

Seminar 1

① Consider training set $S = \{(x_i, y_i)\}_{i=1}^m \subseteq (\mathbb{R}^d \times \{0, 1\})^m$

We consider the classifier from lecture 2, $h_S: \mathbb{R}^d \rightarrow \{0, 1\}$

$$h_S(x) = \begin{cases} y_i, & \text{if } x \in \{x_1, \dots, x_m\} \text{ s.t. } x_i = x \\ 0, & \text{otherwise} \end{cases}$$

We want to show that the classifier h_S can be written as a thresholded polynomial $p_S(x)$, meaning that we want to find a polynomial p_S such that

$$h_S(x) = 1 \Leftrightarrow p_S(x) \geq 0.$$

Let's consider the ~~more complex~~ simpler case, $d=1$ (so x_i is a scalar)

1st try: Consider the polynomial $p_S(x) = -\prod_{i=1}^m (x - x_i)$

If $x = x_i$ for some $i \in \{1, \dots, m\} \Rightarrow p_S(x) = p_S(x_0) = 0$. It will not work if the label of the label $y_0 = 0$. (In this case

$$p_S(x_0) = 0 \Rightarrow h_S(x_0) = 1$$

Also if x doesn't appear in the training data you don't know

if $p_S(x) \geq 0$ or $p_S(x) < 0$

2nd try: Consider the polynomial $p_S(x) = -\prod_{i=1}^m (x - x_i)^2$

If $x = x_i$ for some $i \in \{1, \dots, m\} \Rightarrow p_S(x) = p_S(x_0) = 0 \Rightarrow h_S(x) = 1$.

For points $(x_i, 0) \in S$ it will not work

For all other points it will work fine

3rd try: Consider the polynomial: $p_S(x) = -\prod_{\substack{i=1 \\ y_i=1}}^m (x - x_i)^2$

(in this case if all $y_0 = 0$ then $p_S(x) = -1$)

If $x = x_0$ for some $i \in \{1, \dots, m\} \Rightarrow$ $\begin{cases} y_0 = 1 \Rightarrow p_S(x) = 0 \Rightarrow h_S(x) = 1 \\ y_0 = 0 \Rightarrow p_S(x) < 0 \Rightarrow h_S(x) = 0 \end{cases}$

If $x \neq x_0$ for all $i \in \{1, \dots, m\} \Rightarrow p_S(x) < 0 \Rightarrow h_S(x) = 0$

Otherwise would be $p_S(x) = -\prod_{i=1}^m (x - x_i)^2 y_i$

$$p_S(x) = -\prod_{i=1}^m [(x - x_i)^2 + 1 - y_i]$$

Consider now the general case, d can be > 1

for $d=2$ we have seen that

$$P_S(x) = -\prod_{l=1}^m (x-x_0)^2 \text{ works fine}$$
$$x_0 = L$$

In the general case we consider the L_2 distance (Euclidean distance)

$$P_S(x) = -\prod_{l=1}^m \|x-x_0\|_2^2$$
$$x_0 = L$$

This polynomial will work fine.

(2) $\mathcal{H}_{rec}^2 \rightarrow h_{(a_1, b_1, a_2, b_2)} : \mathbb{R}^2 \rightarrow \{0, 1\}, a_1 \leq b_1 \text{ and } a_2 \leq b_2,$

$$h_{(a_1, b_1, a_2, b_2)}(x_1, x_2) = \begin{cases} 1, & a_1 \leq x_1 \leq b_1 \text{ and } a_2 \leq x_2 \leq b_2 \\ 0, & \text{otherwise} \end{cases}$$

\mathcal{H}_{rec}^2 is an infinite size hypothesis class, it is called the class of all axis aligned rectangles in the plane.

We want to prove that \mathcal{H}_{rec}^2 is PAC-learnable.

From the definition of PAC-learnability we know that

$\mathcal{H} = \mathcal{H}_{rec}^2$ is PAC-learnable if there exists a function $m_H : (0, 1)^2 \rightarrow \mathbb{N}$ and exists a learning algorithm A with the following property: for every $\epsilon, \delta > 0$, for every labeling function $f \in \mathcal{H}_{rec}^2$ (realizability case), for every distribution D on \mathbb{R}^2 when we run the learning algorithm A on a training set S consisting of $m \geq m_H(\epsilon, \delta)$ examples sampled i.i.d from D and labeled by f the algorithm A returns a hypothesis $h_S \in \mathcal{H}$ such that with probability at least $1-\delta$ (over the choice of examples) the risk of h_S is smaller than ϵ :

$$\Pr_{S \sim D^m} (L_{f, D}(h_S) \leq \epsilon) \geq 1-\delta \text{ or otherwise said}$$

$$\Pr_{S \sim D^m} (L_{f, D}(h_S) > \epsilon) \leq \delta.$$

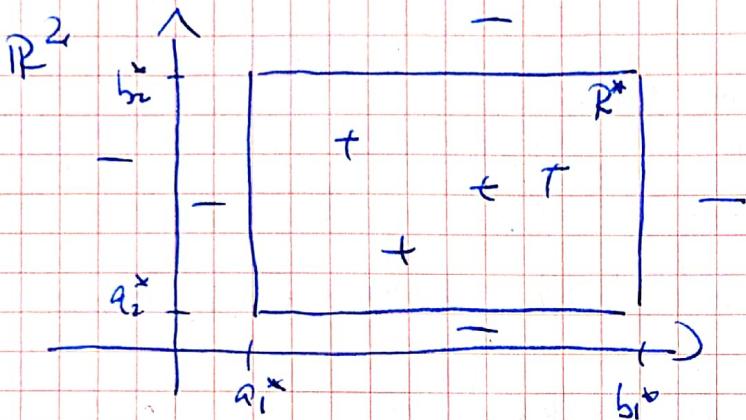
First we need to find the algorithm A.

We are under the realizability assumption, so there exist a labeling function $g \in \mathcal{H}$, $g = h^{*}_{(a_1^*, b_1^*, a_2^*, b_2^*)}$ that labels the training ob-

Consider the training set $S \subset \{(x_i, y_i), (x_2, y_2), \dots, (x_m, y_m)\}$

$$y_i = h^{*}_{(a_1^*, b_1^*, a_2^*, b_2^*)}(x_i), x_i \in \mathbb{R}^2$$

$$x_i = (x_{i1}, x_{i2}) \quad \{$$



h^{*} labels each point drawn from the rectangle $R^* = [a_1^*, b_1^*] \times [a_2^*, b_2^*]$ with label 1, and all other points with label 0.

$$\text{So } h^{*}_{(a_1^*, b_1^*, a_2^*, b_2^*)} = \mathbb{1}_{R^*}$$

Consider the following algorithm A, that takes as input sample the training set S and output h_S .

$$h_S = h(a_{1S}, b_{1S}, a_{2S}, b_{2S}) \text{ where}$$

$$a_{1S} = \min_{\substack{i=1, m \\ y_i=1}} x_{i1}$$

$$a_{2S} = \max_{\substack{i=1, m \\ y_i=1}} x_{i1}$$

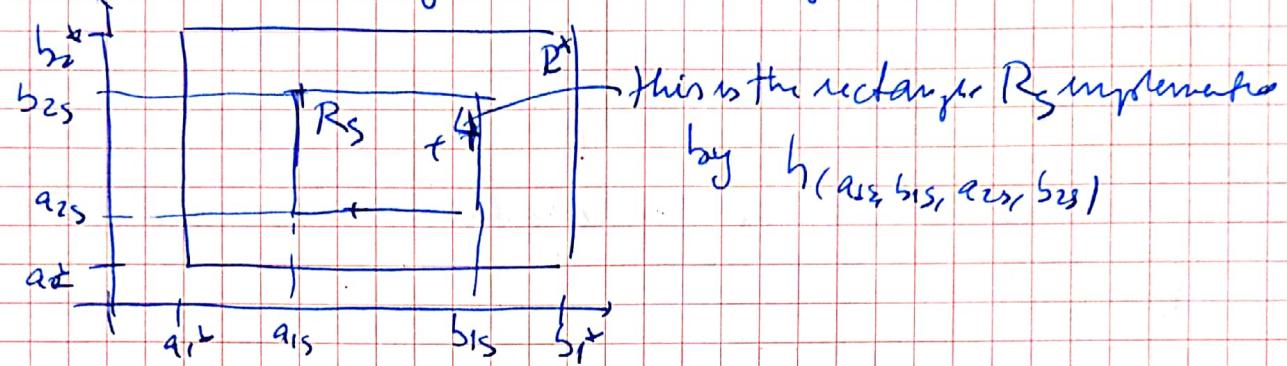
$$b_{1S} = \min_{\substack{i=1, m \\ y_i=1}} x_{i2}$$

$$b_{2S} = \max_{\substack{i=1, m \\ y_i=1}} x_{i2}$$

If all $y_i=0$ then all points x_i have label 0, that is no positive example. In this case choose $z = (z_1, z_2)$ a point that is not in the training set S and take $a_{1S} = b_{1S} = z_1, a_{2S} = b_{2S} = z_2$.

As in the induction $h_S = h_{(a_{1S}, b_{1S}, a_{2S}, b_{2S})} \geq 1$ if $R_S = [a_{1S}, b_{1S}] \times [a_{2S}, b_{2S}]$ is the indicator function of the tightest rectangle enclosing all positive examples.

By construction A is an ERT, meaning that $L_{h^*, D}(h_S) = 0$, h_S doesn't make any errors on the training set S .



Now we want to find the sample complexity $m_{\epsilon, \delta}$ such that

$$\Pr_{S \sim D^m} (L_{h^*, D}(h_S) \leq \epsilon) \geq 1 - \delta \quad \text{when } S \text{ contain } m \geq m_{\epsilon, \delta} \text{ examples.}$$

We make the observation that h_S can make errors in region $R^* - R_S$, labeling points that should get label 1 with label 0. All points $\in R$ will be labeled correctly, all points outside R^* will be labeled correctly.

Let's fix $\epsilon, \delta > 0$ and consider a distribution D over \mathbb{R}^2 .

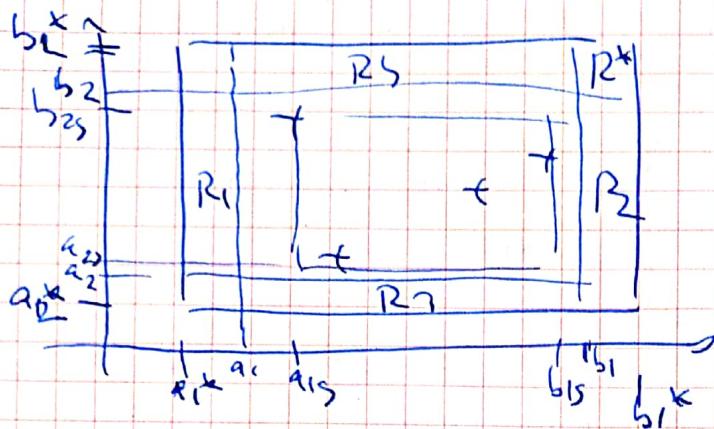
Case 1) if $D(R^*) = \Pr_{x \sim D} (x \in R^*) \leq \epsilon$ then in this case

$$L_{h^*, D}(h_S) \geq \Pr_{x \sim D} (h_S(x) \neq h^*(x)) = \Pr_{x \sim D} (x \in R^* - R_S) \leq$$

$$\leq \Pr_{x \sim D} (x \in R^*) \leq \epsilon \text{ so we have that}$$

$$\Pr_{S \sim D^m} (L_{h^*, D}(h_S) \leq \epsilon) = 1 \quad (\text{this happens all the time})$$

Case 2 $D(R^*) \geq P(x \in R^*) > \varepsilon$



We construct as in induction the rectangles R_1, R_2, R_3, R_4

$$R_1 = [a_1^*, a_1] \times [a_2^*, b_2^*]$$

$$R_2 = [b_3, b_1^*] \times [a_2^*, b_2^*]$$

$$R_3 = [a_1^*, b_1^*] \times [a_2^*, a_2]$$

$$R_4 = [a_1^*, b_1^*] \times [b_2, b_2^*]$$

$$\text{with } D(R_1) \geq P(x \in R_0) = \frac{\varepsilon}{4}.$$

if $R = [a_1, b_1] \times [a_2, b_2]$ (the rectangle returned by A, implemented by h_S) we find each $R_i, i=1 \dots 4$ from

$$L_{h^*, D}(h_S) = P_{x \in D}(h^*(x) \neq h_S(x)) \geq P_{x \in D}(x \in R^* \setminus R_S) \leq$$

$$\leq P_{x \in D}(x \in R_1 \cup R_2 \cup R_3 \cup R_4) \leq \sum_{i=1}^4 P_{x \in D}(x \in R_i)$$

$$\geq \sum_{i=1}^4 D(R_i) \geq \frac{1}{4} \cdot \frac{\varepsilon}{4} = \frac{\varepsilon}{16}$$

So, in this case $P_{S \sim D^m}(L_{h^*, D}(h_S) \leq \varepsilon) \geq 1$ (this happens always)

In order to have $L_{h^*, D}(h_S) > \varepsilon$ we need that R_S will not intersect at least one rectangle R_i .

We denote with ~~$\{S \sim D^m \mid R_S \cap R_i = \emptyset\}$~~

$$F_i = \{S \sim D^m \mid R_S \cap R_i = \emptyset\}$$

This leads to the following:

at least one F_i will happen

$$\begin{aligned} \Pr_{S \sim D^m} (L_{h \times D}(h_S) > \varepsilon) &= \Pr_{S \sim D^m} (F_1 \cup F_2 \cup F_3 \cup F_4) \leq \\ &\leq \sum_{i=1}^4 \Pr_{S \sim D^m} (F_i) \end{aligned}$$

$$\begin{aligned} \text{Now, } \Pr_{S \sim D^m} (F_i) &\approx \text{what is the probability that } R_S \text{ will not intersect } R_i \\ &\approx \text{the probability that no point from } R_i \text{ is sampled in } S \approx \\ &= (1 - \frac{\varepsilon}{4})^m \end{aligned}$$

$$\text{So } \Pr_{S \sim D^m} (L_{h \times D}(h_S) > \varepsilon) \leq \sum_{i=1}^4 \Pr_{S \sim D^m} (F_i) \approx 4 \cdot (1 - \frac{\varepsilon}{4})^m$$

Now, we know from lecture 2 that $1-x \leq e^{-x}$

$$\text{So } 1 - \frac{\varepsilon}{4} \leq e^{-\frac{\varepsilon}{4}}$$

$$\text{So } \Pr_{S \sim D^m} (L_{h \times D}(h_S) > \varepsilon) \leq 4 \cdot (1 - \frac{\varepsilon}{4})^m \leq 4 \cdot e^{-\frac{\varepsilon}{4}m}$$

this is the probability to have
 h_S not accurate $> \varepsilon$

We want to make this probability very small, smaller than δ .

$$4 \cdot e^{-\frac{\varepsilon}{4}m} < \delta$$

$$e^{-\frac{\varepsilon}{4}m} < \frac{\delta}{4} \quad | \cdot \log e$$

$$-\frac{\varepsilon}{4}m < \log \frac{\delta}{4} \quad | \cdot (-\frac{4}{\varepsilon})$$

$$m > -\frac{4}{\varepsilon} \log \frac{\delta}{4} = \frac{4}{\varepsilon} \log \frac{4}{\delta}$$

So if you take $m \geq m_\delta(\varepsilon, \delta) = \frac{4}{\varepsilon} \log \frac{4}{\delta}$ we obtain the desired results.

Repeat the previous question for class of aligned rectangles in \mathbb{R}^d .

in \mathbb{R}^d we have $\mathcal{H}_{\text{rec}}^d = \{ h_{(a_1, b_1, \dots, a_d, b_d)} : \mathbb{R}^d \rightarrow \{0,1\} \mid a_i \leq b_i \}$

$$h_{(a_1, b_1, \dots, a_d, b_d)} = \prod_{[a_i, b_i] \times [a_j, b_j] \cap [c_i, d_i] \neq \emptyset} 1$$

All the arguments used previously will hold, the general result

$$\text{will be that } m_H(\varepsilon, \delta) \geq \frac{2d}{\varepsilon} \cdot \log \frac{2d}{\delta}$$

For $d=2$ we obtain the power result.

The version of algorithm A is given by taking min, max over each dimension, so this mean $O(m \cdot d)$

$$m = \text{number of positive examples} = O\left(\frac{2d}{\varepsilon} \cdot \log \frac{2d}{\delta}\right)$$

$d = \text{number of dimensions}$

So we have the complexity of algorithm A is $O\left(\frac{2d^2}{\varepsilon} \cdot \log \frac{2d}{\delta}\right)$

which is polynomial in $d, \frac{1}{\varepsilon}, \frac{1}{\delta}$.

③ If \mathcal{H} is PAC learnable and $m_H : (0,1)^2 \rightarrow \mathbb{N}$ is its sample complexity

a) Given $\delta \in (0,1)$ and given $0 \leq \varepsilon_1 \leq \varepsilon_2 < 1$ we have that

$$m_H(\varepsilon_1, \delta) \geq m_H(\varepsilon_2, \delta)$$

\mathcal{H} is PAC learnable with sample complexity $m_H(\cdot, \cdot)$ means that there exists a learning algorithm A with the property: for every $\varepsilon, \delta > 0$ when given the algorithm A on a sample set S of m examples, $m \geq m_H(\varepsilon, \delta)$ (samples are labelled by $y \in \mathcal{Y}$ and i.i.d from a distribution D) we have that $h_S = A(S)$ with the risk

$$\Pr_{S \sim D^m} (L_{f, D}(h_S) \leq \varepsilon) \geq 1 - \delta$$

We apply this for ε_2 and $\delta \Rightarrow \Pr_{S \sim D^m} (L_{f, D}(h_S) \leq \varepsilon_1) \geq 1 - \delta$ if

$$m \geq m_H(\varepsilon_1, \delta)$$

We know that $\varepsilon_2 \geq \varepsilon_1$, so we have that

$$\Pr_{S \sim D} (L_{f,D}(h_S) \leq \varepsilon_2) \geq 1 - \delta \text{ if } m \geq m_2(\varepsilon_2, \delta)$$

But $m_2(\varepsilon_2, \delta)$ is the smallest number of examples for which the above inequality holds. So, if it holds for $m \geq m_2(\varepsilon_2, \delta)$ we have that $m_2(\varepsilon_2, \delta) \geq m_2(\varepsilon_1, \delta)$.

b) given $\varepsilon \in (0, 1)$, $0 < \delta_1 \leq \delta_2 < 1$ we have that $m_2(\varepsilon, \delta_1) \geq m_2(\varepsilon, \delta_2)$. Using the same arguments from a) we have that

$$\Pr_{S \sim D} (L_{f,D}(h_S) \leq \varepsilon) \geq 1 - \delta_1 \text{ if } m \geq m_2(\varepsilon, \delta_1)$$

$$\delta_1 \leq \delta_2 \Rightarrow 1 - \delta_1 \geq 1 - \delta_2 \Rightarrow$$

$$\Pr_{S \sim D} (L_{f,D}(h_S) \leq \varepsilon) \geq 1 - \delta_2 \text{ if } m \geq m_2(\varepsilon, \delta_2)$$

But $m_2(\varepsilon, \delta_2)$ is the smallest number of examples for which the above inequality holds (if $m \geq m_2(\varepsilon, \delta_2)$). So, if it holds for $m \geq m_2(\varepsilon, \delta_1)$ we have that $m_2(\varepsilon, \delta_1) \geq m_2(\varepsilon, \delta_2)$.

④ X discrete domain, $H_{\text{singleton}} = \{h_x : x \in X\} \cup \{h^-\}$

$$\forall x \in X, h_x : x \mapsto 1, h_x(x) = \begin{cases} 1, & x = x \\ 0, & x \neq x \end{cases}$$

$$h^- : x \mapsto 0, h^-(x) = 0 \forall x \in X$$

4.1) Describe an algorithm that implements the ERM rule for learning $H_{\text{singleton}}$ in the realizable setup.

Consider $S = \{(x_i, h^*(x_i)), x_0\}$ and for a function over $X\}_{i=1}^m$

The algorithm A is the following:

loop over training examples, if there is an $i \in \{1, \dots, m\}$ such that $y_0 = 1$ then return hypothesis $h_S = A(S) = h_{x_i}$

Otherwise return h^- .

From construction A is ERM meaning that $L_S(h_S) = 0$.

4.2) Show that $\mathcal{H}_{\text{singleton}}$ is PAC-learnable. Provide an upper bound on the sample complexity.

Let $\epsilon, \delta > 0$ and fix a distribution D over \mathcal{X} .

The only case which algorithm fails is the case where $h^* = h_2$ and the sample $S_2 = \{(x_i, y_i)\}_{i=1}^m$, (as sampled i.i.d from $D\}$ doesn't contain any positive example, so all $y_i = 0 \neq h_2(x_i)$.

In this case $h_S = A(S) = h_1$ which is different than h^* . However, we know that if $D(\{z\}) \leq \epsilon$ then everything is ok as we have that:

$$\Pr_{S \sim D} (L_{h^*, D}(h_S) \leq \epsilon) = 1$$

So, we have to upper bound the sample complexity in the case where $D(\{z\}) > \epsilon$ and there is no positive example in the set S (actually, for this problem there is just one positive point training point = 2).

So we have that

$$\Pr_{S \sim D} (L_{h^*, D}(h_S) > \epsilon) = \text{probability that each point in } S$$

is different than z (which has probability $m \epsilon > \epsilon$) $\leq (1 - \epsilon)^m \leq e^{-\epsilon m}$

So if we set $e^{-\epsilon m} < \delta \Leftrightarrow -\epsilon m < \log \delta$

$$m > -\frac{1}{\epsilon} \log \delta \Rightarrow m > \frac{1}{\epsilon} \log \frac{1}{\delta}$$

So if $m \geq \lceil \frac{1}{\epsilon} \log \frac{1}{\delta} \rceil$ we have that $\Pr_{S \sim D} (L_{h^*, D}(h_S) > \epsilon) < \delta$

So the upper bound is $m \chi(\epsilon, \delta)$ is $m \chi(\epsilon, \delta) \leq \lceil \frac{1}{\epsilon} \log \frac{1}{\delta} \rceil$