

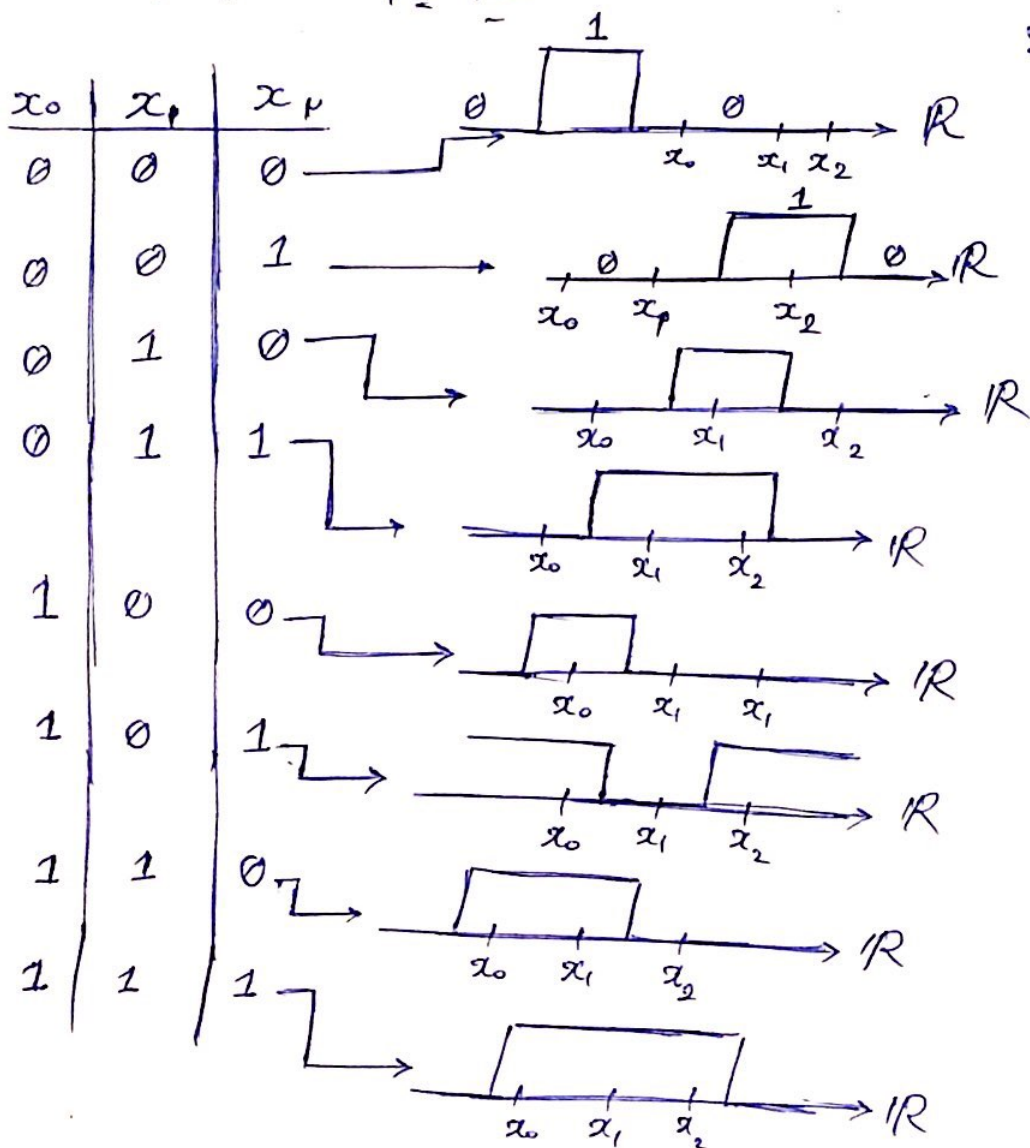
در لایه اول hidden layer - در گره  $v_{i1}$ ، ورودی آن را برابر قرار می دهیم با هفتی  
 ورودی ها  $x^1$  تا  $x^n$  و 1 و وزن این یال ها را برابر قرار می دهیم با  $u^1$  تا  $u^n$  و  
 c. که با این حالت ترکیب خطی  $(u^T x + c)$  محاسبه می شود در ورودی node  $v_{i1}$   
 و با استفاده از  $sign$  مقدار  $sign(u^T x + c)$  در خروجی این گره قرار می گیرد.  
 به همین طور مشابه این کار را برای گره ها دیگر این لایه با وزن ها  $(r_1, \dots, r_n)$   
 و  $(b_1, b_2, \dots, b_n)$  انجام می دهیم. در لایه بعدی جمع 3 گره قبلی را با وزن ها 1  
 محاسبه می کنیم که خروجی  $sign$  برابر با بیت اول خواهد بود. حال جمع دو تایی 3 گره لایه قبل  
 را با وزن یک محاسبه می کنیم:

نهایتاً 3 گره با این اهمیت ندارد

$$H = \{h_{a,b,s} : a \leq b, s \in \{-1, 1\}\}.$$

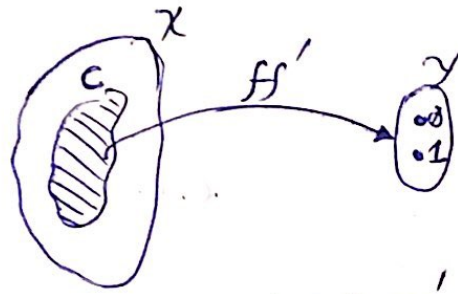
$$h_{a,b,s}(x) = \begin{cases} s & \text{if } x \in [a,b] \\ -s & x \notin [a,b] \end{cases}$$

حل. اگر تعداد نقاط خود را برابر با  $d=3$  در نظر بگیریم داریم:



تکین جدول بالا داریم که  $h$  بر  $C = \{x_0, x_1, x_2\}$  غالب است پس  $d \geq 3$ . اما برای  $(x_0, x_1, x_2, x_3) = (0, 1, 0, 1)$  داریم که حالت مقدار دهی  $C' = \{x_0, x_1, x_2, x_3\}$ ،  $d=3$  را فرض  $x_0 \leq x_1 \leq x_2 \leq x_3$  در اعطای  $h$  وجود ندارد پس  $\overline{d=3}$  خواهد بود

$$H' \subseteq H \rightarrow VCdim(H') \leq VCdim(H)$$



زمن کنید مجموعه  $C$ ، مجموعه‌ای باشد که  $H'$  بر آن غالب است و حال چون  $H' \subseteq H$  است  
 های  $H'$  در  $H$  هستند تمام حالت‌ها مقداردهی  $C$  را پوشش می‌دهند، در  $H$  نیز قرار دارند  
 پس  $H'$  نیز بر  $C$  غالب است و چون  $H$  اعضای بیشتری نسبت به  $H'$  دارد ممکن  $C$  با  
 $|C'| > C$  وجود داشته باشد که  $H$  بر آن غالب باشد پس خواهیم داشت:

$$VCdim(H') \leq VCdim(H)$$

$$X = \mathbb{R}^2, H = \{h_{a_1, b_1, a_2, b_2} : a_1 < a_2, b_1 < b_2\}$$

$$h_{a_1, b_1, a_2, b_2}(x) = \begin{cases} 1 & a_1 \leq x^1 \leq a_2, b_1 \leq x^2 \leq b_2 \\ 0 & \text{o.w} \end{cases}$$

با توجه به اینکه  $d = VCdim(H) = 4$  در قسمت اول به دست آمد و همچنین بازه‌ای نه برای

agnostic PAC learning داریم:

$$C_1 \frac{d + \log \frac{1}{\epsilon}}{\epsilon^2} \leq m_{\text{agg}}(\epsilon, \delta) \leq C_2 \frac{d + \log \frac{1}{\delta}}{\epsilon^2}$$

$$\Rightarrow C_1 \frac{\log \frac{1}{\delta} + 4}{\epsilon^2} \leq m_{\text{agg}}(\epsilon, \delta) \leq C_2 \frac{4 + \log \frac{1}{\delta}}{\epsilon^2}$$



سوال T11 (قسمت ب)

$$\leq C_2 \frac{d + \log \frac{1}{\epsilon}}{\epsilon^2}$$

بالا حد  $\text{bound}$  ای که بدست آوریم یعنی  $m(\delta, \epsilon)$ ، این اطلاعات را به ما

می دهد که برای اینکه در حالت  $\text{infinite class}$ ،  $\text{agnostic PAC learnable}$  باشد

لازم نیست که تعداد نمونه های ما از کل جمعیت به سمت  $\infty$  میل کند و می توان با مقداری حقیقی

که  $m < \infty$  است، این خاصیت در مورد  $f$  داشت و حد پایین عبارت یعنی

$$m(\epsilon, \delta) \gg C_1 \frac{d + \log \frac{1}{\epsilon}}{\epsilon^2}$$

خاصیت مورد نظر را داشت و این تعداد نمونه ها باید از یک حدای که وابسته به  $\epsilon$  و  $\delta$  است بیشتر باشد.

سوال T12

الف (1) با توجه به شرط  $\text{Realizability}$  که می بینیم  $f$  داریم پس کم ترین متوسط

خطا بدرومی کل جمعیت از جمله هر  $\text{Sample Set}$  برابر با صفر خواهد بود. حال با توجه به این که

فرضیه ای که  $A$  از یادگیری بدرومی داده آموزشی داده است، دارای متوسط خطای صفر است و هم

چنین چون در روش  $\text{ERM}$ ، فرضیه ای بدست می آید که دارای کمترین خطا روی  $S$  باشد و

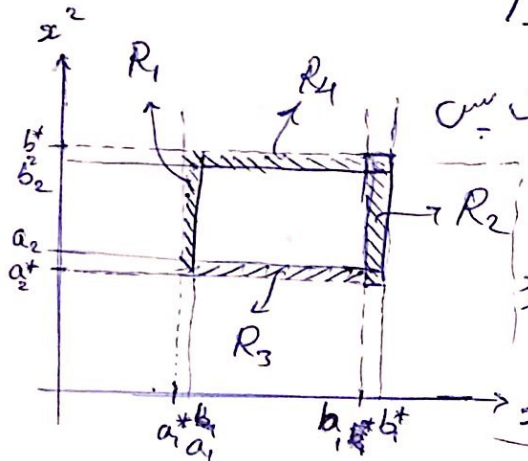
با توجه به این که شرط  $\text{Realizability}$  برقرار است، این خطا برابر با صفر است پس فرضیه

الگوریتم  $A$  می تواند فرضیه الگوریتم  $\text{ERM}$  باشد چون کمترین خطا  $h$  بدرومی  $S$

که قابل دسترسی بوده است، داده است.

الف - 2) 
$$\text{if } m \geq \frac{4 \log \frac{4}{\epsilon}}{\epsilon} \text{ with at least Prob. } 1-\delta \quad |L_D(h_A) - L_S(h_A)| \leq \epsilon.$$

$h_A$ : خروجی الگوریتم  $A$



1) : با توجه به اینکه متوسطی  $R(S)$  برابر با صفر است پس

هیچ نقطه‌ای را اشتباه دسته‌بندی نمی‌کند حال چون در  $R_i$  ها با احتمال  $\frac{\epsilon}{4}$  ممکن است نقطه‌ای وجود داشته باشد که با label آن را نمی‌دانیم پس احتمال

خطا وجود دارد پس نمی‌تواند  $R(S)$  را از  $R_i$  ها برود پس در  $R^*$  قرار دارد.

$$\Rightarrow R(S) \subseteq R^*$$

2) : حال اگر فرض کنیم که  $R(S)$  شامل تمام نقاط  $R_i$  ها باشد، در بدترین حالت ممکن

با احتمال  $\epsilon$ ، نقطه‌ای وجود داشته باشد که negative باشد و در این صورت اشتباه تشخیص داده می‌شود که در این صورت خطای ما برابر  $\epsilon$  خواهد بود:

$$L_S(h_A) \leq \frac{1}{\text{loss prob.}} \times \epsilon = \epsilon$$

3) :

$$\text{Prob}_{S \subset D} \left[ \underbrace{R_i \not\subseteq S}_{A_i} \right] = \text{Prob}_{\substack{x_j \in D \\ j=1,2,\dots,m \\ x_j \in S}} [x_j \notin R_i] = \left(1 - \frac{\epsilon}{4}\right)^m$$

4) :

$$\text{Prob} \left\{ \bigcup_{i=1}^4 A_i \right\} \leq \sum_{i=1}^4 \text{Prob}[A_i] \leq 4 \left(1 - \frac{\epsilon}{4}\right)^m \leq 4e^{-\frac{m\epsilon}{4}}.$$

سین خواهم داشت:

$$\text{Prob} [L_D(h_S) > \epsilon] \leq \text{Prob} \left\{ \bigcup_{i=1}^4 A_i \right\} \leq \frac{4e^{-\frac{m\epsilon}{4}}}{\delta}$$

$$\Rightarrow m \geq \frac{4 \log \frac{4}{\delta}}{\epsilon}$$

ب) مقدار  $m_{ff}(\epsilon, \delta)$  نسبت به  $\epsilon$  و  $\delta$  - نابینا نیست - شرط Realizability بهتر است

محاسبه شده و در نتیجه حد پایین خوبی نسبت به  $\frac{4 \log \frac{4}{\delta}}{\epsilon^2}$  نسبت آوردیم و آن را بگیرد (PAC learning)

بخشنامه ام.

$$l(h(x, y)) = \begin{cases} 1 & h(x) \neq y \\ 0 & h(x) = y \end{cases}$$

$$E_{\text{bayes}} = E_x E_{y|x} [l(h(x, y))]$$

برای گفته کردن مقدار  $E_{\text{bayes}}$  - کافی است که برای هر  $x$  در  $X$  مقدار  $E_{y|x} [l(h(x, y))]$  را بگیریم که در آن صورت مقدار متوسط آن یعنی  $E_{\text{bayes}}$  نیز گفته خواهد شد.

$$\text{Prob} [l(h(x, y)) | x=x] = \begin{cases} \pi(x) & h(x) = 0 \\ 1 - \pi(x) & h(x) = 1 \end{cases}$$

حال برای ایدیه فنی این احتمال داریم:



ادام سوال T/3 صفر

$$E_{y|x} [l(h(x,y)) | x=x] = \begin{cases} \pi(x) & h(x)=0, \text{ صفر} \\ 1-\pi(x) & h(x)=1, \text{ یک} \end{cases}$$

پس متوسط خطای ما برای هر  $x$  یا  $\pi(x)$  است و یا  $1-\pi(x)$ . حال کمترین متوسط خطا برابر است با کمترین مقدار بین دو عبارت  $\pi(x)$  و  $1-\pi(x)$ ، برابر هر  $x \in X$ ، پس خواهیم داشت:

$$\forall x \in X: \min_{y|x} E[l(h(x,y))] = \min(\pi(x), 1-\pi(x))$$

از طرفی برابر تخمین bayes داریم:

$$f(x) = \begin{cases} 1 & \pi(x) \geq \frac{1}{2} \\ 0 & \pi(x) < \frac{1}{2} \end{cases} \Rightarrow l(f_D, (x,y)) = \begin{cases} 1-\pi(x) & \pi(x) \geq \frac{1}{2} \\ \pi(x) & \pi(x) < \frac{1}{2} \end{cases}$$

$$\Rightarrow E_{y|x} [l(f_D, (x,y))] = \begin{cases} 1-\pi(x) & \pi(x) \geq \frac{1}{2} \\ \pi(x) & \pi(x) < \frac{1}{2} \end{cases}$$

$$= \min(\pi(x), 1-\pi(x)).$$

متوسط خطای تخمین کر bayes، دقیقاً برابر با کمترین خطای که است که کل  $h$  ها ممکن می‌توانند داشته باشند.  $\Leftarrow$

$$L_D(f) = \epsilon_{\text{bayes}} \leq L_D(h) : \forall h$$

سوال T/4 الف) با توجه به i.i.d بودن  $l(h(x_i, y_i))$  ها برای  $h$  از  $\epsilon$  و  $\eta$  استفاده می‌کنیم.

$$\text{Prob} \left[ \left| L_D(h) - \underbrace{\frac{1}{m} \sum_{i=1}^m l(h(x_i, y_i))}_{\eta} \right| > \underbrace{\epsilon}_{\eta} \right] \leq 2e^{-2m\epsilon^2} \quad : \epsilon = \eta$$

$$\Rightarrow \text{Prob} [L_D(h) > 2\eta] \leq 2e^{-2m\eta^2}$$

ب) از آن جایی که نمی‌دانیم که آیا الگوریتم  $A$ ،  $h_S$  را بر اساس  $(\mathcal{L}(h_S(x_i)))$  های محدود یا نه یا به بیان دیگر این که  $h_S$ ، با توجه به معادله خطاها بدست می‌آید یا نه، نمی‌توانیم از Hoeffding استفاده کنیم چرا که ممکن است  $(\mathcal{L}(h_S(x_i)))$  به هم وابسته باشند  $\therefore$

حال چون اگر محدود است پس  $\mathcal{L}$  دارد و می‌توانیم از bound آن استفاده کنیم:

$$\forall h \in \mathcal{H} : \text{Prob}[|L_D(h) - L_S(h)| > \epsilon] \leq 2|\mathcal{H}| e^{-2m\epsilon^2}$$

حال با قرار دادن  $L_S(h_S) = \gamma$  و  $\epsilon = \gamma$  داریم:

$$\text{Prob}[|L_D(h) - \gamma| > \gamma] \leq 2|\mathcal{H}| e^{-2m\gamma^2}$$

$$\Rightarrow \text{Prob}[L_D(h) > 2\gamma] \leq 2|\mathcal{H}| e^{-2m\gamma^2}$$

(ج)

$$\epsilon_{\text{Bayes}} = 3\eta \Rightarrow \forall h \in \mathcal{H} \text{ که ممکن است} : L_D(h) \geq 3\eta$$

پس خواهیم داشت که:

$$\text{Prob}[L_D(h) > 2\eta] = 1$$

با توجه به  $\epsilon_{\text{Bayes}}$  ای که در این مسئله داده شده می‌دانیم که نمی‌تواند متوسط خطا برابر هر  $h$  روی توزیع  $\mathcal{D}$  کم‌تر از  $3\eta$  باشد، حال احتمال آنکه برای  $h$  رندم در قسمت الف متوسط خطا بر روی کل distribution بزرگ‌تر از  $2\eta$  باشد  $\frac{1}{2}$  است. ببرگرد  $\mathcal{D}$  صحت است.

(د) مفید به نامسادی Hoeffding بر این اساس که ما اطلاعاتی در مورد توزیع داده نداریم فرض و بر اساس این نامسادی مورد نظر بدست می‌آید، حال اگر اطلاعاتی در مورد توزیع داشته باشیم قطعاً بر روی احتمالات ما ~~تغییر~~ تاثیر خواهد داشت و این را چگونه می‌بخشد.