



Subject:

Date:

No:

## A Formal Learning Model - Exercises

1. Monotonicity of Sample Complexity: Let  $\mathcal{H}$  be a hypothesis class for a binary classification task. Suppose that  $\mathcal{H}$  is PAC learnable and its sample complexity is given by  $m_{\mathcal{H}}(\dots)$ .

Show that  $m_{\mathcal{H}}$  is monotonically nonincreasing in each of its parameters. That is, show that given  $\delta \in (0, 1)$ , and given  $0 < \epsilon_1 \leq \epsilon_2 < 1$ , we have that  $m_{\mathcal{H}}(\epsilon_1, \delta) \geq m_{\mathcal{H}}(\epsilon_2, \delta)$ . Similarly, show that given  $\epsilon \in (0, 1)$ , and given  $0 < \delta_1 \leq \delta_2 < 1$ , we have that  $m_{\mathcal{H}}(\epsilon, \delta_1) \geq m_{\mathcal{H}}(\epsilon, \delta_2)$ .

$D \rightsquigarrow$  unknown distribution over  $\mathcal{X}$

$f \in \mathcal{H} \rightsquigarrow$  target hypothesis

$A \rightsquigarrow$  learning algorithm with sample complexity of  $m_{\mathcal{H}}(\dots)$

Monotonicity of parameter  $\epsilon$ :

fix  $\delta$ , we have to show  $\epsilon_1 \leq \epsilon_2 \Rightarrow m_{\mathcal{H}}(\epsilon_1, \delta) \geq m_{\mathcal{H}}(\epsilon_2, \delta)$

Given an i.i.d. training sequence of size  $m \geq m_{\mathcal{H}}(\epsilon_1, \delta)$

we have  $L_{D,f}(h) \leq \epsilon_1 \leq \epsilon_2 \rightsquigarrow m_{\mathcal{H}}(\epsilon_1, \delta) \geq m_{\mathcal{H}}(\epsilon_2, \delta)$   
by minimality of  $m_{\mathcal{H}}(\epsilon_2, \delta)$

Monotonicity of parameter  $\delta$ :

fix  $\epsilon$ , we have to show  $\delta_1 \leq \delta_2 \Rightarrow m_{\mathcal{H}}(\epsilon, \delta_1) \geq m_{\mathcal{H}}(\epsilon, \delta_2)$

Given an i.i.d. training sequence of size  $m \geq m_{\mathcal{H}}(\epsilon, \delta_1)$

we have  $L_{D,f}(h) \leq$





Subject:

Date:

No:

2. Let  $X$  be a discrete domain, and let  $H_{\text{singleton}} = \{h_z : z \in X\} \cup \{h^-\}$ , where for each  $z \in X$ ,  $h_z$  is the function defined by  $h_z(x) = 1$  if  $x = z$  and  $h_z(x) = 0$  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.

Algorithm

- if a positive instance appears in  $S$  like  $x_+$  then return  $h_{x_+}$
- if  $S$  doesn't have any positive instances then return  $h^-$

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

$\epsilon \in (0, 1)$ ,  $D$  is a distribution over  $X$

i) if  $h^-$  is the true hypothesis  $\forall$

ii) suppose there exists a  $x_+$  instance.

if  $x_+$  appears in  $S \rightarrow \forall$

if  $D(\{x_+\}) \leq \epsilon \rightarrow$  the generalization error is at most  $\epsilon$

if  $x$  doesn't appear in  $S$  and  $D(\{x_+\}) > \epsilon$  then

$$D(\{x_+\} \mid x \text{ doesn't appear in } S) \leq (1 - \epsilon)^m \leq e^{-\epsilon m}$$

$\Rightarrow H_{\text{singleton}}$  is PAC learnable and

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



Date:

No:

3. Let  $X = \mathbb{R}^2$ ,  $Y = \{0, 1\}$ , and let  $H$  be the class of concentric circles in the plane, that is,  $H = \{h_r : r \in \mathbb{R}_+\}$ , where  $h_r(x) = 1[\|x\| \leq r]$ . Prove that  $H$  is PAC learnable (assume realizability) and its sample complexity is bounded by  $m_H(\epsilon, \delta) \leq \left\lceil \frac{\log 1/\delta}{\epsilon} \right\rceil$ . Suppose  $\hat{h}$  be a hypothesis corresponding to the tightest circle that contains all positive instances. Let  $\hat{h}$  has the radius  $\hat{r}$ . Let  $h^*$  be a circle with generalization risk of 0 and its radius is  $r^*$ .

Let  $r$  be a scalar such that  $\mathbb{P}(\{x : r \leq \|x\| \leq r + \epsilon\}) = \epsilon$ . The probability that  $L_D(h_S) \geq \epsilon$  is bounded by the probability that no points belong to  $\{x : r \leq \|x\| \leq r + \epsilon\}$  and it is bounded by  $(1 - \epsilon)^m$  which is less than  $e^{-\epsilon m}$ .  
 $\Rightarrow m_H(\epsilon, \delta) \leq \left\lceil \frac{\log 1/\delta}{\epsilon} \right\rceil$ .

4. In this question, we study the hypothesis class of Boolean conjunctions defined as follows. The instance space is  $X = \{0, 1\}^d$  and the label set is  $Y = \{0, 1\}$ . A literal over the variables  $x_1, \dots, x_d$  is a simple Boolean function that takes the form  $f(x) = x_i$ , for some  $i \in [d]$ , or  $f(x) = 1 - x_i$  for some  $i \in [d]$ . We use the notation  $\bar{x}_i$  as a shorthand for  $1 - x_i$ . A conjunction is any product of literals. In Boolean logic, the product is denoted using the  $\wedge$  sign. For example, the function  $h(x) = x_1(1 - x_2)$  is written as  $x_1 \wedge \bar{x}_2$ .

We consider the hypothesis class of all conjunctions of literals over the  $d$  variables. The empty conjunction is interpreted as the all-positive hypothesis (namely, the function that returns  $h(x) = 1$  for





Subject:

Date:

No:

all  $x$ ). The conjunction  $x_1 \wedge \bar{x}_2$  (and similarly any conjunction involving a literal and its negation) is allowed and interpreted as the all-negative hypothesis (namely, the conjunction that return  $h(x) = 0$  for all  $x$ ). We assume realizability: Namely, we assume that there exists a Boolean conjunction that generates the labels. Thus, each example  $(x, y) \in X \times Y$  consists of an assignment to the  $d$  Boolean variables  $x_1, \dots, x_d$ , and its truth value (0 for false and 1 for true).

For instance, let  $d=3$  and suppose that the true conjunction is  $x_1 \wedge \bar{x}_2$ . Then, the training set  $S$  might contain the following instances:

$((1, 1, 1), 0), ((1, 0, 1), 1), ((0, 1, 0), 0), ((1, 0, 0), 1)$ .

Prove that the hypothesis class of all conjunctions over  $d$  variable is PAC learnable and bound its sample complexity. Propose an algorithm that implements the ERM rule, whose runtime is polynomial in  $d, n$ .

First we compute the size of  $H$

$x_i$  —  $x_i$  appears in the corresponding conjunction  
 $\bar{x}_i$  —  $\bar{x}_i$  appears in the conjunction  
 not  $x_i$  nor  $\bar{x}_i$  appear in the  $n$

$$\Rightarrow |H| = 3^d + 1$$

$H$  is PAC learnable and  $\downarrow$  all-negative hypothesis

$$m_H(\epsilon, \delta) \leq \left\lceil \frac{d \log 3 + \log(1/\delta)}{\epsilon} \right\rceil$$





$h_0 = x_1 \wedge \bar{x}_2 \wedge \dots \wedge x_d \wedge \bar{x}_d \rightsquigarrow h_0$  is always minus hypothesis  
 $((a_1, y_1), \dots, (a_m, y_m)) \rightsquigarrow$  i.i.d. training sequence of size  $m$

For :

negative examples  $\rightsquigarrow$  our algorithm always neglect them

positive examples if  $a_i = 1$  we remove  $\bar{x}_i$  from  $h$

if  $a_i = 0$  we remove  $x_i$  from  $h$

$\rightsquigarrow h_i$  labels positively all the positive examples among  $a_1, a_2, \dots, a_i$   
 and

the set of literals in  $h_i$  contains the set of literals in the target hypothesis

$\rightsquigarrow \Rightarrow h_i$  classifies correctly the negative elements among  $a_1, a_2, \dots, a_m$

$\Rightarrow h$  is ERM

since we have  $m$  examples and for each of them it takes linear time  $d$  to process it the algorithm running time is  $O(m \cdot d)$





Subject:

Date:

No:

5. 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, such that 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,  $\bar{D}_m = (D_1 + D_2 + \dots + D_m) / m$ .

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

$$P[\exists h \in H \text{ s.t. } L_{\bar{D}_m, f}(h) > \epsilon \text{ and } L_{S, f}(h) = 0] \leq |H| e^{-\epsilon m}$$

$$\text{Fix } h \in H \text{ s.t. } L_{\bar{D}_m, f}(h) > \epsilon$$

We have

$$P_{x \sim \bar{D}_m} (h(x) = f(x)) = \frac{P_{x \sim D_1} (h(x) = f(x)) + \dots + P_{x \sim D_m} (h(x) = f(x))}{m}$$

$$< 1 - \epsilon$$

$$P_{\text{Sample Sequence} \sim D_1, \dots, D_m} (L_S(h) = 0) = \prod_{i=1}^m P_{x \sim D_i} (h(x) = f(x))$$

$$= \left( \prod_{i=1}^m P_{x \sim D_i} (h(x) = f(x)) \right)^{\frac{1}{m}})^m$$

$$\text{geometric-} \quad \leq \left( \frac{P_{x \sim D_1} (h(x) = f(x)) + \dots + P_{x \sim D_m} (h(x) = f(x))}{m} \right)^m$$

arithmetic

mean inequality

$$< (1 - \epsilon)^m \leq e^{-\epsilon m}$$

6. Let  $H$  be a hypothesis class of binary classifiers. Show that if  $H$  is agnostic PAC learnable, then  $H$  is PAC learnable as well. Furthermore, if  $A$  is a successful agnostic PAC learner





Date:

No:

For  $H$ , then  $A$  is also a successful PAC learner for  $H$ .

$H$  is agnostic PAC learnable

$A \rightsquigarrow$  a learning algorithm that learns  $H$

$D \rightsquigarrow$  unknown distribution over  $X$

$f \rightsquigarrow$  target function

Assume  $D$  is a joint distribution over  $X \times \{0,1\}$

the conditional probability of  $y$  given  $x$  is determined deterministically by  $f$

$\inf_{h \in H} L_D(h) = 0 \rightsquigarrow$  because of realizability

Let  $A$  consist of a training set  $S$  with  $m$  i.i.d. instances which are labeled by  $f \Rightarrow$  with probability at least  $1-\delta$ ,  $A$  returns an  $h$  with

$$L_D(h) \leq \inf_{h' \in H} L_D(h') + \epsilon \\ = 0 + \epsilon = \epsilon.$$

7. The Bayes Optimal predictor: 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  $\{0,1\}$ ,

$$L_D(f_D) \leq L_D(g).$$

$$\begin{aligned} P(f_D(x) \neq y | X=x) &= 1_{P(y=1|X=x) \geq \frac{1}{2}} \cdot P(y=0|X=x) + \\ &\quad 1_{P(y=1|X=x) < \frac{1}{2}} \cdot P(y=1|X=x) \\ &= 1_{P(y=1|X=x) \geq \frac{1}{2}} (1 - P(y=1|X=x)) + \\ &\quad 1_{P(y=1|X=x) < \frac{1}{2}} (P(y=1|X=x)) \\ &= \min \{ P(y=1|X=x), P(y=0|X=x) \} \end{aligned}$$





Subject:

Date:

No:

$$g: \mathcal{X} \rightarrow \{0, 1\}$$

$$\begin{aligned} P(g(x) \neq y | X=x) &= P(g(x)=0 | X=x) \cdot P(y=1 | X=x) \\ &\quad + P(g(x)=1 | X=x) \cdot P(y=0 | X=x) \\ &\geq P(g(x)=0 | X=x) \cdot \min\{P(y=1 | X=x), P(y=0 | X=x)\} \\ &\quad + P(g(x)=1 | X=x) \cdot \min\{P(y=1 | X=x), P(y=0 | X=x)\} \\ &= \min\{P(Y=1 | X=x), P(Y=0 | X=x)\} \end{aligned}$$

By the law of total expectation:

$$L_D(P_D) = E_{(x, y) \sim D} [1_{P_D(x) \neq y}]$$

$$= E_{x \sim D_x} [E_{y \sim D_{y|x}} (1_{P_D(x) \neq y} | X=x)]$$

$$= E_{x \sim D_x} (P(y=1 | X=x))$$

$$\leq E_{x \sim D_x} [E_{y \sim D_{y|x}} (1_{g(x) \neq y} | X=x)]$$

$$= L_D(g).$$





9. Consider a variant of the PAC model in which there are two example oracles: one that generates positive examples and one that generates negative examples, both according to the underlying distribution  $D$  on  $X$ . Formally, given a target function  $f: X \rightarrow \{0, 1\}$  let  $D^+$  be the distribution over  $X^+ = \{x \in X : f(x) = 1\}$  defined by  $D^+(A) = D(A) / D(X^+)$  for every  $A \subseteq X^+$ . Similarly,  $D^-$  is the distribution over  $X^-$  induced by  $D$ .

The definition of PAC learnability in the two-oracle model is the same as the standard definition of PAC learnability except that here the learner has access to  $m_H^+(\epsilon, \delta)$  i.i.d. examples from  $D^+$  and  $m^-(\epsilon, \delta)$  i.i.d. examples from  $D^-$ . The learner's goal is to output  $h$  s.t. with probability at least  $1 - \delta$  (over the choice of the two training sets, and possibly over the nondeterministic decisions made by the learning algorithm), both  $L_{(D^+, f)}(h) \leq \epsilon$  and  $L_{(D^-, f)}(h) \leq \epsilon$ .

1. (\*) Show that if  $H$  is PAC learnable (in the standard one-oracle model), then  $H$  is PAC learnable in the two-oracle model.

2. (\*\*) Define  $h^+$  to be the always-plus hypothesis and  $h^-$  to be always-minus hypothesis. Assume that  $h^+, h^- \in H$ . Show that if  $H$  is PAC learnable in the -oracle model, then  $H$  is PAC learnable in the standard one-oracle model.

$A \rightsquigarrow$  an algorithm that learns  $H$

$H \rightsquigarrow$  PAC learnable in one-oracle model

We want to prove that  $H$  is PAC learnable in two-oracle model

$D \rightsquigarrow$  distribution over  $X \times \{0, 1\}$

drawing points from the negative and positive oracles with





equal probability  $\equiv$  drawing i.i.d. examples from a distribution  $D'$  with equal probability for positive and negative examples

For every  $A \subseteq X$  we have

$$D'(A) = \frac{1}{2} D^+(A) + \frac{1}{2} D^-(A)$$

$$\Rightarrow D'(\{x: f(x)=1\}) = D'(\{x: f(x)=0\}) = 0.5$$

if we let the algorithm have access to a training set that is drawn i.i.d. from  $D'$  with size  $m_H(\epsilon/2, \delta)$  then with probability at least  $1-\delta$  the algorithm return  $h$  with  $L(D', f)(h) \leq \frac{\epsilon}{2}$

$$L(D', f)(h) = P_{x \sim D'}(f(x) \neq h(x))$$

$$= P_{x \sim D'}[f(x)=1, h(x)=0] + P_{x \sim D'}[f(x)=0, h(x)=1]$$

$$= P_{x \sim D'}(f(x)=1) \cdot P_{x \sim D'}(h(x)=0 | f(x)=1) +$$

$$P_{x \sim D'}(f(x)=0) \cdot P_{x \sim D'}(h(x)=1 | f(x)=0)$$

$$= P_{x \sim D'}(f(x)=1) \cdot P_{x \sim D}(h(x)=0 | f(x)=1)$$

$$+ P_{x \sim D'}(f(x)=0) \cdot P_{x \sim D}(h(x)=1 | f(x)=0)$$

$$= \frac{1}{2} L(D^+, f)(h) + \frac{1}{2} L(D^-, f)(h)$$

it implies that with probability at least  $1-\delta$

$$L(D^+, f)(h) \leq \epsilon$$

$$L(D^-, f)(h) \leq \epsilon$$