

سؤال (1)

$$\|a\|^2 = a^T a$$

الف: اگر  $a$  یک بردار باشد داریم:

$$\Rightarrow \|y - Xw\|^2 = (y - Xw)^T (y - Xw) = (y^T - w^T X^T)(y - Xw)$$

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

$$\frac{\partial \|y - Xw\|^2}{\partial w} = -y^T X - (X^T y)^T + (X^T Xw)^T + w^T X^T X = -2y^T X + 2w^T X^T X = 0$$

$$\Rightarrow 2w^T X^T X = 2y^T X \Rightarrow X^T Xw = X^T y \Rightarrow w = (X^T X)^{-1} X^T y$$

ب:

$$\text{Ridge Regression: } \|y - Xw\|^2 + \lambda \|w\|^2$$

$$= (y^T y - y^T Xw - w^T X^T y + w^T X^T Xw) + (\lambda w^T w)$$

$$\frac{\partial \text{Ridge}}{\partial w} = \underbrace{(-2y^T X + 2w^T X^T X)}_{\text{ارضیت الف}} + 2\lambda w^T = 0 \Rightarrow w^T (X^T X + \lambda I) = y^T X$$

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

الف: مانند اضافه کردن  $m$  داده جدید به مجموعه داده‌ها است که  $m$  تعداد ویژگی‌ها است.

$$X_{n \times m} \Rightarrow \text{مجموعه داده قبلی} \quad X' = \begin{pmatrix} X_{n \times m} \\ \sqrt{\lambda} I_{m \times m} \end{pmatrix} \quad \text{مجموعه داده جدید}$$

$$Y_n \Rightarrow \text{مجموعه تگ قبلی} \quad Y' = \begin{pmatrix} Y_n \\ 0_m \end{pmatrix} \quad \text{مجموعه تگ جدید}$$

$$\|Y' - X'w\|^2 = \|Y' - Xw - \sqrt{\lambda} I_{m \times m} w\|^2 = \underbrace{\|Y - Xw\|^2 + \lambda \|w\|^2}_{L_2 \text{ with regularization}}$$

LSE with adding  
some data points

ب) اگر Correlation بین پارامترهای مرتبط، نامرتبط زیاد باشد،  $\lambda$  خیلی خوب عمل می‌کند. همچنین اگر وابستگی بین پارامترهای ورودی و خروجی زیاد باشد مشکل به وجود می‌آید.  
در حقیقت با  $\lambda$ ، تشخیص اینکه کدام پارامترها مهم‌تر هستند سخت‌تر است. پس اگر حجم داده کم باشد یا وابستگی پارامترها زیاد باشد، تشخیص حتی سخت‌تر هم می‌شود.

ج) کاملاً مناسب قسمت الف عمل می‌کنیم. یعنی  $\sqrt{\lambda_1} I_{m \times m}$  را اضافه کنیم.

$$X' = \begin{pmatrix} X_{n \times m} \\ \sqrt{\lambda_1} I_{m \times m} \end{pmatrix}, \quad Y' = \begin{pmatrix} Y_n \\ 0_m \end{pmatrix}$$

$$\|Y' - X'w\| + \lambda_2 \|w\|_1 = \|Y' - Xw - \sqrt{\lambda_1} I_{m \times m} w\| + \lambda_2 \|w\|_1$$

$$= \|Y - Xw\| + \lambda_1 \|w\|_2 + \lambda_2 \|w\|_1$$

سؤال ۳)  
الف:

$$\Phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \end{bmatrix}, w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix} \Rightarrow y = w^T \Phi(x) + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$\Rightarrow P(y | x_1, x_2) \sim \mathcal{N}(w^T \Phi(x), \sigma^2)$$

$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$  من  $P(y | x)$  است

ب: با فرض استقلال هر داده از داده‌ی دیگر داریم:

$$P(Y | X_1, X_2) = \prod_{i=1}^n P(y^i | x_1^i, x_2^i) = \prod_{i=1}^n \mathcal{N}(w^T \Phi(x^i), \sigma^2)$$

$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y^i - w^T \Phi(x^i))^2}{2\sigma^2}}$$

$$\Rightarrow \log(P(Y | X_1, X_2)) = \sum_{i=1}^n \log e^{-\frac{(y^i - w^T \Phi(x^i))^2}{2\sigma^2}} \times \log \frac{1}{\sqrt{2\pi\sigma^2}}$$

$$\log(P(Y | X_1, X_2)) \propto -\frac{1}{2\sigma^2} \sum_{i=1}^n (y^i - w^T \Phi(x^i))^2$$

با حذف مقادیر constant داریم:

ج: با توجه به منفی بودن عبارت بالا، بهینه‌ترین  $P(Y | X_1, X_2)$  معادل با کمینه کردن حاصل جمع بالا است.  
پس داریم:

$$f(w_0, w_1, w_2, w_3, w_4) = \sum_{i=1}^n (y^i - w^T \Phi(x^i))^2$$

$$\nabla_w f(w) = \nabla_w \left( \sum_{i=1}^n (y_i - w^T \phi(x_i))^2 \right) = -2 \sum_{i=1}^n (y_i - w^T \phi(x_i)) \phi(x_i)^T$$

$$w^{i+1} = w^i - a \nabla_w f(w) = w^i + a \sum_{i=1}^n (y_i - w^T \phi(x_i)) \phi(x_i)^T$$

که  $a$  ضریب یادگیری دهان (rate) است.

سوال ۴

$$Y = Xw + \epsilon \Rightarrow X^T Y = X^T X w + X^T \epsilon \Rightarrow X^T X w = X^T Y - X^T \epsilon$$

$$\Rightarrow w = (X^T X)^{-1} X^T Y - (X^T X)^{-1} X^T \epsilon = (X^T X)^{-1} X^T (Xw + \epsilon) - (X^T X)^{-1} X^T \epsilon$$

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

$$\epsilon \sim \mathcal{N}(0, I) \Rightarrow w \sim (w^*, \|(X^T X)^{-1} X^T\|^2)$$

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

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

$$\Rightarrow \hat{w} \sim \mathcal{N}(w^*, (X^T X)^{-1})$$

$$\text{Var}(Y) = \text{Var}(X\hat{w} + \epsilon) = \text{Var}(X\hat{w}) + \text{Var}(\epsilon)$$

$$= X \text{Var}(\hat{w}) X^T + \text{Var}(\epsilon) = X (X^T X)^{-1} X^T + I$$

$$\text{Var}(a+b) = \text{Var}(a) + \text{Var}(b)$$

$$\text{Var}(Xw) = X \text{Var}(w) X^T$$



$$X = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \Rightarrow A_n = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix} \rightarrow A_n = f(x, n)$$

:ج

پس  $A_n$  - صورت تابعی از  $x_i, n$  و  $x_i$  و  $n$  است

$$\sum_{i=1}^n \phi(x_i) \phi(x_i)^T = \sum_{i=1}^n \begin{bmatrix} 1 & x_i \\ x_i & x_i^2 \end{bmatrix} = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix} = A_n$$

$$\forall s: s^T A_n s = \sum_{i=1}^n s_i^T \phi(x_i) \phi(x_i)^T s_i = \sum_{i=1}^n (s_i^T \phi(x_i)) (s_i^T \phi(x_i))^T$$

$$= \sum_{i=1}^n \|s_i^T \phi(x_i)\|^2 \geq 0 \Rightarrow \text{نیمه مثبت}$$

:ج

$$J(X) = E_{x \sim Q} (\text{Var}(y(x))) = E_{x \sim Q} [x^T (X^T X)^{-1} x]$$

$$X^T X = A_n = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix} \Rightarrow (X^T X)^{-1} = \frac{1}{n \sum x_i^2 - (\sum x_i)^2} \begin{bmatrix} \sum x_i^2 - \sum x_i & -\sum x_i \\ -\sum x_i & n \end{bmatrix}$$

$$\Rightarrow E [x^T (X^T X)^{-1} x] = \frac{1}{n \sum x_i^2 - (\sum x_i)^2} * (\sum x_i^2 + 2 E(x) \sum x_i + E(x^2) n)$$

$$E[x] = 0, \text{Var}(x) = v^2$$

$$\Rightarrow J(X) = \frac{\sum x_i^2 + n v^2}{n \sum x_i^2 - (\sum x_i)^2}$$

الف:

$$Y = f(x, w) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \beta^{-1})$$

$$\Rightarrow P(Y | x, w, \beta) \sim \mathcal{N}(f(x, w), \beta^{-1})$$

با در نظر گرفتن استقلال بین  $\mu$  داده داریم:

$$P(Y | X, w, \beta) = \prod_{i=1}^n \mathcal{N}(f(x_i, w), \beta^{-1})$$

$$= \prod_{i=1}^n \left( \frac{\sqrt{\beta}}{\sqrt{2\pi}} e^{-\frac{(y_i - f(x_i, w))^2}{2\beta^{-1}}} \right) \xrightarrow{\log} -\frac{N}{2} \log(2\pi) + \frac{N}{2} \log(\beta) - \frac{\beta}{2} \sum_{i=1}^n (y_i - f(x_i, w))^2$$

دس بهینه شدن  $P(Y | X, w, \beta)$  معادل شد با کمینه شدن  $\sum_{i=1}^n (y_i - f(x_i, w))^2$  همان  $SSE$  است.

ب:

$$P(w | Y) \propto P(Y | w) P(w)$$

دس برای بهینه کردن  $P(w | Y)$ ، باید  $P(Y | w) P(w)$  را بهینه کنیم. با وگانه کردن از ۲ طرف، نتیجه می شود که  $\log P(Y | w) + \log P(w)$  باید بهینه شود.

$$P(w) \sim \mathcal{N}(0, \alpha^{-1}) = \frac{1}{\sqrt{K\pi}^K} e^{-\frac{1}{2} w^T \alpha I w} \Rightarrow \log P(w) \propto -\sum_{i=1}^K w_i^2$$

$$\log P(Y | w) \propto -\sum_{i=1}^n (y_i - f(x_i, w))^2$$

دس بهینه کردن  $P(w | Y)$  معادل است با کمینه کردن عبارت زیر:

$$\frac{\beta}{2} \sum_{i=1}^n (y_i - f(x_i, w))^2 + \frac{\alpha}{2} \sum_{i=1}^K w_i^2 = \frac{\beta}{2} SSE + \frac{\alpha}{2} \|w\|_2^2 \rightarrow SSE + \text{Ridge-Regularized}$$

$$P(w|Y) \propto P(Y|w)P(w)$$

$$P(w) \propto \frac{\alpha}{r} e^{-\alpha \|w\|} \Rightarrow \log P(w) \sim -\alpha \|w\|$$

ج:

$$MAP \sim -\frac{\beta}{2} \sum_{i=1}^n (y_i - f(x_i; w))^2 - \alpha \|w\|_1$$

از این جا به بعد مسئله قیمت ب محلی کنیم.

پس برای بهینه شدن  $P(w|Y)$  ، باید به است زیر مینه شود:

$$\frac{\beta}{2} \sum_{i=1}^n (y_i - f(x_i; w))^2 + \alpha \|w\|_1 \rightarrow \text{SSE with Lasso regularized}$$

سؤال ۵)

$$P(w) \propto e^{-\frac{1}{2} (w - M_0)^T S_0^{-1} (w - M_0)} = e^{-\frac{1}{2} w^T S_0^{-1} w + M_0^T S_0^{-1} w - \frac{1}{2} M_0^T S_0^{-1} M_0}$$

چون مستقل از  $w$  است  
مغولاً غریب می شود.

الف:

$$\Rightarrow P(w) \propto \exp(J^T w - \frac{1}{2} w^T P w), \quad J = S_0^{-1} M_0, \quad P = S_0^{-1}$$

$$P(w|D) \propto P(D|w) P(w) = \prod_{i=1}^n P(y_i | x_i; w) \times P(w)$$

$$= \left( \prod_{i=1}^n \exp\left(-\frac{\beta}{2} (y_i - w^T x_i)^2\right) \right) \times \exp(J^T w - \frac{1}{2} w^T P w)$$

قیمت های مستقل از  $w$  را در نظر می گیریم چون مغولاً غریب می شود.

$$\Rightarrow P(w|D) \propto \exp \left\{ (J + \beta \sum y_i x_i)^T w - \frac{1}{2} w^T (P + \beta \sum x_i x_i^T) w \right\}$$

$$P_N = S_N^{-1} = P + \beta \sum x_i x_i^T = S_0^{-1} + \beta X^T X \Rightarrow S_N^{-1} = S_0^{-1} + \beta X^T X \quad \checkmark$$

$$J_N = S_N^{-1} M_N = (J + \beta \sum y_i x_i)^T = S_0^{-1} M_0 + \beta X^T Y \Rightarrow M_N = S_N (S_0^{-1} M_0 + \beta X^T Y) \quad \checkmark$$



$$P(w) \sim N(0, \alpha^{-1} I) \Rightarrow M_0 = 0, S_0 = \alpha^{-1} I$$

ب :

$$P(Y|w) \sim N(w^T X, \beta^{-1})$$

طبقی این خواص داشت :

$$S_N^{-1} = S_0^{-1} + \beta X^T X = \alpha I + \beta X^T X$$

$$M_N = S_N (0 + \beta X^T Y) = S_N \times \beta X^T Y$$

$$\Rightarrow P(w|Y) \sim N(M_N, S_N), \quad S_N^{-1} = \alpha I + \beta X^T X, \quad M_N = \beta S_N X^T Y$$

ج : برای بدست آوردن  $M_N$  و  $S_N$  از Prior استفاده کردیم که  $M_0$  بود. ضمیمه کار را برای این قسمت هم انجام می دهیم. و چون  $X$  قطری است، پس می توانیم از  $X$  استفاده کنیم.

$$S_{N+1}^{-1} = S_N^{-1} + \beta x_{n+1} x_{n+1}^T$$

$$M_{N+1} = S_{N+1} (S_N^{-1} M_N + \beta x_{n+1} y_{n+1})$$

پس همان سؤال این است با این تفاوت که در مرحله  $N+1$ ،  $X$  همان  $x_{n+1}$  می شود،  $Y$  همان  $y_{n+1}$  می شود. و به جای  $M_0$  از  $M_N$  استفاده می کنیم چون Prior به روز رسانی می شود. ضمیمه ای ادامه از رابطه بدست آورده استفاده می کنیم :

$$S_{N+1}^{-1} = S_0^{-1} + \beta \sum_{i=1}^{n+1} x_i x_i^T$$

$$M_{N+1} = S_{N+1} (S_N^{-1} M_N + \beta x_{n+1} y_{n+1}) = S_{N+1} (S_N^{-1} (S_N^{-1} M_{N-1} + \beta x_n y_n) + \beta x_{n+1} y_{n+1})$$

$$\Rightarrow M_{N+1} = S_{N+1} (S_0^{-1} M_0 + \beta \sum_{i=1}^{n+1} x_i y_i)$$

پس مرحله به مرحله این چنین می شود

از آن جا که  $X^T Y = \sum x_i y_i$  و  $X^T X = \sum x_i x_i^T$ ، پس جواب معادل این است که صد داده ها یکجا تحلیل می شود.



$$S_{N+1}^{-1} = S_N^{-1} + \beta x_{n+1} x_{n+1}^T$$

$$\Rightarrow S_{N+1} = (S_N^{-1} + \beta x_{n+1} x_{n+1}^T)^{-1} = (S_N^{-1} + \sqrt{\beta} x_{n+1} \sqrt{\beta} x_{n+1}^T)^{-1}$$

$$= S_N - \frac{S_N (\sqrt{\beta} x_{n+1}) (\sqrt{\beta} x_{n+1}^T) S_N}{1 + (\sqrt{\beta} x_{n+1}^T) S_N (\sqrt{\beta} x_{n+1})} = S_N - \frac{\beta S_N x_{n+1} x_{n+1}^T S_N}{1 + \beta x_{n+1}^T S_N x_{n+1}}$$

$$\Rightarrow S_N - S_{N+1} = \frac{\beta S_N x_{n+1} x_{n+1}^T S_N}{1 + \beta x_{n+1}^T S_N x_{n+1}}$$

$$x^T (S_N - S_{N+1}) x = \frac{x^T S_N x_{n+1} x_{n+1}^T S_N x}{\frac{1}{\beta} + x_{n+1}^T S_N x_{n+1}} = \frac{(x^T S_N x_{n+1})^2}{\frac{1}{\beta} + x_{n+1}^T S_N x_{n+1}}$$

چون  $S_N$  مثبت معین است، پس مربع صحت مثبت است. معلوم کردیم که  $S_N$  مثبت معین است.

$$x^T (S_N - S_{N+1}) x \geq 0 \Rightarrow \sigma_N^2 - \sigma_{N+1}^2 \geq 0$$

پس داریم:

$$\Rightarrow \sigma_{N+1}^2 - \sigma_N^2 \leq 0 \Rightarrow \text{واریانس کاهش می یابد.}$$