

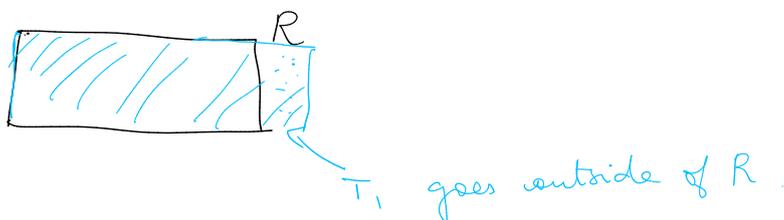
PAC (Take I) Learning: Let \mathcal{C} be a concept class over X . We say that \mathcal{C} is probably approximately learnable if there exists a learning algorithm L with the following property: For every target concept $c \in \mathcal{C}$, for every distribution D over X , for all $0 < \epsilon < 1/2$, $0 < \delta < 1/2$ if L is given access to the example oracle $EX(c, D)$ and inputs ϵ, δ , then with probability at least $1 - \delta$, L outputs a hypothesis $h \in \mathcal{C}$ satisfying $\text{err}(h) \leq \epsilon$.

- The probability is taken over random draws from $EX(c, D)$ as well as any internal randomization of L .

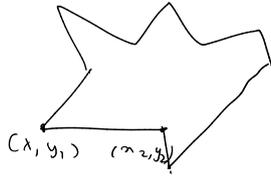
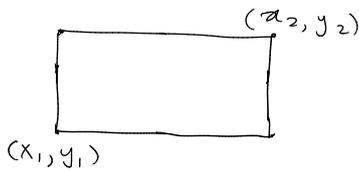
- If L runs in time polynomial in $1/\epsilon, 1/\delta$, we say that \mathcal{C} is efficiently PAC learnable.

Theorem: The concept class \mathcal{C} of axis-aligned rectangles in \mathbb{R}^2 is efficiently PAC learnable.

- We saw a proof of the above last time, with some "missing" details.
- What if D is not "smooth" enough?
- What if $\mathbb{P}_{x \sim D} [x \in R] < \epsilon/4$, where R is the target rectangle, in that case T_i is not contained in R !



Representation size



"complexity of the target function/
concept class!"

"Can store real numbers at unit cost?"

Boolean Functions:

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

$f: X \rightarrow \{0, 1\}$

TRUTH TABLES

2^n entries	00...0	1
	0...1	0
	⋮	

Every boolean function
requires space 2^n .

- Boolean circuits, AND, OR and NOT gates.

- Disjunctive Normal Form (DNF)

$$f = \underbrace{(x_1 \wedge \bar{x}_3 \wedge \dots \wedge x_n)}_{\{0, 1\}^n} \vee (x_3 \wedge \bar{x}_7 \wedge \dots) \vee (x_1 \wedge x_n \wedge \dots)$$

$$f = x_1 \oplus x_2 \oplus \dots \oplus x_n \quad \left. \begin{array}{l} \text{PARITY} \\ \text{FUNCTION} \end{array} \right\} \begin{array}{l} 1 \text{ if odd no. of } x_i \text{ are } 1 \\ 0 \text{ otherwise.} \end{array}$$

- Circuit of depth $\log(n)$ and size $O(n \log n)$

- Need a DNF formula of size 2^{2^n} to represent parity.

- How many boolean functions? 2^{2^n}

- We will use representation to restrict the class of Boolean functions that we will aim to learn?

- Decision trees
- Neural networks

Representation Scheme:

A representation scheme for a concept class C is a mapping $R: \Sigma^* \rightarrow C$, where Σ^* is a finite alphabet.

($R: (\Sigma \cup \mathbb{R})^* \rightarrow C$ if we need real numbers).

$\text{size}: \Sigma^* \rightarrow \mathbb{N}$ (length of the string)

$$\text{size}(c) = \min \{ \text{size}(s) \mid R(s) = c \}$$

Instance Size:

Typically we will consider $X = \{0,1\}^n$, $X = \mathbb{R}^n$, and "n" will be the size parameter with the instance space.

$X = \langle X_n \rangle_{n \geq 1}$, $X = \bigcup_{n \geq 1} X_n$ where X_n represents instances of size n .

C_n is a concept class over X_n ,

$C = \bigcup_{n \geq 1} C_n$ a concept class over X .

PAC Learning (Take II): Let C_n be a representation/concept class over X_n . Let

$C = \bigcup_{n \geq 1} C_n$, & $X = \bigcup_{n \geq 1} X_n$. We say C is PAC learnable if there exists

a learning algorithm L , s.t. $\forall n \in \mathbb{N}$, $\forall C \in C_n$, \forall distribution D over X_n ,

$\forall 0 < \epsilon < 1/2$, $\forall 0 < \delta < 1/2$ if L is given access to the example

oracle $E_{X, C, D}$ and inputs $\epsilon, \delta, \text{size}(C)$, outputs w.p. $\geq 1 - \delta$,

a hypothesis $h \in C_n$ satisfying $\text{err}(h) \leq \epsilon$. (Exercise: Drop this)

(Probability over all sources of randomness)

C is efficiently PAC learnable if the running time of L is

polynomial in \underline{n} , $\underline{\text{size}(C)}$, $\frac{1}{\epsilon}$, $\frac{1}{\delta}$.

Example: Learning Boolean Functions

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

z_i : boolean variable

x

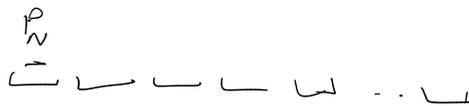
CONJUNCTION: $z_1 \wedge \bar{z}_3 \wedge z_7 \wedge \dots \wedge z_n$

$$f: X \rightarrow \{0,1\}$$

C_n : The class of all conjunctions on n boolean variables!

$$|C_n| = 3^n$$

$$\forall c \in C_n, \text{size}(c) \leq \underline{2n}$$



Q: Can we design an algorithm that learns conjunctions in time polynomial in $n, 1/\epsilon, 1/\delta$?

Examples: $(\underline{x}_1, y_1), (\underline{x}_2, y_2), \dots, (\underline{x}_m, y_m)$

$$\underline{x}_i \in \{0,1\}^n, y_i \in \{0,1\}$$

\exists conjunction c , st. $c(\underline{x}_i) = y_i$

$(11011, 1), (100110, 0), \dots$

ALGORITHM:

$$h = z_1 \wedge \bar{z}_1 \wedge z_2 \wedge \bar{z}_2 \wedge \dots \wedge z_n \wedge \bar{z}_n$$

for $i=1, \dots, m$, {

if $(y_i = 1)$ {

we drop each \bar{z}_j from h if $x_{ij} = 1$

" " " z_j from h if $x_{ij} = 0$

} }

Return h

Remarks: • Algorithm only uses positive examples.

• if $h(\underline{x}) = 1$ then $c(\underline{x}) = 1$

Let's look at a literal l ,

$$p(l) = \sum_{\underline{x} \sim D} \mathbb{P} [c(\underline{x}) = 1 \wedge \text{the literal } l \text{ evaluates to 0 in } \underline{x}]$$

l is "bad" if $p(l) \geq \epsilon/2n$.

A_ϵ : Event that after drawing m examples the literal l is not eliminated from h .

$$P(A_\ell) \leq \left(1 - \frac{\epsilon}{2n}\right)^m$$

$$P(\text{err}(h) \geq \epsilon) \leq P\left(\bigcup_{\substack{\ell \text{ bad} \\ \text{literal}}} A_\ell\right) \leq \sum_{\substack{\ell \text{ bad} \\ \text{literal}}} P(A_\ell) \leq \underbrace{2n \cdot \left(1 - \frac{\epsilon}{2n}\right)^m}_{\text{want}} \leq \delta$$

If we eliminate all the bad literals

$$\text{err}(h) \leq \sum_{\substack{\ell \text{ "good"} \\ \text{literals}}} p(\ell) \leq 2n \cdot \frac{\epsilon}{2n} = \epsilon$$

Confidence Analysis:

$$\text{want } 2n \cdot \left(1 - \frac{\epsilon}{2n}\right)^m \leq \delta$$

$$\text{Sufficient if } 2n \cdot e^{-\frac{m\epsilon}{2n}} \leq \delta, \quad \left(1 - x \leq e^{-x}\right)$$

$$\text{" " } m \geq \frac{2n}{\epsilon} \log\left(\frac{2n}{\delta}\right). \quad \square$$

Theorem: The class of conjunctions is efficiently PAC (Take II)
learnable.