

الف) softmax بصورت smooth، بزرگترین احتمال را برجسته می‌کند، همچنین نسبت به تابع max که فاقد مشتق است، جهت‌دار است زیرا مشتق پذیر است. softmax مقادیر متقی را هم هندل می‌کند زیرا e تابع نمایی است و e بتوان عدد متقی باز هم مثبت است.

ب) واریانس بالا به صورت کلی یعنی مدل پیچیده، وقتی که دیتامون نسبت به پیچیدگی مدل کم باشد هر این طریقی بالا باعث overfitting می‌شود. در واقع وقتی فضای فرضیه H بزرگی داشته باشد یعنی مدل پیچیده واریانس زیاد است. به روش برابر کاهش واریانس در مدل، کم کردن پیچیدگی مدل است.

ب) زیرا رگرسیون Lasso می‌خواهد تعدادی از feature ها را حذف کند و با این کار sparsity را زیاد کند، در صورتی که اگر تمام ویژگی‌ها تا حد خوبی با خروجی مرتبط باشند، باعث خطای بزرگی می‌شود ولی Ridge، ویژگی‌ها را حذف نمی‌کند (مقتول از حذف ویژگی، در نظر گرفتن وزن صفر برابر آن است). بلکه متناسب با اهمیت ^{ویژگی‌ها} به آنها وزن کمتر یا بیشتر می‌دهد و وقتی که همای ویژگی‌ها تا حد خوبی با خروجی مرتبط باشند وزن ^{ویژگی‌ها} می‌آلف صفر به آنها می‌دهد و در نتیجه هیچ‌کدام از ویژگی‌ها آلیکوت نمی‌شوند و مدل جهت عمل می‌کند.

ت) L_2 Regularization باعث می‌شود ضرایب مودن تا حد کوچک شوند. باعث کوچک شدن فضای فرضیه و کمتر شدن variance می‌شود (چون وزن‌ها بزرگ حذف می‌شوند). فضای فرضیه کوچک می‌شود (از صلفی چون فضای فرضیه مودن کوچک شده، ^{مایل} مودن زیاد می‌شود). در کل باعث کاهش واریانس و افزایش بایلیس می‌شود.

ادامه ی ب)

در واقع واریانس بالا نشانه ی اورفیت شدن و وابسته شدن زیاد مدل به دیتا ها است و مدل تعمیم پذیری خوبی نخواهد داشت

۲- تابع loss مان را mse در نظر می گیریم، داریم که:

$$y_k = w^T x_k = x_k^T w \Rightarrow y = X^T w$$

$$Loss = \frac{1}{n} \| y - X^T w \|^2$$

حالا برابر w_j داریم که:

$$SSE(w_j) = \| y - w_j x_j^T \|^2 = (y - w_j x_j^T)^T (y - w_j x_j^T)$$

$$\rightarrow \frac{\partial SSE(w_j)}{\partial w_j} = -2 x_j (y - w_j x_j^T) = 0 \rightarrow x_j y - w_j x_j x_j^T = 0$$

$$\rightarrow w_j = \frac{x_j y}{x_j x_j^T}$$

ب) به کمک بعضی سازی روی یک $Loss = \frac{1}{n} \| y - X^T w \|^2$ باید لیست را کمینه کنیم

$$\nabla_w Loss = \frac{1}{n} \nabla_w ((y - X^T w)^T (y - X^T w)) = 0 \rightarrow X(y - X^T w) = 0$$

$$\rightarrow Xy - XX^T w = 0 \rightarrow Xy = (XX^T)w$$

طبق گفته سوال و کثرت: سطرهای X برهم عمود اند ← داریم که:

$$XX^T = \begin{bmatrix} x_1^T x_1^T & \text{صفر} & \dots & \text{صفر} \\ \vdots & \vdots & \ddots & \vdots \\ \text{صفر} & \dots & \dots & x_n^T x_n^T \end{bmatrix} \rightarrow (XX^T)^{-1} = \begin{bmatrix} \frac{1}{x_1^T x_1^T} & \text{صفر} & \dots & \text{صفر} \\ \vdots & \vdots & \ddots & \vdots \\ \text{صفر} & \dots & \dots & \frac{1}{x_n^T x_n^T} \end{bmatrix}$$

$$\rightarrow w = (XX^T)^{-1} Xy = \begin{bmatrix} \frac{1}{x_1^T x_1^T} & \dots & \dots & \frac{1}{x_n^T x_n^T} \end{bmatrix} \begin{bmatrix} x_1 y \\ \vdots \\ x_n y \end{bmatrix} = \begin{bmatrix} \frac{x_1 y}{x_1^T x_1^T} \\ \vdots \\ \frac{x_n y}{x_n^T x_n^T} \end{bmatrix} \rightarrow w_j = \frac{x_j y}{x_j^T x_j^T}$$

← w_j هار حامل از بجهينه سازي روی تک متغير؛ w_0 هار حامل از بجهينه سازي روی تک برابر است.

(ب) در این مورد باید در نظر بگیریم:

$$\text{Loss}(w_j, w_0) = \text{SSE}(w_j, w_0) = \|y - w_j X_j^T - \tilde{w}_0\|_2^2$$

برابر هر یک از عناصر w_0 با تمام عناصر w_j

$$\rightarrow \frac{\partial \text{Loss}}{\partial w_j} = \frac{\partial}{\partial w_j} ((y - w_j X_j^T - \tilde{w}_0)^T (y - w_j X_j^T - \tilde{w}_0)) =$$

$$X_j (y - w_j X_j^T - \tilde{w}_0) = 0 \rightarrow X_j y - w_j X_j X_j^T - X_j \tilde{w}_0 = 0$$

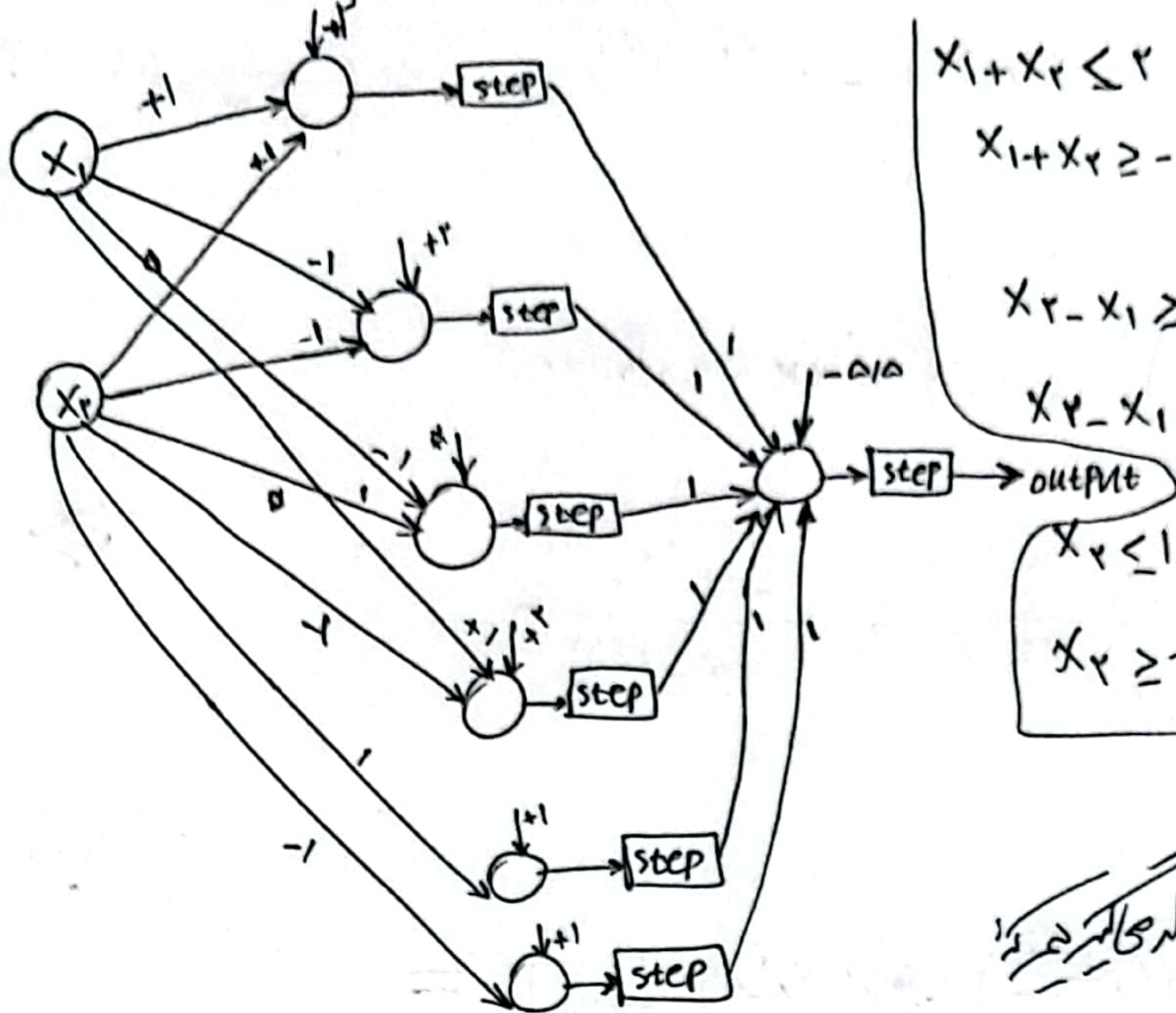
$$\rightarrow w_j X_j X_j^T = X_j y - X_j \tilde{w}_0 \rightarrow w_j = \frac{X_j y - X_j \tilde{w}_0}{X_j X_j^T}$$

$$\frac{\partial \text{Loss}}{\partial w_0} = \frac{\partial}{\partial w_0} ((y - w_j X_j^T - \tilde{w}_0)^T (y - w_j X_j^T - \tilde{w}_0)) = 0 \rightarrow$$

$$[1 \dots 1] (y - w_j X_j^T - \tilde{w}_0) = 0 \rightarrow 1^T y - w_j 1^T X_j^T - n w_0 = 0$$

$$\rightarrow n w_0 = 1^T y - w_j 1^T X_j^T \rightarrow w_0 = \frac{1^T y - w_j 1^T X_j^T}{n}$$

$$\rightarrow w_0 = \frac{\sum y_i - w_j \sum x_{ji}}{n}$$



$$x_1 + x_2 \leq 2$$

$$x_1 + x_2 \geq -2$$

$$x_2 - x_1 \geq -2$$

$$x_2 - x_1 \leq 2$$

$$x_2 \leq 1$$

$$x_2 \geq -1$$

step function مان را اينگونه در نقطه ها تعريف مي كنيم

$$f(x) = \begin{cases} 0 & x < 0 \\ 1 & x \geq 0 \end{cases}$$

$$i=k: \frac{\partial}{\partial z_i} \left(\frac{e^{z_i}}{\sum e^{z_j}} \right) = \frac{e^{z_i} (\sum e^{z_j}) - e^{z_i} e^{z_i}}{(\sum e^{z_j})^2} = \frac{e^{z_i} (\sum_{j \neq i} e^{z_j})}{(\sum e^{z_j})^2} \quad (1)$$

$$\frac{e^{z_i}}{\sum e^{z_j}} - \left(\frac{e^{z_i}}{\sum e^{z_j}} \right)^2 = \text{softmax} - \text{softmax}^2$$

$$i \neq k: \frac{\partial}{\partial z_i} \left(\frac{e^{z_k}}{\sum e^{z_j}} \right) = \frac{-e^{z_k} e^{z_i}}{(\sum e^{z_j})^2}$$

$$\frac{\partial L}{\partial z_i} = -y_i \frac{\partial}{\partial z_i} \log \left(\frac{e^{z_i}}{\sum e^{z_j}} \right) - \sum_{k \neq i} y_k \frac{\partial}{\partial z_i} \log \left(\frac{e^{z_k}}{\sum e^{z_j}} \right) =$$

$$-y_i \left(\frac{\partial}{\partial z_i} (z_i - \log \sum e^{z_j}) \right) - \sum_{k \neq i} y_k \frac{\partial}{\partial z_i} (z_k - \log \sum e^{z_j}) =$$

$$-y_i \left(1 - \frac{e^{z_i}}{\sum e^{z_j}} \right) - \sum_{k \neq i} y_k \left(- \frac{e^{z_i}}{\sum e^{z_j}} \right) = -y_i \left(1 - \frac{e^{z_i}}{\sum e^{z_j}} \right) + \sum_{k \neq i} y_k \left(\frac{e^{z_i}}{\sum e^{z_j}} \right)$$

$$= -y_i + y_i \frac{e^{z_i}}{\sum e^{z_j}} + \sum_{k \neq i} y_k \frac{e^{z_i}}{\sum e^{z_j}} = -y_i + \underbrace{\sum_k y_k}_{=1} \frac{e^{z_i}}{\sum e^{z_j}}$$

$$= -y_i + \frac{e^{z_i}}{\sum e^{z_j}}$$

۵- مدل رگرسیونی با خطا ε را در نظر بگیرید:

$$y = Xw + \varepsilon$$

و $\text{Var}(\varepsilon|X) = \sigma^2 I$, $E(\varepsilon|X) = 0$ ~~$E(\varepsilon) = 0$~~

در این صورت می دانیم که جواب مسئله در حالت حداقلی، ridge عبارت است از:

$$\hat{w}_{ls} = (X^T X)^{-1} X^T y$$

$$\hat{w}_{ridge} = (X^T X + \lambda I)^{-1} X^T y$$

حالا داریم، بایں عبار حساب کنیم:

$$E(\hat{w}_{ls}|X) = E((X^T X)^{-1} X^T (Xw + \varepsilon)|X) = (X^T X)^{-1} X^T X w$$

$$E(\hat{w}_{ridge}|X) = E((X^T X + \lambda I)^{-1} X^T (Xw + \varepsilon)|X) = (X^T X + \lambda I)^{-1} X^T X w$$

ماتریس کوواریانس نیز به صورت زیر است:

~~$$\text{Var}(\hat{w}_{ls}|X) = \text{Var}((X^T X)^{-1} X^T y|X) = E((X^T X)^{-1} X^T \varepsilon \varepsilon^T X (X^T X)^{-1}|X)$$~~

$$\text{Var}(\hat{w}_{ls}|X) = E((\hat{w}_{ls} - E(\hat{w}_{ls}|X))(\hat{w}_{ls} - E(\hat{w}_{ls}|X))^T | X) = E\left(\left((X^T X)^{-1} X^T y - (X^T X)^{-1} X^T X w\right)\left((X^T X)^{-1} X^T y - (X^T X)^{-1} X^T X w\right)^T | X\right)$$

$$= (X^T X)^{-1} X^T E(\varepsilon^2 | X) X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1}$$

$$\text{Var}(\hat{w}_{ridge}|X) = E\left((\hat{w}_{ridge} - E(\hat{w}_{ridge}|X))(\hat{w}_{ridge} - E(\hat{w}_{ridge}|X))^T | X\right) =$$

$$E\left(\left((X^T X + \lambda I)^{-1} X^T y - (X^T X + \lambda I)^{-1} X^T X w\right)\left((X^T X + \lambda I)^{-1} X^T y - (X^T X + \lambda I)^{-1} X^T X w\right)^T | X\right) = \sigma^2 (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1}$$

\downarrow
 $Xw + \varepsilon$

$$\hat{\beta} = X^T X (X^T X + \lambda I)^{-1}$$

حال قرار دهیم

$$\Rightarrow \text{Var}(\hat{w}_{\text{ridge}}|X) = \sigma^2 B^T (X^T X)^{-1} B$$

حالا باید ثابت کنیم $\text{Var}(\hat{w}_{\text{ls}}|X) - \text{Var}(\hat{w}_{\text{ridge}}|X)$ یک ماتریس مثبت معین است.

← اعضاء روی قطار مثبت و در نتیجه $\text{Var}_{\text{ridge}}$ کمتر از Var_{ls} است.

$$\text{Var}(\hat{w}_{\text{ls}}|X) - \text{Var}(\hat{w}_{\text{ridge}}|X) = \sigma^2 (X^T X)^{-1} - \sigma^2 B^T (X^T X)^{-1} B$$

$$= \sigma^2 \left(B^T (B^T)^{-1} (X^T X)^{-1} B^{-1} B - B^T (X^T X)^{-1} B \right) =$$

$$\sigma^2 B^T \left((B^T)^{-1} (X^T X)^{-1} B^{-1} - (X^T X)^{-1} \right) B =$$

$$= \sigma^2 B^T \left(2\lambda (X^T X)^{-2} + \lambda^2 (X^T X)^{-3} \right) B =$$

$$\sigma^2 (X^T X + \lambda I)^{-1} X^T X \left(2\lambda (X^T X)^{-2} + \lambda^2 (X^T X)^{-3} \right) X^T X (X^T X + \lambda I)^{-1} =$$

$$\sigma^2 (X^T X + \lambda I)^{-1} \left(2\lambda I + \lambda^2 (X^T X)^{-1} \right) (X^T X + \lambda I)^{-1}$$

حال بردار نامعقل را از دو طرف در این ماتریس ضرب کنیم، داریم:

$$\sigma^2 Z^T (X^T X + \lambda I)^{-1} (2\lambda I + \lambda^2 (X^T X)^{-1}) (X^T X + \lambda I)^{-1} Z$$

$$\sigma^2 Z^T (X^T X + \lambda I)^{-1} \left(2\lambda I + \lambda^2 (X^T X)^{-1} \right) \underbrace{(X^T X + \lambda I)^{-1} Z}_Z =$$

$$\sigma^2 Z^T (2\lambda I + \lambda^2 (X^T X)^{-1}) Z = 2\sigma^2 Z^T Z + \sigma^2 \lambda^2 Z^T (X^T X)^{-1} Z$$

که $X^T X$ یک ماتریس مثبت معین است $\Rightarrow V^T X^T X V$ نیز مثبت معین است

$\sigma^2 Z^T Z + \sigma^2 \lambda^r Z^T (X^T X)^{-1} Z$ is positive definite

~~Var_{LS} > Var_{ridge}~~

$Var_{LS} > Var_{ridge}$

$$tr\{Var[\hat{Y}(\lambda)]\} = \sigma^2 tr\{X(X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1} X^T\}$$

(ب)

$$X = U D V^T \rightarrow X^T X = V D^2 V^T \rightarrow X^T X + \lambda I = V D^2 V^T + \lambda I = V D^r V^T + \lambda V V^T$$

$$= V(D^r + \lambda I)V^T \rightarrow (X^T X + \lambda I)^{-1} = V(D^r + \lambda I)^{-1}V^T$$

$$tr\{Var[\hat{Y}(\lambda)]\} = \sigma^2 tr\left\{V(D^r + \lambda I)^{-1} \underbrace{V^T V}_I \underbrace{V D^2 V^T}_I (D^r + \lambda I)^{-1} V^T\right\} =$$

$$\sigma^2 tr\left\{V(D^r + \lambda I)^{-1} D^r (D^r + \lambda I)^{-1} V^T\right\} = \sigma^2 tr\left\{\underbrace{V^T V}_I (D^r + \lambda I)^{-1} D^r (D^r + \lambda I)^{-1}\right\}$$

$$= \sigma^2 tr\{(D^r + \lambda I)^{-1} D^r\} = \sigma^2 \sum D_{jj}^r (D_{jj}^r + \lambda)^{-1}$$