## Data-Intensive Distributed Computing CS 451/651 (Fall 2025)



# Data Infrastructure for Machine Learning (v1.00)

Week 6: October 7

Jimmy Lin
David R. Cheriton School of Computer Science
University of Waterloo

These slides are available at http://lintool.github.io/bigdata-2025f/



## Key Questions

What are the key components of an ML solution?

How is the supervised machine learning problem formulated?

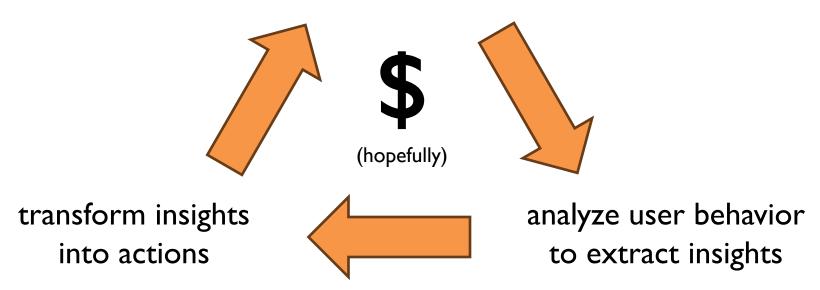
What roles do data platforms and data engineering play?

#### Context...

## The Data Flywheel

(a virtuous cycle)

Build a useful product



Google. Facebook. Twitter. Amazon. Uber.

Context...

#### What's this course about?

The infrastructure that supports the data flywheel.

data platforms + data engineering

# Context... Transform Insights into Actions What does that really mean?

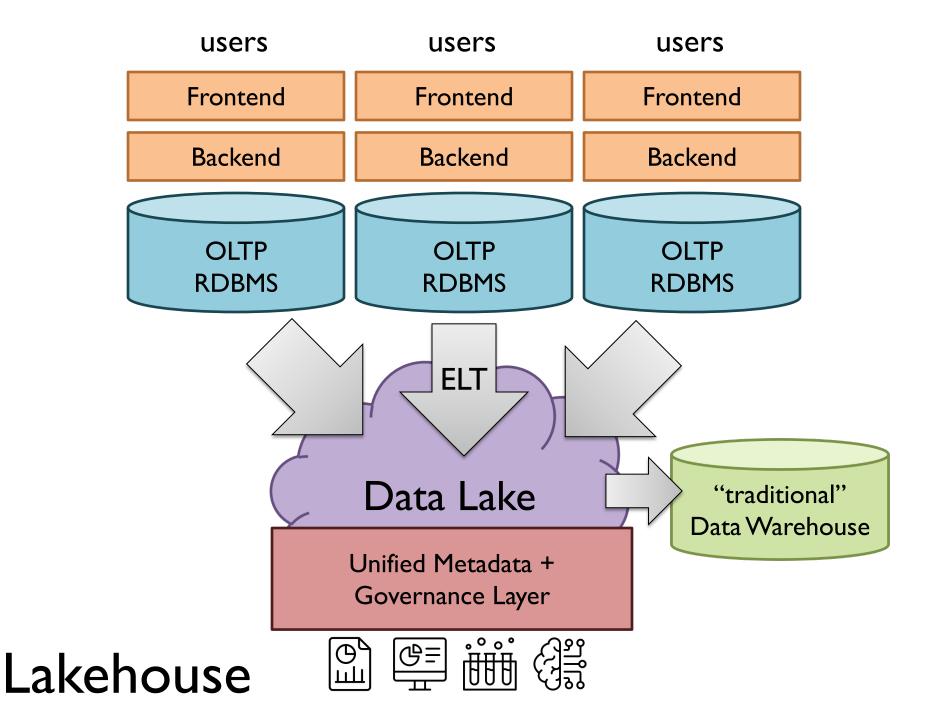
Report generation
Dashboards

Ad hoc analyses
ML models

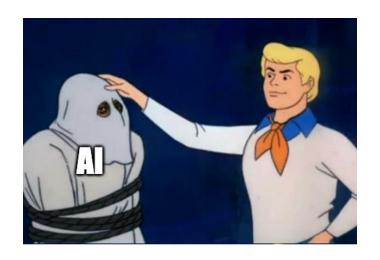
Business Intelligence

Business Intelligence

Data Science









When people talk about AI these days, they really mean ML...

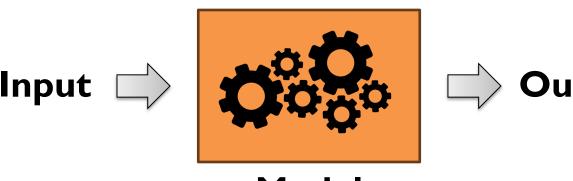




When people talk about ML these days, they really mean supervised ML...



What's supervised machine learning?



**Model** 

(accomplishes some task) (hopefully well)



owl

Input





**Output** 



cat

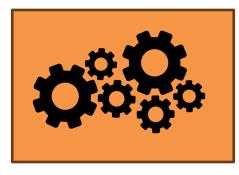
(accomplishes some task) (hopefully well)

amazing spot for good food & a fun time they offer a super unique dine-in experience with their interactive tables! also love that they have innovative weekly feature dishes

#### Input



Worst service I have ever experienced! Waited 25 mins for our waitress to come to our table once seated, no apology! Waited another 30 mins for our drinks, which we had to ask about.











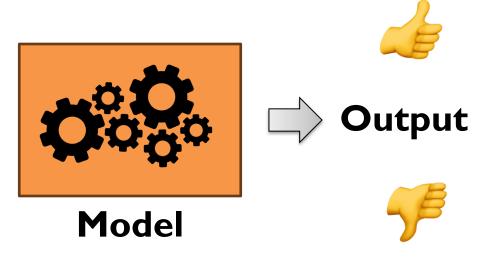
(accomplishes some task) (hopefully well)

amazing spot for good food & a fun time they offer a super unique dine-in experience with their interactive tables! also love that they have innovative weekly feature dishes

#### Input

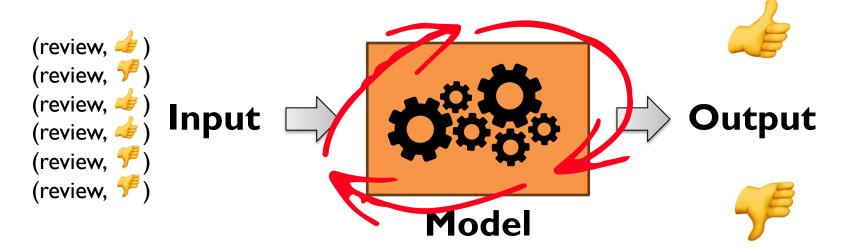


Worst service I have ever experienced! Waited 25 mins for our waitress to come to our table once seated, no apology! Waited another 30 mins for our drinks, which we had to ask about.



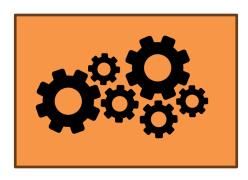
#### How do you build such a model?

#### Model learns from the data



Machine learning algorithm adjusts the model parameters

How do you build such a model?



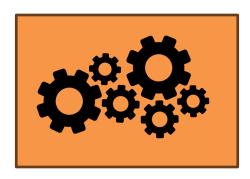
#### Model

Machine learning algorithm adjusts the model parameters

What does that actually mean?

#### Typically, the model takes a parametric form:

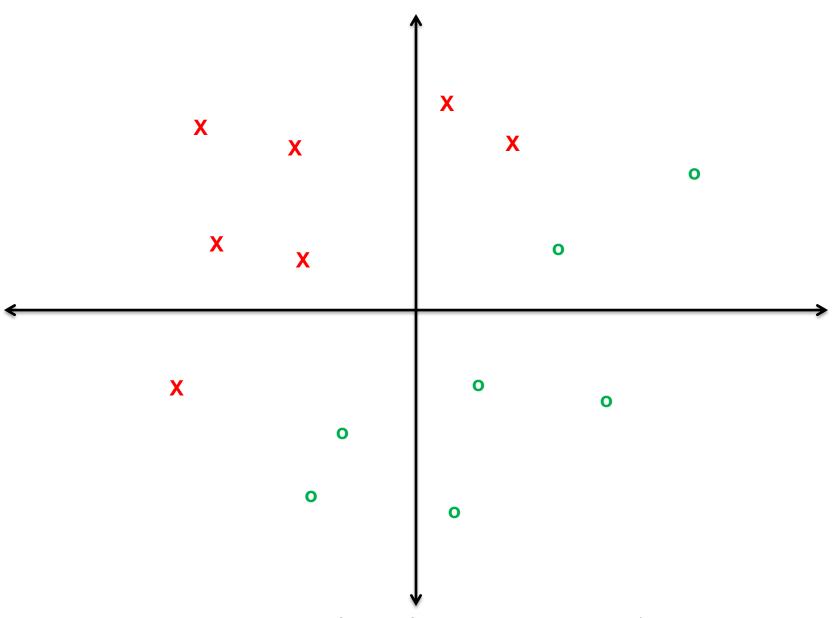
 $f(x_i; \theta)$  model parameters



Model

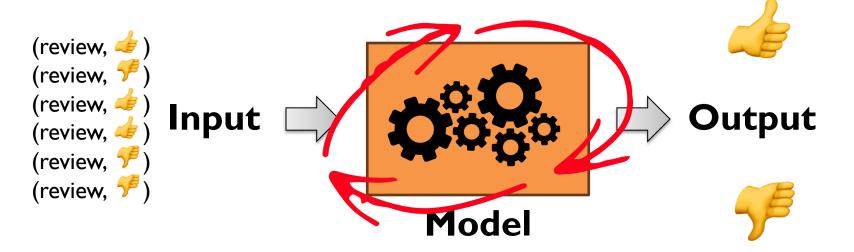
Machine learning algorithm adjusts the model parameters

What does that actually mean?

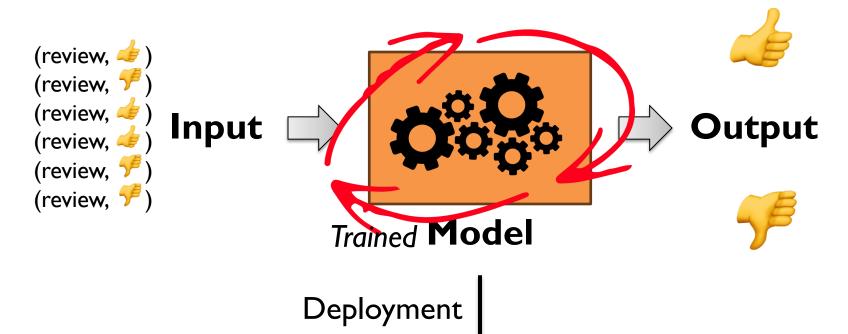


Parametric form for a linear model?

#### Model learns from the data



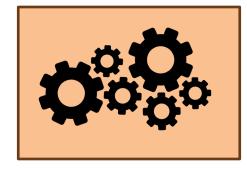
Machine learning algorithm adjusts the model parameters



#### Inference / Prediction

A group of us stopped by yesterday afternoon to enjoy an outdoor lunch. The food was da bomb.









Trained Model



What's supervised machine learning?

Gather training data

Train model

Deploy model

That's it?

#### **Caveats**

"traditional" machine learning (i.e., not LLMs) focus on data infrastructure (i.e., not ML/AI)

## Supervised Machine Learning

The general 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!

## Supervised Binary Classification

Restrict output label to be binary
Yes/No
1/0

Binary classifiers form primitive building blocks for multi-class problems...

## Binary Classifiers as Building Blocks

Example: four-way classification

One vs. rest classifiers

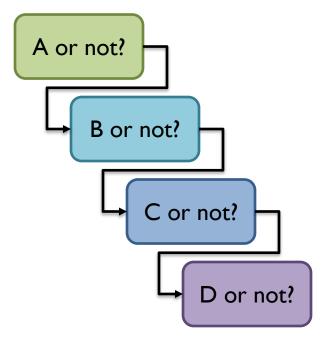
A or not?

B or not?

C or not?

D or not?

Classifier cascades



## Components of an ML Solution

Data

**Features** 

Model

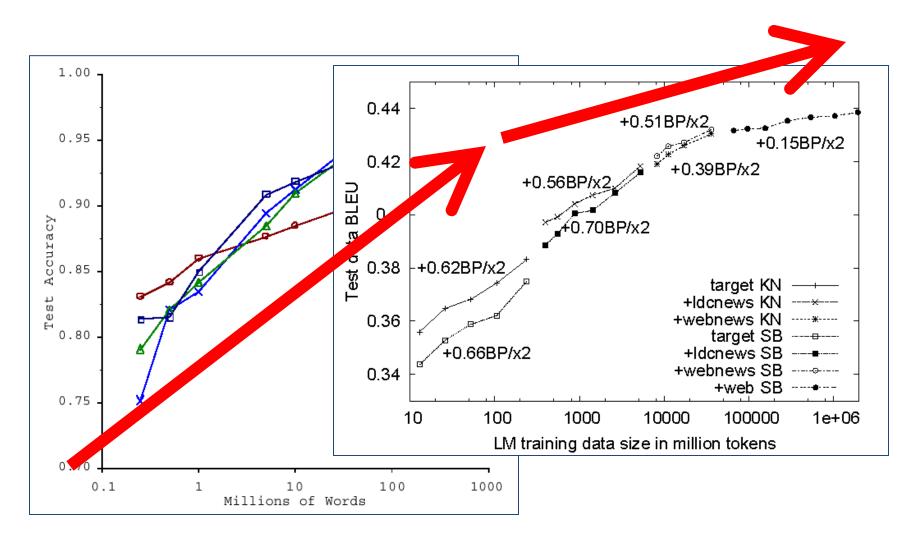
logistic regression, naïve Bayes, SVM, random forests,

neural networks, etc.

gradient descent, stochastic gradient descent, L-BFGS, etc. **Optimization** 

What "matters" the most?

#### No data like more data!



## Components of a ML Solution

Data
Features
Model
Optimization

What's the focus of data engineering?

### Components of a ML Solution

Data
Features
Model
Optimization

What's the focus of data engineering?

#### Instance

amazing spot for good food & a fun time they offer a super unique dine-in experience with their interactive tables! also love that they have innovative weekly feature dishes







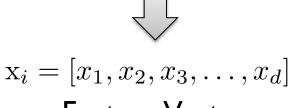
#### **Prediction**



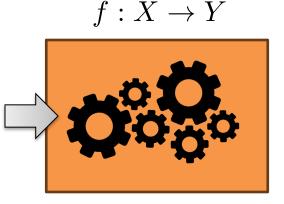


#### Instance

amazing spot for good food & a fun time they offer a super unique dine-in experience with their interactive tables! also love that they have innovative weekly feature dishes



Feature Vector



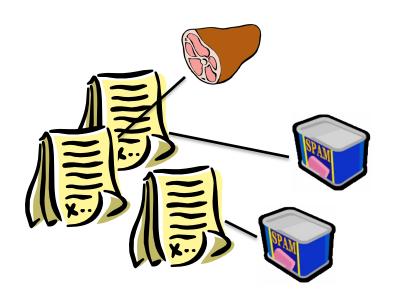
**Prediction** 



Model

Examples of features?
Who comes up with the features? How?

## Features: Spam Classification



#### Instances (emails) 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.

## Limits of Supervised Classification?

Isn't gathering labels a bottleneck?

"Found" data
User behavior logs
Crowdsourcing
Semi-supervised techniques
Self-supervised techniques

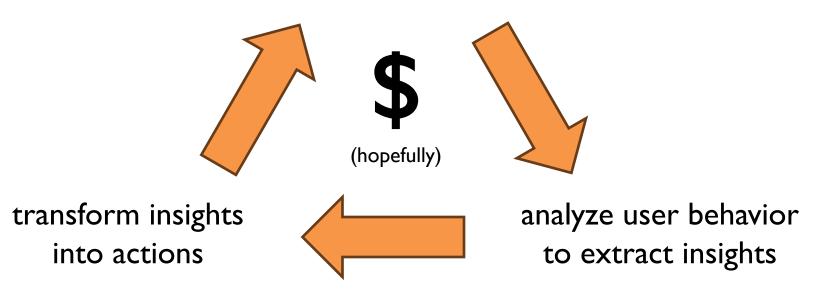
The virtuous cycle of data-driven products

#### This...

## The Data Flywheel

(a virtuous cycle)

Build a useful product



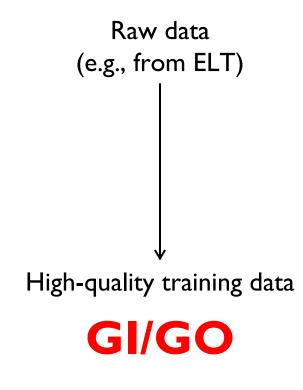
Google. Facebook. Twitter. Amazon. Uber.

Gather training data

Train model

Deploy model

"Data infrastructure for ML": What does that mean?

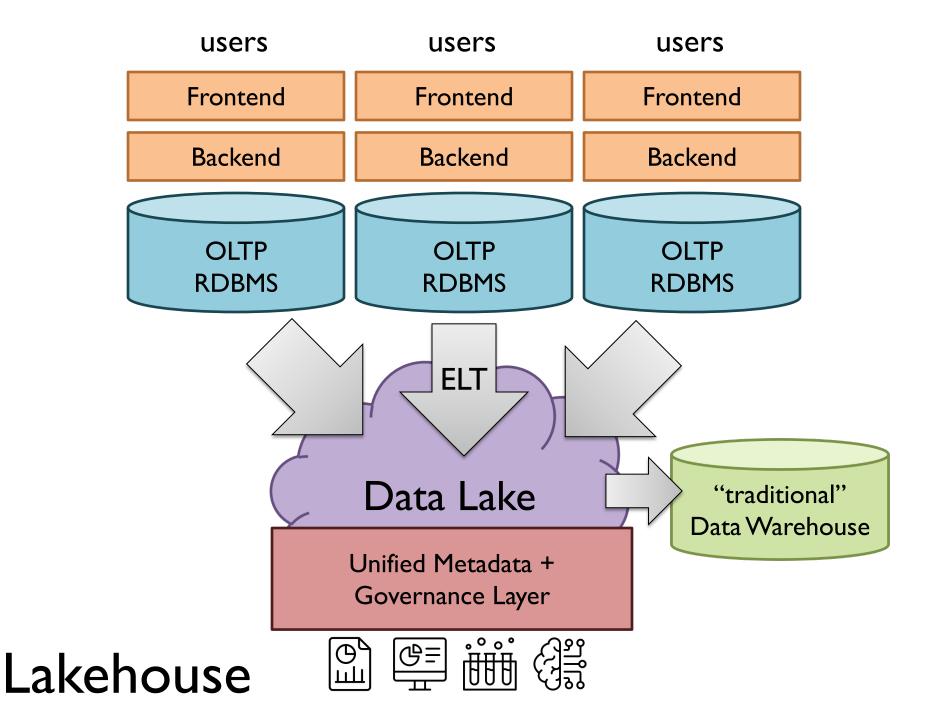


Gather training data

Train model

Deploy model

"Data infrastructure for ML": What does that mean?

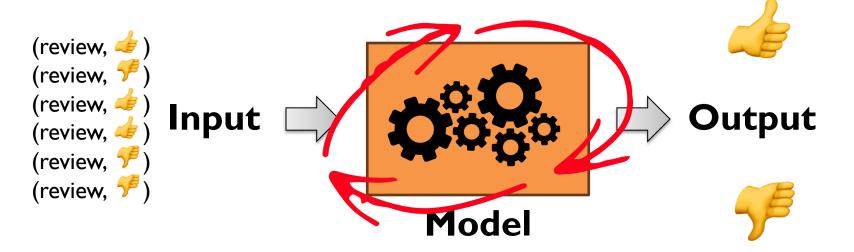


Gather training data

Train model

Deploy model

#### Model learns from the data



Machine learning algorithm adjusts the model parameters

## The Task

Given: 
$$D = \{(\mathbf{x}_i, y_i)\}_i^n$$
 feature vector  $\mathbf{x}_i = [x_1, x_2, x_3, \dots, x_d]$   $y \in \{0, 1\}$ 

Induce:  $f: X \to Y$ 

Such that loss is minimized

$$\frac{1}{n} \sum_{i=0}^{n} \ell(f(\mathbf{x}_i), y_i)$$

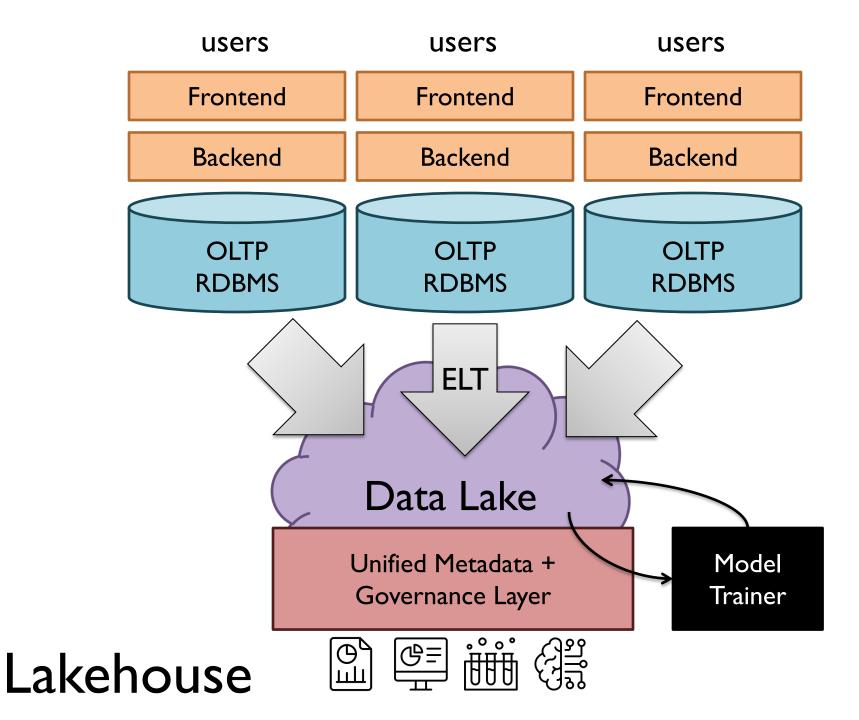
loss function

Typically, we 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

### How do we do it?

Use sklearn: model.fit(X, y)



How do we do it?

Use sklearn: model.fit(X, y)

For many organizations, that (really) is the end of the story! (but let's open up the black box...)

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

## Gradient Descent: Preliminaries

#### Rewrite:

$$\operatorname{arg\,min}_{\theta} \frac{1}{n} \sum_{i=0}^{n} \ell(f(\mathbf{x}_i; \theta), y_i) \qquad \qquad \operatorname{arg\,min}_{\theta} L(\theta)$$

### 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: \*

$$b = a - \gamma \nabla L(a)$$

$$L(a) > L(b)$$

# Gradient Descent: Iterative Update

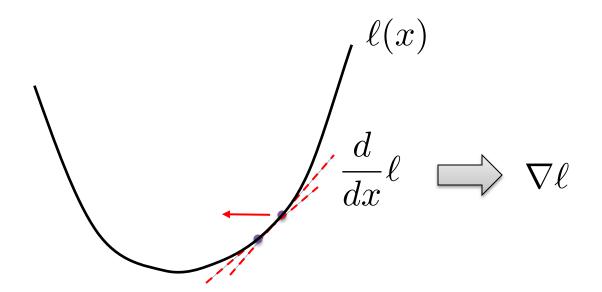
Start at an arbitrary point, iteratively update:

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

We have:

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

## 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

# 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:

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

#### Lots of details:

Figuring out the step size

Getting stuck in local minima

Convergence rate

. . .

### Gradient Descent

Repeat until convergence:

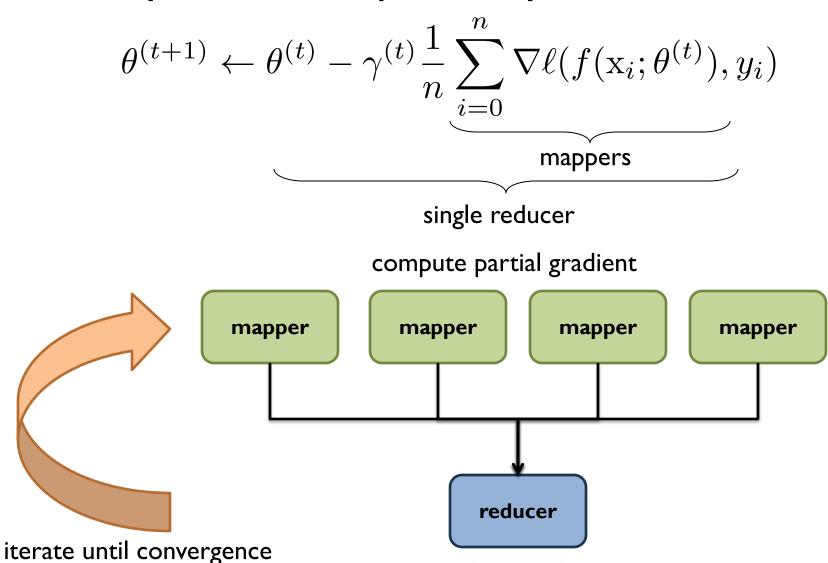
$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \frac{1}{n} \sum_{i=0}^{n} \nabla \ell(f(\mathbf{x}_i; \theta^{(t)}), y_i)$$

Note, sometimes formulated as ascent but entirely equivalent



(Note, this isn't really done anymore, but here's how you would do it...)

# MapReduce / Spark Implementation



update model



# Logistic Regression: Preliminaries

Given: 
$$D = \{(\mathbf{x}_i, y_i)\}_i^n$$
  
 $\mathbf{x}_i = [x_1, x_2, x_3, \dots, x_d]$   
 $y \in \{0, 1\}$ 

**Define:** 
$$f(\mathbf{x}; \mathbf{w}) : \mathbb{R}^d \to \{0, 1\}$$
 
$$f(\mathbf{x}; \mathbf{w}) = \begin{cases} 1 & \text{if } \mathbf{w} \cdot \mathbf{x} \ge t \\ 0 & \text{if } \mathbf{w} \cdot \mathbf{x} < t \end{cases}$$

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

# Relation to the Logistic Function

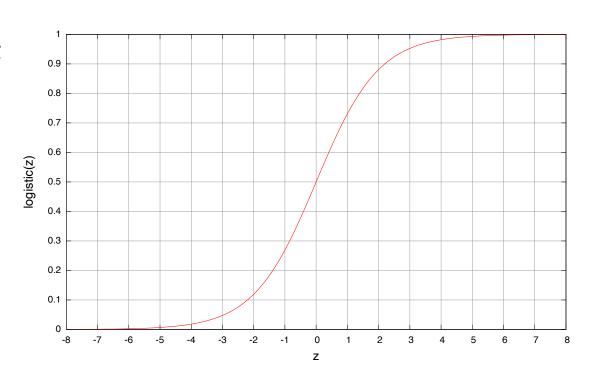
### After some algebra:

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

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

### The logistic function:

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



# Training an LR Classifier

Maximize the conditional likelihood:

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

Define the objective in terms of conditional *log* likelihood:

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

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

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

Substituting:

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

# 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 of 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:

$$\mathbf{w}_{i}^{(t+1)} \leftarrow \mathbf{w}_{i}^{(t)} + \gamma^{(t)} \sum_{j=0}^{n} x_{j,i} \Big( y_{j} - \Pr(y_{j} = 1 | \mathbf{x}_{j}, \mathbf{w}^{(t)}) \Big)$$

### Lots more details...

Want more?

Take a real machine-learning course!

#### Final update rule:

$$\mathbf{w}_{i}^{(t+1)} \leftarrow \mathbf{w}_{i}^{(t)} + \gamma^{(t)} \sum_{j=0}^{n} x_{j,i} \Big( y_{j} - \Pr(y_{j} = 1 | \mathbf{x}_{j}, \mathbf{w}^{(t)}) \Big)$$



## Batch vs. Online

(Batch) Gradient Descent

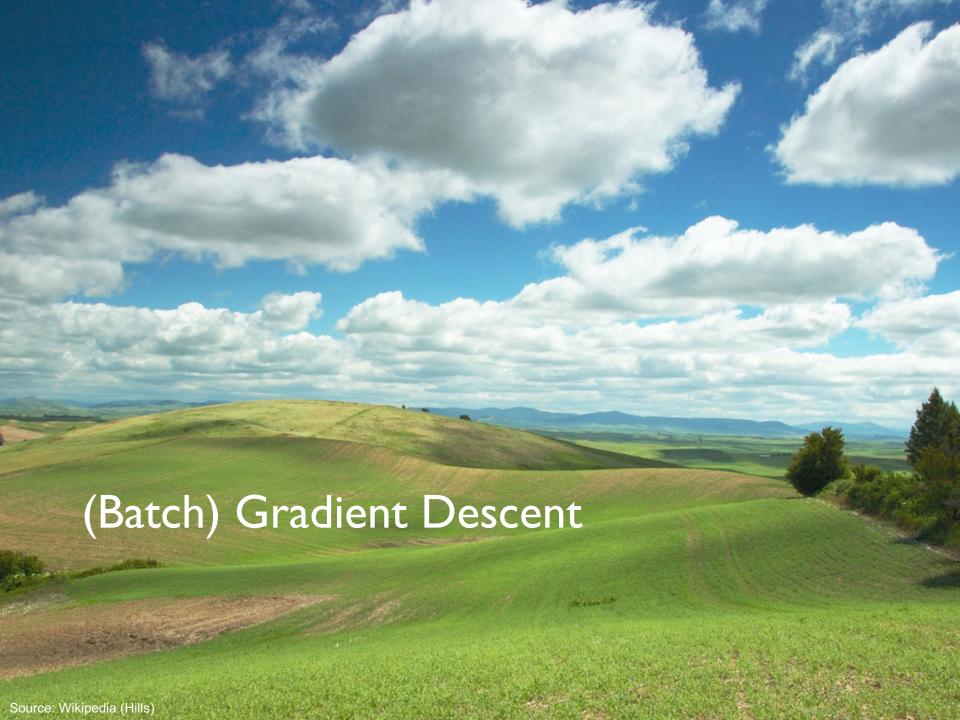
$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \frac{1}{n} \sum_{i=0}^{n} \nabla \ell(f(\mathbf{x}_i; \theta^{(t)}), y_i)$$

"batch" learning: update model after considering <u>all</u> training instances

### Stochastic Gradient Descent (SGD)

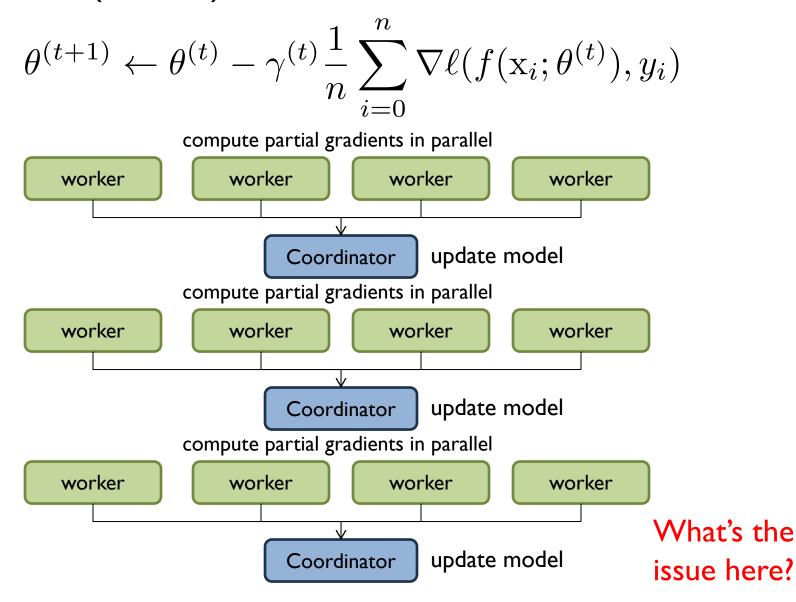
$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \nabla \ell(f(\mathbf{x}; \theta^{(t)}), y)$$

"online" learning: update model after considering <u>each</u> (randomly selected) training instance





# (Batch) Gradient Descent



## Stochastic Gradient Descent

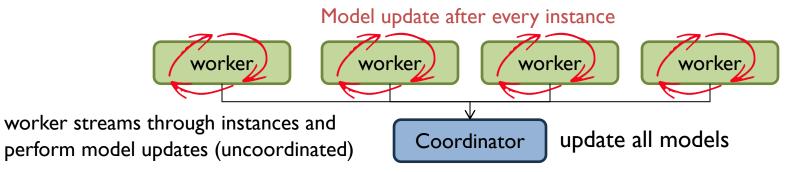
$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \nabla \ell(f(\mathbf{x}; \theta^{(t)}), y)$$

single worker streams through instances and performs model updates



Important: Model update after every instance

### How do you parallelize?



## Stochastic Gradient Descent w/ Mini-Batches

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

single worker streams through instances and performs model updates



Important: Model update after every instance Problem: updates are very noisy

#### Solution: mini-batches

Divide dataset into small batches (e.g., 64)
Perform gradient descent on each mini-batch
Update model after processing each mini-batch

