

به نام خدا

پاسخ تمرین پنجم یادگیری ماشین

امین اسدی

۸۱۰۱۹۶۴۱۰

پاسخ سوالات ۱ تا ۶

سوال اول

$$\alpha = \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix}_{3 \times 1} \quad X = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ -1 & 0 \end{pmatrix}_{3 \times 2} \quad Y = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & -1 \end{pmatrix}_{3 \times 3}$$

$$\hat{\alpha} = \underset{\alpha}{\operatorname{argmax}} \alpha^T \mathbf{1}_n - \frac{1}{2} \alpha^T Y G Y \alpha$$

$$G = X \cdot X^T = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ -1 & 0 \end{pmatrix} \begin{pmatrix} 1 & 0 & -1 \\ 0 & 1 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 & -1 \\ 0 & 1 & 0 \\ -1 & 0 & 1 \end{pmatrix}$$

$$Y G Y = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{pmatrix} \rightarrow \hat{\alpha} = \underset{\alpha = \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix}}{\operatorname{argmax}} \sum_{i=1}^3 \alpha_i - \frac{1}{2} \left(\alpha_1, \alpha_2, \alpha_3 \right) \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix}$$

$$\rightarrow \hat{\alpha} = \underset{\alpha = \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix}}{\operatorname{argmax}} \left[\alpha_1 + \alpha_2 + \alpha_3 - \frac{1}{2} \left(\alpha_1(\alpha_1 + \alpha_3) + \alpha_2^2 + \alpha_3(\alpha_1 + \alpha_3) \right) \right]$$

$$\text{subject to } \sum_{i=1}^3 \alpha_i y_i = 0 \rightarrow \boxed{\alpha_1 + \alpha_2 = \alpha_3}$$

$$\hat{\alpha} = \underset{\alpha}{\operatorname{argmax}} \left[2(\alpha_1 + \alpha_2) - \frac{1}{2} \left(\alpha_1(2\alpha_1 + \alpha_2) + \alpha_2^2 + (\alpha_1 + \alpha_2)(2\alpha_1 + \alpha_2) \right) \right]$$

$$= \underset{\alpha}{\operatorname{argmax}} \left[2\alpha_1 + 2\alpha_2 - \frac{1}{2} \left(2\alpha_1^2 + \alpha_1\alpha_2 + \alpha_2^2 + 2\alpha_1^2 + \alpha_1\alpha_2 + 2\alpha_1\alpha_2 + \alpha_2^2 \right) \right]$$

$$= \underset{\alpha}{\operatorname{argmax}} \left[2\alpha_1 + 2\alpha_2 - \underbrace{2\alpha_1^2 - \alpha_2^2 - 2\alpha_1\alpha_2}_L \right]$$

نقطه بحرانی

$$\begin{cases} \frac{\partial L}{\partial \alpha_1} = 0 \rightarrow 2 - 4\alpha_1 - 2\alpha_2 = 0 \\ \frac{\partial L}{\partial \alpha_2} = 0 \rightarrow 2 - 2\alpha_1 - 2\alpha_2 = 0 \end{cases} \rightarrow \begin{matrix} \alpha_1 = 0 \\ \alpha_2 = 1 \\ \alpha_3 = 0 + 1 = 1 \end{matrix}$$

$$\frac{\partial^2 L}{\partial \alpha_1 \partial \alpha_1} \times \frac{\partial^2 L}{\partial \alpha_2 \partial \alpha_2} - \left(\frac{\partial^2 L}{\partial \alpha_1 \partial \alpha_2} \right)^2 > 0, \quad \frac{\partial^2 L}{\partial \alpha_1 \partial \alpha_1} < 0 \rightarrow \text{maximum}$$

$$\rightarrow \alpha = \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}$$

$$w = \sum_{i=1}^3 \alpha_i y_i x_i = 0(1, 0) + 1(0, 1) + 1(1, 0) = (1, 1)$$

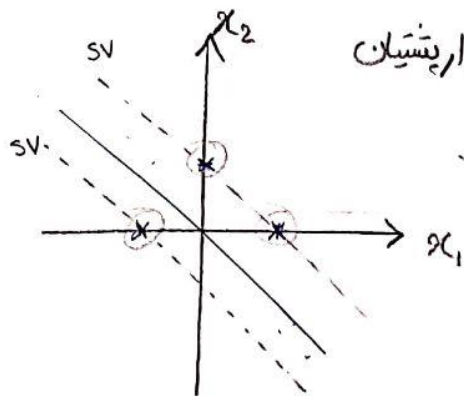
$$b = - \frac{\min_{i: y^{(i)}=1} \omega^T x^{(i)} + \max_{i: y^{(i)}=-1} \omega^T x^{(i)}}{2} \quad \downarrow \quad w = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$= - \frac{\min \left([1, 1] \begin{bmatrix} 1 \\ 0 \end{bmatrix}, [1, 1] \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right) + [1, 1] \begin{bmatrix} -1 \\ 0 \end{bmatrix}}{2}$$

$$= - \frac{1 - 1}{2} = 0$$

$$\rightarrow \text{معادله خط جداساز: } \omega_1' x_1 + \omega_2' x_2 + b' = 0 \rightarrow \boxed{x_2 = -x_1}$$

$$\text{خط های خاصه: } \begin{cases} \omega_1' x_1 + \omega_2' x_2 + b' = 1 \rightarrow \boxed{x_2 = -x_1 + 1} \\ \omega_1' x_1 + \omega_2' x_2 + b' = -1 \rightarrow \boxed{x_2 = -x_1 - 1} \end{cases}$$



در نتیجه هر سه نقطه بردار پشتیبان (SV) می باشند.

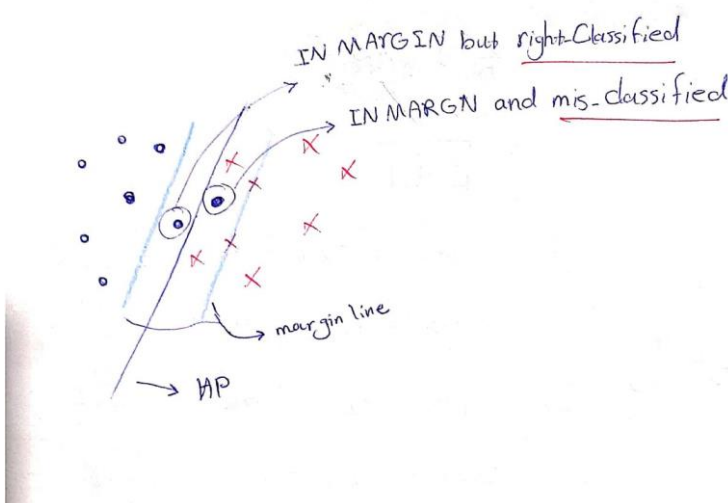
سوال دوم)

(۱)

مساله‌ی ثانویه به ما این امکان را می‌دهد که مفهوم kernel را وارد مساله کنیم تا برای حل مسائلی که به طور خطی جداپذیر نیستند (یا حتی در مسائل جداپذیر خطی برای رفع مشکل واریانس زیاد) از kernel استفاده کنیم. همچنین مساله‌ی ثانویه به ما امکان بهینه انجام دادن محاسبات لازم را می‌دهد زیرا ابتدا آلفاها را می‌یابیم که به جز برای SV ها، برای بقیه نقاط برابر 0 هستند و ترم‌های مربوط به آن‌ها از محاسبات حذف می‌شوند و کافیت برای SV ها محاسبات انجام شوند که به دلیل اینکه تعداد SV ها معمولاً می‌تواند بسیار محدود باشد، محاسبات بسیار ساده تر خواهند شد.

(۲)

از hard margin برای حل مسائلی که کار می‌رود که به صورت خطی جداپذیر هستند که منطقاً در آن‌ها هیچ نقطه‌ای mis-classify نمی‌شود. در دو حالت از soft-margin ها استفاده می‌شود: ۱- در مسائلی که مساله به صورت خطی جداپذیر نیست و بنابراین انعطاف بیشتری به margin می‌دهیم و در نتیجه تعدادی از نقاط به اشتباه طبقه‌بندی خواهند شد. ۲- در حالتی که مساله به صورت خطی جداپذیر است اما حاضر هستیم به خاطر $\text{bias-variance tradeoff}$ تعدادی کمی از نقاط به اشتباه طبقه‌بندی شوند اما در عوض generalization طبقه‌بند را بالا ببریم.



(۳)

این ضریب کنترل‌کننده‌ی trade-off بین خطا و حاشیه است. وقتی این ضریب افزایش می‌یابد در واقع به خطا بهای بیشتری می‌دهیم و در واقع می‌خواهیم خطا را کمتر کنیم بنابراین در این حالت خطا را کمینه می‌کنیم اما حاشیه ممکن است کوچک باشد. اگر این ضریب را کم کنیم به بهای mis-classify شدن برخی نقاط، حاشیه را بیشینه می‌کنیم.

سوال سوم)

کرنل‌ها توابعی هستند که برای تبدیل داده‌ها از یک فضای ویژگی به یک فضای ویژگی با بُعد بالاتر استفاده می‌شوند. علت بردن داده‌ها به فضای بُعد بالاتر این است که ممکن است در بُعد پایین داده‌ها به صورت خطی جداپذیر نباشند و در نتیجه بُعد فضای را افزایش می‌دهیم تا داده‌ها **unfold** شوند و تا جای ممکن به صورت خطی جداپذیر شوند. یک روش این است که از هر نقطه‌ها به فضای جدید منتقل کنیم و سپس محاسبات را در فضای جدید انجام دهیم، اما این روش خصوصا وقتی ابعاد فضای جدید خیلی زیاد است، بسیار هزینه‌بر یا حتی غیر ممکن خواهد بود. اینجاست که از کرنل‌ها استفاده می‌کنیم. با استفاده از کرنل‌ها می‌توانیم به جای انجام محاسبات در بُعد بالا، محاسبات را در بُعد پایین (اولیه) با هزینه کم انجام دهیم و به همان نتیجه دست پیدا کنیم.

اثبات رابطه کاشی-شوارتز:

می‌دانیم کرنل-مارتیس (Kernel Matrix) متناظر با هر زیر مجموعه ناتهی از نقاط باید Positive Semi-Definite باشد. در نتیجه کرنل-مارتیس مربوط به دو نقطه x و y که به صورت زیر است باید PSD باشد:

$$\begin{bmatrix} K(x, x) & K(x, y) \\ K(x, y) & K(y, y) \end{bmatrix}$$

طبق تعریف PSD، باید همه مقادیر ویژه این ماتریس نامنفی باشند و در نتیجه دترمینان آن نیز باید نامنفی باشد پس داریم:

$$K(x, x) K(x, x) - K(x, x)^2 \geq 0 \Rightarrow K(x, x)^2 \leq K(x, x) K(x, x)$$

فرض کنید که $x_{m,i}$ نشان‌دهنده مولفه i ام داده m ام باشد. در این صورت طبق تعریف کرنل داریم:

$$\begin{aligned} \sum_{m=1}^Q \sum_{n=1}^Q K_{\phi}(x_m, x_n) &= \sum_{m=1}^Q \sum_{n=1}^Q \langle \phi(x_m), \phi(x_n) \rangle \\ &= \sum_{m=1}^Q \sum_{n=1}^Q \left\langle \phi \begin{pmatrix} x_{m,0} \\ \vdots \\ x_{m,d_1} \end{pmatrix}, \phi \begin{pmatrix} x_{n,0} \\ \vdots \\ x_{n,d_1} \end{pmatrix} \right\rangle \\ &= \sum_{m=1}^Q \sum_{n=1}^Q \left\langle \begin{pmatrix} x'_{m,0} \\ \vdots \\ x'_{m,d_2} \end{pmatrix}, \begin{pmatrix} x'_{n,0} \\ \vdots \\ x'_{n,d_2} \end{pmatrix} \right\rangle \\ &= [x'_{1,0}{}^2 + x'_{2,0}{}^2 + \dots + x'_{Q,0}{}^2 + 2(x'_{1,0}x'_{2,0} + \dots + x'_{Q-1,0}x'_{Q,0})] + \dots + [x'_{1,d_2}{}^2 \\ &\quad + x'_{2,d_2}{}^2 + \dots + x'_{Q,d_2}{}^2 + 2(x'_{1,d_2}x'_{2,d_2} + \dots + x'_{Q-1,d_2}x'_{Q,d_2})] \end{aligned}$$

$$\begin{aligned}
&= \left(x'_{1,0} + \cdots + x'_{Q,0}\right)^2 + \cdots + \left(x'_{1,d_2} + \cdots + x'_{Q,d_2}\right)^2 \\
\rightarrow \frac{1}{Q} \sqrt{\sum_{m=1}^Q \sum_{n=1}^Q K_{\emptyset}(x_m, x_n)} &= \sqrt{\frac{1}{Q^2} \sum_{m=1}^Q \sum_{n=1}^Q K_{\emptyset}(x_m, x_n)} \\
&= \sqrt{\left(\frac{x'_{1,0} + \cdots + x'_{Q,0}}{Q}\right)^2 + \cdots + \left(\frac{x'_{1,d_2} + \cdots + x'_{Q,d_2}}{Q}\right)^2} = \|\mu_{\emptyset}\|
\end{aligned}$$

و اثبات به پایان می‌رسد.

$$1) K(x, y) = f(x) K_1(x, y) f(y)$$

$$\forall c \in \mathbb{R}^n :$$

$$\begin{aligned} c^T K(x, y) c &= \sum_{i,j} c_i K_1(x_i, y_j) c_j = \sum_{i,j} \underbrace{c_i f(x_i)}_{a_i} K_1(x_i, y_j) \underbrace{f(y_j) c_j}_{a_j} \\ &= \sum_{i,j} a_i K_1(x_i, y_j) a_j = \underbrace{A^T K_1 A}_{\geq 0} \\ &\rightarrow K_1(x, y) \geq 0 \rightarrow \text{Valid} \end{aligned}$$

$$3) K(x, y) = K_1(x, y) + K_2(x, y)$$

$$\forall c \in \mathbb{R}^n :$$

$$\begin{aligned} c^T K(x, y) c &= \sum_{i,j} c_i (K_1(x_i, y_j) + K_2(x_i, y_j)) c_j \\ &= \sum_{i,j} c_i K_1(x_i, y_j) c_j + \sum_{i,j} c_i K_2(x_i, y_j) c_j \\ &= \underbrace{c^T K_1 c}_{\geq 0} + \underbrace{c^T K_2 c}_{\geq 0} \geq 0 \rightarrow \text{Valid} \end{aligned}$$

$$4) K(x, y) = K_1(x, y) K_2(x, y)$$

$$\forall c \in \mathbb{R}^n :$$

$$\begin{aligned} c^T K(x, y) c &= \sum_{i,j} c_i K_1(x_i, y_j) K_2(x_i, y_j) c_j \\ &= \sum_{i,j} K_1(x_i, y_j) \underbrace{c_i K_2(x_i, y_j) c_j}_{\text{طبق بخش قبل} \geq 0} \\ &= \sum_{i,j} K_1(x_i, y_j) K_2'(x_i, y_j) \\ &= \text{trace}(K_1 K_2'^T) \geq 0 \rightarrow \text{Valid} \end{aligned}$$

له حاصل جمع مقادیر ویژه ≥ 0 است

سوال پنجم)

می‌دانیم فاصله دو نقطه \underline{u} و \underline{v} را می‌توان بر حسب ضرب داخلی به صورت زیر بیان کرد:

$$d^2(\underline{u}, \underline{v}) = \|\underline{u} - \underline{v}\|^2 = \langle \underline{u} - \underline{v}, \underline{u} - \underline{v} \rangle = \langle \underline{u}, \underline{u} \rangle + \langle \underline{v}, \underline{v} \rangle - 2\langle \underline{u}, \underline{v} \rangle$$

و نیز طبق تعریف کرنل داریم:

$$K_{\emptyset}(x_1, x_2) = \langle \emptyset(x_1), \emptyset(x_2) \rangle$$

$$\langle \emptyset(x_1), \emptyset(x_2) \rangle = K(x_1, x_2) = K([1, 1]^T, [3, 4]^T) = \exp\left(-\frac{(3-1)^2 + (4-1)^2}{2}\right) = \exp(-6.5)$$

$$\begin{aligned} \langle \emptyset(x_1), \emptyset(x_1) \rangle &= K(x_1, x_1) = K([1, 1]^T, [1, 1]^T) = \exp\left(-\frac{(1-1)^2 + (1-1)^2}{2}\right) \\ &= \exp(0) = 1 \end{aligned}$$

$$\langle \emptyset(x_2), \emptyset(x_2) \rangle = K(x_2, x_2) = K([3, 4]^T, [3, 4]^T) = \exp\left(-\frac{(3-3)^2 + (4-4)^2}{2}\right) = \exp(0) = 1$$

$$\rightarrow d^2(\emptyset(x_1), \emptyset(x_2)) = 1 + 1 - 2 \exp(-6.5) \cong 2 - 2 \times 0.001 = 1.998 \rightarrow$$

$$d = \sqrt{1.998}$$

قسمت دوم سوال ۵:

کافیست تعداد تک‌جمله‌ای‌های حاصل کرنل را بدست آوریم. با توجه به اینکه تعداد ابعاد اولیه برابر d است می‌توان نوشت:

$$x = [x_1, x_2, \dots, x_d]^T$$

$$y = [y_1, y_2, \dots, y_d]^T$$

$$K(x, y) = (x^T y + 1)^2 = (x_1 y_1 + x_2 y_2 + \dots + x_d y_d + 1)^2$$

به وضوح حاصل کرنل برابر مجموع تک‌جمله‌ای‌هایی به فرم $x_1^{c_1} y_1^{c_1} \dots x_d^{c_d} y_d^{c_d} 1^{c_{d+1}}$ خواهد بود که:

$$c_1 + c_2 + \dots + c_d + c_{d+1} = 2 \quad 0 \leq c_1, c_2, \dots, c_d, c_{d+1}$$

پس تعداد این تک‌جمله‌ای‌ها برابر خواهد بود با تعداد جواب‌های معادله بالا که می‌دانیم برابر است با:

$$\binom{d+1+2-1}{2} = \binom{d+2}{2} = \frac{(d+1)(d+2)}{2}$$

سوال ششم)

۱) در رویکرد **generative** ابتدا توزیع **joint** متغیرها بدست آمده و سپس از روی آن توزیع‌های شرطی محاسبه می‌شوند (در واقع تابع **discriminant** را با استفاده از تخمین توزیع‌های مربوط به هر کلاس تخمین می‌زنند) اما در رویکرد **discriminative** به صورت مستقیم تابع توزیع‌های شرطی محاسبه می‌شوند (در واقع تابع **discriminant** را به طور مستقیم تخمین می‌زنند) رویکرد **generative** نیاز به داده‌های زیادی دارد بنابراین رویکرد **discriminative** عملی تر است.

۲) مزیت رویکرد **one-vs-rest** این است که تعداد طبقه‌بند هایی که استفاده می‌کنیم برابر تعداد کلاس‌ها (یا یکی کمتر از آن‌ها) است و در نتیجه با افزایش تعداد کلاس‌ها، تعداد طبقه‌بندهای مورد نیاز به صورت خطی افزایش می‌یابد و بنابراین هزینه کمتری دارد. اگر در هر مرحله یکی از کلاس‌ها را کنار بگذاریم ایراد این روش این است نتیجه به ترتیب کلاس‌ها حساس خواهد بود. اگر در مرحله کلاس را کنار نگذاریم ایراد این خواهد بود که یا بعضی از نواحی به بیش از یک کلاس دسته‌بندی خواهند شد و یا بعضی نواحی به هیچ کلاسی دسته‌بندی نخواهند شد (ambiguous area). رویکرد **one-vs-another** معایب ذکر شده را ندارد اما یک عیب آن این است که برای هر دو کلاس یک طبقه بند نیاز دارد (در کل $\binom{C}{2}$ تا طبقه بند برای C کلاس). معمولاً روش **one-vs-another** عملکرد بهتری نسبت به **one-vs-rest** دارد. در ماشین خطی نیز از همان اول مساله را چند-کلاسه حل می‌کنیم و در نهایت **max** می‌گیریم که این باعث می‌شود مشکلات بالا پیش نیاید اما در حالت کلی دو روش اول نتایج بهتری دارند.

۳)

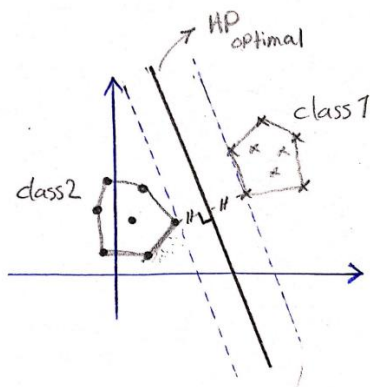
کاهش بُعد یا **dimensionality reduction** با هدف جلوگیری از نحسی ابعاد (curse of dimensionality) انجام می‌شود. نحسی ابعاد وقتی اتفاق می‌افتد که تعداد ویژگی‌ها در مقایسه با تعداد داده‌ها زیاد باشد که باعث **sparse** شدن داده‌ها و در نهایت **overfitting** و بدتر شدن عملکرد طبقه‌بند می‌شود. افزایش بُعد با این هدف انجام می‌شود تا بتوان داده‌هایی که داده‌هایی که در بعد پایین به صورت خطی جداپذیر نیستند به ابعاد بالاتر برد تا به صورت خطی از هم جداپذیر باشند.



(از اسلایدها)

(۴)

در مساله طبقه‌بندی دو کلاسه، ابر صفحه optimal (بیشترین حاشیه) بر کوتاه‌ترین خطی که convex hull مجموعه داده‌های دو کلاس را به هم وصل میکند عمود است و آن را از وسط به دو نصف تقسیم می‌کند.



(۵)

اولیه: minimize $\frac{1}{2} \|w\|^2$ subject to $y_i(w^T x_i + b) \geq 1$

maximize $L(w, b, \alpha) = \frac{1}{2} w^T w + \sum_{i=1}^n \alpha_i (1 - y_i (w^T x_i + b))$

$$\begin{cases} \nabla_w L = 0 \rightarrow w + \sum_{i=1}^n \alpha_i (-y_i) x_i = 0 \rightarrow w = \sum_{i=1}^n \alpha_i y_i x_i \\ \nabla_b L = 0 \rightarrow \sum_{i=1}^n \alpha_i y_i = 0 \end{cases}$$

جاگذاری در L

$$\begin{aligned} \rightarrow L &= \frac{1}{2} \sum_{i=1}^n \alpha_i y_i x_i^T \sum_{j=1}^n \alpha_j y_j x_j + \sum_{i=1}^n \alpha_i (1 - y_i (\sum_{j=1}^n \alpha_j y_j x_j^T x_i + b)) \\ &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i \sum_{j=1}^n \alpha_j y_j x_j^T x_i \\ &\quad - b \sum_{i=1}^n \alpha_i y_i = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j \end{aligned}$$

مساله ثانویه $\rightarrow \begin{cases} \text{maximize } W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j \\ \text{subject to: } \alpha_i \geq 0, \sum_{i=1}^n \alpha_i y_i = 0 \end{cases}$

بله- در **hard margin** بردارهای پشتیبان نقاطی هستند که که نامساوی به حالت تساوی برقرار خواهد بود و در واقع ضریب آلفا متناظر با آنها غیر صفر است. این نقاط در واقع نقاطی هستند که روی خط های حاشیه (margin lines) قرار می گیرند. اما در حالت **soft-margin** علاوه بر این نقاط، نقاطی که در داخل ناحیه **margin** قرار می گیرند و نیز نقاطی که **mis-classify** می شوند نیز بردار پشتیبان هستند.

پایان