

3.5 Exercises

① Monotonicity of Sample Complexity

Let H be a hypothesis class for a binary classification task. Suppose that H is PAC learnable and its sample complexity is given by $m_H(\cdot, \cdot)$.

Show that m_H is monotonically nonincreasing in each of its parameters.

In other words, show that:

$$\text{Given } \begin{cases} \delta \in (0, 1) \\ 0 < \epsilon_1 \leq \epsilon_2 < 1 \end{cases}$$

$$\text{We have } m_H(\epsilon_1, \delta) \geq m_H(\epsilon_2, \delta)$$

Similarly, show that:

$$\text{Given } \begin{cases} \epsilon \in (0, 1) \\ 0 < \delta_1 \leq \delta_2 < 1 \end{cases}$$

$$\text{We have } m_H(\epsilon, \delta_1) \geq m_H(\epsilon, \delta_2)$$

Solution:

• The proof follows from the definition, we restated the PAC learnability definition with realizability assumption (for simplicity), with this

$$\text{adjustments: } \begin{cases} 0 < \epsilon_1 \leq \epsilon_2 < 1 \\ \text{leads to } m_1 \stackrel{\text{def}}{=} m_H(\epsilon_1, \delta) \geq m_H(\epsilon_2, \delta) \stackrel{\text{def}}{=} m_2 \end{cases}$$

• For every $\delta \in (0, 1)$, $0 < \epsilon_1 \leq \epsilon_2 < 1$ and distribution D over X , when running the learning algorithm on $m \geq m_1 \stackrel{\text{def}}{=} m_H(\epsilon_1, \delta) \geq m_H(\epsilon_2, \delta) \stackrel{\text{def}}{=} m_2$ i.i.d. examples generated by D , the algorithm returns a hypothesis h such that, with probability at least $1 - \delta$:

$$L_D(h) \leq \epsilon_1 \leq \epsilon_2$$

By the minimality of m_2 , we conclude that $m_2 \leq m_1$

human language explanation

human language explanation:

- Given sample size $m \geq m_1$, we will have $L_D(h) \leq \epsilon_1$
- Since $\epsilon_1 \leq \epsilon_2$, any hypothesis h that is $L_D(h) \leq \epsilon_1$, also satisfy $L_D(h) \leq \epsilon_2$. And since $m \geq m_1$ is required in the first case, it also satisfy the second (looser) condition, $m \geq m_1 \geq m_2$
- "Minimality of m_2 ": means that m_2 is the smallest possible sample size such that accuracy ϵ_2 is guaranteed.

② Show that a hypothesis class is PAC learnable (with realizability)

Let X be a discrete domain

$$\mathcal{H}_{\text{singleton}} = \underbrace{\{h_z : z \in X\}} \cup \underbrace{\{h^-\}}$$

$$\text{Each } z \in X, h_z(x) = \begin{cases} 1 & \text{if } x = z \\ 0 & \text{otherwise} \end{cases} \quad h^-(x) = 0 \quad \forall x \in X$$

for all x , nothing is special

for each z , there is exactly one h_z
that label it "special"

Realizability assumption holds

that means at least 1 h in $\mathcal{H}_{\text{singleton}}$ satisfy $L_D(h) = 0$.

a) Describe an algorithm that implements the ERM rule for learning $\mathcal{H}_{\text{singleton}}$ in the realizable setup

- Since realizability holds, we need to come up with an algorithm s.t. $L_S(h) = 0$
- Let A be the algorithm that returns hypothesis h_S with the following property:

$$h_S = \begin{cases} h_x & \text{if } \exists x \in S \text{ s.t. } f(x) = 1 \\ h^- & \text{otherwise} \end{cases}$$

Clearly, $L_S(h_S) = 0$, and so A is ERM

b) Show that $H_{\text{singleton}}$ is PAC learnable. Provide an upper bound on the sample complexity

• Let D be distribution over X and $\epsilon \in (0,1)$

• Based on the definition of A in previous section:

• If $f = h^-$, then A return true hypothesis, so $L_{D,f}(h^-) = 0$

• Suppose $\exists x \in X$ s.t. $f(x) = 1$

Let $S|_X = (x_1, \dots, x_m)$ be instances of training set S

We try to upper bound $D^m(\{S|_X : L_{D,f}(h_S) > \epsilon\})$ → How bad training set S represent D

• If $\exists x \in S|_X$, then A returns true hypothesis, so $L_{D,f}(h_S) = 0$

• So the only scenario left for us to upper bound on is when

$\exists x \notin S|_X$

• Also, it can be proven that:

$$D(x) = L_{D,f}(h) \xrightarrow{\text{proof}}$$

$$\Rightarrow D(x) \leq \epsilon \text{ means } L_{D,f}(h) \leq \epsilon$$

$$\Rightarrow D(x) > \epsilon \text{ means } L_{D,f}(h) > \epsilon$$

$$\Leftrightarrow D(x') \leq 1 - \epsilon \text{ means } L_{D,f}(h) > \epsilon$$

$\forall x' \in X|_X$

• Combine 2 points above, we have:

$$\{S|_X : L_{D,f}(h) > \epsilon\} = \{S|_X : \forall x' \in S|_X, D(x') \leq 1 - \epsilon\}$$

And so:

$$D^m(\{S|_X : L_{D,f}(h) > \epsilon\})$$

$$= D^m(\{S|_X : \forall x' \in S|_X, D(x') \leq 1 - \epsilon\})$$

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

$$\leq e^{-\epsilon m}$$

• Let $\delta \in (0,1)$ s.t. $e^{-\epsilon m} \leq \delta$. We can conclude that:

$$m \geq \frac{\log(1/\delta)}{\epsilon}$$

Let's analyze possible values of true error on individual instance x :

• all negative case $h^- : \mathbb{E}_{x \sim D}(h) = 0$

• contain "special" instance:

• Correctly identify "special" instance:
 $\mathbb{E}_{x \sim D}(h_S) = 0$

• Wrongly identify "special" instance:
 $\mathbb{E}_{x \sim D}(h_S) = 1$

So possible values are 0 and 1:

By LOTF and definition of A :

$$L_{D,f}(h) = 1 \cdot D(x) + 0 \cdot D^c(x) = D(x)$$

Therefore, $H_{\text{singleton}}$ is PAC learnable with $m_{H_{\text{singleton}}} \leq \left\lceil \frac{\log(1/\delta)}{\epsilon} \right\rceil$

③ Another prove PAC learnable (with realizability)

Let $X = \mathbb{R}^2$, $Y = \{0, 1\}$

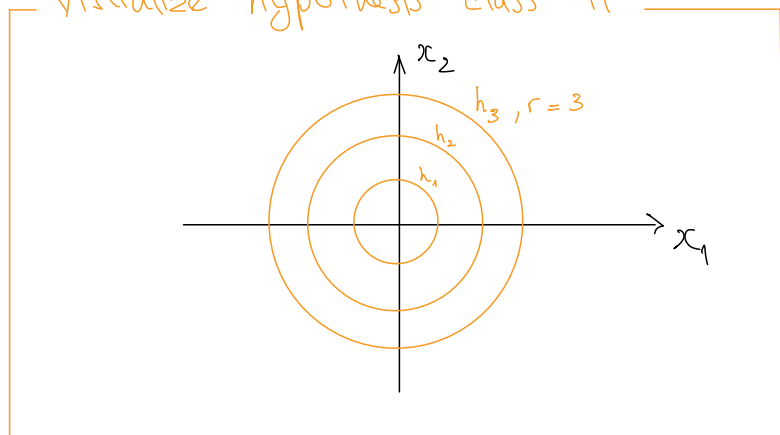
H be hypothesis class of concentric circles in the plane:

$H = \{h_r : r \in \mathbb{R}_+\}$, where $h_r(x) = 1$ if $\|x\| \leq r$

Prove that H is PAC learnable (assume realizability)

H sample complexity is bounded by $m_H(\epsilon, \delta) \leq \left\lceil \frac{\log(1/\delta)}{\epsilon} \right\rceil$

Visualize hypothesis class H



Come up with alg A that is ERM

Solution: Proving PAC learnability has 2 steps

Upper bound sample complexity

- Let D be the distribution over X , $\epsilon \in (0, 1)$ and f be target hypothesis
- Let A be algorithm that returns the smallest circle enclosing all the positive examples from the training set S , where

$C(S)$ be the circle returned
 $r(S)$ be the circle's radius
 $A(S): X \rightarrow Y$ is the hypothesis

We can easily see that $L_S(A(S)) = 0$, so A is ERM

Proof: Since A returns positive examples: $A(S)(x) = 1$

Realizability assumption: $\exists h^* \in H$ s.t. $h^*(x) = 1$

$\Rightarrow A(S) = h^*$, and so $L_S(A(S)) = 0$

conclude A is ERM

• By realizability assumption, $\exists h^* \in \mathcal{H}$ s.t. $L_{D,f}(h^*) = 0$

Let $\begin{cases} C^* & \text{be the circle corresponds to hypothesis } h^* \\ r^* & \text{be the corresponding radius} \end{cases}$

It can be proven that $C(S) \subseteq C^*$, where C^* is the circle enclosing all the positive examples.

Proof: $\begin{cases} C(S) & \text{is the smallest circle enclosing positive examples} \\ C^* & \text{is the circle enclosing positive examples} \end{cases}$

$$\Rightarrow C(S) \subseteq C^*$$

• Since $C(S) \subseteq C^*$, we can prove that $L_{D,f}(A(S)) = D(C^* \setminus C(S))$

Proof: $L_{D,f}(A(S)) = D(\{x \in X : A(S)(x) \neq f(x)\})$ <definition>

$$= D(\{x \in X : x \notin S|_x \text{ and } f(x)=1\})$$

$$= D(C^* \setminus C(S))$$

Why?

$L_{D,f}(A(S))$ consists of:

• $x \in S|_x$ s.t.
 $f(x) = 0$

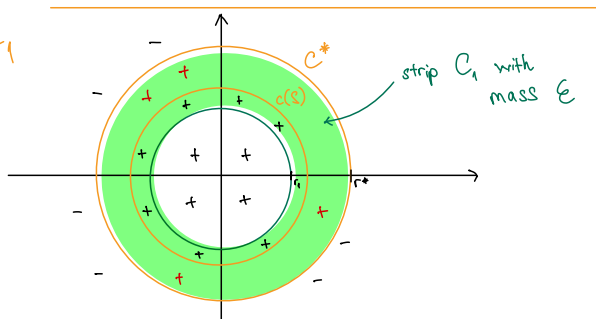
• $x \notin S|_x$ s.t.
 $f(x) = 1$

However, $C(S) \subseteq C^*$
means first case cannot
happen

• Let $r_1 \leq r^*$ be a number such that the probability mass of C_1 is ϵ .
where corresponding strip $C_1 = \{x \in \mathbb{R}^2 : r_1 \leq \|x\| \leq r^*\}$

Visualize C_1

+: contributes to error



From the visualization, we can prove that:

$$L_{D,f}(A(S)) \leq \epsilon$$

Proof:

$$D(C^*) - D(C(S)) \leq D(C_1)$$

$$\Leftrightarrow \underbrace{D(C^* \setminus C(S))}_{\text{positives outside of } C(S)} \leq \underbrace{D(C_1)}_{\text{all positives in strip } C_1}$$

$$\Leftrightarrow L_{D,f}(A(S)) \leq \epsilon \quad \leftarrow \text{previous discussion, subs } D(G)$$

Now, we would like to upper bound $D^m(\{S|_x : L_{D,f}(h_S) > \epsilon\})$.

With the discussion above, we can prove that:

$$\{S|_x : L_{D,f}(h_S) > \epsilon\} = \{S|_x : S|_x \cap C_1 = \emptyset\}$$

→ there are positives in the strip

Proof:

$$\text{Since } \{S|_x : L_{D,f}(h_S) \leq \epsilon\} = \{S|_x : S|_x \cap C_1 \neq \emptyset\}$$

→ there are no positives in the strip

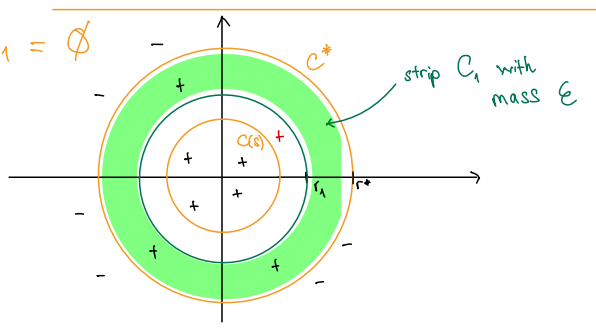
So the opposite holds:

$$\{S|_x : L_{D,f}(h_S) > \epsilon\} = \{S|_x : S|_x \cap C_1 = \emptyset\}$$

→ bad event

Visualize $S|_x \cap C_1 = \emptyset$

+ : contributes to error



Therefore, we can conclude that:

$$\begin{aligned} D^m(\{S|_x : L_{D,f}(h_S) > \epsilon\}) &= D^m(\{S|_x : S|_x \cap C_1 = \emptyset\}) \\ &\leq (1 - \epsilon)^m \\ &\leq e^{-\epsilon m} \end{aligned}$$

Let $\delta \in (0,1)$ such that $e^{-\epsilon m} \leq \delta$, then $m \geq \frac{\log(1/\delta)}{\epsilon}$, and so

\mathcal{H} is PAC learnable with $m_{\mathcal{H}} \leq \left\lceil \frac{\log(1/\delta)}{\epsilon} \right\rceil$