



### 10-601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

# Final Exam Review

Matt Gormley Lecture 29 Apr. 29, 2020

### Reminders

- Homework 9: Learning Paradigms
  - Out: Wed, Apr. 22
  - Due: Wed, Apr. 29 at 11:59pm
  - Can only be submitted up to 3 days late,
     so we can return grades before final exam
- Final Exam Practice Problems
  - Out: Wed, Apr. 29
- Final Exam
  - Mon, May 04 (1pm 4pm)
- Today's In-Class Poll
  - http://poll.mlcourse.org

### **EXAM LOGISTICS**

### Final Exam

#### Time / Location

- Time: Registrar-scheduled Exam
   Mon, May 4th at 1:00pm 4:00pm
- Online Exam: Same format as Midterm Exam 2
- Please watch Piazza carefully for announcements logistics

#### Logistics

- Distribution of Topics: Lectures 19 28 (95%), Lectures 1 18 (5%)
- Format of questions:
  - Multiple choice
  - True / False (with justification)
  - Derivations
  - Short answers
  - Interpreting figures
  - Implementing algorithms on paper
- You are encouraged to bring one 8½ x 11 sheet of notes (front and back)
- Open book according to my definition on Piazza: https://piazza.com/class/k4wzus8w2c11u6?cid=1673

### Final Exam

### How to Prepare

- Attend (or watch) this final exam review session
- Review Practice Problems: Exam 3
  - Disclaimer: the practice problems are somewhere between homework-style problems and exam-style problems
- Review this year's homework problems
- Review the poll questions from each lecture
- Consider whether you have achieved the learning objectives for each lecture / section

### Final Exam

### Advice (for during the exam)

- Solve the easy problems first
   (e.g. multiple choice before derivations)
  - if a problem seems extremely complicated you're likely missing something
- Don't leave any answer blank!
- If you make an assumption, write it down
- If you look at a question and don't know the answer:
  - we probably haven't told you the answer
  - but we've told you enough to work it out
  - imagine arguing for some answer and see if you like it

### Topics for Midterm 1

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

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

### Topics for Midterm 2

- Classification
  - Binary Logistic Regression
  - Multinomial Logistic Regression
- Important Concepts
  - Stochastic Gradient
     Descent
  - Regularization
  - Feature Engineering
- Feature Learning
  - Neural Networks
  - Basic NN Architectures
  - Backpropagation

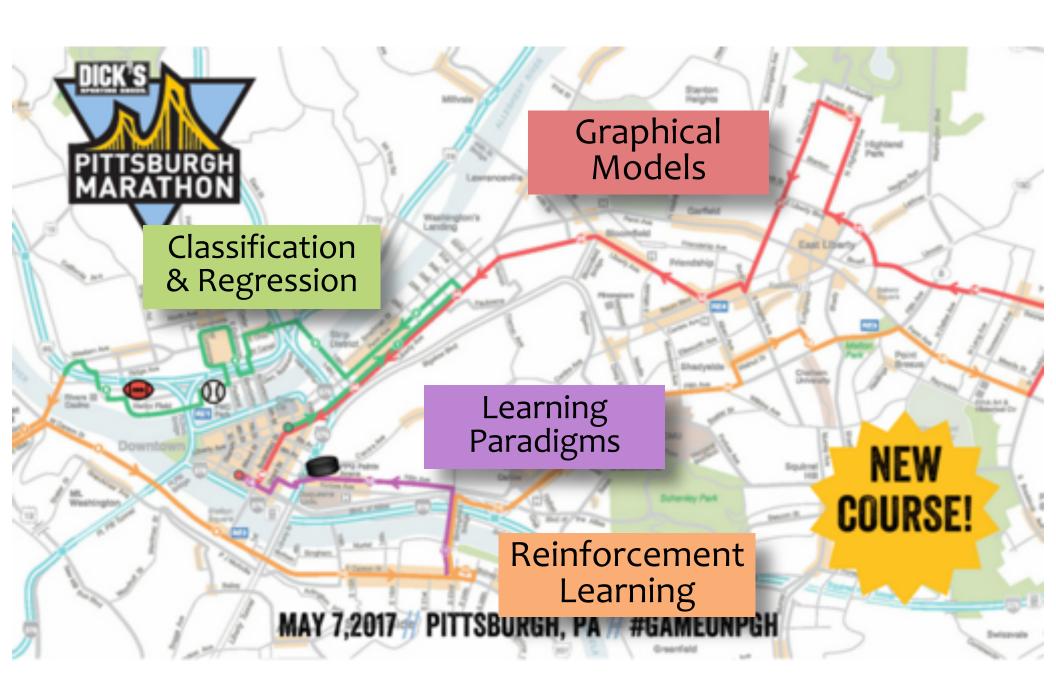
- Learning Theory
  - PAC Learning
- Generative Models
  - Generative vs.
     Discriminative
  - MLE / MAP
  - Naïve Bayes

### **Topics for Final Exam**

- 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
  - SVM (large-margin)
  - Kernels
  - Ensemble Methods
  - Recommender Systems









### Great Race: route and street closing schedule



Learning as Memorization

Learning as
Optimization

Learning from Rewards

Learning and
Structure

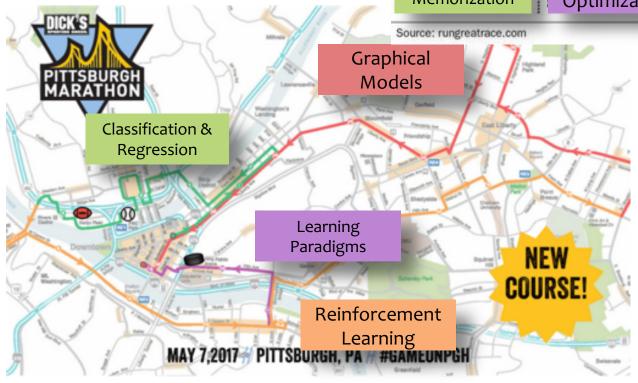
Source: rungreatrace.com Post-Gazette

# A new **combined** course...

... with the best (uphill climbs) from both

#### Great Race: route and street closing schedule





Post-Gazette

Material Covered Before Midterm Exam 2

# **SAMPLE QUESTIONS**

# Matching Game

### Goal: Match the Algorithm to its Update Rule

#### 1. SGD for Logistic Regression

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = p(y|x)$$

#### 2. Least Mean Squares

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}$$

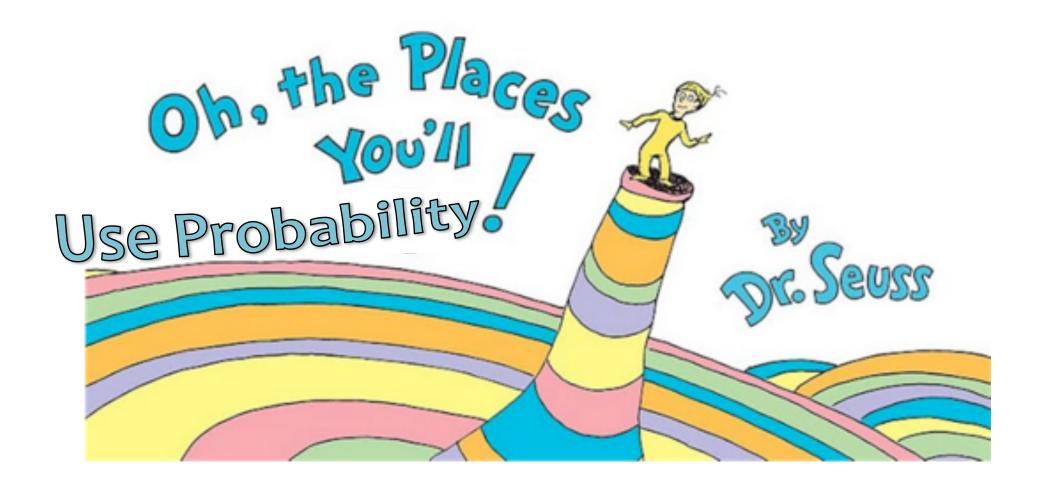
3. Perceptron (next lecture)

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

4. 
$$\theta_k \leftarrow \theta_k + (h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)})$$

$$\theta_k \leftarrow \theta_k + \frac{1}{1 + \exp \lambda (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)})}$$

6. 
$$\theta_k \leftarrow \theta_k + \lambda (h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)}) x_k^{(i)}$$



#### 1.4 Probability

Assume we have a sample space  $\Omega$ . Answer each question with **T** or **F**.

(a) [1 pts.] **T** or **F**: If events A, B, and C are disjoint then they are independent.

(b) [1 pts.] **T** or **F**: 
$$P(A|B) \propto \frac{P(A)P(B|A)}{P(A|B)}$$
. (The sign ' $\propto$ ' means 'is proportional to')

## Medical Diagnosis

#### **Interview Transcript**

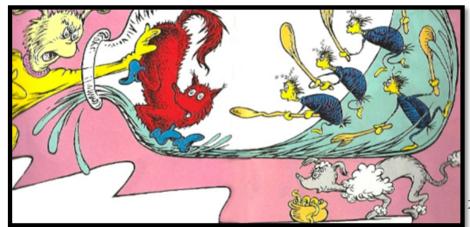
Date: Jan. 15, 2020.

**Parties:** Matt Gormley and Doctor E.

Topic: Medical decision making

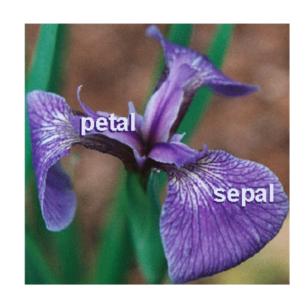
- Matt: Welcome. Thanks for interviewing with me today.
- Dr. E: Interviewing...?
- Matt: Yes. For the record, what type of doctor are you?
- Dr. E: Who said I'm a doctor?
- Matt: I thought when we set up this interview you said—
- Dr. E: I'm a preschooler.
- Matt: Good enough. Today, I'd like to learn how you would determine whether or not your little brother is sick given his symptoms.
- Dr. E: He's not sick.
- Matt: We haven't started yet. Now, suppose he is sneezing. Is he sick?
- Dr. E: No, that's just the sniffles.
- Matt: What if he is coughing; Is he sick?
- Dr. E: No, he just has a cough.
- [Editor's note: preschoolers unilaterally agree that having the sniffles or a cough is not the same as being sick.]

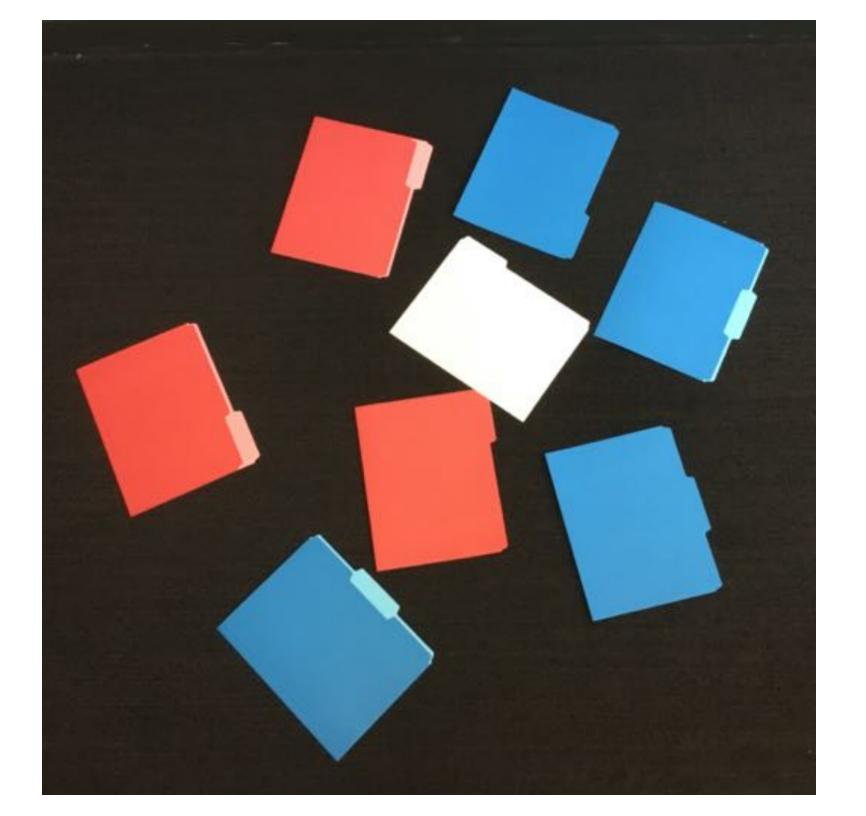
- Matt: What if he's both sneezing and coughing?
- Dr. E: Then he's sick.
- Matt: Got it. What if your little brother is sneezing and coughing, plus he's a doctor.
- Dr. E: Then he's not sick.
- Matt: How do you know?
- Dr. E: Doctors don't get sick.
- Matt: What if he is not sneezing, but is coughing, 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. E: 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 sick.
- Matt: Got it. Can I use this conversation in my lecture?
- Dr. E: Yes





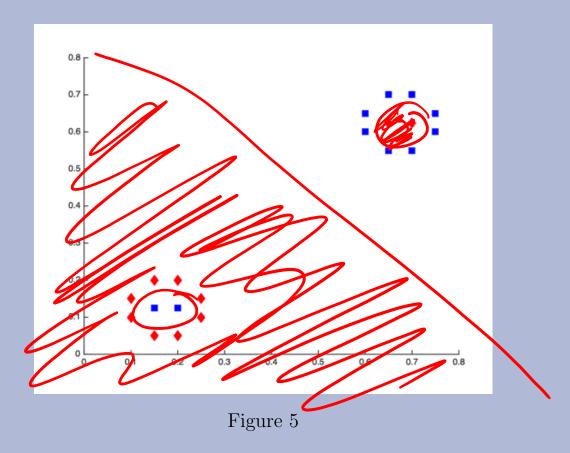
| Species | Sepal<br>Length | Sepal<br>Width | Petal<br>Length | Petal<br>Width |
|---------|-----------------|----------------|-----------------|----------------|
| 0       | 4.3             | 3.0            | 1.1             | 0.1            |
| 0       | 4.9             | 3.6            | 1.4             | 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            |



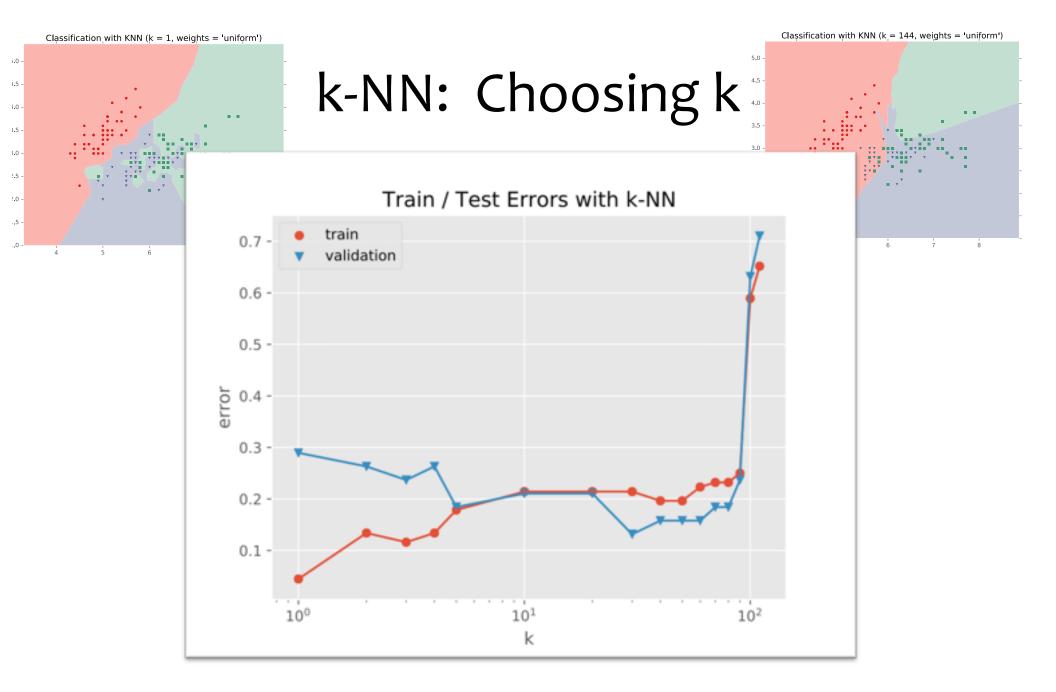


#### 4 K-NN [12 pts]

Now we will apply K-Nearest Neighbors using Euclidean distance to a binary classification task. We assign the class of the test point to be the class of the majority of the k nearest neighbors. A point can be its own neighbor.

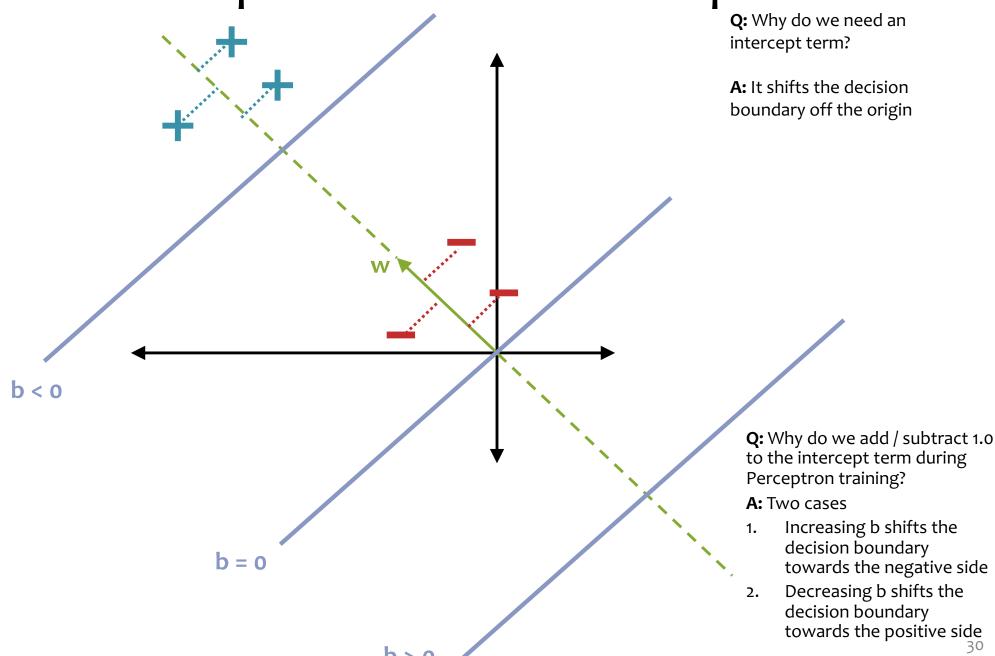


3. [2 pts] What value of k minimizes leave-one-out cross-validation error for the dataset shown in Figure 5? What is the resulting error?

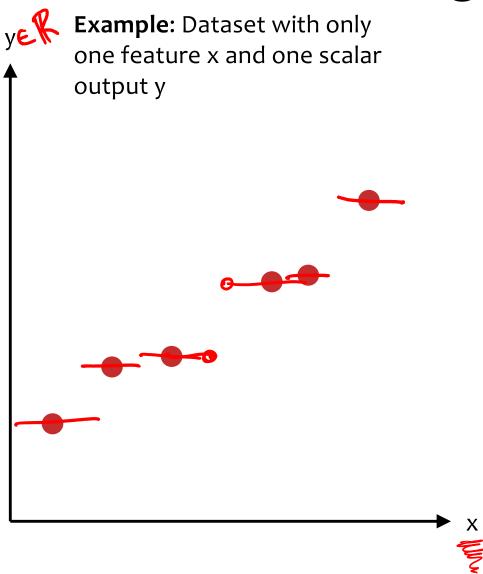


Fisher Iris Data: varying the value of k

### Perceptron & The Intercept Term



# k-NN Regression



#### k=1 Nearest Neighbor Regression

- Train: store all (x, y) pairs
- Predict: pick the nearest x in training data and return its y

#### k=2 Nearest Neighbor Distance Weighted Regression

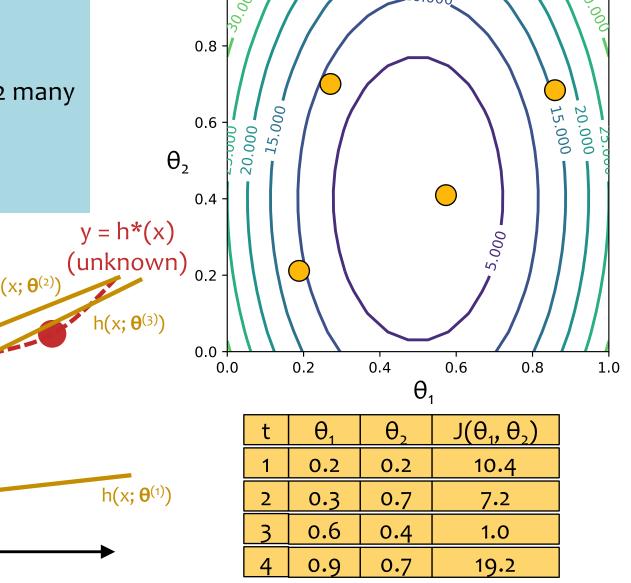
- Train: store all (x, y) pairs
- Predict: pick the nearest two instances x<sup>(n1)</sup> and x<sup>(n2)</sup> in training data and return the weighted average of their y values

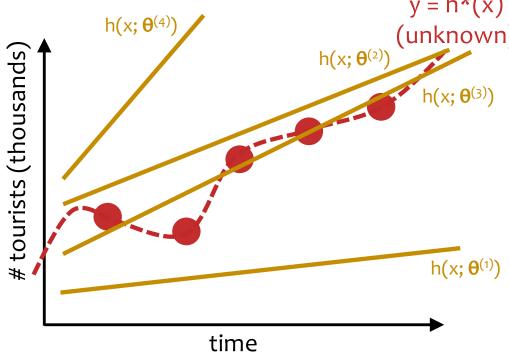


Linear Regression by Rand. Guessing  $J(\theta) = J(\theta_1, \theta_2) = \frac{1}{N} \sum_{i=1}^{N} \left( y^{(i)} - \theta^T \mathbf{x}^{(i)} \right)^2$ 

# 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  $J(\theta)$





#### 3.1 Linear regression

Consider the dataset S plotted in Fig. 1 along with its associated regression line. For each of the altered data sets  $S^{\text{new}}$  plotted in Fig. 3, indicate which regression line (relative to the original one) in Fig. 2 corresponds to the regression line for the new data set. Write your answers in the table below.

| Dataset         | (a) | (b) | (c) | (d) | (e) |
|-----------------|-----|-----|-----|-----|-----|
| Regression line |     |     |     |     |     |

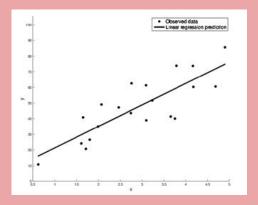


Figure 1: An observed data set and its associated regression line.

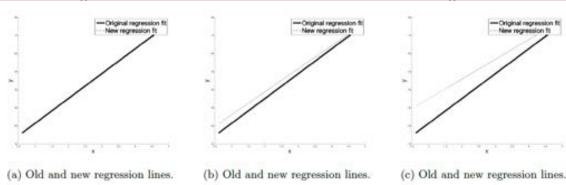
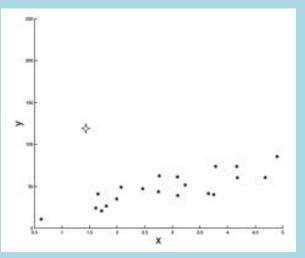


Figure 2: New regression lines for altered data sets  $S^{\text{new}}$ .

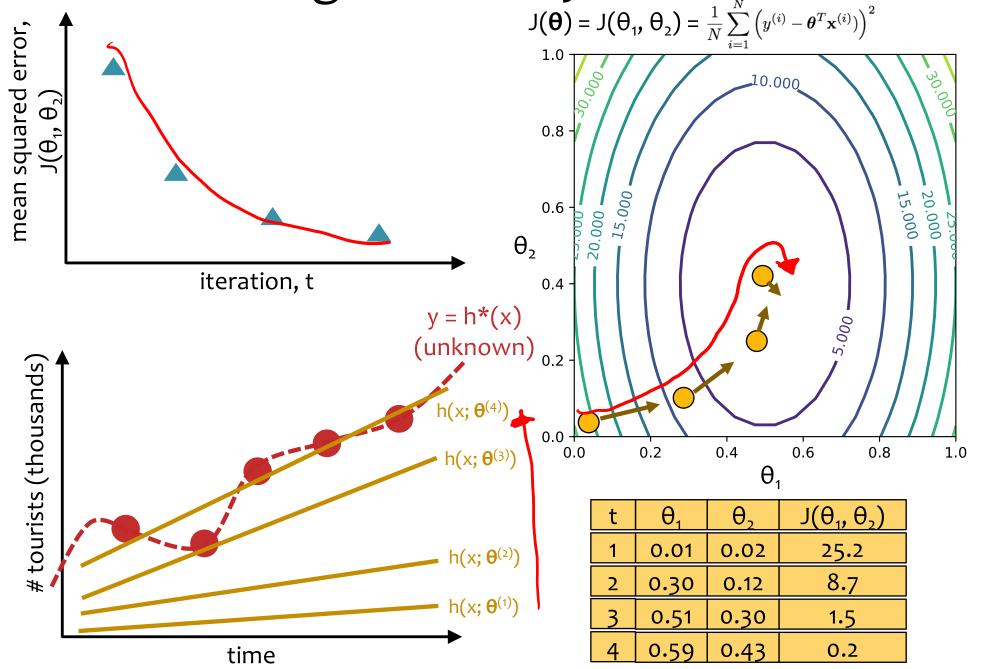


(a) Adding one outlier to the original data set.

# Topographical Maps



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



#### 3.1 Linear regression

Consider the dataset S plotted in Fig. 1 along with its associated regression line. For each of the altered data sets  $S^{\text{new}}$  plotted in Fig. 3, indicate which regression line (relative to the original one) in Fig. 2 corresponds to the regression line for the new data set. Write your answers in the table below.

| Dataset         | (a) | (b) | (c) | (d) | (e) |
|-----------------|-----|-----|-----|-----|-----|
| Regression line |     |     |     |     |     |

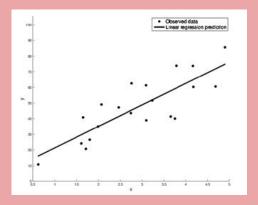


Figure 1: An observed data set and its associated regression line.

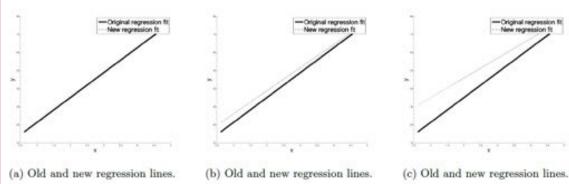
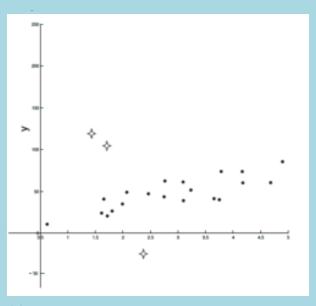


Figure 2: New regression lines for altered data sets  $S^{\text{new}}$ .



(c) Adding three outliers to the original data set. Two on one side and one on the other side.

#### 3.1 Linear regression

Consider the dataset S plotted in Fig. 1 along with its associated regression line. For each of the altered data sets  $S^{\text{new}}$  plotted in Fig. 3, indicate which regression line (relative to the original one) in Fig. 2 corresponds to the regression line for the new data set. Write your answers in the table below.

| Dataset         | (a) | (b) | (c) | (d) | (e) |
|-----------------|-----|-----|-----|-----|-----|
| Regression line |     |     |     |     |     |

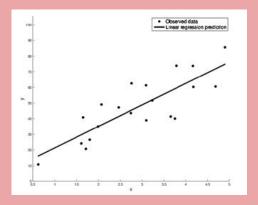


Figure 1: An observed data set and its associated regression line.

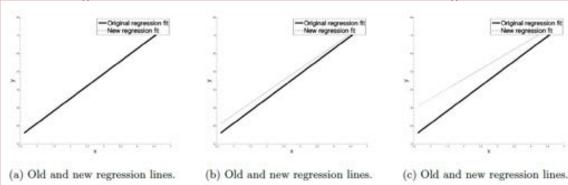
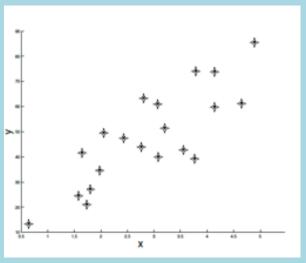


Figure 2: New regression lines for altered data sets  $S^{\text{new}}$ .



(d) Duplicating the original data set.

#### 3.1 Linear regression

Consider the dataset S plotted in Fig. 1 along with its associated regression line. For each of the altered data sets  $S^{\text{new}}$  plotted in Fig. 3, indicate which regression line (relative to the original one) in Fig. 2 corresponds to the regression line for the new data set. Write your answers in the table below.

| Dataset         | (a) | (b) | (c) | (d) | (e) |
|-----------------|-----|-----|-----|-----|-----|
| Regression line |     |     |     |     |     |

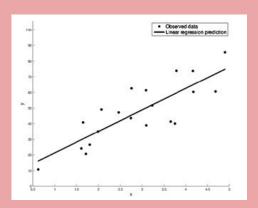


Figure 1: An observed data set and its associated regression line.

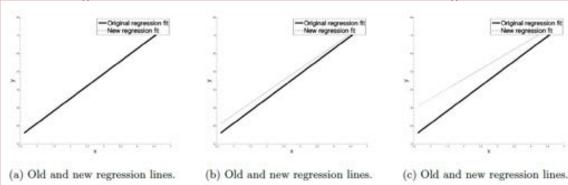
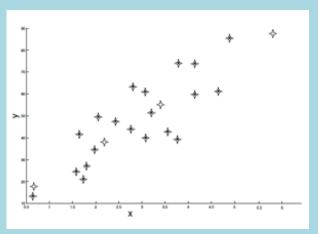
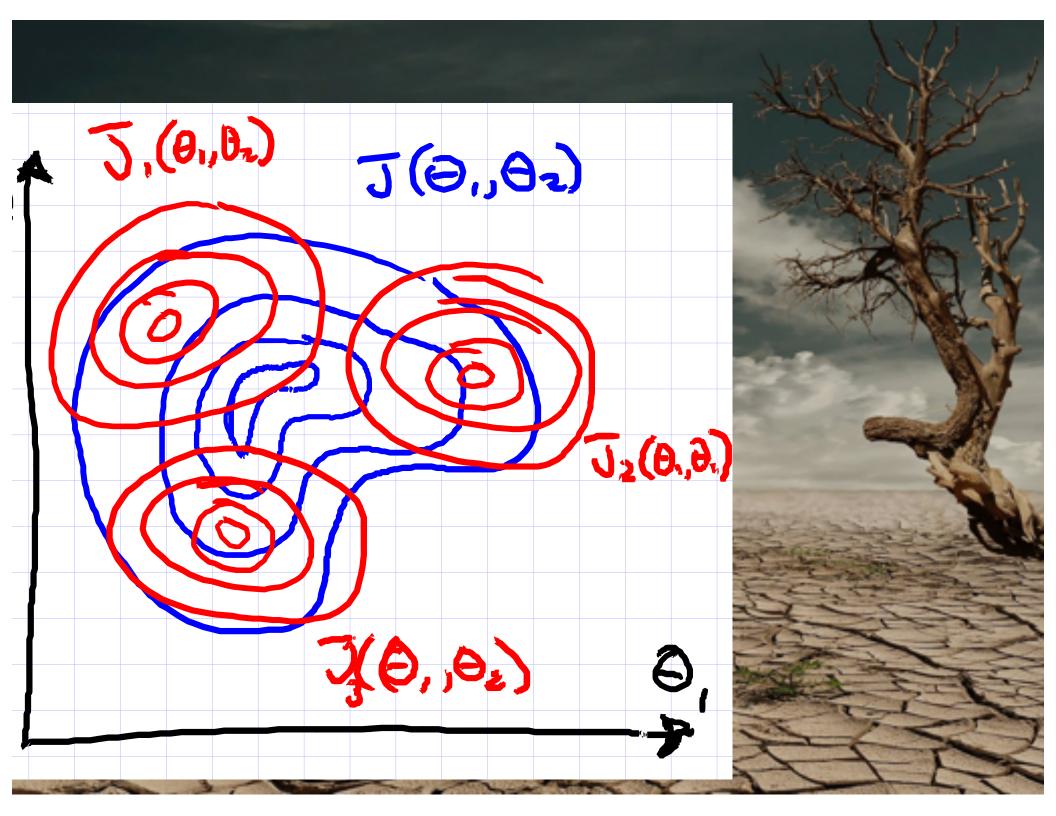


Figure 2: New regression lines for altered data sets  $S^{\text{new}}$ .



(e) Duplicating the original data set and adding four points that lie on the trajectory of the original regression line.



# Robotic Farming

|                                      | Deterministic                               | Probabilistic                         |
|--------------------------------------|---|---------------------------------------|
| Classification (binary output)       | Is this a picture of a wheat kernel?        | Is this plant drought resistant?      |
| Regression<br>(continuous<br>output) | How many wheat kernels are in this picture? | What will the yield of this plant be? |





# Multinomial Logistic Regression



#### 3.2 Logistic regression

Given a training set  $\{(x_i, y_i), i = 1, ..., n\}$  where  $x_i \in \mathbb{R}^d$  is a feature vector and  $y_i \in \{0, 1\}$  is a binary label, we want to find the parameters  $\hat{w}$  that maximize the likelihood for the training set, assuming a parametric model of the form

$$p(y = 1|x; w) = \frac{1}{1 + \exp(-w^T x)}.$$

The conditional log likelihood of the training set is

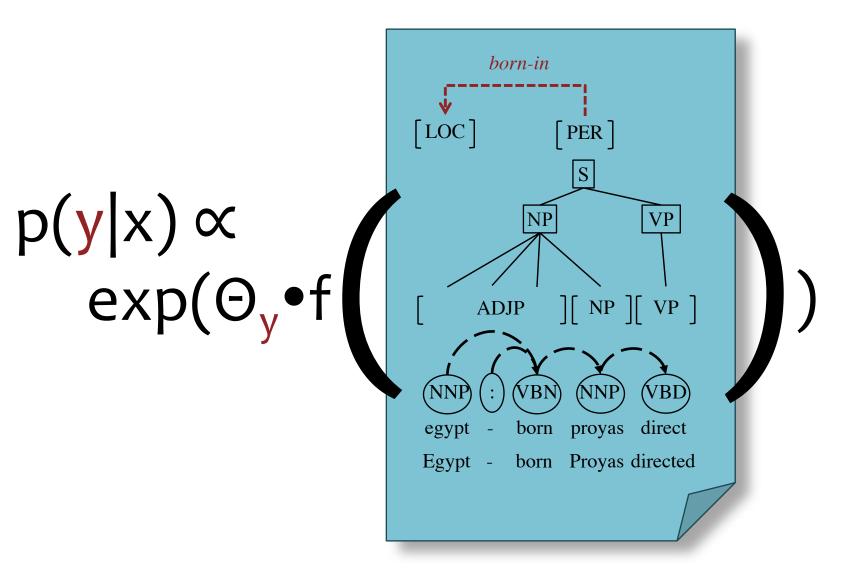
$$\ell(w) = \sum_{i=1}^{n} y_i \log p(y_i, | x_i; w) + (1 - y_i) \log(1 - p(y_i, | x_i; w)),$$

and the gradient is

$$\nabla \ell(w) = \sum_{i=1}^{n} (y_i - p(y_i|x_i; w))x_i.$$

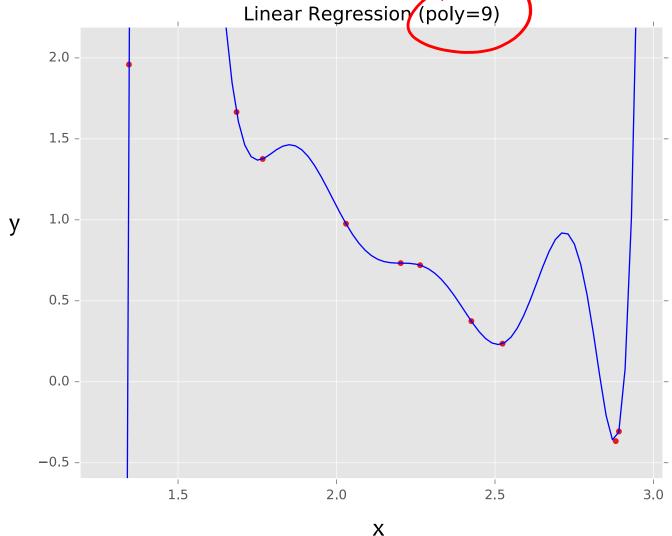
- (b) [5 pts.] What is the form of the classifier output by logistic regression?
- (c) [2 pts.] **Extra Credit:** Consider the case with binary features, i.e,  $x \in \{0,1\}^d \subset \mathbb{R}^d$ , where feature  $x_1$  is rare and happens to appear in the training set with only label 1. What is  $\hat{w}_1$ ? Is the gradient ever zero for any finite w? Why is it important to include a regularization term to control the norm of  $\hat{w}$ ?

#### Handcrafted Features



### Example: Linear Regression

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



true "unknown"
target function is
linear with
negative slope
and gaussian
noise

### Regularization

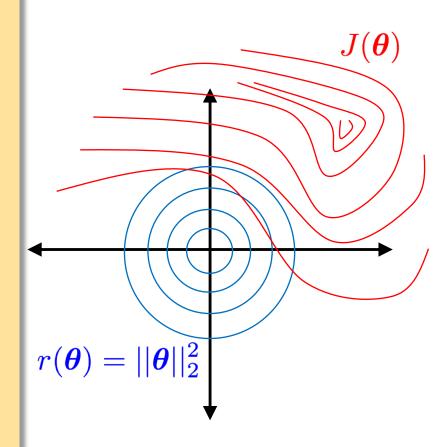
#### **Question:**

Suppose we are minimizing  $J'(\theta)$  where

$$J'(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \lambda r(\boldsymbol{\theta})$$

As  $\lambda$  increases, the minimum of J'( $\theta$ ) will...

- A. ... move towards the midpoint between  $J'(\theta)$  and  $r(\theta)$
- B. ... move towards the minimum of  $J(\theta)$
- C. ... move towards the minimum of  $r(\theta)$
- D. ... move towards a theta vector of positive infinities
- E. ... move towards a theta vector of negative infinities
- F. ... stay the same



#### 2.1 Train and test errors

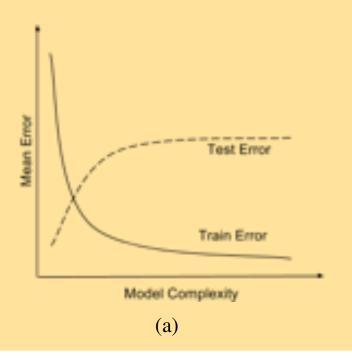
In this problem, we will see how you can debug a classifier by looking at its train and test errors. Consider a classifier trained till convergence on some training data  $\mathcal{D}^{\text{train}}$ , and tested on a separate test set  $\mathcal{D}^{\text{test}}$ . You look at the test error, and find that it is very high. You then compute the training error and find that it is close to 0.

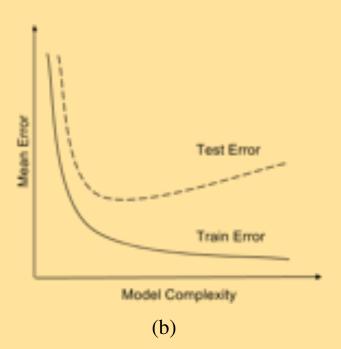
- 1. [4 pts] Which of the following is expected to help? Select all that apply.
  - (a) Increase the training data size.
  - (b) Decrease the training data size.
  - (c) Increase model complexity (For example, if your classifier is an SVM, use a more complex kernel. Or if it is a decision tree, increase the depth).
  - (d) Decrease model complexity.
  - (e) Train on a combination of  $\mathcal{D}^{\text{train}}$  and  $\mathcal{D}^{\text{test}}$  and test on  $\mathcal{D}^{\text{test}}$
  - (f) Conclude that Machine Learning does not work.

#### 2.1 Train and test errors

In this problem, we will see how you can debug a classifier by looking at its train and test errors. Consider a classifier trained till convergence on some training data  $\mathcal{D}^{\text{train}}$ , and tested on a separate test set  $\mathcal{D}^{\text{test}}$ . You look at the test error, and find that it is very high. You then compute the training error and find that it is close to 0.

4. **[1 pts]** Say you plot the train and test errors as a function of the model complexity. Which of the following two plots is your plot expected to look like?





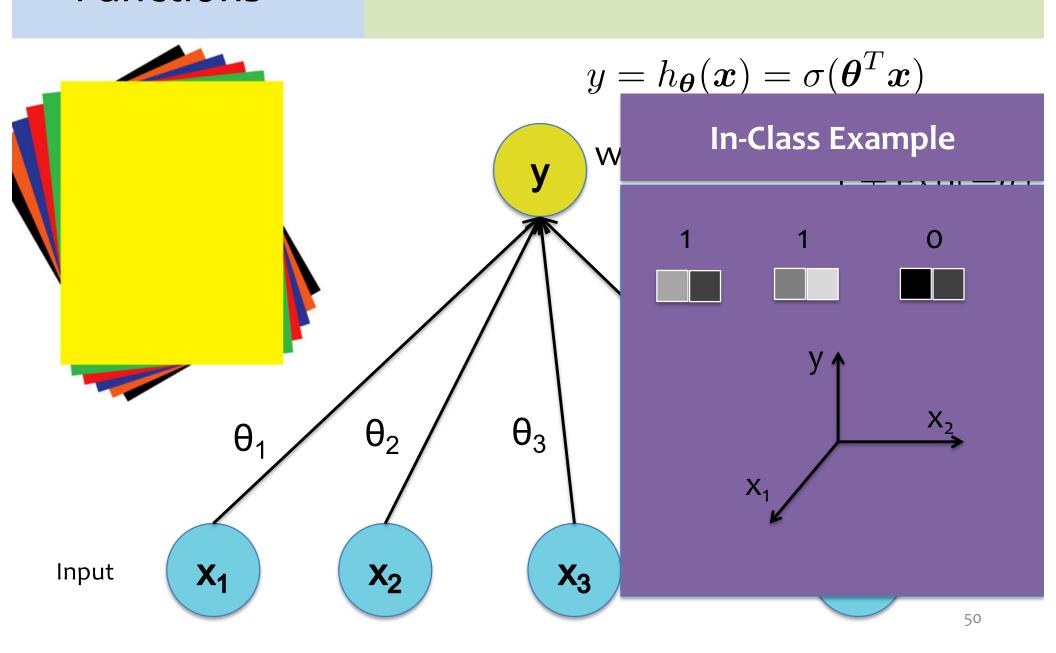
#### 4.1 True or False

Answer each of the following questions with **T** or **F** and **provide a one line justification**.

(a) [2 pts.] Consider two datasets  $D^{(1)}$  and  $D^{(2)}$  where  $D^{(1)} = \{(x_1^{(1)}, y_1^{(1)}), ..., (x_n^{(1)}, y_n^{(1)})\}$  and  $D^{(2)} = \{(x_1^{(2)}, y_1^{(2)}), ..., (x_m^{(2)}, y_m^{(2)})\}$  such that  $x_i^{(1)} \in \mathbb{R}^{d_1}, x_i^{(2)} \in \mathbb{R}^{d_2}$ . Suppose  $d_1 > d_2$  and n > m. Then the maximum number of mistakes a perceptron algorithm will make is higher on dataset  $D^{(1)}$  than on dataset  $D^{(2)}$ .

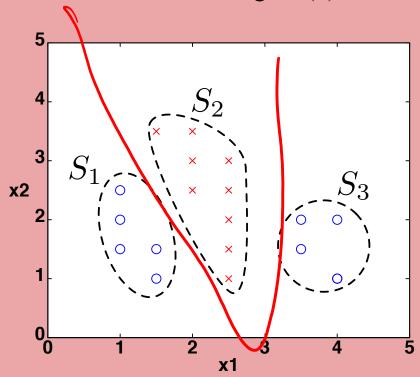
### Decision Functions

### Logistic Regression

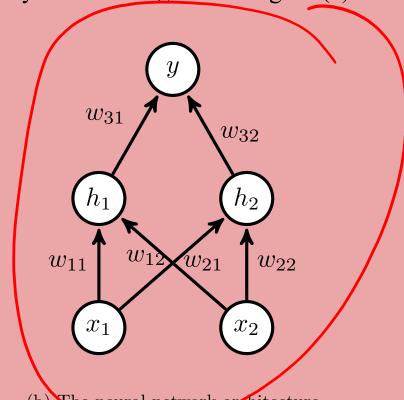


#### **Neural Networks**

Can the neural network in Figure (b) correctly classify the dataset given in Figure (a)?



(a) The dataset with groups  $S_1$ ,  $S_2$ , and  $S_3$ .

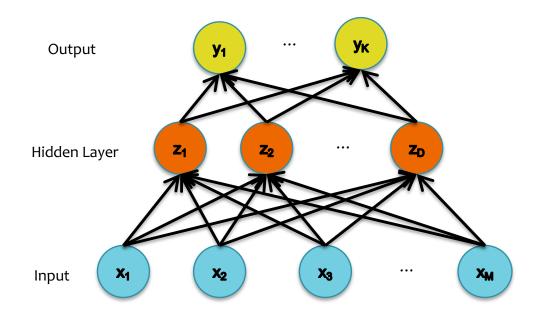


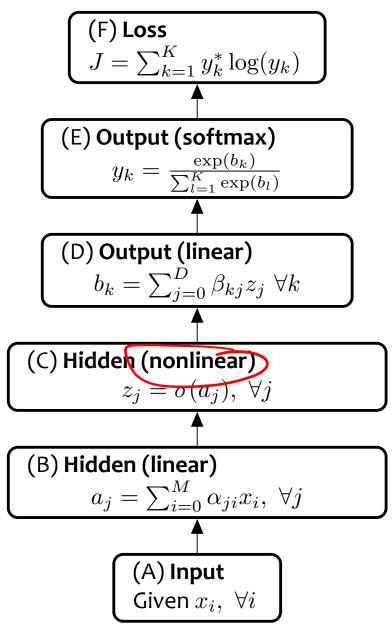
(b) The neural network architecture

### Multi-Class Output

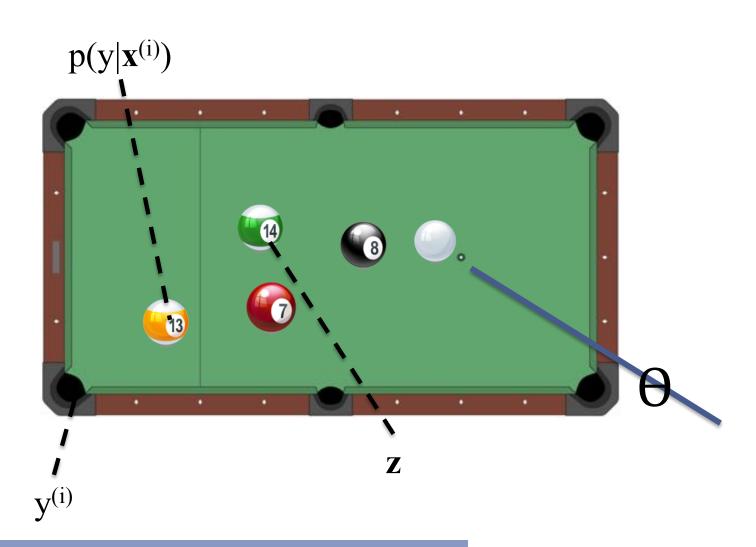
#### Softmax:

$$y_k = \frac{\exp(b_k)}{\sum_{l=1}^K \exp(b_l)}$$



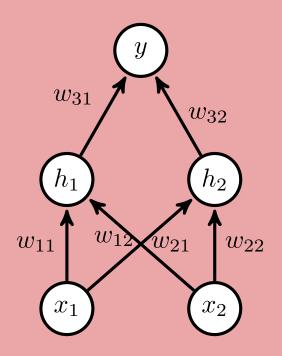


### **Error Back-Propagation**



#### **Neural Networks**

Apply the backpropagation algorithm to obtain the partial derivative of the mean-squared error of y with the true value  $y^*$  with respect to the weight  $w_{22}$  assuming a sigmoid nonlinear activation function for the hidden layer.



(b) The neural network architecture

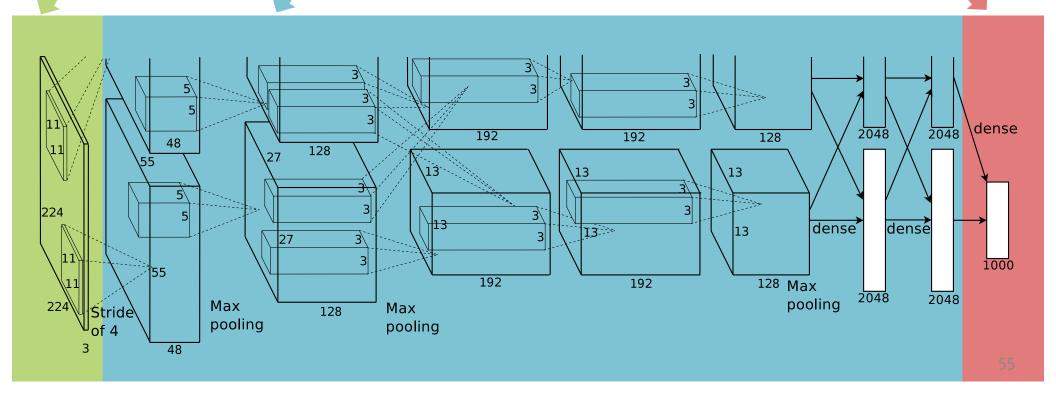
### Architecture #2: AlexNet

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



### **Bidirectional RNN**

inputs:  $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$ 

hidden units:  $\overrightarrow{\mathbf{h}}$  and  $\overleftarrow{\mathbf{h}}$ 

outputs:  $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$ 

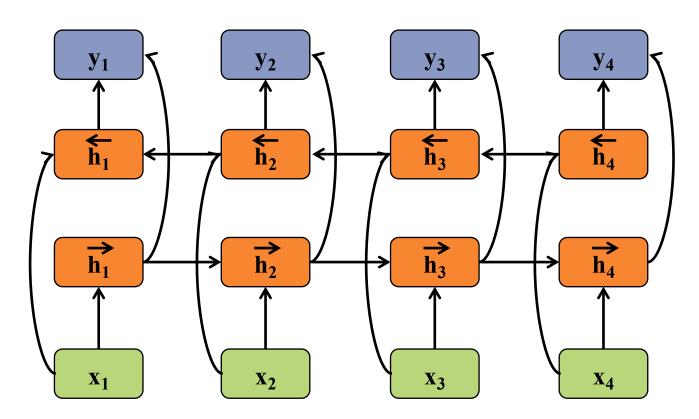
nonlinearity:  $\mathcal{H}$ 

**Recursive Definition:** 

$$\overrightarrow{h}_{t} = \mathcal{H}\left(W_{x\overrightarrow{h}}x_{t} + W_{\overrightarrow{h}}\overrightarrow{h}\overrightarrow{h}\overrightarrow{h}_{t-1} + b_{\overrightarrow{h}}\right)$$

$$\overleftarrow{h}_{t} = \mathcal{H}\left(W_{x\overleftarrow{h}}x_{t} + W_{\overleftarrow{h}}\overleftarrow{h}\overleftarrow{h}_{t+1} + b_{\overleftarrow{h}}\right)$$

$$y_{t} = W_{\overrightarrow{h}}\overrightarrow{h}_{y}\overrightarrow{h}_{t} + W_{\overleftarrow{h}}\overrightarrow{h}_{y}\overleftarrow{h}_{t} + b_{y}$$



## 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?

A. 1-10

B. 11-20

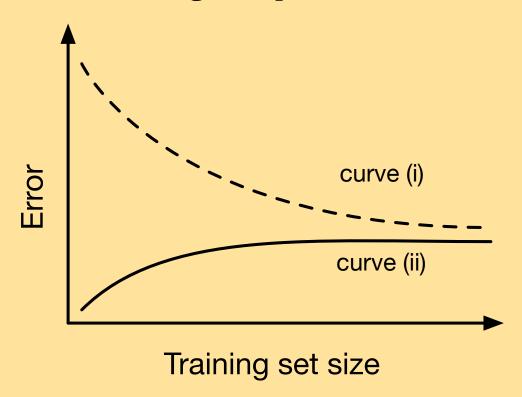
C. 21-30



#### 2.1 True Errors

(b) [4 pts.] **T** or **F**: Learning theory allows us to determine with 100% certainty the true error of a hypothesis to within any  $\epsilon > 0$  error.

#### 2.2 Training Sample Size



- (a) [8 pts.] Which curve represents the training error? Please provide 1–2 sentences of justification.
- (b) [4 pt.] In one word, what does the gap between the two curves represent?

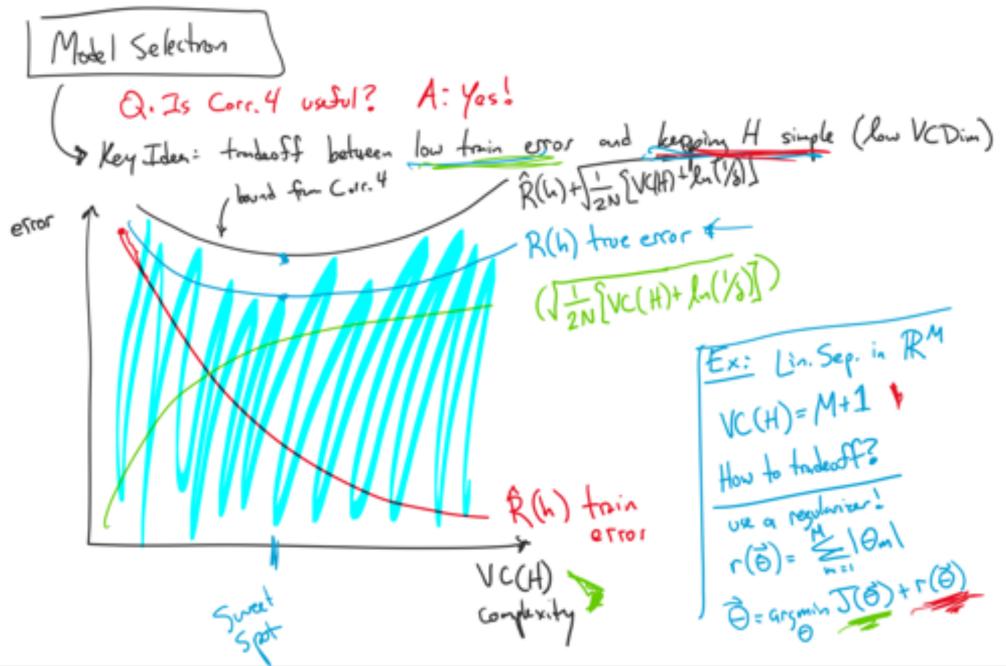
#### 5 Learning Theory [20 pts.]

(a) [3 pts.] **T** or **F**: It is possible to label 4 points in  $\mathbb{R}^2$  in all possible  $2^4$  ways via linear separators in  $\mathbb{R}^2$ .

(d) [3 pts.] **T** or **F**: The VC dimension of a concept class with infinite size is also infinite.

(f) [3 pts.] **T** or **F**: Given a realizable concept class and a set of training instances, a consistent learner will output a concept that achieves 0 error on the training instances.

### PAC Learning & Regularization



#### MLE vs. MAP

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.  $\frac{N}{N}$ 

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

Maximum Likelihood Estimate (MLE)

#### Principle of Maximum a posteriori (MAP) Estimation:

Choose the parameters that maximize the posterior of the parameters given the data.

Prior

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

Maximum a posteriori (MAP) estimate

#### 1.2 Maximum Likelihood Estimation (MLE)

Assume we have a random sample that is Bernoulli distributed  $X_1, \ldots, X_n \sim \text{Bernoulli}(\theta)$ . We are going to derive the MLE for  $\theta$ . Recall that a Bernoulli random variable X takes values in  $\{0,1\}$  and has probability mass function given by

$$P(X;\theta) = \theta^X (1-\theta)^{1-X}.$$

(a) [2 pts.] Derive the likelihood,  $L(\theta; X_1, \ldots, X_n)$ .

(c) **Extra Credit:** [2 pts.] Derive the following formula for the MLE:  $\hat{\theta} = \frac{1}{n} \left( \sum_{i=1}^{n} X_i \right)$ .

#### 1.3 MAP vs MLE

Answer each question with **T** or **F** and **provide a one sentence explanation of your answer:** 

(a) [2 pts.] **T or F:** In the limit, as n (the number of samples) increases, the MAP and MLE estimates become the same.

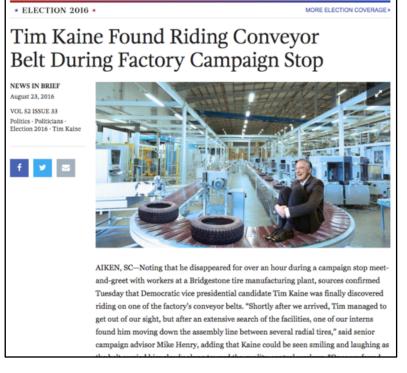
### Fake News Detector

**Today's Goal:** To define a generative model of emails of two different classes (e.g. real vs. fake news)

#### The Economist



#### The Onion



### Model 1: Bernoulli Naïve Bayes

Flip weighted coin



 $x_2$ 

 $\chi_3$ 

 $x_M$ 

 $\mathcal{Y}$ 

 $x_1$ 

If HEADS, flip each red coin



If TAILS, flip each blue coin



We can **generate** data in this fashion. Though in practice we never would since our data is **given**.

Instead, this provides an explanation of **how** the data was generated (albeit a terrible one).

Each red coin corresponds to an  $x_m$ 

#### 1.1 Naive Bayes

You are given a data set of 10,000 students with their sex, height, and hair color. You are trying to build a classifier to predict the sex of a student, so you randomly split the data into a training set and a testing set. Here are the specifications of the data set:

- $sex \in \{male, female\}$
- height  $\in [0,300]$  centimeters
- hair  $\in$  {brown, black, blond, red, green}
- 3240 men in the data set
- 6760 women in the data set

Under the assumptions necessary for Naive Bayes (not the distributional assumptions you might naturally or intuitively make about the dataset) answer each question with **T** or **F** and **provide a one sentence explanation of your answer**:

(a) [2 pts.] **T** or **F**: As height is a continuous valued variable, Naive Bayes is not appropriate since it cannot handle continuous valued variables.

(c) [2 pts.] **T** or **F**: P(height|sex,hair) = P(height|sex).

Material Covered After Midterm Exam 2

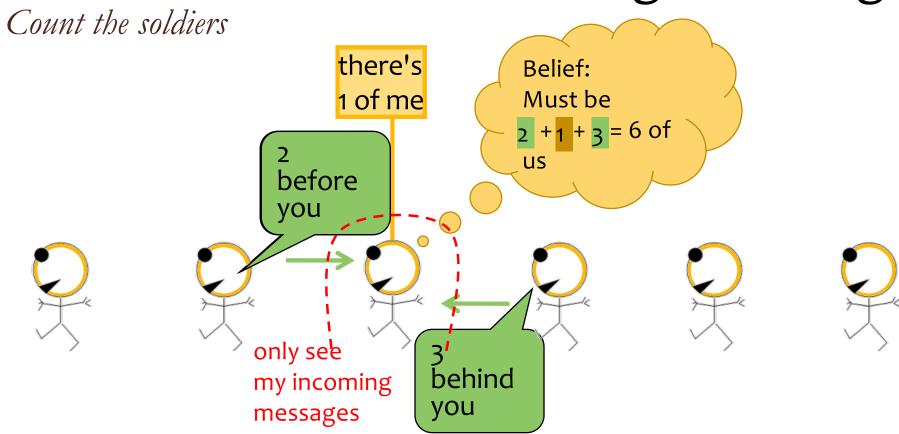
### **SAMPLE QUESTIONS**

### Totoro's Tunnel

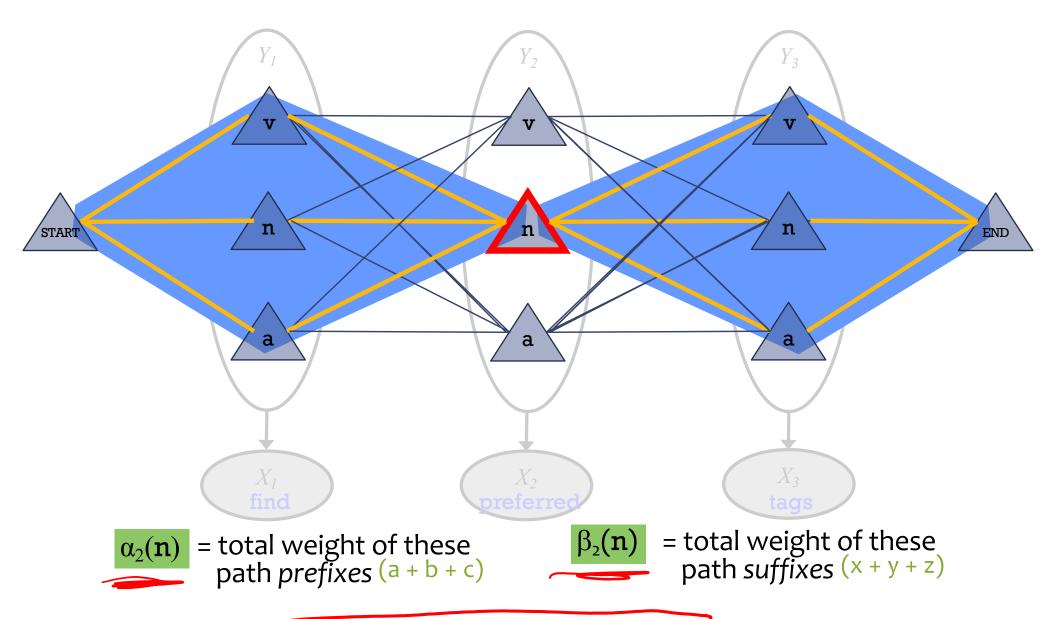




Great Ideas in ML: Message Passing



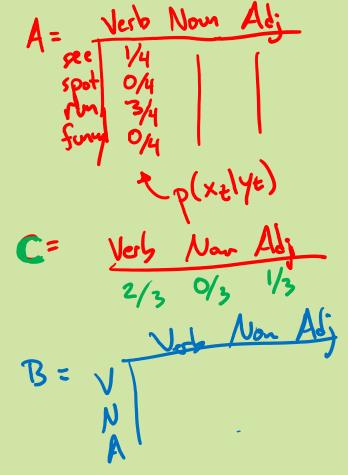
### Forward-Backward Algorithm: Finds Marginals

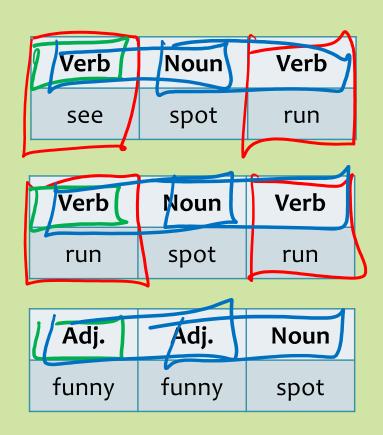


Product gives ax+ay+az+bx+by+bz+cx+cy+cz = total weight of paths

#### 4 Hidden Markov Models

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





#### 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? 

45 POS 

46 POS 

47 POS 

48 P

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: Ryan Reynolds' Voicemail



### 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.,  $R, S, E, A \in \{0, 1\}$ .

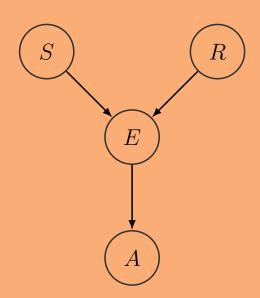


Figure 5: Directed graphical model for problem 5.

# Poll. in Scourse 97

### Sample Questions

(b) [2 nts] How many parameters

(b) [2 pts.] How many parameters, i.e., entries in the CPT tables, 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.,  $R, S, E, A \in \{0, 1\}$ .

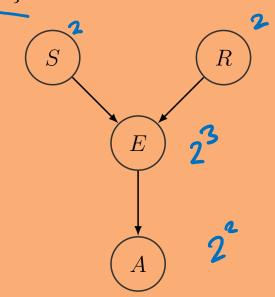


Figure 5: Directed graphical model for problem 5.

(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.

#### 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.,  $R, S, E, A \in \{0, 1\}$ .

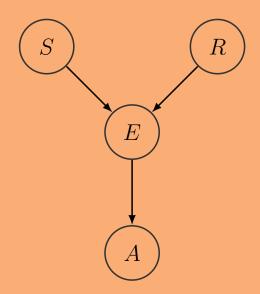
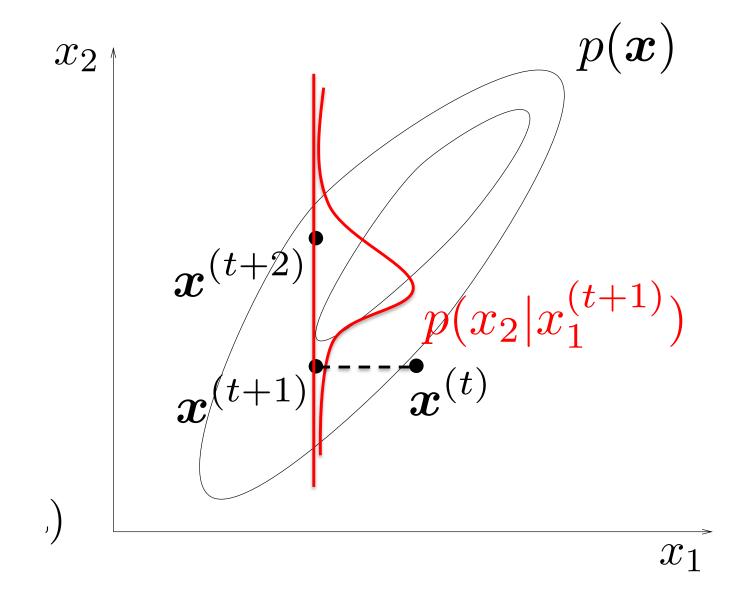


Figure 5: Directed graphical model for problem 5.

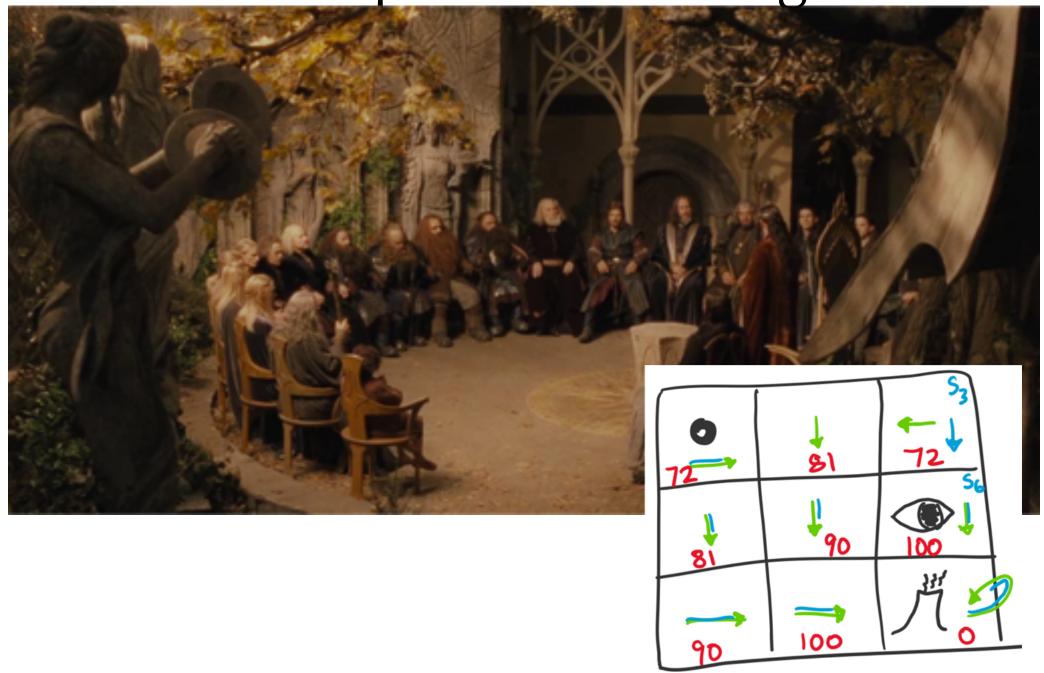
#### 5 Graphical Models

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

## Gibbs Sampling



Example: Path Planning

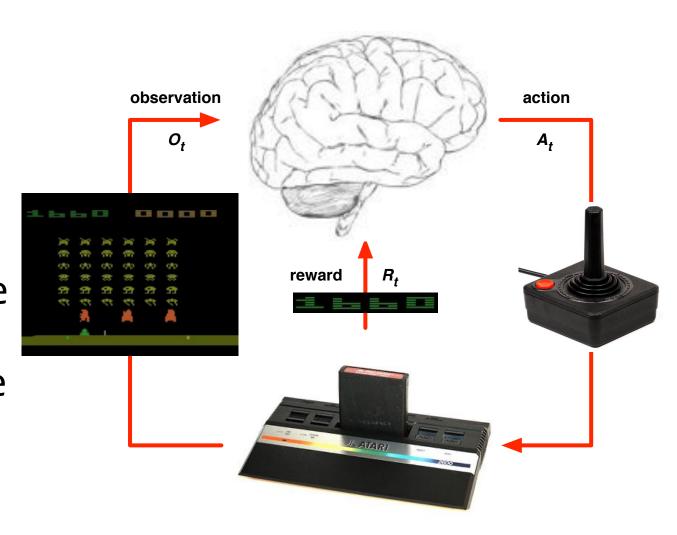


# Today's lecture is brought you by the letter....

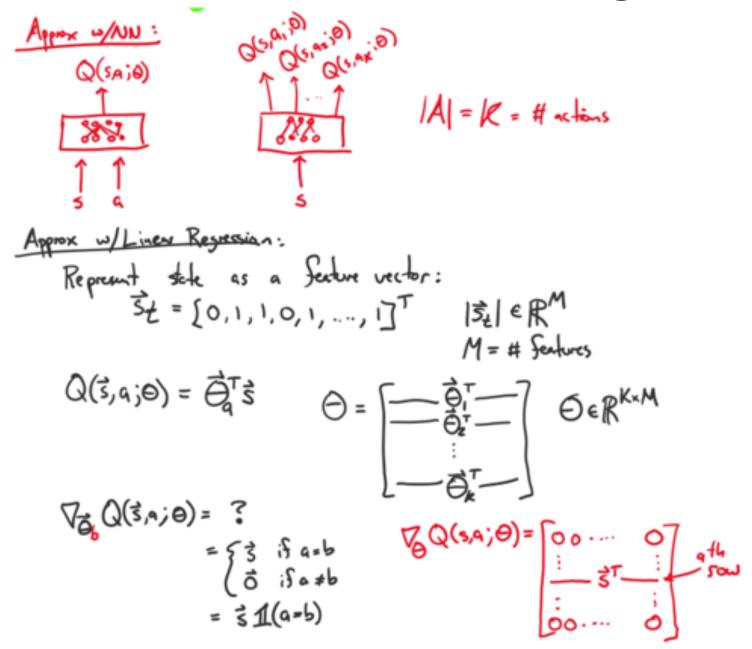


## Playing Atari with Deep RL

- Setup: RL system observes the pixels on the screen
- It receives rewards as the game score
- Actions decide how to move the joystick / buttons



## not-so-Deep Q-Learning



#### 7.1 Reinforcement Learning

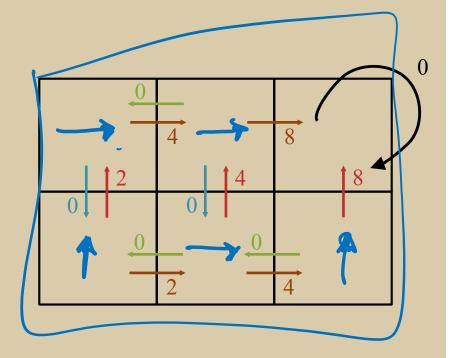
A = calamity

- 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.
    - 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.

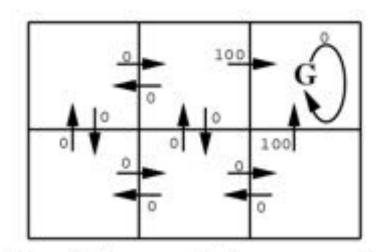
- 4. (1 point) **True or False:** Value iteration is better at balancing exploration and exploitation compared with policy iteration.
  - True
  - ( False

#### 7.1 Reinforcement Learning

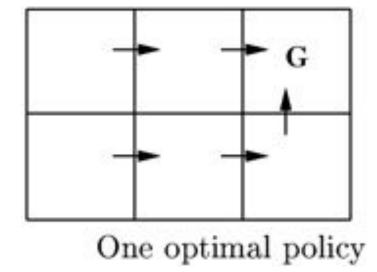
- 1. For the R(s,a) values shown on the arrows below, what is the corresponding optimal policy? Assume the discount factor is 0.1
- 2. For the R(s,a) values shown on the arrows below, which are the corresponding  $V^*(s)$  values? Assume the discount factor is 0.1
- 3. For the R(s,a) values shown on the arrows below, which are the corresponding  $Q^*(s,a)$  values? Assume the discount factor is 0.1

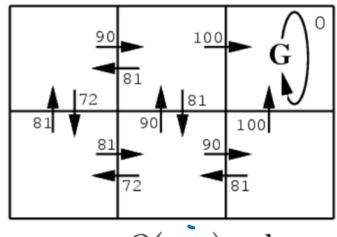


#### Example: Robot Localization

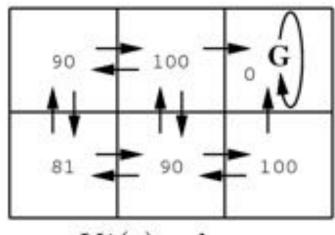


r(s, a) (immediate reward) values





Q(s,a) values



 $V^*(s)$  values

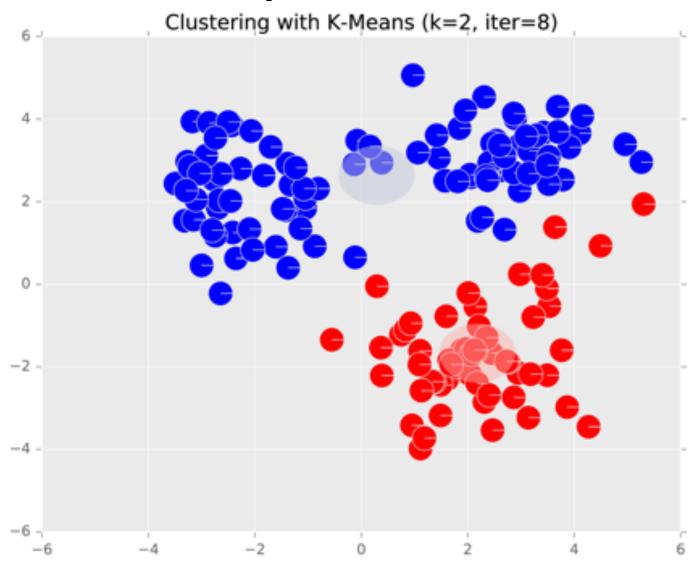
## K-Means Example: A Real-World Dataset



## Example: K-Means



## Example: K-Means

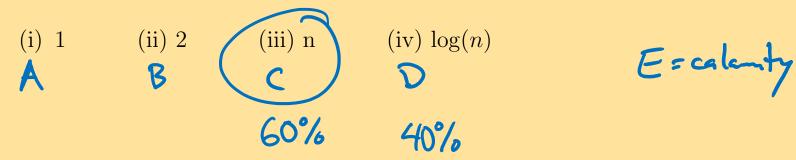


#### Q3

#### Samples Questions

#### 2 K-Means Clustering

(a) [3 pts] We are given n data points,  $x_1, ..., x_n$  and asked to cluster them using K-means. If we choose the value for k to optimize the objective function how many clusters will be used (i.e. what value of k will we choose)? **No justification required.** 



## aka. K-Means Samples Questions

#### 2.2 Lloyd's algorithm

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

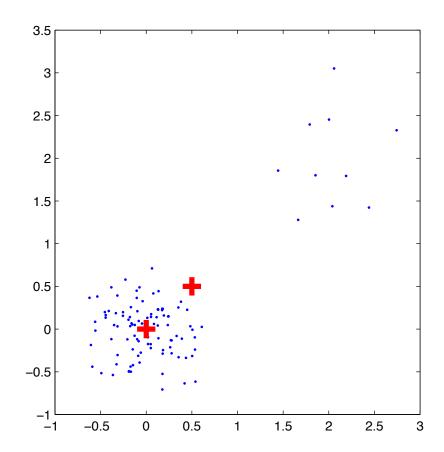


Figure 2: Initial data and cluster centers

#### 2.2 Lloyd's algorithm

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

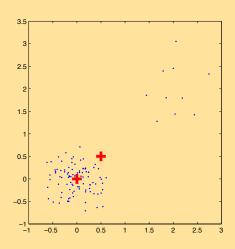
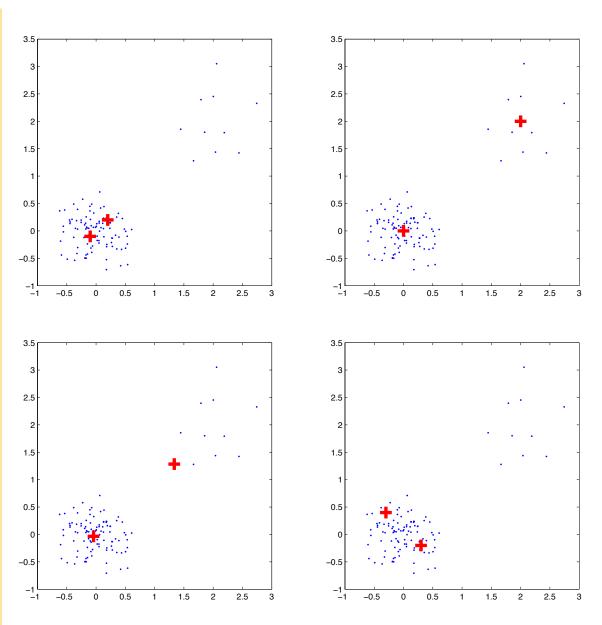


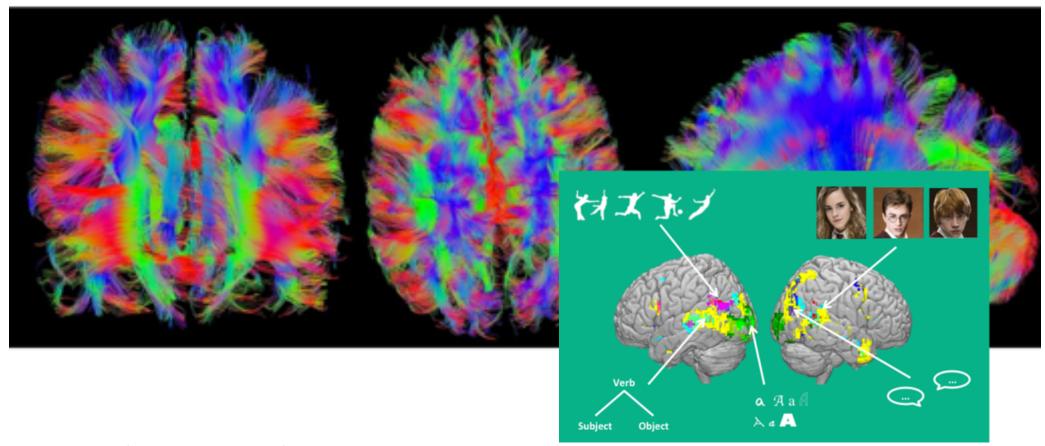
Figure 2: Initial data and cluster centers



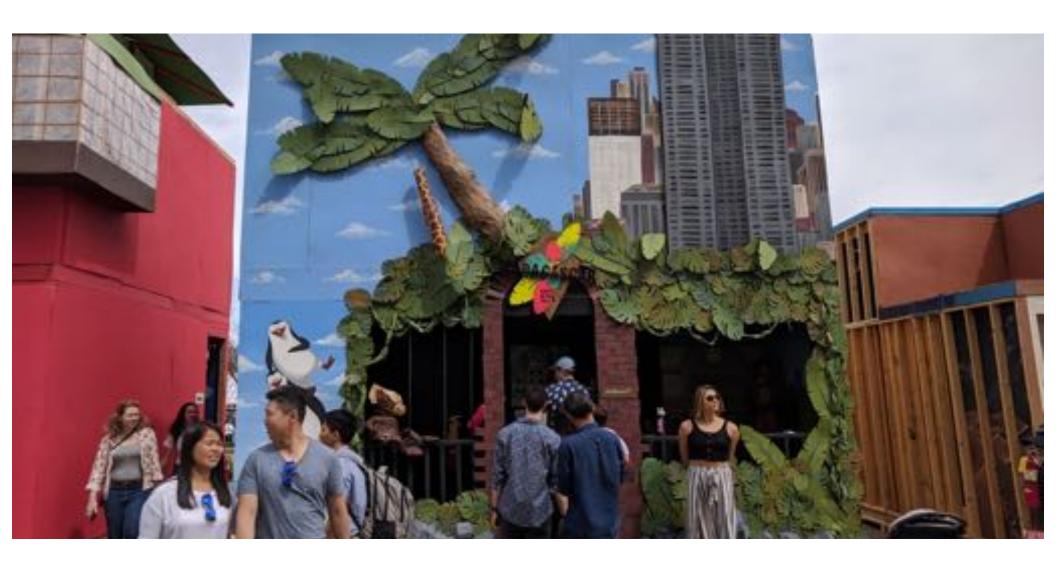
## High Dimension Data

#### Examples of high dimensional data:

- Brain Imaging Data (100s of MBs per scan)



## Shortcut Example

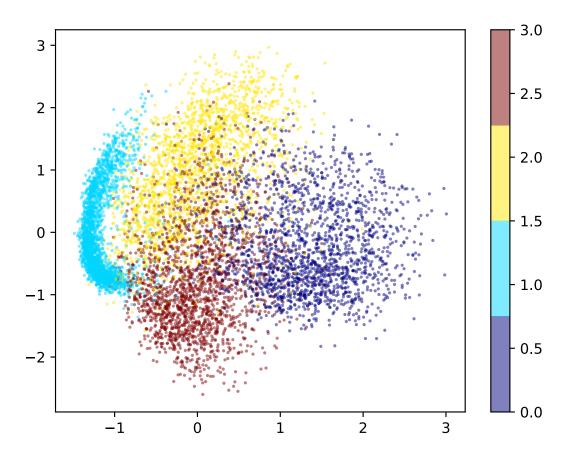


https://www.youtube.com/watch?v=MlJN9pEfPfE

## Projecting MNIST digits

#### **Task Setting:**

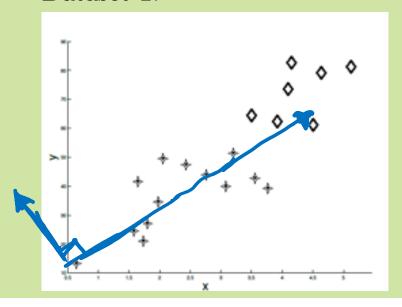
- 1. Take 25x25 images of digits and project them down to 2 components
- 2. Plot the 2 dimensional points



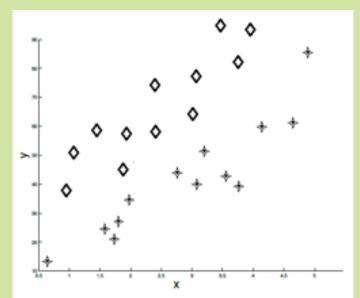
#### 4 Principal Component Analysis [16 pts.]

- (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:



#### 4 Principal Component Analysis

(c) [2 pts.] Assume we apply PCA to a matrix  $X \in \mathbb{R}^{n \times m}$  and obtain a set of PCA features,  $Z \in \mathbb{R}^{m \times n}$ . We divide this set into two,  $Z_1$  and  $Z_2$ . The first set,  $Z_1$ , corresponds to the top principal components. The second set,  $Z_2$ , corresponds to the remaining principal components. Which is more common in the training data:

A: a point with large feature values in  $Z_1$  and small feature values in  $Z_2$ 

B: a point with large feature values in  $Z_2$  and small feature values in  $Z_1$ 

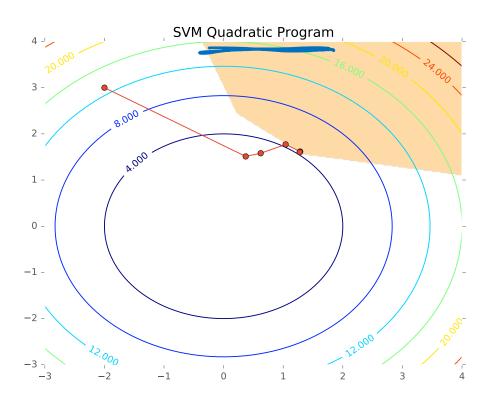
#### 4 Principal Component Analysis

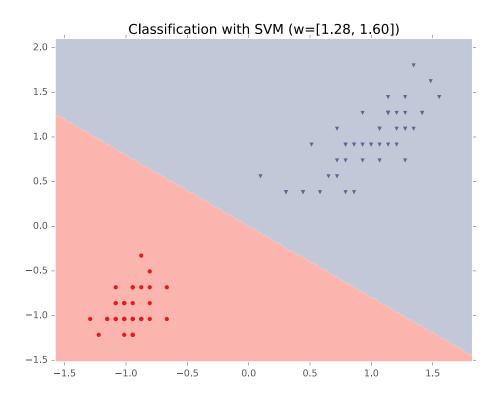
- (i) **T** or **F** The goal of PCA is to interpret the underlying structure of the data in terms of the principal components that are best at predicting the output variable.
- (ii) **T** or **F** The output of PCA is a new representation of the data that is always of lower dimensionality than the original feature representation.
- (iii) **T** or **F** Subsequent principal components are always orthogonal to each other.

## SVM Example:



## SVM QP





## Soft-Margin SVM

#### Hard-margin SVM (Primal)

$$\begin{aligned} & \min_{\mathbf{w},b} \ \frac{1}{2} \|\mathbf{w}\|_2^2 \\ & \text{s.t. } y^{(i)}(\mathbf{w}^T\mathbf{x}^{(i)}+b) \geq 1, \quad \forall i=1,\dots,N \end{aligned}$$

#### Hard-margin SVM (Lagrangian Dual)

$$\max_{\boldsymbol{\alpha}} \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y^{(i)} y^{(j)} \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}$$
s.t.  $\alpha_i \ge 0$ ,  $\forall i = 1, \dots, N$ 

$$\sum_{i=1}^{N} \alpha_i y^{(i)} = 0$$

#### Soft-margin SVM (Primal)

$$\min_{\mathbf{w},b} \frac{1}{2} \|\mathbf{w}\|_2^2 + C\left(\sum_{i=1}^N e_i\right)$$
s.t.  $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \ge 1 - e_i, \quad \forall i = 1, \dots, N$ 

$$e_i \ge 0, \quad \forall i = 1, \dots, N$$

#### Soft-margin SVM (Lagrangian Dual)

$$\max_{\boldsymbol{\alpha}} \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y^{(i)} y^{(j)} \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}$$
s.t.  $0 \le \alpha_i \le C, \quad \forall i = 1, \dots, N$ 

$$\sum_{i=1}^{N} \alpha_i y^{(i)} = 0$$

- (c) [4 pts.] **Extra Credit:** Consider the dataset in Fig. 4. Under the SVM formulation in section 4.2(a),
  - (1) Draw the decision boundary on the graph.
  - (2) What is the size of the margin?
  - (3) Circle all the support vectors on the graph.

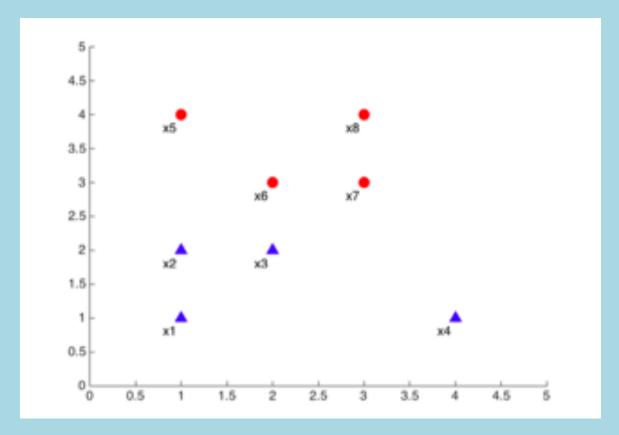
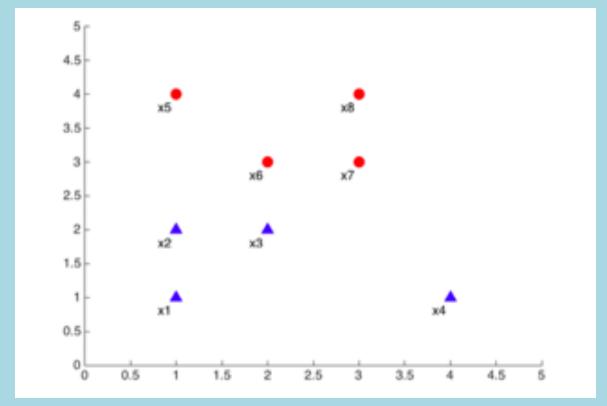


Figure 4: SVM toy dataset

#### 4.2 Multiple Choice

- (a) [3 pt.] If the data is linearly separable, SVM minimizes  $||w||^2$  subject to the constraints  $\forall i, y_i w \cdot x_i \geq 1$ . In the linearly separable case, which of the following may happen to the decision boundary if one of the training samples is removed? Circle all that apply.
  - Shifts toward the point removed
  - Shifts away from the point removed
  - Does not change



3. [Extra Credit: 3 pts.] One formulation of soft-margin SVM optimization problem is:

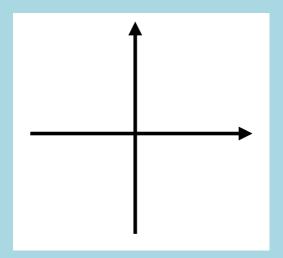
$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_{2}^{2} + C \sum_{i=1}^{N} \xi_{i}$$
s.t.  $y_{i}(\mathbf{w}^{\top} x_{i}) \geq 1 - \xi_{i} \quad \forall i = 1, ..., N$ 

$$\xi_{i} \geq 0 \quad \forall i = 1, ..., N$$

$$C \geq 0$$

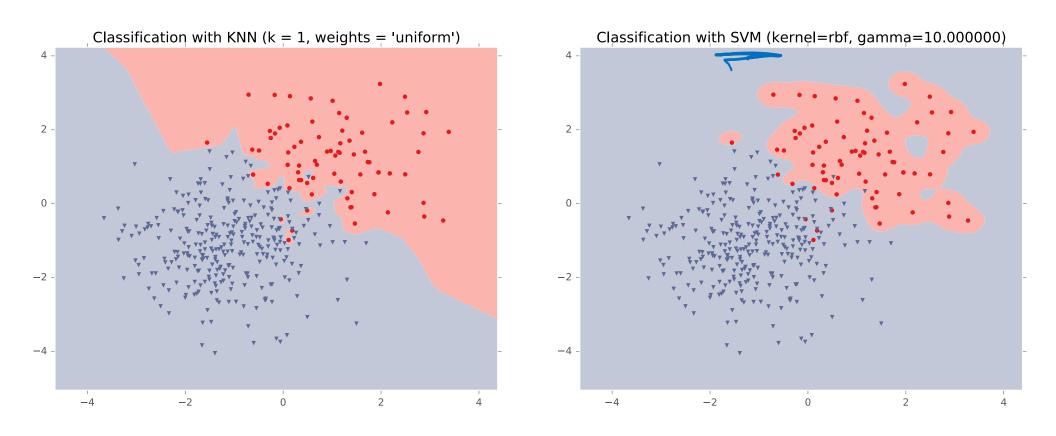
where  $(x_i, y_i)$  are training samples and w defines a linear decision boundary.

Derive a formula for  $\xi_i$  when the objective function achieves its minimum (No steps necessary). Note it is a function of  $y_i \mathbf{w}^\top x_i$ . Sketch a plot of  $\xi_i$  with  $y_i \mathbf{w}^\top x_i$  on the x-axis and value of  $\xi_i$  on the y-axis. What is the name of this function?



#### RBF Kernel Example

#### KNN vs. SVM



RBF Kernel: 
$$K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma ||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$$

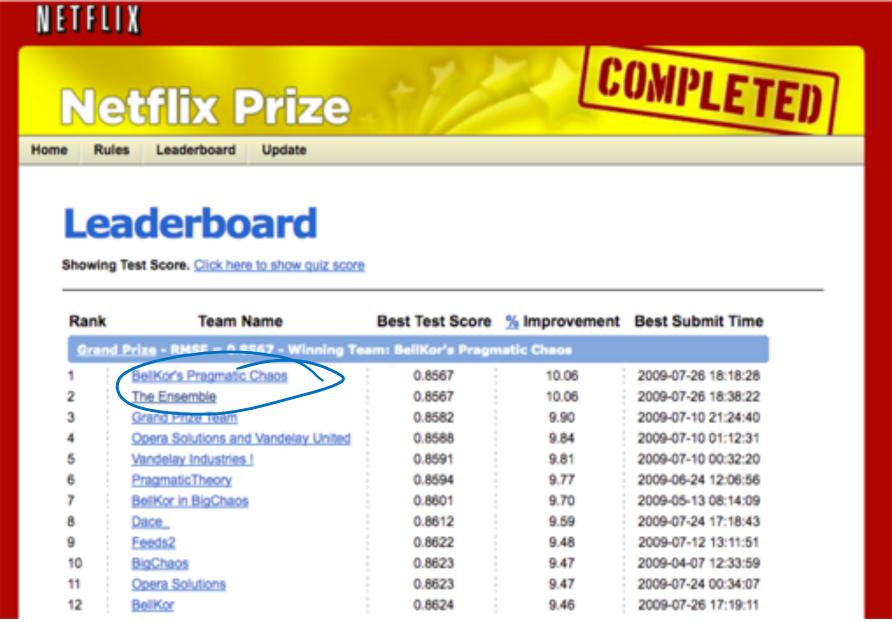
#### 4.3 Analysis

(a) [4 pts.] In one or two sentences, describe the benefit of using the Kernel trick.

(b) [4 pt.] The concept of margin is essential in both SVM and Perceptron. Describe why a large margin separator is desirable for classification.

(e) [2 pts.] **T** or **F**: The function  $K(\mathbf{x}, \mathbf{z}) = -2\mathbf{x}^T\mathbf{z}$  is a valid kernel function.

#### Recommender Systems



Weighted Majority Algorithm

(Littlestone & Warmuth, 1994)

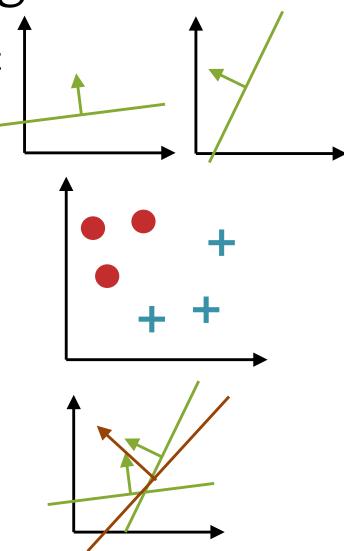
Given: pool A of binary classifiers (that you know nothing about)

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

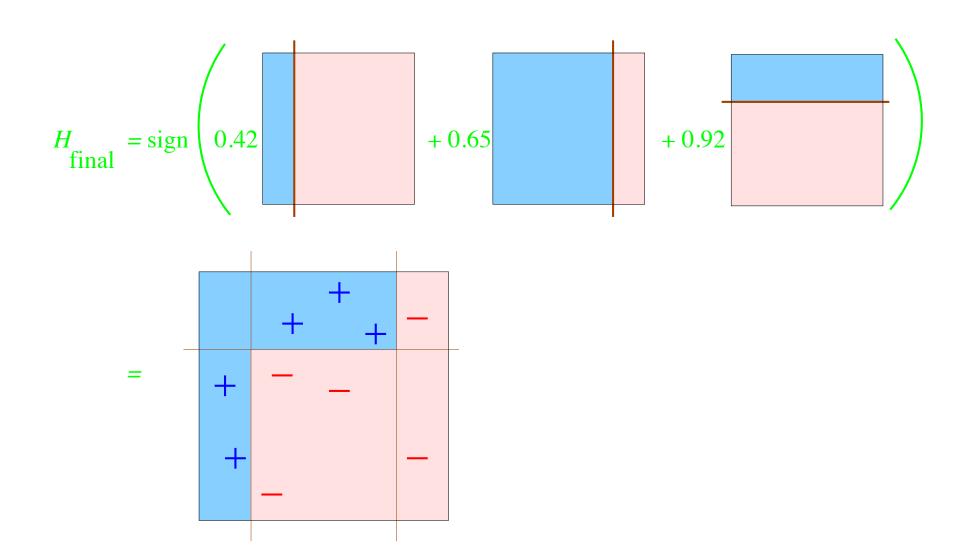
 Goal: design a new learner that uses the predictions of the pool to make new predictions

#### 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  $\beta$



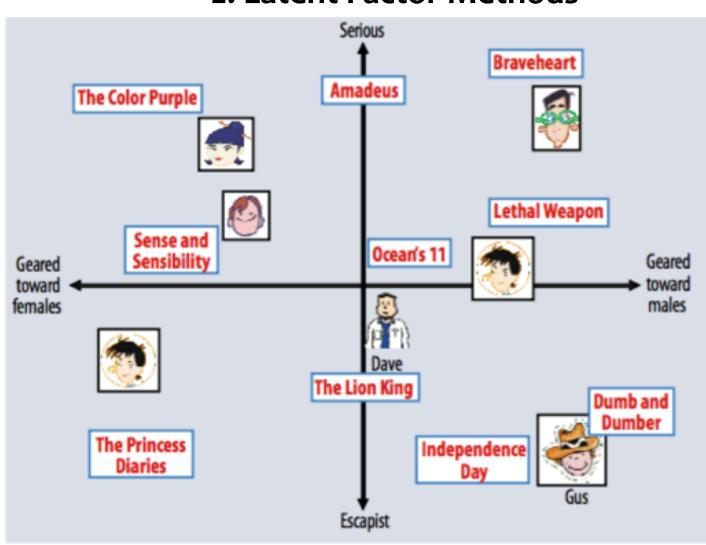
## AdaBoost: Toy Example



## 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



## Crowdsourcing Exam Questions

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

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

#### **Answer Here:**

The Big Picture

#### **MACHINE LEARNING**

| Paradigm                                | Data  |   |
|---|---|---|
| Supervised                              | $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N}$ | $\mathbf{x} \sim p^*(\cdot)$ 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                          |   |

| 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                            | $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N \qquad \mathbf{x} \sim p^*(\cdot)$                                      |

| 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                         | $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N_1} \cup \{\mathbf{x}^{(j)}\}_{j=1}^{N_2}$                     |

| 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                                  | $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), (\mathbf{x}^{(3)}, y^{(3)}), \ldots\}$   |

| 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                         |

| 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                      | $\mathcal{D} = \{(s^{(1)}, a^{(1)}), (s^{(2)}, a^{(2)}), \ldots\}$  |

| 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                  | $\mathcal{D} = \{(s^{(1)}, a^{(1)}, r^{(1)}), (s^{(2)}, a^{(2)}, r^{(2)}), \ldots\}$                                |

# Machine Learning: The Big Picture

#### Whiteboard

- Decision Rules / Models (probabilistic generative, probabilistic discriminative, perceptron, SVM, regression, MDP, graphical models)
- Objective Functions (likelihood, conditional likelihood, hinge loss, mean squared error)
- Regularization (L1, L2, priors for MAP)
- Update Rules (SGD, perceptron)
- Nonlinear Features (preprocessing, kernel trick)

## 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.

#### Big Ideas in ML:

Which are the ideas driving development of the field?

- inductive bias
- generalization / overfitting
- bias-variance decomposition

**Application Area** 

- generative vs. discriminative
- deep nets, graphical models
- PAC learning
- distant rewards

### Facets of Building ML Systems:

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

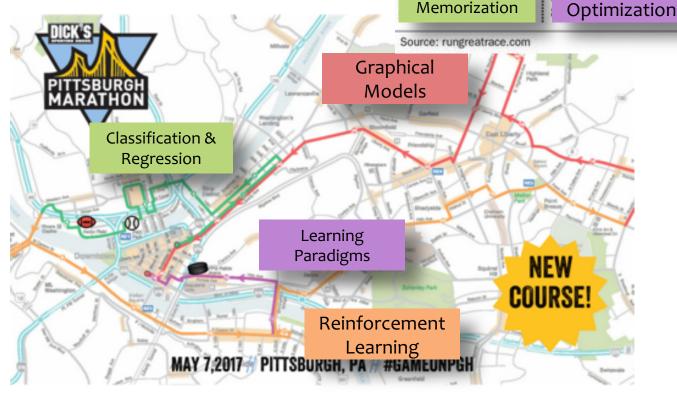
- 1. Data prep
- 2. Model selection
- 3. Training (optimization / search)
- 4. Hyperparameter tuning on validation data
- 5. (Blind) Assessment on test data

## A new **combined** course...

... with the best (uphill climbs) from both

#### Great Race: route and street closing schedule





Post-Gazette

## 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
- 7. 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.

# Q&A