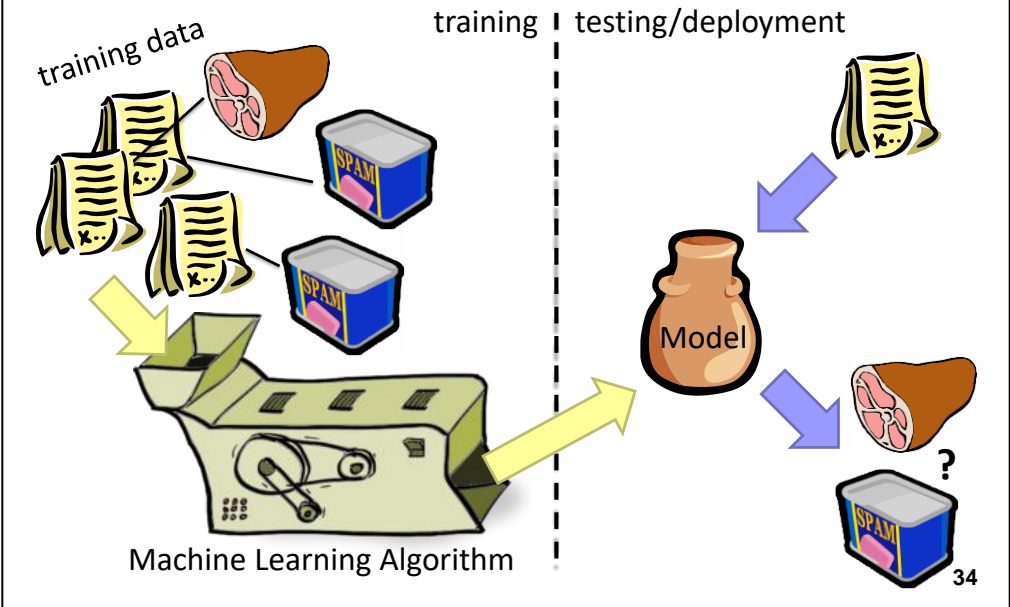
Long vs. short term effects

Multiple, interacting tests

Nosy tech journalists

...

Supervised Machine Learning



Applied ML in Academia

Download interesting dataset (comes with the problem)

Run baseline model

Train/Test

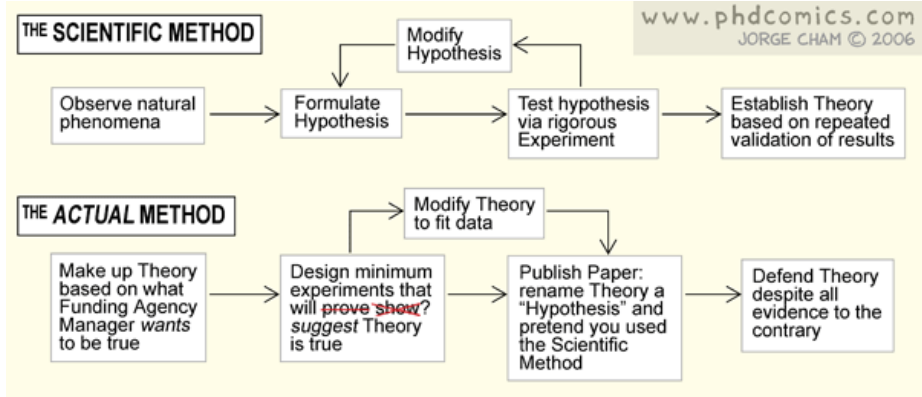
Build better model

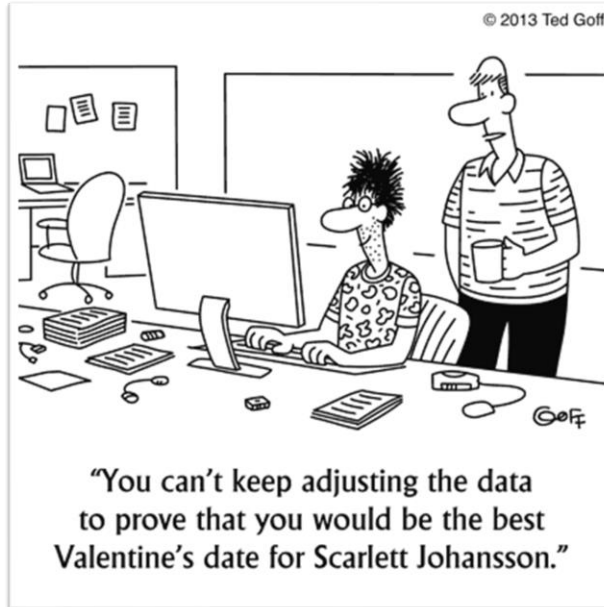
Train/Test

Does new model beat baseline?

Yes: publish a paper!

No: try again!





© 2013 Ted Goff

**“You can’t keep adjusting the data
to prove that you would be the best
Valentine’s date for Scarlett Johansson.”**

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 secret: very little of data science is about machine learning per se!

Data Scientist



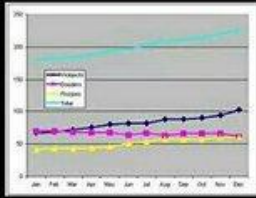
What my friends think I do



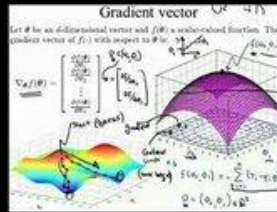
What my mom thinks I do



What society thinks I do



What my boss thinks I do

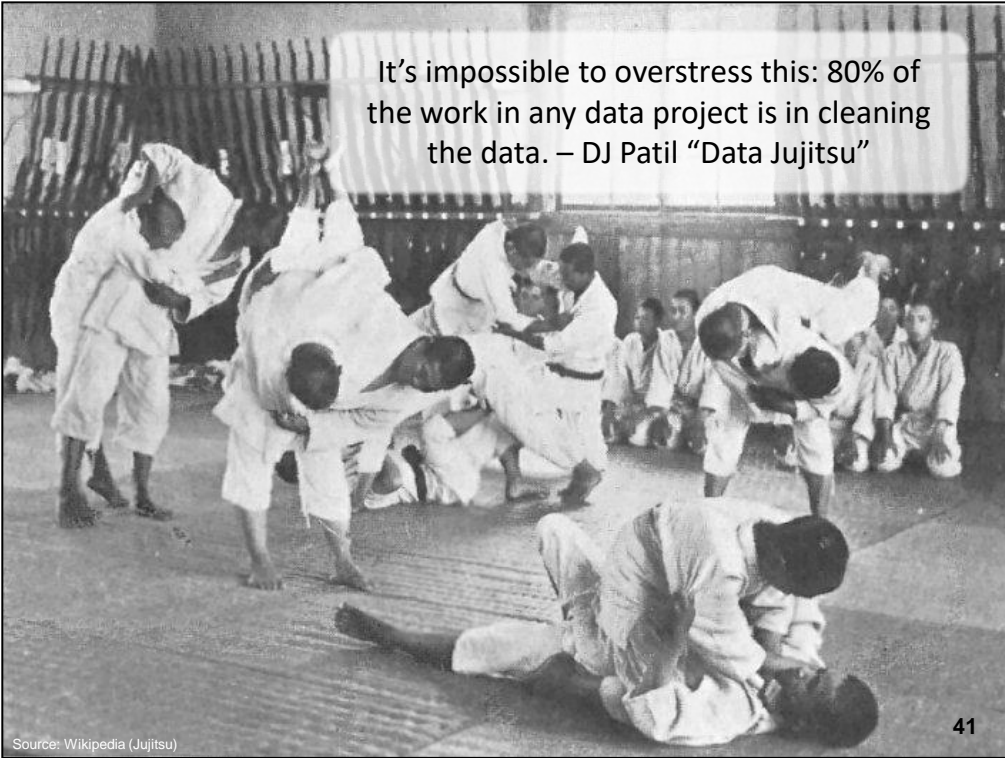


What I think I do



What I actually do

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



Source: Wikipedia (Jujitsu)

41


The New York Times

SECTIONS HOME SEARCH SUBSCRIBE NOW LOG IN

TECHNOLOGY

For 'Big Data' Scientists, Hurdle to Insights Is 'Janitor Work'

By STEVE LOHR AUG. 17, 2014



Monica Rogati, Jawbone's vice president for data science, with Brian Wilt, a senior data scientist.
Peter DaSilva for The New York Times

On finding things...



On naming things...

CamelCase
smallCamelCase
snake_case
camel_Snake
dunder__snake

uid UserId
userId userid
user_id user_id



On feature extraction...

```
^(\w+\s+\d+\s+\d+:\d+:\d+)\s+
(?:@(\S+)\s+(\S+):\s+(\S+)\s+(\S+)
\s+((?:\S+?,\s+)*(?:\S+?))\s+(\S+)\s+(\S+)
\s+\[([^\]]+)\]\s+(\w+)\s+([^\"]\s*\s*
(?:\\\.[^\"]\s*\s*)*\s+(\S+)\s+(\S+)\s+
(\S+)\s+\"([^\"]\s*\s*(?:\\\.[^\"]\s*\s*)*)
\\"\s+\"([^\"]\s*\s*(?:\\\.[^\"]\s*\s*)*)\s*
(\d*-\d-*)?\s*(\d+)?\s*(\d*\.\d\d\.)?
\s+[-\w]+).*
```

An actual Java regular expression used to parse log message at Twitter circa 2010

Friction is cumulative!



[scene: consumer internet company in the Bay Area...]

It's over here...

Well, it wouldn't fit, so we had to shoehorn...

Hang on, I don't remember...

Uh, bad news. Looks like we forgot to log it...

Okay, let's get going... where's the click data?

Well, that's kinda non-intuitive, but okay...

Oh, BTW, where's the timestamp of the click?

[grumble, grumble, grumble]

Frontend Engineer

Develops new feature, adds logging code to capture clicks

Data Scientist

Analyze user behavior, extract insights to improve feature

46

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...

Finally works!

Congratulations, you're halfway there...



Source: Wikipedia (Hills)

48

Congratulations, you're halfway there...

Does it actually work?

A/B testing

Is it fast enough?

Good, you're two thirds there...

Productionize



Source: Wikipedia (Oil refinery)

Productionize

What are your jobs' dependencies?
How/when are your jobs scheduled?
Are there enough resources?
How do you know if it's working?
Who do you call if it stops working?

Infrastructure is critical here!
(plumbing)



Takeaway lesson:
Most of data science isn't glamorous!

Source: Wikipedia (Plumbing)

52