

«به نام خدا»

تقریب سدی بنجم

۹۵۱۰۵۳۰۲

حسین ابراهیمی

سوال T20

symmetric

الف (۱). برای ضرب داخلی که در فضای F صورت می گیرد، خاصیت

داریم پس :

$$\langle \psi(x), \psi(x') \rangle = \langle \psi(x'), \psi(x) \rangle$$

$$\Rightarrow K(x, x') = K(x', x) \quad \square$$

ب (۲). می دانیم جواب بهینه مسئله در فضای F برابر است با :

$$w^* = \sum_{i=1}^m d_i \psi(x_i) \quad i=1, \dots, m, d_i \in \mathbb{R}$$

پس داریم :

$$\|w^*\|^2 \geq 0 \Rightarrow \langle w^*, w^* \rangle = \left\langle \sum_{i=1}^m d_i \psi(x_i), \sum_{j=1}^m d_j \psi(x_j) \right\rangle \geq 0$$

$$\Rightarrow \sum_{i=1}^m \sum_{j=1}^m d_i d_j \langle \psi(x_i), \psi(x_j) \rangle \geq 0$$

$$\Rightarrow \sum_{i=1}^m \sum_{j=1}^m d_i d_j K(x_i, x_j) \geq 0 \quad \square$$

$\Rightarrow K$ is positive semi-definite function.

(ب)

$$\sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j K(x_i, x_j) \stackrel{(*)}{=} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j G_{ij}$$

$$G_{ij} = K(x_i, x_j) \quad (*)$$

$$\Rightarrow \sum_{i=1}^m \alpha_i \sum_{j=1}^m \alpha_j G_{ij} = \underline{\underline{d^T G d \geq 0}}, \quad d \in \mathbb{R}^m / 0$$

G is real matrix \Rightarrow G is $\underbrace{\text{semi-definite matrix}}_{\text{Positive}}$.

سوال T21

$$l \leq d \Rightarrow l = 0, 1, \dots, d$$

(الف)

$$\Rightarrow l=0 \rightarrow 1$$

$$\Rightarrow l=1 \rightarrow K$$

$$\Rightarrow |X_d| = 1 + K + \dots + K^d$$

مقدار هندسی

$$l=d \rightarrow K^d \Rightarrow |X_d| = \frac{K^{d+1} - 1}{K - 1}$$

همان طور که می بینیم با افزایش d مقدار $|X_d|$ بصورت فزاینده با K^{d+1} افزایش می یابد و نتیجه خواهم داشت که:

$$|X_d| = \underbrace{\frac{1}{1 - 1/K}}_{\text{Constant}} K^d - \underbrace{\frac{1}{K-1}}_{\text{مهم}} \Rightarrow |X_d| = O(K^d)$$

(ب)

$$w^u = \begin{cases} 1 & u=v \\ 0 & \text{o.w} \end{cases} \Rightarrow w = (0, 0, \dots, 1, \dots, 0, 0, \dots)$$

مؤلفه مربوط به ویژگی v در v feature Set

$$\begin{aligned} \Rightarrow \|w\|^2 &= \langle w, w \rangle = \langle (0, \dots, 1, \dots, 0), (0, \dots, 0, 1, \dots, 0) \rangle \\ &= 1 \Rightarrow \|w\| = 1 \end{aligned}$$

حال برای $\gamma = \frac{1}{2}$ margin کافی است تا برابر با $-\frac{1}{2}$ قرار دهیم در این صورت

خواهیم داشت: اگر x دارای زیربنا v باشد در این صورت: $\psi^v(x) = 1$
 اگر x' ویژگی v را در زیربنای خود نداشته باشد: $\psi^v(x') = 0$

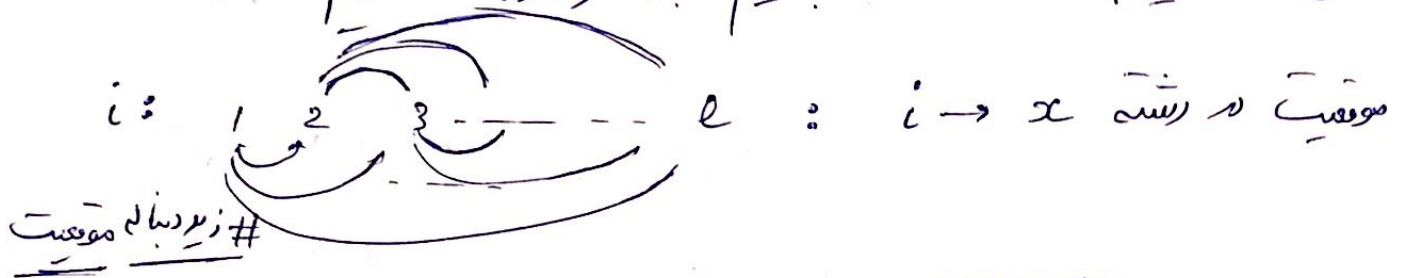
$$x: \quad \underbrace{\langle w, \psi(x) \rangle}_1 - \frac{1}{2} = \frac{1}{2} \geq \frac{1}{2} \quad \checkmark \quad \text{در این صورت داریم:}$$

$$x': \quad \underbrace{\langle w, \psi(x') \rangle}_0 - \frac{1}{2} = -\frac{1}{2} \leq -\frac{1}{2} \quad \checkmark$$

پس $(w, -\frac{1}{2})$ تعریف شده $\frac{1}{2}$ margin خواسته شده را دارند.

$$\|\psi(x)\|^2 = \langle \psi(x), \psi(x) \rangle = \# \text{ of substrings of } x. \quad (ج)$$

کافی است زیربالات x را محاسبه کنیم. برای هر i در x داریم:



$$1 : l$$

$$2 : l-1$$

⋮

$$l : 0$$

$$\Rightarrow l + (l-1) + \dots + 0 = \left\lfloor \frac{l(l-1)}{2} \right\rfloor < l^2$$

$$\Rightarrow \|\psi(x)\|^2 \leq l^2 \leq d^2 \Rightarrow \|\psi(x)\| \leq d$$

$$\Rightarrow \|\psi(x)\| = O(d)$$

⑤ باتوجه به تابع ψ ، فایلهای که شامل ویرگول ψ نباشند به نقطه صفر محور مربوط ψ در F نگاشت میشوند فایلهای که دارای ψ نباشند به نقطه 1 این محور پس مسئله حداقل با (w, b) معرفی شده در قسمت (ب) $Seperable$ است. حال چون مسئله $Seperable$ است از $Hard-SVM$ استفاده نمیکنیم تا بیشترین $margin$ درست ببار.

$$K(x, x') = \langle \psi(x), \psi(x') \rangle \quad (5)$$

حال مقدار $\psi(x)$ براساس ψ باید زیر دنباله d در x قرار دارد یا نه، حال مولفه های

$$\langle \psi(x), \psi(x') \rangle = \psi_1(x) \cdot \psi_1(x') + \dots + \psi_s(x) \cdot \psi_s(x') \quad \text{داریم:}$$

حال مقدار هر کدام از عبارت بالا وقتی 1 است که دنباله ψ در هر دوی آن ها وجود داشته باشد حال جمع آن ها برابر خواهد بود با تعداد زیر دنباله ها مشترکی که هر دو فایلهای x دارند.

$$K(x, x') = \langle \psi(x), \psi(x') \rangle = \# \text{ of common substrings}$$

⑥ کافی است چک کنیم که آیا زیر دنباله های فایلهای x در زیر دنباله های فایلهای x' وجود دارند یا نه و تعداد آن ها را بشماریم. تعداد زیر دنباله های x از $O(d^2)$ میباشد.

حال کافی است هر کدام از آن ها را با $O(d^2)$ تا زیر دنباله های فایلهای x' از $O(1)$

مقایسه کنیم که هزینه این عمل از $O(d^2)$ است. در کل حداکثر d^2 بار این کار را انجام میدهم پس پیچیدگی الگوریتم حداکثر از $O(d^4)$ است که برای اعضای مختلف x

① بله، بایستی به (w, b) ای که در قسمت (ب) تعریف کردیم داریم که برای هر $(x, y) \in D$ ، $\psi(x) \in \mathbb{R}^n$ را به محور مربوط v نگاشت می‌کنیم. در آن محور x های دارای ویرس به 1 و x های نه ویرس ندارند به 0 نگاشت داده می‌شود که بایستی به آن *seperable* هستند.

$$w^* = w_{(b)} = \begin{cases} w^u = 1 & : u = v \\ w^u = 0 & : 0 \cdot w \end{cases}, \quad b^* = -\frac{1}{2}$$

$$\Rightarrow \min_{\psi(x)} y_i (\langle w^*, \psi(x) \rangle + b^*) = \boxed{\frac{1}{2} = \gamma}$$

$$\max_{x \in \mathbb{R}^d} \|\psi(x)\|^2 = \max_{x \in \mathbb{R}^d} \langle \psi(x), \psi(x) \rangle = d^2$$

«قسمت ج»

$$\Rightarrow \|\psi(x)\| \leq d \Rightarrow \boxed{\rho = d}$$

(ع) باتوجه به اینکه مسئله به صورت خطی در فضای F قابل جداسازی است، خطای آموزش Hard-SVM بر روی داده آموزش S برابر با صفر خواهد بود:

$$\Rightarrow L_S(h) = \frac{1}{m} \sum_{i=1}^m d^{0-1}((w, b), (x_i, y_i))$$

طبق قضیه 15.4 در کتاب (w^*, b^*) مربوط به این فزیم هموژن داشته باشند، بین منظور (w, b) که در قسمت (ز) بدست آمده با بردن در فضای \mathbb{R}^{d+1} و سپس Normalize کردن برای اینکه $\|w\|=1$ شود به فزیم هموژن تبدیل

یکنیم:

$$w^* = (0, 0, \dots, 1, 0, \dots, 0, -\frac{1}{2}) \Rightarrow \|w^*\| = \frac{\sqrt{5}}{2}$$

← مربوط به γ

$$\Rightarrow \tilde{w} = \frac{w^*}{\|w^*\|} \Rightarrow \tilde{w} = (0, 0, \dots, 0, \frac{2}{\sqrt{5}}, 0, \dots, 0, -\frac{1}{\sqrt{5}})$$

$$15.4 \Rightarrow L_D(\tilde{w}) \leq \sqrt{\frac{4 \left(\frac{d}{\frac{1}{\sqrt{5}}} \right)^2}{m}} + \sqrt{\frac{2 \log(2/\delta)}{m}}$$

$$\leq \frac{2d}{\sqrt{5}m} + \sqrt{\frac{2 \log(2/\delta)}{m}} \quad (\gamma = \frac{1}{\sqrt{5}}, \rho = d)$$

(ب) باتوجه به قسمت قبل، از ناهمبستگی بالا استفاده می‌کنیم:

با اتصال حاصل 0.99 داریم:

$$\frac{2 \times 10^3}{\sqrt{5} \sqrt{m}} + \frac{\sqrt{2 \times \ln \frac{2}{10^{-2}}}}{\sqrt{m}} \leq \frac{1}{100}$$

$$\Rightarrow \sqrt{m} \geq \left(\frac{2 \times 10^3}{\sqrt{5}} + \sqrt{2 \ln \frac{2}{10^{-2}}} \right) \times 10^2 \approx 900$$

$$\Rightarrow \boxed{m \geq 900^2 = 810000}$$

کجایی). از حد بالای بدست آمده در PAC استعاره می‌کنیم. چون الگوریتم یادگیرنده $learning$

که استعاره کردیم به صورت خطی دارد. ما را حساب سازی کرده است پس خطای تجربی آن صفر است که حداقل مقدار ممکن است پس ~~الگوریتم~~ الگوریتم را می‌توان خوبی ERM در نظر گرفت و $bound$ های بدست آمده را اعمال کرد:

$$m \geq \frac{\log \frac{1}{\epsilon} + VCdim(f)}{\epsilon} \cdot C$$

$$VCdim(f) = |X_d| + 1 = O(K^d) + 1 \leq 256^{1000} + 1$$

$$\Rightarrow m \geq \frac{\log 10^2 + 256^{1000} + 1}{10^{-2}} \cdot C \Rightarrow \boxed{m \geq 10^2 (256^{1000} + 6)}$$

$$h_v(x) = \begin{cases} 1 & v \text{ is a substring of } x. \\ 0 & \text{otherwise} \end{cases} \quad (ک)$$

نویس می‌کنیم برای $D \sim (x, y)$ بخواهیم ویدئوی را تشخیص دهیم. حال چون برای هر v h_v ای در f وجود دارد که به درستی ویدئوس در هر فایل تشخیص می‌دهد پس برای هر x از D نیز، فرضیه‌ای در f داریم که بدون خطا وجود ویدئوس را تشخیص می‌دهد پس داریم که $L_D(h_v) = 0$ پس شرط $Realizability$ برقرار است.

حال چون f محدود است $Uniform$ داریم پس اگر الگوریتم یادگیری خود را روش $Convergence$

ERM قرار دهیم، فرضیه PAC خواهیم داشت و از آن $learning$

رابطه زیر برقرار نیست حتی اگر:

$$|f| = |x_d| \stackrel{\text{تابع نسبت}}{=} O(K^d) \Rightarrow |f| \leq K^d$$

الف

حال با اعداد حلقه 1-8 داریم:

$$m_{ff}(\epsilon, \delta) \leq \frac{\log \frac{|f|}{\delta}}{\epsilon} \leq \frac{\log \frac{K^d}{\delta}}{\epsilon}$$

$$\Rightarrow \left\lceil m_{ff}(\epsilon, \delta) \leq \frac{\log \frac{K^d}{\delta}}{\epsilon} \right\rceil$$

اگر داشته باشیم: $\epsilon, \delta = 10^{-2}, 10^{-2}$, $d = 1000$, $K = 256$

$$m \geq m_{ff}(\epsilon, \delta) \Rightarrow m \geq \frac{\log \frac{256^{1000}}{10^{-2}}}{10^{-2}} = \frac{\log 10^2 + \log 256^{1000}}{10^{-2}}$$

$$\Rightarrow m \geq \frac{2 \log 10 + 1000 \log 256}{10^{-2}} \approx 56 \times 10^4$$

نتیجه برقرار است از حد پایین بروی m ، بند «ط» بسیار بزرگ است اما نسبت به Sample complexity که در قسمت «ی» برقرار آوریدیم کمبود چشمگیری داشته است.

$$L_D(h) = \mathbb{E}_x \mathbb{E}_y [(y - h(x))^2 | x] \quad (\text{الف})$$

از طرفی داریم:

$$\mathbb{E}_x \mathbb{E}_y [(y - h(x))^2 | x] = \sum_x \mathbb{E}_y [(y - h(x))^2 | x] \mathbb{P}(x) \quad (*)$$

برای کمینه کردن $L_D(h)$ ، کافی است مینیم عبارت (*) را نسبت آوریم. حال چون ما توزیع D مربوط به داده‌ها را می‌دانیم پس مقادیر $\mathbb{P}(x)$ برای x های مختلف از D را داریم پس در نتیجه برای مینیم کردن $L_D(h)$ ، عبارت $\mathbb{E}_y [(y - h(x))^2 | x]$ برای x های مختلف مینیم شود. (در واقع $L_D(h)$ به صورت ترکیب خطی از مقادیر $\mathbb{E}_y [(y - h(x))^2 | x]$ می‌شود که برای کمینه کردن $L_D(h)$ باید هر کدام از این term ها کمینه شوند)

$$\begin{aligned} L &= \mathbb{E}_y [(y - h(x))^2 | x] = \mathbb{E}_y [y^2 - 2yh(x) + h^2(x) | x] \\ &= \mathbb{E}_y [y^2 | x] - 2h(x) \mathbb{E}_y [y | x] + h^2(x) \end{aligned} \quad (\text{ب})$$

$$\Rightarrow \nabla_h L = 0 \Rightarrow \frac{\partial L}{\partial h} = 0 \Rightarrow -2 \mathbb{E}_y [y | x] + 2h(x) = 0$$

$$\Rightarrow h(x) = \mathbb{E}_y [y | x]$$

$$\Rightarrow \boxed{h_D(x) = \mathbb{E}_y [y | x]}$$

$$h_D(x) = \mathbb{E}_y[y|x] \quad (2)$$

$$\begin{aligned} \Rightarrow \mathbb{E}_y[(y - h_D(x))^2|x] &= \mathbb{E}_y[(y - \mathbb{E}_y[y|x])^2|x] \\ &= \mathbb{E}_y[y^2 - 2y\mathbb{E}_y[y|x] + (\mathbb{E}_y[y|x])^2|x] \\ &= \mathbb{E}_y[y^2|x] - 2\mathbb{E}_y[y|x]\mathbb{E}_y[y|x] + \mathbb{E}_x^2[y|x] \\ &= \mathbb{E}_y[y^2|x] - \mathbb{E}_y^2[y|x] = \text{Variance}(y|x) = \sigma_{y|x}^2 \end{aligned}$$

$$\Rightarrow L_D(h_D(x)) = \mathbb{E}_x[\sigma_{y|x}^2] = \min_{h \in \mathcal{H}} L_D(h) = \epsilon_{\text{Bayes}}$$

در بین تمام h ها بهترین

$$h_D(x) = \mathbb{E}_y[y|x] = y \times \underbrace{\mathbb{P}[y|x]}_{\text{deterministic رابط} = 1} = y = f(x)$$

$$\Rightarrow \mathbb{E}_y[(\overbrace{y}^{f(x)} - h_D(x))^2|x] = \mathbb{E}[(y - y)^2|x] = 0 = \sigma_{y|x}^2$$

$$\Rightarrow L_D(h_D) = \mathbb{E}_x[\underbrace{\sigma_{y|x}^2}_0] = 0$$

$$h_D(x) = \mathbb{E}_y[y|x] = \sum_y y \cdot \mathbb{P}[y|x] = \sum_y y \cdot \frac{\mathbb{P}(x, y)}{\mathbb{P}(x)} \quad (3)$$

همه x, y را در نظر بگیریم \rightarrow

$$h_D(x) = \sum_y y \cdot \frac{\mathbb{P}(x) \mathbb{P}(y)}{\mathbb{P}(x)} = \sum_y y \cdot \mathbb{P}(y) = \mathbb{E}_y[y]$$

$$E_y[(y - h(x))^2 / x] = E_y[(y - E[y])^2 / x]$$

$$= \frac{E_y[y^2 / x]}{E[y^2]} - 2 \frac{E[y] E[y / x]}{E[y]} + \frac{E^2[y / x]}{E^2[y]}$$

$$= E[y^2] - E^2[y] = \sigma_y^2$$

$$\Rightarrow L_D(h_D) = E_x[\sigma_y^2] = \sigma_y^2$$

$$L_D(h) = \mathbb{E}_x \mathbb{E}_y [(h(x) - y)^2 | x], \quad h(x) = c \quad (9)$$

$$\Rightarrow L_D(h) = \mathbb{E}_x \mathbb{E}_y [(c - y)^2 | x] = \mathbb{E}_x \mathbb{E}_y [y^2 - 2cy + c^2 | x]$$

$$= \mathbb{E}_x [\mathbb{E}_y [y^2 | x] - 2c \mathbb{E}_y [y | x] + c^2]$$

$$= \mathbb{E}_x \mathbb{E}_y [y^2 | x] - 2c \mathbb{E}_x \mathbb{E}_y [y | x] + c^2$$

$$\frac{dL_D}{dc} = 0 \Rightarrow -2 \mathbb{E}_x \mathbb{E}_y [y | x] + 2c = 0$$

$$\Rightarrow \boxed{c = \mathbb{E}_x \mathbb{E}_y [y | x] = \mathbb{E}_{x,y} [y]}$$

$$\begin{aligned} \rightarrow L_D\left(\frac{c}{\mathbb{E}_{x,y} [y]}\right) &= \mathbb{E}_x \mathbb{E}_y \left[\left(\mathbb{E}_x \mathbb{E}_y [y | x] - y \right)^2 | x \right] \\ &= \mathbb{E}_x \mathbb{E}_y [y^2 | x] - \left(\mathbb{E}_x \mathbb{E}_y [y | x] \right)^2 \quad \left| \begin{array}{l} \text{مقدار خطای اِزا} \\ \text{c بهینه} \end{array} \right. \end{aligned}$$

$$\mathbb{E}_x [\sigma_{y|x}^2] \leq \sigma_y^2 \quad (10) \quad \text{رابطه صحیح به صورت مابین است:}$$

وقتی حالت تساوی رخ می‌دهد، y و x از هم مستقل باشند اما وقتی رابطه استقلال ندارند بدین معناست که به ازای هر D به x ، به یک توزیع وابسته به x ، برای مقادیر y آن می‌رسیم یعنی $\mathbb{P}(Y/x)$ و سپس از روی این توزیع وابسته به x ، تخمین خود را که برابر

است با $\mathbb{E}_y [y | x]$ بدست می‌آوریم. اما وقتی استقلال برقرار باشد، بهر هر x ، توزیع y

ها وابسته به آن را نداریم و باید از روی خود توزیع y تخمین را بدست آوریم که برابر با $\mathbb{E}_y [y]$ است. به نظر می‌رسد که در نتیجه خطای بیشتری خواهیم داشت.

$$\min_{h \in H} L_D(h) \geq \epsilon_{\text{bayes}} = \mathbb{E}_x [\sigma_{y|x}^2] \quad (5)$$

بله، زیرا مقدار ϵ_{bayes} کمترین خطای است که به ازای بعضی h های ممکن درست می‌باشد.

که شامل مجموعه فرضیه h نیز می‌شود، پس کمترین انتخاب ممکن از بین H



حداقل کمترین خطای برابر با ϵ_{bayes} خواهد داشت.

(6) اگر شرط Realizability بر روی h داشته باشیم بین معنایست که کمترین خطای برابر با

h برابر با صفر است:

$$\underbrace{\min_{h \in H} L_D(h)}_0 \geq \epsilon_{\text{bayes}} \geq 0 \Rightarrow \boxed{\epsilon_{\text{bayes}} = 0}$$

(5) وقتی شرط Realizability بر روی H برقرار است بین معنایست h ای در H وجود

دارد که رابطه قطعی بین x و y در D برقرار می‌کند: $y = f(x)$

$$E_y[y|x] = y \iff \mathbb{P}(y|x) = 1$$

که با توجه به آن

$$E_x \left[\underbrace{E_y \left(\underbrace{E_y[y|x]}_y - y \right)^2}_0 \right] = E_x[0] = 0$$

$$\Rightarrow \epsilon_{\text{bayes}} = 0$$

در واقع با شرط Realizability حالت احتمالاتی معادله y برای x از بین می‌رود و برای هر

x یک y خواهیم داشت که با انتخاب این، خطای روی D صفر می‌شود.

(ک) ابتدا: $(\epsilon_{app} - \epsilon_{bayes} \geq 0)$

ϵ_{bayes} خطای است که می‌توان از بین تمام h های ممکن بدست آورد که مجموع فرضیه H کمترین

تیر عضو آن است. حال خطای بدست آمده برای بهترین h از H ، کوچکتر یا مساوی با

$$\epsilon_{app} = \min_{h \in H} L_D(h) \geq \epsilon_{bayes} = \min_{h: \text{تمام } h \text{ های ممکن}} L_D(h)$$

$$\Rightarrow (\epsilon_{app} - \epsilon_{bayes} \geq 0)$$

• ϵ_{bayes} : بدین معناست که بهترین خطای h می‌توان h ها می‌توان برای مدل بدست آورد چقدر است. (\hat{h})

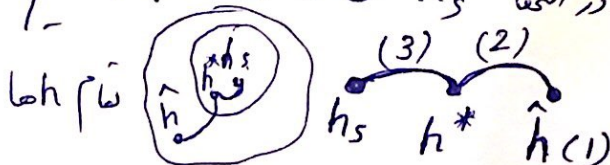
• $\epsilon_{bayes} - \epsilon_{app}$: اگر ما بهترین انتخاب h را از بین H انجام دهیم - چقدر با بهترین h ممکن برای مدل فاصله داریم.

• ϵ_{est} : h ای که بر اساس آموزش یا آزمون داده آموزشی S بدست می‌آوریم، مقدار خطای آن چقدر با مقدار خطای بهترین h ای که می‌توانیم از H انتخاب کنیم فاصله دارد.

• جمله اول (ϵ_{bayes}) در این جمع بیان می‌کند که حداقل چقدر خطای هر h ممکن

داریم (شامل تمام h ها از جمله h_s)، حال بهترین h از H چقدر خطا نسبت

به بهترین h ممکن دارد (جمله دوم) و در انتها h_s ای که انتخاب کرده‌ایم چقدر خطا نسبت به بهترین h در H دارد.



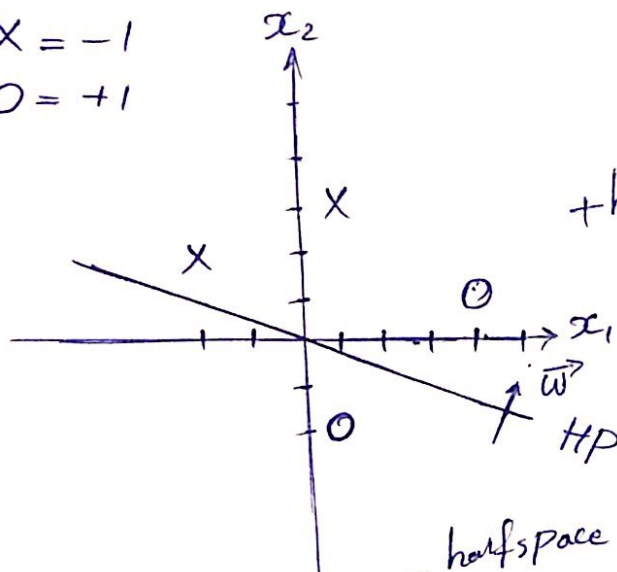
$$X = \mathbb{R}^2, Y = \{-1, +1\}.$$

① T23 مس

$$S = \{((-2, 2), -1), ((1, 3), -1), ((1, -2), +1), ((4, 1), +1)\}.$$

$$X = -1$$

$$O = +1$$



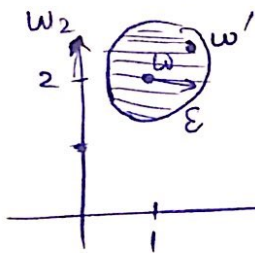
+ halfspace

$$HP: \langle (1, 2), (x_1, x_2) \rangle = 0$$

$$h = \text{sign}[x_1 + 2x_2]$$

$$\text{- halfspace} \Rightarrow w = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, b = 0$$

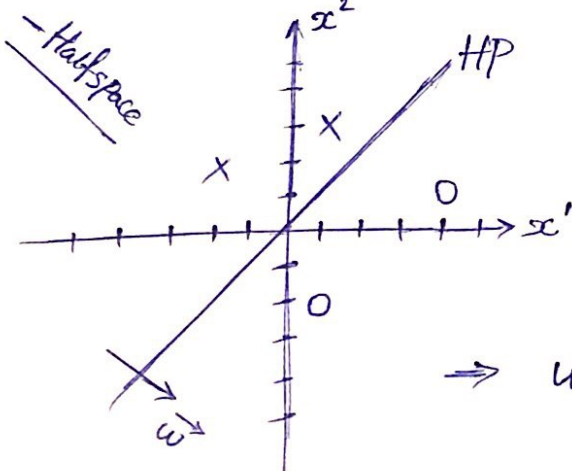
$$\Rightarrow L_S(w) = \frac{1}{4} \sum_{i=1}^4 \ell(w, (x_i, y_i)) = \frac{1}{4} [+1 + 1 + 0 + 1] = \frac{3}{4}$$



$$\forall w': |w' - w| < \epsilon, \text{ for some } \epsilon > 0 :$$

$$L_S(w') = L_S(w) = \frac{3}{4}$$

$\Rightarrow w$ is local minimum of $L_S(w)$.



$$HP: \langle (1, -1), (x_1, x_2) \rangle = 0$$

$$h = \text{sign}[x_1 - x_2]$$

$$\Rightarrow w^* = \begin{bmatrix} 1 \\ -1 \end{bmatrix} \Rightarrow \boxed{L_S(w^*) = 0}$$

\Rightarrow چون خطای کمتر از 0 نداریم پس w^* global minimum است و $L_S(w)$ global min است.

②