

PAC LEARNING

- WHAT IS THE "TRUE ERROR" OR GENERALIZATION ERROR OF A CLASSIFIER?
- DECISION TREES: FIX T

D ON NEW EXAMPLES

(x, y)
CHALLENGE
LABEL L

$$\Pr_{(x,y) \sim D} [T(x) \neq y]$$

TRUE ERROR
GENERALIZATION ERROR

OF T

"		
	x^1	y^1
	:	:
	:	:
	:	:
	x^m	y^m
"		

S TRAINING SET

LEARNER IS GIVEN S

X	Y
x^1	y^1
:	:
:	:
x^m	y^m

- EACH
 - $x^i \in \{0,1\}^n$
 - $y^i \in \{0,1\}$

LET'S BUILD A DECISION TREE

x^i ARE DISTINCT

YOU CAN BUILD A DECISION TREE (SIZE $\geq |S|$)
 THAT IS CONSISTENT WITH ALL THE POINTS
 IN S.

QUESTION: HOW WELL DOES THIS TREE GENERALIZE
 WHAT IS THE TRUE ERROR OF THIS
 TREE?

X

0 5 C

How can we estimate the true error of a decision tree?

"HOLD-OUT" OR A "VALIDATION SET"

$S \leftarrow$ TRAINING SET

$H \leftarrow$ HOLD-OUT

CROSS-VALIDATION

1. USE S
TO BUILD A DECISION
TREE

2. ESTIMATE TREE'S
TRUE ERROR
VIA ITS ERROR ON
 H .

ANOTHER APPROACH:

TRADE-OFF TRAINING ERROR WITH
"MODEL COMPLEXITY"

DEFINE ANOTHER POTENTIAL FUNCTION ϕ

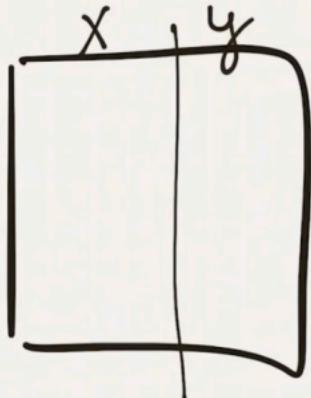
$\phi: \text{TREES} \rightarrow \mathbb{R}$ GIVEN A TRAINING SET S

$$\phi(T) = \underbrace{\text{TRAINING ERROR ON } S}_{\text{MINIMIZE}} + \alpha \cdot \frac{\text{SIZE}(T)}{|S|}$$

MINIMIZE ϕ

HYPERPARAMETER

ANOTHER APPROACH: MDL "MINIMUM DESCRIPTION LENGTH PRINCIPLE"



SOME NUMBER OF BITS NEEDED
TO ENCODE S.
 $m \cdot (n+1)$

BUILD A TREE T

LET'S SAY T IS CORRECT ON 90% OF S
AND INCORRECT ON 10%.

WE CAN ENCODE S USING #BITS(T) + #BITS TO
ENCODE THAT 10%. WE GOT
WRONG.

PAC MODEL OF LEARNING

THERE IS A DISTRIBUTION \mathcal{D} ON $\{0,1\}^n$ (\mathbb{R}^n)

FUNCTION CLASS $\mathcal{C} = \{$ DECISION TREES OF SIZE $S\}$

LEARNER (RUNS IN POLYNOMIAL-TIME)

Fix $c \in \mathcal{C}$

c IS THE UNKNOWN DEC. TREE
WE WANT TO LEARN

RECEIVES

(x', y') $x \text{ and } y = c(x)$

(x'', y'') $y'' = c(x'')$

\vdots
 (x^m, y^m) GOAL: OUTPUT $h \in \mathcal{C}$

LEARNER
SHOULD BE
EFFICIENT
(n, S)

$\Pr_{x \sim \mathcal{D}} [h(x) \neq c(x)] \leq \epsilon = .01$

X

"OVER THE DRAWS FROM Ω "

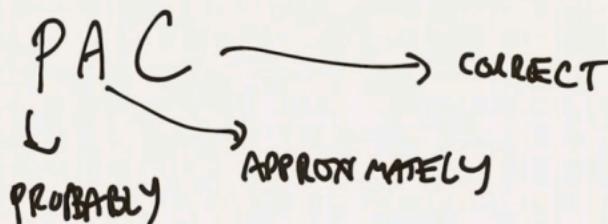
0 5 10

WITH PROBABILITY AT LEAST $1 - \delta$

THE LEARNER SHOULD OUTPUT A HYPOTHESIS

$$h \text{ s.t. } \Pr_{x \sim \Omega} [h(x) \neq c(x)] \leq \underline{\epsilon}.$$
RUN-TIME = polynomial $(\frac{1}{\epsilon}, \frac{1}{\delta}, n, s)$.IMAGINE LEARNER REQUESTS x^1, \dots, x^m
 $x^1 = \dots = x^m$

PAC



CORRECT L. VALIANT [1984]

PROBABLY APPROXIMATELY

W^XHEN CAN WE LEARN A FUNCTION CLASS?
WHAT FUNCTION CLASSES CAN WE PAC LEARN?

GIVE LEARNER AN ALGORITHM A

A : TRAINING SETS \rightarrow DECISION TREES

A(S) OUTPUT A TREE T THAT IS CONSISTENT WITH S. SIZE OF T

IS GOING TO BE AT MOST S.

A ALWAYS OUTPUTS A CONSISTENT HYPOTHESIS FROM C

GIVEN ANY TRAINING SET (ASSUMING THERE IS ONE)

Q: GIVEN ALGORITHM A, HOW CAN WE PAC LEARN C?

- DRAW SUFFICIENTLY MANY TRAINING POINTS = "S" $\begin{pmatrix} x \\ y \end{pmatrix}$
- USE A TO FIND \hat{c}_0 CONSISTENT WITH S.
- OUTPUT \hat{c}_0 .

Q: How large should S BE?

MARBLE GAME

TWO JARS

JAR 1

ALL BLUE
MARBLES

JAR 2

90% RED MARBLES
10% BLUE MARBLES.

GOAL: FIGURE OUT IF YOU'VE BEEN GIVEN
JAR 1 OR JAR 2

You RECEIVE A RANDOM ELEMENT OF JAR
ANY-TIME you want.

→ PICK RANDOM MARBLE FROM JAR

CASE 1: MARBLE IS RED JAR 2

ELSE 2: MARBLE IS BLUE.

PICK AT MOST 100 MARBLES.

WHAT IS THE PROBABILITY OF FAILURE?

PROB OF FAILURE is $(.1)^{100}$ ← FAILURE
PROB
"/S- PARAMETER"

Let's RETURN TO PAC LEARNING

- DRAW MANY SAMPLES S
- RUN A
- OUTPUT CLASSIFIER C THAT IS
CONSISTENT WITH S GIVEN FROM A.

WHAT IS THE PROBABILITY THIS PROCEDURE
FAILS? $\leq S$

BAD EVENT:

OUTPUT c THAT IS CONSISTENT WITH S
BUT HAS TRUE ERROR $> \epsilon$.

$\Pr[\text{BAD EVENT}] ?$

ENUMERATED ALL FUNCTIONS IN $\mathcal{C} = \{c_1, \dots, c_N\}$

FIX c_1 ASSUME c_1 HAS TRUE ERROR $> \epsilon$

WHAT IS $\Pr_S [c_1 \text{ is consistent with } S] ? \leq (1-\epsilon)^{|S|}$

FIX c_2 " c_2 HAS " " $> \epsilon$

$\Pr_S [c_2 \text{ is consistent with } S] ? \leq (1-\epsilon)^{|S|}$

o o c

FOR EVERY c_i (WITH ERROR $> \epsilon$)
 $\Pr_{S} [c_i \text{ is consistent on } S] \leq (1-\epsilon)^{|S|}$ ↴ BAD

Q: RANDOMLY FORM S , WHAT IS THE PROB
 THERE EXISTS A FUNCTION $c \in C$ WHOSE ERROR
 $> \epsilon$ AND IS CONSISTENT WITH S ?

UNION BOUND) $A, B \quad \Pr[A \cup B] \leq \Pr[A] + \Pr[B]$

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

$$\text{SOLVE FOR } |S|: \quad (1-x) \approx e^{-x}$$

x

o ↪ ↵

$$|C| \cdot (1-\epsilon)^{|S|} \leq \delta$$

$$|C| \cdot e^{-\epsilon|S|} \leq \delta \quad \underline{(1+xze^x)}$$

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

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

IF YOU CHOOSE
TRAINING POINTS
LARGER THAN $\log\left(\frac{|C|}{\delta}\right)$

$$|S| > \frac{\log\left(\frac{|C|}{\delta}\right)}{\epsilon}$$

ITEM WITH PROB? δ ,
FUNCTION OUTPUT
 c IS $1-\epsilon$
ACCURATE.

x

o s c

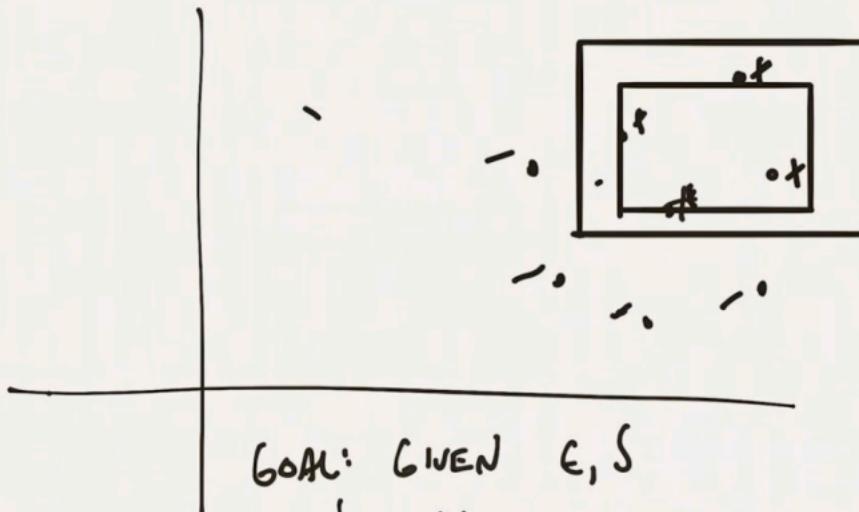
SUGGESTS A

"CONSISTENT HYPOTHESIS"

APPROACH TO LEARNING.

PAC - LEARNING AXIS - PARALLEL RECTANGLES.

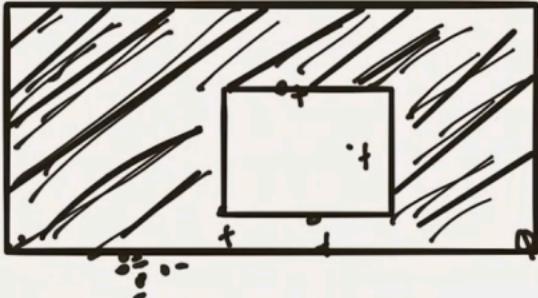
WE ARE WALKING IN 2-DIMENSIONS



GOAL: GIVEN ϵ, δ
 OUTPUT h THAT IS ϵ -ACCURATE
 WITH PROB $> 1 - \delta$

- Labeled + if the point is INSIDE c_1 , axis-p rect.

- if the point is OUTSIDE c_1 , axis-p rect.



CLAIM: TIGHTEST FITTING RECTANGLE WORKS!

QUESTION: HOW LARGE TO CHOOSE IS /

BAD EVENT: TIGHTEST FITTING RECTANGLE IS SMALL i.e. LOTS OF PROB MASS EXISTS

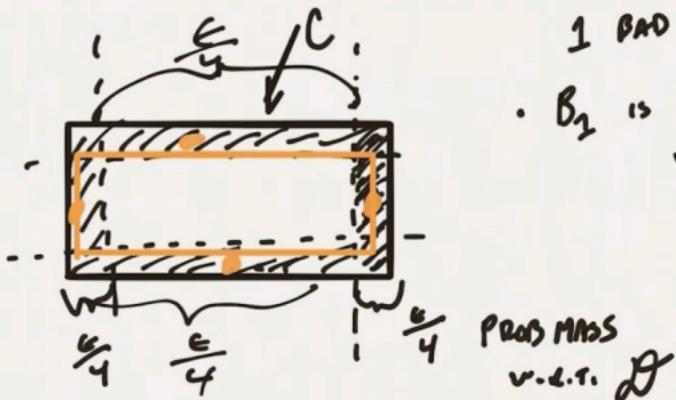
OUTSIDE OF h .
HOW DO WE bound THE PROBABILITY THIS HAPPENS?

ANALYZE $h^{\text{tightest fitting}}$ AND SAY SOMETHING ABOUT h VS ALL RECTANGLES THAT ARE LARGE BUT CONTAIN h .

X

O S C

- B_3 - WE SEE NO POINT IN LEFT STRIP



- 1 BAD EVENT
- B_3 IS THE EVENT WE SEE NO POINT IN RIGHT STRIP

- B_4 IS EVENT WE SEE NO POINT IN TOP STRIP.

- B_2 IS THE EVENT WE SEE NO POINT IN BOTTOM STRIP.

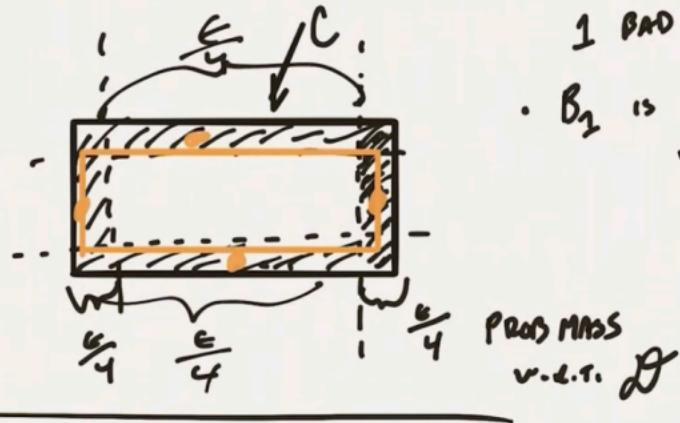
CLAIM • IF NEITHER B_1, B_2, B_3, B_4 OCCUR THEN

L, RIGHTEST FITTING RECT IS E-ACCURATE.

X

0 5 0

- B_3 - WE SEE NO POINT IN LEFT STRIP



- B_3 IS THE EVENT WE SEE NO POINT IN LEFT STRIP
- B_2 IS THE EVENT WE SEE NO POINT IN RIGHT STRIP

- B_1 IS EVENT WE SEE NO POINT IN TOP STRIP.

- B_1 IS THE EVENT WE SEE NO POINT IN BOTTOM STRIP.

CLAIM. IF CHOOSE m RANDOM SAMPLES $\Pr[B_1]$

$$\Pr[B_1] \leq \left(1 - \frac{1}{4}\right)^m$$

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

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

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

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

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

$$\begin{aligned}-\frac{\epsilon m}{4} &\leq \log\left(\frac{\delta}{4}\right) \\ \Rightarrow m &\geq \frac{4 \cdot \log\left(\frac{\delta}{4}\right)}{\epsilon}\end{aligned}$$

h) TRAPEZOID RULE, WILL BE ϵ -ACCURATE w.r.t. $\int - f$

x

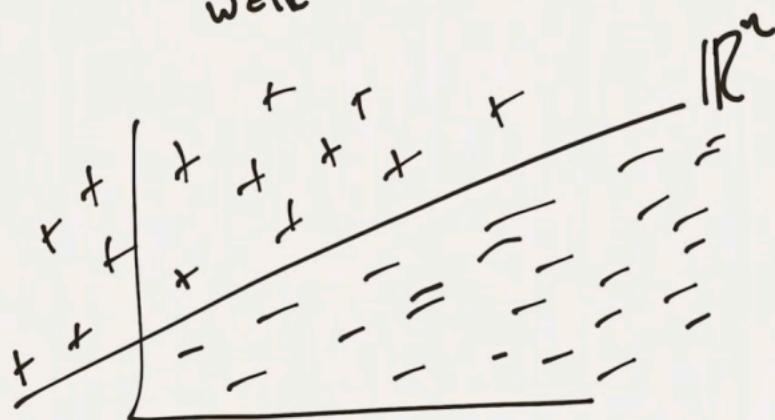
o s c

$\mathcal{C} = \{\text{HALFSPACES}\}$

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

$w \in \mathbb{R}^n$ $x \in \mathbb{R}^n$

f is BOOLEAN
0 or 1



PAC LEARNING
ALGORITHMS
FOR HALFSPACES.

ONE APPROACH FOR LEARNING HALFSPACES

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

IS UNKNOWN

$\theta \in \mathbb{R}$

UNKNOWN

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

GIVEN DRAWS FROM \mathcal{D} $(x_i, f(x))$

$$(01010, \text{pos}) \rightarrow w_2 + w_4 > \theta \quad w_i \in \mathbb{Z}_{\geq 0}$$

$$(0110, \text{neg}) \rightarrow w_2 + w_3 \leq \theta \quad \text{SOME BOUNDED RANGE.}$$

EACH LABELED EXAMPLE \rightarrow LINEAR INEQUALITY.

SYSTEM OF LINEAR INEQUALITIES.

CAN WE FIND A CONSISTENT HYPOTHESIS?

GENERAL-PURPOSE TOOL CALLED LINEAR PROGRAMMING

CROSS-VALIDATION

- HOLD-OUT APPROACH FOR TESTING/APPROXIMATING THE TRUE ERROR OF A CLASSIFIER
- LET'S ASSUME CLASSIFICATION; SO HYPOTHESIS h IS GOING TO OUTPUT $\{0, 1\}$
 $\{-1, +1\}$ VALUES.
"HOLD-OUT"
1. LEAVE SOME PART OF TRAINING SET OUT DURING TRAIN TIME.
2. TEST CLASSIFIER ON THIS HELD OUT SET.

X 0 5 C

1st INEQUALITY MARKOV'S INEQUALITY

let X be R.V. THAT TAKES ON ONLY POS VALUES.

$$\Pr [X \geq k \cdot \mathbb{E}[x]] \leq \frac{1}{k}$$

CHEBYSHEV'S INEQUALITY

$$\text{VAR}(x) = \mathbb{E}[(x - \mathbb{E}[x])^2] \quad \underline{\mathbb{E}[x] = \mu}$$

$$\sqrt{\text{VAR}(x)} = \text{STANDARD DEVIATION}(x) = \underline{\sigma}$$

$$\Pr [|x - \mu| > t \cdot \sigma] \leq \frac{1}{t^2}$$

CHERNOFF BOUND

$$X_1, X_2, \dots, X_n \quad \mathbb{E}[X_i] = p \quad X_i \in \{0, 1\}$$

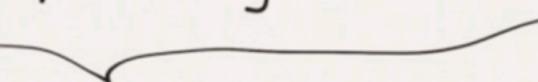
$$S = \sum_{i=1}^n X_i \quad \mu = \mathbb{E}[S] = p \cdot n$$

$$\mathbb{E}[X_1 + \dots + X_n] = p \cdot n$$

$$\Pr[S > \mu + \delta_n] \leq e^{-2n\delta_n^2}$$

$$\Pr[S < \mu - \delta_n] \leq e^{-2n\delta_n^2}$$

$$\Rightarrow \Pr[|S - \mu| > \delta_n] \leq 2 \cdot e^{-2n\delta_n^2}$$



X 0 5 C

APPLY THE CHERNOFF BOUND TO THE CASE
OF ESTIMATING THE TRUE ERROR OF A
CLASSIFIER

HOLD-OUT SET S WE'LL SAY $|S|=n$

FIX h (GENERATED USING SOME WD TRAINING SET)

RECALL Ω , S IS A SAMPLE DRAWN FROM Ω
IND OF TRAINING SET

$$z = \Pr_{x \in \Omega} [h(x) \neq c(x)]$$

UNKNOWN FUNCTION TRYING TO LEARN

Let X_i be r.v. EQUALS 1 if h is INCORRECT
ON THE i^{th} element
& S .

0 if h is correct on the i^{th}
element of S .

X_1, \dots, X_n

$X_i = \begin{cases} 1 & \text{if } h \text{ is incorrect on } i^{\text{th}} \\ 0 & \text{otherwise.} \end{cases}$

$$S = \sum_{i=1}^n X_i$$

$$\mathbb{E}[S] = n \cdot p \leftarrow \text{true error of } h.$$

$$p = \mathbb{E}[X_i] = \mathbb{E}[X_1] = \mathbb{E}[X_n]$$

$$\Pr[|S - n \cdot p| > \delta_n] \leq 2 e^{-2n\delta^2}$$

(RECALL p is TRUE ERROR OF CLASSIFIER h).

$$S = .1$$

$$\Pr[|S - n \cdot p| > .1n] \leq 2 e^{-\frac{2n}{100}}$$

HOW LARGE TO CHOOSE
 n BEFORE THIS

QUANTITY BECOMES SMALL?

$$e^{-\frac{2n}{100}} < \frac{\alpha}{2}$$

$$\alpha$$

$$\text{IF } |S - n \cdot p| \leq .1n$$

$$\frac{-2n}{100} < \log\left(\frac{\alpha}{2}\right) \Rightarrow n > 50 \cdot \log\left(\frac{\alpha}{2}\right)$$

\Rightarrow ERROR RATE
ON S IS
WITHIN .1 OF
TRUE ERROR RATE.

HOLD-OUT SET IS SOMEWHAT EXPENSIVE

- DATA IS EXPENSIVE
- IF WE WANT TO TRY OUT MULTIPLE METHODS FOR GENERATING CLASSIFIERS, WE QUICKLY LOSE CONFIDENCE IN OUR ESTIMATES.

CROSS-VALIDATION

- TRAIN USING FOLDS 2...FOLD K
- TEST ON FOLD 1



x

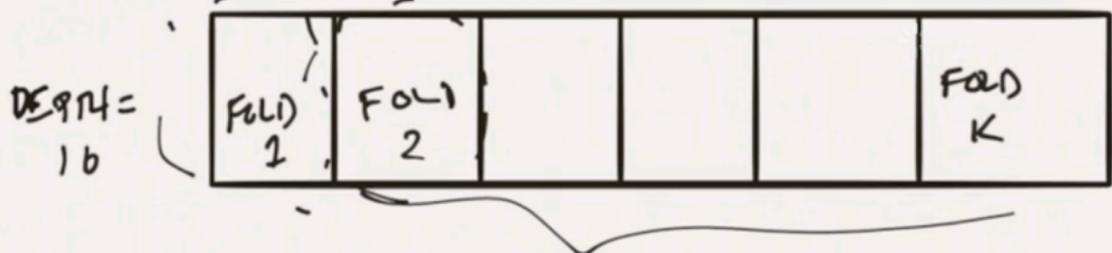
0 5 c

LET'S GO BACK TO DECISION TREES

TRAINING SET S

- SHOULD I BUILD A DECISION TREE OF DEPTH 10
 " " " " " " 15

WE DECIDE USING CROSS VALIDATION



X

O D C

LET'S GO BACK TO DECISION TREES

TRAINING SET S

- SHOULD I BUILD A DECISION TREE OF DEPTH 10
 " " " " " " 15

WE DECIDE USING CROSS VALIDATION

