

1) An axis aligned rectangle classifier in the plane is a classifier that assigns the value 1 to a point iff it is inside a certain rectangle. Formally, given real numbers  $a_1 \leq b_1$ ,  $a_2 \leq b_2$ , define the classifier  $h(a_1, b_1, a_2, b_2)$

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

The class of all axis aligned rectangles in the plane is defined as  $H_{\text{rec}}^2 = \{ h(a_1, b_1, a_2, b_2) : a_1 \leq b_1 \text{ and } a_2 \leq b_2 \}$

We Assume realizability assumption.

1. Let  $A$  be the algorithm that returns the smallest rectangle enclosing all positive examples in the training set. Show that  $A$  is an ECM.

$A$  is an ERM if

$$A = \underset{h(a_1, b_1, a_2, b_2) \in H_{\text{rec}}^2}{\operatorname{argmin}} L_S(h(a_1, b_1, a_2, b_2))$$

Since  $A$  contains all positive examples, it will correctly classify all positive examples due to the definition of the classifier.

Due to the realizability assumption, we have a classifier / Algorithm  $A^*$  s.t.  $L_S(A^*) = 0$

Thus, for a classifier  $A^*$  to have  $L_S(A^*) = 0$  it must contain all positive examples and no negative examples.

The assumption assumes existence of such an  $A^*$ .

Consider the case where  $a_1, b_1, a_2$ , or  $b_2$  of  $A$  realize a new boundary  $a_1', b_1', a_2'$ , or  $b_2'$  s.t.

$$a \leq a_1', b_1' \leq b, a_2 \leq a_2', b_2' \leq b_2$$

in other words, the area of  $A$  is smaller.

Since  $A$  is the tightest/smallest rectangle that contains all positive examples it must have a positive example on the boundary of the rectangle. If this wasn't the case, there would exist a smaller  $A$ .

Therefore, decreasing the area of  $A$  by translating a boundary would result in a new classifier  $A'$  that would misclassify the example that was on the boundary before translation.

We also know  $A \subseteq A^*$  (shown later), so we do not need to consider negative examples now being correctly classified.

$$\Rightarrow L_S(A) < L_S(A')$$

The other case is when the area of  $A$  gets larger by translating a boundary in the opposite direction. Since  $A$  contains all positive examples increasing the area of  $A$  will only allow the possibility that the new classifier  $A''$  will incorrectly classify a negative example,

$$\Rightarrow L_S(A) \leq L_S(A'')$$

In conclusion,  $L_S(A) \leq L_S(A')$

$$L_S(A) \leq L_S(A'')$$

and the set of all rectangle classifiers

$$H_{\text{rec}}^2 = \{ A, \{A'\}, \{A''\} \}$$

← area more than  $A$   
 The set of all classifiers w/ area less than  $A$ .

thus  $A = \arg\min_{h(a_1, b_1, a_2, b_2) \in H_{\text{rec}}^2} L_S(h(a_1, b_1, a_2, b_2))$

⇒ and  $A$  is an ERM.

2. Show that if  $A$  receives a training set of size  $\geq \frac{4 \log(4/\delta)}{\epsilon}$  then, w/ probability of at least  $1 - \delta$  it returns a hypothesis with error at most  $\epsilon$ .

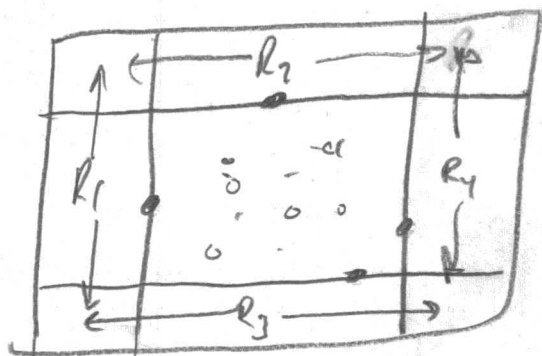
This will be done in steps following the outline of the problem.

if  $R_1, \dots, R_4$  have  $\epsilon/4$  error each.

and samples lie in

$$R_1 \text{ to } R_4 \quad L(h) \leq \epsilon$$

← this is worst case scenario with error  $= \epsilon$  because samples lie on boundaries of  $R_1, R_2, R_3$ , and  $R_4$ .



Show that  $R(S) \subseteq R^*$

$R(S)$  is rectangle returned by  $A$ .

$R^*$  is the rectangle that generates labels.

A realization of  $R^*, r^*$ , induces a realization  $r(S)$  of  $R(S)$ . Since  $R(S)$  is defined by the smallest rectangle by containing all positive instances, and  $R^*$  contains all positive instances as it is the true probability distribution,  $\Rightarrow R(S) \subseteq R^*$  (this proves  $A \in A^*$  used earlier)

Show that if  $S$  contains positive examples in all of the rectangles  $R_1, R_2, R_3, R_4$ , then the hypothesis returned by  $A$  has error at most  $\epsilon$ .

Due to realizability assumption, we have a perfect classifying rectangle,  $R(S)$ . If a positive example appears in  $R_1$  (as defined in the problem statement), then we would incorrectly classify at most  $\epsilon/4$  data points under  $D$ .

Similarly for rectangles  $R_2, R_3$ , and  $R_4$  contributing  $\epsilon/4$  to  $L_D(h)$ . Therefore, due to geometry of the rectangles, the contribution  $R$  to  $L_D(h)$  by  $A$  is

at most  $\epsilon$  (as  $R_1, R_2, R_3, R_4$  can overlap). and the error is at most the union of the area of these 4 rectangles.



For each  $i \in \{1, \dots, 4\}$ , upper bound the probability that  $S$  does not contain an example from  $R_i$

$\Pr \{ S \text{ does not contain an example from } R_i \}$

$$= \Pr \{ \forall (x_i, y_i) \in S, : x_i \notin R_i \} \leq \prod_{i=1}^{|S|} (1 - \frac{\epsilon}{4}) = (1 - \frac{\epsilon}{4})^{|S|}$$

$\epsilon/4$  probability  $x_i \in R_i$   $\Rightarrow$   $(1 - \epsilon/4)$  probability  $x_i \notin R_i$

Use the union bound to conclude the argument

$$\Pr \{ \forall (x_i, y_i) \in S : (x_i \notin R_1) \cup x_i \notin R_2 \cup x_i \notin R_3 \cup x_i \notin R_4 \}$$

$$\leq \Pr \{ \forall (x_i, y_i) \in S : x_i \notin R_1 \} + \Pr \{ \forall (x_i, y_i) \in S : x_i \notin R_2 \}$$

$$+ \Pr \{ \forall (x_i, y_i) \in S : x_i \notin R_3 \} + \Pr \{ \forall (x_i, y_i) \in S : x_i \notin R_4 \}$$

$$= 4 (1 - \epsilon/4)^{|S|} \leq 4 (e^{-\epsilon/4})^{|S|} = 4 e^{-\frac{\epsilon |S|}{4}} \leq 8$$

simplifying,

$$\Rightarrow \frac{e^{-\frac{\epsilon |S|}{4}}}{4} \leq 8/4$$

$$\ln \left( e^{-\frac{\epsilon |S|}{4}} \right) \leq \ln(8/4)$$

$$-\frac{\epsilon |S|}{4} \leq \ln(2)$$

conclusion,

$$\Rightarrow |S| \geq \left( \frac{4}{\epsilon} \right) \ln(2)$$

training set size must be greater than  $(4/\epsilon) \ln(2)$

Let  $X$  be a discrete domain, and let  $H_{\text{singleton}} = \{h_z: z \in X\}$

$H_{\text{singleton}} = \{h_z: z \in X\} \cup \{h^-\}$ , where for each  $z \in X$ ,  $h_z$  is the function defined by  $h_z(x) = 1$  if  $x = z$ ,  $h_z(x) = 0$  if  $x \neq z$

$h^-$  is simply the all-negative hypothesis, namely,  $\forall x \in X, h^-(x) = 0$ .

The realizability assumption here implies that the true hypothesis  $f$  labels negatively all examples in the domain, perhaps except one.

1. Describe an algorithm that implements the ERM rule for learning

$H_{\text{singleton}}$  in the realizable setup.

- Due to the realizability assumption, only one example

can be positive (1). Thus, the algorithm that follows is ERM.

Iterate over  $y_i = 1$ :

if  $y_i = 1$ :

output  $h_z$

terminate

else:

continue

output  $h^-$  (if no example has  $y_i = 1$ )

Now, since we know from the realizability assumption the true hypothesis  $f$ , there is two scenarios 1) all negative 2)  $(n-1)$  negative, 1 positive.

The algorithm handles both cases and gives  $L_S(h) = 0$

So it has to be ERM.

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

Recall PAC learnable if there exists a function  $m_H : (0,1)^2 \rightarrow \mathbb{N}$  and a learning algorithm w/ the following properties:

- For every  $\epsilon, \delta \in (0,1)$ , for every distribution  $D$  over  $X$ , and for every labelling function  $f: X \rightarrow \{0,1\}$ , if the realizability assumption holds wrt  $H, D, f$ , when running the learning algorithm on  $m \geq m_H(\epsilon, \delta)$  iid examples generated by  $D$  and labelled by  $f$ , the algorithm returns a hypothesis  $h$  s.t. w/ probability of at least  $1-\delta$  (over the choice of the examples),  $L_{(D,f)}(h) \leq \epsilon$

WTS  $\Pr \{ S \mid L_{(D,f)}(h_s) > \epsilon \} \leq \delta$

Previously I mentioned that there is two possibilities for  $f$

1) All negative

2) All negative except 1 sample

- however in a random sample this  $f$  could generate all negative.

Consider the first  $f$  (all negative). Then from part (a), then the algorithm designed in part (a) would return  $h^-$  which will always have  $L_{(D,f)}(h^-) = 0$ , thus for any  $\epsilon$ ,  $L_{(D,f)}(h_s) > \epsilon$  never occurs and  $\Pr \{ S \mid L_{(D,f)}(h_s) > \epsilon \} = 0$  and this is less than  $\delta$  for any  $\delta \in (0,1)$

If  $f$  is all negative (except for possibly one sample) the algorithm will output  $h_+$  or  $h_-$ . If it outputs  $h_+$ , then  $L_{(D,f)}(h_+) = 0$  and for similar reasoning,  $\Pr\{S \mid L_{(D,f)}(h_+) > \epsilon\}$  can't occur because  $L_{(D,f)}(h_+) = 0$  and  $\epsilon > 0$ , so  $\Pr\{S \mid L_{(D,f)}(h_+) > \epsilon\} \leq \delta$  for  $\delta > 0$ .

If the algorithm outputs  $h_-$ , well,  $L_{(D,f)}(h_-)$  is non-zero.

Let event  $E$  represent the event that there is a positive example.  $\bar{E}$  when there is no positive example.

Shown in book in 23.1 ( $\{S \mid x_i L_{(D,f)}(h_-) > \epsilon\} \subseteq M$ )

$$\Pr\{S \mid L_{(D,f)}(h_-) > \epsilon\} \leq \Pr\{S \mid L_S(h_-) = 0\} = \Pr\{\bar{E}\}$$

when  $\epsilon < L_{(D,f)}(h_-)$  by definition

$$\Pr\{h_-(x) \neq f(x)\} = \Pr\{E\}$$

$$\Pr\{\bar{E}\} = (1 - \Pr\{E\})^m$$

(due to iid assumption with  $m$  samples)

Since  $\epsilon < \Pr\{E\}$  (from (x))

$$\Pr\{S \mid L_{(D,f)}(h_-) > \epsilon\} \leq \Pr\{\bar{E}\} = \Pr\{1 - \Pr\{E\}\} \leq (1 - \epsilon)^m \leq \delta$$

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

$$m \log(1 - \epsilon) \leq \log \delta$$

$$m(-\log(1 - \epsilon)) \geq -\log \delta$$

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

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

$$\log(1 - \epsilon)$$

thus if we have enough

data,  $\Pr\{S \mid L_{(D,f)}(h_-) > \epsilon\} \leq \delta$  and we can conclude that this class,  $H_{\text{singleton}}$ , is PAC learnable.



Let  $X = \mathbb{R}^2$ ,  $Y = \{0, 1\}$ , and let  $\mathcal{H}$  be the class of concentric circles in the plane, that is,  $\mathcal{H} = \{h_r : r \in \mathbb{R}^+\}$  where  $h_r(x) = 1 \{ \|x\| \leq r \}$ . Prove that  $\mathcal{H}$  is PAC learnable, (assume realizability), and its sample complexity is bounded by

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \frac{\log(1/\delta)}{\epsilon}$$

With the realizability assumption for class  $\mathcal{H} = \{h_r : r \in \mathbb{R}^+\}$  we can have a learning algorithm that outputs a circle w/ radius  $r_h$  s.t.  $L_S(h) = 0$  (choose highest circle that fits  $S$  s.t.  $L_S(h) = 0$ ).  
 (similar to first rectangle problem) if we choose  $r_h$  to satisfy the above condition, then because  $L_{D,F}(h) \geq \epsilon \Rightarrow$  area in between  $f$  and  $h \geq \epsilon$   
 Similarly, the  $\Pr\{S : L_{D,F}(h) > \epsilon\} \leq \Pr\{S : L_S(h) = 0\}$

$$\leq \sum_{h: L_{D,F}(h) > \epsilon} \Pr\{S : L_S(h) = 0\} = \Pr\{h(x_j) = f(x_j), \forall j\}$$

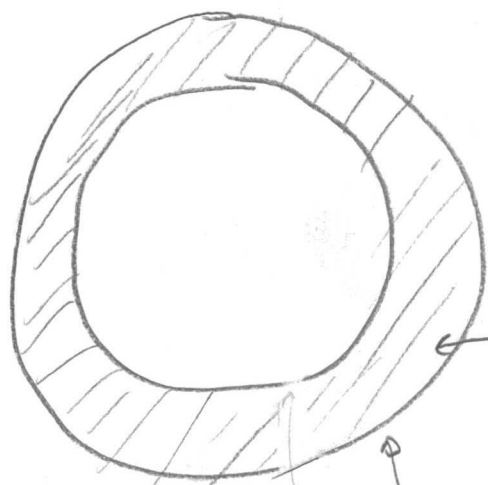
$$= \prod_{j=1}^m (1 - L_{D,F}(h)) \leq \prod_{j=1}^m (1 - \epsilon) \leq (1 - \epsilon)^m \leq e^{-m\epsilon} \leq \delta$$

$$\log(e^{-m\epsilon}) \leq \log \delta \Rightarrow -m\epsilon \leq \log \delta$$

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

and thus  $m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{\log(1/\delta)}{\epsilon} \right\rceil$ , (the ceiling preserves bound while giving natural number)

for a given  $\epsilon, \delta$  it is possible to achieve  $L_{D,F}(h) < \epsilon$   
 and if  $m \geq \frac{\log(1/\delta)}{\epsilon}$  (the ceiling preserves bound while giving natural number)



mass of ring  $= \epsilon$  under  $D$

true boundary  $f$  generating positive labels

$$\Pr \{ x_i \notin R \} = (1-\epsilon)$$

$$\Pr \{ (x_1, x_2, \dots, x_n) \notin R \} = (1-\epsilon)^m = \Pr \{ S \notin R \}$$

$$\Pr \{ S : L_{(D,f)}(h) > \epsilon \}$$

is the probability  $(x_1, x_2, \dots, x_n) \notin R$

$$\leq \Pr \{ S \notin R \} = (1-\epsilon)^m \leq e^{-m\epsilon} \leq \delta$$

If no sample point  $x_i$  is in  $R$ , our highest circle

(learned classifier) will have  $L_D(h) > \epsilon$  because

the mass of ring between learned classifier and boundary

$f > \epsilon$ .

$$e^{-m\epsilon} \leq \delta \Rightarrow \log(e^{-m\epsilon}) \leq \log(\delta) \Rightarrow -m\epsilon \leq \log(\delta)$$

$$\Rightarrow m\epsilon \geq \log(1/\delta) \Rightarrow m \geq \log(1/\delta)/\epsilon$$

$$m_H(\epsilon, \delta) \geq \left\lceil \frac{\log(1/\delta)}{\epsilon} \right\rceil$$

Let  $X$  be a domain and let  $D_1, D_2, \dots, D_m$  be a sequence of distributions over  $X$ . Let  $H$  be a finite class of binary classifiers over  $X$  and let  $f \in H$ . Suppose we are getting a sample  $S$  of  $m$  examples, s.t. the instances are independent but are not identically distributed. The  $i$ th instance is sampled from  $D_i$  and then  $y_i$  is set to be  $f(x_i)$ . Let  $\bar{D}_m$  denote the average, that is,

$$\bar{D}_m = \frac{D_1 + \dots + D_m}{m}$$

Fix an accuracy parameter  $\epsilon \in (0, 1)$ . Show that

$$\Pr \{ \exists h \in H \text{ s.t. } L(\bar{D}_m, f)(h) > \epsilon \text{ and } L_S(f)(h) = 0 \} \leq |H| e^{-\epsilon m}$$

$$L(\bar{D}_m, f)(h) = \frac{1}{m} \sum_i L_{(D_i, f)}(h) \quad (\text{from piazza and book})$$

$$\Pr \{ \{ \exists h \in H \text{ s.t. } L(\bar{D}_m, f)(h) > \epsilon \text{ and } L_S(f)(h) = 0 \} \}$$

$$\leq \Pr \{ \bigcup_{h \in H} \{ L(\bar{D}_m, f)(h) > \epsilon \text{ and } L_S(f)(h) = 0 \} \}$$

$H_B$  is the set of "bad hypothesis"  $h$  that satisfy  $L(\bar{D}_m, f)(h) > \epsilon$  and  $L_S(f)(h) = 0$ . This inequality is described in chapter 2 of the textbook. Note  $H_B \subseteq H$

$$\Pr \left\{ \bigcup_{h \in H_B} (L(\bar{D}, F) > \epsilon \text{ and } L_S(h) = 0) \right\}$$

$$\leq \sum_{i=1}^{|H_B|} \Pr \left\{ L(\bar{D}, F) > \epsilon \text{ and } L_S(h_i) = 0 \right\} \quad \text{by union bound.}$$

$$= \sum_{j=1}^{|H_B|} \Pr \{ L_S(h_j) = 0 \} \Pr \{ L(\bar{D}, F) > \epsilon \mid L_S(h_j) = 0 \} \quad \text{(card bound)}$$

by joint probability factorization. shown in book

$$= \sum_{j=1}^{|H_B|} \left( \prod_{i=1}^m (1 - L_{(D_i, F)}(h_j)) \right) \Pr \{ L(\bar{D}, F) > \epsilon \mid L_S(h_j) = 0 \}$$

because  $\Pr \{ L_S(h) = 0 \} = \prod_{i=1}^m (1 - L_{(D_i, F)}(h))$

and  $L_{(D_i, F)}(h)$  shown in book

By geometric-arithmetic mean inequality

$$\leq \sum_{i=1}^{|H_B|} \left( \frac{1}{m} \sum_{i=1}^m (1 - L_{(D_i, F)}(h)) \right)^m \Pr \{ L(\bar{D}, F) > \epsilon \mid L_S(h) = 0 \}$$

$$= \sum_{i=1}^{|H_B|} \left( 1 - L(\bar{D}, F)(h) \right)^m \Pr \{ L(\bar{D}, F) > \epsilon \mid L_S(h) = 0 \}$$

this term is at most 1.

$$\leq \sum_{i=1}^{|H_B|} (1 - \epsilon)^m \Pr \{ L(\bar{D}, F) > \epsilon \mid L_S(h) = 0 \}$$

$$\leq |H_B| (1 - \epsilon)^m \leq |H| (1 - \epsilon)^m \leq |H| e^{-\epsilon m}$$

Which  $H_B \subseteq H$  concludes the proof

Show that for every probability distribution  $D$ , the Bayes optimal predictor  $f_D$  is optimal, in the sense that for every classifier  $g$  from  $X$  to  $\{0, 1\}$   $L_D(f_D) \leq L_D(g)$

Bayes  $L_D(h) = \mathbb{E}_{(x,y) \sim D} [1(h(x) \neq y)]$  by definition

$= \mathbb{E}_x \left[ \mathbb{E}_{y|x}^{\text{minimize}} [1(h(x) \neq y)] \right]$  by law of iterative expectation

$= \mathbb{E}_x [P(Y=1|x) \cdot 1(h(x) \neq 1) + P(Y=0|x) \cdot 1(h(x) \neq 0)]$

for any prediction  $h(x) = 0$  contributes  $P(Y=1|x)$  to error

and  $h(x) = 1$  contributes  $P(Y=0|x)$  to error, thus to

Bayes  $\min_h \mathbb{E}_x [P(Y=1|x) \cdot 1(h(x) \neq 1) + P(Y=0|x) \cdot 1(h(x) \neq 0)]$

which  $L_D(f_D) = \mathbb{E}_x [\min(P(Y=1|x), P(Y=0|x))]$  rule to minimize above expression.

↑ minimum possible error.

any other classifier  $L_D(g)$  will have larger error.

We notice that the Bayes optimum classifier is the decision

rule  $L_D(f_D)$ , thus  $L_D(f_D) = \text{Bayes optimum} \leq L_D(g) = \text{any other}$

because  $L_D(f_D)$  is the minimum possible loss w.r.t  $D$