

سوال ①

الف) میزان تخلفی مثبت بودن ξ_i را ξ_i ←

$$\langle w, x_i \rangle - y_i - \epsilon \leq \xi_i$$

$$y_i - \langle w, x_i \rangle - \epsilon \leq \xi_i^* \quad \text{where } \xi_i, \xi_i^* \geq 0$$

پس می توانیم تابع هزینه کلی را بصورت زیر بنویسیم:

$$\min_{w \in \mathbb{R}^m, \xi_i, \xi_i^* \in \mathbb{R}^n} \frac{1}{2} \|w\|^2 + c \sum_{i=1}^n (\xi_i + \xi_i^*)$$

ب) تابع لاگرانژین بصورت زیر نوشته می شود:

$$\begin{aligned} L(w, \xi_i, \xi_i^*, \alpha, \alpha^*, \beta, \beta^*) &= \frac{1}{2} \|w\|^2 + c \sum_{i=1}^n (\xi_i + \xi_i^*) \\ &- \sum_{i=1}^n (\beta_i \xi_i + \beta_i^* \xi_i^*) - \sum_{i=1}^n \alpha_i (\epsilon + \xi_i - y_i + \langle w, x_i \rangle) \\ &- \sum_{i=1}^n \alpha_i^* (\epsilon + \xi_i^* + y_i - \langle w, x_i \rangle) \end{aligned}$$

چون باید شرط مثبت بودن را ارضاء کند $\alpha_i, \alpha_i^*, \beta_i, \beta_i^* \geq 0$

$$\min_{w, \xi_i, \xi_i^*} \max_{\alpha, \alpha^*, \beta, \beta^*} L(w, \xi, \xi^*, \alpha, \alpha^*, \beta, \beta^*)$$

$$= \max_{\alpha, \alpha^*, \beta, \beta^*} \min_{w, \xi, \xi^*} L(w, \xi, \xi^*, \alpha, \alpha^*, \beta, \beta^*)$$

حل مسبقاً: w, ξ, ξ^* مستقيم

$$\frac{\partial L}{\partial w} = w - \sum_{i=1}^n (\alpha_i - \alpha_i^*) x_i = 0$$

$$\frac{\partial L}{\partial \xi_i} = C - \beta_i - \alpha_i = 0 \Rightarrow \beta_i = C - \alpha_i \geq 0$$

$$\frac{\partial L}{\partial \xi_i^*} = C - \beta_i^* - \alpha_i^* = 0 \Rightarrow \beta_i^* = C - \alpha_i^* \geq 0$$

$$\begin{aligned} L &= \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*) - \sum_{i=1}^n \beta_i \xi_i + \beta_i^* \xi_i^* \\ &\quad - \sum_{i=1}^n \alpha_i (\epsilon + \xi_i - y_i + \langle w, x_i \rangle) - \sum_{i=1}^n \alpha_i^* (\epsilon + \xi_i^* + y_i - \langle w, x_i \rangle) \\ &= \frac{1}{2} \left\| \sum_{i=1}^n (\alpha_i^* - \alpha_i) x_i \right\|^2 + \sum_{i=1}^n \xi_i (C - \beta_i - \alpha_i) + \sum_{i=1}^n \xi_i^* (C - \beta_i^* - \alpha_i^*) \\ &\quad - \epsilon \sum_{i=1}^n (\alpha_i + \alpha_i^*) + \sum_{i=1}^n y_i (\alpha_i - \alpha_i^*) + \sum_{i=1}^n (\alpha_i^* - \alpha_i) \langle w, x_i \rangle \\ &\quad \quad \quad \left\langle \sum_{j=1}^n (\alpha_j - \alpha_j^*) x_j, x_i \right\rangle \end{aligned}$$

$$= -\frac{1}{2} \sum_{i,j=1}^n (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) \langle x_i, x_j \rangle -$$

$$\in \sum_{i=1}^n (\alpha_i + \alpha_i^*) + \sum_{i=1}^n y_i (\alpha_i - \alpha_i^*)$$

صورت دومین:

$$\max_{\alpha, \alpha^*} -\frac{1}{2} \sum_{i,j=1}^n (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) \langle x_i, x_j \rangle$$

$$- \in \sum_{i=1}^n (\alpha_i + \alpha_i^*) + \sum_{i=1}^n y_i (\alpha_i - \alpha_i^*)$$

$$\alpha, \alpha^* \in [0, c]$$

ب) مسئله دارای هدف درجه دوم است و محدودیت‌های خطی است پس می‌توان با یک برنامه Quadratic Programming آن را حل کرد.

ت) در پاسخ optimal (استیم): $\alpha_i (\epsilon + \xi_i - y_i + \langle w, x_i \rangle) = 0$

$$\alpha_i^* (\epsilon + \xi_i^* + y_i - \langle w, x_i \rangle) = 0$$

$$\alpha_i, \alpha_i^* > 0 \Rightarrow \epsilon + \xi_i - y_i + \langle w, x_i \rangle = 0$$

$$\alpha_i^*, \xi_i^* = 0$$

اگر $\xi_i = 0$ باشد نشان می‌دهد که α_i روی مرزهای SVM قرار دارد پس α_i یک بردار margin support است. اگر $\xi_i > 0$ باشد نشان می‌دهد α_i خارج از منطقه حاشیه ϵ است پس α_i یک بردار non-margin support است.

اذاً n برای α_i^* و α_i^* نیز این موقع برقرار است.

$$f(\alpha) = \langle w, \alpha \rangle \quad (c)$$

$$\Rightarrow f(\alpha) = \sum_{i=1}^n (\alpha_i - \alpha_i^*) \langle \alpha_i, \alpha \rangle$$

$$w = \sum_{i=1}^n (\alpha_i - \alpha_i^*) \alpha_i$$

می‌توان عبارت بالا را به فرم کرنل نوشت:

$$f(\alpha) = \sum_{i=1}^n (\alpha_i - \alpha_i^*) K(\alpha_i, \alpha)$$

(ج)

نقص ϵ و C برعکس می‌باشد.
 هرچه ϵ کوچکتر باشد SVM سعی می‌کند روی خط‌های کمتری fit شود.
 مدل تولید SVM پیچیده‌تر و واریانس آن زیادتر و بایاس آن کمتر می‌شود.
 هرچه C بزرگتر باشد، مدل تولید ساده‌تر، واریانس آن پایین‌تر اما بایاس آن بیشتر خواهد بود.

C معیاری است که چقدر برای استیلاهای جریمه در نقطه می‌گیریم.
 هرچه قدر C را بزرگتر در نقطه بگیریم، مدل می‌خواهد جریمه بیشتری برای خطا در نقطه بگیرد پس روی داده‌ها $over fit$ می‌شود و بایاس کم اما واریانس زیاد دارد.
 اما هرچه قدر C کوچک باشد، مدل ساده‌تر با بایاس زیاد و واریانس کم خواهیم داشت.

سوال ②

برای بدست آوردن حداقل تعداد نمونه m از فرمول زیر استفاده می کنیم:

$$m \geq \frac{1}{\epsilon} \left(\ln |H| + \ln \frac{1}{\delta} \right)$$

$$m \geq \frac{1}{0.05} \left(\ln 1000 + \ln \frac{1}{0.05} \right)$$

$$m \geq 198.06 \Rightarrow \boxed{m \geq 199} \quad \text{چون باید طبیعی باشد.}$$

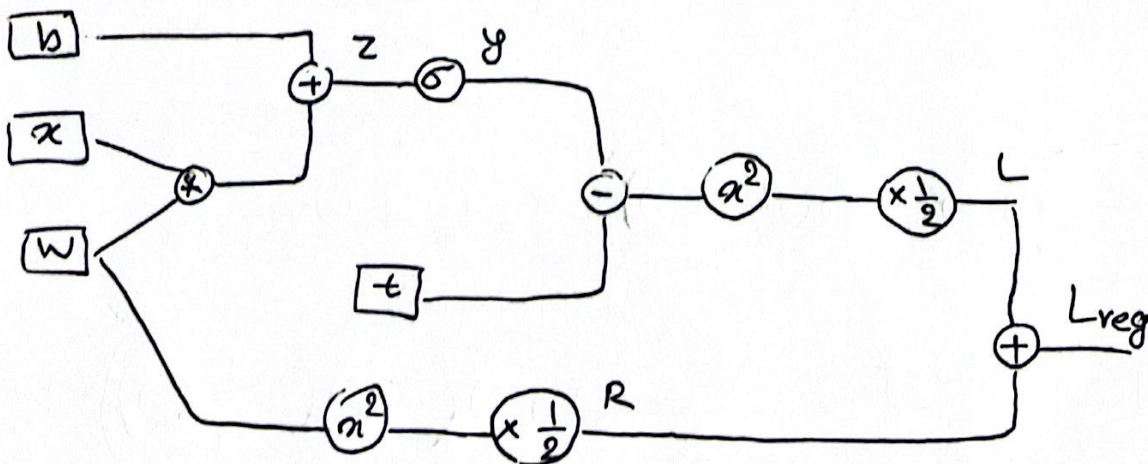
$$z = wx + b$$

$$y = \sigma(z)$$

سوال ③

الف)

$$L_{\text{reg}} = \underbrace{\frac{1}{2} (y - t)^2}_L + \underbrace{\lambda \frac{1}{2} w^2}_R$$



$$\frac{\partial L_{\text{reg}}}{\partial L} = 1, \quad \frac{\partial L_{\text{reg}}}{\partial R} = \lambda$$

$$\frac{\partial L}{\partial y} = (y - t)$$

$$\frac{\partial L_{\text{reg}}}{\partial y} = \frac{\partial L_{\text{reg}}}{\partial L} \times \frac{\partial L}{\partial y} = (y - t)$$

$$\frac{\partial L_{\text{reg}}}{\partial z} = \frac{\partial L_{\text{reg}}}{\partial y} \times \frac{\partial y}{\partial z} = (y - t) \left(\sigma(z) (1 - \sigma(z)) \right)$$

↓
 $y = \sigma(z)$

$$\frac{\partial L_{\text{reg}}}{\partial b} = \frac{\partial L_{\text{reg}}}{\partial z} \times \frac{\partial z}{\partial b} = (y - t) \left(\sigma(z) - (1 - \sigma(z)) \right)$$

$$\frac{\partial L_{\text{reg}}}{\partial w} = \frac{\partial L_{\text{reg}}}{\partial L} \times \frac{\partial L}{\partial z} \times \frac{\partial z}{\partial w} + \frac{\partial L_{\text{reg}}}{\partial R} \times \frac{\partial R}{\partial w}$$

$$\frac{\partial L_{\text{reg}}}{\partial w} = (y - t) \left(\sigma(z) (1 - \sigma(z)) \right) \times \alpha + \lambda w$$

$$\frac{\partial L_{\text{reg}}}{\partial \alpha} = \frac{\partial L_{\text{reg}}}{\partial z} \times \frac{\partial z}{\partial \alpha} = (y - t) \left(\sigma(z) (1 - \sigma(z)) \right) \cdot w$$

(ب) وزن‌ها بصورت زرد در ابتدا انتخاب می‌شوند چون شروع از هر وزنی ما را می‌تواند به نقطه‌ی بهینه‌ی دیگری برساند و برخی از این نقاط بهتر از نقاط دیگر باشند. اما اگر همواره از یک نقطه مشخص فرایند آموزش را انجام دهیم، همواره به یک نقطه بهینه.

local/global minima می‌رسیم و این دلخواه نخواهد بود.

وزن‌ها بصورت مقادیر کوچک نزدیک صفر انتخاب می‌شوند. چون در این نقاط گزاین تابع activation غیر صفر است و کمک می‌کند تا با الگوریتم گزاین کاهشی به نقاط بهینه‌ی درست برسیم. اما اگر وزن‌ها مقادیر بزرگی باشند مثلاً در تابع فعال سگ sigmoid دارای گزاین صفر خواهیم بود و باعث جا به جایی درستی در الگوریتم gradient descent نخواهد شد و در یک جا ممکن است متوقف شویم که بهینه محلی است و نتوانیم از آن فرار کنیم.

$$\alpha = 0.1 \rightarrow \text{learning_rate}$$

(ج)

$$\text{initial value : } \begin{cases} a=1 \\ b=1 \\ w=2 \\ t=1 \\ \lambda=0.01 \end{cases}$$

Forward propagation

$$z = w\alpha + b = 2 + 1 = 3$$

$$y = \frac{1}{1 + e^{-z}} = \frac{1}{1.0497} = 0.952$$

$$L_{\text{reg}} = \frac{1}{2} (1 - 0.952)^2 + 0.01 \times \frac{1}{2} 2^2$$

backward propagation

$$w' = w - \alpha \frac{\partial L_{\text{reg}}}{\partial w} = 2 - 0.1 \times \left[(0.952 - 1)(0.952(1 - 0.952)) \cdot 1 + 0.01 \times 2 \right]$$

$$\Rightarrow w' = 2 - 0.1 \times 0.0178 \approx 1.99822$$

$$b' = b - \alpha \frac{\partial L_{\text{reg}}}{\partial b} = 1 - 0.1 \left[(0.952 - 1)(0.952(1 - 0.952)) \right]$$

$$\Rightarrow b' = 1 - 0.1 \times (-0.0021) \approx 1.00021$$