



# Big Data Infrastructure

CS 489/698 Big Data Infrastructure (Winter 2016)

Week 8: Data Mining (1/4)

March 1, 2016

Jimmy Lin

David R. Cheriton School of Computer Science  
University of Waterloo

These slides are available at <http://lintool.github.io/bigdata-2016w/>

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

# Structure of the Course

Analyzing Text

Analyzing Graphs

Analyzing  
Relational Data

Data Mining

“Core” framework features  
and algorithm design

# Supervised Machine Learning

The generic problem of function induction given sample instances of input and output

## Focus today

Classification: output draws from finite discrete labels

Regression: output is a continuous value

This is not meant to be an exhaustive treatment of machine learning!

# Classification



# Applications

Spam detection

Sentiment analysis

Content (e.g., genre) classification

Link prediction

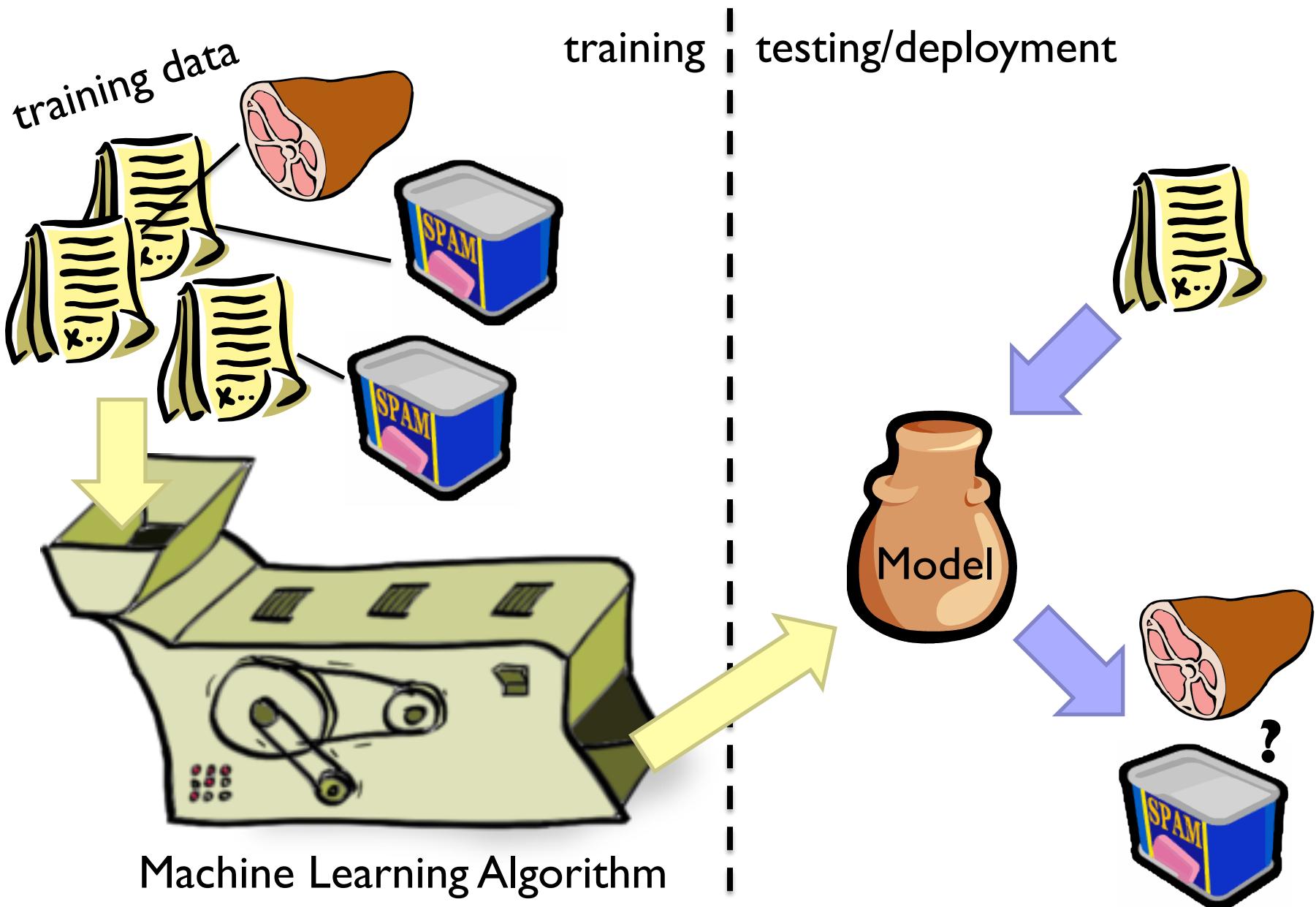
Document ranking

Object recognition

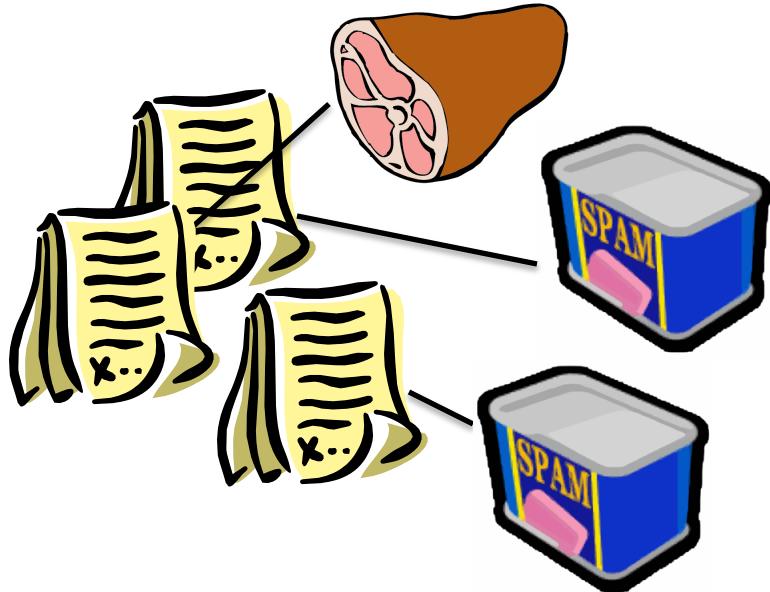
Fraud detection

And much much more!

# Supervised Machine Learning



# Feature Representations



Who comes up with the features?  
How?

Objects are represented in terms of features:

“Dense” features: sender IP, timestamp, # of recipients, length of message, etc.

“Sparse” features: contains the term “viagra” in message, contains “URGENT” in subject, etc.

# Applications

Spam detection

Sentiment analysis

Content (e.g., genre) classification

Link prediction

Document ranking

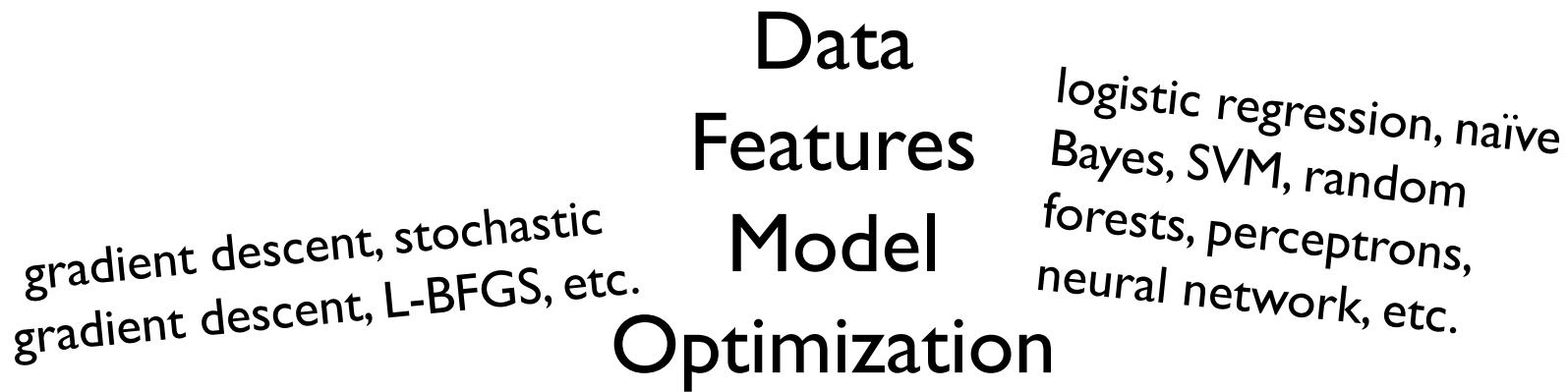
Object recognition

Fraud detection

And much much more!

Features are highly  
application-specific!

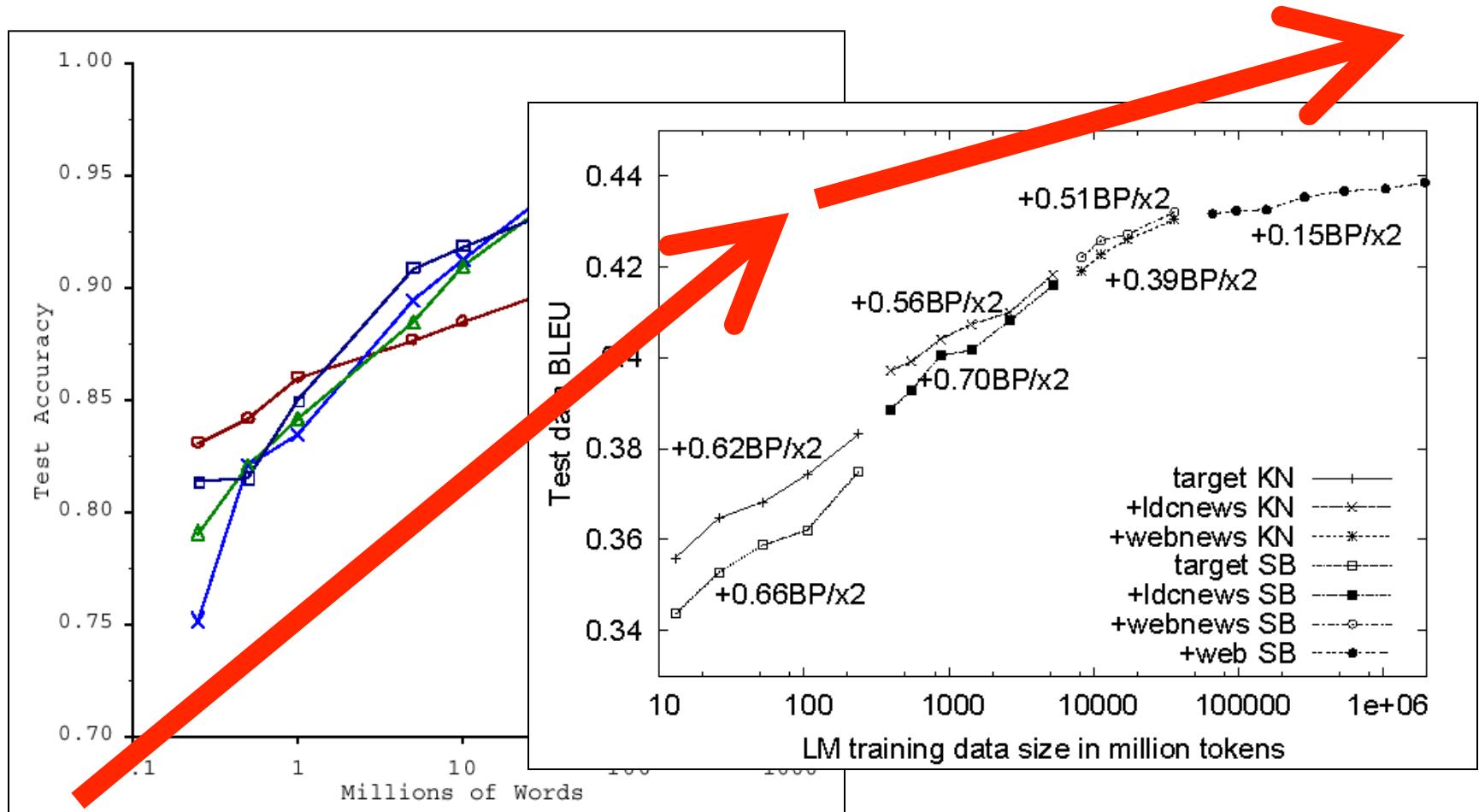
# Components of a ML Solution



What “matters” the most?

# No data like more data!

s/knowledge/data/g;



# Limits of Supervised Classification?

- Why is this a big data problem?
  - Isn't gathering labels a serious bottleneck?
- Solution: crowdsourcing
- Solution: user behavior logs
  - Learning to rank
  - Computational advertising
  - Link recommendation
- The virtuous cycle of data-driven products

# Supervised Binary Classification

- Restrict output label to be *binary*
  - Yes/No
  - 1/0
- Binary classifiers form a primitive building block for multi-class problems
  - One vs. rest classifier ensembles
  - Classifier cascades

# The Task

- Given  $D = \{(x_i, y_i)\}_i^n$   


↓ label  
↑ (sparse) feature vector

$$x_i = [x_1, x_2, x_3, \dots, x_d]$$

$$y \in \{0, 1\}$$

- Induce  $f : X \rightarrow Y$ 
  - Such that loss is minimized
- $$\frac{1}{n} \sum_{i=0}^n \ell(f(x_i), y_i)$$


↑ loss function
- Typically, consider functions of a parametric form:

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


↑ model parameters

**Key insight: machine learning as an optimization problem!**  
**(closed form solutions generally not possible)**

# Gradient Descent: Preliminaries

- Rewrite:

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=0}^n \ell(f(\mathbf{x}_i; \theta), y_i) \quad \longrightarrow \quad \arg \min_{\theta} L(\theta)$$

- Compute gradient:

- “Points” to fastest increasing “direction”

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

- So, at any point: \*

$$\mathbf{b} = \mathbf{a} - \gamma \nabla L(\mathbf{a})$$

$$L(\mathbf{a}) \geq L(\mathbf{b})$$

\* caveats

# Gradient Descent: Iterative Update

- Start at an arbitrary point, iteratively update:

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \nabla L(\theta^{(t)})$$

- We have:

$$L(\theta^{(0)}) \geq L(\theta^{(1)}) \geq L(\theta^{(2)}) \dots$$

- Lots of details:

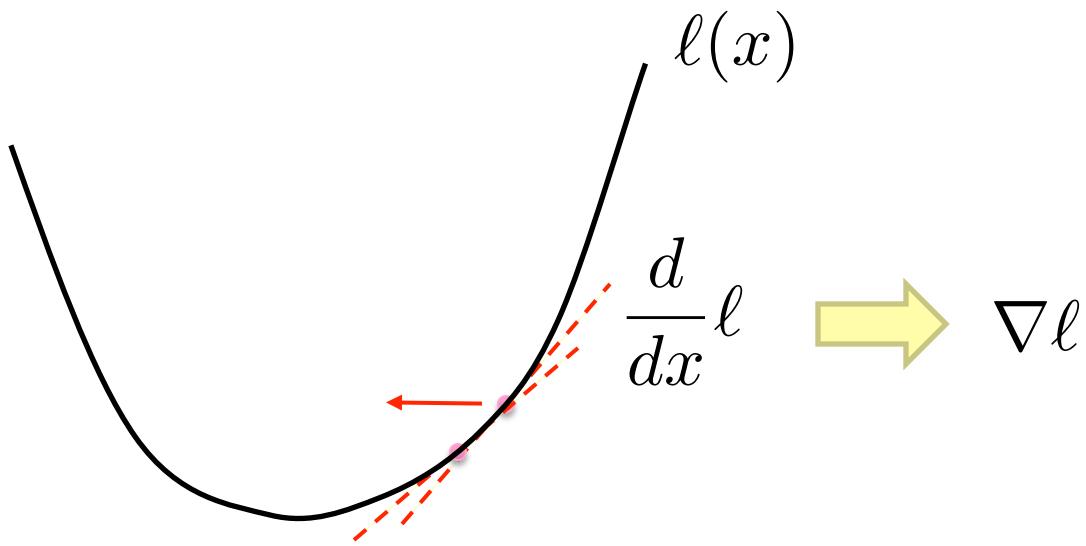
- Figuring out the step size
- Getting stuck in local minima
- Convergence rate
- ...

# Gradient Descent

Repeat until convergence:

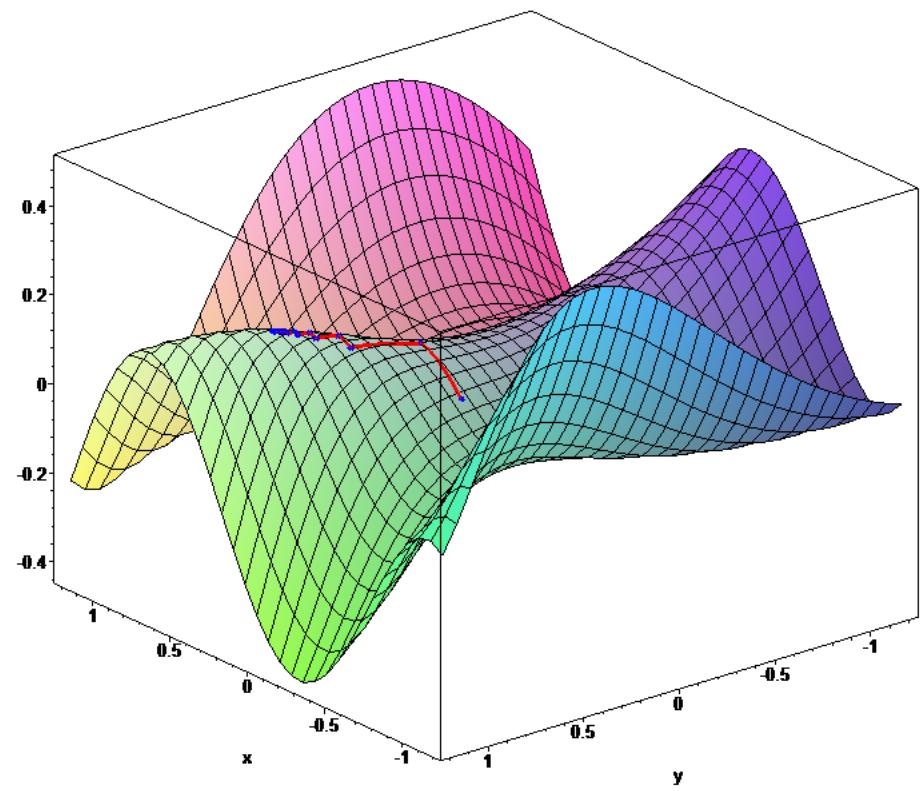
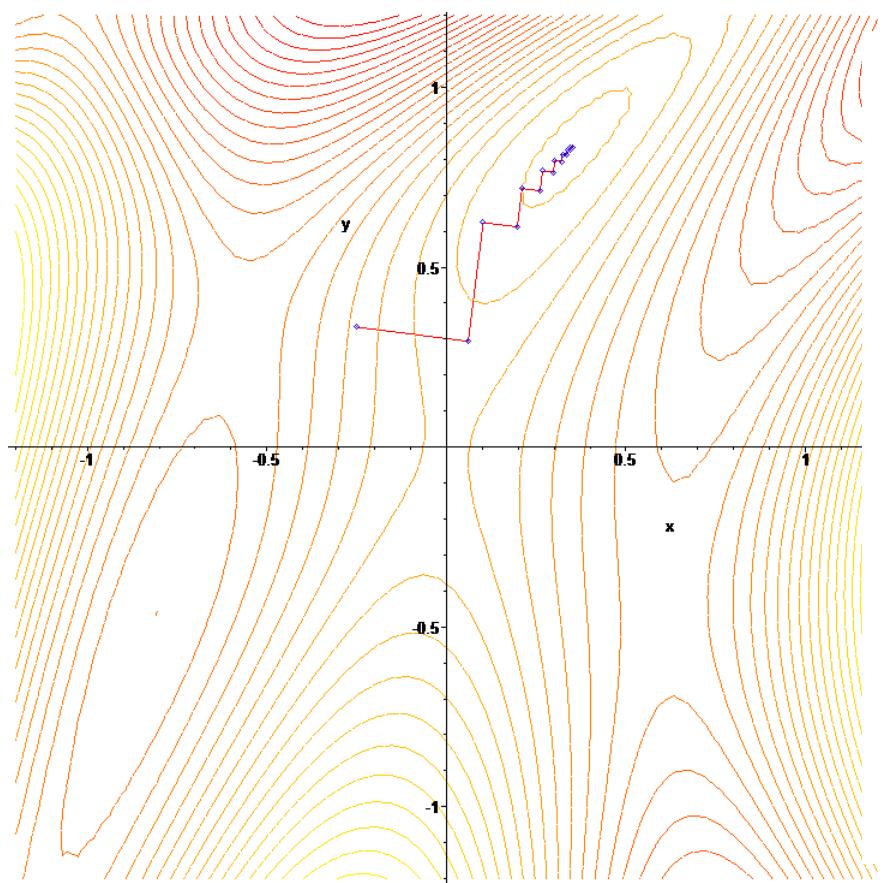
$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \frac{1}{n} \sum_{i=0}^n \nabla \ell(f(\mathbf{x}_i; \theta^{(t)}), y_i)$$

# Intuition behind the math...



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

New weights    Old weights                      Update based on gradient



The background image shows a wide, open landscape with rolling green hills. The sky above is a vibrant blue, filled with large, white, fluffy clouds. The foreground is a mix of green grass and some brown, possibly dry, areas. In the distance, more hills and mountains are visible under the same cloudy sky.

# Gradient Descent

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

# Lots More Details...

- Gradient descent is a “first order” optimization technique
  - Often, slow convergence
  - Conjugate techniques accelerate convergence
- Newton and quasi-Newton methods:
  - Intuition: Taylor expansion

$$f(x + \Delta x) = f(x) + f'(x)\Delta x + \frac{1}{2}f''(x)\Delta x^2$$

- Requires the Hessian (square matrix of second order partial derivatives): impractical to fully compute

# Logistic Regression



# Logistic Regression: Preliminaries

- Given  $D = \{(x_i, y_i)\}_i^n$

$$x_i = [x_1, x_2, x_3, \dots, x_d]$$
$$y \in \{0, 1\}$$

- Let's define:

$$f(x; w) : \mathbb{R}^d \rightarrow \{0, 1\}$$

$$f(x; w) = \begin{cases} 1 & \text{if } w \cdot x \geq t \\ 0 & \text{if } w \cdot x < t \end{cases}$$

- Interpretation:

$$\ln \left[ \frac{\Pr(y=1|x)}{\Pr(y=0|x)} \right] = w \cdot x$$

$$\ln \left[ \frac{\Pr(y=1|x)}{1 - \Pr(y=1|x)} \right] = w \cdot x$$

# Relation to the Logistic Function

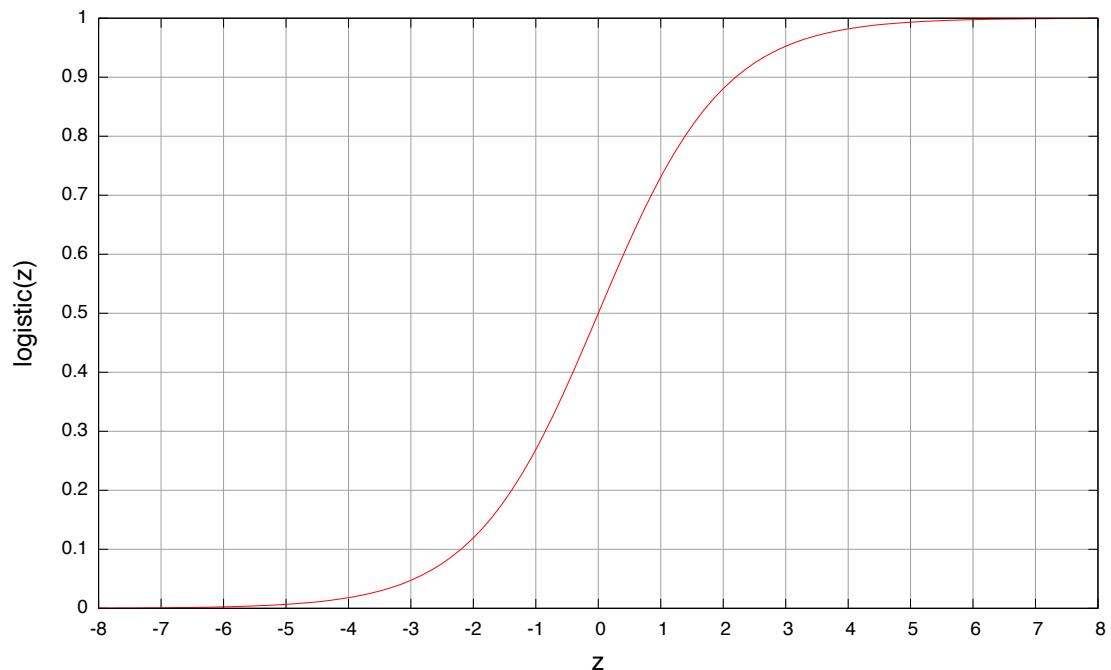
- After some algebra:

$$\Pr(y = 1|x) = \frac{e^{w \cdot x}}{1 + e^{w \cdot x}}$$

$$\Pr(y = 0|x) = \frac{1}{1 + e^{w \cdot x}}$$

- The logistic function:

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



# Training an LR Classifier

- Maximize the conditional likelihood:

$$\arg \max_w \prod_{i=1}^n \Pr(y_i | \mathbf{x}_i, w)$$

- Define the objective in terms of conditional log likelihood:

$$L(w) = \sum_{i=1}^n \ln \Pr(y_i | \mathbf{x}_i, w)$$

- We know  $y \in \{0, 1\}$  so:

$$\Pr(y | \mathbf{x}, w) = \Pr(y = 1 | \mathbf{x}, w)^y \Pr(y = 0 | \mathbf{x}, w)^{(1-y)}$$

- Substituting:

$$L(w) = \sum_{i=1}^n \left( y_i \ln \Pr(y_i = 1 | \mathbf{x}_i, w) + (1 - y_i) \ln \Pr(y_i = 0 | \mathbf{x}_i, w) \right)$$

# LR Classifier Update Rule

- Take the derivative:

$$L(\mathbf{w}) = \sum_{i=1}^n \left( y_i \ln \Pr(y_i = 1 | \mathbf{x}_i, \mathbf{w}) + (1 - y_i) \ln \Pr(y_i = 0 | \mathbf{x}_i, \mathbf{w}) \right)$$

$$\frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}) = \sum_{i=0}^n \mathbf{x}_i \left( y_i - \Pr(y_i = 1 | \mathbf{x}_i, \mathbf{w}) \right)$$

- General form for update rule:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \gamma^{(t)} \nabla_{\mathbf{w}} L(\mathbf{w}^{(t)})$$

$$\nabla L(\mathbf{w}) = \left[ \frac{\partial L(\mathbf{w})}{\partial w_0}, \frac{\partial L(\mathbf{w})}{\partial w_1}, \dots, \frac{\partial L(\mathbf{w})}{\partial w_d} \right]$$

- Final update rule:

$$\mathbf{w}_i^{(t+1)} \leftarrow \mathbf{w}_i^{(t)} + \gamma^{(t)} \sum_{j=0}^n x_{j,i} \left( y_j - \Pr(y_j = 1 | \mathbf{x}_j, \mathbf{w}^{(t)}) \right)$$

# Lots more details...

- Regularization
- Different loss functions
- ...

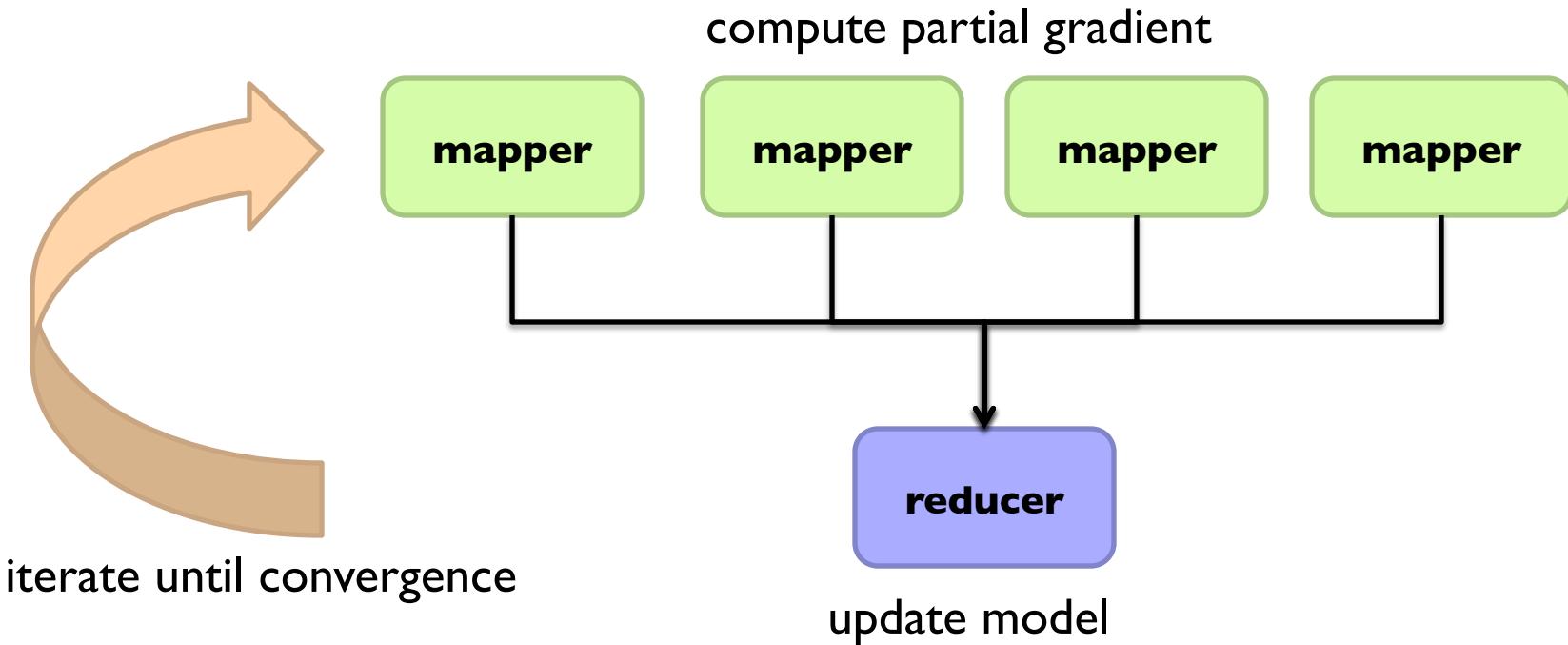
Want more details?  
Take a real machine-learning course!

# MapReduce Implementation

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

mappers

single reducer



# Shortcomings

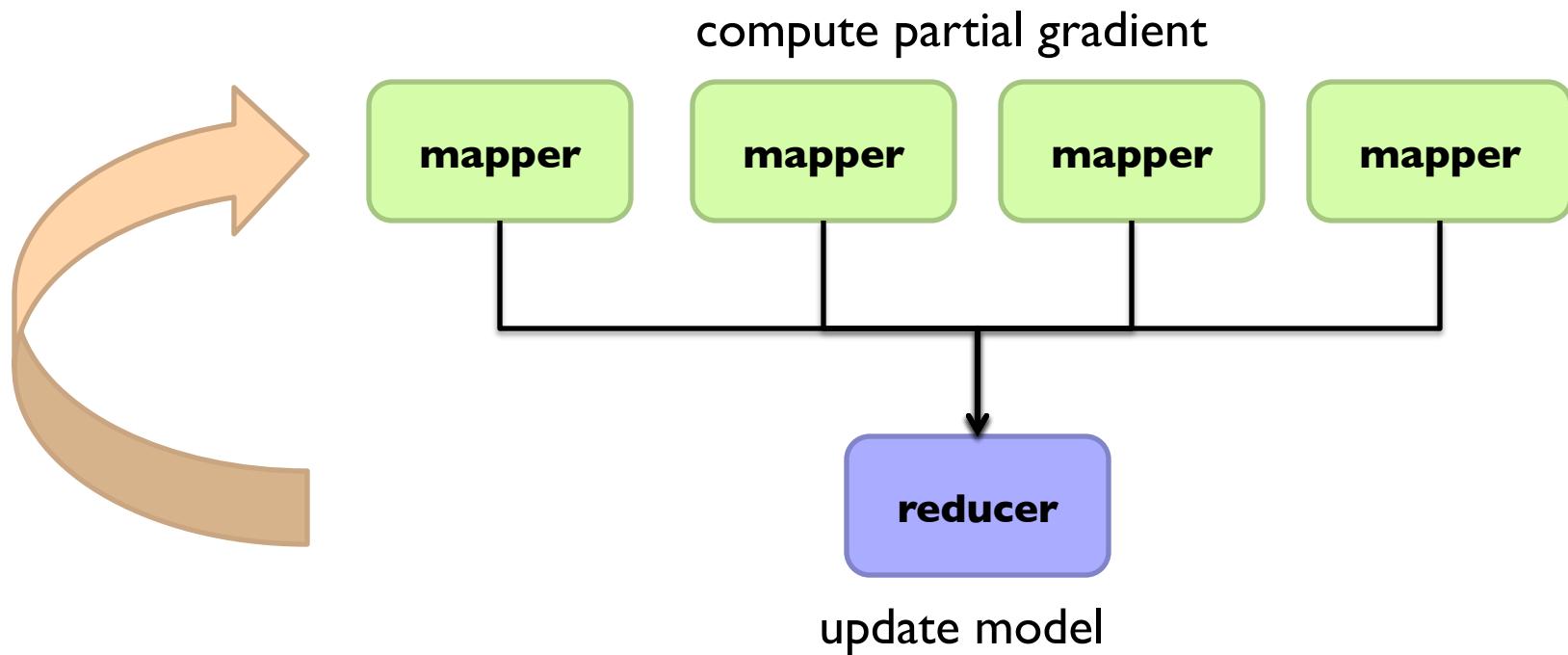
- Hadoop is bad at iterative algorithms
  - High job startup costs
  - Awkward to retain state across iterations
- High sensitivity to skew
  - Iteration speed bounded by slowest task
- Potentially poor cluster utilization
  - Must shuffle all data to a single reducer
- Some possible tradeoffs
  - Number of iterations vs. complexity of computation per iteration
  - E.g., L-BFGS: faster convergence, but more to compute

# Spark Implementation

```
val points = spark.textFile(...).map(parsePoint).persist()
```

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

What's the difference?



# Gradient Descent

A photograph of a vibrant water park featuring a complex of multi-colored water slides (yellow, blue, green, orange, purple) winding through a steel frame structure. In the foreground, a large yellow splash pool is visible, with several people in swimwear playing in the water. The background shows a clear blue sky with scattered white clouds and some green trees.

# Stochastic Gradient Descent

# Batch vs. Online

## Gradient Descent

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

“batch” learning: update model after considering all training instances

## Stochastic Gradient Descent (SGD)

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

“online” learning: update model after considering each (randomly-selected) training instance

In practice... just as good!

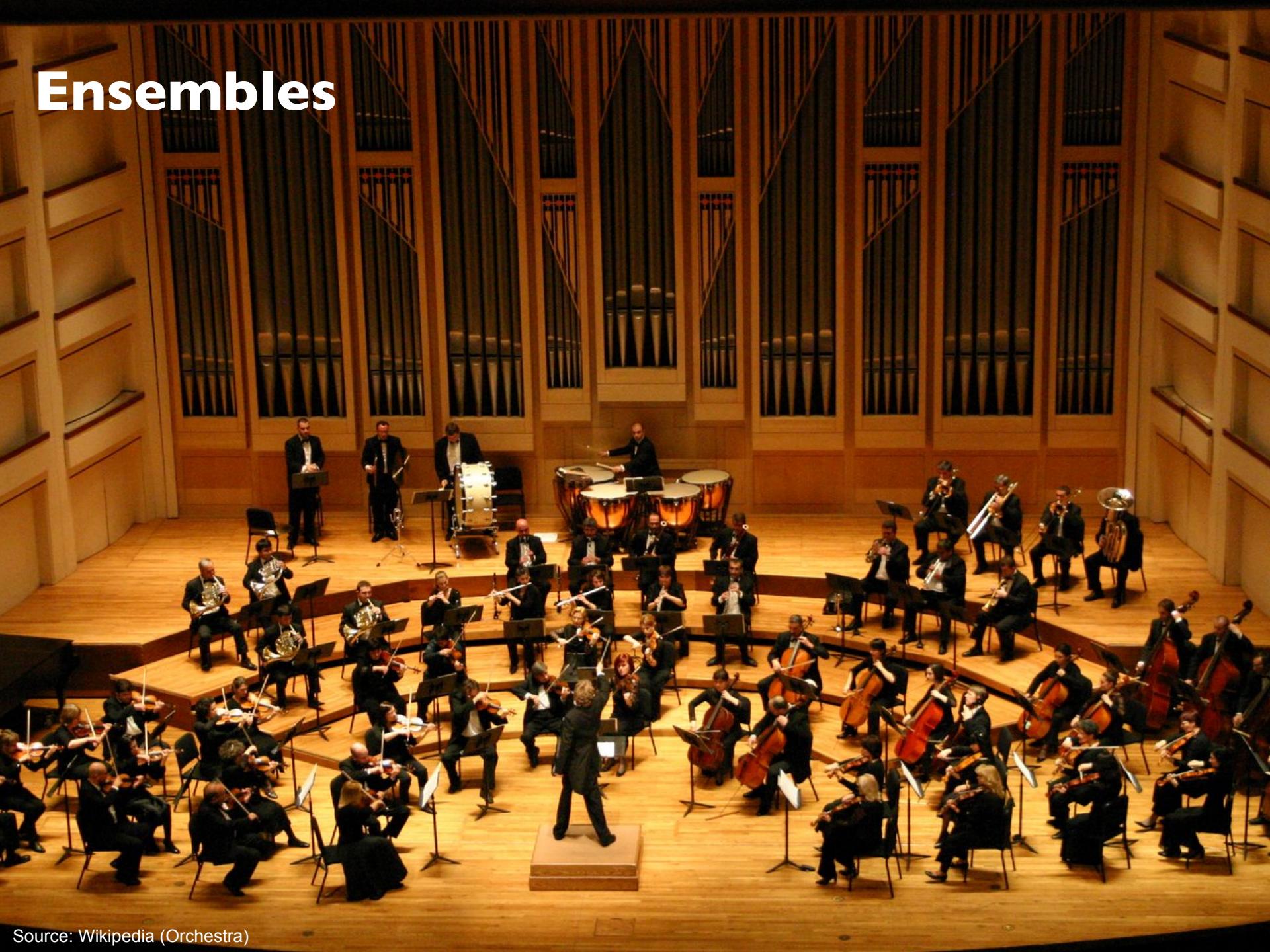
# Practical Notes

- Most common implementation:
  - Randomly shuffle training instances
  - Stream instances through learner
- Single vs. multi-pass approaches
- “Mini-batching” as a middle ground between batch and stochastic gradient descent

We've solved the iteration problem!

What about the single reducer problem?

# Ensembles



# Ensemble Learning

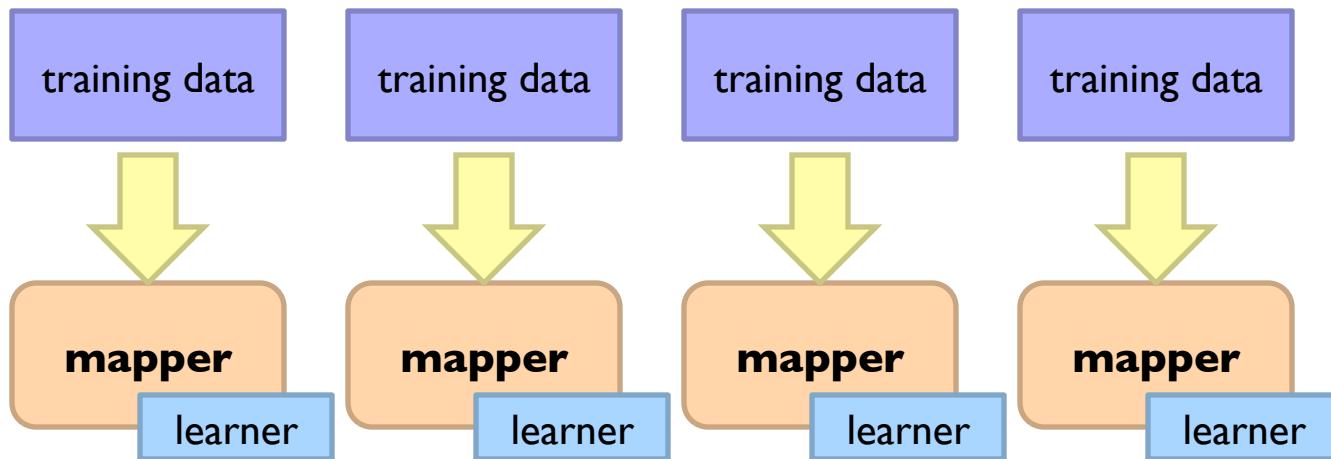
- Learn multiple models, combine results from different models to make prediction
- Why does it work?
  - If errors uncorrelated, multiple classifiers being wrong is less likely
  - Reduces the variance component of error
- A variety of different techniques:
  - Majority voting
  - Simple weighted voting:
$$y = \arg \max_{y \in Y} \sum_{k=1}^n \alpha_k p_k(y|x)$$
  - Model averaging
  - ...

# Practical Notes

- Common implementation:
  - Train classifiers on different input partitions of the data
  - Embarrassingly parallel!
- Contrast with other ensemble techniques, e.g., boosting

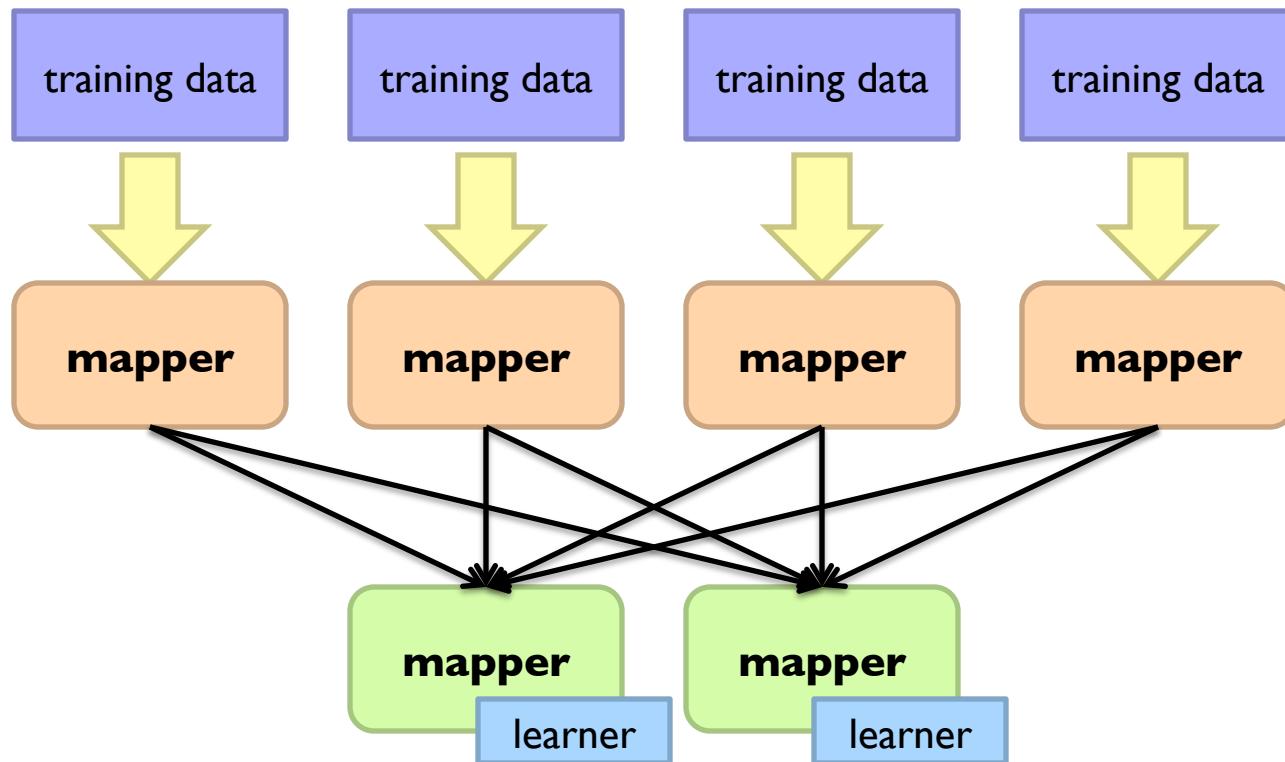
# MapReduce Implementation

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



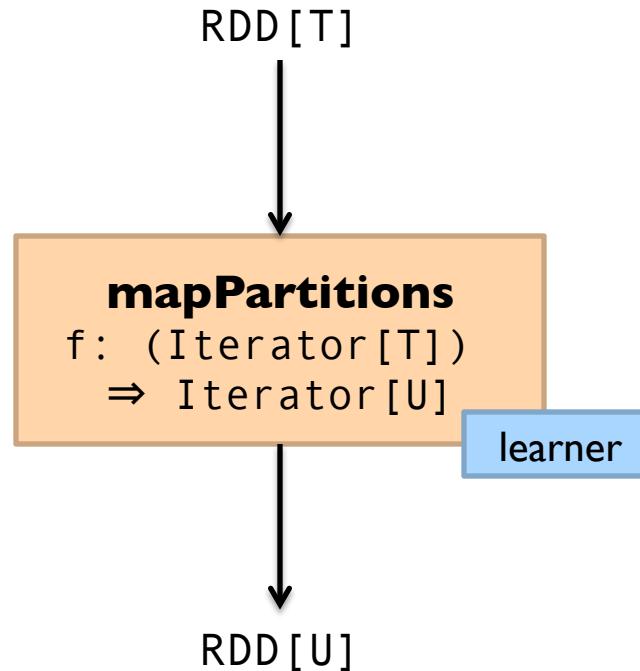
# MapReduce Implementation

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



# What about Spark?

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



# MapReduce Implementation: Details

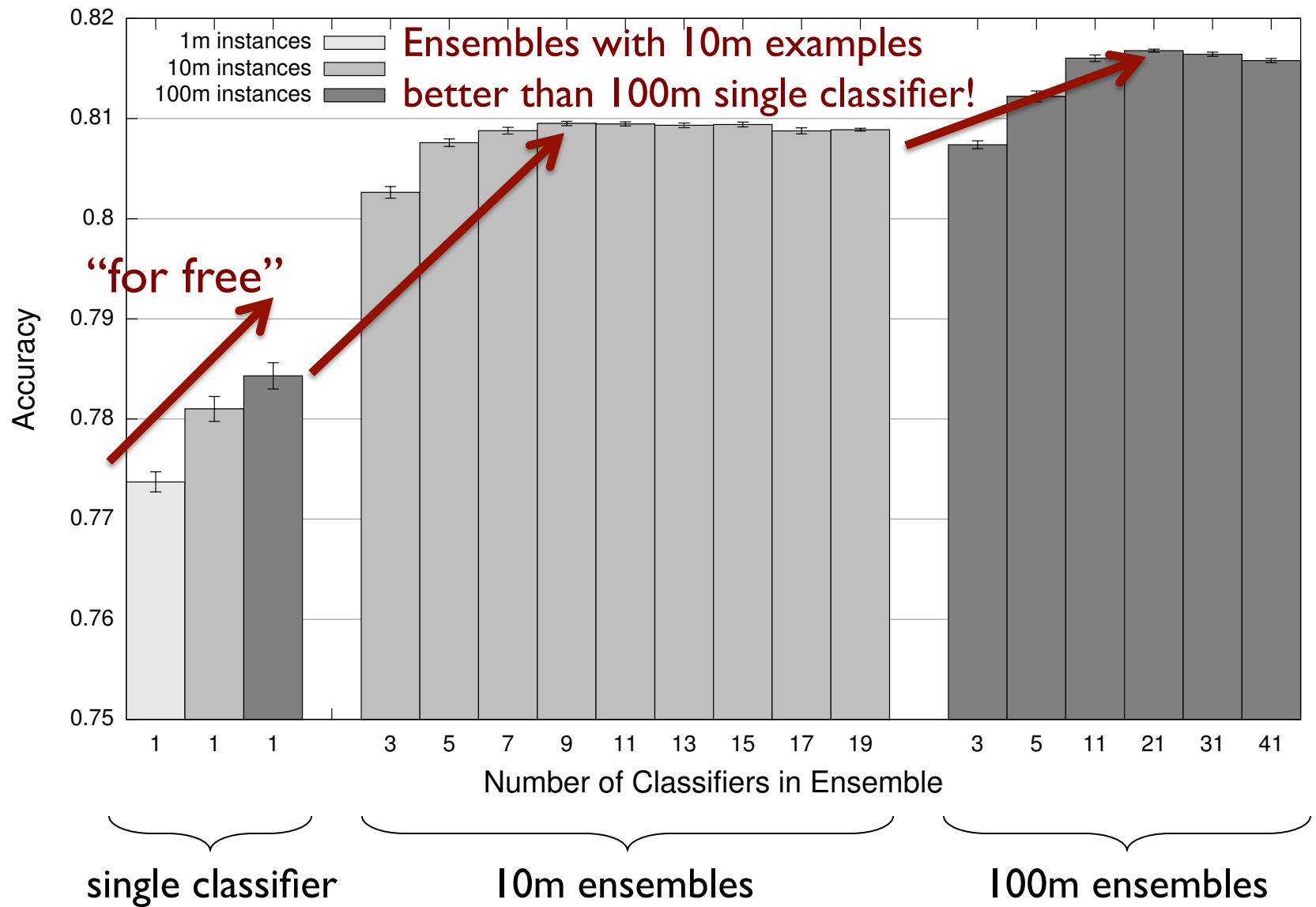
- Two possible implementations:
  - Write model out as “side data”
  - Emit model as intermediate output

# Sentiment Analysis Case Study

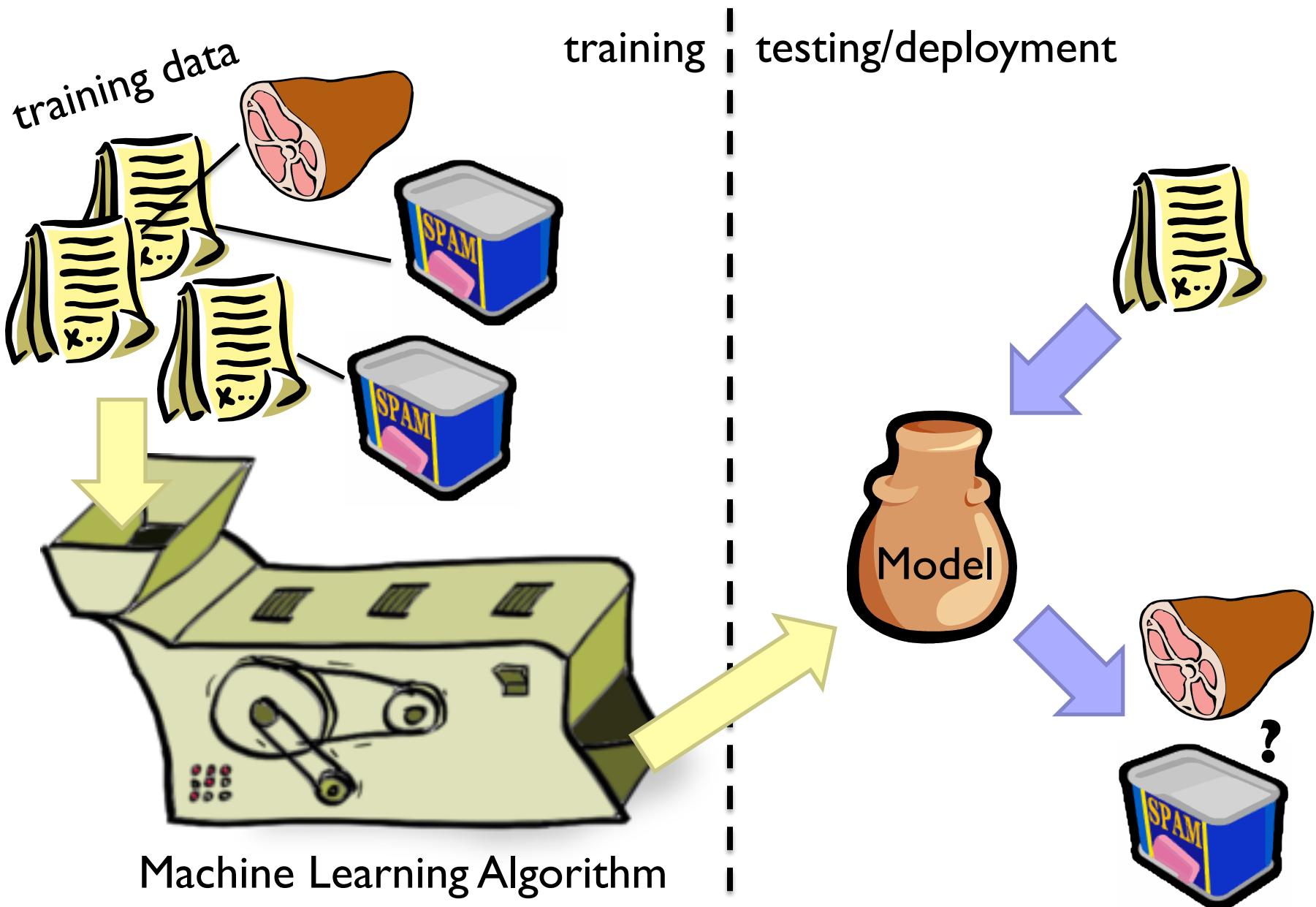
Lin and Kolcz, SIGMOD 2012

- Binary polarity classification: {positive, negative} sentiment
  - Independently interesting task
  - Illustrates end-to-end flow
  - 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:
  - Logistic regression with SGD (L2 regularization)
  - Ensembles of various sizes (simple weighted voting)

Diminishing returns...



# Supervised Machine Learning



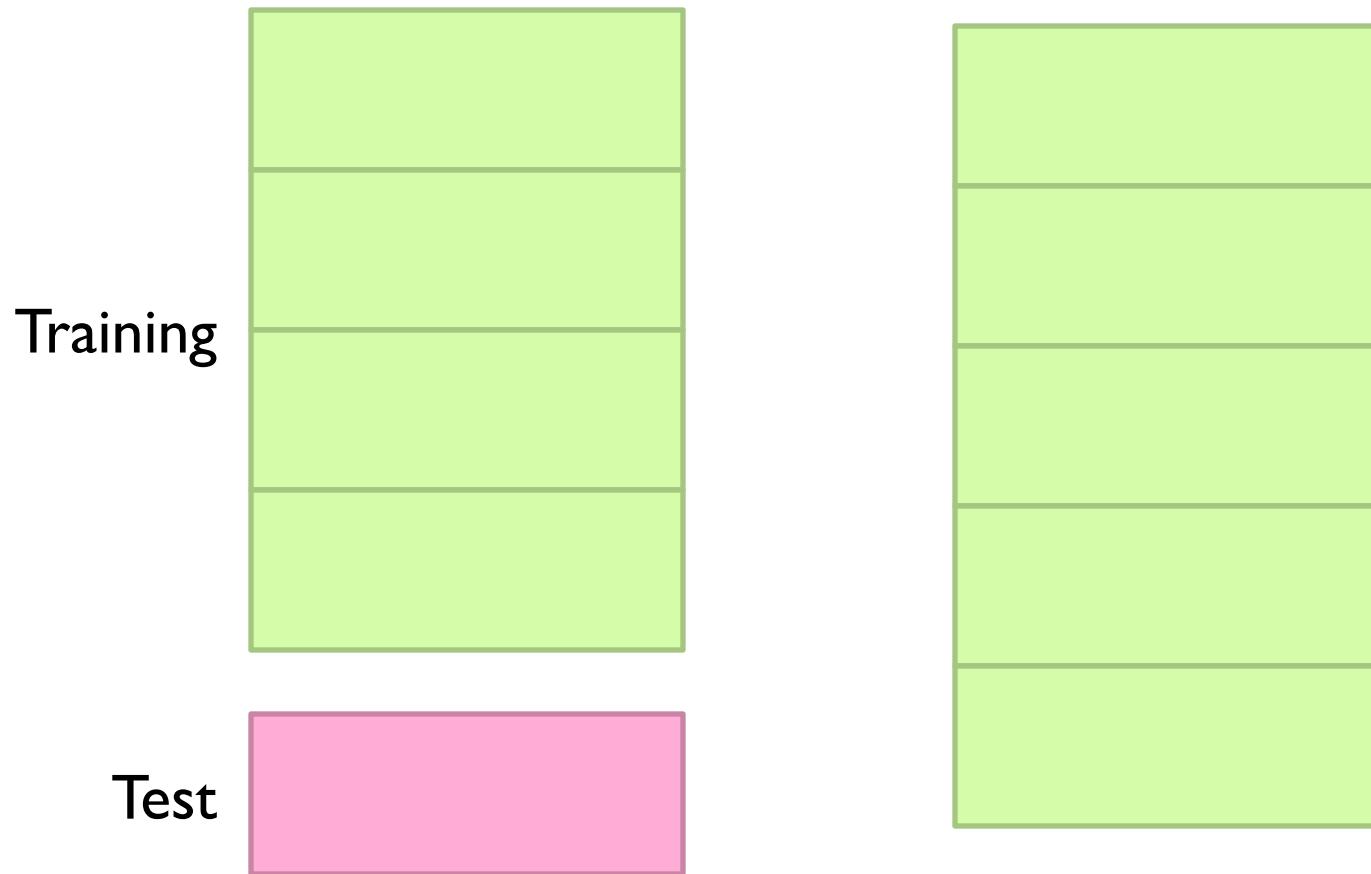
# Applied ML in Academia

- 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!

# Three Commandants of Machine Learning

Thou shalt not mix training and testing data  
Thou shalt not mix training and testing data  
Thou shalt not mix training and testing data

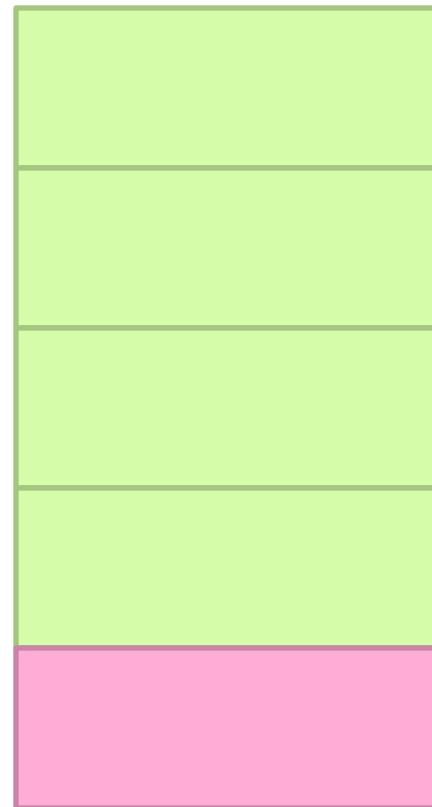
# Training/Testing Splits



What happens if you need more?

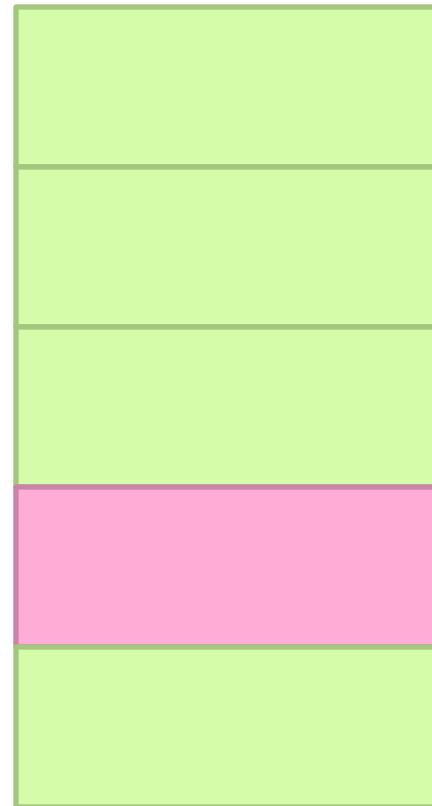
Cross-Validation

# Training/Testing Splits



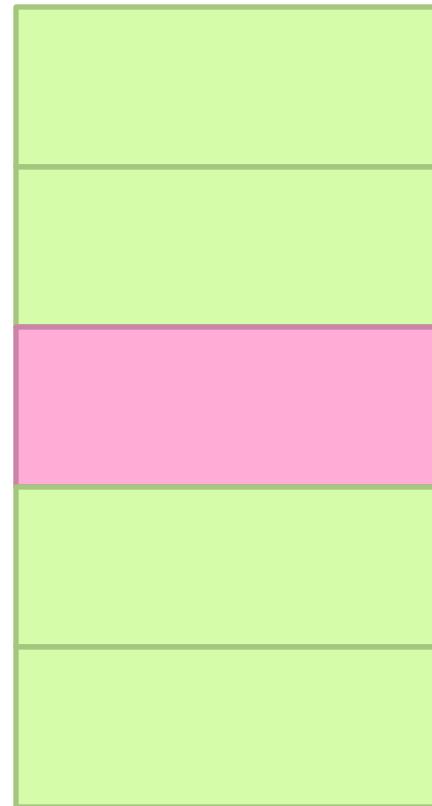
Cross-Validation

# Training/Testing Splits



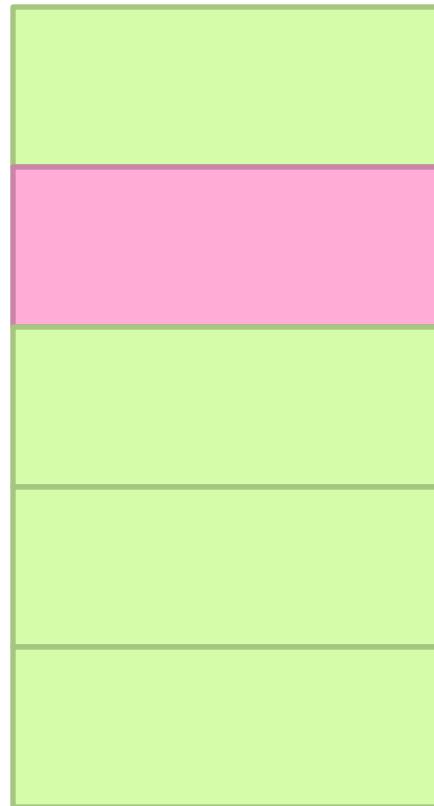
Cross-Validation

# Training/Testing Splits



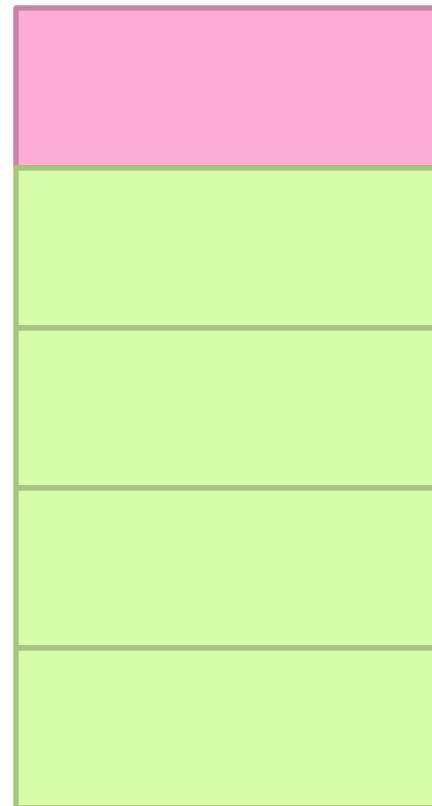
Cross-Validation

# Training/Testing Splits



Cross-Validation

# Training/Testing Splits



Cross-Validation

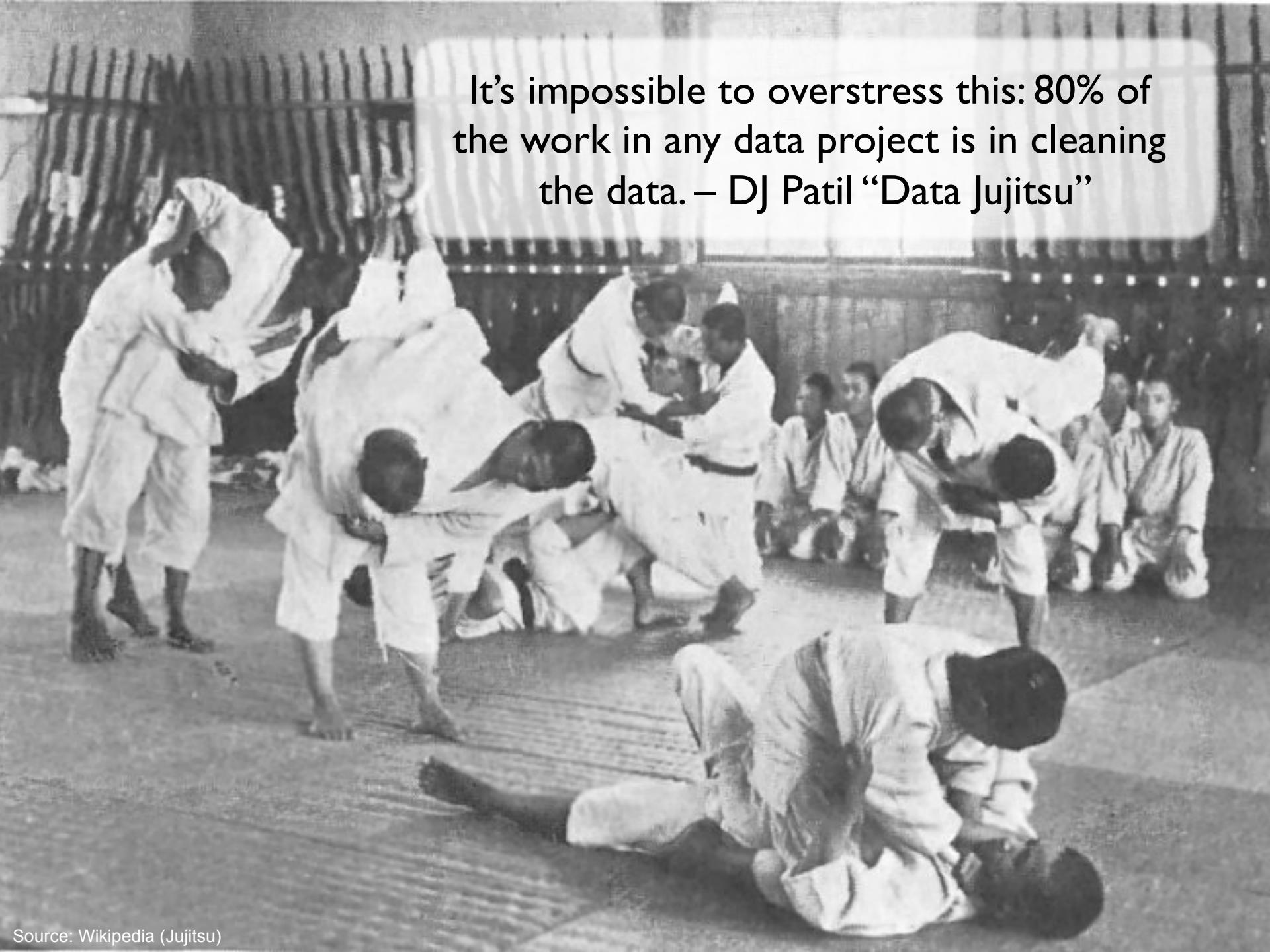
# Applied ML in Academia

- 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!

# Applied ML in Industry

- Formulate the problem (from some vague business goal)
- Find the data
- Prepare and clean the data
- Run data through machine learning library
  - Train/test on retrospective data
- Does the new model beat baseline?
  - Yes: Run live A/B test
  - No: try again!
- Does the new model work in A/B testing?
  - Yes: productionize the model
  - No: try again!

Where's most of the time spent?



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

SUBSCRIBE NOW

LOG IN



The New York Times

≡ SECTIONS

HOME

SEARCH

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

# Components of a ML Solution

Data  
Features  
Model  
Optimization

What “matters” the most?

A photograph of a traditional Japanese rock garden. In the foreground, a gravel path is raked into fine, parallel lines. Several large, dark, irregular stones are scattered across the garden. A small, shallow pond is visible in the middle ground, surrounded by more stones and low-lying green plants. In the background, there are more trees and shrubs, and the wooden buildings of a residence are visible behind the garden wall.

Questions?