

$$1. Y = f(x) + \varepsilon$$

$$(a) \text{ Show } E[(Y - \hat{f}(x; D))^2] = \text{Var}(\varepsilon) + E[(f(x) - E[\hat{f}(x; D)|x])^2] + E[(\hat{f}(x; D) - E[\hat{f}(x; D)|x])^2]$$

$$\begin{aligned} \text{So } E[(Y - \hat{f}(x; D))^2] &= E[(Y - f(x) + f(x) - \hat{f}(x; D))^2] \\ &\stackrel{(1)}{=} E[(Y - f(x))^2] + \stackrel{(2)}{=} E[(f(x) - \hat{f}(x; D))^2] \\ &\quad + \stackrel{(3)}{=} 2E[(Y - f(x))(f(x) - \hat{f}(x; D))] \end{aligned}$$

$$\textcircled{1} E[(Y - f(x))^2] = \text{Var}(\varepsilon)$$

$$\begin{aligned} \textcircled{2} E[(f(x) - \hat{f}(x; D))^2] &= E[(f(x) - E_b[\hat{f}(x; D)|x] + E_b[\hat{f}(x; D)|x] - \hat{f}(x; D))^2] \\ &= E[(f(x) - E_b[f(x; D)|x])^2] \quad \text{--- bias} \\ &\quad + E[(E_b[\hat{f}(x; D)|x] - \hat{f}(x; D))^2] \quad \text{--- variance} \end{aligned}$$

$$\begin{aligned} \textcircled{3} &= E[(Y - f(x))(f(x) - \hat{f}(x; D))] \\ &= E[E[(Y - f(x))(f(x) - \hat{f}(x; D)) | x, D]] \\ &= E[(f(x) - f(x))(f(x) - \hat{f}(x; D))] \\ &= 0 \end{aligned}$$

$$\begin{aligned} \textcircled{1} + \textcircled{2} + \textcircled{3} &\Rightarrow E[(Y - \hat{f}(x; D))^2] = \text{Var}(\varepsilon) + E[(f(x) - E[\hat{f}(x; D)|x])^2] \\ &\quad + E[(\hat{f}(x; D) - E[\hat{f}(x; D)|x])^2] \end{aligned}$$

(b) As the flexibility increase, the model has small bias and high variance.

2. (a)

$$L = \|y - \underline{X}\beta\|^2 + \lambda \|\beta\|^2 = (y - \underline{X}\beta)^T (y - \underline{X}\beta) + \lambda \beta^T \beta$$

$$\frac{\partial L}{\partial \beta} = 0 \Rightarrow -2\underline{X}^T y + 2\underline{X}^T \underline{X} \beta + 2\lambda \beta = 0$$

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

$$(b) SE(\hat{\beta})^2 = E[(\hat{\beta} - \beta)(\hat{\beta} - \beta)^T]$$

$$= E\left[\left((\underline{X}^T \underline{X} + \lambda \underline{I})^{-1} \underline{X}^T y - \beta\right)\left((\underline{X}^T \underline{X} + \lambda \underline{I})^{-1} \underline{X}^T y - \beta\right)^T\right]$$

$$= E\left[\left((\underline{X}^T \underline{X} + \lambda \underline{I})^{-1} \underline{X}^T (\underline{X}\beta + \varepsilon) - \beta\right)\left((\underline{X}^T \underline{X} + \lambda \underline{I})^{-1} \underline{X}^T (\underline{X}\beta + \varepsilon) - \beta\right)^T\right]$$

$$= E\left[\left((\underline{X}^T \underline{X} + \lambda \underline{I})^{-1} \underline{X}^T \underline{X} \beta + (\underline{X}^T \underline{X} + \lambda \underline{I})^{-1} \underline{X}^T \varepsilon - \beta\right)\left((\underline{X}^T \underline{X} + \lambda \underline{I})^{-1} \underline{X}^T \underline{X} \beta + (\underline{X}^T \underline{X} + \lambda \underline{I})^{-1} \underline{X}^T \varepsilon - \beta\right)^T\right]$$

$$\because \underline{X}^T \underline{X} = a \underline{I}_p$$

$$= E\left[\left(\frac{a}{a+\lambda} \beta + \frac{1}{a+\lambda} \underline{X}^T \varepsilon - \beta\right)\left(\frac{a}{a+\lambda} \beta + \frac{1}{a+\lambda} \underline{X}^T \varepsilon - \beta\right)^T\right]$$

$$= E\left[\left(\frac{1}{a+\lambda} (-\lambda \beta + \underline{X}^T \varepsilon)\right)\left(\frac{1}{a+\lambda} (-\lambda \beta + \underline{X}^T \varepsilon)\right)^T\right]$$

$$= \left(\frac{1}{a+\lambda}\right)^2 (\lambda^2 \beta \beta^T + a \sigma^2 \underline{I}_p)$$

$$\rightarrow SE(\beta_{jj}) = \frac{1}{a+\lambda} (\lambda^2 \{\beta \beta^T\}_{jj} + a \sigma^2)$$

Q3:-

Given dataset $D = \{(x_1, y_1), (x_2, y_2), (x_3, y_3)\}$
 $(1, 2), (3, 7), (5, 8)$

3-fold Cv:-

Divide "D" into three groups

D_1, D_2, D_3 of size $n/k = 3/3 = 1$

Let $D_1 = \{(x_1, y_1)\} \rightarrow$ validation set; $D/D_1 = \{(x_2, y_2), (x_3, y_3)\}$
 as training set.

$$\hat{\beta}_1 = 0.5; \hat{\beta}_0 = 5.50$$

$$\therefore MSE_1 = (y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1)^2$$

$$= (2 - 5.50 - 0.5 \times 1)^2 = 16$$

$$MSE_j = \frac{1}{n_j} \sum_{i \in (x_i, y_i) \in D_j} (y_i - \hat{f}(x_i | D_j))^2$$

Let $D_2 = \{(x_2, y_2)\}$; $D/D_2 = \{(x_1, y_1), (x_3, y_3)\}$

$$\hat{\beta}_1 = 1.5 \quad \hat{\beta}_0 = 0.5$$

$$\therefore MSE_2 = (7 - 0.5 - 1.5 \times 3)^2 = 4$$

Let $D_3 = \{(x_3, y_3)\}$; $D/D_3 = \{(x_1, y_1), (x_2, y_2)\}$

$$\hat{\beta}_1 = 2.5 \quad \hat{\beta}_0 = -0.5$$

$$\therefore MSE_3 = (8 + 0.5 - 2.5 \times 5)^2 = 16$$

$$\therefore CV_{(3)} = \sum_{j=1}^3 \frac{n_j}{n} MSE_j \Rightarrow \frac{1}{3} [16 + 4 + 16] \Rightarrow 12 \quad \star$$

4. (a)

$$l(\beta_0, \beta_1) = \prod_{i=1}^n P(Y=y_i | X=x_i; \beta_0, \beta_1)$$

$$= \prod_{i: y_i=1}^n P(x_i; \beta_0, \beta_1) \prod_{i: y_i=0}^n (1 - P(x_i; \beta_0, \beta_1))$$

$$\Rightarrow \log l(\beta_0, \beta_1) = \sum_{i=1}^n \left[\frac{1}{2}(y_i+1) \log P(x_i; \beta_0, \beta_1) + \frac{1}{2}(1-y_i) \log (1 - P(x_i; \beta_0, \beta_1)) \right]$$

$$= \sum_{i=1}^n \left[\frac{1}{2}(y_i+1) \log \frac{1}{1 + e^{-y_i(\beta_0 + \beta_1 x_i)}} + \frac{1}{2}(1-y_i) \log \frac{1}{1 + e^{-y_i(\beta_0 + \beta_1 x_i)}} \right]$$

$$= \sum_{i=1}^n \log \frac{1}{1 + e^{-y_i(\beta_0 + \beta_1 x_i)}}$$

$$\Rightarrow \hat{\beta} = \arg \max_{\beta} \sum_{i=1}^n \log \frac{1}{1 + e^{-y_i(\beta_0 + \beta_1 x_i)}}$$

$$= \arg \min_{\beta} \sum_{i=1}^n \log (1 + e^{-y_i(\beta_0 + \beta_1 x_i)})$$

$$\Rightarrow \hat{\beta} = \arg \min_{\beta} \sum_{i=1}^n \log (1 + e^{-y_i(\beta_0 + \beta_1 x_i)})$$

(b) $\hat{Y} = \arg \max_{k \in \{1, \dots, K\}} P(Y=k | X=x)$

$$= \arg \max_{k \in \{1, \dots, K\}} \frac{\pi_k f_k(x)}{\sum_{l=1}^K \pi_l f_l(x)}$$

$$= \arg \max_{k \in \{1, \dots, K\}} \frac{\pi_k \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left(-\frac{1}{2\sigma_k^2} (x - \mu_k)^2\right)}{\sum_{l=1}^K \pi_l \frac{1}{\sqrt{2\pi\sigma_l^2}} \exp\left(-\frac{1}{2\sigma_l^2} (x - \mu_l)^2\right)}$$

$$= \arg \max_{k \in \{1, \dots, K\}} \log \pi_k - \frac{1}{2\sigma_k^2} (x - \mu_k)^2 - \frac{1}{2} \log \sigma_k^2$$

Q.5(a) given Predictors (age, study, hno)

Backward stepwise selection algorithm

1) $p=3$; M_3 (full model) = (age + study + hno) \underline{RSS} 14.4

2) for M_2 : choose the best (i.e. with smallest RSS)
among (age + study), (age + hno), (study + hno)
 $RSS = 18.2$ 17.4 17.6

$\therefore M_2 = (\text{age} + \text{hno})$ $\underline{RSS} : 17.4$

3) for M_1 : choose the best among age, hno, study
 $M_1 = (\text{hno})$ $\underline{RSS} : 22.3$

Q.5(b) $C_p = \frac{1}{n} [RSS + 2(p+1)\hat{\sigma}^2]$, $n=10$

where

$\hat{\sigma}^2$ = Estimate of variance of error ϵ (with full model)

$$= \frac{RSS}{n-p-1} = \frac{14.4}{10-4} = \frac{14.4}{6} = 2.4$$

↑
"3"

for M_3 :

$$C_p = \frac{1}{10} [14 \cdot 4 + 2(3+1) \times 2 \cdot 4] = 3.66$$

for M_2 :

$$C_p = \frac{1}{10} [17 \cdot 4 + 2(2+1) \times 2 \cdot 4] = 3.18$$

for M_1 :

$$C_p = \frac{1}{10} [22 \cdot 3 + 2(1+1) \times 2 \cdot 4] = 3.19$$

Choose the model that has smallest " C_p "
from above, Model $M_2 = (\text{age} + \text{hwr})$ has least C_p
 $\therefore M_2$ is the best model \times

Question 6:

a) Cross validation has less bias and tends not to over estimate the test error rate as much as validation set does, because CV repeatedly fit the statistical learning method using entire training data, whereas in validation set approach uses part of the size of original dataset.

b) By using Linear model, $p(x) = \beta_0 + \beta_1 x$ may go beyond $[0, 1]$

\therefore We must model $p(x)$ using a function that gives outputs between 0 and 1 $\forall x$.

logistic function is one as such.

c) $t_j = \frac{\hat{\beta}_j}{\hat{SE}(\hat{\beta}_j)}$ should be large;

which measures how far $\hat{\beta}_j$ is to 0 relative to $\hat{SE}(\hat{\beta}_j)$.