

1 Ex.6.2

Define the vector-valued function $b(x)^T = (1, x, x^2, \dots, x^k)$ for $k \geq 0$. Let \mathbf{B} be the $N \times (k+1)$ regression matrix with i th row $b(x_i)^T$, and $\mathbf{W}(x_0)$ the $N \times N$ diagonal matrix with i th diagonal element $K_\lambda(x_0, x_i)$. Then we have

$$b(x_0)^T = b(x_0)^T (\mathbf{B}^T \mathbf{W}(x_0) \mathbf{B})^{-1} \mathbf{B}^T \mathbf{W}(x_0) \mathbf{B}. \quad (1)$$

Note the definition of $l_i(x_0)$ in (6.9) in text, from (1), we have

$$\begin{aligned} 1 &= b(x_0)^T (\mathbf{B}^T \mathbf{W}(x_0) \mathbf{B})^{-1} \mathbf{B}^T \mathbf{W}(x_0) \mathbf{1} = \sum_{i=1}^N l_i(x_0) \\ x_0 &= b(x_0)^T (\mathbf{B}^T \mathbf{W}(x_0) \mathbf{B})^{-1} \mathbf{B}^T \mathbf{W}(x_0) \mathbf{B}_2 = \sum_{i=1}^N l_i(x_0) x_i \\ &\dots \\ x_0^k &= b(x_0)^T (\mathbf{B}^T \mathbf{W}(x_0) \mathbf{B})^{-1} \mathbf{B}^T \mathbf{W}(x_0) \mathbf{B}_{k+1} = \sum_{i=1}^N l_i(x_0) x_i^k \end{aligned}$$

where \mathbf{B}_i is the i th column of \mathbf{B} (note that $\mathbf{B}_1 = \mathbf{1}$). Therefore we have $b_0(x_0) = \sum_{i=1}^N l_i(x_0) = 1$ and

$$b_1(x_0) = \sum_{i=1}^N (x_i - x_0) l_i(x_0) = \sum_{i=1}^N l_i(x_0) x_i - x_0 \sum_{i=1}^N l_i(x_0) = x_0 - x_0 \cdot 1 = 0.$$

For $j \geq 2$, we have

$$\begin{aligned} b_j(x_0) &= \sum_{i=1}^N (x_i - x_0)^j l_i(x_0) \\ &= \sum_{i=1}^N \left(\sum_{b=0}^j C_j^b (-1)^b x_i^{j-b} x_0^b \right) l_i(x_0) \\ &= \sum_{b=0}^j C_j^b (-1)^b x_0^b \left(\sum_{i=1}^N l_i(x_0) x_i^{j-b} \right) \\ &= \sum_{b=0}^j C_j^b (-1)^b x_0^b x_0^{j-b} \\ &= \sum_{b=0}^j C_j^b (-1)^b x_0^j \\ &= (1-1)^j x_0^j \\ &= 0. \end{aligned}$$

By Taylor expansion we have

$$\begin{aligned} E \left[\hat{f}(x_0) \right] - f(x_0) &= \sum_{i=1}^N l_i(x_0) f(x_i) - f(x_0) \\ &= f(x_0) \sum_{i=1}^N l_i(x_0) - f(x_0) + f'(x_0) \sum_{i=1}^N (x_i - x_0) l_i(x_0) \\ &\quad + \frac{f''(x_0)}{2} \sum_{i=1}^N (x_i - x_0)^2 l_i(x_0) \\ &\quad + \dots \\ &\quad + (-1)^k \frac{f^{(k)}(x_0)}{k!} \sum_{i=1}^N (x_i - x_0)^k l_i(x_0) \\ &\quad + R \\ &= R, \end{aligned}$$

where the remainder term R involves $(k+1)$ th and higher-order derivatives of f , on which the bias only depends.

2 Ex.6.3

Let's first introduce notations. Define the vector-valued function $b(x)^T = (1, x, x^2, \dots, x^d)$ for $d \geq 1$. Let \mathbf{B} be the $N \times (d+1)$ regression matrix with i th row $b(x_i)^T$.

$$\mathbf{B} = \begin{pmatrix} 1 & x_1 & \cdots & x_1^d \\ 1 & x_2 & \cdots & x_2^d \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_N & \cdots & x_N^d \end{pmatrix} = \begin{pmatrix} b(x_1)^T \\ b(x_2)^T \\ \vdots \\ b(x_N)^T \end{pmatrix} \in R^{N \times (d+1)}$$

and

Let $\mathbf{W}(x)$ the $N \times N$ diagonal matrix with i th diagonal element $K_\lambda(x, x_i)$, that is,

$$\mathbf{W}(x) = \begin{pmatrix} K_\lambda(x, x_1) & 0 & \cdots & 0 \\ 0 & K_\lambda(x, x_2) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & K_\lambda(x, x_N) \end{pmatrix} \in R^{N \times N}$$

Note that $\mathbf{W}(x) = \mathbf{W}^T(x)$ By definition of $l(x)$ (see, e.g... (6.9) in the text), we have

$$l(x_0)^T = b(x_0)^T (\mathbf{B}^T \mathbf{W}(x_0) \mathbf{B})^{-1} \mathbf{B}^T \mathbf{W}(x_0)$$

Denote $b = b(x_0)$ and $\mathbf{W} = \mathbf{W}(x_0)$ to simplify the notations from now on, we have

$$\begin{aligned} \|l(x_0)\|^2 &= l(x_0)^T l(x_0) \\ &= b^T (\mathbf{B}^T \mathbf{W} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{W} \mathbf{W}^T \mathbf{B} (\mathbf{B}^T \mathbf{W} \mathbf{B})^{-1} b \end{aligned} \quad (1)$$

We need to show $\|l(x_0)\|^2$ is increasing in d . The expression involves with the weighted kernel matrix \mathbf{W} , however it turns out $\|l(x_0)\|^2$ does not depend on \mathbf{W} . Note that we could plug $\mathbf{I} = \mathbf{B} \mathbf{B}^T (\mathbf{B} \mathbf{B}^T)^{-1} = (\mathbf{B} \mathbf{B}^T)^{-1} \mathbf{B} \mathbf{B}^T$ between \mathbf{W} and \mathbf{W}^T in (1), we obtain

$$\begin{aligned} &\|l(x_0)\|^2 \\ &= b^T (\mathbf{B}^T \mathbf{W} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{W} \mathbf{B} \mathbf{B}^T (\mathbf{B} \mathbf{B}^T)^{-1} (\mathbf{B} \mathbf{B}^T)^{-1} \mathbf{B} \mathbf{B}^T \mathbf{W}^T \mathbf{B} (\mathbf{B}^T \mathbf{W} \mathbf{B})^{-1} b \\ &= b^T \mathbf{B}^T (\mathbf{B} \mathbf{B}^T)^{-1} (\mathbf{B} \mathbf{B}^T)^{-1} \mathbf{B} b, \end{aligned}$$

therefore we see $\|l(x_0)\|^2$ is independent of \mathbf{W} So we can take $\mathbf{W} = \mathbf{I}$ in (1), which gives

$$\|l(x_0)\|^2 = b^T (\mathbf{B}^T \mathbf{B})^{-1} b \quad (2)$$

Now consider the case for $d+1$, we denote

$$\hat{\mathbf{B}} = \begin{pmatrix} 1 & x_1 & \cdots & x_1^d & x_1^{d+1} \\ 1 & x_2 & \cdots & x_2^d & x_2^{d+1} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & x_N & \cdots & x_N^d & x_N^{d+1} \end{pmatrix} = (\mathbf{B} \quad c) \in R^{N \times (d+2)}$$

where

$$c = \begin{pmatrix} x_1^{d+1} \\ x_2^{d+1} \\ \vdots \\ x_N^{d+1} \end{pmatrix} \in R^{N \times 1}$$

Similarly, denote $\hat{b}^T = (b^T, x_0^{d+1}) = (1, x_0, x_0^2, \dots, x_0^d, x_0^{d+1}) \in R^{1 \times (d+2)}$. In that case, (2) becomes

$$\left\| \hat{l}(x_0) \right\|^2 = \hat{b}^T \left(\hat{\mathbf{B}}^T \hat{\mathbf{B}} \right)^{-1} \hat{b} \quad (3)$$

Further, we have

$$\hat{\mathbf{B}}^T \hat{\mathbf{B}} = \begin{pmatrix} \mathbf{B}^T \\ c^T \end{pmatrix} \begin{pmatrix} \mathbf{B} & c \end{pmatrix} = \begin{pmatrix} \mathbf{B}^T \mathbf{B} & \mathbf{B}^T c \\ c^T \mathbf{B} & c^T c \end{pmatrix} \in R^{(d+2) \times (d+2)}.$$

Note that $c^T c \in R^1$ is a scalar. Recall the formula for block matrix inverse, (e.g., Schur complement), we have

$$\begin{aligned} & \left(\hat{\mathbf{B}}^T \hat{\mathbf{B}} \right)^{-1} \\ &= \begin{pmatrix} (\mathbf{B}^T \mathbf{B})^{-1} + \frac{1}{k} (\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T c c^T \mathbf{B} (\mathbf{B}^T \mathbf{B})^{-1} & -\frac{1}{k} (\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T c \\ -\frac{1}{k} c^T \mathbf{B} (\mathbf{B}^T \mathbf{B})^{-1} & \frac{1}{k} \end{pmatrix} \end{aligned}$$

where

$$k = c^T c - c^T \mathbf{B} (\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T c \quad (4)$$

Denote $\beta = (\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T c \in R^{(d+2) \times 1}$, plug $\left(\hat{\mathbf{B}}^T \hat{\mathbf{B}} \right)^{-1}$ into (3), we get

$$\begin{aligned} & \left\| \hat{l}(x_0) \right\|^2 \\ &= \begin{pmatrix} b^T & x_0^{d+1} \end{pmatrix} \begin{pmatrix} (\mathbf{B}^T \mathbf{B})^{-1} + \frac{1}{k} \beta \beta^T & -\frac{1}{k} \beta \\ -\frac{1}{k} \beta^T & \frac{1}{k} \end{pmatrix} \begin{pmatrix} b \\ x_0^{d+1} \end{pmatrix} \\ &= b^T (\mathbf{B}^T \mathbf{B})^{-1} b + \frac{1}{k} \left[b^T \beta \beta^T b - x_0^{d+1} \beta^T b - x_0^{d+1} b^T \beta + (x_0^{d+1})^2 \right] \\ &= b^T (\mathbf{B}^T \mathbf{B})^{-1} b + \frac{1}{k} (x_0^{d+1} - b^T \beta)^2 \quad (\text{note } b^T \beta \in R) \\ &= \|l(x_0)\|^2 + \frac{1}{k} (x_0^{d+1} - b^T \beta)^2. \end{aligned}$$

Therefore, it suffices to show that $k > 0$ for k defined in (4). To do that, we only need to show

$$\mathbf{B} (\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T \preceq \mathbf{I}_N.$$

Consider the QR decomposition of \mathbf{B}

$$\mathbf{B} = \mathbf{Q} \mathbf{R}$$

where \mathbf{Q} is an $N \times (d+1)$ orthogonal matrix, $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}_N$, and $\mathbf{R} \in R^{(d+1) \times (d+1)}$ is an upper triangular matrix. Then

$$\mathbf{B} (\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T = \mathbf{Q} \mathbf{R} (\mathbf{R}^T \mathbf{R})^{-1} \mathbf{R}^T \mathbf{Q}^T = \mathbf{Q} \mathbf{Q}^T. \quad (5)$$

Let $(\mathbf{Q} \mathbf{Q}_1)$ be an $N \times N$ orthogonal matrix, we have

$$\mathbf{I}_N = \begin{pmatrix} \mathbf{Q} & \mathbf{Q}_1 \end{pmatrix} \begin{pmatrix} \mathbf{Q}^T \\ \mathbf{Q}_1^T \end{pmatrix} = \mathbf{Q} \mathbf{Q}^T + \mathbf{Q}_1 \mathbf{Q}_1^T$$

The result (5) follows by noting $\mathbf{Q}_1 \mathbf{Q}_1^T$ is positive semi-definite. The proof is now complete.

3 Ex.6.7

Note that local regression smoothers are linear estimators, and we can write

$$\hat{\mathbf{f}} = \mathbf{S}_{\lambda \mathbf{y}}$$

where $\{\mathbf{S}_{\lambda}\}_{ij} = l_i(x_j)$ for $l_i(x)$ defined by (6.8) in the text. Then we know

$$y_i - \hat{f}^{-i}(x_i) = \frac{y_i - \hat{f}(x_i)}{1 - \{\mathbf{S}_{\lambda}\}_{ii}}$$

4 Ex.6.10

Consider the in-sample prediction error (7.18) and the training error $\overline{\text{err}}$ in the case of squared-error loss:

$$\begin{aligned}\text{Err}_{\text{in}} &= \frac{1}{N} \sum_{i=1}^N E_{Y^0} \left(Y_i^0 - \hat{f}(x_i) \right)^2 \\ \overline{\text{err}} &= \frac{1}{N} \sum_{i=1}^N \left(y_i - \hat{f}(x_i) \right)^2.\end{aligned}$$

Add and subtract $f(x_i)$ and $E\hat{f}(x_i)$ in each expression and expand. Hence establish that the average optimism in the training error is

$$\frac{2}{N} \sum_{i=1}^N \text{Cov}(\hat{y}_i, y_i),$$

. So we know that

$$\text{PE}(\lambda) = \text{ASR}(\lambda) + \frac{2}{N} \sum_{i=1}^N \text{Cov}(\hat{y}_i, y_i)$$

and we have

$$\begin{aligned}\sum_{i=1}^N \text{Cov}(\hat{y}_i, y_i) &= \text{trace}(\text{Cov}(\hat{\mathbf{y}}, \mathbf{y})) \\ &= \text{trace}(\text{Cov}(\mathbf{S}\mathbf{y}, \mathbf{y})) \\ &= \text{trace}(\mathbf{S}\mathbf{Cov}(\mathbf{y}, \mathbf{y})) \\ &= \text{trace}(\mathbf{S}\mathbf{Var}(\mathbf{y})) \\ &= \text{trace}(\mathbf{S})\sigma_\epsilon^2.\end{aligned}$$

Then the proof is straightforward.