

# PAC-learning, VC Dimension, *cont.*

Machine Learning – 10701/15781  
Carlos Guestrin  
Carnegie Mellon University

October 24<sup>th</sup>, 2007

©2005-2007 Carlos Guestrin

1

## A simple setting...

- Classification
  - $m$  data points
  - Finite number of possible hypothesis (e.g., dec. trees of depth  $d$ )
- A learner finds a hypothesis  $h$  that is **consistent** with training data
  - Gets zero error in training  $\rightarrow$   $\text{error}_{\text{train}}(h) = 0$
- What is the probability that  $h$  has more than  $\varepsilon$  true error?
  - $\text{error}_{\text{true}}(h) \geq \varepsilon$

©2005-2007 Carlos Guestrin

2

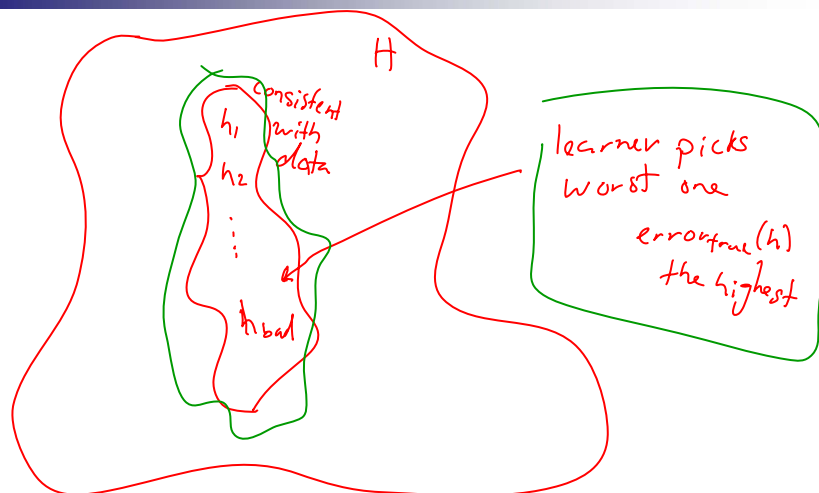
## How likely is a bad hypothesis to get $m$ data points right?

- Hypothesis  $h$  that is **consistent** with training data  $\rightarrow$  got  $m$  i.i.d. points right
  - $h$  "bad" if it gets all this data right, but has high true error
- Prob.  $h$  with  $\text{error}_{\text{true}}(h) \geq \epsilon$  gets one data point right  
 $P(h \text{ gets one point right}) \leq 1 - \epsilon$
- Prob.  $h$  with  $\text{error}_{\text{true}}(h) \geq \epsilon$  gets  $m$  data points right  
 $P(h \text{ gets } m \text{ i.i.d. points right}) \leq (1 - \epsilon)^m$   
exponentially small (as  $m$  increases)

©2005-2007 Carlos Guestrin

3

## But there are many possible hypothesis that are consistent with training data



©2005-2007 Carlos Guestrin

4

## How likely is learner to pick a bad hypothesis

- Prob.  $h$  with  $\text{error}_{\text{true}}(h) \geq \epsilon$  gets  $m$  data points right

$$P(h_{\text{bad}} \text{ consistent with data}) \leq (1-\epsilon)^m$$

- There are  $k$  hypothesis consistent with data

- How likely is learner to pick a bad one?

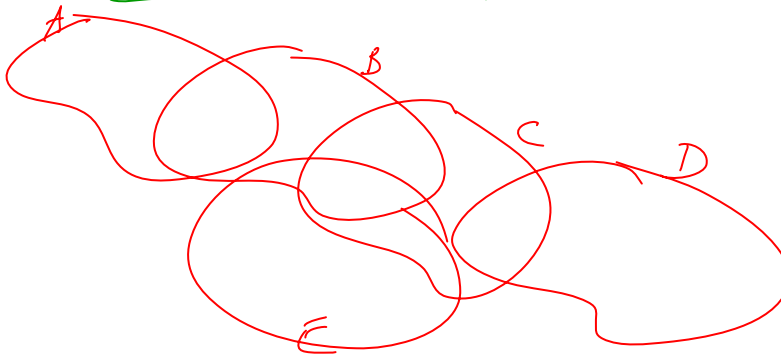
$$P(\exists h \text{ that is bad and consistent with data}) \\ = P(h_1 \text{ bad consistent} \vee h_2 \text{ bad consistent} \vee \dots \vee h_k \text{ bad consistent})$$

©2005-2007 Carlos Guestrin

5

## Union bound

- $P(A \text{ or } B \text{ or } C \text{ or } D \text{ or } \dots) \leq P(A) + P(B) + P(C) + \dots$



©2005-2007 Carlos Guestrin

6

## How likely is learner to pick a bad hypothesis

- Prob.  $h$  with  $\text{error}_{\text{true}}(h) \geq \epsilon$  gets  $m$  data points right

$$P(h_{\text{bad, consistent}}) \leq (1-\epsilon)^m$$

- There are  $k$  hypothesis consistent with data

- How likely is learner to pick a bad one?

$$P(\exists \text{ bad } h \text{ consistent with data}) \leq k (1-\epsilon)^m$$

$$\leq |H| (1-\epsilon)^m$$

$$\leq |H| e^{-m\epsilon}$$

$$(1-\epsilon)^m \leq e^{-m\epsilon}$$

what's  $k$ ?

$$k \leq |H|$$

↑ # of hypotheses

©2005-2007 Carlos Guestrin

7

## Review: Generalization error in finite hypothesis spaces [Haussler '88]

- **Theorem:** Hypothesis space  $H$  finite, dataset  $D$  with  $m$  i.i.d. samples,  $0 < \epsilon < 1$ : <sup>you give me  $\epsilon$</sup>  for any learned hypothesis  $h$  that is consistent on the training data:

$$P(\text{error}_{\text{true}}(h) \geq \epsilon) \leq |H| e^{-m\epsilon}$$



©2005-2007 Carlos Guestrin

8

# Using a PAC bound

Typically, 2 use cases:  $P(\text{error}_{\text{true}}(h) > \epsilon) \leq |H|e^{-m\epsilon}$

□ 1: Pick  $\epsilon$  and  $\delta$ , give you  $m$  *pick  $m, \delta$ , equation holds*

□ 2: Pick  $m$  and  $\delta$ , give you  $\epsilon$  *vs:  $\epsilon \geq \frac{1}{m} (\ln|H| + \ln \frac{1}{\delta})$*

*1,  $\epsilon = 0.1$   $\delta = 0.05$*

$|H|e^{-m\epsilon} \leq \delta$

$\ln(|H|e^{-m\epsilon}) \leq \ln \delta$  *small  $\epsilon$  requires many more points*

$m \geq \frac{1}{\epsilon} (\ln|H| + \ln \frac{1}{\delta})$

*$|H|$  may be big, only depends on  $|H|$   
 $\ln \delta$  not hard to achieve  
1 by  $\delta$*

$\epsilon \geq \frac{1}{m} (\ln|H| + \ln \frac{1}{\delta})$

$\downarrow$  with prob at least  $1 - \delta$

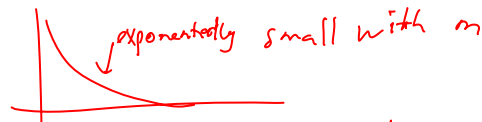
$\text{error}_{\text{true}}(h) \leq \frac{1}{m} (\ln|H| + \ln \frac{1}{\delta})$

©2005-2007 Carlos Guestrin 9

## Review: Generalization error in finite hypothesis spaces [Haussler '88]

- **Theorem:** Hypothesis space  $H$  finite, dataset  $D$  with  $m$  i.i.d. samples,  $0 < \epsilon < 1$ : for any learned hypothesis  $h$  that is consistent on the training data:

$$P(\text{error}_{\text{true}}(h) > \epsilon) \leq |H|e^{-m\epsilon}$$



\* result about  $h$  consistent with train data  
no guarantee your alg. will find such  $h$ .

Even if  $h$  makes zero errors in training data, may make errors in test

## Limitations of Haussler '88 bound

$$P(\text{error}_{\text{true}}(h) > \epsilon) \leq |H|e^{-m\epsilon}$$

- Consistent classifier

- Size of hypothesis space

$$\ln |H|$$

©2005-2007 Carlos Guestrin

11

## What if our classifier does not have zero error on the training data?

- A learner with **zero** training errors may make mistakes in test set
- What about a learner with  $\text{error}_{\text{train}}(h)$  in training set?

no longer assume  $\text{error}_{\text{train}}(h) = 0$

©2005-2007 Carlos Guestrin

12

## Simpler question: What's the expected error of a hypothesis?

- The error of a hypothesis is like estimating the parameter of a coin!

*data*  
  
*true  $\theta$  of coin*

- Chernoff bound: for  $m$  i.i.d. coin flips,  $x_1, \dots, x_m$ , where  $x_i \in \{0, 1\}$ . For  $0 < \epsilon < 1$ :

$$P\left(\theta - \frac{1}{m} \sum_i x_i > \epsilon\right) \leq e^{-2m\epsilon^2}$$

*true parameter* (pointing to  $\theta$ )  
*estimate* (pointing to  $\frac{1}{m} \sum_i x_i$ )

$$\theta_{MLE} = \frac{1}{m} \sum_i x_i$$

©2005-2007 Carlos Guestrin

13

## Using Chernoff bound to estimate error of a single hypothesis

$$P\left(\theta - \frac{1}{m} \sum_i x_i > \epsilon\right) \leq e^{-2m\epsilon^2}$$

*some hypothesis  $h$*

$$\text{error}_{\text{train}}(h) = \frac{1}{m} \sum_i \mathbb{1}(h(x_i) \neq y_i)$$

$$\theta \leftarrow \text{error}_{\text{true}}(h) = E_x[\mathbb{1}(h(x) \neq y)]$$

$$P(\text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h) > \epsilon) \leq e^{-2m\epsilon^2}$$

©2005-2007 Carlos Guestrin

14

## But we are comparing many hypothesis: **Union bound**

For each hypothesis  $h_i$ :

$$P(\text{error}_{\text{true}}(h_i) - \text{error}_{\text{train}}(h_i) > \epsilon) \leq e^{-2m\epsilon^2}$$

What if I am comparing two hypothesis,  $h_1$  and  $h_2$ ?

better?  $h_1$  or  $h_2$

$$\left. \begin{array}{l} \text{error}_{\text{train}}(h_1) \stackrel{?}{>} \text{error}_{\text{train}}(h_2) \\ \text{error}_{\text{true}}(h_1) \stackrel{?}{>} \text{error}_{\text{true}}(h_2) \end{array} \right\} \text{, want } \text{error}_{\text{true}}(h_1) > \text{error}_{\text{true}}(h_2)$$

$$P(\text{error}_{\text{true}}(h_1) - \text{error}_{\text{train}}(h_1) > \epsilon \vee \text{error}_{\text{true}}(h_2) - \text{error}_{\text{train}}(h_2) > \epsilon)$$

$$\leq P(\text{error}_{\text{true}}(h_1) - \text{error}_{\text{train}}(h_1) > \epsilon) + P(\text{error}_{\text{true}}(h_2) - \text{error}_{\text{train}}(h_2) > \epsilon)$$

$$\leq 2e^{-2m\epsilon^2}$$

©2005-2007 Carlos Guestrin

15

## Generalization bound for $|H|$ hypothesis

- Theorem:** Hypothesis space  $H$  finite, dataset  $D$  with  $m$  i.i.d. samples,  $0 < \epsilon < 1$  : for any learned hypothesis  $h$ :

$$P(\text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h) > \epsilon) \leq |H|e^{-2m\epsilon^2} \leq \sigma$$

holds w.p.  $1 - \sigma$ :

$$\epsilon \geq \sqrt{\frac{\ln |H| + \ln \frac{1}{\sigma}}{2m}} \quad \text{with prob. } 1 - \sigma$$

$$\text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h) \leq \sqrt{\frac{\ln |H| + \ln \frac{1}{\sigma}}{2m}}$$

©2005-2007 Carlos Guestrin

16



## PAC bound and Bias-Variance tradeoff

$$P(\text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h) > \epsilon) \leq |H|e^{-2m\epsilon^2}$$

or, after moving some terms around,  
with probability at least  $1-\delta$ :

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{\ln |H| + \ln \frac{1}{\delta}}{2m}}$$

Handwritten notes and table:

- as  $m \rightarrow \infty$ , "variance"  $\rightarrow 0$ ,  $\text{error}_{\text{train}}(h)$  could increase
- $|H| \rightarrow \infty$  increases as  $m$  increases in the hope that  $\text{error}_{\text{train}}$  decreases
- Want small  $\text{error}_{\text{true}}(h)$
- fixed  $m$ 

complex hypothesis space	small	large, because $ H $ is big
simple hyp. space	large	small because $\ln(H)$ is small

- Important: PAC bound holds for all  $h$ , but doesn't guarantee that algorithm finds best  $h$ !!!

©2005-2007 Carlos Guestrin

17

## What about the size of the hypothesis space?

$$m \geq \frac{1}{2\epsilon^2} \left( \ln |H| + \ln \frac{1}{\delta} \right)$$

- How large is the hypothesis space?

$\ln |H|$

©2005-2007 Carlos Guestrin

18

# Boolean formulas with $n$ binary features

$n$  binary attributes

$$m \geq \frac{1}{2\epsilon^2} \left( \ln |H| + \ln \frac{1}{\delta} \right)$$

$H$ : any boolean formula

$x_1$	$x_2$	...	$x_n$	$y$
0	0	0	0	0
0	0	0	1	0
...	...	...	...	...
1	1	1	1	1

$2^n$  rows  
 $|H| = 2^n$   
 $\ln |H| = 2^n \ln 2$   
 really really big  
 need every row to learn

$H$ : all conjunctions of literals

$x_1 \wedge x_2 \wedge x_7$   
 $x_3 \wedge x_8 \wedge \neg x_{12}$   
 ;

for each  $x_i$ : 3 choices  
 - not included  
 - included positive  
 - included negative

$$|H| = 3^n$$

$$\ln |H| = n \ln 3$$

©2005-2007 Carlos Guestrin

19

# Number of decision trees of depth $k$

$$m \geq \frac{1}{2\epsilon^2} \left( \ln |H| + \ln \frac{1}{\delta} \right)$$

Recursive solution

Given  $n$  attributes

$H_k$  = Number of decision trees of depth  $k$

$$H_0 = 2$$

$$H_{k+1} = (\text{\#choices of root attribute}) \cdot$$

(# possible left subtrees) \*

(# possible right subtrees)

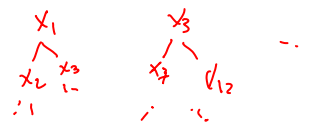
$$H_{k+1} = n \cdot H_k \cdot H_k$$

Write  $L_k = \log_2 H_k$

$$L_0 = 1$$

$$L_{k+1} = \log_2 n + 2L_k$$

$$\text{So } L_k = (2^k - 1)(1 + \log_2 n) + 1$$



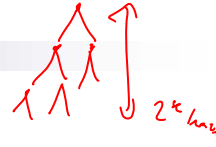
$$\ln |H| \leq (2^k - 1)(1 + \log_2 n) + 1$$

©2005-2007 Carlos Guestrin

20

## PAC bound for decision trees of depth k

$$m \geq \frac{\ln 2}{2\epsilon^2} \left( (2^k - 1)(1 + \log_2 n) + 1 + \ln \frac{1}{\delta} \right)$$



### ■ Bad!!!

- Number of points is exponential in depth!

### ■ But, for m data points, decision tree can't get too big...

*no more than m leaves*

**Number of leaves never more than number data points**

©2005-2007 Carlos Guestrin

21

## Number of decision trees with k leaves

$$m \geq \frac{1}{2\epsilon^2} \left( \ln |H| + \ln \frac{1}{\delta} \right)$$

$H_k$  = Number of decision trees with k leaves

$H_0 = 2$

$$H_{k+1} = n \sum_{i=1}^k H_i H_{k+1-i}$$

**Loose bound:**

$$H_k = n^{k-1} (k+1)^{2k-1}$$

*really big*

**Reminder:**

$$|\text{DTs depth } k| = 2 * (2n)^{2k-1}$$

*really really big*

©2005-2007 Carlos Guestrin

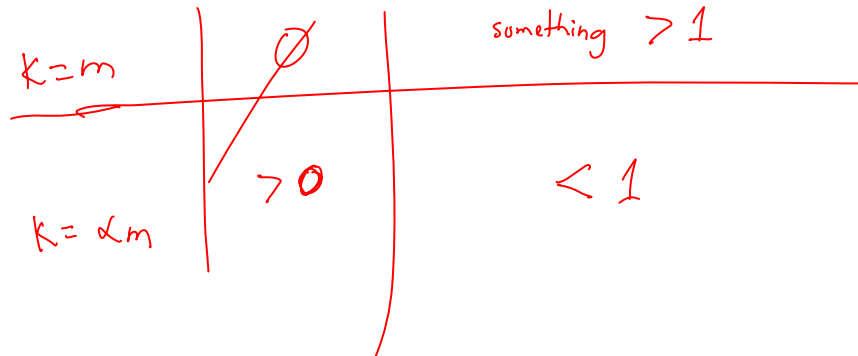
22

## PAC bound for decision trees with k leaves – Bias-Variance revisited

$$H_k = n^{k-1}(k+1)^{2k-1} \quad \text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{\ln |H| + \ln \frac{1}{\delta}}{2m}}$$

Number of features

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{(k-1) \ln n + (2k-1) \ln(k+1) + \ln \frac{1}{\delta}}{2m}}$$



## Announcements

### ■ Midterm:

- Thursday Oct. 25th, Thursday 5-6:30pm, MM A14
  - All content up to, and including SVMs and Kernels
    - Not learning theory
  - any book, class notes, your printouts of class materials that are on the class website, including my annotated slides and relevant readings, and Andrew Moore's tutorials. You cannot use materials brought by other students.
  - Calculators are not necessary.
  - No laptops, PDAs or cellphones.

## What did we learn from decision trees?

- Bias-Variance tradeoff formalized

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{(k-1) \ln n + (2k-1) \ln(k+1) + \ln \frac{1}{\delta}}{2m}}$$

- Moral of the story:

Complexity of learning not measured in terms of size hypothesis space, but in maximum number of points that allows consistent classification

- Complexity  $m$  – no bias, lots of variance
- Lower than  $m$  – some bias, less variance

©2005-2007 Carlos Guestrin

25

## What about continuous hypothesis spaces?

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{\ln |H| + \ln \frac{1}{\delta}}{2m}}$$

- Continuous hypothesis space:

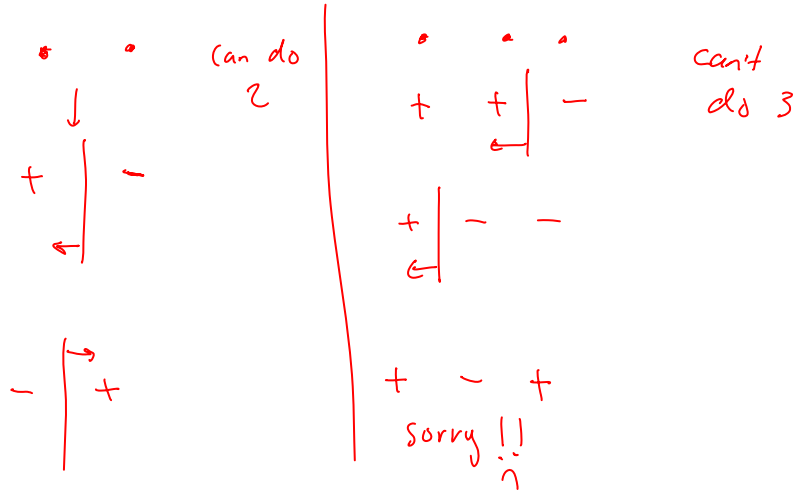
- $|H| = \infty$
- Infinite variance???

- **As with decision trees, only care about the maximum number of points that can be classified exactly!**

©2005-2007 Carlos Guestrin

26

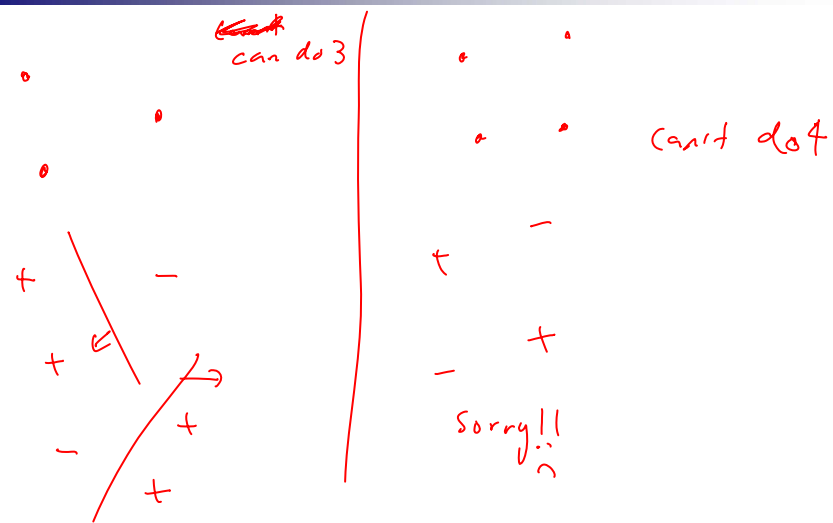
# How many points can a linear boundary classify exactly? (1-D)



©2005-2007 Carlos Guestrin

27

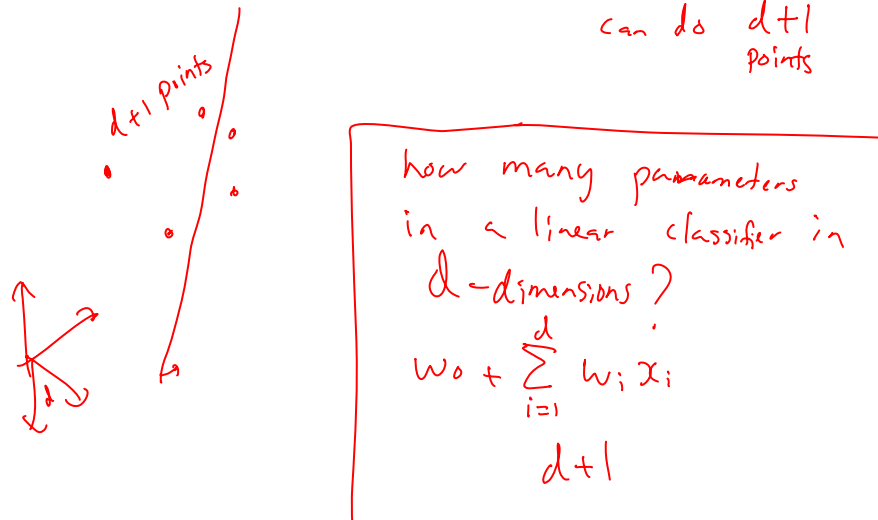
# How many points can a linear boundary classify exactly? (2-D)



©2005-2007 Carlos Guestrin

28

## How many points can a linear boundary classify exactly? (d-D)



©2005-2007 Carlos Guestrin

29

## PAC bound using VC dimension

- Number of training points that can be classified exactly is VC dimension!!! *e.g. linear classifiers*
- Measures relevant size of hypothesis space, as *may be continuous*  
with decision trees with  $k$  leaves

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{VC(H) \left( \ln \frac{2m}{VC(H)} + 1 \right) + \ln \frac{4}{\delta}}{m}}$$

only depend on  $VC(H)$   
not on  $|H|$

©2005-2007 Carlos Guestrin

30

# Shattering a set of points

*Definition:* a **dichotomy** of a set  $S$  is a partition of  $S$  into two disjoint subsets.

*Definition:* a set of instances  $S$  is **shattered** by hypothesis space  $H$  if and only if for every dichotomy of  $S$  there exists some hypothesis in  $H$  consistent with this dichotomy.

©2005-2007 Carlos Guestrin

31

# VC dimension

*Definition:* The **Vapnik-Chervonenkis dimension**,  $VC(H)$ , of hypothesis space  $H$  defined over instance space  $X$  is the size of the largest finite subset of  $X$  shattered by  $H$ . If arbitrarily large finite sets of  $X$  can be shattered by  $H$ , then  $VC(H) \equiv \infty$ .

©2005-2007 Carlos Guestrin

32



## PAC bound using VC dimension

- Number of training points that can be classified exactly is VC dimension!!!
  - Measures relevant size of hypothesis space, as with decision trees with  $k$  leaves
  - Bound for infinite dimension hypothesis spaces:

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{VC(H) \left( \ln \frac{2m}{VC(H)} + 1 \right) + \ln \frac{4}{\delta}}{m}}$$

©2005-2007 Carlos Guestrin

33

## Examples of VC dimension

- Linear classifiers:
  - $VC(H) = d+1$ , for  $d$  features plus constant term  $b$
- Neural networks
  - $VC(H) = \text{\#parameters}$
  - Local minima means NNs will probably not find best parameters
- 1-Nearest neighbor?

©2005-2007 Carlos Guestrin

34

## Another VC dim. example - What can we shatter?

- What's the VC dim. of decision stumps in 2d?

©2005-2007 Carlos Guestrin

35

## Another VC dim. example - What can't we shatter?

- What's the VC dim. of decision stumps in 2d?

©2005-2007 Carlos Guestrin

36

## What you need to know

- Finite hypothesis space
  - Derive results
  - Counting number of hypothesis
  - Mistakes on Training data
- Complexity of the classifier depends on number of points that can be classified exactly
  - Finite case – decision trees
  - Infinite case – VC dimension
- Bias-Variance tradeoff in learning theory
- Remember: will your algorithm find best classifier?

©2005-2007 Carlos Guestrin

37

# Big Picture

Machine Learning – 10701/15781  
Carlos Guestrin  
Carnegie Mellon University

October 24<sup>th</sup>, 2007

©2005-2007 Carlos Guestrin

38

## What you have learned thus far

- Learning is function approximation
- Point estimation
- Regression
- Naïve Bayes
- Logistic regression
- Bias-Variance tradeoff
- Neural nets
- Decision trees
- Cross validation
- Boosting
- Instance-based learning
- SVMs
- Kernel trick
- PAC learning
- VC dimension
- *margin* bounds
- Mistake bounds



©2005-2007 Carlos Guestrin

39

## Review material in terms of...

- Types of learning problems
- Hypothesis spaces  
*what they can represent*
- Loss functions
- Optimization algorithms

©2005-2007 Carlos Guestrin

40

# BIG PICTURE

(a few points of comparison)

