

PAC LEARNING

11

PROBABLY APPROXIMATELY CORRECT

CLASSIFICATION:

INPUT

m data points

$(x_1, y_1), \dots, (x_m, y_m)$

x_i chosen i.i.d. from some distribution

$y_i = c(x_i)$

UNKNOWN CONCEPT
 $c: X \rightarrow Y$

GOAL

learn c

APPROACH

- consider a finite set of hypotheses \mathcal{H}
(decision trees of depth d)
- find h with

$$\text{err}_{\text{train}}(h) = 0$$

(we say that h is CONSISTENT with data)

WHAT IS THE PROBABILITY THAT

$$\text{err}_{\text{true}}(h) > \varepsilon ?$$

$$\Pr[\text{err}(h) \neq c(x)]$$

PROBABILITY OF A MISTAKE

[2]

BAD HYPOTHESES: w_1, \dots, w_K
 (TRUE ERROR $> \varepsilon$)

STEP I:

What is the probability that a fixed BAD hypothesis gets all data right?
 say w_3 :

$$\Pr[w_3 \text{ makes an error}] > \varepsilon \quad [\text{BAD HYPOTHESIS}]$$

$$\Pr[w_3 \text{ gets a single point right}] \leq 1 - \varepsilon$$

m points

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

STEP II:

Prob. that any bad hypothesis gets
all m points right

$$\leq K \cdot (1 - \varepsilon)^m \quad [\text{UNION BOUND}]$$

$$\leq |\mathcal{H}| \cdot (1 - \varepsilon)^m \leq 12e \cdot e^{-m\varepsilon} \quad [\text{by } 1 - x \leq e^{-x}]$$

WITH ENOUGH DATA, WE MANAGE
 TO EXCLUDE ALL BAD HYPOTHESES
 WITH HIGH PROBABILITY

THEOREM [Haussler 1988]

For a finite hypothesis space \mathcal{H} , m i.i.d.
 training examples, $0 < \varepsilon < 1$, and any w
 consistent w/ data:

$$\Pr[\text{err}_{\text{true}}(w) > \varepsilon] \leq 12e \cdot e^{-m\varepsilon}$$

3)

SUPPOSE WE WANT THIS PROBABILITY $< \delta$:

- How many examples suffice?

$$\text{WANT: } |\mathcal{H}| e^{-m\epsilon} \leq \delta$$

$$\text{REQUIRE: } m \geq \frac{\ln |\mathcal{H}| + \ln(1/\delta)}{\epsilon}$$

SAMPLE COMPLEXITY

NOTE:

- grows gently with $|\mathcal{H}|$ and amount of confidence $1-\delta$
- can be large if we want ϵ to be small

- What is the largest error we get after perfectly fitting m samples?

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

REWRITE ERROR BOUND:

$$\text{err}_{\text{true}}(h) \leq \frac{1}{\log_2 e} \cdot \frac{\underbrace{\# \text{bits to identify } h}_{\log_2 |\mathcal{H}|} + \log_2(1/\delta)}{m}$$

If \mathcal{H} finite, encode each h by its index: $1, \dots, |\mathcal{H}|$ - requires $\lceil \log_2 |\mathcal{H}| \rceil$ bits for each h .

Other coding schemes possible and the result still holds (\mathcal{H} not necessarily finite):

$$\text{err}_{\text{true}}(h) \leq \frac{1}{\log_2 e} \cdot \frac{\text{size}(h) + \log_2(1/\delta)}{m}$$

OCCAM'S RAZOR: SHORTER HYPOTHESES THAT EXPLAIN THE DATA HAVE BETTER GENERALIZATION

EXAMPLE:

41

DECISION TREES FOR n BINARY ATTRIBUTES x_0, \dots, x_{n-1}

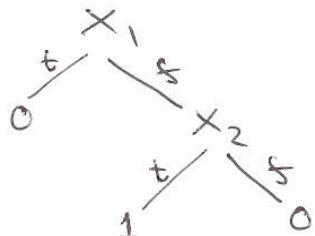
- define encoding recursively:

$$\text{ENCODE}(\text{leaf}) = \begin{cases} 00 & \text{if predicts 0} \\ 01 & \text{if predicts 1} \end{cases}$$

$$\text{ENCODE}\left(\begin{array}{c} x_i \\ t \diagup \diagdown s \\ \text{LEFT} \quad \text{RIGHT} \end{array}\right) = 1 \cdot \langle i \text{ in binary} \rangle$$

- $\text{ENCODE}(\text{LEFT})$
- $\text{ENCODE}(\text{RIGHT})$

SAY: ATTRIBUTES x_0, x_1, x_2, x_3



$\underbrace{101}_{x_1 \text{ pred. } 0} \underbrace{00}_{0} \underbrace{110}_{x_2 \text{ pred. } 1} \underbrace{01}_{1} \underbrace{00}_{0}$

k leaves, $k-1$ inner nodes: $\text{size}(h) = 2k + (k-1) \log n$
BITS

IF WE CAN FIT DATA WITH A TREE
WHICH HAS k leaves:

$$\text{err}_{\text{true}}(h) \leq \frac{1}{\log_2 e} \cdot \frac{2k + (k-1) \log_2 n + \log_2(V\delta)}{m}$$

x
 < 1 (BOUND NON-TRIVIAL)
 if $k < \frac{m}{\text{err}_{\text{true}}(V\delta)}$

LIMITATIONS OF PAC LEARNING

5

- consistency requirement:
 - what if concept $\notin \mathcal{H}$
 - or there is inherent noise in data?
 - coding can be cumbersome
 - e.g.: how to encode half-spaces?

DROP CONSISTENCY REQUIREMENT:

- m i.i.d. data points,
but y_i is not necessarily
a deterministic function of x_i

APPROACH

- find $w \in \mathbb{R}^d$ that makes fewest mistakes on the data

NOW WE CAN SHOW:

PROVED USING THE UNION BOUND (ON "BAD" EVENTS)
AND HOEFFDING'S INEQUALITY

[6]

Let $\mathcal{H} = \{h_1, \dots, h_{|\mathcal{H}|}\}$

BAD EVENT j : $j = 1, \dots, |\mathcal{H}|$

$$\text{err}_{\text{true}}(h_j) > \text{err}_{\text{train}}(h_j) + \varepsilon$$

$$\Pr[\text{err}_{\text{true}}(h_j) - \text{err}_{\text{train}}(h_j) > \varepsilon] \leq 33?$$

=

DETOUR: CLASSIFICATION AS COIN FLIPS

INCORRECT $\hat{\triangleq}$ heads

CORRECT $\hat{\triangleq}$ tails

$\text{err}_{\text{true}}(h_j)$ $\hat{\triangleq}$ PROBABILITY OF HEADS, θ

$\text{err}_{\text{train}}(h_j)$ $\hat{\triangleq}$ MAXIMUM LIKELIHOOD ESTIMATE
 $\hat{\theta} = \frac{\# \text{misclassifications}}{m}$

HOEFFDING INEQUALITY:

$$\Pr[\theta - \hat{\theta} > \varepsilon] \leq e^{-2m\varepsilon^2}$$

\nwarrow RANDOM VARIABLE

GENERAL HOEFFDING BOUND:

z_1, z_2, \dots, z_m indep. samples from
a distribution over $[0, 1]$
with mean μ

$$\Pr[\mu - \underbrace{\frac{1}{m} \sum_{i=1}^m z_i}_{\text{empirical average}} > \varepsilon] \leq e^{-2m\varepsilon^2}$$

4

SUMMARIZING

INCONSISTENT HYPOTHESIS MODEL:

With prob $1-\delta$, for all $w \in \mathcal{H}$:

$$\text{err}_{\text{true}}(w) \leq \text{err}_{\text{train}}(w) + \sqrt{\frac{\ln(2\ell) + \ln(1/\delta)}{2m}}$$

or similarly as before

$$\text{err}_{\text{true}}(w) \leq \text{err}_{\text{train}}(w) + \sqrt{\frac{1}{\log_2} \cdot \frac{\text{size}(w) + \log_2(\ell/\delta)}{2m}}$$

want small \downarrow \downarrow

fixed m, δ	for complex hypotheses	small	large
	for simpler hypotheses	large	small

“bias”

“variance”

DECREASES TO ZERO
w/ MORE DATA

NOTE: a worse dependence

on # samples (need more samples)
than the CONSISTENT CASE