



#### 10-301/601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

# Significance Testing + Exam 3 Review

Matt Gormley Lecture 27 Dec. 7, 2022

#### Reminders

- Homework 9: Learning Paradigms
  - Out: Fri, Dec 2
  - Due: Fri, Dec 9 at 11:59pm
     (only two grace/late days permitted)
- Exam 3 Practice Problems
  - Out: Wed, Dec 7
- Exam 3
  - Thu, Dec 15 (9:30am 11:30am)

## Crowdsourcing Exam Questions

#### **In-Class Exercise**

- Select one of lecture-level learning objectives
- 2. Write a question that assesses that objective
- Adjust to avoid 'trivia style' question

#### **Answer Here:**

#### SIGNIFICANCE TESTING

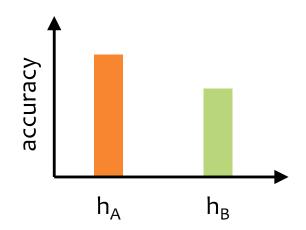
#### Which classifier is better?

**Goal**: Given two classifiers:  $h_A(x)$  and  $h_B(x)$  which is better?

h<sub>B</sub>(x)

 $h_A(x)$ 

Common Approach: Evaluate each classifier on a test set and report which has higher accuracy.



#### Two Sources of Variance

- 1. Randomness in training
- 2. Randomness in our test data

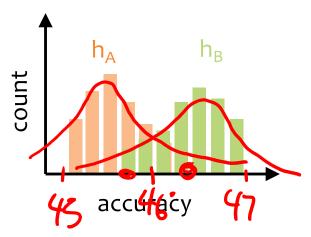
## 1. Randomness in training

Example: Assume we are training a deep neural network with a nonconvex objective function via random restarts

We collect a sequence of classifiers for R random restarts:

- $h_B(x)^{(1)} \leftarrow train(D, seed = time in ms)$
- $h_B(x)^{(2)} \leftarrow train(D, seed = time in ms)$
- **\*** ...
- $h_B(x)^{(R)} \leftarrow train(D, seed = time in ms)$

#### Solution: histogram



#### Solution: confidence interval

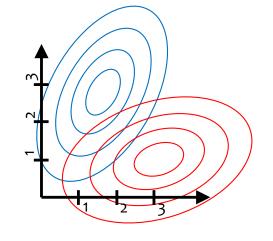
report variance of h<sub>A</sub> and h<sub>B</sub>

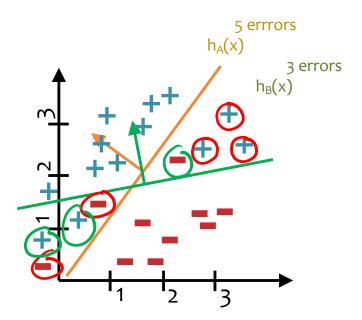
- Ex:
- h<sub>A</sub> 45% +/- 5%
- h<sub>B</sub> 47% +/- 8%

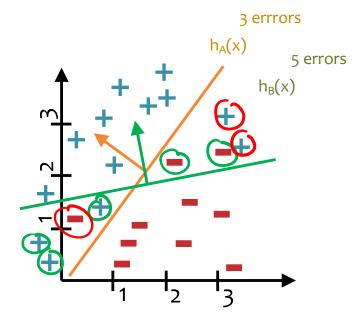
#### 2. Randomness in our test data

**Recall:** we assume  $\underline{x^{(i)}} \sim p^*(\cdot)$  and  $\underline{y^{(i)}} = c^*(\underline{x^{(i)}})$  or  $(x^{(i)}, y^{(i)}) \sim p^*(\cdot, \cdot)$ 

**Data:** Assume the data is drawn from a generative distribution p\*(x|y)p\*(y) where p\*(y) is an even coin flip and p\*(x|y=red) is the red Gaussian and p\*(x|y=blue) is the blue Gaussian.





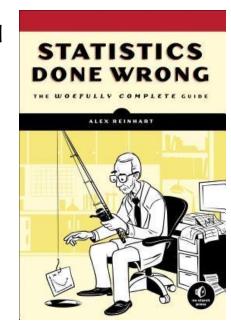


Solution: significance testing

## Significance Testing in ML

"And because any medication or intervention usually has some real effect, you can always get a statistically significant result by collecting so much data that you detect extremely tiny but relatively unimportant differences. As Bruce Thompson wrote, Statistical significance testing can involve a tautological logic in which tired researchers, having collected data on hundreds of subjects, then conduct a statistical test to evaluate whether there were a lot of subjects, which the researchers already know, because they collected the data and know they are tired. This tautology has created considerable damage as regards the cumulation of knowledge."

Alex Reinhart
 Statistics Done Wrong: The Woefully Complete Guide



For machine learning, significance testing is usually still answering an important question:

Did we evaluate our model on enough test data to conclude that our improvement over the baseline is surprising?

## Significance Testing in ML

#### Paired Bootstrap Test

**Key Idea:** simulate the resampling of many test sets **Algorithm:** 

- 1. Draw B bootstrap samples: where  $\mathbf{v} \leq \mathcal{N}'$   $S^{(b)} = \{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)})(\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots, (\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}$ with replacement from test data  $D_{\text{test}}$ ,  $|\mathcal{D}_{\text{test}}| = \mathcal{N}'$
- 2. Let v = 0
- 3. For b = 1,..., B if  $\delta(S^{(b)}) > 2\delta(D_{test})$ : V = V + 1  $\delta(D') = \text{difference in accuracy between } h_A \text{ and } h_B \text{ on } D'$
- 4. Return p-value as v/B

Ho = null hypothesis = performance of  $h_A$  and  $h_B$  is the same

#### **EXAM LOGISTICS**

#### Exam 3

- Time / Location
  - Time: Thu, Dec 15 at 8:30 9:30am 11:30am
  - Location & Seats: You have all been split across multiple rooms.
     Everyone has an assigned seat in one of these room.
  - Please watch Piazza carefully for announcements.
- Logistics
  - Covered material: Lectures 18 26
    - (only K-Means from Lecture 26)
  - Format of questions:
    - Multiple choice
    - True / False (with justification)
    - Derivations
    - Short answers
    - Interpreting figures
    - Implementing algorithms on paper
  - No electronic devices
  - You are allowed to **bring** one 8½ x 11 sheet of notes (front and back)

#### Exam 3

#### How to Prepare

- Attend (or watch) this exam review session
- Review practice problems
- Review homework problems
- Review the **poll questions** from each lecture
- Consider whether you have achieved the learning objectives for each lecture / section
- Write your cheat sheets

## Topics for Exam 1

- Foundations
  - Probability, Linear
     Algebra, Geometry,
     Calculus
  - Optimization
- Important Concepts
  - Overfitting
  - Experimental Design

- Classification
  - Decision Tree
  - KNN
  - Perceptron
- Regression
  - Linear Regression

### Topics for Exam 2

- Classification
  - Binary LogisticRegression
- Important Concepts
  - Stochastic GradientDescent
  - Regularization
  - Feature Engineering
- Feature Learning
  - Neural Networks
  - Basic NN Architectures
  - Backpropagation

- Learning Theory
  - PAC Learning
- Generative Models
  - Generative vs.
     Discriminative
  - MLE / MAP
  - Naïve Bayes

- Regression
  - Linear Regression

## Topics for Exam 3

- Graphical Models
  - HMMs
  - Learning and Inference
  - Bayesian Networks
- Reinforcement Learning
  - Value Iteration
  - Policy Iteration
  - Q-Learning
  - Deep Q-Learning

- Other Learning Paradigms
  - K-Means
  - PCA
  - Ensemble Methods
  - Recommender Systems

#### **MATERIAL COVERED ON EXAM 1**

### Supervised Binary Classification

- Step 1: training
  - Given: labeled training dataset

Goal: learn a classifier from the training dataset

Step 2: prediction

Given: unlabeled test date: learned classifier

- Goal: predict a label for e instance
- Step 3: evaluation
  - Given: predictions from : labeled test datas
  - Goal: compute the test e rate (i.e. error rate on th dataset)

Key question in Machine Learning:

How do we learn the classifier from data?

## Medical Diagnosis

#### **Interview Transcript**

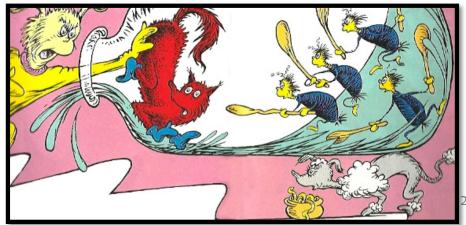
**Date**: Jan. 15, 2022

Parties: Matt Gormley and Doctor S.

**Topic:** Medical decision making

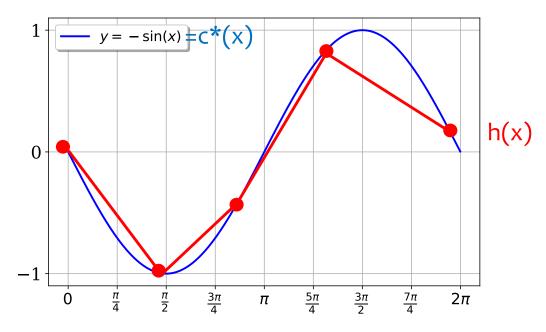
- Matt: Welcome. Thanks for interviewing with me today.
- Dr. S: Interviewing...?
- Matt: Yes. For the record, what type of doctor are you?
- Dr. S: Who said I'm a doctor?
- Matt: I thought when we set up this interview you said—
- Dr. S: I'm a preschooler.
- Matt: Good enough. Today, I'd like to learn how you would determine whether or not your little brother is allergic to cats given his symptoms.
- Dr. S: He's not allergic.
- Matt: We haven't started yet. Now, suppose he is sneezing. Does he have allergies to cats?
- Dr. S: Well, we don't even have a cat, so that doesn't make any sense.
- Matt: What if he is itchy; Does he have allergies?
- Dr. S: No, that's just a mosquito.
- [Editor's note: preschoolers unilaterally agree that itchiness is always caused by mosquitos, regardless of whether mosquitos were/are present.]

- Matt: What if he's both sneezing and itchy?
- Dr. S: Then he's allergic.
- Matt: Got it. What if your little brother is sneezing and itchy, plus he's a doctor.
- Dr. S: Then, thumbs down, he's not allergic.
- Matt: How do you know?
- Dr. S: Doctors don't get allergies.
- Matt: What if he is not sneezing, but is itchy, and he is a fox....
- Matt: ... and the fox is in the bottle where the tweetle beetles battle with their paddles in a puddle on a noodle-eating poodle.
- Dr. S: Then he is must be a tweetle beetle noodle poodle bottled paddled muddled duddled fuddled wuddled fox in socks, sir. That means he's definitely allergic.
- Matt: Got it. Can I use this conversation in my lecture?
- Dr. S: Yes



### **Function Approximation**

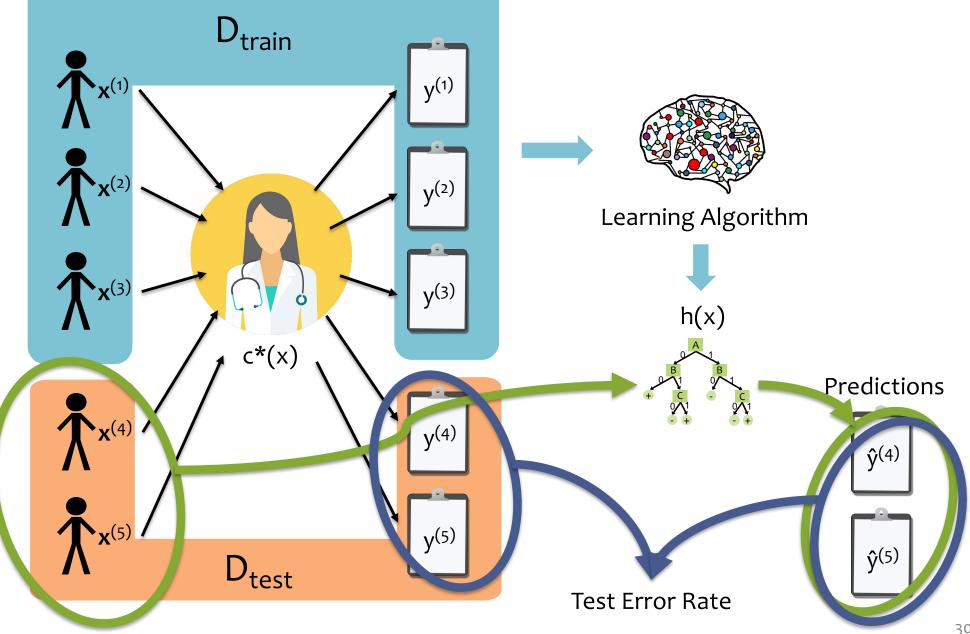
**Quiz:** Implement a simple function which returns  $-\sin(x)$ .



#### A few constraints are imposed:

- 1. You can't call any other trigonometric functions
- You can call an existing implementation of sin(x) a few times (e.g. 100) to test your solution
- 3. You only need to evaluate it for x in [0, 2\*pi]

## Supervised Machine Learning



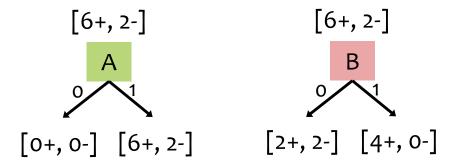


## Decision Tree Learning Example

#### **Dataset:**

Output Y, Attributes A and B

Y	Α	В
-	1	0
-	1	0
+	1	0
+	1	0
+	1	1
+	1	1
+	1	1
+	1	1



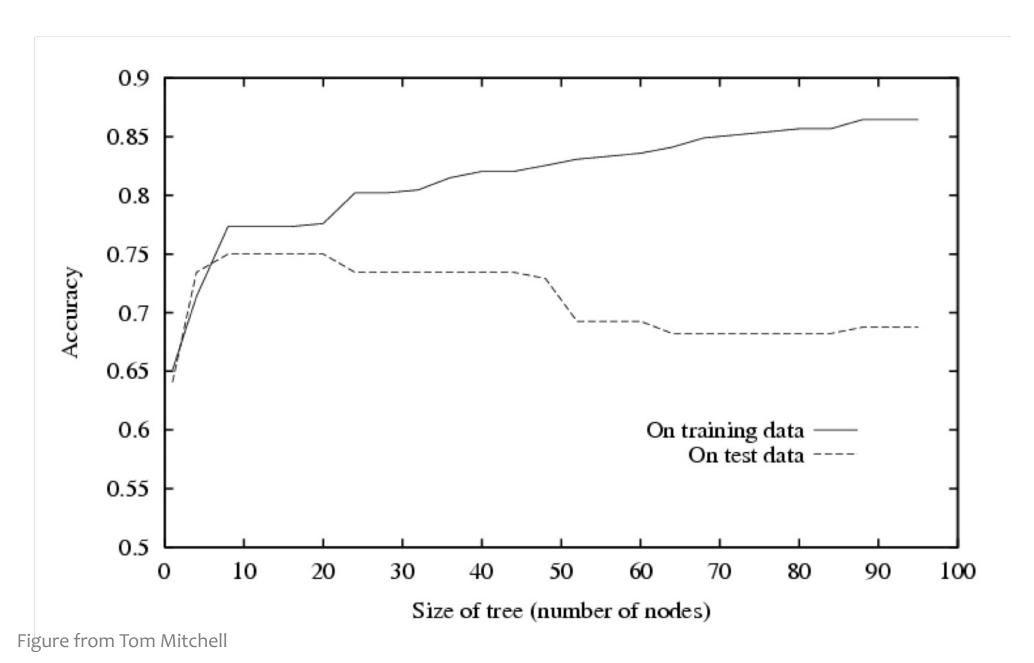
#### **Mutual Information**

$$H(Y) = -2/8 \log(2/8) - 6/8 \log(6/8)$$

$$H(Y|A=0) =$$
 "undefined"  
 $H(Y|A=1) = -2/8 \log(2/8) - 6/8 \log(6/8)$   
 $= H(Y)$   
 $H(Y|A) = P(A=0)H(Y|A=0) + P(A=1)H(Y|A=1)$   
 $= 0 + H(Y|A=1) = H(Y)$   
 $I(Y; A) = H(Y) - H(Y|A=1) = 0$ 

$$H(Y|B=0) = -2/4 \log(2/4) - 2/4 \log(2/4)$$
  
 $H(Y|B=1) = -0 \log(0) - 1 \log(1) = 0$   
 $H(Y|B) = 4/8(0) + 4/8(H(Y|B=0))$   
 $I(Y;B) = H(Y) - 4/8 H(Y|B=0) > 0$ 

## Overfitting in Decision Tree Learning





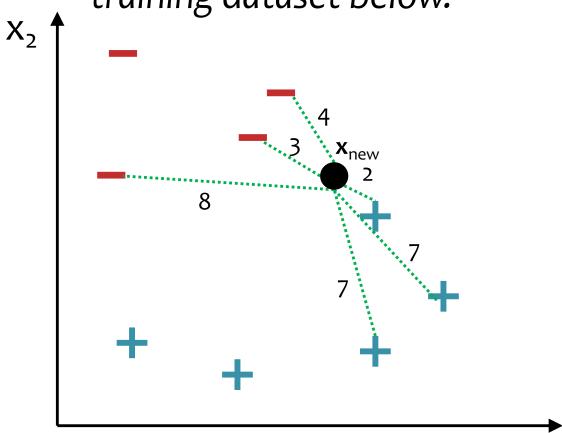
Species	Sepal Length	Sepal Width	Petal Length	Petal Width
0	4.3	3.0	1.1	0.1
0	4.9	3.6	1.4	0.1
0	5.3	3.7	1.5	0.2
1	4.9	2.4	3.3	1.0
1	5.7	2.8	4.1	1.3
1	6.3	3.3	4.7	1.6
1	6.7	3.0	5.0	1.7





## k-Nearest Neighbors

Suppose we have the training dataset below.



How should we label the new point?

It depends on k:

if 
$$k=1$$
,  $h(x_{new}) = +1$ 

if 
$$k=3$$
,  $h(x_{new}) = -1$ 

if 
$$k=5$$
,  $h(x_{new}) = +1$ 







### Hyperparameter Optimization

#### **Question:**

True or False: given a finite amount of computation time, grid search is more likely to find good values for hyperparameters than random search.

## Answer: Grid Layout Random Layout Random Layout Random Layout

Important parameter

Figure 1: Grid and random search of nine trials for optimizing a function  $f(x,y) = g(x) + h(y) \approx g(x)$  with low effective dimensionality. Above each square g(x) is shown in green, and left of each square h(y) is shown in yellow. With grid search, nine trials only test g(x) in three distinct places. With random search, all nine trials explore distinct values of g. This failure of grid search is the rule rather than the exception in high dimensional hyper-parameter optimization.

Important parameter

## Linear Models for Classification

Key idea: Try to learn this hyperplane directly

#### Looking ahead:

- We'll see a number of commonly used Linear Classifiers
- These include:
  - Perceptron
  - Logistic Regression
  - Naïve Bayes (under certain conditions)
  - Support Vector
     Machines

Directly modeling the hyperplane would use a decision function:

$$h(\mathbf{x}) = \operatorname{sign}(\boldsymbol{\theta}^T \mathbf{x})$$

for:

$$y \in \{-1, +1\}$$

## Perceptron Mistake Bound

**Guarantee:** if some data has margin  $\gamma$  and all points lie inside a ball of radius R, then the online Perceptron algorithm makes  $\leq (R/\gamma)^2$  mistakes

(Normalized margin: multiplying all points by 100, or dividing all points by 100, doesn't change the number of mistakes! The algorithm is invariant to scaling.)



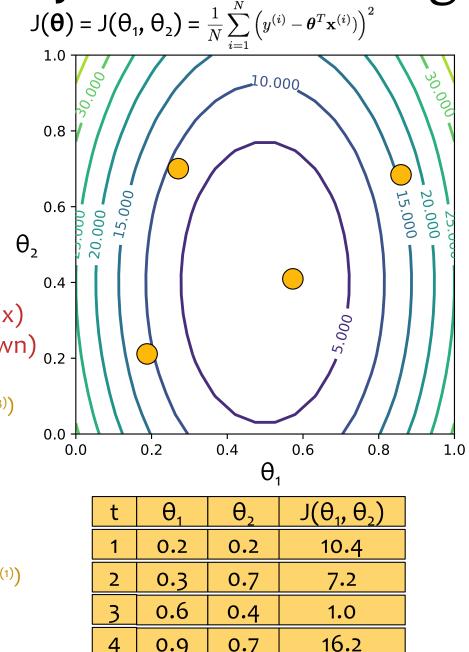
**Def:** We say that the (batch) perceptron algorithm has **converged** if it stops making mistakes on the training data (perfectly classifies the training data).

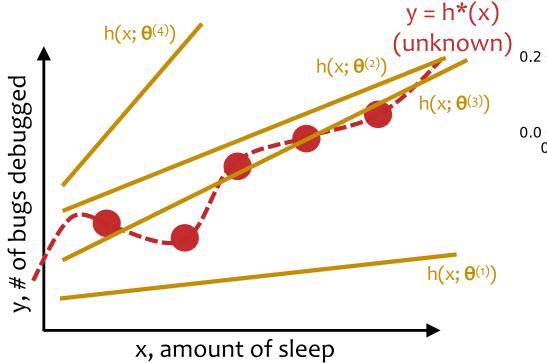
Main Takeaway: For linearly separable data, if the perceptron algorithm cycles repeatedly through the data, it will converge in a finite # of steps.

Linear Regression by Rand. Guessing

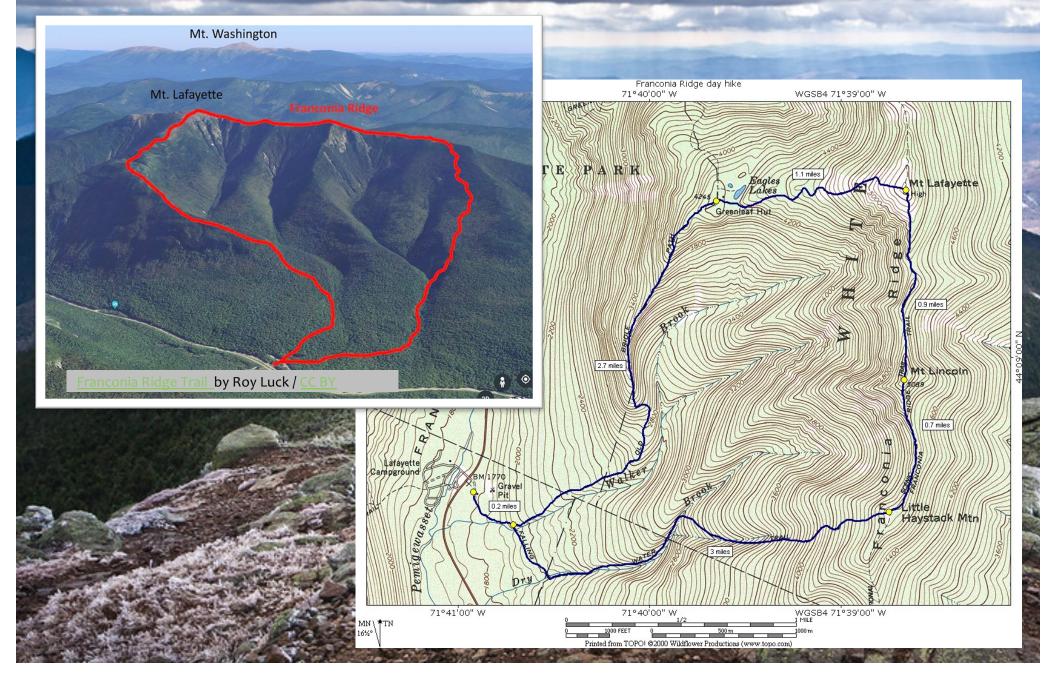
## Optimization Method #0: Random Guessing

- 1. Pick a random  $\theta$
- 2. Evaluate  $J(\theta)$
- 3. Repeat steps 1 and 2 many times
- 4. Return  $\theta$  that gives smallest  $J(\theta)$

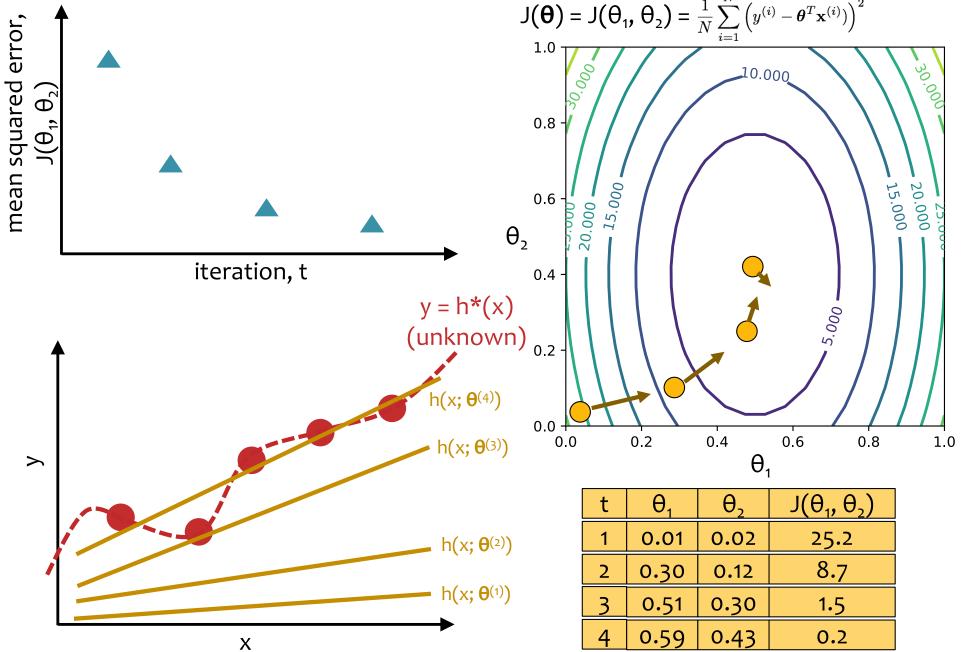




## Topographical Maps



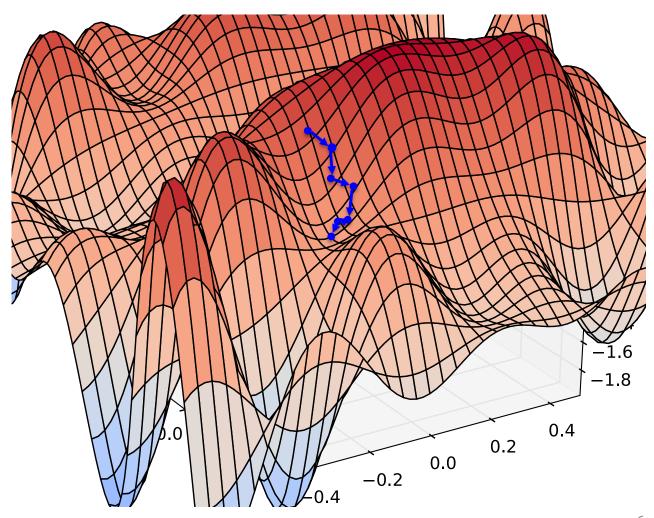
## Linear Regression by Gradient Desc. $J(\theta) = J(\theta_1, \theta_2) = \frac{1}{N} \sum_{i=1}^{N} (y^{(i)} - \theta^T \mathbf{x}^{(i)})^2$

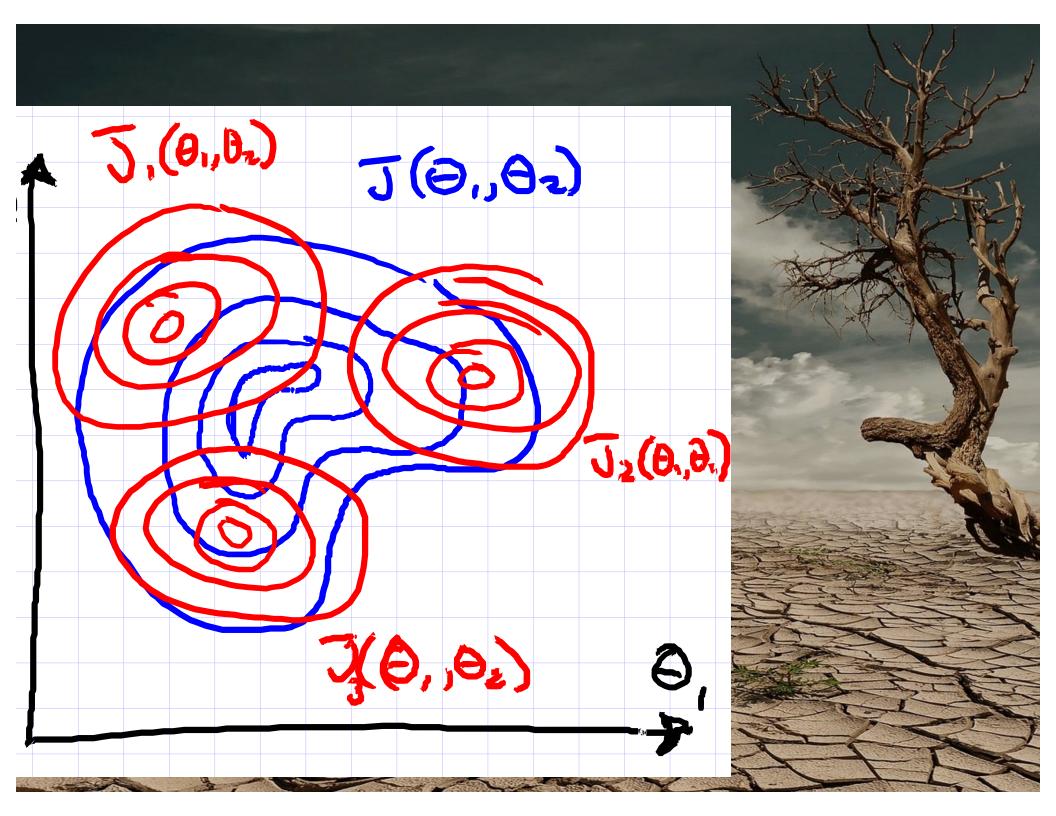


#### MATERIAL COVERED ON EXAM 2

### **Gradient Descent & Convexity**

- Gradient
   descent is a
   local
   optimization
   algorithm
- If the function is nonconvex, it will find a local minimum, not necessarily a global minimum
- If the function is convex, it will find a global minimum





# Probabilistic Learning

#### **Function Approximation**

Previously, we assumed that our output was generated using a deterministic target function:

$$\mathbf{x}^{(i)} \sim p^*(\cdot)$$

$$y^{(i)} = c^*(\mathbf{x}^{(i)})$$

Our goal was to learn a hypothesis h(x) that best approximates c\*(x)

#### **Probabilistic Learning**

Today, we assume that our output is **sampled** from a conditional **probability distribution**:

$$\mathbf{x}^{(i)} \sim p^*(\cdot)$$

$$y^{(i)} \sim p^*(\cdot|\mathbf{x}^{(i)})$$

Our goal is to learn a probability distribution p(y|x) that best approximates  $p^*(y|x)$ 

#### MLE

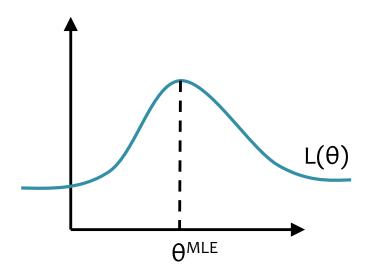
Suppose we have data  $\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$ 

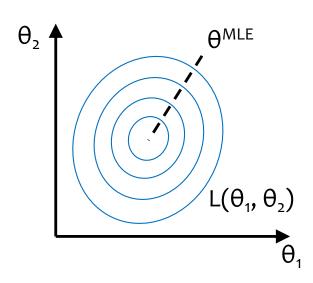
#### Principle of Maximum Likelihood Estimation:

Choose the parameters that maximize the likelihood of the data. N

$$\boldsymbol{\theta}^{\mathsf{MLE}} = \operatorname*{argmax}_{\boldsymbol{\theta}} \prod_{i=1} p(\mathbf{x}^{(i)} | \boldsymbol{\theta})$$

Maximum Likelihood Estimate (MLE)





# Logistic Regression

**Data:** Inputs are continuous vectors of length M. Outputs are discrete.

$$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N$$
 where  $\mathbf{x} \in \mathbb{R}^M$  and  $y \in \{0, 1\}$ 

**Model:** Logistic function applied to dot product of parameters with input vector.

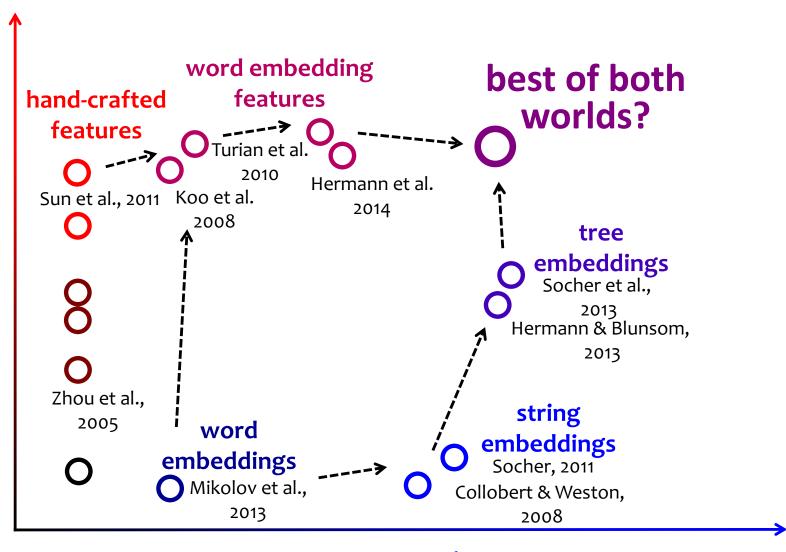
$$p_{\boldsymbol{\theta}}(y=1|\mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})}$$

**Learning:** finds the parameters that minimize some objective function.  ${m heta}^* = \mathop{\rm argmin}_{m heta} J({m heta})$ 

Prediction: Output is the most probable class.

$$\hat{y} = \operatorname*{argmax} p_{\boldsymbol{\theta}}(y|\mathbf{x})$$
$$y \in \{0,1\}$$

### Where do features come from?



Feature Learning

# Example: Linear Regression

**Goal:** Learn  $y = \mathbf{w}^T f(\mathbf{x}) + b$  where f(.) is a polynomial basis function

i	у	х	•••	<b>x</b> <sup>9</sup>	
1	2.0	1.2	•••	(1.2)9	
2	1.3	1.7	•••	(1.7)9	
	•••	•••	•••	•••	у
10	1.1	1.9	•••	(1.9)9	



X

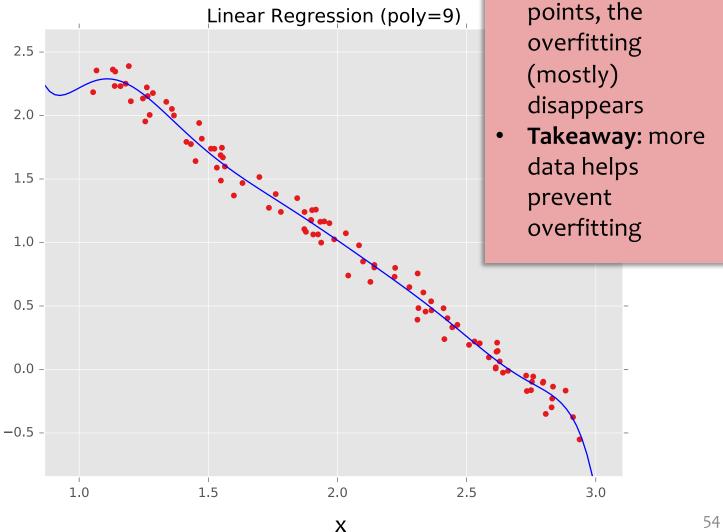
- But with N = 100
   points, the
   overfitting
   (mostly)
   disappears
- Takeaway: more data helps prevent overfitting

3.0

# Example: Linear Regression

**Goal:** Learn  $y = \mathbf{w}^T f(\mathbf{x}) + b$  where f(.) is a polynomial basis function

i	у	х	•••	<b>x</b> <sup>9</sup>	
1	2.0	1.2	•••	(1.2)9	
2	1.3	1.7	•••	(1.7)9	
3	0.1	2.7	•••	(2.7)9	)
4	1.1	1.9	•••	(1.9)9	
•••	•••	•••	•••	•••	
	•••	•••	•••		
•••	•••	•••	•••	•••	
98	•••	•••	•••	•••	
99	•••	•••	•••	•••	
100	0.9	1.5	•••	(1.5)9	



With just N = 10

But with N = 100

points we overfit!

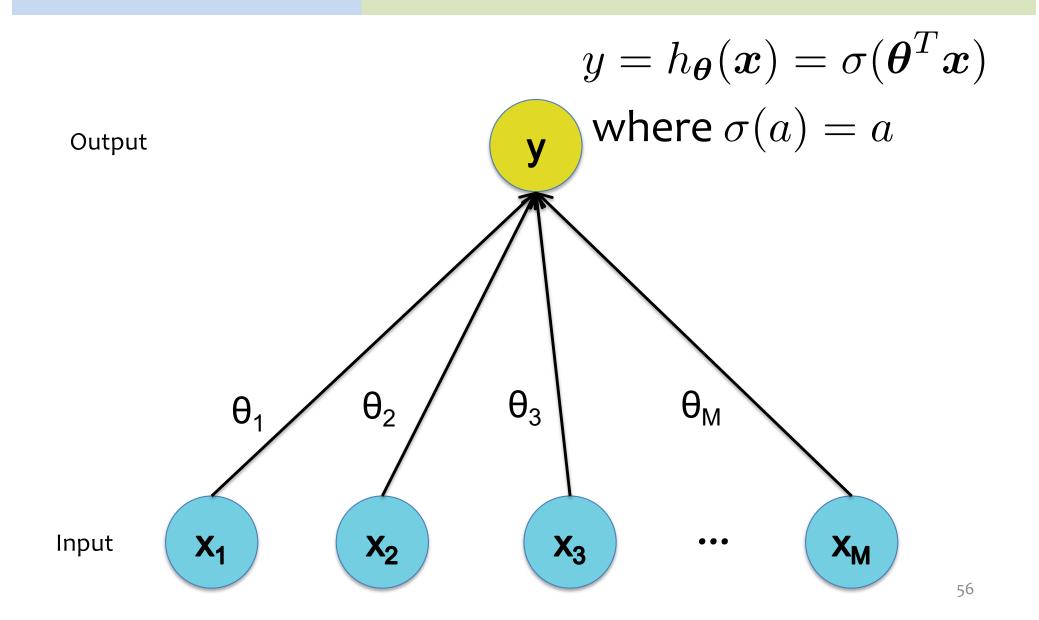
# Regularization

- **Given** objective function:  $J(\theta)$
- Goal is to find:  $\hat{\boldsymbol{\theta}} = \operatorname{argmin} J(\boldsymbol{\theta}) + \lambda r(\boldsymbol{\theta})$
- **Key idea:** Define regularizer  $r(\theta)$  s.t. we tradeoff between fitting the data and keeping the model simple
- Choose form of  $\mathbf{r}(\boldsymbol{\theta})$ :

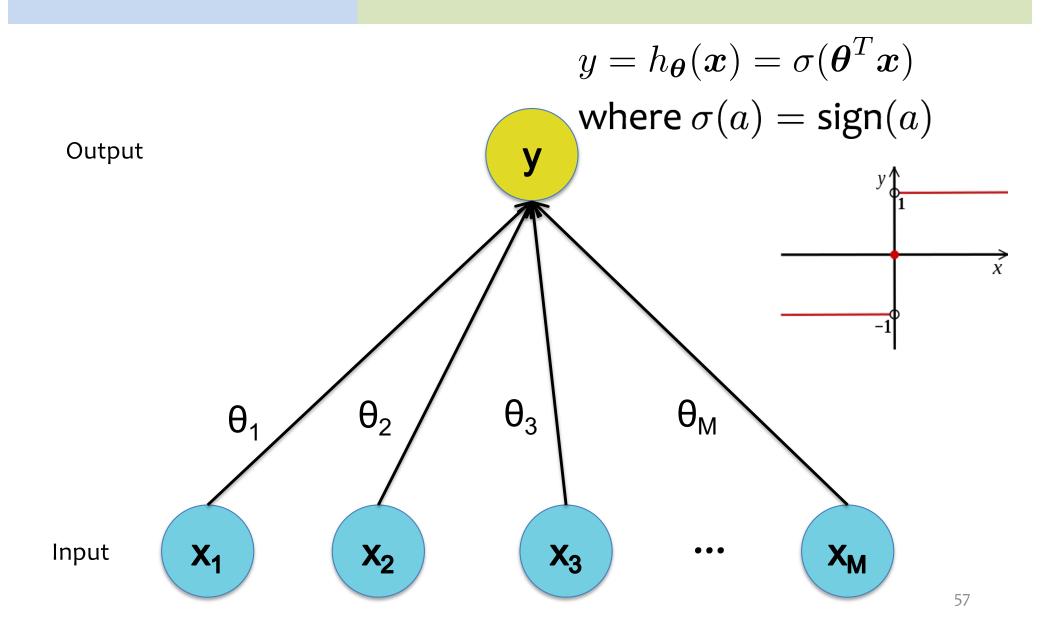
   Example: q-norm (usually p-norm):  $\|\boldsymbol{\theta}\|_q = \left(\sum_{m=1}^M |\theta_m|\right)^{\overline{q}}$

q	$r(oldsymbol{ heta})$	yields parame- ters that are	name	optimization notes
0	$  \boldsymbol{\theta}  _0 = \sum \mathbb{1}(\theta_m \neq 0)$	zero values	Lo reg.	no good computa- tional solutions
$\frac{1}{2}$	$  oldsymbol{ heta}  _1 = \sum   heta_m  \ (  oldsymbol{ heta}  _2)^2 = \sum  heta_m^2$	zero values small values	L1 reg. L2 reg.	subdifferentiable differentiable

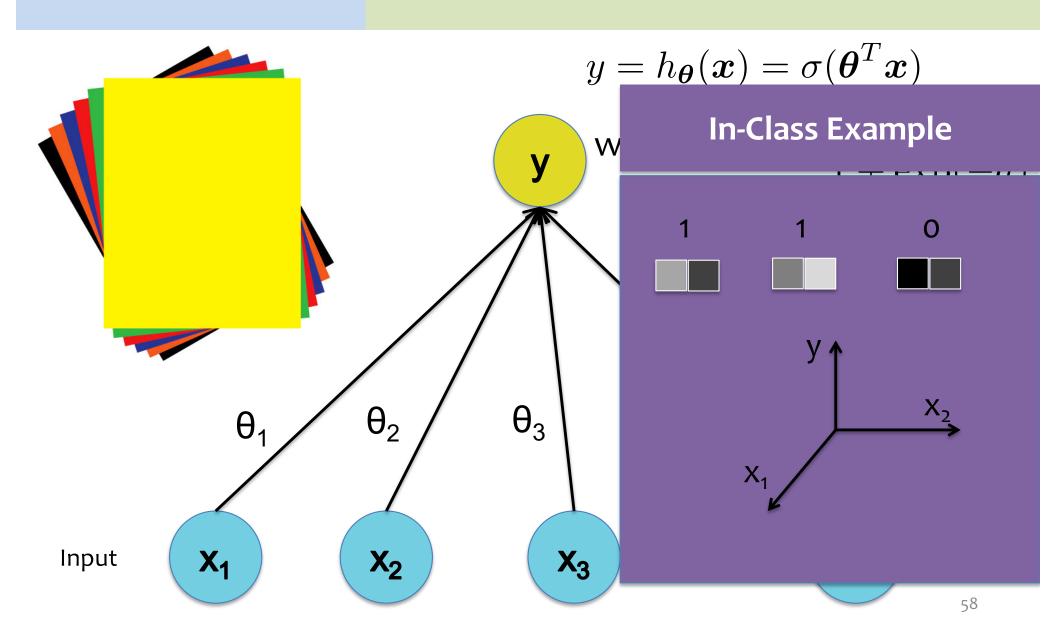
# Linear Regression



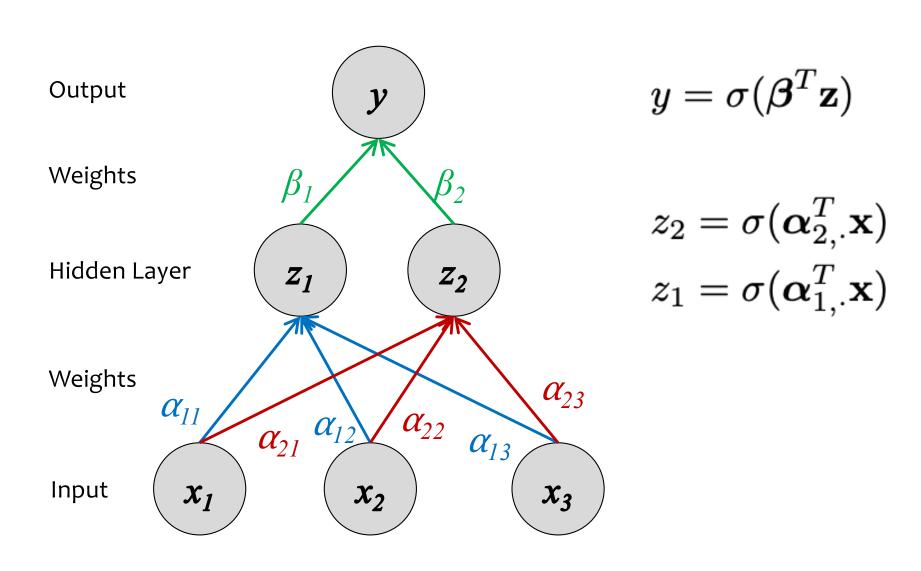
### Perceptron



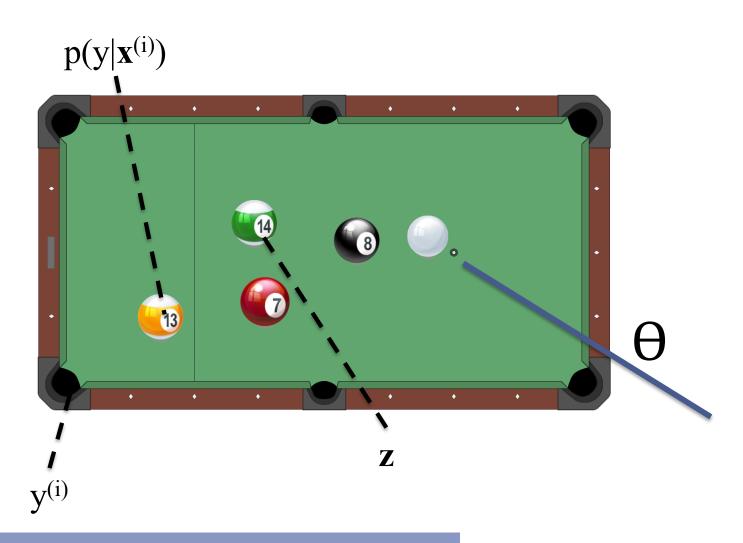
# Logistic Regression



### Neural Network



# **Error Back-Propagation**



### **Training**

# Differentiation Quiz

#### Differentiation Quiz #1:

Suppose x = 2 and z = 3, what are dy/dx and dy/dz for the function below? Round your answer to the nearest integer.

$$y = \exp(xz) + \frac{xz}{\log(x)} + \frac{\sin(\log(x))}{xz}$$

**Answer:** Answers below are in the fo

```
A. [42, -72]
```

```
rom math import *
# Define function
   return exp(x*z) + x*z/log(x) + sin(log(x)) / (x*z)
```

$$(=2; \mathbf{z} = 3; \mathbf{e} = 1e-8)$$

### Architecture #2: AlexNet

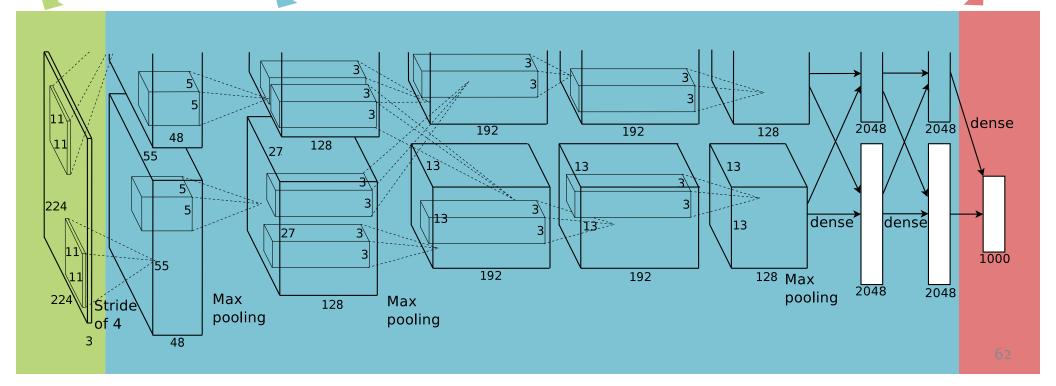
#### **CNN for Image Classification**

(Krizhevsky, Sutskever & Hinton, 2012) 15.3% error on ImageNet LSVRC-2012 contest

Input image (pixels)

- Five convolutional layers (w/max-pooling)
- Three fully connected layers

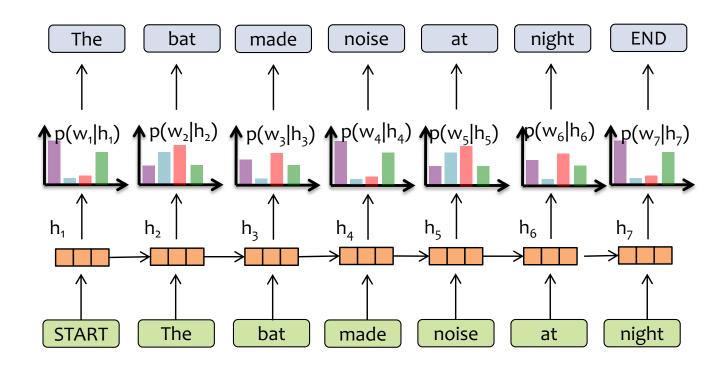
1000-way softmax





# ID-601 course staff

# RNN Language Model



#### **Key Idea:**

- (1) convert all previous words to a fixed length vector
- (2) define distribution  $p(w_t | f_{\theta}(w_{t-1}, ..., w_1))$  that conditions on the vector  $\mathbf{h}_t = f_{\theta}(w_{t-1}, ..., w_1)$

# Sampling from an RNN-LM

#### ??

VIOLA: Why, Salisbury must find his flesh and thought That which I am not aps, not a man and in fire, To show the reining of the raven and the wars To grace my hand reproach within, and not a fair are hand, That Caesar and my goodly father's world; When I was heaven of presence and our fleets, We spare with hours, but cut thy council I am great, Murdered a master's ready there My powe so much as hell: Some service i

bondman here, Would show hi

KING LEAR: O, if you we feeble sight, the courtesy of your law, Your sight and several breath, will wear the gods With his heads, and my hands are wonder'd at the deeds, So drop upon your lordship's head, and your opinion Shall be against your honour.

#### ??

my will.

CHARLES: Marry, do I, sir; and I came to acquaint you with a matter. I am given, sir, secretly to understand that your younger brother Orlando hath a disposition to come in disguised against me to try a fall. To-morrow, sir, I wrestle for my credit; and he that escapes me without some broken limb shall acquit him Which is the real is but young and tender; and, uld be loath to foil him, as I Shakespeare?! honour, if he come in: nx love to you, I came hither to acquaint you wi that either you might stay him from his int ent or brook such

TOUCHSTONE: For my part, I had rather bear with you than bear you; yet I should bear no cross if I did bear you, for I think you have no money in your purse.

disgrace well as he shadn into, in that it is a

thing of his own search and altogether against

# PAC-MAN Learning For some hypothesis $h \in \mathcal{H}$ :

1. True Error

2. Training Error

$$\hat{R}(h)$$

#### Question 2:

What is the expected number of PAC-MAN levels Matt will complete before a Game-

#### Over?

- Α. 1-10
- B. 11-20
- 21-30



# Sample Complexity Results

**Definition 0.1.** The **sample complexity** of a learning algorithm is the number of examples required to achieve arbitrarily small error (with respect to the optimal hypothesis) with high probability (i.e. close to 1).

#### Four Cases we care about...

Real	lizabl	le

Agnostic

Finite  $|\mathcal{H}|$ 

**Thm.** 1  $N \geq \frac{1}{\epsilon} \left[ \log(|\mathcal{H}|) + \log(\frac{1}{\delta}) \right]$  labeled examples are sufficient so that with probability  $(1-\delta)$  all  $h \in \mathcal{H}$  with  $\hat{R}(h) = 0$  have  $R(h) \leq \epsilon$ .

**Thm.** 2  $N \geq \frac{1}{2\epsilon^2} \left[ \log(|\mathcal{H}|) + \log(\frac{2}{\delta}) \right]$  labeled examples are sufficient so that with probability  $(1-\delta)$  for all  $h \in \mathcal{H}$  we have that  $|R(h) - \hat{R}(h)| \leq \epsilon$ .

Infinite  $|\mathcal{H}|$ 

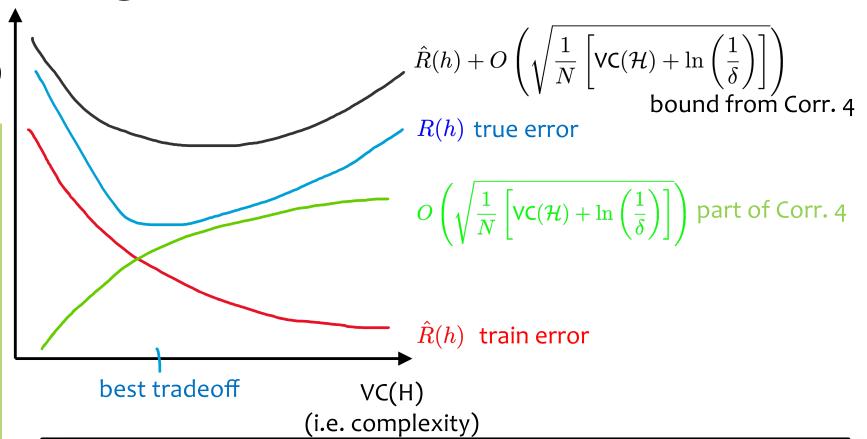
**Thm. 3**  $N = O(\frac{1}{\epsilon} \left[ \text{VC}(\mathcal{H}) \log(\frac{1}{\epsilon}) + \log(\frac{1}{\delta}) \right])$  labeled examples are sufficient so that with probability  $(1 - \delta)$  all  $h \in \mathcal{H}$  with  $\hat{R}(h) = 0$  have  $R(h) \leq \epsilon$ .

**Thm.** 4  $N = O(\frac{1}{\epsilon^2} \left[ \text{VC}(\mathcal{H}) + \log(\frac{1}{\delta}) \right])$  labeled examples are sufficient so that with probability  $(1 - \delta)$  for all  $h \in \mathcal{H}$  we have that  $|R(h) - \hat{R}(h)| \leq \epsilon$ .

# Learning Theory & Model Selection

error
(i.e. lower →
good data fit)

Key Point:
we want
to tradeoff
between
low
training
error and
keeping H
simple
(low VCDim)



Q: Is Corollary 4 useful? A: Yes! Ex: H = Linear Separators in R<sup>M</sup>

VC(H) = M+1

Q: In practice, how do we tradeoff between error and VC(H)?

A: Use a regularizer! That is, reducing the number of (effective) features reduces the VC dimension. More features usually leads to a better fit to the data.

# Text Data

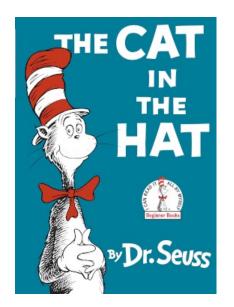


10/31/22

$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	у
("hat")	("cat")	("dog")	("fish")	("mom")	("dad")	(Dr. Seuss)

$x_1$ ("hat")	x <sub>2</sub> ("cat")	x <sub>3</sub> ("dog")	x <sub>4</sub> ("fish")	x <sub>5</sub> ("mom")	x <sub>6</sub> ("dad")	y (Dr. Seuss)
1	1	0	0	0	0	1

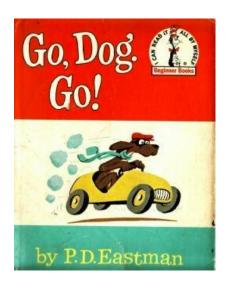
The Cat in the Hat (by Dr. Seuss)



Source: https://en.wikipedia.org/wiki/The Cat in the Hat#/media/File:The Cat in the Hat.png

$x_1$ ("hat")	x <sub>2</sub> ("cat")	x <sub>3</sub> ("dog")	x <sub>4</sub> ("fish")	x <sub>5</sub> ("mom")	x <sub>6</sub> ("dad")	<i>y</i> (Dr. Seuss)
1	1	0	0	0	0	1
0	0	1	0	0	0	0

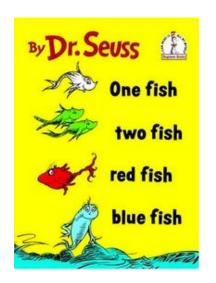
Go, Dog. Go! (by P. D. Eastman)



Source: https://en.wikipedia.org/wiki/Go, Dog. Go!#/media/File:Go Dog Go.jpg

x <sub>1</sub> ("hat")	x <sub>2</sub> ("cat")	x <sub>3</sub> ("dog")	x <sub>4</sub> ("fish")	$x_5$ ("mom")	x <sub>6</sub> ("dad")	<i>y</i> (Dr. Seuss)
1	1	0	0	0	0	1
0	0	1	0	0	0	0
0	0	0	1	0	0	1

One Fish, Two Fish, Red Fish, Blue Fish (by Dr. Seuss)

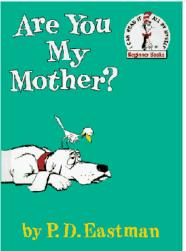


Source

https://en.wikipedia.org/wiki/One Fish, Two Fish, Red Fish, Blue Fish#/media/File:One Fish Two Fish Red Fish Blue Fish (cover art).jpg

x <sub>1</sub> ("hat")	x <sub>2</sub> ("cat")	x <sub>3</sub> ("dog")	x <sub>4</sub> ("fish")	x <sub>5</sub> ("mom")	x <sub>6</sub> ("dad")	<i>y</i> (Dr. Seuss)
1	1	0	0	0	0	1
0	0	1	0	0	0	0
0	0	0	1	0	0	1
0	0	0	0	1	0	0

Are You My Mother? (by P. D. Eastman)



Source: https://en.wikipedia.org/wiki/Are You My Mother%3F#/media/File:Areyoumymother.gif

### Recipe for Closed-form MLE

- 1. Assume data was generated i.i.d. from some model (i.e. write the generative story)  $x^{(i)} \sim p(x|\theta)$
- 2. Write log-likelihood

$$\ell(\boldsymbol{\theta}) = \log p(\mathbf{x}^{(1)}|\boldsymbol{\theta}) + \dots + \log p(\mathbf{x}^{(N)}|\boldsymbol{\theta})$$

3. Compute partial derivatives (i.e. gradient)

$$\frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_1} = \dots$$
$$\frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_2} = \dots$$
$$\frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_M} = \dots$$

4. Set derivatives to zero and solve for  $\theta$ 

$$\partial \ell(\theta)/\partial \theta_{\rm m} = {\rm o \ for \ all \ m} \in \{1, ..., M\}$$
  
 $\theta^{\rm MLE} = {\rm solution \ to \ system \ of \ M \ equations \ and \ M \ variables}$ 

5. Compute the second derivative and check that  $\ell(\theta)$  is concave down at  $\theta^{\text{MLE}}$ 

# Recipe for Closed-form MAP Estimation

- 1. Assume data was generated i.i.d. from some model (i.e. write the generative story)  $\theta \sim p(\theta)$  and then for all i:  $x^{(i)} \sim p(x|\theta)$
- 2. Write log-likelihood

$$\ell_{MAP}(\theta) = \log p(\theta) + \log p(x^{(1)}|\theta) + \dots + \log p(x^{(N)}|\theta)$$

3. Compute partial derivatives (i.e. gradient)

$$\partial \ell_{MAP}(\boldsymbol{\theta})/\partial \theta_1 = \dots$$
  
 $\partial \ell_{MAP}(\boldsymbol{\theta})/\partial \theta_2 = \dots$   
 $\dots$   
 $\partial \ell_{MAP}(\boldsymbol{\theta})/\partial \theta_M = \dots$ 

4. Set derivatives to zero and solve for  $\theta$ 

$$\partial \ell_{MAP}(\theta)/\partial \theta_{m} = 0$$
 for all  $m \in \{1, ..., M\}$   
 $\theta^{MAP} = \text{solution to system of } M \text{ equations and } M \text{ variables}$ 

5. Compute the second derivative and check that  $\ell(\theta)$  is concave down at  $\theta^{\text{MAP}}$ 

#### Classification and Regression: The Big Picture

#### **Recipe for Machine Learning**

- 1. Given data  $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N$
- 2. (a) Choose a decision function  $h_{\theta}(\mathbf{x}) = \cdots$  (parameterized by  $\theta$ )
  - (b) Choose an objective function  $J_{\mathcal{D}}(\boldsymbol{\theta}) = \cdots$  (relies on data)
- 3. Learn by choosing parameters that optimize the objective  $J_{\mathcal{D}}(\boldsymbol{\theta})$

$$\hat{oldsymbol{ heta}} pprox rgmin_{oldsymbol{ heta}} J_{\mathcal{D}}(oldsymbol{ heta})$$

4. Predict on new test example  $\mathbf{x}_{\mathsf{new}}$  using  $h_{\boldsymbol{\theta}}(\cdot)$ 

$$\hat{y} = h_{\boldsymbol{\theta}}(\mathbf{x}_{\mathsf{new}})$$

#### **Optimization Method**

- Gradient Descent:  $\theta \to \theta \gamma \nabla_{\theta} J(\theta)$
- $$\begin{split} \bullet \; \; \mathsf{SGD:} \; & \boldsymbol{\theta} \to \boldsymbol{\theta} \gamma \nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta}) \\ \; \mathsf{for} \; & i \sim \mathsf{Uniform}(1, \dots, N) \\ \; \mathsf{where} \; & J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} J^{(i)}(\boldsymbol{\theta}) \end{split}$$
- mini-batch SGD
- closed form
  - 1. compute partial derivatives
  - 2. set equal to zero and solve

#### **Decision Functions**

- Perceptron:  $h_{\theta}(\mathbf{x}) = \operatorname{sign}(\boldsymbol{\theta}^T \mathbf{x})$
- Linear Regression:  $h_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}$
- Discriminative Models:  $h_{\theta}(\mathbf{x}) = \operatorname*{argmax}_{y} p_{\theta}(y \mid \mathbf{x})$ 
  - Logistic Regression:  $p_{\theta}(y = 1 \mid \mathbf{x}) = \sigma(\boldsymbol{\theta}^T \mathbf{x})$
  - $\circ$  Neural Net (classification):  $p_{\theta}(y=1\mid \mathbf{x}) = \sigma((\mathbf{W}^{(2)})^T \sigma((\mathbf{W}^{(1)})^T \mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)})$
- Generative Models:  $h_{\theta}(\mathbf{x}) = \operatorname*{argmax}_{y} p_{\theta}(\mathbf{x}, y)$ 
  - $\circ$  Naive Bayes:  $p_{m{ heta}}(\mathbf{x},y) = p_{m{ heta}}(y) \prod_{m=1}^M p_{m{ heta}}(x_m \mid y)$

#### **Objective Function**

- MLE:  $J(oldsymbol{ heta}) = -\sum_{i=1}^N \log p(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$
- MCLE:  $J(oldsymbol{ heta}) = -\sum_{i=1}^N \log p(\mathbf{y}^{(i)} \mid \mathbf{x}^{(i)})$
- L2 Regularized:  $J'(\theta) = J(\theta) + \lambda ||\theta||_2^2$  (same as Gaussian prior  $p(\theta)$  over parameters)
- L1 Regularized:  $J'(\theta) = J(\theta) + \lambda ||\theta||_1$  (same as Laplace prior  $p(\theta)$  over parameters)

### MATERIAL COVERED ON EXAM 3

### Totoro's Tunnel

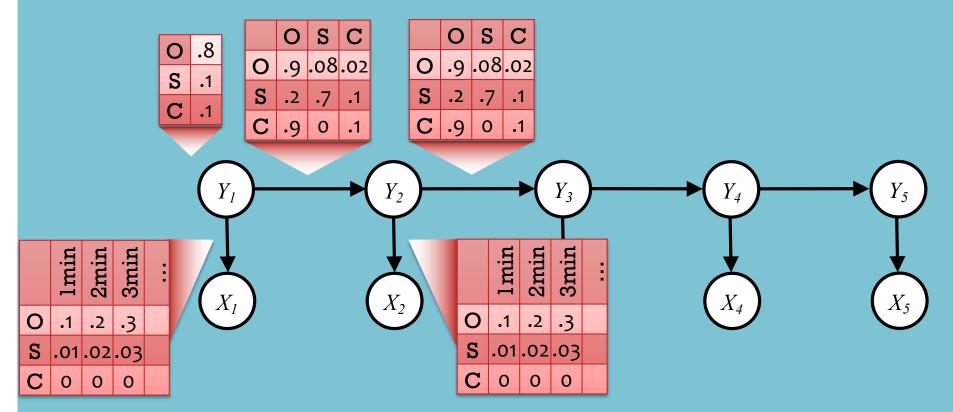




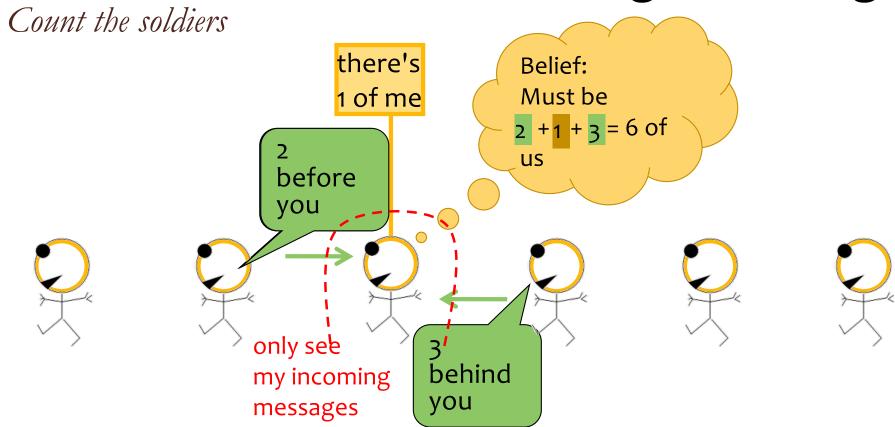
### Hidden Markov Model

#### **HMM Parameters:**

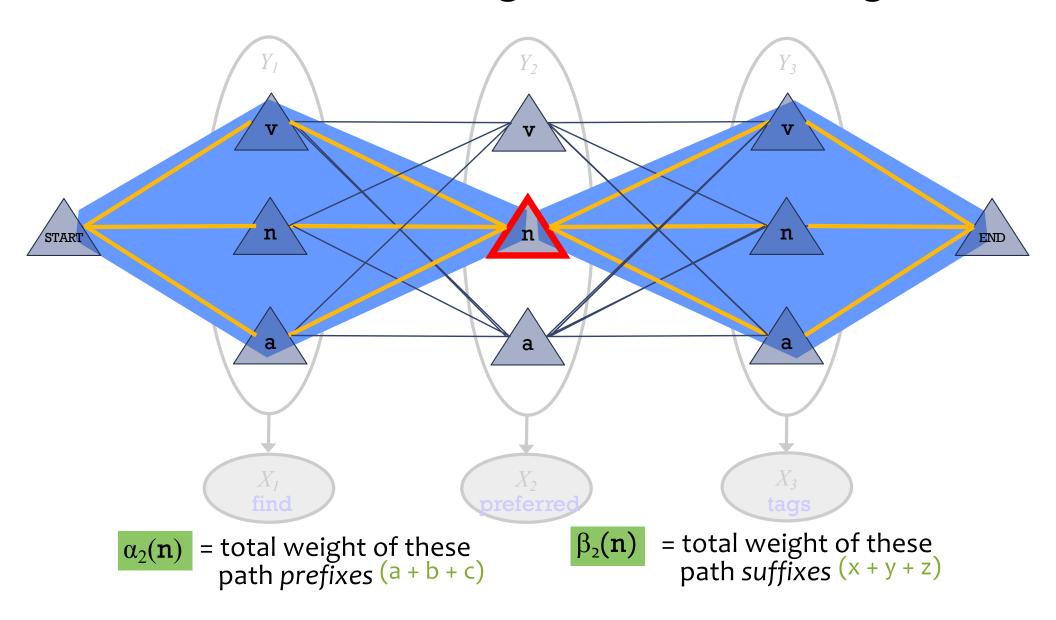
Emission matrix, **A**, where  $P(X_t = k | Y_t = j) = A_{j,k}, \forall t, k$ Transition matrix, **B**, where  $P(Y_t = k | Y_{t-1} = j) = B_{j,k}, \forall t, k$ Initial probs, **C**, where  $P(Y_1 = k) = C_k, \forall k$ 



Great Ideas in ML: Message Passing

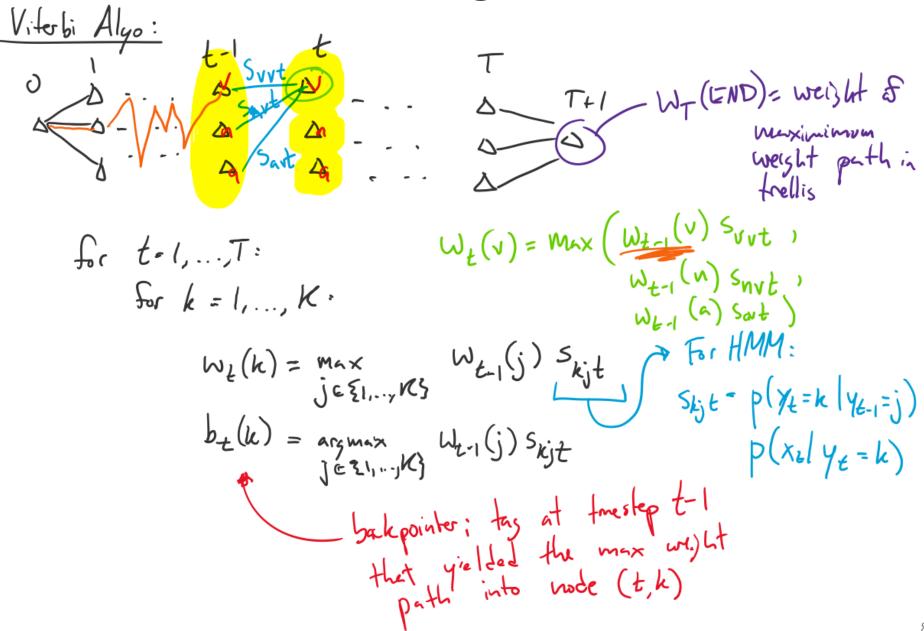


### Forward-Backward Algorithm: Finds Marginals



Product gives ax+ay+az+bx+by+bz+cx+cy+cz = total weight of paths

# Viterbi Algorithm



#### 4 Hidden Markov Models

1. Given the POS tagging data shown, what are the parameter values learned by an HMM?

Verb	Noun	Verb
see	spot	run

Verb	Noun	Verb
run	spot	run

Adj.	Adj.	Noun
funny	funny	spot

#### 4 Hidden Markov Models

- 1. Given the POS tagging data shown, what are the parameter values learned by an HMM?
- 2. Suppose you a learning an HMM POS Tagger, how many POS tag sequences of length 23 are there?
- 3. How does an HMM efficiently search for the most probable tag sequence given a 23-word sentence?

Verb	Noun	Verb
see	spot	run

Verb	Noun	Verb
run	spot	run

Adj.	Adj.	Noun
funny	funny	spot

# Example: Voting for PA Senate Seat

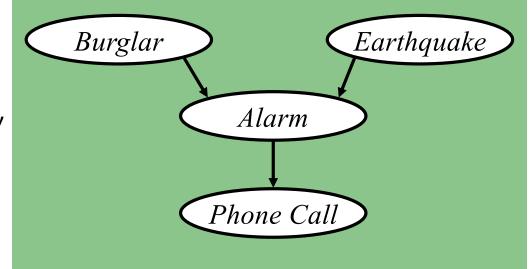


## The "Burglar Alarm" example

- After you get this phone call, suppose you learn that there was a medium-sized earthquake in your neighborhood. Oh, whew! Probably not a burglar after all.
- Earthquake "explains away" the hypothetical burglar.
- But then it must **not** be the case that

 $Burglar \perp\!\!\!\perp Earthquake \mid PhoneCall$  even though

 $Burglar \perp \!\!\! \perp Earthquake$ 



## Example: Tornado Alarms

### Hacking Attack Woke Up Dallas With Emergency Sirens, Officials Say

By ELI ROSENBERG and MAYA SALAM APRIL 8, 2017



Warning sirens in Dallas, meant to alert the public to emergencies like severe weather, started sounding around 11:40 p.m. Friday, and were not shut off until 1:20 a.m. Rex C. Curry for The New York Times

- Imagine that you work at the 911 call center in Dallas
- 2. You receive six calls informing you that the Emergency Weather Sirens are going off
- 3. What do you conclude?

(a) [2 pts.] Write the expression for the joint distribution.

#### 5 Graphical Models [16 pts.]

We use the following Bayesian network to model the relationship between studying (S), being well-rested (R), doing well on the exam (E), and getting an A grade (A). All nodes are binary, i.e.,  $R, S, E, A \in \{0, 1\}$ .

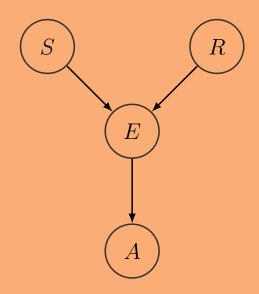


Figure 5: Directed graphical model for problem 5.

(b) [2 pts.] How many parameters are necessary to describe the joint distribution?



### 5 Graphical Models [16 pts.]

We use the following Bayesian network to model the relationship between studying (S), being well-rested (R), doing well on the exam (E), and getting an A grade (A). All nodes are binary, i.e.,  $R, S, E, A \in \{0, 1\}$ .

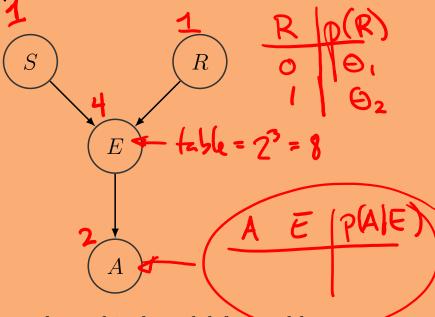


Figure 5: Directed graphical model for problem 5.

P0//

(d) [2 pts.] Is S marginally independent of R? Is S conditionally independent of R given E? Answer yes or no to each questions and provide a brief explanation why.

A = bric B = Yes C= No

### 5 Graphical Models [16 pts.]

We use the following Bayesian network to model the relationship between studying (S), being well-rested (R), doing well on the exam (E), and getting an A grade (A). All nodes are binary, i.e.,  $R, S, E, A \in \{0, 1\}$ .

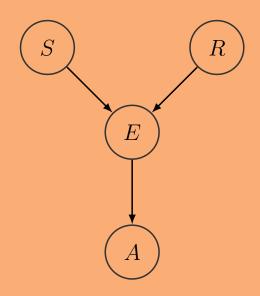


Figure 5: Directed graphical model for problem 5.

#### 5 Graphical Models

(f) [3 pts.] Give two reasons why the graphical models formalism is convenient when compared to learning a full joint distribution.

## A Few Problems for Bayes Nets

Suppose we already have the parameters of a Bayesian Network...

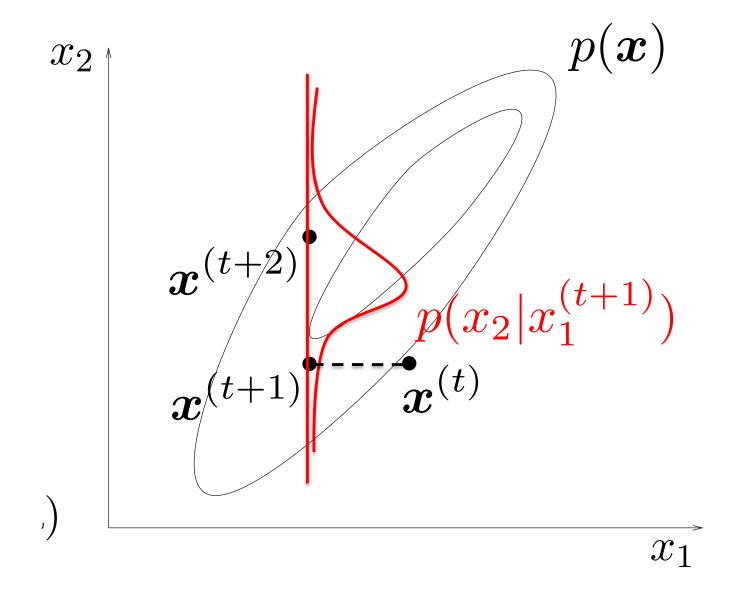
- How do we compute the probability of a specific assignment to the variables?
   P(T=t, H=h, A=a, C=c)
- 2. How do we draw a sample from the joint distribution?  $t,h,a,c \sim P(T, H, A, C)$
- 3. How do we compute marginal probabilities? P(A) = ...



- 4. How do we draw samples from a conditional distribution?  $t,h,a \sim P(T, H, A \mid C = c)$
- 5. How do we compute conditional marginal probabilities?  $P(H \mid C = c) = ...$

Can we use samples ?

# Gibbs Sampling



# MDP Example: Multi-armed bar

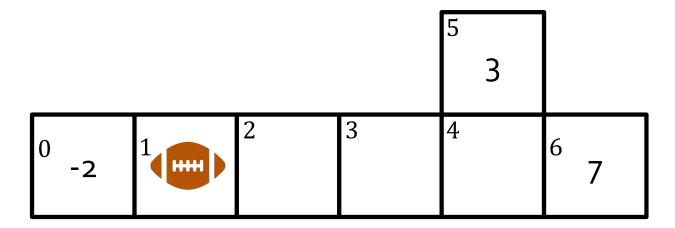
• Single state: |S| = 1

• Three actions:  $\mathcal{A} = \{1, 2, 3\}$ 

Rewards are stochas



## RL: Value Function Example



$$R(s,a) = \begin{cases} -2 & \text{if entering state 0 (safety)} \\ 3 & \text{if entering state 5 (field goal)} \\ 7 & \text{if entering state 6 (touch down)} \\ 0 & \text{otherwise} \end{cases}$$

$$\gamma = 0.9$$

Today's lecture is brought to you by the letter Q



# Today's lecture is brought to you by the letter Q

- Inputs: reward function R(s, a), transition probabilities  $p(s' \mid s, a)$
- Initialize  $V(s) = 0 \ \forall \ s \in \mathcal{S}$  (or randomly)
- While not converged, do:
  - For  $s \in S$ 
    - For  $a \in \mathcal{A}$

$$Q(s,a) = R(s,a) + \gamma \sum_{s' \in S} p(s' \mid s,a) V(s')$$

•  $V(s) \leftarrow \max_{a \in \mathcal{A}} Q(s, a)$ 

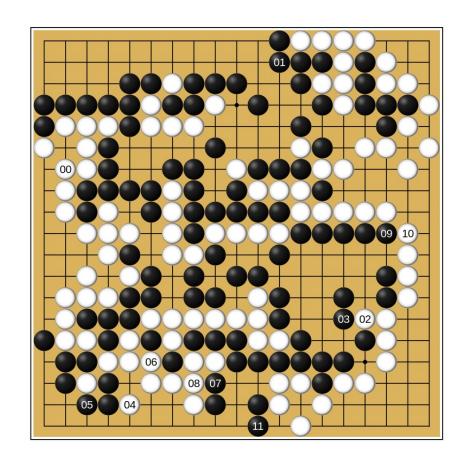
• For  $s \in S$ 

$$\pi^*(s) \leftarrow \underset{a \in \mathcal{A}}{\operatorname{argmax}} R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' \mid s, a) V(s')$$

• Return  $\pi^*$ 

## Playing Go

- 19-by-19 board
- Players alternate placing black and white stones
- The goal is claim more territory than the opponent



The number of legal Go board states is ~10<sup>170</sup> (<a href="https://en.wikipedia.org/wiki/Go\_and\_mathematics">https://en.wikipedia.org/wiki/Go\_and\_mathematics</a>) compared to the number of possible games of chess, ~10<sup>120</sup>

#### 7.1 Reinforcement Learning





3. (1 point) Please select one statement that is true for reinforcement learning and supervised learning.

- O Reinforcement learning is a kind of supervised learning problem because you can treat the reward and next state as the label and each state, action pair as the training data.
- 80% C
- O Reinforcement learning differs from supervised learning because it has a temporal structure in the learning process, whereas, in supervised learning, the prediction of a data point does not affect the data you would see in the future.

#### 7.1 Reinforcement Learning

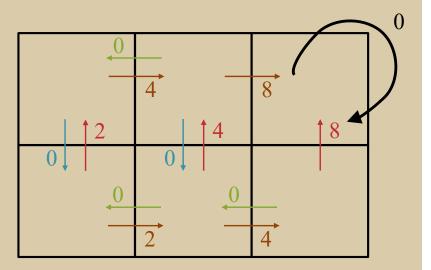
- 3. (1 point) Please select one statement that is true for reinforcement learning and supervised learning.
  - Reinforcement learning is a kind of supervised learning problem because you can treat the reward and next state as the label and each state, action pair as the training data.
  - O Reinforcement learning differs from supervised learning because it has a temporal structure in the learning process, whereas, in supervised learning, the prediction of a data point does not affect the data you would see in the future.

Q3

- 4. (1 point) **True or False:** Value iteration is better at balancing exploration and exploitation compared with policy iteration.
  - A O True 35%
  - C | False | 5%
  - B = toxic

#### 7.1 Reinforcement Learning

- 1. For the R(s,a) values shown on the arrows below, what is the corresponding optimal policy? Assume the discount factor is 0.1
- 2. For the R(s,a) values shown on the arrows below, which are the corresponding  $V^*(s)$  values? Assume the discount factor is 0.1
- 3. For the R(s,a) values shown on the arrows below, which are the corresponding  $Q^*(s,a)$  values? Assume the discount factor is 0.1
- 4. Could we change R(s,a) such that all the  $V^*(s)$  values change but the optimal policy stays the same? If so, show how and if not, briefly explain why not.





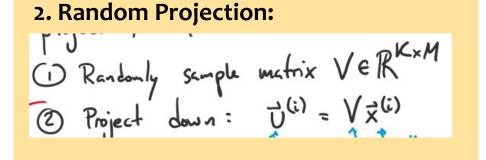
## Shortcut Example



https://www.youtube.com/watch?v=MIJN9pEfPfE

## PCA section in one slide...

# 1. Dimensionality reduction:



#### 3. Definition of PCA:

Choose the matrix V that either...

- 1. minimizes reconstruction error
- 2. consists of the K eigenvectors with largest eigenvalue

The above are equivalent definitions.

#### 4. Algorithm for PCA:

The option we'll focus on:

Run Singular Value
Decomposition (SVD) to
obtain all the eigenvectors.
Keep just the top-K to form V.
Play some tricks to keep
things efficient.

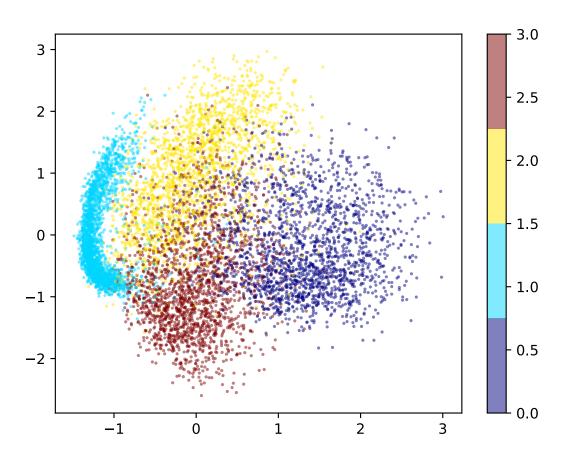
#### 5. An Example



# Projecting MNIST digits

#### **Task Setting:**

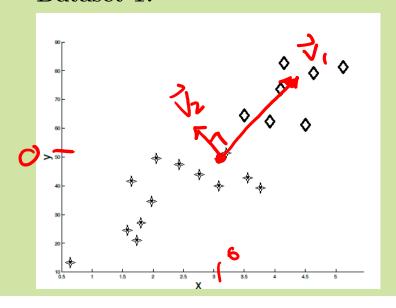
- 1. Take 25x25 images of digits and project them down to 2 components
- 2. Plot the 2 dimensional points



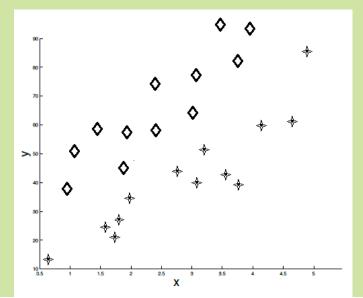
## 4 Principal Component Analysis [16 pts.]

- (a) In the following plots, a train set of data points X belonging to two classes on  $\mathbb{R}^2$  are given, where the original features are the coordinates (x, y). For each, answer the following questions:
  - (i) [3 pt.] Draw all the principal components.
  - (ii) [6 pts.] Can we correctly classify this dataset by using a threshold function after projecting onto one of the principal components? If so, which principal component should we project onto? If not, explain in 1–2 sentences why it is not possible.

#### Dataset 1:



#### Dataset 2:

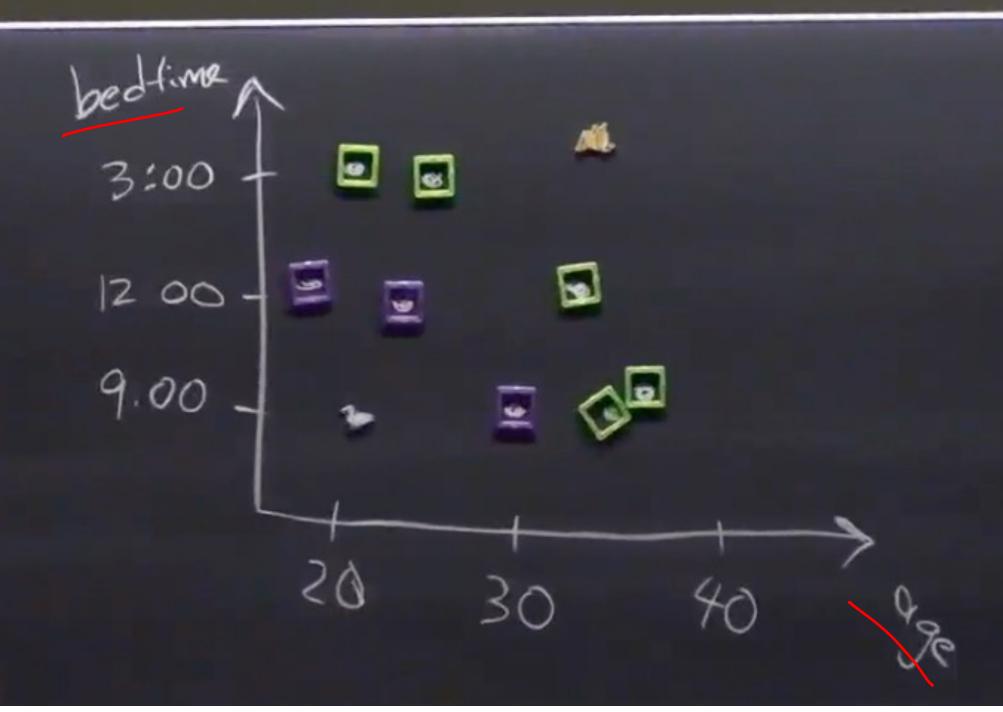


## K-Means Algorithm

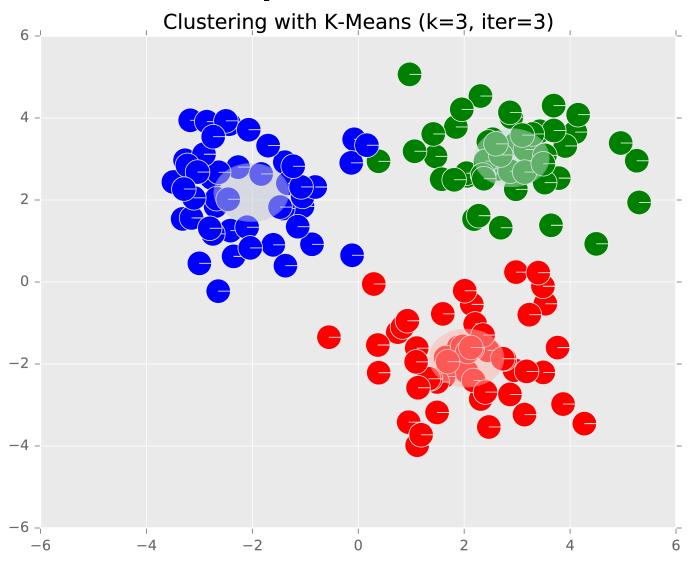
Given unlabeled feature vectors

$$D = \{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$$

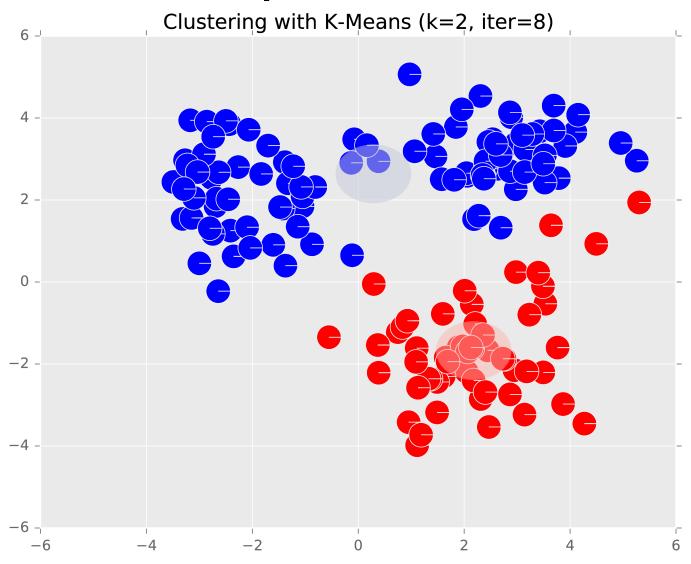
- Initialize cluster centers  $c = \{c^{(1)}, \dots, c^{(K)}\}$
- Repeat until convergence:
  - for i in {1,..., N}
     z<sup>(i)</sup> ← index j of cluster center nearest to x<sup>(i)</sup>
     for j in {1,...,K}
     c<sup>(j)</sup> ← mean of all points assigned to cluster j



# Example: K-Means



## Example: K-Means



#### 2.2 Lloyd's algorithm

Q4

Circle the image which depicts the cluster center positions after 1 iteration of Lloyd's algorithm.

K-means

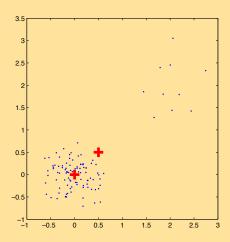
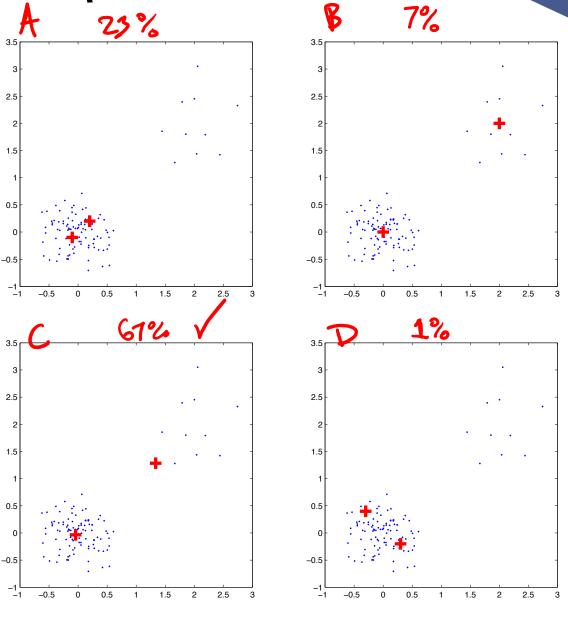
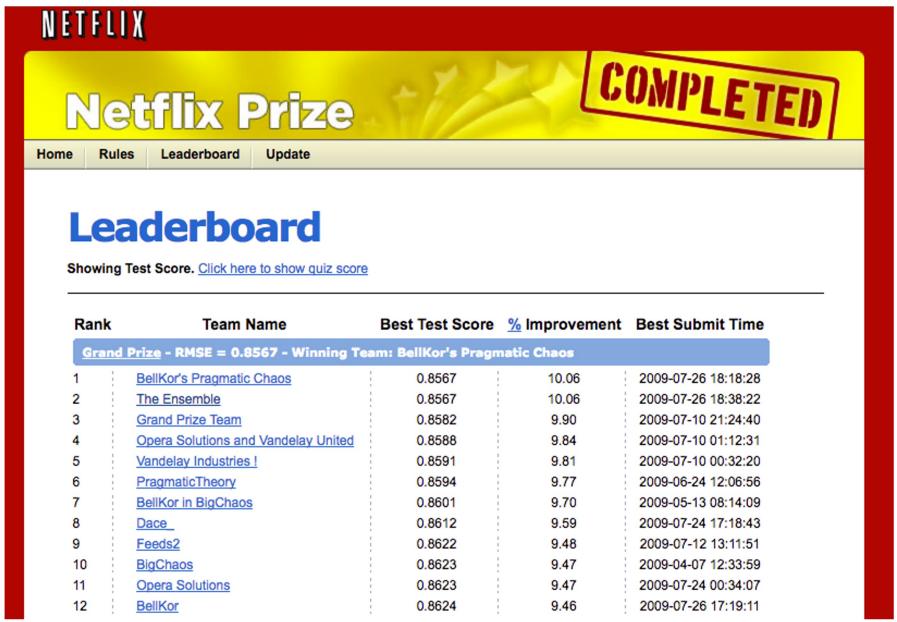


Figure 2: Initial data and cluster centers





## Recommender Systems



Weighted Majority Algorithm

(Littlestone & Warmuth, 1994)

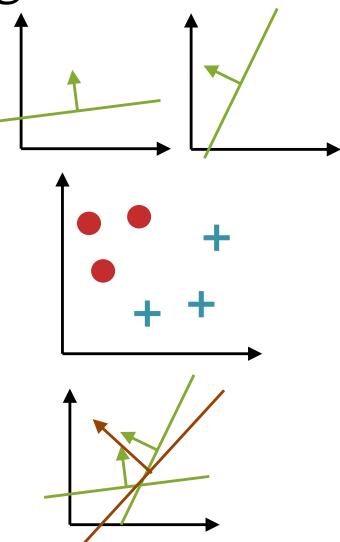
Given: pool A of binary classifiers (that you know nothing about)

 Data: stream of examples (i.e. online learning setting)

 Goal: design a new learner that uses the predictions of the pool to make new predictions

## • Algorithm:

- Initially weight all classifiers equally
- Receive a training example and predict the (weighted) majority vote of the classifiers in the pool
- Down-weight classifiers that contribute to a mistake by a factor of  $\boldsymbol{\beta}$



# Weighted Majority Algorithm

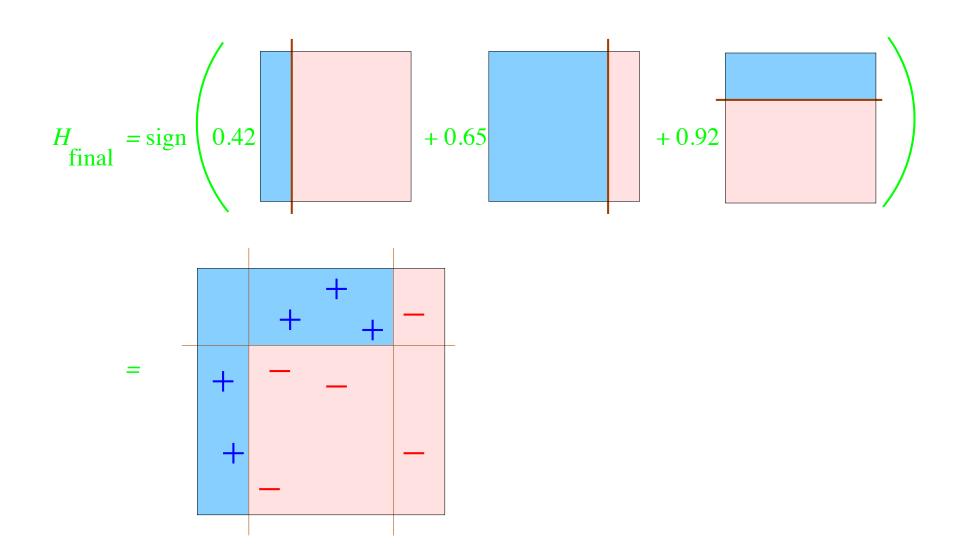
Theorems (Littlestone & Warmuth, 1994)

For the general case where WM is applied to a pool  $\mathcal{A}$  of algorithms we show the following upper bounds on the number of mistakes made in a given sequence of trials:

- 1.  $O(\log |\mathcal{A}| + m)$ , if one algorithm of  $\mathcal{A}$  makes at most m mistakes.
- 2.  $O(\log \frac{|A|}{k} + m)$ , if each of a subpool of k algorithms of A makes at most m mistakes.
- 3.  $O(\log \frac{|A|}{k} + \frac{m}{k})$ , if the total number of mistakes of a subpool of k algorithms of A is at most m.

These are
"mistake
bounds" of the
variety we saw
for the
Perceptron
algorithm

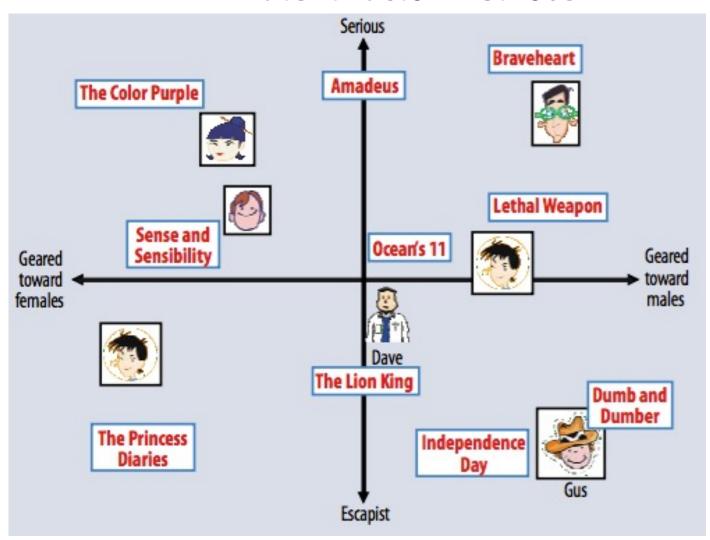
## AdaBoost: Toy Example



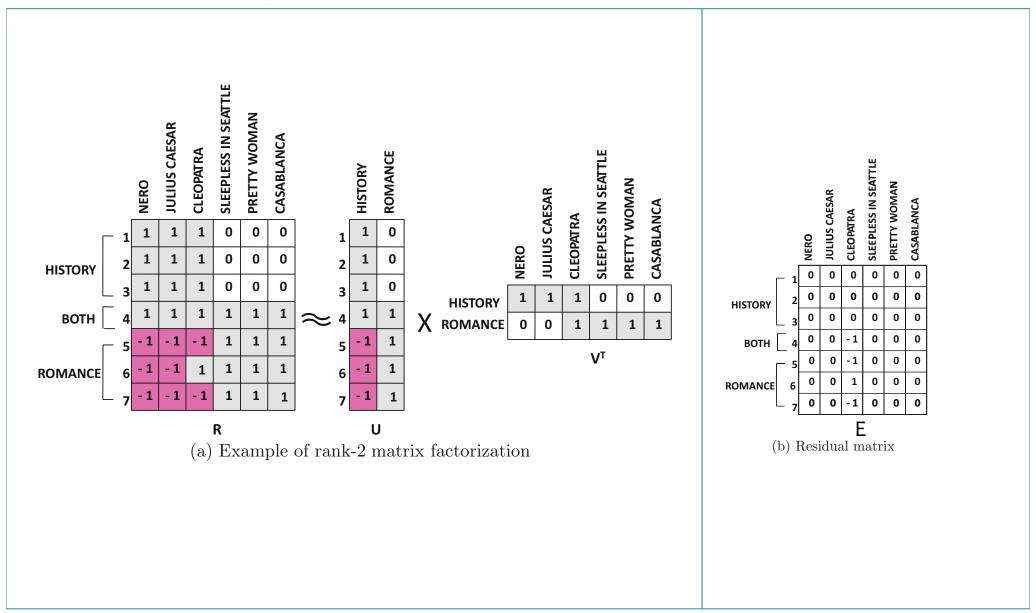
## Two Types of Collaborative Filtering

#### 2. Latent Factor Methods

- Assume that both movies and users live in some lowdimensional space describing their properties
- Recommend a
   movie based on its
   proximity to the
   user in the latent
   space
- Example Algorithm:
   Matrix Factorization



## Example: MF for Netflix Problem



P0//

## Recommending Movies

## Question: 05

Which of the following pieces of information about user behavior could be used to improve a collaborative filtering system?

## Select all that apply

- A. # of times a user watched a given movie
- B. Total # of movies a user has watched
- C. How often a user turns on subtitles
- D. # of times a user paused a given movie
  - E. How many accounts a user is associated with
  - F. # of DVDs a user can rent at a time
  - G. None of the above

### Classification and Regression: The Big Picture

#### **Recipe for Machine Learning**

- 1. Given data  $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N$
- 2. (a) Choose a decision function  $h_{m{ heta}}(\mathbf{x}) = \cdots$  (parameterized by  $m{ heta}$ )
  - (b) Choose an objective function  $J_{\mathcal{D}}(\boldsymbol{\theta}) = \cdots$  (relies on data)
- 3. Learn by choosing parameters that optimize the objective  $J_{\mathcal{D}}(\boldsymbol{\theta})$

$$\hat{oldsymbol{ heta}} pprox rgmin_{oldsymbol{ heta}} J_{\mathcal{D}}(oldsymbol{ heta})$$

4. Predict on new test example  $\mathbf{x}_{\mathsf{new}}$  using  $h_{\boldsymbol{\theta}}(\cdot)$ 

$$\hat{y} = h_{\boldsymbol{\theta}}(\mathbf{x}_{\mathsf{new}})$$

#### **Optimization Method**

- Gradient Descent:  $\theta \to \theta \gamma \nabla_{\theta} J(\theta)$
- $$\begin{split} \bullet \; & \mathsf{SGD:} \; \boldsymbol{\theta} \to \boldsymbol{\theta} \gamma \nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta}) \\ & \mathsf{for} \; i \sim \mathsf{Uniform}(1,\ldots,N) \\ & \mathsf{where} \; J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} J^{(i)}(\boldsymbol{\theta}) \end{split}$$
- mini-batch SGD
- closed form
  - 1. compute partial derivatives
  - 2. set equal to zero and solve

#### **Decision Functions**

- Perceptron:  $h_{\theta}(\mathbf{x}) = \operatorname{sign}(\boldsymbol{\theta}^T \mathbf{x})$
- Linear Regression:  $h_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}$
- Discriminative Models:  $h_{\theta}(\mathbf{x}) = \operatorname*{argmax}_{y} p_{\theta}(y \mid \mathbf{x})$ 
  - Logistic Regression:  $p_{\theta}(y = 1 \mid \mathbf{x}) = \sigma(\boldsymbol{\theta}^T \mathbf{x})$
  - o Neural Net (classification):  $p_{\theta}(y=1\mid \mathbf{x}) = \sigma((\mathbf{W}^{(2)})^T \sigma((\mathbf{W}^{(1)})^T \mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)})$
- Generative Models:  $h_{\theta}(\mathbf{x}) = \operatorname*{argmax}_{y} p_{\theta}(\mathbf{x}, y)$ 
  - $\circ$  Naive Bayes:  $p_{m{ heta}}(\mathbf{x},y) = p_{m{ heta}}(y) \prod_{m=1}^M p_{m{ heta}}(x_m \mid y)$

#### **Objective Function**

- ullet MLE:  $J(oldsymbol{ heta}) = -\sum_{i=1}^N \log p(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$
- ullet MCLE:  $J(oldsymbol{ heta}) = -\sum_{i=1}^N \log p(\mathbf{y}^{(i)} \mid \mathbf{x}^{(i)})$
- L2 Regularized:  $J'(\theta) = J(\theta) + \lambda ||\theta||_2^2$  (same as Gaussian prior  $p(\theta)$  over parameters)
- L1 Regularized:  $J'(\theta) = J(\theta) + \lambda ||\theta||_1$  (same as Laplace prior  $p(\theta)$  over parameters)

# Learning Paradigms

Paradigm	Data
Supervised	$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N \qquad \mathbf{x} \sim p^*(\cdot) \text{ and } y = c^*(\cdot)$
$\hookrightarrow$ Regression	$y^{(i)} \in \mathbb{R}$
$\hookrightarrow$ Classification	$y^{(i)} \in \{1, \dots, K\}$
$\hookrightarrow$ Binary classification	$y^{(i)} \in \{+1, -1\}$
$\hookrightarrow$ Structured Prediction	$\mathbf{y}^{(i)}$ is a vector
Unsupervised	$\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^{N} \qquad \mathbf{x} \sim p^*(\cdot)$
Semi-supervised	$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N_1} \cup \{\mathbf{x}^{(j)}\}_{j=1}^{N_2}$
Online	$\mathcal{D} = \{ (\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), (\mathbf{x}^{(3)}, y^{(3)}), \ldots \}$
Active Learning	$\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ and can query $y^{(i)} = c^*(\cdot)$ at a cost
Imitation Learning	$\mathcal{D} = \{ (s^{(1)}, a^{(1)}), (s^{(2)}, a^{(2)}), \ldots \}$
Reinforcement Learning	$\mathcal{D} = \{ (s^{(1)}, a^{(1)}, r^{(1)}), (s^{(2)}, a^{(2)}, r^{(2)}), \ldots \}$

## ML Big Picture

#### **Learning Paradigms:**

What data is available and when? What form of prediction?

- supervised learning
- unsupervised learning
- semi-supervised learning
- reinforcement learning
- active learning
- imitation learning
- domain adaptation
- online learning
- density estimation
- recommender systems
- feature learning
- manifold learning
- dimensionality reduction
- ensemble learning
- distant supervision
- hyperparameter optimization

#### **Theoretical Foundations:**

What principles guide learning?

- probabilistic
- ☐ information theoretic
- evolutionary search
- ☐ ML as optimization

#### **Problem Formulation:**

What is the structure of our output prediction?

boolean Binary Classification

categorical Multiclass Classification

ordinal Ordinal Classification

real Regression

ordering Ranking

multiple discrete Structured Prediction

multiple continuous (e.g. dynamical systems)

both discrete & (e.g. mixed graphical models)

cont.

Application Areas

Key challenges?

NLP, Speech, Computer
Vision, Robotics, Medicine,

## Facets of Building ML Systems:

How to build systems that are robust, efficient, adaptive, effective?

- 1. Data prep
- 2. Model selection
- 3. Training (optimization / search)
- 4. Hyperparameter tuning on validation data
- 5. (Blind) Assessment on test data

#### Big Ideas in ML:

Which are the ideas driving development of the field?

- inductive bias
- generalization / overfitting
- bias-variance decomposition
- generative vs. discriminative
- deep nets, graphical models
- PAC learning
- distant rewards

## Course Level Objectives

#### You should be able to...

- Implement and analyze existing learning algorithms, including well-studied methods for classification, regression, structured prediction, clustering, and representation learning
- 2. Integrate multiple facets of practical machine learning in a single system: data preprocessing, learning, regularization and model selection
- 3. Describe the the formal properties of models and algorithms for learning and explain the practical implications of those results
- 4. Compare and contrast different paradigms for learning (supervised, unsupervised, etc.)
- 5. Design experiments to evaluate and compare different machine learning techniques on real-world problems
- 6. Employ probability, statistics, calculus, linear algebra, and optimization in order to develop new predictive models or learning methods
- Given a description of a ML technique, analyze it to identify (1) the expressive power of the formalism; (2) the inductive bias implicit in the algorithm; (3) the size and complexity of the search space; (4) the computational properties of the algorithm: (5) any guarantees (or lack thereof) regarding termination, convergence, correctness, accuracy or generalization power.

## **Course Staff**





#### Team A (HW2, HW6)



Hayden Kim



Aditi Sharma



Sami Kale



Hang Shu

#### Team B (HW3, HW7)



Chu Weng



Kalvin Chang



**Brandon Wang** 



Monica Geng



Shubham Virmani

#### Team C (HW4, HW8)



Abhi Vijayakumar



**Lulu Ricketts** 



Jack Lyu



Jeneel Mashru



Chongling Zhu

**Pranay Gundam** 

## Team D (HW5, HW9)



Tara Lakdawala



Qiuyi Yin



Alex Xie



Yuchen Xu



Lavanya Gupta

