1] Sales = 2.939 + 0.046×TV + 0.189×radio − 0.001 ×newspaper

   isn't associated with sales

   in the presence of TV & radio

---

3] $Y = 50 + 20 x_1 + 0.07 x_2 + 35 x_3 + 0.01 x_4 - 10 x_5$

a) i and iii can be true.

b) $50 + 20 \times 4 + 0.07 \times 110 + 35 \times 1 + 0.01 \times 4 \times 110 - 10 \times 4 \times 1 = 137$

   (female ↑ over the $35 \times 1$ term)

c) small coefficient for interaction GPA×IQ ($x_4$) means it has a bit association with sallary in the presence of other features.

---

5] $\hat{y}_i = \hat{\beta} x_i \Rightarrow \hat{y}_i = \dfrac{\sum\limits_{i=1}^{n} x_i y_i}{\sum x_i^2} \times x_i$

   داریم ضریب مشترک بین‌شان جبرآکره

   $= \sum \left( \dfrac{x_i' x_i}{\sum x_i'^2} \right) y_i'$

   $a_i'$

   $\Downarrow$

   $\boxed{\hat{y}_i = \sum_i a_i' y_i'}$

1] a) between different models, a model which has better fit to training data, has smallest training RSS.

b) if the data have large P (features), may overfit and. in this situation forward stepwise selection has smaller. test RSS.

c) i) True      ii) True      iii) false

    iv) false         v) True

2] a) iii → lasso regression can remove some features by zero coefficient and decrease flexibility (complexity)

b) iii → ridge is also act like lasso.

c) ii → non-linear methods can lead to increase in flexibility.

5] a) $(y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{12})^2 + (y_2 - \hat{\beta}_1 x_{21} - \hat{\beta}_2 x_{22})^2 + \lambda (\hat{\beta}_1^2 + \hat{\beta}_2^2)$ → ridge regression

b) $\hat{\beta}_1 = \hat{\beta}_2 \ (y_1^2 + \hat{\beta}_1^2 x_{11}^2 + \hat{\beta}_2^2 x_{12}^2 - 2\hat{\beta}_1 x_{11} y_1 - 2\hat{\beta}_2 x_{12} y_1 + 2\hat{\beta}_1 \hat{\beta}_2 x_{11} x_{12})$

$+ (y_2^2 + \hat{\beta}_1^2 x_{21}^2 + \hat{\beta}_2^2 x_{22}^2 - 2\hat{\beta}_1 x_{21} y_2 - 2\hat{\beta}_2 x_{22} y_2 + 2\hat{\beta}_1 \hat{\beta}_2 x_{21} x_{22})$

$+ \lambda \hat{\beta}_1^2 + \lambda \hat{\beta}_2^2$

$\xrightarrow{\text{minimization}}$ $\dfrac{\partial L(\beta)}{\partial \beta_1} = 0 \rightarrow (2\hat{\beta}_1 x_{11}^2 - 2x_{11}y_1 + 2\hat{\beta}_2 x_{11}x_{12}) +$

$$(2\hat{\beta}_1 x_{21}^2 - 2x_{21}y_2 + 2\hat{\beta}_2 x_{21}x_{22}) + 2\lambda\hat{\beta}_1 = 0$$

$\left.\begin{array}{l} x_{11} = x_{12} = x_0 \\[2mm] x_{21} = x_{22} = x_1 \end{array}\right\} \rightarrow (\hat{\beta}_1 x_0^2 - x_0 y_1 + \hat{\beta}_2 x_0^2) + (\hat{\beta}_1 x_1^2 - x_1 y_2 + \hat{\beta}_2 x_1^2) + \lambda\hat{\beta}_1 = 0$

$$\Rightarrow \hat{\beta}_1(x_0^2 + x_1^2) - x_0 y_1 - x_1 y_2 + \hat{\beta}_2(x_0^2 + x_1^2) + \lambda\hat{\beta}_1 = 0$$

$$\lambda\hat{\beta}_1 + \hat{\beta}_1(x_0^2 + x_1^2) + \hat{\beta}_2(x_0^2 + x_1^2) = x_0 y_1 + x_1 y_2$$

$\begin{array}{c} x_0 \quad x_1 \\ \boxed{x_{11}} \; \boxed{x_{21}} = 0 \end{array}$

$\hat{\beta}_1(x_0^2 + x_1^2 + \overset{(x_0+x_1)^2}{2x_0x_1} - 2x_0x_1) + \hat{\beta}_2\underbrace{(x_0^2 + x_1^2 + 2x_0x_1}_{} - 2x_0x_1) + \lambda\hat{\beta}_1$

$$\lambda\hat{\beta}_1 - 2\hat{\beta}_1 x_0 x_1 - 2\hat{\beta}_2 x_0 x_1 = x_0 y_1 + x_1 y_2$$

$$\Rightarrow \lambda\hat{\beta}_1 = 2\hat{\beta}_1 x_0 x_1 + 2\hat{\beta}_2 x_0 x_1 + x_0 y_1 + x_1 y_2$$

---

$\dfrac{\partial L(\beta)}{\partial \beta_2} \rightarrow (2\hat{\beta}_2 x_{12}^2 - 2x_{12}y_1 + 2\hat{\beta}_1 x_{11}x_{12}) +$

$$(2\hat{\beta}_2 x_{22}^2 - 2x_{22}y_2 + 2\hat{\beta}_1 x_{21}x_{22}) + 2\lambda\hat{\beta}_2 = 0$$

$\left.\begin{array}{l} x_{11} = x_{12} = x_0 \\[2mm] x_{21} = x_{22} = x_1 \end{array}\right\} \rightarrow (\hat{\beta}_2 x_0^2 - x_0 y_1 + \hat{\beta}_1 x_1^2) + (\hat{\beta}_2 x_1^2 - x_1 y_2 + \hat{\beta}_1 x_1^2) + \lambda\hat{\beta}_2 = 0$

$$\Rightarrow \lambda\hat{\beta}_2 = 2\hat{\beta}_1 x_0 x_1 + 2\hat{\beta}_2 x_0 x_1 + x_0 y_1 + x_1 y_2$$

$$\Rightarrow \lambda\beta_1 = \lambda\beta_2 \rightarrow \boxed{\beta_1 = \beta_2}$$

c) lasso:

should be minimize

$$(y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{12})^2 + (y_2 - \hat{\beta}_1 x_{21} - \hat{\beta}_2 x_{22})^2 + \lambda (|\hat{\beta}_1| + |\hat{\beta}_2|)$$

d) in lasso regression we have absolute for $\beta$ and can't do partial derivation. so other possible solutions like convex optimization may be useful.

a)

7] $D = \{x[1], \ldots, x[n]\} \rightarrow L(\sigma^2) = P(x[1]) \times \cdots \times P(x[n])$

likelihood

$$= \prod_{i=1}^{n} P(x[i])$$

$$= \left(\beta_0 + \sum_{d=1}^{p} x_{ij} \beta_j + \varepsilon_i\right)^n \rightarrow \ell(x) = \log L(M)$$

b) prior $\rightarrow P(\beta) = \frac{1}{2b} \exp(-|\beta|/b)$

$\rightarrow$ posterior = prior $\times$ Liklihood $= \frac{1}{2b} \exp(-|\beta|/b) \times \log\left(\beta_0 + \sum_{j=1}^{p} x_{ij} \beta_j\right)^n$

c) I can't solve this section.

d) posterior = prior $\times$ likelihood $= \left(\frac{1}{\log \sqrt{2\pi c}}\right)^n \exp\left(-\frac{1}{2c^2} \sum_{i=1}^{n} (x[i] - \hat{\mu})^2\right)$

$\mathcal{N}(0, c)$

$$\times \log\left(\beta_0 + \sum_{j=1}^{p} x_{ij} \beta_j\right)^n$$

e) I can't solve this.

**Problem 3 :** show that the ridge hat matrix is not a projection matrix.

projection matrix : → 1) idenpotent → $H^2 = H$

- the ridge hat matrix → proof : $H^2 \neq H$   $(H = x(x^Tx + \lambda I)^{-1}x^T)$

$$H^2 - M \neq 0$$

$$'X = UDV^T$$

$$H(M - I) \neq 0$$

$$M = UDV^T(VD^TU^T \cdot UDV^T + \lambda I)^{-1} VD^TU^T$$

$$UU^T \cdot DD^T(D^TD + \lambda I)^{-1}\left(\left[UU^T \cdot DD^T(D^TD + \lambda I)^{-1}\right] - I\right) \neq 0$$

$$UU^T \cdot DD^T(D^TD + \lambda I)^{-1} \cdot UU^T \cdot DD^T(D^TD + \lambda I)^{-1} - UU^T \cdot DD^T(D^TD + \lambda I)^{-1}$$

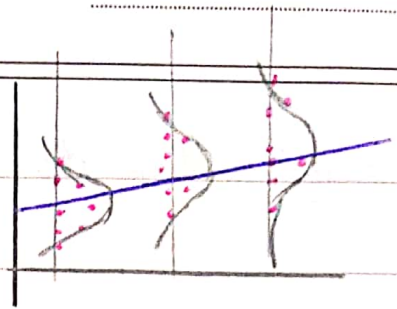$$UU^T \cdot \underbrace{DD^T(D^TD)^{-1}}_{I}(I + \lambda I)^{-1} \cdot \underbrace{DD^T(D^TD)^{-1}}_{I}(I + \lambda I)^{-1} - UU^T(I + \lambda I)^{-1}$$

$$= UU^T(I + \lambda)^2 - UU^T(I + \lambda) \neq 0 \longrightarrow (I + \lambda)^2 \neq (I + \lambda)$$

$$\boxed{H^2 \neq H}$$

# problem 4 : weighted linear regression →

_derive the optimal solution $\hat{\beta}_n$ for the weighted loss function:

$$L(\beta) = \frac{1}{2}\sum_{i=1}^{n} w[i] \times (y[i] - \beta^T x[i])^2$$

$$\Rightarrow L(\beta) = W(Y - x\beta)^T (Y - x\beta)$$

$$= w^T(Y^T Y - 2\beta^T x^T Y + \beta^T x^T x \beta)$$

$$= w^T Y^T Y - 2w^T \beta^T x^T Y + w^T \beta^T x^T x \beta$$

$$\Rightarrow \hat{\beta} = \operatorname*{argmin}_{\beta} L(\beta) \to \frac{\partial L(\beta)}{\partial \beta} = -2w^T x^T Y + 2w^T x^T x \beta$$

$$\Rightarrow \hat{\beta} = \frac{w^T x^T Y}{w^T x^T x} = \boxed{(w^T x^T x)^{-1} w^T x^T Y}$$

$$\Rightarrow \text{Hessian} \to \nabla^2 L(\beta) = 2w^T x^T x$$

## problem 5 : Maximum likelihood estimation of multinomial distribution.

$x \to$ random variable $\to x_1, x_2, \ldots, x_K$

have (k parameters) $\to \Pi = (\pi_1, \ldots, \pi_K)$, $P(x = x_K) = \pi_K$

subject to $\to \sum_K \pi_K = 1 \to \sum_K \pi_K - 1 = 0$

observed data $\to (n_1, \ldots, n_K) \to n_K$: the number of times the value $x_K$ appears in the data

prove the MLE for $\pi_K$ is $\dfrac{n_K}{n}$ where $n = \sum n_K$

$\underset{\downarrow}{\overbrace{\dfrac{P(D|\theta)}{}}}$

$$P(n_1, \ldots, n_K | \pi_1, \ldots, \pi_K) = \binom{n}{n_1, n_2, n_K} \pi_1^{n_1} \pi_2^{n_2} \cdots \pi_K^{n_K} = \dfrac{n!}{\prod n_i!} \prod \pi_i^{n_i}$$

$$\ell(\overbrace{\pi_1, \ldots, \pi_K}^{\theta}) = \log n! - \sum_{i=1}^{k} \log n_i! + \sum_{i=1}^{k} n_i \log \pi_i$$

$$\underset{\text{lagrangian}}{L(\pi_1, \ldots, \pi_K, \lambda)} = \ell(\pi_1, \ldots, \pi_K) + \lambda\left(\sum_{i=1}^{k} \pi_i - 1\right)$$

$$= \log n! - \sum_{i=1}^{k} \log n_i! + \sum_{i=1}^{k} n_i \log \pi_i + \lambda\left(\sum_{i=1}^{k} \pi_i - 1\right)$$

$$\dfrac{\partial L(\pi_1, \ldots, \pi_K, \lambda)}{\partial \pi} = 0 \to \sum_{i=1}^{k} n_i \dfrac{1}{\pi_i \ln 10} + \lambda = 0 \Rightarrow \lambda = -\sum \dfrac{n_i}{\pi_i \ln 10}$$

$$\dfrac{\partial L(\pi_1, \ldots, \pi_K, \lambda)}{\partial \lambda} = 0 \to \sum_{i=1}^{k} \pi_i - 1 = 0 \to \sum \pi_i = 1$$