### CIS520: Machine Learning Lyle Ungar



Poll Everywhere Poll Everywhere, Inc. Communication

E Everyone

This app is compatible with all of your devices.

Install *Poll Everywhere* from app store or go to https://pollev.com/lyleungar251 What's your favorite word?

happy

Start the presentation to see live content. Still no live content? Install the app or get help at PollEv.com/app

### Administrivia

#### ♦ Always look at the wiki lectures page

- Readings, worksheets, quiz, survey
- Remember the resources page on the course wiki
- ◆ Ed you should be receiving announcements!!!
- Waitlist done
- ♦ HW0, HW1, worksheets...
- ◆ Office hours: wiki: "people/office hours"
- Pods sign up!!!

### **Course Cadence**

- ♦ MW: new material in lectures
- ◆ F: review material DRLB A2
- ♦ WTh Pods in person
  - Meet people; think about material
- Th-T Worksheets, HW, quizzes
  - Do everything by midnight Tuesday

### **Overview of ML**



#### Rashidi, Betts & Green 2019

### Nonparametric Learning Lyle Ungar

K-NN Norms, Distance Overfitting and Model Complexity **Decision Trees** Entropy, Information gain

### k-Nearest Neighbors (kNN)

#### • To predict y at a point x

- Find the k nearest neighbors
- $\hat{y}(\mathbf{x})$  = the majority label or

the average of the y's of those points

https://www.youtube.com/watch?v=PB4qATziTlQ

### **Norms and Distances**

### Norms

#### For all $a \in R$ and all $u, v \in V$ ,

- $L_p(av) = |a| L_p(v)$
- $L_{\rho}(\boldsymbol{u} + \boldsymbol{v}) \leq L_{\rho}(\boldsymbol{u}) + L_{\rho}(\boldsymbol{v})$ 
  - triangle inequality or subadditivity
- If  $L_{\rho}(\mathbf{v}) = 0$  then  $\mathbf{v}$  is the zero vector
  - implies //v// = 0 iff v is the zero vector

### **L**<sub>p</sub> **norm**, $||\mathbf{x}||_{p}$ : $(\Sigma_{j} |\mathbf{x}_{j}|^{p})^{1/p}$

What is

||(1,2,3)||<sub>1</sub> ? A) 1 B) 3 C) sqrt(14) D) sqrt(14/3) E) none of the above



What is

||(1,2,3)||<sub>2</sub> ? A) 1 B) 3 C) sqrt(14) D) sqrt(14/3) E) none of the above



## What is

||(1,2,3)||<sub>1/2</sub> ? A) 1 B) 3 C) sqrt(14) D) sqrt(14/3) E) none of the above



### L<sub>0</sub> pseudo-norm

 $||\mathbf{x}||_0$  = number of elements  $x_j \neq 0$ 

How is this not a real norm?

What is  $||(1,2,3)||_0$ ?

A) 1 B) 3

C) sqrt(14)

- D) sqrt(14/3)
- E) none of the above



### Distance

#### • Every norm generates a distance

$$d_{p}(\mathbf{x},\mathbf{y}) = ||\mathbf{x}-\mathbf{y}||_{p}$$

### **Distance function (metric)**

- 1.  $d(x, y) \ge 0$  (*non-negativity*, or separation axiom)
- 2. d(x, y) = 0 if and only if x = y (coincidence axiom)
- 3. d(x, y) = d(y, x) (*symmetry*)
- 4.  $d(x, z) \le d(x, y) + d(y, z)$  (subadditivity / triangle inequality).

https://en.wikipedia.org/wiki/Metric\_(mathematics)

### Lines of equal distance from (0,0)









*L*<sub>1</sub>*norm* 

*L*<sub>inf</sub>*norm* 

Convexity Is  $||\mathbf{x}||_{1/2}$  convex? Concave Convex

Image credit: https://writingexplained.org/concave-vs-convex-difference

A figure is convex if any line segment connecting two points on the surface of the figure lies entirely inside the figure; otherwise it is concave



# Different norms give different decision boundaries







L<sub>2</sub>

L

### **Components of ML - kNN**

Representation: nonparametric

•  $\hat{y} = f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^{T} \mathbf{x}$ 

Loss function (with L2 distance)

•  $L(\mathbf{y}, \, \hat{\mathbf{y}}) = ||\mathbf{y} - \hat{\mathbf{y}}||_2$ 

Optimization method: not required

- $\operatorname{argmin}_{w} L(y, \hat{y}(w))$
- gradient descent

### How to pick k?

• What loss function are we trying to minimize?  $\|y - \hat{y}(x)\|_{p}$ 

### Linear regression on 3 data sets



### 1-NN on 3 data sets (L<sub>1</sub>)









### 9-NN on 3 data sets $(L_1)$







# In high dimensions most points are equally close to each other.

#### Consider a 100-dimensional cube.

• A vertex represented is a "one hot encoding" or "indicator" function, a vector with 99 zeros and one 1.

#### • What is the distance between any two vertices?

• 0,1,2, more, it varies

#### • Generate points at random with half 0's and half 1's.

• How far away (on average) are two such points?

Half the coordinates will be the same, so sqrt(50)

## **Decision Trees** and Information Theory

Lyle Ungar University of Pennsylvania

### **Decision Trees** Recursive partition trees, ID3, C4.5, CART, CHAID



https://www.nytimes.com/interactive/2019/08/08/op inion/sunday/party-polarization-quiz.html

### What symptom tells you most about the disease? S1 S2 S3 D A) S1

- y n n y n y y y
- n y n n
- n n n n
- y y n y

A) S1 B) S2 C) S3

Why?



#### What symptom tells you most about the disease? **S2/D S1/D** s3/D У n y y n n **y** 2 1 **y** 1 2 0 0 У A) S1 B) S2 **n** 2 2 **n** 1 **n** 1 2 1 **S**3

Why?

#### If you know S1=n, what symptom tells you most about the disease? S1 S2 **S**3 D **A) S1** B) S2 C) S3 n n y y У У У n Why? n У n n A, B, or C? n n n n AA

BB

c **c** 

tion to see live content. Still on live content? Install the app or get bein at Polic

y y n y

### **Resulting decision tree**

S1 y/ \n D S3 y/ \n D ~ D

The key question: what criterion to use do decide which question to ask?

## Entropy and Information Gain

### Andrew W. Moore Carnegie Mellon University

www.cs.cmu.edu/~awm awm@cs.cmu.edu

412-268-7599

modified by Lyle Ungar

### **Bits**

You observe a set of independent random samples of X

#### You see that X has four possible values

P(X=A) = 1/4	P(X=B) = 1/4	P(X=C) = 1/4	P(X=D) = 1/4

So you might see: BAACBADCDADDDA...

You transmit data over a binary serial link. You can encode each reading with two bits (e.g. A = 00, B = 01, C = 10, D = 11) 010000100100111011001111100...

### **Fewer Bits**

#### Someone tells you that the probabilities are not equal

P(X=A) = 1/2P(X=B) = 1/4P(X=C) = 1/8P(X=D) = 1/8It is possible to invent a coding for your transmission that only<br/>uses 1.75 bits on average per symbol. How?

А	0
В	10
С	110
D	111

(This is just one of several ways)

### **Fewer Bits**

#### Suppose there are three equally likely values...

P(X=A) = 1/3 P(X=B) = 1/3 P(X=C) = 1/3

Here's a naïve coding, costing 2 bits per symbol



Can you think of a coding that only needs 1.6 bits per symbol on average?

In theory, it can in fact be done with 1.58496 bits per symbol.

General Case: EntropySuppose X can have one of m values...  $V_1 V_2 V_m$  $P(X=V_1) = p_1$  $P(X=V_2) = p_2$ .... $P(X=V_m) = p_m$ 

What's the smallest possible number of bits, on average, per symbol, needed to transmit a stream of symbols drawn from X's distribution?

It is

$$H(X) = -p_1 \log_2 p_1 - p_2 \log_2 p_2 - \dots - p_m \log_2 p_m$$
$$= -\sum_{j=1}^m p_j \log_2 p_j$$



H(X) = The entropy of X

• "High Entropy" means X is from a uniform (boring) distribution

• "Low Entropy" means X is from varied (peaks and valleys) distribution Copyright © 2001, 2003, Andrew W. Moore



### **Entropy in a nut-shell**



Low Entropy



**High Entropy** 

### **Entropy in a nut-shell**



### Why does entropy have this form?

$$H(X) = -p_1 \log_2 p_1 - p_2 \log_2 p_2 - \dots - p_m \log_2 p_m$$
$$= -\sum_{j=1}^m p_j \log_2 p_j$$

Entropy is the expected value of the information content (surprise) of the message  $log_2p_i$ 

If an event is certain, the entropy is A) 0 B) between 0 and  $\frac{1}{2}$ C)  $\frac{1}{2}$ D) between  $\frac{1}{2}$  and 1 E) 1



### Why does entropy have this form?

$$H(X) = -p_1 \log_2 p_1 - p_2 \log_2 p_2 - \dots - p_m \log_2 p_m$$
$$= -\sum_{j=1}^m p_j \log_2 p_j$$

If two events are equally likely, the entropy is

A) 0 B) between 0 and  $\frac{1}{2}$ C)  $\frac{1}{2}$ D) between  $\frac{1}{2}$  and 1 E) 1

•		A, B, C, D or E	
A	A		
В	В		
С	с		
D	D		
E	E		
e.		Start the presentation to see live content. Still no live content! Install the app or get help at PollEv.com/app	

### Specific Conditional Entropy H(Y|X=v)

#### Suppose I'm trying to predict output Y and I have input X

- X = College Major
- Y = Likes "Gladiator"

X	Y
Math	Yes
History	Νο
CS	Yes
Math	Νο
Math	Νο
CS	Yes
History	Νο
စ <b>)၊ ဆုံးဂြ</b> င်္ 2001, 20(	3, <b>YARS</b> rew W. Moore

Assume this reflects the true probabilities

#### e.g. From this data we estimate

- *P*(*LikeG* = Yes) = 0.5
- *P*(*Major* = *Math* & *LikeG* = *No*) = 0.25
- *P*(*Major* = *Math*) = 0.5
- P(LikeG = Yes | Major = History) = 0

#### Note:

• *H*(X) = 1.5 •*H*(Y) = 1

### Specific Conditional Entropy H(Y|X=v)

X = College Major Y = Likes "Gladiator"

X	Y
Math	Yes
History	Νο
CS	Yes
Math	Νο
Math	No
CS	Yes
History	Νο
Math	Yes

**Definition of Specific Conditional Entropy:** H(Y | X = v) = The entropy of Y among only those records in which X has value v

#### **Example:**

- *H*(*Y*|*X*=*Math*) = 1
- H(Y|X=History) = 0
- H(Y|X=CS) = 0

### Conditional Entropy H(Y|X)

X = College Major Y = Likes "Gladiator"

X	Y
Math	Yes
History	Νο
CS	Yes
Math	Νο
Math	No
CS	Yes
History	Νο
Math	Yes

Copyright © 2001, 2003, Andrew W. Moore

#### **Definition of Conditional Entropy:**

H(Y|X) = The average specific conditional entropy of Y

If you choose a record at random what will be the conditional entropy of *Y*, conditioned on that row's value of *X* 

= Expected number of bits to transmit Y if both sides will know the value of X

$$= \Sigma_j \operatorname{Prob}(X=v_j) H(Y \mid X=v_j)$$

Conditional EntropyX = College MajorY = Likes "Gladiator"

### Definition of Conditional Entropy:

X Y Math Yes **History** No CS Yes Math No Math No CS Yes **History** No Math Yes

Copyright © 2001, 2003, Andrew W. Moore

H(Y|X) = The average conditional entropy of Y =  $\Sigma_i Prob(X=v_i) H(Y | X = v_i)$ 

#### **Example:**

$v_j$	Prob(X=v <sub>j</sub> )	$H(Y \mid X = v_j)$
Math	0.5	1
History	0.25	0
CS	0.25	0

H(Y|X) = 0.5 \* 1 + 0.25 \* 0 + 0.25 \* 0 = 0.5

### Information Gain C = College Major Definition of Information Gain:

X = College Major Y = Likes "Gladiator"

#### Y X Math Yes **History** No CS Yes Math No Math No CS Yes **History** No Math Yes

Copyright © 2001, 2003, Andrew W. Moore

IG(Y|X) = I must transmit Y. How many bits on average would it save me if both ends of the line knew X?

IG(Y|X) = H(Y) - H(Y|X)

#### Example:

- H(Y) = 1
- H(Y|X) = 0.5
- Thus IG(Y|X) = 1 0.5 = 0.5

### **Information Gain Example**



### **Another example**



### What is Information Gain used for?

If you are going to collect information from someone (e.g. asking questions sequentially in a decision tree), the "best" question is the one with the highest information gain.

Information gain is useful for model selection

# What question did we not ask (or answer) about decision trees?

### What you should know

- ♦ K-NN
  - hyperparameter k controls model complexity
- Norm, distance
- Convexity
- Entropy, information gain
- The standard decision tree algorithm
  - Recursive partition to maximize information gain



Start the presentation to see live content. Still no live content? Install the app or get help at PollEv.com/app

#### What questions do you have on today's class?

Тор

Start the presentation to see live content. For screen share software, share the entire screen. Get help at pollev.com/app