

## PAC Model of Learning:

There is a distribution  $D$  on  $\{0, 1\}^n$  ( $\mathbb{R}^n$ )  
with function class  $C = \{\text{decision trees of size } s\}$

A Learner (runs in polynomial time) receives:

Fix  $c \in C$ ,  $c$  is the unknown decision tree  
script that we want to learn.

Receives  $(x^1, y^1)$ ,  $x \sim D$ , and  $y = c(x)$ .  
 $(x^2, y^2)$  drawn from probability Distribution  $D$   
 $\vdots$   
 $(x^n, y^n)$  and  $y^i = c(x^i)$

Goal: Learner output  $h \in C$   $\nwarrow$  classifier/hypothesis.

$$\Pr_{x \sim D} [h(x) \neq c(x)] \leq \epsilon$$

$\uparrow$  something small.

Learner should be efficient:  $(n, s)$  runtime.

## Goal of PAC

with probability <sup>"Drawn over distribution D"</sup> at least  $1 - \delta$ , the learner should output a hypothesis  $h$  s.t.

$$P_{x \sim D} [h(x) \neq \underbrace{c(x)}_{\substack{\text{theoretical best model drawn from } C}}] \leq \epsilon$$

$\uparrow$   
our model

Runtime = polynomial  $(\frac{1}{\epsilon}, \frac{1}{\delta}, n, S)$ .

Why need  $1 - \delta$ ? <sup>p + f success</sup> Imagine learner requests

$x^1, \dots, x^m$ , but receives the same training example every time, so no way the error will be small. So we have to allow for some probability of failure. The probability of seeing the same example over and over again is very small.

PAC = Probably Approximately Correct, L. Valiant [1984]

$\delta$  = failure probability. If you want an algorithm that has smaller + smaller probability of inaccuracy,  $\epsilon$ , you get to run in more time + use samples.

This works for any binary classification (boolean).

$\uparrow$  PAC Learning.

Probably: confidence  $1 - \delta$

Approximately correct:  $1 - \epsilon$

When can we PAC learn a function class?

What function classes can we PAC learn?

Give the learner an algorithm  $A$ .

$A$ : training sets  $\xrightarrow{\text{map}}$  decision trees.

$A(S)$  output a tree  $T$  that is consistent with  $S$ .

Size of  $T$  is going to be at most  $S$ .  
 $\uparrow$  # of nodes.  $\uparrow$  little  $S$

$A$  always outputs a consistent hypothesis from  $C$  given any training set (assuming there is one),  
 $\uparrow$  script  $C$

Question: Given algorithm  $A$ : how can we PAC learn  $C$ ?

- Draw sufficiently many training points.
- Use  $A$  to find  $c \in C$  consistent with  $S$ .  
 $\uparrow$  function class/script  $C$ .
- output  $c$ .  $\leftarrow$  little  $c$

Q: How large should  $S$  (training set) be?

Example: Jar 1: all blue marbles. | Jar 2: 90% red, 10% blue.

Goal: Figure out if given Jar 1 or 2.

Case 1: red marble  $\rightarrow$  Jar 2.

Case 2: blue marble  $\rightarrow$  probably Jar 1.

draw 100 total.

$$\Pr[\text{failure}] = (0.1)^{100} = \text{"}\delta\text{-parameter."}$$

Let's return to PAC learning:

- Draw many samples  $S$
- Run  $A$  (algorithm)
- Output classifier  $c$  that is consistent with  $S$  given from  $A$ .

What is the probability this procedure fails?  $\leq \delta$

Bad Event: we output  $c \in C$  that is consistent with  $S$  but the true error is greater than  $\epsilon$ .

$P[\text{Bad Event}]$ ? Let's imagine that we have enumerated all functions in  $C = \{c_1, c_2, \dots, c_n\}$

Fix  $c_1$ . Assume  $c_1$  has true error  $> \epsilon$ .

What is  $\Pr_S[c_1 \text{ is consistent with } S]$ ?

Answer:  $\leq \underbrace{(1-\epsilon)}_{\text{success}}^{|S|}$  (same as marble example: draw the wrong example)  
size of  $S$

Fix  $c_2$ . " $c_2$  has "  $> \epsilon$ .

"  $\Pr_S[c_2 \text{ "}]$ ? Also  $\leq (1-\epsilon)^{|S|}$ .

For every  $c_i$  (with error  $> \epsilon$ )

$\Pr_S[c_i \text{ is consistent on } S] \leq (1-\epsilon)^{|S|}$ .

Q: Randomly from  $S$ , what is the prob there exists a function  $c \in C$  whose error  $> \epsilon$  and is consistent with  $S$ ?

Use the union bound:  $P(A \cup B) \leq P(A) + P(B)$ . want the bound to be less than  $\delta$ .

$$Pr[\text{Bad Event}] \leq |C| \cdot (1 - \epsilon)^{|S|} \leq \delta$$

There is  $c \in C$  that is consistent with  $S$  but has true error  $> \epsilon$ .  
 at most script  $C$  functions to consider.  
 so we can just add up the failure probabilities. This is worse case scenario.

Solve for  $|S|$ : use property  $(1+x) = e^x, (1-x) = e^{-x}$

$$P[\text{Bad Event}] = |C| \cdot (1 - \epsilon)^{|S|} \leq \delta$$

$$|C| \cdot e^{-\epsilon|S|} \leq \delta$$

$$e^{-\epsilon|S|} \leq \frac{\delta}{|C|}$$

$$-\epsilon|S| \leq \ln\left(\frac{\delta}{|C|}\right)$$

$$\boxed{|S| \geq \frac{\ln\left(\frac{|C|}{\delta}\right)}{\epsilon}}$$

∴ If you choose # training points larger than  $\frac{\ln\left(\frac{|C|}{\delta}\right)}{\epsilon}$ , then with probability  $\geq 1 - \delta$ , function output

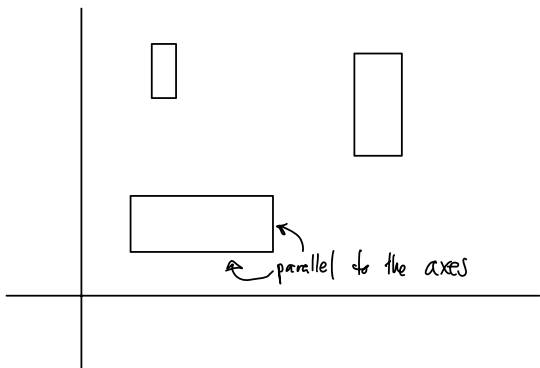
$C$  is at least  $(1 - \epsilon)$  accurate.

⇒ suggest polynomial size bounds.

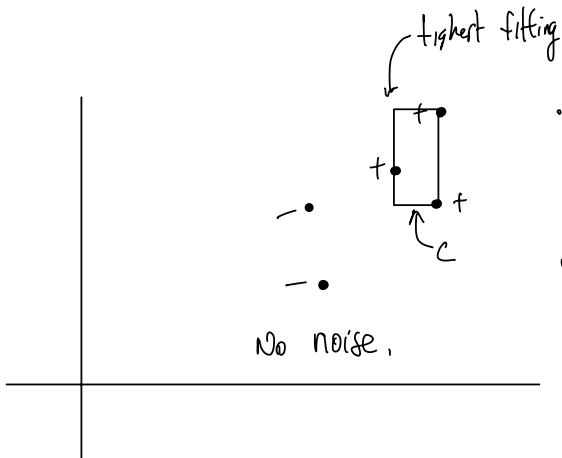
⇒ Suggest a "consistent hypothesis" approach to learning.

## Infinite Function

PAC-learning. Axis-parallel rectangles.



Infinite many axis-parallel rectangles.

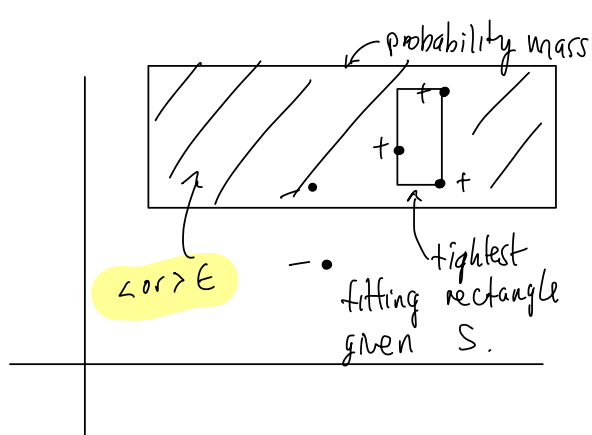


- labelled + if the point is inside  $C$ , unknown axis-parallel rectangle.
- labelled - if the point is outside  $C$ .

Goal : Given  $\epsilon, \delta$ , output  $h$  that is  $\epsilon$ -accurate with prob  $\geq (1 - \delta)$ .

claim: the tightest fitting rectangle works. Tightest fitting rectangle containing all the positive points in the set  $S$ .

Question: How large to choose  $|S|$ ?  
 $\uparrow$   
 training set

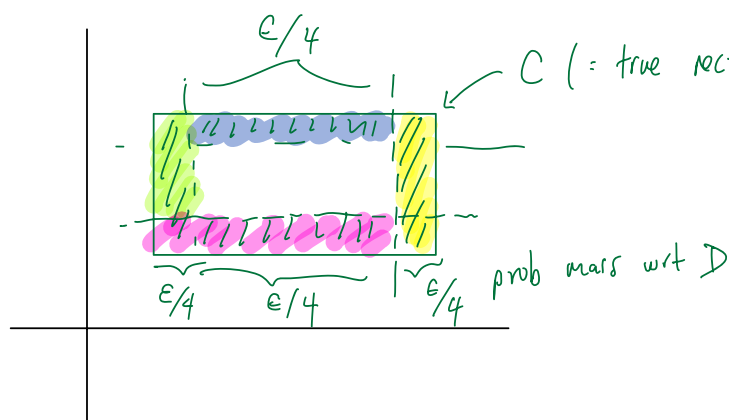


Bad Event: tightest fitting rectangle is small, i.e.

lots of probability mass of the true rectangle exists outside of  $h$ , the tightest fitting rectangle.

$P(\text{Bad Event}) = \text{shaded region}$ . If  $P(\text{Bad Event}) \geq \epsilon$ , then trouble.  
How do we bound the probability this happens?

analyze  $h$  and say something about  $h$  vs all rectangles that are large and contain  $h$ .



Intuition: taking more examples/training data decreases the probability for Bad Event.

$$\text{Bad Events} = B_1 + B_2 + B_3 + B_4.$$

- $B_1$  is the event we see no points in the right strip.
  - $B_2$  is event we see no point in the bottom strip.
  - $B_3$  " " left " "
  - $B_4$  " " top " "
- } jar example.

Claim: If neither  $B_1, B_2, B_3, B_4$  occur,  
then  $h$ , which is the highest fitting rectangle,  
is  $\epsilon$ -accurate.

What does it mean?

True error of the highest fitting rectangle is  
going to be at most  $\epsilon$  as long as  
none of the bad events occur.

Claim: If choose  $m$  random samples

$$\Pr[B_1] \leq \left(1 - \frac{\epsilon}{4}\right)^m \quad \begin{matrix} \uparrow \text{no point in right strip.} \\ \text{each draw independent.} \end{matrix} \quad (\text{marble jar example}).$$

$$\Pr[B_1 \cup B_2 \cup B_3 \cup B_4] \leq 4 \cdot \left(1 - \frac{\epsilon}{4}\right)^m$$

$$\text{want it} \leq \delta : 4 \left(1 - \frac{\epsilon}{4}\right)^m \leq \delta$$

$$\left(1 - \frac{\epsilon}{4}\right)^m \leq \frac{\delta}{4}$$

$$\begin{matrix} 1+x \approx e^x \\ 1-x \approx e^{-x} \end{matrix}$$

$1-\delta$  accurate.

$$e^{-\epsilon m/4} \leq \frac{\delta}{4}$$

$$-\frac{\epsilon m}{4} \leq \ln \left[ \frac{\delta}{4} \right]$$

$$\Rightarrow \boxed{m \geq \frac{4 \ln \left[ \frac{4}{\delta} \right]}{\epsilon}}$$

$h$ , the highest triangle, will be  $\epsilon$ -accurate with  
probability  $\geq (1-\delta)$ .

---



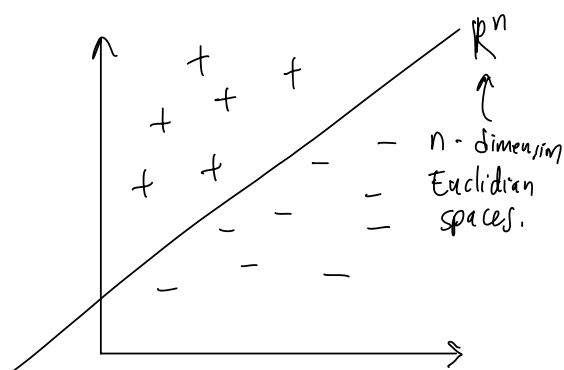
← function class  
 $C = \text{Half Spaces}$

•  $C = \{ \text{half spaces} \} = \text{dividing } n\text{-dimension Euclidean spaces into half-spaces,}$   
 ↑ infinite.

$$f = \text{sign}(w \cdot x - \theta)$$

$w \in \mathbb{R}^n$  (unknown),  $x \in \mathbb{R}^n$  (vector),  $\theta \in \mathbb{R}$  (scalar, unknown).

$f$  is boolean. 0 if output is negative,  
 1 if output is positive.



Want:  
 come up with PAC learning algorithms for half spaces.

one Approach for Learning Halfspaces:

$w \in \mathbb{R}^n$  is unknown

$\theta \in \mathbb{R}$  is unknown

$$f = \text{sign}\left(\sum_{i=1}^n w_i x_i - \theta\right)$$

Given draws from  $\mathcal{D}(x, f(x))$   
 $x \sim$

$(01010, \text{pos}) \rightarrow w_2 + w_4 > \theta$   
 $(0110, \text{neg}) \rightarrow w_2 + w_3 \leq \theta$

Assume  $w_i$  is integer in some bounded range.

Each labeled example  $\rightarrow$  linear inequalities,  
system of linear inequalities,  
Can we find a consistent hypothesis?

Use: General-purpose tool called linear programming.

Perceptron: algorithm for learning half-spaces.

Read: Ch 3 of textbook.