

1.1  $g$  is a valid kernel

$$h(n, y) = \frac{1}{4} \left( g(n+y, n+y) - g(n-y, n-y) \right) \stackrel{(*)}{\stackrel{(**)}}{=} \frac{1}{4} (4g(n, y)) = g(n, y)$$

$$g(n+y, n+y) = \underbrace{g(n, n+y)}_{g(n, n) + g(y, n)} + \underbrace{g(y, n+y)}_{g(n, y) + g(y, y)} = g(n, n) + g(y, y) + 2g(n, y) \quad (*)$$

$$g(n-y, n-y) = \underbrace{g(n, n-y)}_{g(n, n) - g(y, n)} - \underbrace{g(y, n-y)}_{g(n, y) - g(y, y)} = g(n, n) + g(y, y) - 2g(n, y) \quad (**)$$

بنابراین نشان دادیم:  $h(n, y) = g(n, y)$  و  $h$  هم یک کرنل معتبر است

1.2.a

یک راه برای اثبات معتبر بودن  $k_3$  آن است که نشان دهیم  $\Phi_3$  ای وجود دارد که  $k_3(n_1, n_2) = \langle \Phi_3(n_1), \Phi_3(n_2) \rangle$

$$k_1 \text{ is valid} \rightarrow k_1(n_1, n_2) = \langle \Phi_1(n_1), \Phi_1(n_2) \rangle = \begin{bmatrix} \varphi_1(n_1) \\ \vdots \\ \varphi_m(n_1) \end{bmatrix} \cdot \begin{bmatrix} \varphi_1(n_2) \\ \vdots \\ \varphi_m(n_2) \end{bmatrix} = \sum_{i=1}^m \varphi_i(n_1) \varphi_i(n_2)$$

$$k_2 \text{ is valid} \rightarrow k_2(n_1, n_2) = \langle \Phi_2(n_1), \Phi_2(n_2) \rangle = \begin{bmatrix} \varphi'_1(n_1) \\ \vdots \\ \varphi'_{m'}(n_1) \end{bmatrix} \cdot \begin{bmatrix} \varphi'_1(n_2) \\ \vdots \\ \varphi'_{m'}(n_2) \end{bmatrix} = \sum_{i=1}^{m'} \varphi'_i(n_1) \varphi'_i(n_2)$$

$$\Phi_3(n) = \begin{bmatrix} \varphi_1(n) \\ \vdots \\ \varphi_m(n) \\ \varphi'_1(n) \\ \vdots \\ \varphi'_{m'}(n) \end{bmatrix} \rightarrow k_3(n_1, n_2) = \sum_{i=1}^m \varphi_i(n_1) \varphi_i(n_2) + \sum_{i=1}^{m'} \varphi'_i(n_1) \varphi'_i(n_2) = k_1(n_1, n_2) + k_2(n_1, n_2)$$

حال  $\Phi_3$  را به صورت زیر تعریف می‌کنیم

بنابراین  $k_3$  معتبر است

1.2.b

مانند سرنال قبل نشان می دهیم  $\Phi_4$  ای وجود دارد که

$$k_1 \text{ is valid} \rightarrow k_1(n_1, n_2) = \dots = \sum_{i=1}^m \varphi_i(n_1) \varphi_i(n_2)$$

$$k_2 \text{ is valid} \rightarrow k_2(n_1, n_2) = \dots = \sum_{i=1}^{m'} \varphi'_i(n_1) \varphi'_i(n_2)$$

$$k_4(n_1, n_2) = k_1(n_1, n_2) k_2(n_1, n_2) = \sum_{i=1}^m \sum_{j=1}^{m'} \underbrace{\varphi_i(n_1)} \underbrace{\varphi_i(n_2)} \underbrace{\varphi'_j(n_1)} \underbrace{\varphi'_j(n_2)}$$

$$\Phi_4(n) = \begin{bmatrix} \varphi''_{11}(n) \\ \vdots \\ \varphi''_{mm'}(n) \end{bmatrix} \quad \varphi''_{ij}(n) = \varphi_i(n) \varphi'_j(n)$$

حال  $\Phi_4$  را به صورت مقابل تعریف می کنیم

$$\rightarrow k_4(n_1, n_2) = \sum_{i=1}^m \sum_{j=1}^{m'} \varphi''_{ij}(n_1) \varphi''_{ij}(n_2) = \langle \Phi_4(n_1), \Phi_4(n_2) \rangle$$

بنابراین  $k_4$  معتبر است

1.2.c

$$e^{k_1(n_1, n_2)} = \lim_{N \rightarrow \infty} \sum_{i=0}^N \frac{k_1(n_1, n_2)^i}{i!} = \sum_{i=0}^{\infty} \frac{k_1(n_1, n_2)^i}{i!}$$

اگر بسط تیلور  $e^{k_1(n_1, n_2)}$  را در نقطه 0 بنویسیم داریم

(\*) طبق بخش 1.2.b می دانیم که ضرب 2 کرنل معتبر همچنان معتبر است پس این کار را می توان به تعداد دلخواه انجام داد ، یعنی لام استقرآ آن است که

$k^n \times k$  خود عامل ضرب 2 کرنل معتبر است پس  $k_1(n_1, n_2)^i$  کرنل معتبر است

(\*) همچنین ضرب عدد مثبت در کرنل معتبر نیز یک کرنل معتبر است زیرا :

$$k \text{ is valid} \rightarrow k(n_1, n_2) = \sum_{i=1}^m \varphi_i(n_1) \varphi_i(n_2) \rightarrow \forall \lambda > 0 \quad \lambda k(n_1, n_2) = \sum_{i=1}^m \frac{\sqrt{\lambda} \varphi_i(n_1)}{\varphi'_i(n_1)} \frac{\sqrt{\lambda} \varphi_i(n_2)}{\varphi_i(n_2)}$$

(\*) طبق بخش 1.2.a جمع کرنل های معتبر نیز معتبر است

$$\rightarrow e^{k_1(n_1, n_2)} \text{ is valid}$$

1.2.d

مسابه سوال قبل بسط نیاید.  $\frac{1}{1-x_1^T x_2}$  را در نقطه 0 می نویسیم

$$\frac{1}{1-x_1^T x_2} = \lim_{N \rightarrow \infty} \sum_{i=0}^N (x_1^T x_2)^i = \lim_{N \rightarrow \infty} \frac{1 - (x_1^T x_2)^{N+1}}{1 - x_1^T x_2}$$

یعنی اگر خواصیم به  $\frac{1}{1-x_1^T x_2}$  برسیم باید شرط

تصادف هندی را داشته باشیم.

if  $|x_1^T x_2| < 1 \rightarrow \frac{1}{1-x_1^T x_2} = \sum_{i=0}^{\infty} (x_1^T x_2)^i \rightarrow k_\phi$  is valid

اگر این شرط را نداشته باشیم لزوماً  $k_\phi$  معتبر نیست. به عنوان مثال، مثال نقض زیر را در نظر بگیرید.

2 samples in  $\mathbb{R}^3$   $x_1 = \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix}$ ,  $x_2 = \begin{bmatrix} 2 \\ 4 \\ 6 \end{bmatrix}$

$\rightarrow$  kernel matrix =  $\begin{bmatrix} \frac{1}{1-35} & \frac{1}{1-49} \\ \frac{1}{1-49} & \frac{1}{1-56} \end{bmatrix} \rightarrow \exists x = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} -\frac{1}{34} & -\frac{1}{43} \\ -\frac{1}{43} & -\frac{1}{56} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = -\frac{1}{34} - \frac{2}{43} - \frac{1}{56} < 0$

$k_\phi$  is not valid.  $\leftarrow$  kernel matrix is not positive semi-definite

1.3

$k$  is valid  $\rightarrow k(x, x') = \phi(x) \cdot \phi(x')$

$$\hat{k}(x^{(i)}, y^{(i)}) = \sum_{x^{(j)} \in A} \sum_{y^{(j)} \in B} \phi(x^{(i)}) \cdot \phi(y^{(i)}) = \sum_{x^{(j)} \in A} \sum_{y^{(j)} \in B} \sum_{j=1}^d \varphi(x_j^{(i)}) \varphi(y_j^{(i)}) = \otimes$$

$$\underbrace{\sum_{x^{(i)} \in A} \phi(x^{(i)})}_{\hat{\Phi}(x^{(i)})} \cdot \underbrace{\sum_{y^{(i)} \in B} \phi(y^{(i)})}_{\hat{\Phi}(y^{(i)})} = \begin{bmatrix} \sum_{x^{(i)} \in A} \varphi(x_1^{(i)}) \\ \vdots \\ \sum_{x^{(i)} \in A} \varphi(x_d^{(i)}) \end{bmatrix} \cdot \begin{bmatrix} \sum_{y^{(i)} \in B} \varphi(y_1^{(i)}) \\ \vdots \\ \sum_{y^{(i)} \in B} \varphi(y_d^{(i)}) \end{bmatrix} = \otimes$$

بنابراین نشان دادیم  $\hat{k}(x^{(i)}, y^{(i)}) = \hat{\Phi}(x^{(i)}) \cdot \hat{\Phi}(y^{(i)})$  پس  $\hat{k}$  کرنل معتبر است

2.1

پس از دست آمدن  $\alpha$  به دست آوردن  $b$  داریم:

$$y^{(n)} (w^T x^{(n)} + b) = 1$$

به ازای هر  $x^{(n)}$  که  $SV$  باشد

حال اگر فرض را در  $\alpha_i$  ضرب کرده و sum بگیریم رابطه زیر برقرار است زیرا طبق KKT condition می دانیم که  $x^{(i)}$  های که  $SV$  نباشند،  $\alpha_i$  متناظرشان 0 است

$$\sum_{i=1}^N \alpha_i y^{(i)} (w^T x^{(i)} + b) = \sum_{i=1}^N \alpha_i \rightarrow w^T w = \|w\|^2 = \sum_{i=1}^N \alpha_i$$

$$w^T \underbrace{\sum_{i=1}^N \alpha_i y^{(i)} x^{(i)}}_w + b \underbrace{\sum_{i=1}^N \alpha_i y^{(i)}}_0$$

$$\text{margin} = \frac{1}{\|w\|} = \frac{1}{\sqrt{\sum_{i=1}^N \alpha_i}}$$

2.2

می دانیم که اگر بعد جدایی افشان کنیم فاصله نقاط تنها می تواند افزایش یابد. با توجه به اینکه معادله سوال زمین کرده که margin افزایش پیدا نمی کند پس تفسیری در margin داریم. پس ابرصفحه جدا کننده همان ابرصفحه ی قبلی می ماند و فقط یک نفر به آن افشان می شود اگر داده ها از قبل linearly separable باشند بعد از افزایش بعد نیز قابل جدا کردن اند و همین ابرصفحه این کار را انجام می دهد پس در فاز training ایده یی نمی آید و مدل robust است اما برای generalization مشکلاتی نظیر spurious relationship ایجاد می شود

2.3

دو داده با لیبیل متفاوت linearly separable است بنابراین طبق نرمول که برای SVM hard margin داریم، ماتریس زیر را به Q ورودی می دهیم:

$$Q = \begin{bmatrix} y^{(1)} y^{(1)} x^{(1)T} x^{(1)} & y^{(1)} y^{(2)} x^{(1)T} x^{(2)} \\ y^{(2)} y^{(1)} x^{(2)T} x^{(1)} & y^{(2)} y^{(2)} x^{(2)T} x^{(2)} \end{bmatrix}$$

سوال لازم برای اینکه بهینه سازی convex شود و جواب را بتوانیم از Q بگیریم آن است که این ماتریس PSD باشد

لیبل متفاوت

$$\begin{bmatrix} x^{(1)T} x^{(1)} - x^{(1)T} x^{(2)} \\ -x^{(2)T} x^{(1)} & x^{(2)T} x^{(2)} \end{bmatrix}$$

برای ماتریس های 2x2 می دانیم که ماتریس PSD است اگر و تنها اگر شرایط زیر برقرار باشند:

- Q is symmetric. ← مشاهده می شود که برقرار است ✓
- trace(Q) > 0. ← trace(Q) =  $\|x^{(1)}\|^2 + \|x^{(2)}\|^2 > 0$  ✓
- det(Q) ≥ 0. ← det(Q) =  $\|x^{(1)}\|^2 \|x^{(2)}\|^2 - \langle x^{(1)}, x^{(2)} \rangle^2 \geq 0$  ✓

طبق کوشی شوارتز

بنابراین Q جواب (X) را به دست می آورد و طبق نرمول های درس w و b را به دست می آوریم و معادله ابرصفحه را به صورت کامل داریم و بنابراین فاصله آن تا مبدأ را که برابر  $\frac{|b|}{\|w\|}$  است را نیز داریم. اگر لیبیل ها متفاوت نبود ابرصفحه بین آنها نمی افتاد و margin مقدار می خورد نمی داشت



3.1

متغیر  $\epsilon_i^*$  را میزان تخلفی مربوط به کمترین بودن  $f(x^{(i)})$  از  $y^{(i)}$  در نظر می‌گیریم پس می‌توان نوشت

$$\epsilon_i^* = \max(y^{(i)} - f(x^{(i)}) - \epsilon, 0)$$

متغیر  $\epsilon_i$  را میزان تخلفی مربوط به بیشترین بودن  $f(x^{(i)})$  از  $y^{(i)}$  در نظر می‌گیریم پس می‌توان نوشت

$$\epsilon_i = \max(f(x^{(i)}) - y^{(i)} - \epsilon, 0)$$

متغیر  $\epsilon$  را برای پارامتری کردن جریمه تخلفی به فرمول افاده کردیم

حال دافع است که در آن واحد تنها یکی از جریمه‌ها ( $\epsilon_i^*$  یا  $\epsilon_i$ ) اعمال خواهد شد زیرا شرط مثبت بودن آنها استراک ندارد

$$\epsilon_i^* \geq 0 \rightarrow y^{(i)} \geq f(x^{(i)}) + \epsilon$$

استراک ندارد (\*)

$$\epsilon_i \geq 0 \rightarrow y^{(i)} \leq f(x^{(i)}) - \epsilon$$

$$L_{\epsilon}(x^{(i)}, y^{(i)}, f) = \begin{cases} \max(y^{(i)} - f(x^{(i)}) - \epsilon, 0) & y^{(i)} \geq f(x^{(i)}) \\ \max(f(x^{(i)}) - y^{(i)} - \epsilon, 0) & y^{(i)} < f(x^{(i)}) \end{cases} \stackrel{(*)}{=} \epsilon_i + \epsilon_i^*$$

$$① \epsilon_i^* \geq 0$$

$$③ \epsilon_i \geq 0$$

$$② \epsilon_i^* \geq y^{(i)} - f(x^{(i)}) - \epsilon \quad ④ \epsilon_i \geq f(x^{(i)}) - y^{(i)} - \epsilon$$

بنابراین صورت primal مسئله به صورت گفته شده در می‌آید و شرط دوم او را داریم

3.2

primal

$$\min_{w, b, \epsilon_i, \epsilon_i^*} \max_{\alpha_i, \alpha_i^*, \beta_i, \beta_i^*} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N (\epsilon_i + \epsilon_i^*) - \sum_{i=1}^N \alpha_i^* \epsilon_i^* - \sum_{i=1}^N \beta_i^* (\epsilon_i^* - y^{(i)} + \overbrace{f(x^{(i)})}^{w^T x^{(i)} + b} + \epsilon) - \sum_{i=1}^N \alpha_i \epsilon_i - \sum_{i=1}^N \beta_i (\epsilon_i - \underbrace{f(x^{(i)})}_{w^T x^{(i)} + b} + y^{(i)} + \epsilon)$$

مسئله dual جای  $\min$  و  $\max$  عوض می‌کند. با توجه به اینکه برای مسئله  $primal = dual$  پس ابتدا نسبت به  $w, b, \epsilon_i$  و  $\epsilon_i^*$  کمینه می‌کنیم

$$\nabla_w L = w - \sum_{i=1}^N \beta_i^* x^{(i)} + \sum_{i=1}^N \beta_i x^{(i)} = 0 \rightarrow w = \sum_{i=1}^N (\beta_i^* - \beta_i) x^{(i)}$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^N -\beta_i^* + \sum_{i=1}^N \beta_i = 0 \rightarrow \sum_{i=1}^N (\beta_i - \beta_i^*) = 0$$

$$\frac{\partial L}{\partial \epsilon_i} = C - \alpha_i - \beta_i = 0 \rightarrow \alpha_i = C - \beta_i$$

$$\frac{\partial L}{\partial \epsilon_i^*} = C - \alpha_i^* - \beta_i^* = 0 \rightarrow \alpha_i^* = C - \beta_i^*$$

KKT conditions

$$\alpha_i \epsilon_i = (C - \beta_i) \epsilon_i = 0$$

$$\alpha_i^* \epsilon_i^* = (C - \beta_i^*) \epsilon_i^* = 0$$

$$\beta_i (\epsilon_i - f(x^{(i)}) + y^{(i)} + \epsilon) = 0$$

$$\beta_i^* (\epsilon_i^* - y^{(i)} + f(x^{(i)}) + \epsilon) = 0$$

### 3.2 ادامه

$$\begin{aligned}
 \rightarrow L &= \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N (\beta_i^* - \beta_i) (\beta_j^* - \beta_j) x^{(i)T} x^{(j)} + C \sum_{i=1}^N (\cancel{\varepsilon_i} + \cancel{\varepsilon_i^*}) \\
 &\quad - \sum_{i=1}^N (\cancel{\beta_i^*}) \varepsilon_i^* - \sum_{i=1}^N \beta_i^* (\varepsilon_i^* - y^{(i)}) + \sum_{j=1}^N ((\beta_j^* - \beta_j) x^{(j)T} x^{(i)}) + \cancel{\beta_i} + \varepsilon \\
 &\quad - \sum_{i=1}^N (\cancel{\beta_i}) \varepsilon_i - \sum_{i=1}^N \beta_i (\varepsilon_i - \sum_{j=1}^N ((\beta_j^* - \beta_j) x^{(j)T} x^{(i)})) + \cancel{\beta_i} + y^{(i)} + \varepsilon \\
 &= \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N (\beta_i^* - \beta_i) (\beta_j^* - \beta_j) x^{(i)T} x^{(j)} + \sum_{i=1}^N \cancel{\beta_i^*} \varepsilon_i^* + \sum_{i=1}^N \cancel{\beta_i} \varepsilon_i \\
 &\quad - \varepsilon \sum_{i=1}^N (\beta_i + \beta_i^*) + \sum_{i=1}^N y^{(i)} (\beta_i^* - \beta_i) - \sum_{i=1}^N \cancel{\beta_i^*} \varepsilon_i^* - \sum_{i=1}^N \cancel{\beta_i} \varepsilon_i \\
 &\quad - \sum_{i=1}^N \sum_{j=1}^N (\beta_i^* - \beta_i) (\beta_j^* - \beta_j) x^{(j)T} x^{(i)} \\
 &= \boxed{-\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N (\beta_i^* - \beta_i) (\beta_j^* - \beta_j) x^{(i)T} x^{(j)} - \varepsilon \sum_{i=1}^N (\beta_i + \beta_i^*) + \sum_{i=1}^N y^{(i)} (\beta_i^* - \beta_i)}
 \end{aligned}$$

$$\left. \begin{aligned} \beta_i &\geq 0 \\ \alpha_i &\geq 0 \rightarrow C - \beta_i \geq 0 \rightarrow \beta_i \leq C \end{aligned} \right\} \rightarrow 0 \leq \beta_i \leq C$$

حال باید عبارت به دست آمده  $\max$  شود، با شرط رویه در

نام گذاری من با صورت سوال فرق دارد و به جای  $\alpha$ ،  $\beta$  گذاشته ام

### 3.3

مسئله دسته بندی به تقاری به حسب یک متغیر لاگرانژ رسیدیم و می توان آن را توسط QP حل کرد. کافی است فرم ماتریسی عبارت به دست آمده را بنویسیم و سپس ماتریس ضرایب را به دست آوریم. البته چون عبارت به دست آمده با مسئله دسته بندی متفاوت است، روزه متفاوتی باید برای به دست آوردن ماتریس ضرایب و ورودی دادن آن به QP انجام شود

### 3.4

مساوی قبل داده هایی که  $\beta_i$  یا  $\beta_i^*$  شان مثبت باشند SV فعالند بعد

$$\left. \begin{aligned} \beta_i > 0 &\xrightarrow{KKT} \varepsilon_i = f(x^{(i)}) - y^{(i)} - \varepsilon \\ \varepsilon_i &= \max(f(x^{(i)}) - y^{(i)} - \varepsilon, 0) \end{aligned} \right\} \rightarrow y^{(i)} \leq f(x^{(i)}) - \varepsilon$$

⊗ بنابراین داده هایی که خارج از رنج  $F \pm \varepsilon$  هستند SV هستند

$$\left. \begin{aligned} \beta_i^* > 0 &\xrightarrow{KKT} \varepsilon_i^* = y^{(i)} - f(x^{(i)}) - \varepsilon \\ \varepsilon_i^* &= \max(y^{(i)} - f(x^{(i)}) - \varepsilon, 0) \end{aligned} \right\} \rightarrow y^{(i)} \geq f(x^{(i)}) + \varepsilon$$

3.5

$$W = \sum_{i=1}^N (\beta_i^* - \beta_i) x^{(i)} \rightarrow F(x^{(k)}) = \sum_{i=1}^N \left( (\beta_i^* - \beta_i) x^{(i)T} x^{(k)} \right) + b$$

از رابطه دوبار می توان برای بیشینه استفاده کرد

همچنین می توان از کرنل هم استفاده کرد و در تمامی فرمول ها به جای جواب ضرب داخلی  $\phi(x^{(i)})$  ها نیاز است و نمی خواهد فقط را ملاقات کنیم

$$W = \sum_{i=1}^N (\beta_i^* - \beta_i) \phi(x^{(i)}) \rightarrow F(x^{(k)}) = \sum_{i=1}^N \left( (\beta_i^* - \beta_i) \underbrace{\phi(x^{(i)})^T \phi(x^{(k)})}_{k(x^{(i)}, x^{(k)})} \right) + b$$

3.6

اگر  $F(x)$  را معادل خط / ابر صفحه به دست آمده بدانیم ، بازه ای  $\pm \epsilon$  بازه ای است که خطای به آن اجازه tolerate می شود  
 همپایه هایی که داخل این بازه قرار دارند اردو صفر دارند - بنابراین افزایش  $\epsilon$  باعث کاهش پیچیدگی مدل و آسان گرفتن می شود  
 و کاهش  $\epsilon$  نیز باعث افزایش پیچیدگی مدل و سخت گیری می شود  $\left( \frac{1}{\epsilon} \right)$  ضریب regularization عمل می کند

$\epsilon$  بزرگ یعنی  $\epsilon_1$  و  $\epsilon_2$  باید کوچک باشند و اجازه violation کمی می دهیم و به همین صورت  $\epsilon$  کوچک یعنی اجازه violation بیشتری می دهیم  
 مانند قبل افزایش  $\epsilon$  به نوعی باعث کاهش پیچیدگی مدل و کاهش  $\epsilon$  باعث افزایش پیچیدگی مدل می شود  $\left( \frac{1}{\epsilon} \right)$  ضریب regularization عمل می کند

4.1

اگر سمبل که حرف می شود جزء SV ها باشد هیچ تغییری رخ نمی دهد زیرا فقط SV ها روی margin تأثیر می گذارند

اگر سمبل از SV ها باشد اما به جز آن حداقل 2 تا SV دیگر با همان لپل وجود داشته باشد باز تغییری ایجاد نمی شود زیرا 2 نقطه برای نگه داشتن خط کافی است (ما برای خط صرف رژیم اما به طور کلی برای ابرصفحه n بعدی حداقل n-1 سمبل SV نیاز است تا مرز جابجا نشود) در غیر این صورت مرز رژیم به سمت سمبل حرف شده می رود تا margin بزرگتری به دست آید

logistic regression مرز رژیم تغییر می کند زیرا تابع هزینه ای داریم که به ازای هر سمبل error آن را حساب کرده و بر اساس آن مرز رژیم را جابجا می کند پس حتی اگر batch هم داشته باشیم، هر سمبل در loss تأثیر گذاشته و هدف سمبل تماماً مرز رژیم را جابجا می کند

4.2

می دانیم که  $\epsilon_i > 1$  معادل است با لپل خوردن سمبل نام است. یعنی به طور کلی داریم

$\epsilon_i = 0 \rightarrow \text{correctly classified}$

$0 < \epsilon_i \leq 1 \rightarrow \text{correctly classified but inside margin}$

$\epsilon_i > 1 \rightarrow \text{misclassified}$

پس واضح است که  $\sum_{i=1}^N \epsilon_i$  یک کران بالا برای  $\sum_{i=1}^N \text{error}$  است که استنباط لپل زده شده اند

4.3

$$\text{minimize } \frac{1}{2} \|w\|^2 + c \sum_{i=1}^N \epsilon_i$$

$c \rightarrow \infty$

یعنی باید  $\epsilon_i$  ها کوچک باشند و اجازه violation نمی دهیم (به حالت hard margin برگشتیم)

$c \rightarrow 0$

یعنی  $\epsilon_i$  ها می توانند بسیار بزرگ باشند و اجازه violation زیادی می دهیم. با توجه به این حالت وزن بیشتری برای  $\|w\|^2$  است پس  $\|w\|$  بیشتر کمینه می شود و بنابراین margin بزرگتری داریم

$\frac{1}{c}$  را می توان به عنوان ضریب regularization دید زیرا بین پیچیدگی مدل، generalization یک trade-off ایجاد می کند اگر  $c$  را بزرگ قرار دهیم اجازه پیچیدگی بیشتری به مدل می دهیم، و اگر  $c$  را کم انتخاب کنیم جلوی پیچیدگی را می گیریم



4.4

مدل logistic regression به دلیل داشتن تابع هزینه بین ابرصفحه‌های جداگشته که به درستی دنیای آموزش را جدا کرده اند نیز تفاوت قائل می‌شود و به طور کلی عملکرد خوبی در پیدا کردن مرز که بهتر generalize کند دارد. اما SVM ابرصفحه با margin بیشتر را انتخاب می‌کند و از این نظر از عیب generalization از LR بهتر است

بنابراین اگر دنیا جدایی پذیر خطی باشد هر 2 مدل جواب را پیدا می‌کنند و روی داده آموزش به ارور صفر می‌رسند. تفاوت آن است که SVM در زمان تست بهتر generalize می‌کند (هم از نظر ابرصفحه‌ای که margin بزرگتری دارد و هم از نظر robust بودن نسبت به outlier ها در داده آموزش)

4.5

در این حالت عملکرد 2 مدل شبیه هم می‌شود (به جز robust بودن SVM نسبت به outlier ها) و هر 2 مدل جوابی sub optimal برمی‌گردانند که تا حدی توانسته داده را دسته‌بندی کند.

ضمناً هر دو مدل می‌توانند از تبدیل غیرخطی  $\Phi$  استفاده کنند اما SVM از این نظر بهتر است چون نیازی ندارد که فضای  $\Phi$  را ملاقات کند و بنابراین سریع‌تر است. حتی می‌تواند داده‌ها را به فضای  $\infty$  بعدی ببرد که LR نمی‌تواند و از این نظر قابلیت دسته‌بندی بالاتری دارد