

سوال اول

پارت اول: گزیده C صحیح است.

دقت کنید که زیاد بودن ضریب همبستگی از $feature$ ها به معنای مهم بودن آن است. (دقت کنید) اما با توجه به اینکه درباره توزیع داده ها $correlation$ آن ها را اطلاع ایم نمی توان درباره توزیع نظریات در $feature$ ها با یکدیگر $correlation$ زیاد می دارد. آن یکی فیتیله دارای ضریب بزرگ با علامت مخالف باشد. و این دو را از هم جدا می کند. با بر این درباره $retained$ کردن آن توان نظریات $(a, b) \Leftarrow$ بنابراین C صحیح است.

پارت دوم:

(a) درست است.

زمانی که Bias زیاد است به معنای آن است که مدل $underfitting$ دارد و می تواند توزیع را به خوبی مدل کند و افزایش تعداد $training$ و $test$ ها تأثیری در آن ندارد. مثلاً $f(x) = x^2 + x^3$ را بنویسیم. $Poly$ از درجه ۲ مدل کنیم. $\hat{f}(x) = ax + b \Leftarrow$ هر چه تعداد $data$ زیاد شود باز هم $\hat{f}(x) = ax + b$ می ماند $f(x) = x^2 + x^3$ را مدل کند.

(b) غلط است.

زیرا اگر مدل $overfit$ شود. آنگاه می توان آورده های $train$ را کم کرد اما نمی توان $test$ error را کاهش داد. $overfitted \text{ model} \rightarrow \text{increase complexity}$
 \downarrow
 $\text{still overfitted with lower training error but higher test error.}$

(c) غلط است.

اگر مدل $underfit$ باشد نیاز داریم تا $Complexity$ مدل را افزایش دهیم. و این افزایش به کاهش $training$ و $test$ error منجر می شود. $underfitted \text{ model} \rightarrow \text{increase complexity} \rightarrow \begin{cases} \text{decrease training error} \\ \text{decrease test error} \end{cases}$

(d) غلط است. زمانی که $training$ error زیاد است. بدان معناست که باید $complexity$ مدل افزایش یابد. $Bias$ کاهش یابد باید از یک $Polynomial$ regression پیچیده تر (باز درجه بیشتر) استفاده کرد. $Linear$ regression (دقت کنید که $\Phi(x)$ را می بینیم به $Poly$ با درجه بیشتر می رسم تا ثابت زمان بگذرد).

a) Residuals are $\epsilon = Y - X\beta$ وقت کثیره β همان سالت

$$\text{error} = \sum_N \epsilon_i^2 = [\epsilon_1 \dots \epsilon_n] \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix} = \epsilon^T \epsilon$$

We want to minimize $\epsilon^T \epsilon = (Y - X\beta)^T (Y - X\beta)$

$$\frac{\partial (Y - X\beta)^T (Y - X\beta)}{\partial \beta} = -2X^T(Y - X\beta) \Rightarrow X^T(Y - X\beta) = 0 \Rightarrow X^T X \beta = X^T Y$$

$$\Rightarrow \beta = (X^T X)^{-1} X^T Y$$

b) with L_2 regularization σ error = $\sum_N \epsilon_i^2 + \lambda \sum \beta_i^2$

$$\text{error} = \epsilon^T \epsilon + \lambda \beta^T \beta \Rightarrow (Y - X\beta)^T (Y - X\beta) + \frac{\lambda}{2} \beta^T \beta$$

$$\frac{\partial \text{error}}{\partial \beta} = -2X^T(Y - X\beta) + 2\beta \frac{\lambda}{2} = 0 \Rightarrow X^T X \beta + \lambda \beta = X^T Y$$

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

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

$$\text{if } \lambda = 0 \rightarrow w^* = (X^T X)^{-1} X^T Y$$

زمان کثیره $X^T X = \Sigma^T X$ است حال باریم \leftarrow

$$c) \Sigma X = XF \Rightarrow X^T \Sigma^T = F^T X^T \Rightarrow X^T \Sigma = F^T X^T$$

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

حال در w_{new}^* عبارت بالا را $X^T \Sigma^{-1}$ با X^T جایگزین می کنیم.

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

$$\Rightarrow \text{if } \Sigma X = XF \text{ and } F \text{ not singular} \rightarrow w_{\text{new}}^* = w_{\text{old}}^* \rightarrow$$

$$\text{if } \Sigma X = XF \Rightarrow w_{\text{new}}^* = w_{\text{old}}^* \checkmark$$

$$W_{new}^* = W_{old}^*$$

ادامہ سوال ۲ C

غرض لکھو ←

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

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

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

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

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

ماتریس

$(X^T \Sigma^{-1} X)$ حاکم و پوزیٹو

$(X^T X)^{-1}$ پوزیٹو و حاکم

$[(X^T \Sigma^{-1} X) (X^T X)^{-1}]^{-1}$ حاکم و پوزیٹو

←

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

←

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

F ایک حاکم و پوزیٹو non singular ہے

if $W_{new}^* = W_{old}^* \Rightarrow$ we found F which is non singular matrix and

بنا براس

$$\Sigma X = X F \checkmark$$

if $\Sigma X = X F$
F non singular

$$\Rightarrow W_{new}^* = W_{old}^* \Rightarrow$$

$$W_{old}^* = W_{new}^* \iff \Sigma X = X F$$

F is ^{non} singular.

Proof Done ✓

در صفحہ بعدی ایک توضیح
(فائدہ ی دھم)

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

دست افند گفته شده:

$$\underbrace{[(X^T X)^{-1} X^T - (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1}]}_C Y = 0 \Rightarrow CY = 0 \quad \leftarrow$$

دست کنند به ازای هر Y ، رابطه بالا برقرار است (زیرا Y یک توزیع آمار دارد)

$$\forall Y \quad CY = 0 \Rightarrow \boxed{C = 0} \quad \leftarrow$$

مسئله سوال ۱۲

آزادی و دارای ایجاد λ_1, λ_2

باید به loss یک عبارت $\lambda_1 w_i^2$ اضافه کرد.

برای اینکار \leftarrow به Y سطریهای $0, 1$ اضافه
چک کنیم و آنرا معرفی کردیم.

برای ستون X هر بار:

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

$$X_w = \begin{bmatrix} \sqrt{\lambda_1} & & & \\ & \sqrt{\lambda_2} & & \\ & & \ddots & \\ & & & \sqrt{\lambda_n} \end{bmatrix}_{D \times D}$$

$$\Rightarrow X' = \begin{bmatrix} X_w \\ X \end{bmatrix} = \begin{bmatrix} \sqrt{\lambda_1} & \dots & \sqrt{\lambda_n} \\ & & X \end{bmatrix}_{D \times D}$$

$$\text{error} = \sum |Y_i' - X_i' w|^2 + \lambda_2 \sum |w_i|$$

حال داریم:

$$= \text{error} = \sum_i |Y_i' - X_i' w|^2 + \sum_{D+1}^{\infty} |Y_i' - X_i' w|^2 + \lambda_2 \sum |w_i|$$

$$\text{error} = \underbrace{\sum |X_i' w|^2}_{\sum (\sqrt{\lambda_i} w_i)^2} + \sum |Y_i - X_i w|^2 + \lambda_2 \sum |w_i| \Rightarrow \sum \lambda_i w_i^2 + \sum |Y_i - X_i w|^2 + \lambda_2 \sum |w_i|$$

$$\text{error} = (Y - Xw)^T (Y - Xw) + \lambda_1 w^T w + \lambda_2 \sum |w_i|$$

$$\text{error} = (Y' - X'w)^T (Y' - X'w) + \lambda_2 \sum |w_i|$$

$$f(x|\mu, \Sigma) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right\}$$

$$\text{minimize KL} = \text{minimize} \int p(x) \log p(x) dx - \int p(x) \log q(x) dx$$

independent of μ and Σ

$$\text{minimization of KL} = \text{maximization} \int p(x) \log q(x) dx = \text{maximize}$$

$$\Psi = \int p(x) \left[\log\left(\frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}}\right) - \frac{1}{2} \exp\left((x-\mu)^T \Sigma^{-1} (x-\mu)\right) \right] dx = \Psi$$

$$\text{maximize for } \mu \Rightarrow \frac{\partial \Psi}{\partial \mu} = -\frac{1}{2} \int 2 \Sigma^{-1} (x-\mu) p(x) dx = 0$$

$$\Rightarrow \int (x-\mu) p(x) dx = 0 \Rightarrow \int x p(x) dx = \mu \int p(x) dx \Rightarrow \boxed{\mu = E_p(x)}$$



$$\Sigma^T = \Sigma, \frac{\partial}{\partial A} Y^T A Y = Y Y^T, \frac{1}{|A|} = |A^{-1}| \leftarrow \text{اینها تو جبهه کنی از مرتبه اولی دانستم}$$

$$\frac{\partial}{\partial A} \log |A| = (A^{-1})^T \leftarrow \text{اینها = آخری سوال ۳}$$

کتاب اول
مست دوم

$$\text{KL maximization} = \text{maximize} \int p(x) \log\left(\frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}}\right) - \frac{1}{2} \exp\left((x-\mu)^T \Sigma^{-1} (x-\mu)\right) dx$$

$$\Rightarrow \frac{\partial \Psi}{\partial \Sigma} = \int p(x) \log(2\pi) dx + \frac{1}{2} \int p(x) \log |\Sigma^{-1}| - \frac{1}{2} \int p(x) (x-\mu)^T \Sigma^{-1} (x-\mu) dx$$

$$\frac{\partial \Psi}{\partial \Sigma^{-1}} = \int \frac{\partial}{\partial \Sigma^{-1}} p(x) \log p(x) dx + \frac{1}{2} \int p(x) ((\Sigma^{-1})^{-1})^T dx - \frac{1}{2} \int p(x) (x-\mu) (x-\mu)^T dx$$

$$\Rightarrow \Sigma \int p(x) dx = \int p(x) (x-\mu) (x-\mu)^T dx \Rightarrow \Sigma = \int (x-\mu) (x-\mu)^T p(x) dx$$

$$\Rightarrow \Sigma = E\{(x-\mu)(x-\mu)^T\} = \text{Var}_p(x)$$

ادرس دوم
به صورت کلی

است. maximum likelihood
ایستادگی که minimize کردن KL برای

اینده برای یک تابع convex حالت \rightarrow داریم \leftarrow
 $f(g(x)) \leq f(g(x_1))$

می دانیم \rightarrow کانوکس است بنابراین

$$KL(P||Q) = - \int P(x) \ln \frac{Q(x)}{P(x)} dx \geq - \ln \int P(x) \frac{Q(x)}{P(x)} dx \geq - \ln 1$$

$$KL(P||Q) = - \int P(x) \ln \frac{Q(x)}{P(x)} dx \geq 0$$

تابع convex است $\rightarrow KL(P||Q) \geq 0$ است \Leftarrow تابع $Q(x)$ - تدار اصفی است.
عبارت این تبارسانی KL صفر است که $Q(x) = P(x)$ باشد.

$$P(x) \approx Q(x|\theta)$$

حال چون توزیع $P(x)$ را نمی دانیم داریم
به تبارسانی \rightarrow است.

$$KL(P||Q) = \sum_N - \ln \left(\frac{Q(x_i|\theta)}{P(x_i)} \right) P(x_i) \Rightarrow \frac{\partial KL}{\partial \theta} = \sum_N \frac{\partial}{\partial P(x_i)} (-\ln Q(x_i|\theta) + \ln P(x_i))$$

$$\text{minimize } \sum_N - \ln Q(x|\theta) = \text{maximize } \sum_N \ln(Q|\theta) \Rightarrow \text{maximize } \prod_N Q(x|\theta)$$

که معادل \Leftarrow است \Leftarrow maximum likelihood
است \Leftarrow maximum likelihood

$$\mu = E_P(x)$$

$$\Sigma = \text{var}(x)$$

* از تخمین ML بایستی آگیزد.

$$\det(A) = \sum_j (-1)^{ij} A_{ij} M_{ij}$$

$$\frac{\partial \det(A)}{\partial A} = (\bar{A}^{-1})^T = C$$

$$\frac{\partial}{\partial A_{ij}} \log |A| = \frac{C_{ij}}{|A|} \leadsto \frac{\partial \log |A|}{\partial A_{ij}} = \frac{C}{|A|}$$

$$\bar{A}^{-1} = \frac{C^T}{|A|} \Rightarrow (\bar{A}^{-1})^T = \left(\frac{C}{|A|} \right) = (\bar{A}^{-1})^T \quad \checkmark$$

we apply only on col j^{th}

$$y_i = w_i x_{ij} + \varepsilon_i \sim \mathcal{N}(0, \Sigma) \quad \Sigma = \begin{bmatrix} \varepsilon_0 \\ \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix} \quad (a)$$

\rightarrow shows the col j^{th}

$$Y = w_i x_j + \varepsilon \Rightarrow \text{error} = \sum \varepsilon_i^2 = \varepsilon^T \varepsilon$$

$$\text{minimize } (Y - w_j x_j)^T (Y - w_j x_j) \Rightarrow \frac{\partial}{\partial w_j} = -2 X^T (Y - w_j x_j) = 0$$

$$X_j^T Y = w_j X_j^T X_j \Rightarrow \boxed{w_j = \frac{X_j^T Y}{X_j^T X_j}}$$

مشکلی که X_j نشان دهنده ستون j ام X است.

$$w^* = (X^T X)^{-1} X^T Y \quad (b)$$

$$\boxed{X^T X = I} \quad \text{orthogonal} \quad \text{از آنجا که در هر سوال گفته که ستون های } X \text{ orthogonal هستند}$$

$$w^* = X^T Y \Rightarrow w_j^* = (X^T Y)_j \Rightarrow w_j^* = X_j^T Y \quad \text{ستون } j \text{ ام} \quad (I)$$

$$w_j^* = \frac{X_j^T Y}{X_j^T X_j} \quad \text{train هر feature} \quad \text{حالت از قبیل استاندارد می کنیم} \leftarrow \text{دیده} \quad (II)$$

$$(I), (II) \Rightarrow w_j^* = X_j^T Y$$

اگر X ستون های orthogonal باشد نتیجه حاصل از train روی هر feature

باجت train روی کل feature ها یکی است.

در واقع چون X orthogonal است \Leftarrow Projection هر ستون روی ستون دیگر صفر است.

~~در واقع چون X orthogonal است \Leftarrow Projection هر ستون روی ستون دیگر صفر است.~~

$$\text{error} = \sum_{i=1}^N (y_i - w_j x_{ij} - w_0)^2 \Rightarrow \frac{\partial \text{error}}{\partial w_0} = 0$$

$$\frac{\partial \text{error}}{\partial w_0} = 0 \Rightarrow \sum_{i=1}^N (y_i - w_j x_{ij} - w_0) = 0 \Rightarrow \sum_{i=1}^N w_0 = \sum_{i=1}^N y_i - \sum_{i=1}^N w_j x_{ij}$$

$$\Rightarrow N w_0 = \sum_{i=1}^N y_i - w_j \sum_{i=1}^N x_{ij} \Rightarrow w_0 = \frac{\sum_{i=1}^N y_i}{N} - w_j \frac{\sum_{i=1}^N x_{ij}}{N} \Rightarrow$$

$$w_0 = E(y) - w_j E[X_j] = E(y) - w_j E(x_j) \checkmark \rightarrow \textcircled{I}$$

سَوْن زامات X_j

مَنَالِي سَوْن زام

$$\frac{\partial \text{error}}{\partial w_j} = \sum_{i=1}^N 2 x_{ij} (y_i - w_j x_{ij} - w_0) = 0 \Rightarrow \sum_{i=1}^N w_j x_{ij}^2 = \sum_{i=1}^N x_{ij} y_i - \sum_{i=1}^N w_0 x_{ij}$$

$$\Rightarrow w_j \sum_{i=1}^N x_{ij}^2 = \sum_{i=1}^N x_{ij} y_i - w_0 \times N \times E(x_j) \Rightarrow \textcircled{I} \text{ substed the } w_0 \text{ in this}$$

$$\Rightarrow w_j \sum_{i=1}^N x_{ij}^2 = \sum_{i=1}^N x_{ij} y_i - [E(y) - w_j E(x_j)] N E(x_j) \Rightarrow$$

$$w_j \left(\sum_{i=1}^N x_{ij}^2 - N E(x_j)^2 \right) = \sum_{i=1}^N x_{ij} y_i - N E(y) E(x_j) \Rightarrow$$

$$w_j \left(\frac{1}{N} \sum_{i=1}^N x_{ij}^2 - E(x_j)^2 \right) = \frac{1}{N} \sum_{i=1}^N x_{ij} y_i - E(y) E(x_j) \Rightarrow$$

$$\Rightarrow \frac{1}{N} \sum_{i=1}^N x_{ij}^2 = E(x_j^2) \Rightarrow E(x_j^2) - E(x_j)^2 = \text{Var}(x_j)$$

$$\frac{1}{N} \sum_{i=1}^N x_{ij} y_i = E(x_j y) \Rightarrow E(x_j y) - E(y) E(x_j) = \text{Cov}(x_j, y)$$

$$\Rightarrow w_j = \frac{\text{Cov}(x_j, y)}{\text{Var}(x_j)} \checkmark$$

Done.