

$$y_i = x_i^T \beta_0 + \varepsilon_i \quad (1)$$

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

$$R_{tr}(\hat{\beta}) = \frac{1}{N} \sum_{i=1}^N (y_i - x_i^T \hat{\beta})^2 = \frac{1}{N} (y - X\hat{\beta})^T (y - X\hat{\beta})$$

$$H = X (X^T X)^{-1} X^T$$

$$y - X\hat{\beta} = y - Hy$$

$$R_{tr}(\hat{\beta}) = \frac{1}{N} (y - Hy)^T (y - Hy)$$

$$H^T = X (X^T X)^{-1} X^T = H$$

$$H^T = X \underbrace{(X^T X)^{-1} X^T X}_{I} (X^T X)^{-1} X^T = X (X^T X)^{-1} X^T = H$$

$$H^T v = \lambda^T v \xRightarrow{H^T=H} H v = \lambda^T v \Rightarrow \lambda v = \lambda^T v \Rightarrow$$

$$\lambda v (\lambda - 1) = 0 \Rightarrow \lambda = 0 \quad \vee \quad \lambda = 1$$

به تعداد رتبه  
 $P = H$

$$\text{tr}(H) = \text{rank}(H) = p$$



$$E[R_{tr}(\hat{\beta})] = \frac{1}{N} E[\text{tr}((y - Hy)(y - Hy)^T)] \quad (1)$$

$$= \frac{1}{N} E[\text{tr}((I - H)y y^T (I - H)^T)] = \frac{1}{N} \text{tr}[(I - H)E[yy^T](I - H)]$$

$$= \frac{1}{N} \text{tr}[(I - H)[\text{Var}(y) + E[y]E[y]^T](I - H)]$$

$$= \frac{1}{N} \text{tr}[(I - H)[\sigma^2 I + X\beta\beta^T X^T](I - H)]$$

$$= \frac{1}{N} \text{tr}[\sigma^2(I - H) + (I - H)X\beta\beta^T X^T(I - H)]$$

$$= \frac{1}{N} \text{tr}[\sigma^2(I - H)] = \frac{1}{N} \text{tr}[\sigma^2(I - H)] = \frac{\sigma^2}{N} (N - p)$$

$$= \sigma^2 \left(1 - \frac{p}{N}\right) \leq \sigma^2 \leq E[R_{te}(\hat{\beta})]$$

$$* (y - Hy)^T (y - Hy) = \text{tr}[(y - Hy)^T (y - Hy)]$$

$$= \text{tr}[(y - Hy)(y - Hy)^T]$$

$$E[(y - \bar{x}^T \hat{\beta})^T] = E[(\bar{x}^T \beta + \bar{\varepsilon} - \bar{x}^T \hat{\beta})^T] = E[\bar{x}^T (\beta - \hat{\beta})] + E[\bar{\varepsilon}] + \bar{x}^T E[\hat{\beta} - \beta]$$

$$E[R_{te}(\hat{\beta})] = \frac{1}{N} \sum_{i=1}^N E[(y_i - \bar{x}_i^T \hat{\beta})^2] \geq \sigma^2$$



$$Y = \beta_0 + \beta_1 X + \varepsilon \Rightarrow \hat{\beta}^{(lin)} = \arg \min_{\beta} \|y - X_{linear} \beta\|_2^2$$

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \varepsilon \Rightarrow \hat{\beta}^{(poly)} = \arg \min_{\beta} \|y - X_{poly} \beta\|_2^2$$

$$\min_{\beta} \|y - X_{linear} \beta\|_2^2 = \|y - X_{linear} \hat{\beta}^{(lin)}\|_2^2$$

$$\beta^* = (\hat{\beta}_0^{(lin)}, \hat{\beta}_1^{(lin)}, \dots, 0) :$$

$$\|y - X_{poly} \beta^*\|_2^2 = \|y - X_{linear} \hat{\beta}^{(lin)}\|_2^2$$

$$\|y - X_{poly} \hat{\beta}^{(poly)}\|_2^2 \leq \|y - X_{poly} \beta^*\|_2^2 = \|y - X_{linear} \hat{\beta}^{(lin)}\|_2^2$$

$$\Rightarrow \min_{\beta} \|y - X_{poly} \beta\|_2^2 \leq \min_{\beta} \|y - X_{linear} \beta\|_2^2$$



(۲)

$$y = f(x) + \varepsilon$$

$$\hat{y}_1 = f_1(x)$$

$$\hat{y}_p = f_p(x)$$

$$E_1 = \frac{1}{n} \sum_{i=1}^n (f(x_i) + \varepsilon_i - f_1(x_i))^2$$

$$E_p = \frac{1}{n} \sum_{i=1}^n (f(x_i) + \varepsilon_i - f_p(x_i))^2$$

$$E = E[(f(x) + \varepsilon - \hat{f}(x))^2] = \int \int p(x, \varepsilon) (f(x) + \varepsilon - \hat{f}(x))^2 dx d\varepsilon$$

در تعداد داده زیاد خطای هر دو تقریباً برابر است اما در تعداد داده کم تابع (مدل) بیشترین خطای گسترده روی داده های train دارد و overfit رخ میدهد.

برعکس در داده های test مدل درجه ۳ (که بیشترین است و overfit کرده) خطای

بیشتری دارد.

$$n \rightarrow \infty : \int \int (f(x) + \varepsilon - \hat{f}_1(x))^2 dx d\varepsilon = \int \int (f(x) + \varepsilon - \hat{f}_p(x))^2 dx d\varepsilon = E$$

$$\left. \begin{aligned} |E_1 - E| &< \varepsilon \\ |E_p - E| &< \varepsilon \end{aligned} \right\} \Rightarrow |E_1 - E_p| = |E_1 - E - E_p + E| \leq |E_1 - E| + |E_p - E| < 2\varepsilon$$



$$J_{\lambda}(w) = \frac{1}{r} (y - Xw)^T (y - Xw) + \lambda \|w\|_1 \quad \text{الف (۳)}$$

$$= \frac{1}{r} (y^T y - y^T X w - w^T X^T y + \underbrace{w^T X^T X w}_I) + \lambda \|w\|_1$$

$$= \frac{1}{r} y^T y + (-y^T X w + \frac{1}{r} w^T w + \lambda \|w\|_1)$$

$$= \underbrace{\frac{1}{r} y^T y}_{g(y)} + \sum_{i=1}^d \underbrace{\left( \frac{1}{r} w_i^2 + \lambda |w_i| - y^T X_{:,i} w_i \right)}_{f(X_{:,i}, y_i, w_i, \lambda)}$$

$$w_i > 0 \Rightarrow J_{\lambda}(w) = \frac{1}{r} y^T y + \sum_{i=1}^d \frac{1}{r} w_i^2 + \lambda w_i - y^T X_{:,i} w_i$$

$$\frac{\partial J_{\lambda}(w)}{\partial w_i} = 0 \Rightarrow w_i + \lambda - y^T X_{:,i} = 0 \Rightarrow$$

$$w_i = y^T X_{:,i} - \lambda$$

$$w_i < 0 \Rightarrow \text{به طور مشابه} \Rightarrow w_i = y^T X_{:,i} + \lambda \quad \text{ب}$$



۳) ت - در صورتی که  $w_i$  صفر می شود که در هر دو رابطه بالا صدق کند یعنی:

$$y^T x_{:,i} - \lambda = y^T x_{:,i} + \lambda \Rightarrow 2\lambda = 0 \Rightarrow \lambda = 0$$

$$y^T x_{:,i} - \lambda > 0 \Rightarrow y^T x_{:,i} > \lambda \text{ و } y^T x_{:,i} < -\lambda \Rightarrow |y^T x_{:,i}| < \lambda \Rightarrow w_i = 0$$

$$J_{\lambda}^{(w)} = \frac{1}{2} \|y - Xw\|_2^2 + \frac{\lambda}{2} \|w\|_2^2$$

$$= \frac{1}{2} (y - Xw)^T (y - Xw) + \frac{\lambda}{2} w^T w$$

$$= \frac{1}{2} y^T y - y^T Xw + \frac{1}{2} w^T X^T X w + \frac{\lambda}{2} w^T w$$

$$= \frac{1}{2} y^T y + \sum_{i=1}^d \frac{1}{2} (\lambda + 1) w_i^2 - y^T x_{:,i} w_i$$

$$\frac{\partial J_{\lambda}(w)}{\partial w_i} = 0 \Rightarrow (\lambda + 1) w_i - y^T x_{:,i} = 0 \Rightarrow$$

$$w_i = \frac{y^T x_{:,i}}{\lambda + 1}$$

زمانی که  $y^T x_{:,i} \neq 0$  یا  $\lambda \rightarrow \infty$ ،  $w_i$  صفر می شود.  $y^T x_{:,i}$  درست مانیت  $\lambda$  نیز

در عمل به نسبت بزرگ این در این روش احتمال آنکه  $w_i$  ها صفر (sparse) نشوند، کمتر است.



(۳) الف -  $\epsilon_i < 0 \Rightarrow y^{(i)} (w^T x^{(i)} + b) \geq 1 - \epsilon_i$

معادله است با

$$y^{(i)} (w^T x^{(i)} + b) \geq 1$$

در نتیجه چون می‌خواهیم  $\min \frac{1}{2} \|w\|^2 + \frac{c}{2} \sum_{i=1}^m \epsilon_i^2$  را حل کنیم این عبارت سعی می‌کند  $\epsilon$  ها را صفر کند و به ازای  $\epsilon$  ها منفی آن ها را دقیقاً صفر می‌کند و همان طور که نشأت را داریم با constraint ها مسئله در نا قضا نمی‌شود بنابراین  $\epsilon$  را نگذاریم خود به خود  $\epsilon$  می‌شود.

ب -

$$L(w, \alpha, \epsilon, \alpha) = \frac{1}{2} \|w\|^2 + \frac{c}{2} \sum_{i=1}^m \epsilon_i^2 - \sum_{i=1}^m \alpha_i [y^{(i)} (w^T x^{(i)} + b) - 1 + \epsilon_i]$$

ج -

شرایط KKT:

$$\nabla_w L = w - \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} = 0 \Rightarrow w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$$

$$\nabla_b L = - \sum_{i=1}^m \alpha_i y^{(i)} = 0$$

$$\nabla_{\epsilon} L = c \epsilon - \alpha = 0 \Rightarrow c \epsilon_i = \alpha_i$$

$$\alpha_i [y^{(i)} (w^T x^{(i)} + b) - 1 + \epsilon_i] = 0$$

$$y^{(i)} (w^T x^{(i)} + b) \geq 1 - \epsilon_i$$

$$\alpha_i \geq 0$$



به نام خدا

پارسا بالینیان - ۱۳۰۱۰۰۸۶۵

$$L = \frac{1}{\mu} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y^{(i)} y^{(j)} x^{(i)T} x^{(j)} + \quad (۴) \rightarrow$$

$$\frac{1}{\mu} \sum_{i=1}^m \alpha_i \xi_i - \sum_{i=1}^m \alpha_i \xi_i + \sum_{i=1}^m \alpha_i - \sum_{i=1}^m \alpha_i y^{(i)} b$$

$$- \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y^{(i)} y^{(j)} x^{(i)T} x^{(j)} =$$

$$\sum_{i=1}^m \alpha_i - \frac{1}{\mu} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y^{(i)} y^{(j)} x^{(i)T} x^{(j)} - \frac{1}{\mu} \sum_{i=1}^m \frac{\alpha_i^2}{C}$$

مسئله dual به صورت زیر تعریف می شود :

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{\mu} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y^{(i)} y^{(j)} x^{(i)T} x^{(j)} - \frac{1}{\mu} \sum_{i=1}^m \frac{\alpha_i^2}{C}$$

$$\text{s.t.} \quad \begin{cases} \alpha_i \geq 0 \\ \sum_{i=1}^m \alpha_i y^{(i)} = 0 \end{cases}$$



$$\nabla_{\beta} J = 0 \Rightarrow -A^T \gamma w (\gamma - A\beta) = 0 \Rightarrow$$

$$A^T \omega Y - A^T \omega A \beta = 0 \Rightarrow \beta = (A^T \omega A)^{-1} A^T \omega Y$$

$$A^T \sim A \in \mathbb{R}^{p \times p}$$

جواب منصرم  $\Rightarrow$  اگر همه  $\lambda$ ها با هم برابر نباشند  
بفرمایید زیرا  $A^T w A$  دارد و ندارد.  
 $\Rightarrow$  " " " " " "  
نسب " " " " " "  
ندارد.

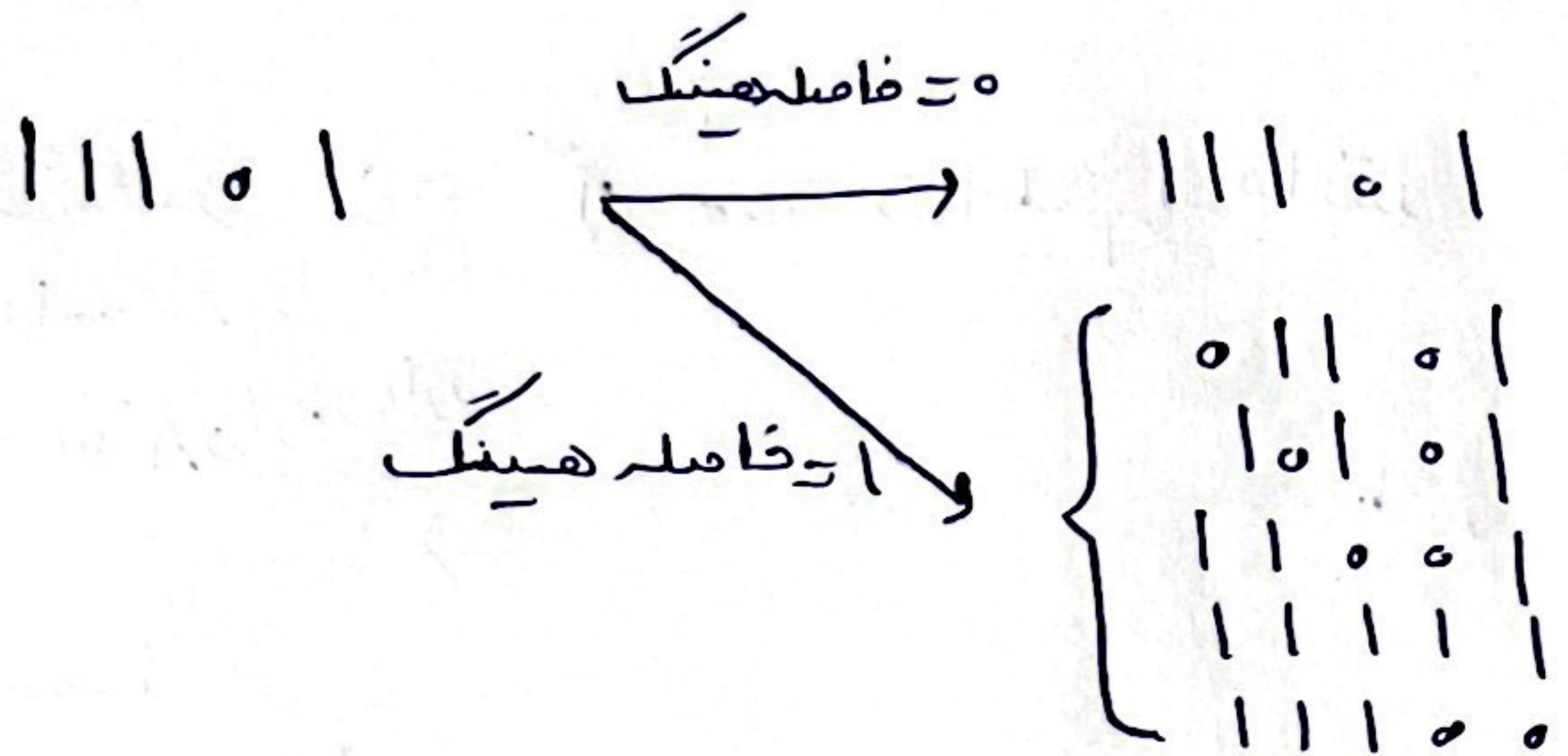
$$\text{Rank}(A^T w A) = \begin{cases} 2 & \text{اگر همه } \lambda\text{ها با هم برابر نباشند} \\ 1 & \text{وگرنه باشد} \end{cases}$$

$$\beta(t+1) = \beta(t) - \alpha \nabla_{\beta} J(t) = \beta(t) + \alpha (A^T y_w (y - \bar{A}_{\beta(t)}))$$

$\alpha$ : step-size

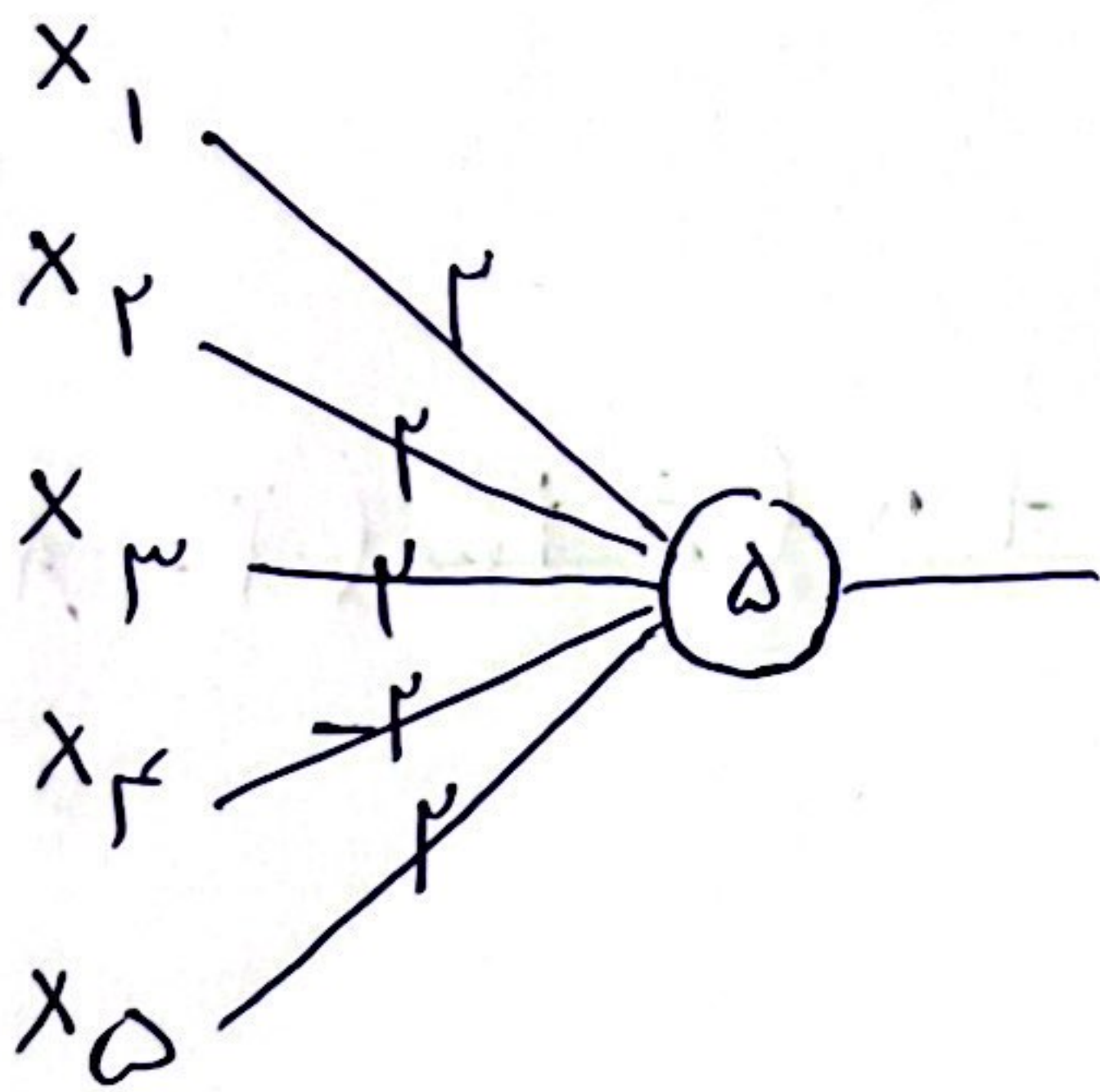


(۶) الف - ۵ نورون ورودی (برای ۵ بیت) و ۱ نورون خروجی (اگر بخواهیم خروجی را با بیز کد کنیم در غیر این صورت می توانیم ۲ نورون خروجی داشته باشیم و خروجی را به صورت one-hot کد کنیم)



ب - بله می توان - زیرا می توان مسئله را به صورت روبه رو مدل کرد و شبکه یک لایه MLP آن در صفت بعدی رسم شده است.

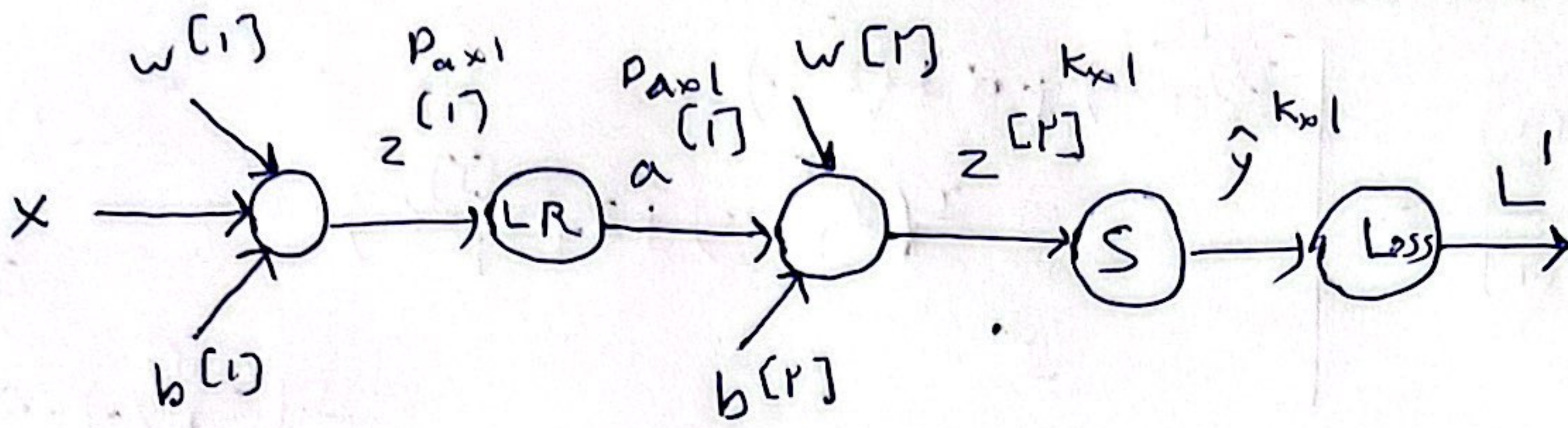
ج -





به نام خدا

پارسا یا لفرمان - ۱۴۰۱/۰۸/۰۵



الف -

$$w^{[2]} \in \mathbb{R}^{K \times D_a}, \quad b^{[2]} \in \mathbb{R}^{K \times 1}$$

$$z^{[1]} \in \mathbb{R}^{D_a \times m}$$

$$\frac{\partial \hat{y}_K}{\partial z_K^{[2]}} = \frac{-e^{z_K^{[2]}} e^{z_K^{[2]}} + e^{z_K^{[2]}} \sum_{j=1}^K e^{z_j^{[2]}}}{\left( \sum_{j=1}^K e^{z_j^{[2]}} \right)^2} = \hat{y}_K - \hat{y}_K^2 \quad \text{ب-}$$

$$\hat{y}_K = \frac{e^{z_K^{[2]}}}{\sum_{j=1}^K e^{z_j^{[2]}}}$$

$$\frac{\partial \hat{y}_K}{\partial z_i^{[2]}} = \frac{-e^{z_i^{[2]}} e^{z_K^{[2]}}}{\left( \sum_{j=1}^K e^{z_j^{[2]}} \right)^2} = -\hat{y}_i \hat{y}_K \quad \text{ج-}$$

$$\frac{\partial L}{\partial z_i^{[2]}} = \frac{\partial \hat{y}_i}{\partial z_i^{[2]}} \frac{\partial L}{\partial \hat{y}_i} = \hat{y}_i (1 - \hat{y}_i) \frac{-y_i}{\hat{y}_i} = \hat{y}_i - 1 \quad \text{د-}$$

$$\frac{\partial L}{\partial z_i^{[2]}} = \frac{\partial \hat{y}_K}{\partial z_i^{[2]}} \frac{\partial L}{\partial \hat{y}_K} = -\hat{y}_i \hat{y}_K \frac{-y_K}{\hat{y}_K} = \hat{y}_i$$

$$\delta_i = \frac{\partial z^{[2]}}{\partial a^{[1]}} = w^{[2]T} \quad \text{ه-}$$



به نام خدا

پارسا بالیزیان - ۴۰۰۱۰۰۸۲۵

$$\delta_2 = \frac{\partial a^{(1)}}{\partial z^{(1)}} = \begin{bmatrix} I(z_1^{(1)} > 0) + \alpha I(z_1^{(1)} < 0) \\ I(z_2^{(1)} > 0) + \alpha I(z_2^{(1)} < 0) \\ \vdots \\ I(z_n^{(1)} > 0) + \alpha I(z_n^{(1)} < 0) \end{bmatrix}$$

(۷) و -

$$\frac{\partial L}{\partial x^{(1)}} = \underbrace{\frac{\partial z^{(1)}}{\partial w^{(1)}}}_{x^T} \underbrace{\frac{\partial a^{(1)}}{\partial z^{(1)}}}_{\delta_2} \underbrace{\frac{\partial z^{(2)}}{\partial a^{(2)}}}_{\delta_1} \underbrace{\frac{\partial \hat{y}}{\partial z^{(2)}}}_{\delta}$$

$$= \delta_2 \delta_1 \delta \cdot x^T$$

ح - مشکل اولیه Softmax در عمل این است که چون ممکن است  $z_j^{(2)}$  ها خیلی بزرگ باشند در نتیجه  $e^{z_j^{(2)}}$  خیلی خیلی بزرگ می شود و overflow رخ می دهد در نتیجه یک مقدار DC که برابر بزرگ مقدار  $z_j^{(2)}$  ها است از همه مقدار کم می شود تا  $\max_j e^{z_j^{(2)}} = 1$  شود و overflow اتفاق نیفتد.

دقت کنید:

$$\hat{y}_i = \frac{e^{z_i^{(2)} - m}}{\sum_{j=1}^K (e^{z_j^{(2)} - m})} = \frac{e^{z_i^{(2)}}}{e^{-m} \sum_{j=1}^K e^{z_j^{(2)}}}$$

در نتیجه ما ساختار مکرر را تغییر ندادیم و می توانیم از  $\hat{y}_i$  استفاده کنیم.

$$= \frac{e^{z_i^{(2)}}}{\sum_{j=1}^K e^{z_j^{(2)}} = \hat{y}_i$$

قبل