

# Introduction to Learning Theory Part 2

CS 760@UW-Madison



# Goals for the lecture



you should understand the following concepts

- consistent learners and version spaces
- PAC learnability and sample complexity
- VC-dimension

# PAC Learning Theory



# Probably Approximately Correct (PAC) learning



[Valiant, CACM 1984]

- Consider a class  $C$  of possible target concepts defined over a set of instances  $\mathcal{X}$  of length  $n$ , and a learner  $L$  using hypothesis space  $H$
- $C$  is PAC learnable by  $L$  using  $H$  if, for all
  - $c \in C$
  - distributions  $\mathcal{D}$  over  $\mathcal{X}$
  - $\varepsilon$  such that  $0 < \varepsilon < 0.5$
  - $\delta$  such that  $0 < \delta < 0.5$
- learner  $L$  will, with probability at least  $(1-\delta)$ , output a hypothesis  $h \in H$  such that  $error_{\mathcal{D}}(h) \leq \varepsilon$  in time that is polynomial in
  - $1/\varepsilon$
  - $1/\delta$
  - $n$
  - $size(c)$

# PAC learning and consistency



- Suppose we can find hypotheses that are consistent with  $m$  training instances.
- We can analyze PAC learnability by determining whether
  1. The needed  $m$  grows polynomially in the relevant parameters
  2. the processing time per training example is polynomial

# Version spaces



- A hypothesis  $h$  is *consistent* with a set of training examples  $D$  of target concept if and only if  $h(\mathbf{x}) = c(\mathbf{x})$  for each training example  $\langle \mathbf{x}, c(\mathbf{x}) \rangle$  in  $D$

$$\text{consistent}(h, D) \equiv (\forall \langle \mathbf{x}, c(\mathbf{x}) \rangle \in D) h(\mathbf{x}) = c(\mathbf{x})$$

- The version space  $VS_{H,D}$  with respect to hypothesis space  $H$  and training set  $D$ , is the subset of hypotheses from  $H$  consistent with all training examples in  $D$

$$VS_{H,D} \equiv \{h \in H \mid \text{consistent}(h, D)\}$$

# Exhausting the version space



- The version space  $VS_{H,D}$  is  $\varepsilon$ -exhausted with respect to  $c$  and  $D$  if every hypothesis  $h \in VS_{H,D}$  has true error  $< \varepsilon$

$$\left( \forall h \in VS_{H,D} \right) error_D(h) < \varepsilon$$

# Exhausting the version space



- Suppose that every  $h$  in our version space  $VS_{H,D}$  is consistent with  $m$  training examples
- The probability that  $VS_{H,D}$  is not  $\varepsilon$ -exhausted (i.e. that it contains some hypotheses that are not accurate enough)

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

Proof:

$$(1 - \varepsilon)^m$$

probability that some hypothesis with error  $> \varepsilon$  is consistent with  $m$  training instances

$$k(1 - \varepsilon)^m$$

there might be  $k$  such hypotheses

$$|H|(1 - \varepsilon)^m$$

$k$  is bounded by  $|H|$

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

$(1 - \varepsilon) \leq e^{-\varepsilon}$  when  $0 \leq \varepsilon \leq 1$





# Sample complexity for finite hypothesis spaces

[Blumer et al., *Information Processing Letters* 1987]

- we want to reduce this probability below  $\delta$

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

- solving for  $m$  we get

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

log dependence on  $H$

$\varepsilon$  has stronger influence than  $\delta$

# PAC analysis example: learning conjunctions of Boolean literals



- each instance has  $n$  Boolean features
- learned hypotheses are of the form  $Y = X_1 \wedge X_2 \wedge \neg X_5$

How many training examples suffice to ensure that with prob  $\geq 0.99$ , a consistent learner will return a hypothesis with error  $\leq 0.05$  ?

there are  $3^n$  hypotheses (each variable can be present and unnegated, present and negated, or absent) in  $H$

$$m \geq \frac{1}{.05} \left( \ln(3^n) + \ln\left(\frac{1}{.01}\right) \right)$$

for  $n=10$ ,  $m \geq 312$

for  $n=100$ ,  $m \geq 2290$

# PAC analysis example: learning conjunctions of Boolean literals



- we've shown that the sample complexity is polynomial in relevant parameters:  $1/\epsilon$ ,  $1/\delta$ ,  $n$
- to prove that Boolean conjunctions are PAC learnable, need to also show that we can find a consistent hypothesis in polynomial time (the FIND-S algorithm in Mitchell, Chapter 2 does this)

FIND-S:

initialize  $h$  to the most specific hypothesis  $x_1 \wedge \neg x_1 \wedge x_2 \wedge \neg x_2 \dots x_n \wedge \neg x_n$

for each positive training instance  $x$

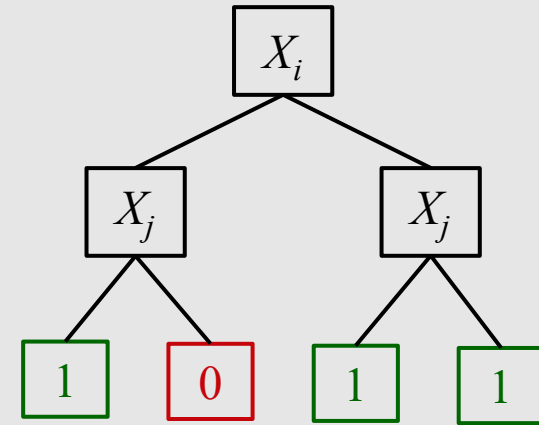
    remove from  $h$  any literal that is not satisfied by  $x$

output hypothesis  $h$



# PAC analysis example: learning decision trees of depth 2

- each instance has  $n$  Boolean features
- learned hypotheses are DTs of depth 2 using only 2 variables



$$|H| = \binom{n}{2} \times 16 = \frac{n(n-1)}{2} \times 16 = 8n(n-1)$$

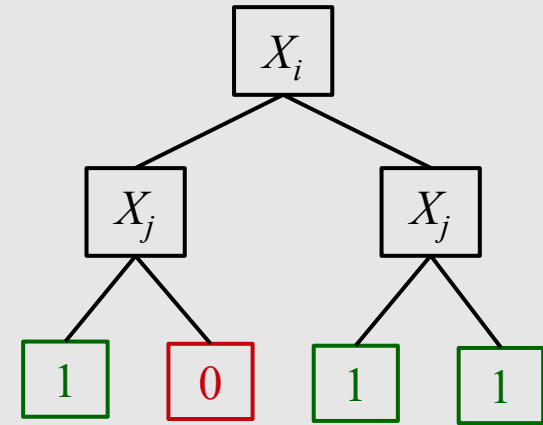
# possible split choices

# possible leaf labelings

# PAC analysis example: learning decision trees of depth 2



- each instance has  $n$  Boolean features
- learned hypotheses are DTs of depth 2 using only 2 variables



How many training examples suffice to ensure that with prob  $\geq 0.99$ , a consistent learner will return a hypothesis with error  $\leq 0.05$  ?

$$m \geq \frac{1}{.05} \left( \ln(8n^2 - 8n) + \ln\left(\frac{1}{.01}\right) \right)$$

for  $n=10$ ,  $m \geq 224$

for  $n=100$ ,  $m \geq 318$

# PAC analysis example: $K$ -term DNF is not PAC learnable



- each instance has  $n$  Boolean features
- learned hypotheses are of the form  $Y = T_1 \vee T_2 \vee \dots \vee T_k$  where each  $T_i$  is a conjunction of  $n$  Boolean features or their negations

$|H| \leq 3^{nk}$ , so sample complexity is polynomial in the relevant parameters

$$m \geq \frac{1}{\varepsilon} \left( nk \ln(3) + \ln \left( \frac{1}{\delta} \right) \right)$$

however, the computational complexity (time to find consistent  $h$ ) is not polynomial in  $m$  (e.g. graph 3-coloring, an NP-complete problem, can be reduced to learning 3-term DNF)

# Comments on PAC learning



- PAC analysis formalizes the learning task and allows for non-perfect learning (indicated by  $\varepsilon$  and  $\delta$ )
  - Requires polynomial computational time
- finding a consistent hypothesis is sometimes easier for larger concept classes
  - e.g. although  $k$ -term DNF is not PAC learnable, the more general class  $k$ -CNF is
- PAC analysis has been extended to explore a wide range of cases
  - the target concept not in our hypothesis class
  - infinite hypothesis class (VC-dimension theory)
  - noisy training data
  - learner allowed to ask queries
  - restricted distributions (e.g. uniform) over  $\mathcal{D}$
  - etc.
- most analyses are worst case
- sample complexity bounds are generally not tight

# The Agnostic Case





# What if the target concept is not in our hypothesis space?



- so far, we've been assuming that the target concept  $c$  is in our hypothesis space; this is not a very realistic assumption
- *agnostic learning* setting
  - don't assume  $c \in H$
  - learner returns hypothesis  $h$  that makes fewest errors on training data

# Hoeffding bound



- we can approach the agnostic setting by using the Hoeffding bound
- let  $Z_1 \dots Z_m$  be a sequence of  $m$  independent Bernoulli trials (e.g. coin flips), each with probability of success  $E[Z_i] = p$
- let  $S = Z_1 + \dots + Z_m$

$$P[S < (p - \varepsilon)m] \leq e^{-2m\varepsilon^2}$$

# Agnostic PAC learning



- applying the Hoeffding bound to characterize the error rate of a given hypothesis

$$P[\text{error}_{\mathcal{D}}(h) > \text{error}_{\mathcal{D}}(h) + \varepsilon] \leq e^{-2m\varepsilon^2}$$

- but our learner searches hypothesis space to find  $h_{best}$

$$P[\text{error}_{\mathcal{D}}(h_{best}) > \text{error}_{\mathcal{D}}(h_{best}) + \varepsilon] \leq |H|e^{-2m\varepsilon^2}$$

- solving for the sample complexity when this probability is limited to  $\delta$

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

# VC-dimension



# What if the hypothesis space is not finite?



- **Q:** If  $H$  is infinite (e.g. the class of perceptrons), what measure of hypothesis-space complexity can we use in place of  $|H|$  ?
  
  
  
  
  
  
  
  
  
  
- **A:** the largest subset of  $\mathcal{X}$  for which  $H$  can guarantee zero training error, regardless of the target function.

this is known as the *Vapnik-Chervonenkis dimension* (VC-dimension)

# Shattering and the VC dimension

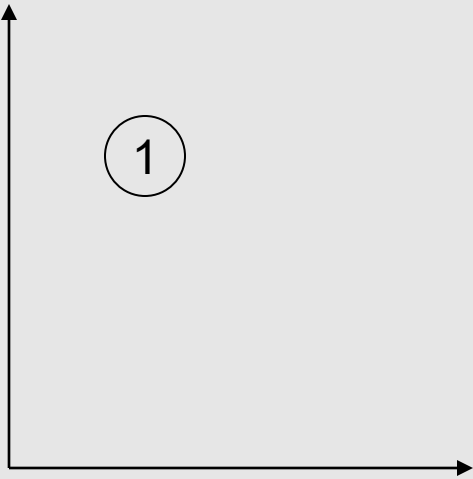


- a set of instances  $D$  is *shattered* by a hypothesis space  $H$  iff for every dichotomy of  $D$  there is a hypothesis in  $H$  consistent with this dichotomy
- the *VC dimension* of  $H$  is the size of the largest set of instances that is shattered by  $H$

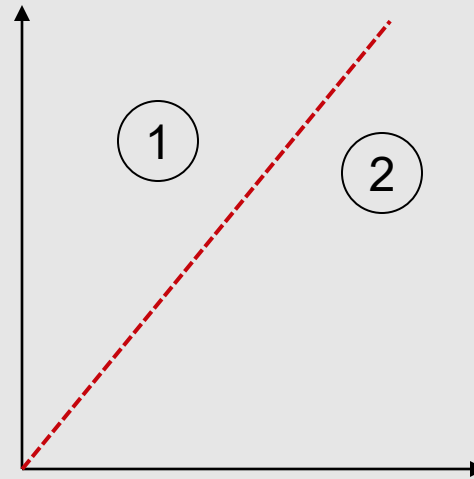
# Infinite hypothesis space with a finite VC dimension

consider:  $H$  is set of lines in 2D (i.e. perceptrons in 2D feature space)

can find an  $h$  consistent with 1 instance  
no matter how it's labeled



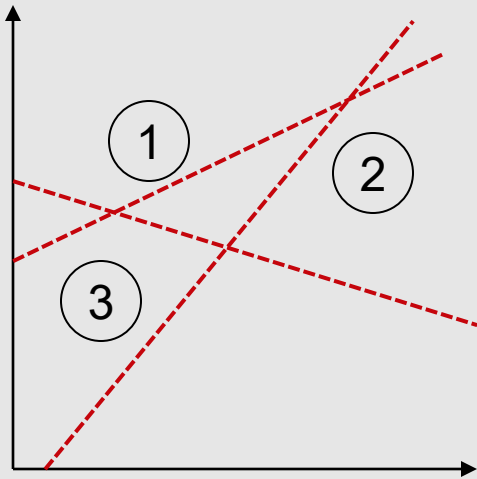
can find an  $h$  consistent with 2  
instances no matter labeling



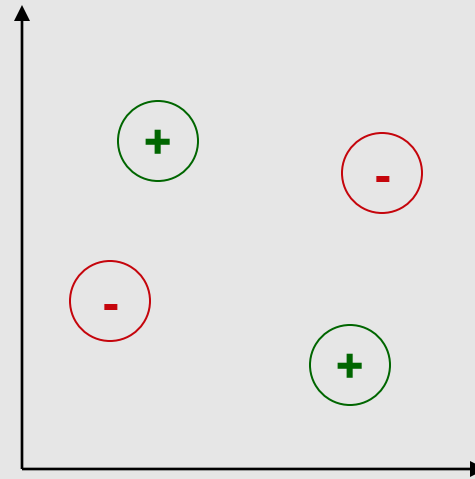
# Infinite hypothesis space with a finite VC dimension

consider:  $H$  is set of lines in 2D

can find an  $h$  consistent with 3 instances no matter labeling (assuming they're not colinear)



cannot find an  $h$  consistent with 4 instances for some labelings



can shatter 3 instances, but not 4, so the  $VC\text{-dim}(H) = 3$

more generally, the  $VC\text{-dim}$  of hyperplanes in  $n$  dimensions =  $n+1$



# VC dimension for finite hypothesis spaces



for finite  $H$ ,  $\text{VC-dim}(H) \leq \log_2 |H|$

Proof:

suppose  $\text{VC-dim}(H) = d$

for  $d$  instances,  $2^d$  different labelings possible

therefore  $H$  must be able to represent  $2^d$  hypotheses

$$2^d \leq |H|$$

$$d = \text{VC-dim}(H) \leq \log_2 |H|$$

# Sample complexity and the VC dimension

- using  $\text{VC-dim}(H)$  as a measure of complexity of  $H$ , we can derive the following bound [Blumer et al., *JACM* 1989]

$$m \geq \frac{1}{\varepsilon} \left( 4 \log_2 \left( \frac{2}{\delta} \right) + 8 \text{VC-dim}(H) \log_2 \left( \frac{13}{\varepsilon} \right) \right)$$

$m$  grows  $\log \times$  linear in  $\varepsilon$  (better than earlier bound)

can be used for both finite and infinite hypothesis spaces

# Lower bound on sample complexity

[Ehrenfeucht et al., *Information & Computation* 1989]



- there exists a distribution  $\mathcal{D}$  and target concept in  $C$  such that if the number of training instances given to  $L$

$$m < \max \left[ \frac{1}{\varepsilon} \log \left( \frac{1}{\delta} \right), \frac{\text{VC-dim}(C) - 1}{32\varepsilon} \right]$$

then with probability at least  $\delta$ ,  $L$  outputs  $h$  such that  $error_{\mathcal{D}}(h) > \varepsilon$



# THANK YOU

Some of the slides in these lectures have been adapted/borrowed from materials developed by Mark Craven, David Page, Jude Shavlik, Tom Mitchell, Nina Balcan, Elad Hazan, Tom Dietterich, and Pedro Domingos.

