

# Data-Intensive Distributed Computing

CS 431/631 451/651 (Fall 2021)

## Part 7: Data Mining (1/4)

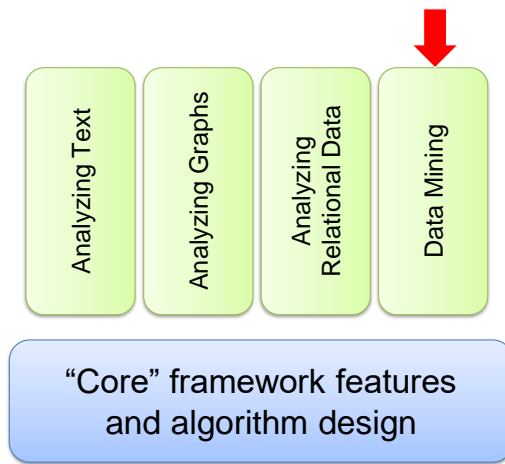
Ali Abedi

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

## Structure of the Course



# 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!

3

# Classification



Source: Wikipedia (Sorting)

# Applications

Spam detection

Sentiment analysis

Content (e.g., topic) classification

Link prediction

Document ranking

Object recognition

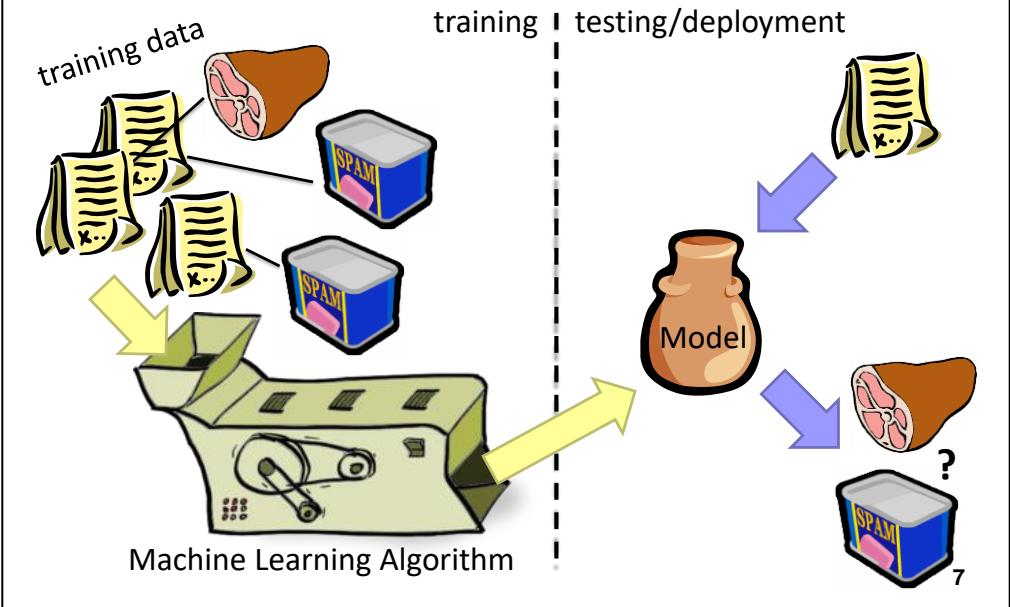
Fraud detection

And much much more!

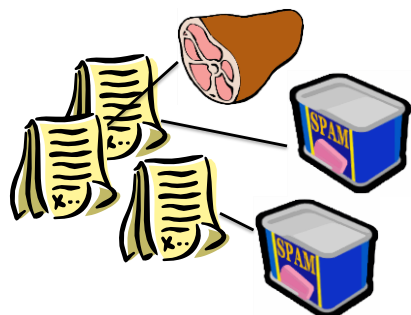
## Applications



# 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

Data  
Features  
Model  
Optimization

*logistic regression, naive Bayes, SVM, random forests, perceptrons, neural networks, etc.*

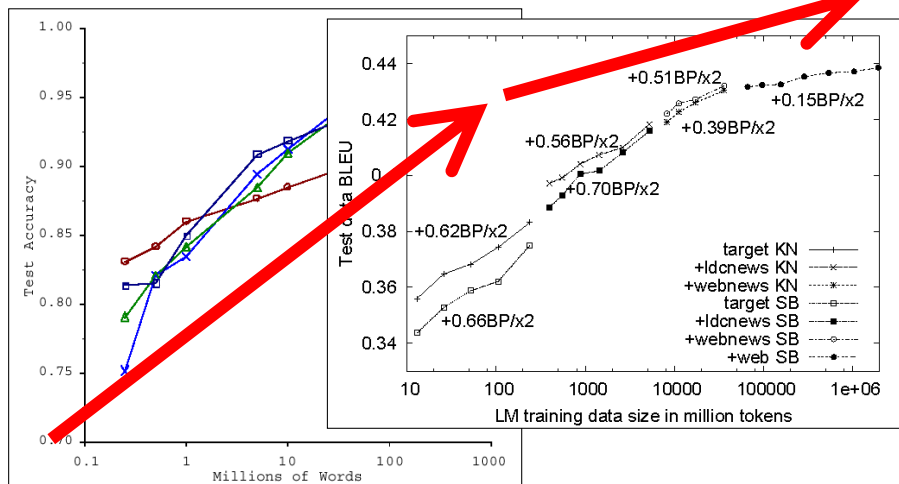
*gradient descent, stochastic gradient descent, L-BFGS, etc.*

What “matters” the most?

10

Data matters the most. If we throw a lot of data at almost any algorithm it performs good.

# No data like more data!



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

11

We have seen these examples before. For example stupid backoff outperforms other algorithms when it's trained on a lot of data.

# Limits of Supervised Classification?

Why is this a big data problem?

Isn't gathering labels a serious bottleneck?

Solutions

Crowdsourcing

Bootstrapping, semi-supervised techniques

Exploiting user behavior logs



12

## Supervised *Binary* Classification

Restrict output label to be *binary*

Yes/No

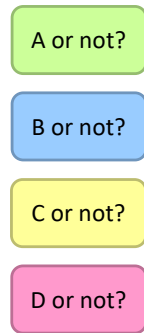
1/0

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

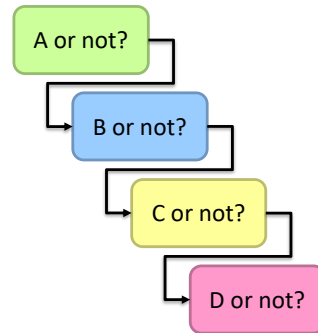
# Binary Classifiers as Building Blocks

Example: four-way classification

One vs. rest classifiers





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

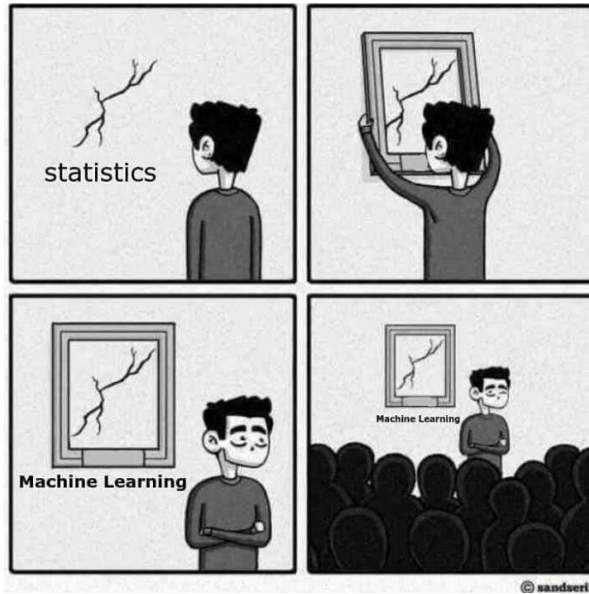
Typically, we consider functions of a parametric form:

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

 model parameters

15

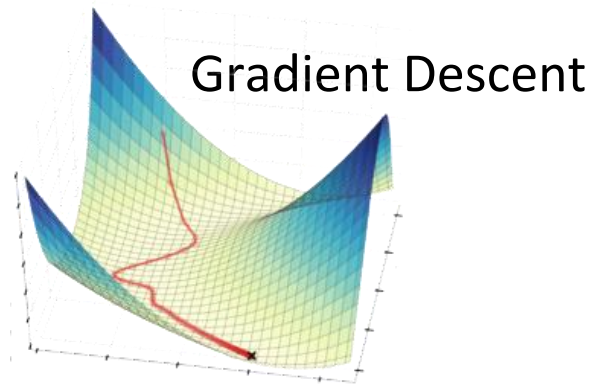
Key insight: machine learning as an optimization problem!





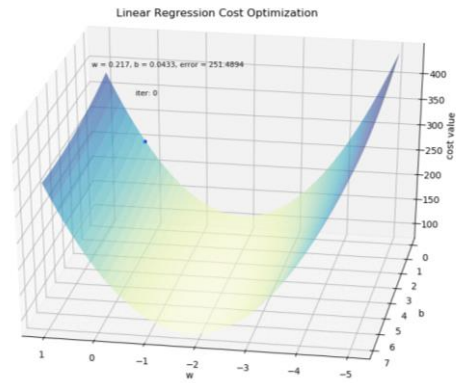
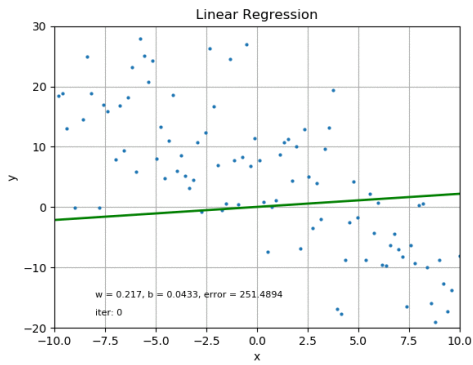
$$\arg \min_{\theta} \frac{1}{n} \sum_{i=0}^n \ell(f(x_i; \theta), y_i)$$

Because the closed form solutions  
generally not possible ...





# Linear Regression + Gradient Descent



## Gradient Descent: Preliminaries

Rewrite:

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

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:

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

$$L(a) \geq L(b)$$

## 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:

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

# 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:

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

Note, sometimes formulated as *ascent* but entirely equivalent

23

## Even More Details...

Gradient descent is a “first order” optimization technique

Often, slow convergence

Newton and quasi-Newton methods:

Intuition: Taylor expansion

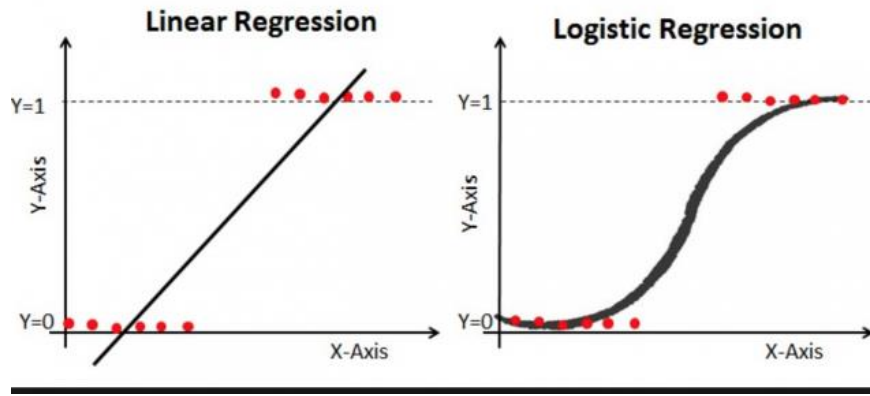
$$f(x + \Delta x) = f(x) + f'(x)\Delta x + \frac{1}{2}f''(x)\Delta x^2$$

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





## Logistic Regression: Preliminaries



## Logistic Regression: Preliminaries

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

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

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

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

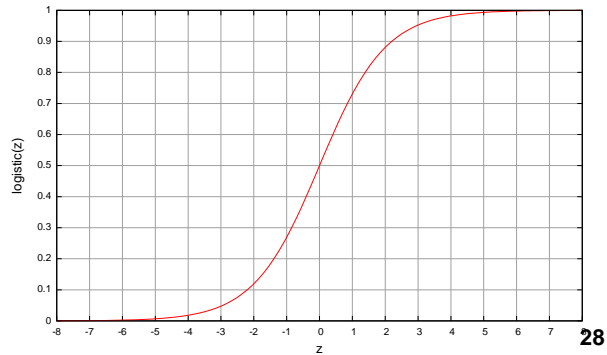
After some algebra:

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

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

The logistic function:

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



28

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

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:

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

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:

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

30

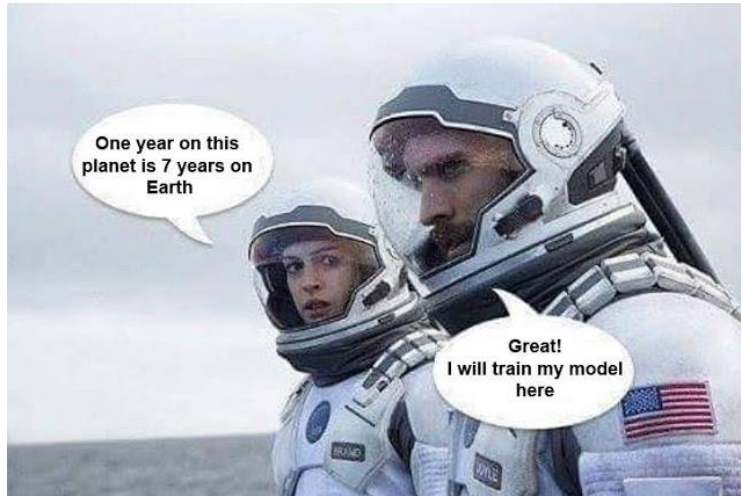
Lots more details...

Regularization  
Different loss functions

...

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

31



$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \gamma^{(t)} \sum_{j=0}^n x_{j,i} (y_j - \Pr(y_j = 1 | x_j, w^{(t)}))$$

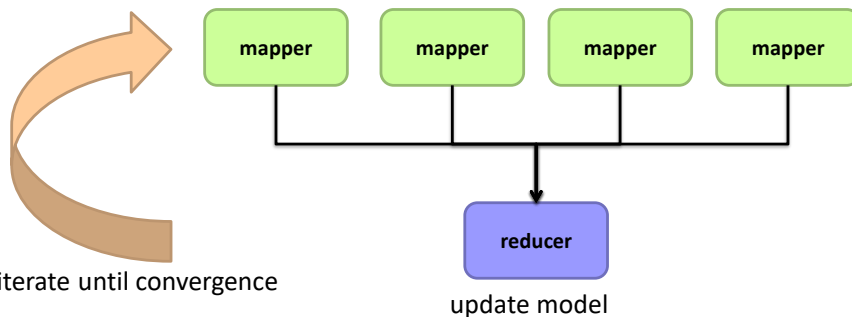


# MapReduce Implementation

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

mappers  
single reducer

compute partial gradient



# 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

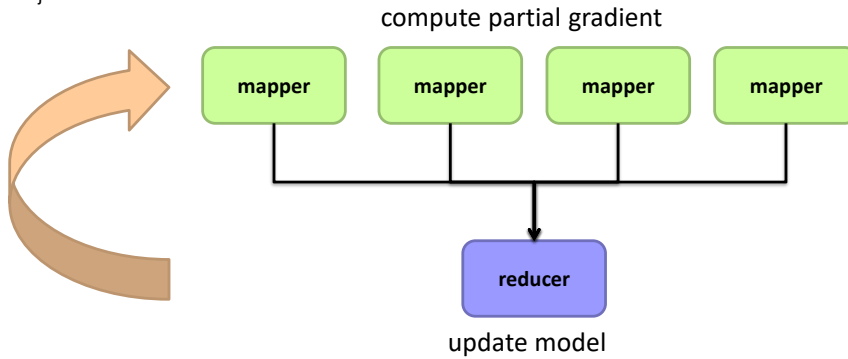
E.g., L-BFGS: faster convergence, but more to compute

# 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?*



35