

RECAP

• PAC Learning:

- all target concepts and distributions, accuracy & confidence parameter.
- allow learning algorithm to output a hypothesis from a different (typically larger) hypothesis class.
[If we insist on the output hypothesis to be from the concept class being learned, we will call it PROPER LEARNING]
- efficient if the running time is polynomial in $\gamma_\varepsilon, \gamma_\delta, n$ and size(C).
- notion of "Polynomially evaluable" hypothesis classes

RESULTS:

- (Proper) PAC-learning algorithms for learning CONJUNCTIONS and axis-aligned RECTANGLES (in \mathbb{R}^2).
- Hardness of proper PAC-learning of 3-term DNF.
- PAC-learning algorithm (improper) for learning 3-term-DNF by using a (proper) PAC learning algorithm for 3-CNF.

Consistent Learning

Explanatory model:

Given $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$ $x_i \in X$ (instance space)
 $y_i \in \{0, 1\}$.

"Can we find a hypothesis that is consistent with all of the given data?"

if $\underline{x} = \underline{x}_1$, then output y_1 .	
else if $\underline{x} = \underline{x}_2$ " " y_2	
?	
- - $\underline{x} = \underline{x}_m$ - - y_m	
else output 0.	

$$X = \{0, 1\}^n$$

Exercise: Show that you can find a consistent hypothesis that is a DNF formula with at most m terms.

Occam's Razor: Find a simple an explanation as possible.

"Shortest" hypothesis that is consistent with the given data.

For learning, "short enough" hypothesis will suffice.

"Explanation is shorter than the give data!"

Looking for succinct explanations.

[Connections to Kolmogorov complexity; "Minimum Description Length Principle"; Bayesian interpretation]

Consistent Learner: We say that there exists a consistent learner for a concept class C using a hypothesis class H , if $\forall n \geq 1, \forall c \in C_n, \forall m \geq 1$, given $(x_1, c(x_1)), (x_2, c(x_2)), \dots, (x_m, c(x_m))$, the learner outputs $h \in H_n$, such that $h(x_i) = c(x_i)$ for $i = 1, \dots, m$.

We say that there is an efficient consistent learner if the algorithm runs in time polynomial in $n, \text{size}(c), m$.

(EMPIRICAL RISK MINIMIZATION - using the zero-one loss)

(Occam's Razor Theorem)

Theorem: Let C be a concept class and H a hypothesis class. Let L be a consistent learner for C using h , and furthermore L is efficient. Then $\forall n \geq 1, \forall D$ over $X_n, \forall c \in C_n$, if L is given a sample of size m drawn from $EX(c, D)$ such that

$$m \geq \frac{1}{\varepsilon} (\log |H_n| + \log \frac{1}{\delta})$$

then L is guaranteed to output a hypothesis $h \in H_n$, that with probability at least $1 - \delta$, satisfies $\underline{\text{err}}(h) \leq \varepsilon$. (If $\log |H_n|$ is polynomial in n and $\text{size}(c)$ then C is efficiently PAC learnable using H .)

Proof: Call hypothesis $h \in H_n$ 'bad' if $\text{err}(h) \geq \varepsilon$.

A_n event that m random examples from $EX(C, D)$ are consistent with h .

if h is bad, then $P(A_n) \leq (1 - \varepsilon)^m$.

$$\mathcal{E} = \bigcup_{h \text{ bad}} A_h$$

$$P(\mathcal{E}) \leq \sum_{h \text{ bad}} P(A_h) \leq |H_n| \cdot (1 - \varepsilon)^m. \quad \begin{matrix} \text{Want this to} \\ \text{be} \leq \delta \end{matrix}$$

union bound

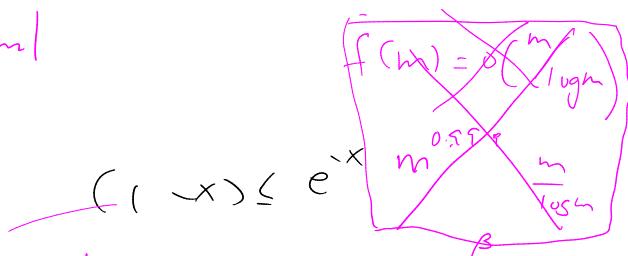
want

$$|H_n| \cdot (1 - \varepsilon)^m \leq \delta$$

$|H_{n,m}|$

It suffices to show that

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



equivalent to

$$m \geq \frac{1}{\varepsilon} \left(\underbrace{\log |H_n|}_{\log |H_{n,m}|} + \log \frac{1}{\delta} \right)$$

$$\log |H_{n,m}| \leq \text{poly}(n) \cdot m$$

If ε does not occur, then the consistent learner must output $h \in H_n$, s.t. $\text{err}(h) \leq \varepsilon$. \square

For CONJUNCTIONS: Provided $m \geq \frac{2n}{\varepsilon} \log \left(\frac{2n}{\delta} \right)$ the algorithm "works". [Proved last week].

$$= \frac{2n \log(2n)}{\varepsilon} + \frac{2n \log(\frac{1}{\delta})}{\varepsilon}$$

H_n CONJUNCTIONS, $|H_n| = 3^n$, $\log |H_n| = n \log 3$.

If $m \geq \frac{1}{\varepsilon} (n \log 3 + \log \frac{1}{\delta})$ the consistent learner gives a PAC learner.

3-term DNF

3-term DNF of size n : $2^{\binom{n}{2}}$

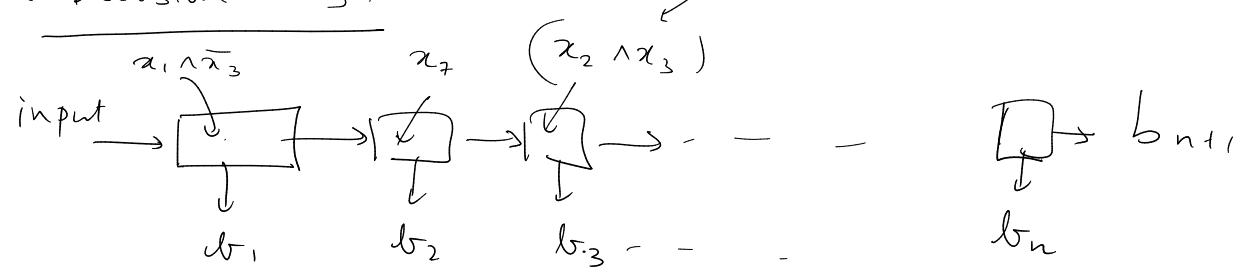
$\log - \Theta(n)$

3-CNF formulae on n variables: $2^{\frac{8}{3}n^3}$ $\rightarrow \Theta(n^3)$.

[Without worrying about computational complexity, any 3-term DNF formula consistent with $\frac{1}{\varepsilon} (c \cdot n + \log \frac{1}{\delta})$ examples has a PAC-guarantee.]

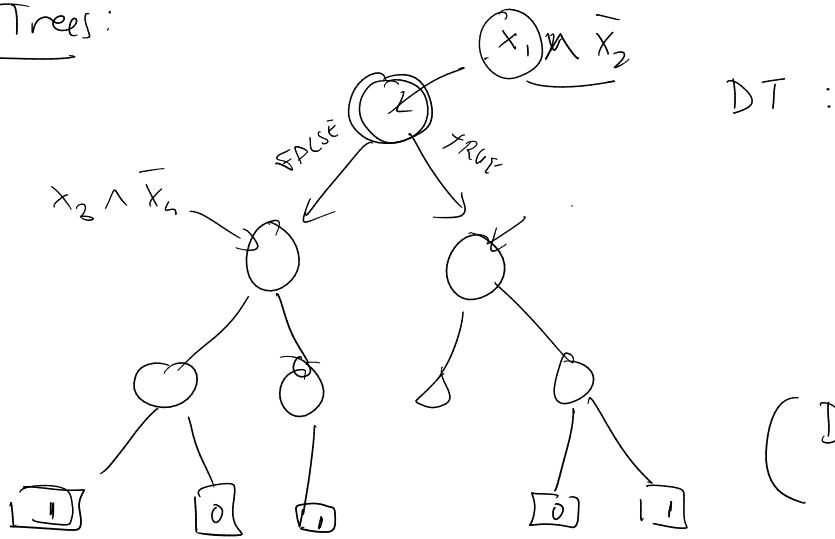
(If we output a 3-CNF formula, Occam's Razor Theorem will tell us that $\frac{1}{\varepsilon} (c \cdot n^3 + \log \frac{1}{\delta})$ examples are sufficient)

k -Decision Lists:



$b_i \in \{0, 1\}$ (PAC-Learnable)

Decision Trees:



(Decision trees
of poly size
not PAC-learnable)