

# Mathematical Foundations of Machine Learning

Spring 2022

Tsinghua University

---

Lecturer: Yuan Zhou

**Homework 1**

Posted: Feb 25, 2022

Due: Mar 14 23:59, 2022

Name: 董浚哲

ID: 2019011985

---

*There are 6 problems in this homework. The total amount of available points is 100. Throughout this exercise, we use  $\log$  to denote the natural logarithm (i.e., the  $\ln$  function).*

1. [10 pts.] **VC-dimension of axis aligned rectangles in  $\mathbb{R}^d$ .** Let  $\mathcal{H}_{\text{rec}}^d$  be the class of axis aligned rectangles in  $\mathbb{R}^d$ . In Lecture 2, we have seen that  $\text{VCdim}(\mathcal{H}_{\text{rec}}^2) = 4$ . Prove that for general  $d \in \mathbb{N}_+$ , we have that  $\text{VCdim}(\mathcal{H}_{\text{rec}}^d) = 2d$ .

*Solution.*

On one hand,  $\mathcal{H}_{\text{rec}}^d$  shatters the set  $\{\pm e_1, \pm e_2, \dots, \pm e_n\}$ , where  $e_i = (0, \dots, 0, 1, 0, \dots, 0)$  with the only “1” at the  $i$ -th place. Denote the rectangle as  $[a_1, b_1] \times \dots \times [a_n, b_n]$ , where  $a_i = -\frac{1}{2} - h(-e_i)$ ,  $b_i = \frac{1}{2} + h(e_i)$

On the other hand, consider the set of points  $\mathbf{x}_1, \dots, \mathbf{x}_{2n+1}$ . Now we construct a scheme where at least one point labeled 0 is inevitably covered by the box.  $\forall i \in [1, n]$ , consider  $\mathbf{x}_{i_{\text{floor}}}, \mathbf{x}_{i_{\text{ceil}}}$  s.t.  $x_{i_{\text{floor}}}^i = \min\{x_k^i : 1 \leq k \leq 2n+1\}$ ,  $x_{i_{\text{ceil}}}^i = \max\{x_k^i : 1 \leq k \leq 2n+1\}$ . Label both points with “1”. If a point is not labeled after this process, label it with “0”. In this way, at most  $2n$  points are labeled “1” while at least 1 point is labeled with “0”. Meanwhile, to predict the selected  $2n$  points right, all points are covered by the rectangle since  $a_i \leq x_{i_{\text{floor}}}^i, b_i \geq x_{i_{\text{ceil}}}^i$ , which results in predicting the point labeled “0” wrong. So  $\mathcal{H}_{\text{rec}}^d$  cannot shatter any set of  $2n+1$  points.

□

2. [10 pts.] **VC-dimension of a vector space of real functions.** Let  $F$  be a finite-dimensional vector space of real functions on  $\mathbb{R}^n$ ,  $\dim(F) = r < \infty$ . Let  $\mathcal{H}$  be the set of hypotheses:

$$\mathcal{H} = \{\{x : f(x) \geq 0\} : f \in F\}.$$

Show that  $d$ , the VC-dimension of  $\mathcal{H}$ , is finite and that  $d \leq r$ . (**Hint:** select an arbitrary set of  $m = r + 1$  points and consider  $\{(f(x_1), \dots, f(x_m)) : f \in F\}$ )

*Solution.* Suppose the contrary is true s.t. for  $m = r + 1$ ,  $\forall C \subset \mathbb{R}^n, |C| = m$ ,  $\mathcal{H}$  shatters  $C = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$ . Select a basis of  $F$ :  $\mathbb{B} = \{f_1, f_2, \dots, f_r\}$ . Give points in  $C$  an arbitrary label. Since  $\mathcal{H}$  shatters  $C$ ,  $\exists f(x) = \sum_{i=1}^r \alpha_i f_i(x)$  and a partition  $I \cup J = \{1, 2, \dots, m\}, f(x_i) \geq 0 \ \forall i \in I, f(x_j) < 0 \ \forall j \in J$ . This results in the equation:

$$\begin{bmatrix} f_1(x_1) & \cdots & f_r(x_1) \\ \vdots & \ddots & \vdots \\ f_1(x_m) & \cdots & f_r(x_m) \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_r \end{bmatrix} = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_m \end{bmatrix}$$

where  $\forall i \in I, \beta_i \geq 0; \forall j \in J, \beta_j < 0$ .

This is an overdetermined equation, so unless rows of the matrix are linear related there's no solution (WLOG assume the last row can be linearly represented by others), but in this case  $\beta_m$  is no longer free to be positive or negative, which is a contradiction.  $\square$

3. [20 pts.] **The weighted majority learner.** Recall the “The Halving Learner” in Lecture 1. We reached the conclusion that the halving learner makes at most  $\log_2 |\mathcal{H}|$  mistakes under the assumption that there exists one hypothesis  $h \in \mathcal{H}$  which consistently provides the correct label. In this exercise, we will explore the case when such an assumption does not hold, and we will design a learner trying to match up with the best hypothesis in  $\mathcal{H}$  (i.e., the one that proves the least number of wrong labels).

Now, let us assume that there exists one hypothesis  $h \in \mathcal{H}$  that makes at most  $M$  mistakes. The goal of our learner is to upper bound the number of mistakes by a function of  $M$ .

- (a) **The simple method.** We run the halving learner until every hypothesis in  $\mathcal{H}$  makes at least one mistake (so that all hypotheses are eliminated). We call such a process an *epoch*. In our simple method, we repeatedly run the epochs – whenever an epoch ends, we start a new one with all the (eliminated) hypotheses “alive”. Please upper bound the number of epochs and the number of mistakes the simple method would make.

- (b) **The weighted majority approach.** Let us assign non-negative weights  $w_h$  to each hypothesis  $h \in \mathcal{H}$ . Instead of elimination (which can be viewed as a hard penalty), we impose a soft penalty to a hypothesis function  $h$  whenever it makes a mistake – we reduce the weight  $w_h$  by a pre-determined factor  $C \in [0, 1)$ . Note that the elimination (hard penalty) is equivalent to set  $C = 0$  so that the weight of the hypothesis function  $h$  is set to 0.

In Algorithm 1 we describe the details of the weighted majority approach. The name “weighted majority” comes from the prediction rule that is based on the weighted majority of the hypotheses, where in contrast, the halving learner admits either 0 or 1 weight and predicts based on the majority of the weight-1 hypotheses.

To analyze the weighted majority learner, let us consider the potential function (at time  $t$ ) that is defined as  $w^{(t)} = \sum_{h \in \mathcal{H}} w_h^{(t)}$ . Please determine the constant  $C$  in Algorithm 1 and prove that  $w^{(t+1)} \leq \frac{3}{4}w^{(t)}$  if the learner makes a mistake at time  $t$ .

---

**Algorithm 1** The Weighted Majority Learner

---

```

1: Initialization: Set the initial weight for each  $h \in \mathcal{H}$  as  $w_h^{(1)} = 1$ 
2: for  $t = 1, 2, \dots$  do
3:   Get  $x_t$ 
4:   Predict  $\hat{y}_t = \text{WeightedMajority}\{h(x_t) : h \in \mathcal{H}\}$  i.e.  $\hat{y}_t = \text{sign}(\sum_{h \in \mathcal{H}} w_h^{(t)} h(x_t))$ 
5:   Get  $y_t$  and update weights as follows –
6:   if  $y_t \neq \hat{y}_t$  then
7:     Let  $w_h^{(t+1)} = C w_h^{(t)}$  for all  $h \in \mathcal{H}$ 
8:   else
9:     For all  $h \in \mathcal{H}$  such that  $h(x_t) = \hat{y}_t$ , let  $w_h^{t+1} = C w_h^t$ 
10:    For all  $h \in \mathcal{H}$  such that  $h(x_t) \neq \hat{y}_t$ , let  $w_h^{(t+1)} = w_h^{(t)}$ 
11:   end if
12: end for

```

---

- (c) Under the assumption that there exists one hypothesis in  $\mathcal{H}$  that makes at most  $M$  mistakes, prove a lower bound of the potential function at any time  $t$ .
- (d) Derive the mistake bound for Algorithm 1 by combining the above conclusions.

*Solution.*

- (a) Denote the wise hypothesis who makes at most  $M$  mistakes as  $\hat{h}$ . In the worst case, exactly one (different) mistake is made in the first  $M$  epoch, and in the  $M+1$ th epoch,  $\hat{h}$  never makes an mistake. In each epoch, at most  $\log_2 |\mathcal{H}|$  mistakes are made so in total no more than  $(M+1)\log_2 |\mathcal{H}| - 1$  mistakes are made.
- (b) In the worst case, every prediction reaches a tie:  $\hat{y}_t = 0$  and is regarded as a mistake. Here we propose  $C = \frac{1}{2}$ . Whenever a tie is made, hypotheses that amounts to exactly half the total weight (e.g.  $\sum_{h \text{ wrong}} w_h^{(t)} = \frac{1}{2}w^{(t)}$ ), and their weight are reduced to half of it, so

$$w^{(t+1)} \leq \frac{1}{2} \sum_{h \text{ wrong}} w_h^{(t)} + \sum_{h \text{ correct}} w_h^{(t)} = \frac{1}{2} \cdot \frac{1}{2}w^{(t)} + \frac{1}{2}w^{(t)} = \frac{3}{4}w^{(t)}$$

- (c) In the worst case, every hypothesis apart from  $\hat{h}$  makes an mistake in every test, and  $\hat{h}$  makes an mistake in the first  $M$  tests. So (observe that no more mistakes can be made after  $2^{-M} \geq (\frac{1}{2})^{-t}(|\mathcal{H}| - 1)$  in the worst case, which means  $t \leq M + \log(|\mathcal{H}| - 1)$ )

$$w^{(t)} \geq \begin{cases} (\frac{1}{2})^{-t}|\mathcal{H}| & t \leq M \\ 2^{-M} + (\frac{1}{2})^{-t}(|\mathcal{H}| - 1) & M < t \leq M + \log_2(|\mathcal{H}| - 1) \end{cases}$$

- (d) Above we've proved that  $t \leq M + \log_2(|\mathcal{H}| - 1)$ , which is an lower bound. Combine (b) and we see that the upper bound is given by  $2^{-M} \geq \frac{1}{2}(\frac{3}{4})^t|\mathcal{H}|$ , which results in  $t \leq (2 - \log_2 3)^{-1}(M - 1 + \log_2(|\mathcal{H}|))$ .

□

4. [10 pts.] **Learning concentric circles.** Let  $\mathcal{X} = \mathbb{R}^2$  and consider the set of concepts of the form  $c = \{(x, y) : x^2 + y^2 \leq r^2\}$  for some real number  $r$ . Show that this class can be  $(\epsilon, \delta)$ -PAC-learned from training data of size  $m \geq (1/\epsilon) \log(1/\delta)$ .

(**Hint:** The proof of Lemma 6.1 in the textbook *Understanding Machine Learning* might serve as a reference.)

*Solution.*

Denote the (concentric) circle with radius  $r$  to be  $C_r = \{(x, y) : x^2 + y^2 \leq r^2\}$ . Suppose  $r^*$  is the radius s.t.  $L_D(C_{r^*}) = 0$ . Let  $D_x$  be the marginal distribution over the domain  $\mathbb{X}$  and denote  $r_1 < r^* < r_2$  the radii s.t.

$$P_{\mathbf{x} \sim D_x}(\mathbf{x} \in C_{r^*} \setminus C_{r_1}) = P_{\mathbf{x} \sim D_x}(\mathbf{x} \in C_{r_2} \setminus C_{r^*}) = \varepsilon$$

Given a training data set, denote  $b_1 = \max\{|\mathbf{x}| : \mathbf{x} \in C_{r^*}\}$ ,  $b_2 = \min\{|\mathbf{x}| : \mathbf{x} \notin C_{r^*}\}$ , and denote  $C_{r_S}$  the corresponding ERM hypothesis. By definition  $b_1 < r_S < b_2$ , so

$$\begin{aligned} & P_{\mathbf{x} \sim D^m}(L_D(C_{r_S}) > \varepsilon) \\ & \leq P_{\mathbf{x} \sim D^m}(r_2 < b_2 \vee r_1 > b_1) \\ & \leq P_{\mathbf{x} \sim D^m}(r_2 < b_2) + P_{\mathbf{x} \sim D^m}(r_1 > b_1) \\ & = (1 - \varepsilon)^m + (1 - \varepsilon)^m \leq 2e^{-\varepsilon m} \end{aligned}$$

Take the estimate of  $m$  into the inequality above and we get the desired statement.  $\square$

5. [20 pts.] **The robust ellipsoid learner in adversarial setting.** In this problem we consider a special adversarial setting of the online learning game. Specifically, when we wish to predict the label with the parameter  $\mathbf{w}_t$  given the instance  $\mathbf{x}_t$ , the adversary will perturb the model parameter to  $\tilde{\mathbf{w}}_t$  (see Line 7 in Algorithm 2). Assuming  $\|\tilde{\mathbf{w}}_t - \mathbf{w}_t\| \leq \epsilon_t \leq \sqrt{\lambda_{\min}(A_t)}/(10d^3)$  where  $\lambda_{\min}(A_t)$  is the smallest eigenvalue of  $A_t$ , we aim to prove that, if we slightly adjust the update process of  $A_t$  (see the red constant in Line 10 in Algorithm 2), our ellipsoid learner will still work with a similar mistake bound.

Fix any time  $t$ , let us write  $A_t$  in the form of  $A_t = UD^2U^\top$  where  $U$  is orthonormal and  $D$  is diagonal with non-negative entries.

- (a) Show that if  $y_t \langle UD\mathbf{z} + \mathbf{w}_t, \mathbf{x}_t \rangle > 0$  and  $y_t \langle \tilde{\mathbf{w}}_t, \mathbf{x}_t \rangle < 0$ , we have  $y_t \langle UD\mathbf{z}, \mathbf{x}_t \rangle > -\epsilon \|\mathbf{x}_t\|$ .
- (b) Show that if  $y_t \langle UD\mathbf{z}, \mathbf{x}_t \rangle > -\epsilon \|\mathbf{x}_t\|$ , we have  $\langle \mathbf{z}, y_t DU^\top \mathbf{x}_t \rangle / \|DU^\top \mathbf{x}_t\| > -1/(10d^3)$ .
- (c) Prove the Enclosing Ellipsoid Lemma: when  $\hat{y}_t \neq y_t$ , we have that  $\mathcal{E}_t \cap \{\mathbf{w} : y_t \langle \mathbf{w}, \mathbf{x}_t \rangle > 0\} \subseteq \mathcal{E}_{t+1}$ .

---

**Algorithm 2** The Robust Ellipsoid Learner.

---

```
1: Initialization:  $\mathbf{w}_1 = \mathbf{0}, \mathbf{A}_1 = I$ 
2: for  $t = 1, 2, \dots$  do
3:   Get  $\mathbf{x}_t$ 
4:   if  $|\mathcal{E}_t \cap G^d| = 1$  then
5:     Choose the only element in  $\mathcal{E}_t \cap G^d$  and predict  $\hat{y}_t = \text{sign}(\mathbf{w}_t^\top \mathbf{x}_t)$ 
6:   else
7:     Predict  $\hat{y}_t = \text{sign}(\tilde{\mathbf{w}}_t^\top \mathbf{x}_t)$ , where  $\tilde{\mathbf{w}}_t$  is the perturbed parameter selected by the
       adversary such that  $\|\tilde{\mathbf{w}}_t - \mathbf{w}_t\| \leq \epsilon_t$ 
8:   end if
9:   if  $y_t = \hat{y}_t$  then
10:    Update:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \frac{y_t}{d+1} \frac{A_t \mathbf{x}_t}{\sqrt{\mathbf{x}_t^\top A_t \mathbf{x}_t}} = \mathbf{w}_t + \mathbf{u}_t$$

$$A_{t+1} = \frac{d^2}{d^2 - 1} \left( A_t - \frac{2}{d+1} \frac{A_t \mathbf{x}_t \mathbf{x}_t^\top A_t}{\mathbf{x}_t^\top A_t \mathbf{x}_t} \right) = \left( 1 + \frac{1}{2d^2} \right) \frac{d^2}{d^2 - 1} (A_t - 2(d+1) \mathbf{u}_t \mathbf{u}_t^\top)$$

11:   else
12:      $\mathbf{w}_{t+1} = \mathbf{w}_t$  and  $A_{t+1} = A_t$ 
13:   end if
14: end for
```

---

(d) Prove the Volume Reduction Lemma:  $\text{Vol}(\mathcal{E}_{t+1}) \leq \text{Vol}(\mathcal{E}_t) e^{-1/(4d+4)}$  and also prove the mistake bound for the robust learner under the presence of the adversary.

*Solution.*

(a)  $y_t \langle UD\mathbf{z} + \mathbf{w}_t, \mathbf{x}_t \rangle = y_t \langle UD\mathbf{z}, \mathbf{x}_t \rangle + y_t \langle \mathbf{w}_t - \tilde{\mathbf{w}}_t, \mathbf{x}_t \rangle + y_t \langle \tilde{\mathbf{w}}_t, \mathbf{x}_t \rangle$ , so

$$y_t \langle UD\mathbf{z}, \mathbf{x}_t \rangle = y_t \langle UD\mathbf{z} + \mathbf{w}_t, \mathbf{x}_t \rangle - y_t \langle \tilde{\mathbf{w}}_t, \mathbf{x}_t \rangle - y_t \langle \mathbf{w}_t - \tilde{\mathbf{w}}_t, \mathbf{x}_t \rangle > 0 + 0 - \|\tilde{\mathbf{w}}_t - \mathbf{w}_t\| \|\mathbf{x}_t\| \geq -\epsilon \|\mathbf{x}_t\|$$

(b) Observe that  $y_t \langle UD\mathbf{z}, \mathbf{x}_t \rangle = \langle \mathbf{z}, y_t DU^T \mathbf{x}_t \rangle$ , so  $LHS \geq -\frac{\sqrt{\lambda_{\min}(A_t)}}{(10d^3)} \frac{\|\mathbf{x}_t\|}{\|DU^T \mathbf{x}_t\|}$ . Take it into the inequality and simplify, and we see that it's equivalent to prove:

$$\sqrt{\lambda_{\min}(A_t)} \|\mathbf{x}_t\| = \sqrt{\lambda_{\min}(A_t)} \|U^T \mathbf{x}_t\| \leq \|DU^T \mathbf{x}_t\|$$

which is trivially by the definition of SVD. The equality holds when  $U^T \mathbf{x}_t$  is a right singular vector of the smallest singular value  $\sigma_n = \sqrt{\lambda_{\min}(A)}$

(c) When  $\hat{y}_t \neq y_t$ , then  $y_t = -\text{sign}(\tilde{\mathbf{w}}_t^T \mathbf{x}_t) \Rightarrow y_t \langle \tilde{\mathbf{w}}_t, \mathbf{x}_t \rangle = -|\langle \tilde{\mathbf{w}}_t, \mathbf{x}_t \rangle| < 0$

The goal is:

$$\forall \mathbf{w} \text{ s.t. } y_t \langle \mathbf{w}, \mathbf{x}_t \rangle > 0 \wedge (\mathbf{w} - \mathbf{w}_t)^T U D^{-2} U^T (\mathbf{w} - \mathbf{w}_t) \leq 1 \Rightarrow \mathbf{w} \in \mathcal{E}_{t+1}$$

Denote  $\mathbf{z} = D^{-1} U^T (\mathbf{w} - \mathbf{w}_t)$ , then the goal is equivalent to:

$$\forall \mathbf{z} \text{ s.t. } y_t \langle U D \mathbf{z} + \mathbf{w}_t, \mathbf{x}_t \rangle > 0 \wedge \mathbf{z}^T \mathbf{z} \leq 1 \Rightarrow U D \mathbf{z} \in \mathcal{E}(A_{t+1}^{\frac{1}{2}}, t)$$

In this case, by (a) we have  $y_t \langle U D \mathbf{z}, \mathbf{x}_t \rangle \geq -\epsilon \|\mathbf{x}_t\|$ . So it suffices to prove that:

$$\forall \mathbf{z} \text{ s.t. } \langle \mathbf{z}, y_t D U^T \mathbf{x}_t \rangle \geq -\epsilon \|\mathbf{x}_t\| \wedge \mathbf{z}^T \mathbf{z} \leq 1 \Rightarrow U D \mathbf{z} \in \mathcal{E}(A_t^{\frac{1}{2}}, u_t)$$

Denote  $V$  as the Householder matrix (which is symmetric and orthogonal, so we do not bother distinguishing  $V$  and  $V^T$ ) s.t.  $V y_t D U^T \mathbf{x}_t = \|D U^T \mathbf{x}_t\| e_1$ , and denote  $\tilde{\mathbf{z}} = V \mathbf{z}$ , then the new goal is equivalent to (denote the first dimension of  $\tilde{\mathbf{z}}$  as  $\tilde{z}_1$ ):

$$\forall \tilde{\mathbf{z}} \text{ s.t. } \|D U^T \mathbf{x}_t\| \tilde{z}_1 \geq -\epsilon \|\mathbf{x}_t\| \wedge \tilde{\mathbf{z}}^T \tilde{\mathbf{z}} \leq 1 \Rightarrow U D^T V \tilde{\mathbf{z}} \in \mathcal{E}(A_{t+1}^{\frac{1}{2}}, u_t)$$

by (b), it suffices to prove that:

$$\forall \tilde{\mathbf{z}} \text{ s.t. } \tilde{z}_1 \geq -\frac{1}{10d^3} \wedge \tilde{\mathbf{z}}^T \tilde{\mathbf{z}} \leq 1 \Rightarrow U D^T V \tilde{\mathbf{z}} \in \mathcal{E}(A_{t+1}^{\frac{1}{2}}, u_t)$$

**Claim:**  $\forall \tilde{\mathbf{z}} \text{ s.t. } \tilde{z}_1 > -\frac{1}{10d^3}, \tilde{\mathbf{z}}^T \tilde{\mathbf{z}} \leq 1 \Rightarrow \tilde{\mathbf{z}} \in \mathcal{E}(A^{\frac{1}{2}}, \tilde{u})$ , where

$$\tilde{u} = (\frac{1}{d+1}, 0, \dots, 0)^T \quad A = (1 + \frac{1}{2d^2}) \frac{d^2}{d^2 - 1} (\text{diag}(\frac{d-1}{d+1}, 1, \dots, 1))$$

*proof of the claim.* Denote  $a = \tilde{z}_1 \in [-\frac{1}{10d^3}]$ ,  $b = \sqrt{\sum_{i=2}^n \tilde{z}_i^2} \leq \sqrt{1 - a^2}$ .

It suffices to prove that  $(\tilde{\mathbf{z}} - \tilde{u})^T A^{-1} (\tilde{\mathbf{z}} - \tilde{u}) \leq 1$ :

$$\begin{aligned} & (\tilde{\mathbf{z}} - \tilde{u})^T A^{-1} (\tilde{\mathbf{z}} - \tilde{u}) \\ &= (1 + \frac{1}{2d^2})^{-1} \frac{d^2 - 1}{d^2} [\frac{d+1}{d-1} (a - \frac{1}{d+1})^2 + b^2] \\ &\leq (1 + \frac{1}{2d^2})^{-1} \frac{d^2 - 1}{d^2} [\frac{d+1}{d-1} (a - \frac{1}{d+1})^2 + (1 - a^2)] \\ &= (1 + a(a-1) \frac{2(d+1)}{d^2}) (1 + \frac{1}{2d^2})^{-1} \end{aligned}$$

So (after simplification) it suffices to prove that

$$a(a-1) \leq \frac{1}{4d^2(d+1)}$$

LHS takes its maximum at  $a = -\frac{1}{10d^3}$ , so it suffices to prove that

$$\frac{1}{10d^3}(1 + \frac{1}{10d^3}) \leq \frac{1}{4d^2(d+1)}$$

which is true for (calculated by Wolfram Alpha)  $d > 0.92809$ . Thus the claim is proved.  $\square$

$\forall \tilde{z} : \tilde{z}^\top \tilde{z} \leq 1$  and  $\tilde{z}_1 > 0$ , we have  $UDV^\top \tilde{z} \in \mathcal{E}(UDV^\top A^{1/2}, UDV^\top \tilde{\mathbf{u}})$  By defn.:  
 $V^\top = \begin{bmatrix} \frac{y_t DU^\top \mathbf{x}_t}{\|DU^\top \mathbf{x}_t\|} & \dots & \dots \end{bmatrix}$ . Therefore,  $UDV^\top \tilde{\mathbf{u}} = \frac{y_t UD^2 U^\top \mathbf{x}_t}{(d+1)\|DU^\top \mathbf{x}_t\|} = \frac{y_t}{d+1} \frac{A_t \mathbf{x}_t}{\sqrt{\mathbf{x}_t^\top A_t \mathbf{x}_t}} =$   
 $\mathbf{u}_t \in \mathcal{E}(UDV^\top A^{1/2}, UDV^\top \tilde{\mathbf{u}}) = \mathcal{E}((UDV^\top A V D U^\top)^{1/2}, \mathbf{u}_t) = \mathcal{E}(A_{t+1}^{1/2}, \mathbf{u}_t)$

$$\begin{aligned} V^\top A V &= (1 + \frac{1}{2d^2}) \frac{d^2}{d^2 - 1} V^\top \left( I - \text{diag} \left( \frac{2}{d+1}, 0, 0, \dots, 0 \right) \right) V \\ &= (1 + \frac{1}{2d^2}) \frac{d^2}{d^2 - 1} \left( I - \frac{2}{d+1} \frac{DU^\top \mathbf{x}_t \mathbf{x}_t^\top U D}{\|DU^\top \mathbf{x}_t\|^2} \right) \\ UDV^\top A V D U^\top &= (1 + \frac{1}{2d^2}) \frac{d^2}{d^2 - 1} \left( U D^2 U^\top - \frac{2}{d+1} \frac{U D^2 U^\top \mathbf{x}_t \mathbf{x}_t^\top U D^2 U^\top}{\|DU^\top \mathbf{x}_t\|^2} \right) \\ &= (1 + \frac{1}{2d^2}) \frac{d^2}{d^2 - 1} \left( A_t - \frac{2}{d+1} \frac{A_t \mathbf{x}_t \mathbf{x}_t^\top A_t}{\|DU^\top \mathbf{x}_t\|^2} \right) = A_{t+1} \end{aligned}$$

(d) Repeat what's written on lecture notes word by word, and we see that

$$\begin{aligned} & \left( \frac{\text{Vol}(\mathcal{E}_{t+1})}{\text{Vol}(\mathcal{E}_t)} \right)^2 \\ &= \frac{\det(A_{t+1})}{\det(A_t)} = \left( 1 + \frac{1}{d^2} \right)^d \left( \frac{d^2}{d^2 - 1} \right)^d \left( 1 - \frac{2}{d+1} \right) \\ &= \left( 1 + \frac{1}{2d^2} \right)^d \left( 1 + \frac{1}{d^2 - 1} \right)^{d-1} \left( \frac{d-1}{d+1} \cdot \frac{d^2}{d^2 - 1} \right) = \left( 1 + \frac{1}{2d^2} \right)^d \left( 1 + \frac{1}{d^2 - 1} \right)^{d-1} \left( 1 - \frac{1}{d+1} \right)^2 \\ &\leq \exp\left(\frac{d}{2d^2}\right) \exp\left(\frac{d-1}{d^2-1}\right) \cdot \exp\left(-\frac{2}{d+1}\right) \leq \exp\left(-\frac{1}{2(d+1)}\right) \end{aligned}$$

The lower bound of the volume remains unchanged, so the mistake bound is doubled:  $d(4d+4) \log(2n)$  mistakes.



□

6. [30 pts.] **Learning in the presence of noises.** We consider a finite hypothesis set  $\mathcal{H}$ , assume that the target concept is in  $\mathcal{H}$ , and adopt the following noise model: the label of a training point received by the learner is randomly changed with probability  $\eta \in (0, \frac{1}{2})$ . The exact value of the noise rate  $\eta$  is not known to the learner but an upper bound  $\eta'$  is revealed to him with  $\eta \leq \eta' < 1/2$ .

- (a) For any  $h \in \mathcal{H}$ , let  $d(h)$  denote the probability that the label of a training point received by the learner disagrees with the one given by  $h$ . Let  $h^*$  be the target hypothesis, show that  $d(h^*) = \eta$ .
- (b) More generally, show that for any  $h \in \mathcal{H}$ ,  $d(h) = \eta + (1 - 2\eta)L_{\mathcal{D}}(h)$ , where  $L_{\mathcal{D}}(h)$  denotes the generalization error of  $h$ .
- (c) Fix  $\epsilon > 0$  for this and all the following questions. Use the previous questions to show that if  $L_{\mathcal{D}}(h) > \epsilon$ , then  $d(h) - d(h^*) \geq \epsilon'$ , where  $\epsilon' = \epsilon(1 - 2\eta')$ .
- (d) For any hypothesis  $h \in \mathcal{H}$  and sample  $S$  of size  $m$ , let  $\hat{d}(h)$  denote the fraction of the points in  $S$  whose labels disagree with  $h$ . We will consider the algorithm  $L$  that, after receiving  $S$ , returns the hypothesis  $h_S$  with the smallest number of disagreements (thus  $\hat{d}(h_S)$  is minimized). To show the PAC-learning guarantee for  $L$ , we will show that for any  $h$ , if  $L_{\mathcal{D}}(h) > \epsilon$ , then with high probability  $\hat{d}(h) \geq \hat{d}(h^*)$ . First, show that for any  $\delta > 0$ , with probability at least  $1 - \delta/2$ , for  $m \geq \frac{2}{\epsilon'^2} \log \frac{2}{\delta}$ , the following holds:

$$\hat{d}(h^*) - d(h^*) \leq \epsilon'/2.$$

- (e) Second, show that for any  $\delta > 0$ , with probability at least  $1 - \delta/2$ , for  $m \geq \frac{2}{\epsilon'^2} (\log |\mathcal{H}| + \log \frac{2}{\delta})$ , the following holds for all  $h \in \mathcal{H}$ :

$$d(h) - \hat{d}(h) \leq \epsilon'/2.$$

- (f) Finally, show that for any  $\delta > 0$ , with probability at least  $1 - \delta$ , for  $m \geq \frac{2}{\epsilon'^2(1-2\eta')^2} (\log |\mathcal{H}| + \log \frac{2}{\delta})$ , the following holds for all  $h \in \mathcal{H}$  with  $L_{\mathcal{D}}(h) > \epsilon$ :

$$\hat{d}(h) - \hat{d}(h^*) \geq 0.$$

(**Hint:** use  $\widehat{d}(h) - \widehat{d}(h^*) = [\widehat{d}(h) - d(h)] + [d(h) - d(h^*)] + [d(h^*) - \widehat{d}(h^*)]$  and use previous questions to lower bound each of these three terms.)

*Solution.*

- (a) By definition,  $d(h^*)$  is the probability that the label received by the learner is different from that predicted by  $h^*$ . Since  $h^*$  always makes correct predictions, it  $d(h^*)$  is the probability that the received label is different from the original one. So by the definition of  $\eta$ ,  $d(h^*) = \eta$ .
- (b) By definition,  $L_D(h) = P_{x \sim D}\{h(x) \neq h^*(x)\}$ . Denote the label given to the learner as  $\tilde{y}$ . Then
- With probability  $1 - \eta$ ,  $h^*(x) = \tilde{y}$ . In this case,  $h(x)$  has to be different from  $h^*(x)$  to be wrong.
  - With probability  $\eta$ ,  $h^*(x) \neq \tilde{y}$ . In this case,  $h(x)$  has to be the same with  $h^*(x)$  to be wrong.

So  $d(h) = (1 - \eta)L_D(h) + \eta(1 - L_D(h)) = \eta + (1 - 2\eta)L_D(h)$ .

- (c) Since  $d(h^*) \equiv \eta$ , so  $d(h) - d(h^*) = (1 - 2\eta)L_D > (1 - 2\eta')L_D > (1 - 2\eta')\varepsilon = \varepsilon'$
- (d) Now that  $d(h^*) = \eta$ , it's equivalent to show that  $\hat{d}(h^*)$  holds for the given conditions. Since  $h^*$  always gives the correct prediction,  $\hat{d}(h^*)$  is the fraction of points whose labels have been changed. Given  $m$ , the number of the changed points  $X \sim B(m, \eta)$ . So

$$P(\hat{d}(h^*) < \eta + \frac{\varepsilon'}{2}) = P(X < m \cdot (\eta + \frac{\varepsilon'}{2}))$$

By central limit theorem,

$$P(\hat{d}(h^*) < \frac{\varepsilon'}{2} + \eta) \geq \Phi(\frac{\varepsilon'}{2} \sqrt{\frac{m}{\eta(1-\eta)}}) > \Phi(\varepsilon' \sqrt{m}) \geq \Phi(\sqrt{2 \log(\frac{2}{\delta})})$$

Meanwhile, as long as  $\delta < 2e^{-\frac{1}{2}} \approx 1.21 > 1$  (which is always satisfied since  $\delta < 1$ ), we have

$$1 - \Phi(\sqrt{2 \log(\frac{2}{\delta})}) < \int_{\sqrt{2 \log(\frac{\delta}{2})}}^{\infty} t e^{-t^2} dt = \frac{\delta}{2}$$

so  $P(\hat{d}(h^*) - d(h^*) < \frac{\varepsilon'}{2}) > 1 - \frac{\delta}{2}$

(e) With similar approach as above, we see that for a fixed  $h \in \mathcal{H}$ ,  $P((d(h) - \hat{d}(h)) > \frac{\varepsilon'}{2}) \leq \Phi(-\varepsilon' \sqrt{m}) > \frac{1}{2} e^{-\varepsilon'^2 m} \geq \frac{\delta}{H}$ . Now that there're  $|H|$  hypothesis, we add these errors and get:  $\forall h \in \mathcal{H}$

$$P(d(h) - \hat{d}(h) \leq \frac{\varepsilon'}{2}) \geq 1 - \frac{\delta}{2}$$

(f) Observe that

$$\hat{d}(h) - \hat{d}(h^*) = [\hat{d}(h) - d(h)] + [d(h) - d(h^*)] + [d(h^*) - \hat{d}(h^*)]$$

By previous arguments:

- $\hat{d}(h) - d(h) \geq -\varepsilon'/2$
- $d(h) - d(h^*) \geq \varepsilon'$
- $d(h^*) - \hat{d}(h^*) \geq -\varepsilon'/2$

Sum these terms and we reach the desired argument.

□