

سؤال ۱:

Part 1: گزینه C. بزرگ بودن ضریب به معنای مهم بودن و تأثیر زیاد آن روی مدل است ولی ممکن است با یک ویژگی دیگر وابستگی بالایی داشته باشد و مثلاً آن ویژگی بتواند این را خنثی کند و تأثیر آن را از مدل از بین ببرد. بنابراین در مورد اهمیت در مدل و امکان نگهداری یا حذف آن نمی توان نظر قطعی داد.

Part 2:

- (a) درست، مثلاً زمانی که مدل ما ساده است و under fit شده است با افزایش داده نمی توان این مشکل را برطرف نمود.
- (b) نادرست، مثلاً ممکن است مدل پیچیده باشد و ما به شدت overfit شویم و خطای train کم شود ولی خطای تست به شدت زیاد شود.
- (c) نادرست، خطای train کم می شود ولی اگر داده برای train زیاد باشد و حتی اگر خود مدل هم پیچیده باشد خطای test نیز کم خواهد شد.
- (d) نادرست، ممکن است مدل پیچیده تر باشد و نیاز شود از یک چندجمله ای با درجه ی بالاتر استفاده نمود.

$$L(w) = \frac{1}{p} \sum_{i=1}^n (x^{(i)T} w - y^{(i)})^2 = (Xw - Y)^T (Xw - Y) \quad \text{سوال ۲:} \quad (a)$$

$$\Rightarrow L(w) = w^T X^T X w - w^T X^T Y - Y^T X w - Y^T Y$$

$$\frac{\partial L}{\partial w} = Y X^T X w - X^T Y - X^T Y = 0 \Rightarrow Y X^T X w = Y X^T Y \Rightarrow$$

$$\Rightarrow X^T X w = X^T Y \Rightarrow \underline{w^* = (X^T X)^{-1} X^T Y}$$

$$L(w) = (Xw - Y)^T (Xw - Y) + \lambda w^T w \quad (b)$$

$$\frac{\partial L}{\partial w} = Y X^T X w - Y X^T Y + 2\lambda w = 0 \Rightarrow (X^T X + \lambda I) w = X^T Y$$

$$\Rightarrow w^* = (X^T X + \lambda I)^{-1} X^T Y$$

(c) ① ابتدا رفت را اثبات می کنیم یعنی اگر $\sum X = XF$ و F وارون پذیر:

$$\sum X = XF \xRightarrow{T} X^T \Sigma^T = F^T X^T \xRightarrow{\Sigma^T = \Sigma} X^T \Sigma = F^T X^T \xRightarrow{\Sigma^{-1}} X^T \Sigma^{-1}$$

$$X^T = F^T X^T \Sigma^{-1} \xRightarrow{F^{-T} X^T} F^{-T} X^T = X^T \Sigma^{-1} *$$

در w_{new}^* به جای $X^T \Sigma^{-1}$ عبارت $F^{-T} X^T$ می گذاریم:

$$w_{new}^* = (F^{-T} X^T X)^{-1} F^{-T} X^T Y = (X^T X)^{-1} F^T X F^{-T} X^T Y = \underbrace{(X^T X)^{-1} X^T Y}_{w_{opt}^*}$$

$$\text{فرض: } (X^T X)^{-1} X^T Y = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y$$

② حال برگشت را اثبات می کنیم:

$$\Rightarrow [(X^T X)^{-1} X^T Y - (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y] Y = 0$$

چون Y لاینل داده ها است و برای هر دیتاست این تساوی برقرار است Y می تواند در فضای بوج (null space)

ماتریس سمت چپ باشد بنابراین ماتریس سمت چپ صفر است یعنی:

$$(X^T X)^{-1} X^T = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1}$$

ادامه اثبات صفحه بعد:

$$(X^T X)^{-1} X^T = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1}$$

$$\underbrace{(X^T \Sigma^{-1} X)}_A (X^T X)^{-1} X^T = X^T \Sigma^{-1}$$

ماتریس A حاصل ضرب ۲ ماتریس وارون پذیر است پس وارون پذیر است.

$$A X^T = X^T \Sigma^{-1} \xrightarrow{\times \Sigma} A X^T \Sigma = X^T \xrightarrow{\text{ترنسپوز}} \Sigma^T X A^T = X \Rightarrow$$

Σ متقارن

$$\Rightarrow \Sigma X A^T = X \Rightarrow \Sigma X = X \underbrace{A^{-T}}_F = X F$$

A وارون پذیر

پس ثابت کردیم F وجود دارد و همان A^{-T} است.

(d) در سؤال خواسته شده نشان دهیم که $\lambda_1 \|w\|_2^2$ را می توان با اضافه کردن چند داده حذف کرد.

$$Y' = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}_{m \times 1}$$

Y' و X' به شکل زیر تعریف می کنیم با فرض m بعدی بودن w:

$$X' = \begin{bmatrix} \sqrt{\lambda_1} & 0 & \dots & 0 \\ 0 & \sqrt{\lambda_1} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & \sqrt{\lambda_1} \end{bmatrix}_{m \times m} = \sqrt{\lambda_1} I \quad \Rightarrow \quad \|Y' - X'w\|_2^2 = \left\| \begin{bmatrix} \sqrt{\lambda_1} w_1 \\ \sqrt{\lambda_1} w_2 \\ \vdots \\ \sqrt{\lambda_1} w_m \end{bmatrix} \right\|_2^2 =$$

$$= \lambda_1 \|w\|_2^2$$

بنابراین با اضافه کردن X' و Y' به X و Y داریم:

$$\text{new } X = \begin{bmatrix} X' \\ X \end{bmatrix}_{(m+n) \times m} \quad \text{new } Y = \begin{bmatrix} Y' \\ Y \end{bmatrix}_{(m+n) \times 1}$$

$$\begin{aligned} |(new X)w - new Y|^2 &= \left| \begin{bmatrix} X' \\ X \end{bmatrix} w - \begin{bmatrix} Y' \\ Y \end{bmatrix} \right|^2 = \|Xw - Y\|_2^2 + \|X'w - Y'\|_2^2 \\ &= \|Xw - Y\|_2^2 + \lambda_1 \|w\|_2^2 \end{aligned}$$

$$|Xw - Y|^2 + \lambda_1 \|w\|_2^2 + \lambda_2 \|w\|_1 = |new X w - new Y|^2 + \lambda_2 \|w\|_1 \quad \text{بنابراین:}$$

$$\text{minimize } \int P(n) \log \left(\frac{P(n)}{Q(n)} \right) dn = \int P(n) (\log P(n) - \log Q(n)) dn = \text{سوال ۳:}$$

$$KL = \int P(n) \left(\log P(n) - \log \frac{1}{(2\pi)^n} |\Sigma|^{-\frac{1}{2}} \exp \left(-\frac{1}{2} (n-\mu)^T \Sigma^{-1} (n-\mu) \right) \right) dn$$

$$\frac{\partial KL}{\partial \mu} \stackrel{\text{لینینز}}{=} \int \frac{\partial}{\partial \mu} P(n) \log \left(\frac{P(n)}{Q(n)} \right) dn = \int P(n) \left(\frac{\partial}{\partial \mu} \left(-\frac{1}{2} (n-\mu)^T \Sigma^{-1} (n-\mu) \right) \right) dn$$

$$= \int P(n) \times -\frac{1}{2} \times (-2) \Sigma^{-1} (n-\mu) dn = \int P(n) \Sigma^{-1} (n-\mu) dn = 0$$

$$\Rightarrow \Sigma^{-1} \int n P(n) dn = \Sigma^{-1} \int \mu P(n) dn = \Sigma^{-1} \mu \underbrace{\int P(n) dn}_1 = \Sigma^{-1} \mu$$

$$\Rightarrow \int n P(n) dn = \mu \Rightarrow \mu = E_P[X]$$

$$\frac{\partial KL}{\partial \Sigma} = \int P(n) \frac{\partial}{\partial \Sigma} \left(-\frac{1}{2} \log |\Sigma| + \frac{1}{2} (n-\mu)^T \Sigma^{-1} (n-\mu) \right)$$

$$|\Sigma| = \sum_{j=1}^n \underbrace{\sigma_{jj}}_{\text{عناصر اویان}} \underbrace{|\Sigma_{jj}|}_{\text{ماتریس مینور}} \quad \frac{\partial |\Sigma|}{\partial \Sigma} = \begin{bmatrix} \frac{\partial \Sigma_{11} |\Sigma_{jj}|}{\partial \sigma_{11}} & \dots & \frac{\partial \Sigma_{nj} |\Sigma_{jj}|}{\partial \sigma_{n1}} \\ \vdots & & \vdots \\ \frac{\partial \Sigma_{1n} |\Sigma_{jj}|}{\partial \sigma_{1n}} & \dots & \frac{\partial \Sigma_{nn} |\Sigma_{jj}|}{\partial \sigma_{nn}} \end{bmatrix}$$

$$= \begin{bmatrix} |\Sigma_{11}| & \dots & |\Sigma_{n1}| \\ \vdots & & \vdots \\ |\Sigma_{1n}| & \dots & |\Sigma_{nn}| \end{bmatrix} = \text{Adj}(\Sigma) = |\Sigma| \Sigma^{-1}$$

$$\frac{\partial \log |\Sigma|}{\partial \Sigma} = \frac{1}{|\Sigma|} \begin{bmatrix} |\Sigma_{11}| & \dots & |\Sigma_{n1}| \\ \vdots & & \vdots \\ |\Sigma_{1n}| & \dots & |\Sigma_{nn}| \end{bmatrix} = \Sigma^{-1} \Rightarrow$$

$$\frac{\partial}{\partial \Sigma} \log |\Sigma| = \Sigma^{-1}$$

در صفحہ بعد ہم مناسب مشتق بعدی می درازیم:

$$\frac{\partial}{\partial \Sigma} (x-\mu)^T \Sigma^{-1} (x-\mu) = -\Sigma^{-1} (x-\mu) (x-\mu)^T \Sigma^{-1}$$

$$\Rightarrow \int P(x) \frac{\partial}{\partial \Sigma} \log \frac{P(x)}{Q(x)} dx = \int P(x) \left(-\frac{1}{\gamma} \Sigma^{-1} + \frac{1}{\gamma} \Sigma^{-1} (x-\mu) (x-\mu)^T \Sigma^{-1} \right) dx = 0$$

$$\Rightarrow \frac{1}{\gamma} \Sigma^{-1} \underbrace{\int P(x) dx}_1 = \frac{1}{\gamma} \Sigma^{-1} \left(\int P(x) (x-\mu) (x-\mu)^T dx \right) \Sigma^{-1}$$

$$\Rightarrow \Sigma = \int (x-\mu) (x-\mu)^T P(x) dx = \text{Var}_P(X)$$

سوال ۳:
(a)

$$\frac{1}{2} \sum_{i=1}^n (w_j x_j^{(i)} - y^{(i)})^2 = L$$

$$\frac{\partial L}{\partial w_j} = \sum_{i=1}^n x_j^{(i)} (w_j x_j^{(i)} - y^{(i)}) = 0 \Rightarrow \sum_{i=1}^n w_j x_j^{(i)} = \sum_{i=1}^n y^{(i)} x_j^{(i)}$$

$$\Rightarrow w_j \times \sum_{i=1}^n x_j^{(i)} = \sum_{i=1}^n y^{(i)} x_j^{(i)} \Rightarrow w_j \times x_j^T x_j = x_j^T y$$

$$w_j = \frac{x_j^T y}{x_j^T x_j} \quad \text{که } x_j = \begin{bmatrix} x_j^{(1)} \\ \vdots \\ x_j^{(n)} \end{bmatrix} \text{ بنا بر این داریم}$$

(b) طبق سوال ۲ داریم که

$$W = (X^T X)^{-1} X^T y$$

$$w_{m \times 1} = [w_1, \dots, w_m]^T$$

$$X = \begin{bmatrix} x^{(1)} & \dots & x^{(n)} \end{bmatrix}^T_{n \times m}$$

حالا اگر X ماتریس orthonormal باشد داریم

$$W = X y \quad \Leftarrow \text{بنا بر این } X^T X = I$$

$$W = X y \Rightarrow \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix} = \begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix} \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} \Rightarrow w_j = x_j^T y$$

طبق بحثی قبلی داشتیم که $w_j = \frac{x_j^T y}{x_j^T x_j}$ چون ستون ها برهم عمودند $x_j^T x_j = 1$ بنا بر این

w_j درست آمده در این حالت برابر حالتی است که به صورت مستقل یادگیری می شد.

$$W := \begin{bmatrix} w_0 \\ w_j \end{bmatrix} \quad X_j := \begin{bmatrix} 1 \\ x_j^{(i)} \end{bmatrix}$$

feature نام از داده ها

(c)

$$X := \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_j^{(1)} & \dots & x_j^{(n)} \end{bmatrix}_{n \times h}$$

حال تابع L را به فرم ماتریسی بیان می کنیم:

$$L = \sum_{i=1}^n (w^T x_j^{(i)} - y^{(i)})^2 = (w^T X - y^T)^T (w^T X - y^T) = (x^T w - y)(w^T X - y^T) =$$

$$= x^T w w^T X - x^T w y^T - y w^T X + y y^T$$

$$\frac{\partial L}{\partial w} = 2 w^T X x^T - y^T x^T - y^T x^T = 0 \Rightarrow w^T X x^T = y^T x^T \Rightarrow$$

$$\Rightarrow X x^T w = X y \Rightarrow w = (X x^T)^{-1} X y$$

$$X X^T = \begin{bmatrix} 1 & \dots & 1 \\ x_j^{(1)} & \dots & x_j^{(n)} \end{bmatrix} \begin{bmatrix} 1 & x_j^{(1)} \\ \vdots & \vdots \\ 1 & x_j^{(n)} \end{bmatrix} = \begin{bmatrix} n & \sum x_j^{(i)} \\ \sum x_j^{(i)} & \sum x_j^{(i)2} \end{bmatrix} =$$

$$= n \times \begin{bmatrix} 1 & E[x_j] \\ E[x_j] & E[x_j^2] \end{bmatrix} \Rightarrow |X X^T| = n^2 (E[x_j^2] - E[x_j]^2) = n^2 \text{Var}(x_j)$$

$$\Rightarrow (X X^T)^{-1} \stackrel{\textcircled{1}}{=} \frac{1}{n \text{Var}(x_j)} \begin{bmatrix} E[x_j^2] & -E[x_j] \\ -E[x_j] & 1 \end{bmatrix}$$

$$\Rightarrow X Y = \begin{bmatrix} 1 & \dots & 1 \\ x_j^{(1)} & \dots & x_j^{(n)} \end{bmatrix} \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} = \begin{bmatrix} \sum y^{(i)} \\ \sum y^{(i)} x_j^{(i)} \end{bmatrix} \stackrel{\textcircled{2}}{=} n \times \begin{bmatrix} E[y] \\ E[y x_j] \end{bmatrix}$$

$$\stackrel{\textcircled{1} \times \textcircled{2}}{\Rightarrow} \begin{bmatrix} w_0 \\ w_j \end{bmatrix} = (X X^T)^{-1} X Y = \frac{1}{\text{Var}(x_j)} \begin{bmatrix} E[x_j^2] E[y] - E[x_j] E[x_j y] \\ -E[x_j] E[y] + E[x_j y] \end{bmatrix} =$$

$$\Rightarrow w_j = \frac{E[x_j y] - E[x_j] E[y]}{\text{Var}(x_j)} = \frac{\text{Cov}(x_j, y)}{\text{Var}(x_j)}$$

$$\frac{\partial L}{\partial w} = \sum_{i=1}^n [1 \ x_j^{(i)}] (w_j x_j^{(i)} + w_0 - y^{(i)}) = 0 \Rightarrow$$

$$\sum_{i=1}^n -w_j x_j^{(i)} + y^{(i)} = \sum_{i=1}^n w_0 = n w_0 \Rightarrow$$

$$\Rightarrow w_0 = \frac{1}{n} \sum y^{(i)} - \frac{1}{n} \sum x_j^{(i)} w_j = E[y] - w_j E[x_j]$$

$$\Rightarrow w_0 = E[y] - w_j E[x_j]$$