

## سؤال ۱

در VC-dim، VC-theory یک اندازه گیری از ظرفیت (پیچیدگی) مجموعه ای از توابع است که می تواند توسط یک الگوریتم طبقه بندی باینری آموخته شود. برای آن که VC-dim را به صورت عمیق تر بررسی کنیم باید بدانیم که:

ما می گوئیم  $\mathcal{H}$  یک مجموعه می رود نقاط  $C$  داخل domain را shatter می کند اگر  $\mathcal{H}$  تمامی توابعی که از این مجموعه می  $C$  به صفر یا یک می روند را بسازد یعنی تمامی حالات را داشته باشیم.  $(\mathcal{H}_C = \{h(c_1), h(c_2), \dots, h(c_m) : h \in \mathcal{H}\})$

حال برای VC-dim داریم: بیش ترین عضو مجموعه  $C$  به صورتی که  $\mathcal{H}$  می تواند  $C$  را shatter کند. برای آن که نشان دهیم  $VC-dim(\mathcal{H}) = d$  می باشد داریم:

- ①  $\exists$  set  $C$  of size  $d$  that is shattered by  $\mathcal{H}$
  - ② VC of size  $d+1$  that is not shattered by  $\mathcal{H}$
- حداکثر مقدار VC-dim زمانی است که برای هر  $d$  یکی از حالات output را داشته باشیم بنا بر این خواهیم داشت

$$VC-dim(\mathcal{H}) \leq \lceil \log_2(|\mathcal{H}|) \rceil$$

حال که با موضوع VC-dim آشنا شدیم می خواهیم برویم سراغ تئوری اصلی یادگیری ماشین و ارتباط آن با VC-dim. بر اساس قضیه ی اساسی یادگیری آماری داریم:

let  $\mathcal{H}$  be a hypothesis class of function from  $X$  to  $\{0,1\}$  and the loss function be 0-1 loss then the followings are equivalent:

- ①  $\mathcal{H}$  is Agnostic PAC learnable
- ②  $\mathcal{H}$  is PAC learnable

③ Any ERM rule is PAC-learnable for  $\mathcal{H}$

④  $\mathcal{H}$  has a finite VC-dim

همچنین از (quantitative V) theory of learning داریم:

Assume that  $\text{VC-dim}(\mathcal{H}) = d < \infty$  then there are constants

$c_1, c_2, c_1', c_2'$  such that

①  $\mathcal{H}$  is Agnostic PAC-learnable with sample complexity

$$c_1 \frac{d + \log(\frac{1}{\delta})}{\epsilon^2} \leq m_{\mathcal{H}}(\epsilon, \delta) \leq c_2 \frac{d + \log(\frac{1}{\delta})}{\epsilon^2}$$

②  $\mathcal{H}$  is PAC-learnable with sample complexity:

$$c_1' \frac{d + \log(\frac{1}{\delta})}{\epsilon} \leq m_{\mathcal{H}}(\epsilon, \delta) \leq c_2' \frac{d + \log(\frac{1}{\delta})}{\epsilon}$$

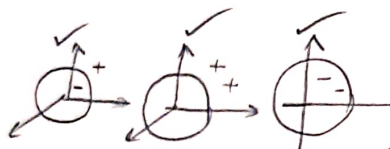
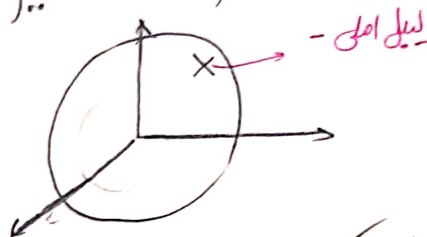
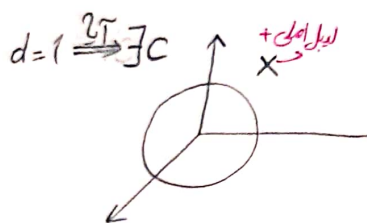
در واقع اگر به من یک  $\mathcal{H}$  ای دهند اول باید چک کنم که Agnostic PAC یا PAC قابل یادگیری هست یا نه. کلید این کار VC-dim است. اگر یود بود ✓ اوکی و سپس می فهمیم از چ اروری sample نیاز داریم. حلال که با اصل قضیه آشنا شدیم داریم

**الف** این classifier در واقع مدل کشنده‌ی یک کره است. همچنین از بخش قبل می‌دانیم برای آن که بگوییم  $d = \dim(V_C) = ?$  - (اینجا  $V_C$  به مجموعه‌ی داده‌ها اشاره دارد) - می‌تواند و با کلاس

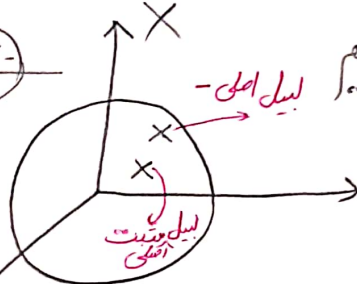
set  $C$  of size  $d$  that is shattered by  $H$

$\forall C$  of size  $d+1$  " " not " " " "

بنابراین از شروع می‌کنیم و جابجایی روییم



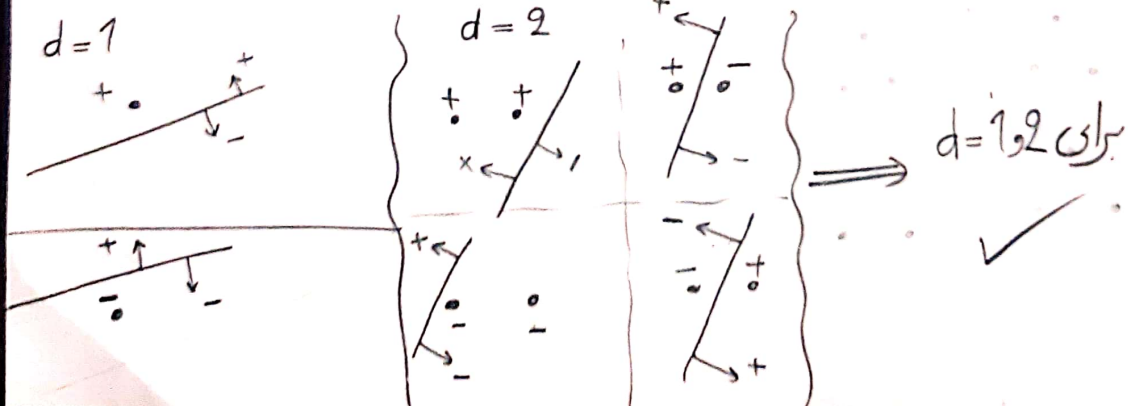
$d=2 \Rightarrow \exists C$



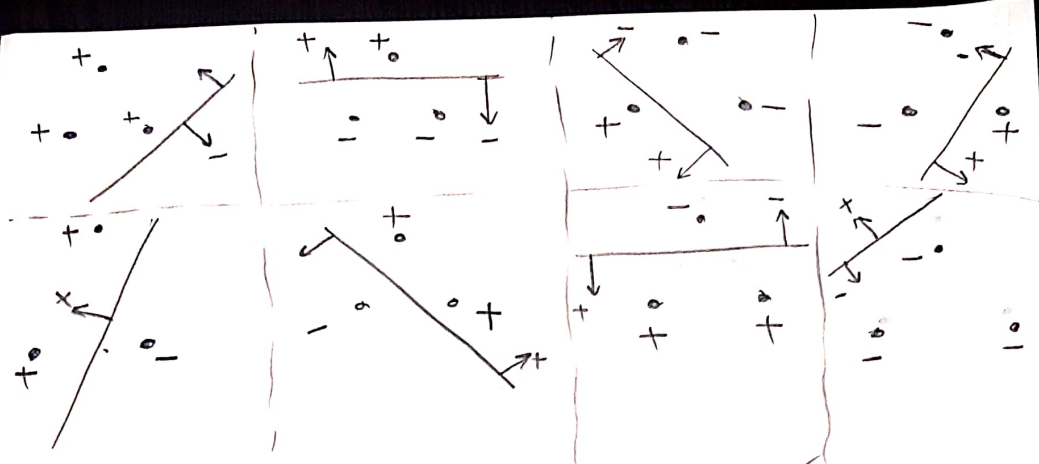
$\forall C$  ✓ حال برای  $d=2$  داریم (اینجا  $V_C$  داریم یعنی)

classifier داده‌های بالا را نمی‌تواند از هم جدا کند بنابراین  $V_C - \dim = 1$

**ب** برای  $\text{sgn}(x_1\theta_1 + x_2\theta_2 + \theta)$  داریم:







بر  $d=3$  ✓  
اما برای  $d=4$  داریم:

به هیچ روشی نمی تواند

$\checkmark$   $\text{vc-dim}=3$

ج] برای حالت کلی  $N$  داریم  $N+1$

اثبات: فرض می کنیم  $d+1$  نقطه فضای  $d$  بعدی داریم

$$x_i = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{id} \end{pmatrix} \text{ for } i=1, \dots, d+1$$

حالتی خواهیم نشان دهیم آرایش داخلی از این  $d+1$  نقطه و همواره  $d$  خط جدا کننده (shattering) پذیر باشد. می دانیم  $h(x)$  در Half-space ها این گونه predict می کند

$$h(x_i) = \text{sgn}(w_1 x_{i1} + \dots + w_d x_{id} + b)$$

فرض کنید زوجی هائی هدف  $y_i$  باشد،  $i = 1, \dots, d+1$   
 فرض کنید  $y_1$  تا  $y_{d+1}$  یک labeling راغوان برای  $x_1$  تا  $x_{d+1}$  باشد  
 $d+1$  معادله خواهیم داشت

$$\text{sgn}(w_1 x_{11} + \dots + w_d x_{1d} + b) = y_1$$

$$\text{sgn}(w_1 x_{(d+1)1} + \dots + w_d x_{(d+1)d} + b) = y_{d+1}$$

اگر معادله را به صورت زیر بنویسیم و جواب داشته باشد معادله ای بالا هم جواب خواهد داشت:

$$w_1 x_{11} + \dots + w_d x_{1d} + b = y_1$$

$$w_1 x_{(d+1)1} + \dots + w_d x_{(d+1)d} + b = y_{d+1}$$

که معادله ای؟ لا را می توان به صورت  $AX = Y$  نوشت. اگر آرایش نمونه ها طوری باشد

که ماتریس  $(d+1) \times (d+1)$   $A$  وارون پذیر باشد  $w_1$  تا  $w_d$  و  $b$  به صورت

یکسان به دست خواهند آمد

برای  $d+1$  نقطه آرایش خاصی وجود دارد که؟  $\text{shattering}$  امکان

پذیر است

حال فرض کنید  $d+1$  نقطه داشته باشیم. طبق قضیه Radon هر مجموعه  $d+2$  تایی از نقاط در  $\mathbb{R}^d$  می تواند به دو مجموعه  $d$  تایی تقسیم شود به طوری که convex hull آن ها تهی باشد. حال اگر یکی از این مجموعه ها را  $label$  + دیگری را  $label$  منفی بنویسیم نمی توان با همیچ  $half$  space ای نقاط داخل اشتراک را درست  $label$  زد

$$\implies VC\text{-dim}(\mathcal{H}) < d+1$$

طبق توضیحات بالا  $VC\text{-dim}(\mathcal{H}) = d+1$

## سوال ۲:

الف) برای به دست آوردن رگرسیون خطی بر روی داده‌ها داریم:

$$\hat{y}^{(i)} = \theta^T x^{(i)} \rightarrow y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}, \text{ where } \epsilon^{(i)} \sim N(0, \sigma^2), \quad \epsilon^{(i)} \text{ are i.i.d.}$$

با توجه به قاعده maximum likelihood داریم:

$$P(Y|X, \theta) = \prod_{i=1}^n P(y^{(i)}|x^{(i)}, \theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (y^{(i)} - \theta^T x^{(i)})^2\right)$$

باید  $Y$  به گونه‌ای انتخاب شود که این احتمال را ماکزیم کند؛ یا به طور معادل عبارت  $-\ln(P(Y|X, \theta))$  را مینیمم کند:

$$-\ln(P(Y|X, \theta)) = -\sum_{i=1}^n -\frac{1}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$$

با توجه به این‌که برای همه مقادیر  $Y$  مقدار  $\sigma$  یکسان است؛ کفایست عبارت زیر را مینیمم کنیم که همان مینیمم کردن خطای MSE در رگرسیون خطی است:

$$\min \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$$

## سوال ۳:

$$f_X(x, \mu, \sigma^2) = \prod_{i=1}^n f(x_i, \mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (x_i - \mu)^2\right)$$

$$\rightarrow \ln f_X(x, \mu, \sigma^2) = -\sum_{i=1}^n -\frac{1}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

## سوال ۴:

الف) در رگرسیون خطی باید تابع هزینه که به صورت MSE تعریف می‌شود را کمینه کنیم: یعنی:

$$\min J(\theta) = \sum_{i=1}^n (h_{\theta}(x_i) - y_i)^2$$

می‌توانیم  $h_{\theta}(x_i)$  را برای رگرسیون خطی به صورت زیر بنویسیم:

$$h_{\theta}(x_i) = \theta^T x^{(i)} \rightarrow \min \sum_{i=1}^n (\theta^T x^{(i)} - y_i)^2$$

$$\rightarrow J(\theta) = (X\theta - y)^T (X\theta - y)$$

$$\frac{\partial J(\theta)}{\partial \theta} = 2X^T X\theta - 2X^T y = 0 \rightarrow \theta = (X^T X)^{-1} X^T y$$

ب) با اضافه کردن L2 Regularization در تابع هزینه داریم:

$$J(\theta) = \sum_{i=1}^n (\theta^T x^{(i)} - y_i)^2 + \lambda ||\theta||_2^2 = (X\theta - y)^T (X\theta - y) + \lambda \theta^T \theta$$

$$\frac{\partial J(\theta)}{\partial \theta} = 0 \rightarrow 2X^T X\theta - 2X^T y + 2\lambda \theta = 0 \rightarrow (X^T X + \lambda I)\theta = X^T y$$

$$\rightarrow \theta = (X^T X + \lambda I)^{-1} X^T y$$