

① VC dim (Perception)  $\geq d+1$

$$f(x) = \begin{cases} 1 & w^T x + b \geq 0 \\ -1 & \text{else} \end{cases} \Rightarrow x^{(0)} = (0, \dots, 0)^T, x^{(i)} = (0, \dots, 1, 0, \dots, 0)^T, X = \{x^{(i)} \mid i \geq 0\}$$

$$\textcircled{2} y = (y_0, \dots, y_d)^T \in \{1, -1\}^{d+1}$$

$$\Rightarrow \text{[scribbled out]} \quad b = \frac{1}{r} y_0, \quad w = (w_1, \dots, w_d), \quad w_i = y_i \quad (1 \leq i \leq d) \Rightarrow$$

$$w^T x + b = \frac{1}{r} y_0 + \sum_{j=1}^d y_j x_j \quad x = x^{(i)} \Rightarrow w^T x^{(i)} + b = \begin{cases} \frac{1}{r} y_0 & i=0 \\ y_i & i \neq 0 \end{cases}$$

← ادا نقطہ ایش میں

(II)  $\text{vc dim}(\text{Perception}) < d+1 \Rightarrow w^T x + b = w^T x$

$\Rightarrow \bigcup_i x^{(i)}, \dots, x^{(d+r)} \in \mathbb{R}^d \Rightarrow X^T = (x^T, 1), W^T = (w^T, b), x^{(i)} \leftrightarrow x^{(i)}$   
 $X \in \mathbb{R}^{d+1}$

$$\Rightarrow x^{(1)}, \dots, x^{(d+1)} \in \mathbb{R}^d \Rightarrow \exists i: x^{(i)} = \sum_{j \neq i} a_j x^{(j)}, \quad a_j \neq 1$$

$$S = \{j \mid j \neq i, a_j \neq 0\} \Rightarrow \forall j \in S: \text{sign}(a_j) = y^{(j)}, y^{(i)} = -1$$

$$\Rightarrow h(x) = \begin{cases} 1 & f(x) > 0 \\ -1 & \text{else} \end{cases}, \quad \begin{array}{l} \text{طبق فرض خفن} \\ \text{وجود دارد که } \alpha \text{ را نقطه} \\ \text{را بدیش می دهد} \end{array} \quad \left| \quad \begin{array}{l} \forall j \in S: \alpha_j \cdot w^T x^{(j)} > 0 \\ w^T x^{(i)} \leq 0 \quad (*) \end{array} \right.$$

DATA OFFICE  $w^T x^{(i)} = w^T \left( \sum_{j \neq i} \alpha_j x^{(j)} \right) = \sum_{j \in S} \alpha_j \cdot w^T x^{(j)}$   $\rightarrow$  ~~✗~~



✓  
 پس  $W$  وجود ندارد ✓  
 در  $d_+$  نقطه  $d_+$  برآید

①  $\Rightarrow \forall \epsilon \dim(\text{Perception}) = d_+ \Rightarrow \text{lowest break point} = d_+$   
 ②

← جواب نهی  $\delta$



(۲) الف) ادعا کنیم  $VCdim(H) = n$

(I)  $VCdim(H) \geq n$  :

مجموعه  $\{e_1, e_2, \dots, e_n\}$  را در نظر بگیریم.  $e_i = \langle 0 \dots 0 1 0 \dots 0 \rangle$  اندکس  $i$

به ازای هر مقدار  $b_i$  به این  $n$  بردار می‌خواهیم یک  $h$  معین کنیم به شکلی باشد.

$$\Rightarrow b_1 \leftrightarrow e_1, b_2 \leftrightarrow e_2, \dots, b_n \leftrightarrow e_n$$

$$\Rightarrow S = \{i \mid b_i = 1\} \Rightarrow h_S(e_i) = \begin{cases} 1 & i \in S \\ 0 & i \notin S \end{cases} = \begin{cases} 1 & b_i = 1 \\ 0 & b_i = 0 \end{cases} \Rightarrow VCdim(H) \geq n \checkmark$$

(II)  $VCdim(H) \leq n+1$  :

به ازای هر  $n+1$  بردار  $e_1, \dots, e_{n+1}$  وجود خواهد داشت به شکلی خطی وابسته بردارهاست:

$$e_j = e_1 \oplus e_2 \oplus \dots \oplus e_k \Rightarrow b_1 = b_2 = \dots = b_k = 1, b_j = 0 \Rightarrow VCdim(H) \leq n+1$$

غیر ممکن است.

(I)  $\Rightarrow VCdim(H) = n$

(II)

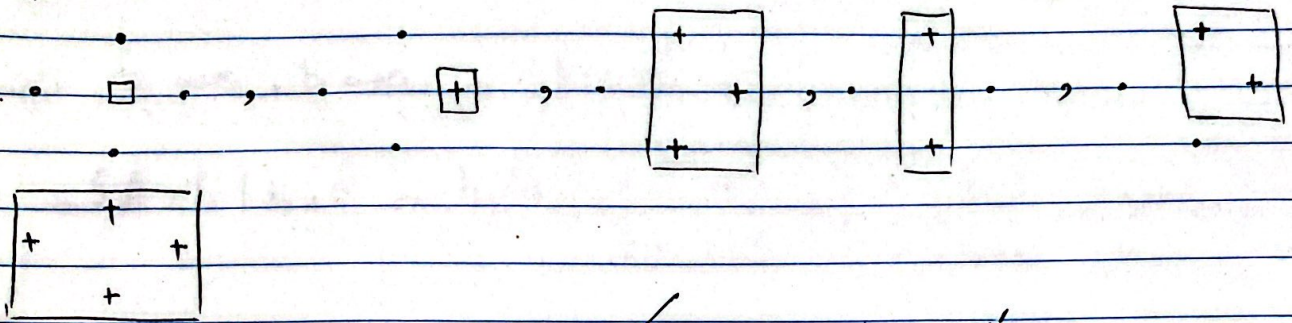


بسیار ادعای کنیم  $VCdim(H) = 2d$

(I)  $VCdim(H) \geq 2d$

$e_i^+ = [0 \dots 0 \ 1 \ 0 \dots 0]^T$  و  $e_i^- = [0 \dots 0 \ -1 \ 0 \dots 0]^T \Rightarrow E^+ = \{e_i^+\}, E^- = \{e_i^-\} \Rightarrow S = E^+ \cup E^-$   
 ← اندیس  $i$  ← اندیس  $i$

از  $d=2$  آغاز میکنیم:



الفن برای هر بعد جدید  $d'$ ،  $e_{d'}^+$  و  $e_{d'}^-$  را اضافه میکنیم.  
 موجه کنیم به همان تابع فرض برت آمده برای  $d'$ ، افلاص لازم، اضافه میکنیم تا  $d'$  را نیز دربرگیرد.

از اینجا  $e_{d'}^+$  و  $e_{d'}^-$  در لایحه کافی  $\{e_1^+, \dots, e_{d'}^+, e_1^-, \dots, e_{d'}^-\}$  هستند، افلاص کردن افلاص به مستقل جهت دربر گرفتن نقاط جدید ممکن است. (بسته به مقادیرهای مختلف به  $e_{d'}^+$  و  $e_{d'}^-$ ، حالات مختلف پیدا میکنند) - همراه

برای  $d=2$  به قدر و اگر برای  $d$  به قدر باشد، برای  $d+1$  هم به قدر است  $\Rightarrow VCdim(H) \geq 2d$

(II)  $VCdim(H) \leq 2d$

برای هر  $2d+1$  نقطه در فضا، ما یکم و نیمم در ایجاد متعلق را محاسبه میکنیم.

$\min_{x \in X} \{m_1, m_2, \dots, m_d\}$  و  $\max_{x \in X} \{M_1, M_2, \dots, M_d\}$ ،  $m_i = \min_{x \in X} x_i$ ،  $M_i = \max_{x \in X} x_i$

$\Rightarrow \exists x \forall k : m_k \leq x_k \leq M_k \Rightarrow$  اگر  $x_k$  برابر (ا-) مقادیر شود و یا  $x_k$  ها برابر (ا+)، آنگاه امکان پذیر نیست چون  $x_k$  داخل مستقل نخواهد بود.

(I)  $VCdim(H) = 2d$

(II) DATA OFFICE

به ازای هر  $n$  نقطه در فضای  $\mathbb{R}^d$ ، وجود خواهد داشت  $p$  به گونه‌ای که  $H_1 = \{ \text{sign}(p) \mid p: \mathbb{R}^d \rightarrow \mathbb{R} \}$   $n$  نقطه را پوشش دهد (توجه: در نظر بگیرید که دقیقاً همان نقاط را به  $\{ \pm 1 \}$  می‌تابانند).

$H = \{ \text{sign}(f \cdot g) \mid g: \mathbb{R}^d \rightarrow \mathbb{R}, f \in F \}$ 
 $\Rightarrow H_i \subseteq H$

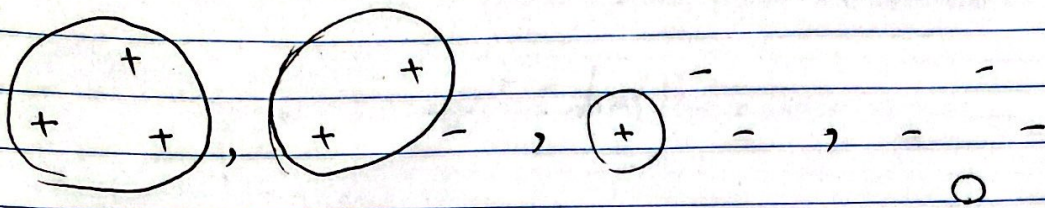
Scanned with CamScanner



$$vc \dim(H) = 3$$

(> ادعا می کنیم)

①  $vc \dim(H) \geq 3$



پس برای هر نقطه بالا، فضای فرض  $H$  را پوشش می دهد.

②  $vc \dim(H) < 4$

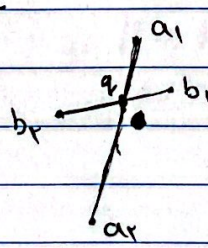
به ازای هر 4 نقطه در فضا :

۱) روی یک خط باشند :  $+ - + -$

۲) ۳ تا از نقاط روی یک خط باشند :  $-$   
 $+ - +$

۳) هیچکدام روی یک خط نباشند :  $+$   
 $-$   
 $+ +$

۴) هیچکدام روی یک خط نباشند :  $+$   
و convex باشند



ادعا می کنیم برای دو حالت  $\{a_1: +, a_2: +, b_1: -, b_2: -\}$  و  $\{a_1: -, a_2: -, b_1: +, b_2: +\}$  عمل نیست.

برهان خلف  $\rightarrow$  اگر هر دوی اینها عمل باشند  $C_1$  دایره مربوط به حالت ۱ و  $C_2$  دایره مربوط به حالت ۲  $\Rightarrow S = (C_1 / C_2) \cup (C_2 / C_1)$  را در نظر بگیرید.

$S$  شامل 4 ناحیه disjoint است. (و نهایت دو دایره نیست هم).

و در هیچکدام 4 ناحیه غیر مشترک وجود ندارد  $\Rightarrow vc \dim(H) < 4$

## مسئله‌ی ۲.

پاسخ.

الف) فرض کنید که  $\mathcal{E}(h_*) \leq \mathcal{E}(h^*) + \eta$  در آن صورت طبق نامساوی هافدینگ برای  $\gamma = \sqrt{\frac{1}{\gamma m} \log \frac{\gamma k}{\delta}}$  با احتمال حداقل  $1 - \delta$  داریم که

$$\begin{aligned}\hat{\mathcal{E}}(h_*) &\leq \mathcal{E}(h_*) + \gamma \leq \mathcal{E}(h^*) + \gamma + \eta \leq \mathcal{E}(\hat{h}) + \gamma + \eta \\ &\leq \hat{\mathcal{E}}(\hat{h}) + 2\gamma + \eta\end{aligned}$$

نامساوی های اول و آخر بر این اساس هستند که تمامی فرض ها در  $H$  باید در بازه خطای  $\gamma$  از خطای تعمیم اصلی خود باشند. نامساوی دوم بر اساس فرض سوال است. نامساوی سوم بر این فرض است که  $h^*$  کمینه کننده خطای تعمیم واقعی است و در نتیجه  $\mathcal{E}(h^*) \leq \mathcal{E}(\hat{h})$ .

با توجه به این موارد در نتیجه شرط برعکس یعنی

$$\hat{\mathcal{E}}(h_*) > \hat{\mathcal{E}}(\hat{h}) + 2\gamma + \eta$$

با احتمال حداکثر  $\delta$  روی می دهد.

(ب)

فرض کنید که  $\mathcal{E}(h_*) > \mathcal{E}(h^*) + \eta$

در آن صورت طبق نامساوی هافدینگ برای  $\gamma = \sqrt{\frac{1}{\gamma m} \log \frac{\gamma k}{\delta}}$  با احتمال حداقل  $1 - \delta$  داریم که

$$\begin{aligned}\hat{\mathcal{E}}(h_*) &\geq \mathcal{E}(h_*) - \gamma \geq \mathcal{E}(h^*) - \gamma + \eta \geq \mathcal{E}(h^*) - 2\gamma + \eta \\ &\geq \hat{\mathcal{E}}(\hat{h}) - 2\gamma + \eta\end{aligned}$$

نامساوی های اول و آخر بر این اساس هستند که تمامی فرض ها در  $H$  باید در بازه خطای  $\gamma$  از خطای تعمیم اصلی خود باشند. نامساوی دوم بر اساس فرض سوال است. نامساوی سوم بر این فرض است که  $\hat{h}$  کمینه کننده خطای empirical است.

با توجه به این موارد در نتیجه شرط برعکس یعنی

$$\hat{\mathcal{E}}(h_*) < \hat{\mathcal{E}}(\hat{h}) - 2\gamma + \eta$$

با احتمال حداکثر  $\delta$  روی می دهد.

ج) فرض کنید که  $h_* = h^*$ . با استفاده از نامساوی هافدینگ که  $\gamma = \sqrt{\frac{1}{\gamma m} \log \frac{\gamma k}{\delta}}$  با احتمال حداقل  $1 - \delta$  داریم که

$$\begin{aligned}\hat{\mathcal{E}}(h_*) &\leq \mathcal{E}(h_*) + \gamma = \mathcal{E}(h^*) + \gamma \\ &\leq \mathcal{E}(\hat{h}) + \gamma \leq \hat{\mathcal{E}}(\hat{h}) + 2\gamma\end{aligned}$$

نامساوی های اول و آخر بر این اساس هستند که تمامی فرض ها در  $H$  باید در بازه خطای  $\gamma$  از خطای تعمیم اصلی خود باشند. تساوی دوم بر اساس فرض سوال است. سومین نامساوی بر این اساس است که  $h^*$  کمینه کننده خطای تعمیم واقعی است.

حال توجه کنید که برای  $\eta$  و  $\delta$  ثابت، اگر  $m \rightarrow \infty$  داریم  $\gamma = \sqrt{\frac{1}{\gamma m} \log \frac{\gamma k}{\delta}} \rightarrow 0$

این بدین معنی است که برای  $m$  به اندازه کافی بزرگ  $4\gamma < \eta$  و معادلا  $\eta - 2\gamma < 2\gamma$ . در نتیجه با احتمال حداقل  $1 - \delta$  اگر  $m$  به اندازه کافی بزرگ باشد  $\hat{\mathcal{E}}(h_*) \leq \hat{\mathcal{E}}(\hat{h}) + \eta - 2\gamma$  در نتیجه الگوریتم جواب YES بر می‌گرداند.



6. [10 points] **Learning theory: Relaxed generalization bounds**

Let  $Z_1, Z_2, \dots, Z_m$  be independent and identically distributed random variables drawn from a Bernoulli( $\phi$ ) distribution where  $P(Z_i = 1) = \phi$  and  $P(Z_i = 0) = 1 - \phi$ . Let  $\hat{\phi} = (1/m) \sum_{i=1}^m Z_i$ , and let any  $\gamma > 0$  be fixed. Hoeffding's inequality, as we saw in class, states

$$\mathbb{P}(|\phi - \hat{\phi}| > \gamma) \leq 2 \exp(-2\gamma^2 m)$$

However, this relies on the assumption that the random variables  $Z_1, \dots, Z_m$  are all *jointly* independent. In this problem we will relax this assumption by only assuming *pairwise* independence among the  $Z_i$ . In this case we cannot apply Hoeffding's inequality, but the following inequality (Chebyshev's inequality) holds:

$$P(|\phi - \hat{\phi}| > \gamma) \leq \frac{\text{Var}(Z_i)}{m\gamma^2}$$

where  $\text{Var}(Z_i)$  denotes the variance of the random variable  $Z_i$  and for  $Z_i \sim \text{Bernoulli}(\phi)$  we have  $\text{Var}(Z_i) = \phi(1 - \phi)$ .

Given our hypothesis set  $\mathcal{H} = \{h_1, \dots, h_k\}$  and  $m$  pairwise but not necessarily jointly independent data samples  $(x, y) \sim \mathcal{D}$ , we now derive guarantees on the generalization error of our best hypothesis

$$\hat{h} = \underset{h \in \mathcal{H}}{\text{argmin}} \hat{\varepsilon}(h)$$

where as usual we define  $\hat{\varepsilon}(h) = \frac{1}{m} \sum_{i=1}^m 1\{h(x^{(i)}) \neq y^{(i)}\}$ , where  $(x^{(i)}, y^{(i)})$  are examples from the training set.

- (a) [2 points] What is the maximum possible value of  $\text{Var}(Z_i) = \phi(1 - \phi)$ ? From now on we will instead use this maximal value such that the bounds we derive hold for all possible  $\phi$ .

**Answer:** We find the maximum value by using the first and second order conditions. Differentiating and setting to 0 gives  $\phi = 1/2$ . By finding the second derivative ( $-2$ ), we confirm that this point is a maximum. Hence we substitute  $\text{Var}(Z_i)$  with  $1/4$  for the remainder of the question.

(b) [4 points] Let  $\gamma > 0$ .

- i. [2 points] Give a non-trivial (i.e. not the constant 1) upper bound on the probability that  $|\hat{\varepsilon}(\hat{h}) - \varepsilon(\hat{h})| > \gamma$ .
- ii. [1 points] Fix  $\delta \in (0, 1)$ . Using your upper bound, how large must the sample size  $m$  be before you can guarantee that

$$\mathbb{P}(|\hat{\varepsilon}(\hat{h}) - \varepsilon(\hat{h})| > \gamma) \leq \delta,$$

that is, that the training error and generalization error are within  $\gamma$  of one another with probability at least  $1 - \delta$ ?

- iii. [1 points] How does this sample size compare to what is achievable using Hoeffding's inequality?

**Answer:** We first use the Union bound to find

$$\begin{aligned} P(\exists h \in \mathcal{H}, |\varepsilon(h) - \hat{\varepsilon}(h)| > \gamma) &\leq \sum_{i=1}^k P(|\varepsilon(h_i) - \hat{\varepsilon}(h_i)| > \gamma) \\ &\leq \sum_{i=1}^k \frac{1}{4m\gamma^2} \text{ (using Chebyshev's inequality)} \\ &= \frac{k}{4m\gamma^2}. \end{aligned}$$

(Note that applying Chebyshev's inequality to  $\hat{h}$  does *not* work.) Setting this equal to  $\delta$  and solving for  $m$ , we find the solution:

$$m = \frac{k}{4\delta\gamma^2}$$

Hence the number of training examples required to make this guarantee is linear in  $k$  instead of logarithmic as when we used Hoeffding's inequality.



- (c) [4 points] Show that with probability at least  $1 - \delta$ , the difference between the generalization error of  $\hat{h}$  and the generalization error of the best hypothesis in  $\mathcal{H}$  (i.e. the hypothesis  $h^* = \operatorname{argmin}_{h \in \mathcal{H}} \varepsilon(h)$ ) is bounded by  $\sqrt{k/(m\delta)}$ .

**Answer:** First we solve for  $\gamma$  in the bound we found in (b):

$$\gamma = \sqrt{\frac{k}{4m\delta}}$$

Let  $h^* = \operatorname{argmin}_{h \in \mathcal{H}} \varepsilon(h)$ . By uniform convergence and the definition of  $\hat{h}$  (see Lecture Notes 4, page 7),

$$\varepsilon(\hat{h}) \leq \varepsilon(h^*) + 2\gamma$$

Hence  $|\varepsilon(\hat{h}) - \varepsilon(h^*)| \leq 2\gamma = \sqrt{k/(m\delta)}$  as desired.