

Textbook:

4.7)

Non-parametric model

$$Y_i = f(x_{i1}, \dots, x_{ip}) + \epsilon_i$$

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

(2) ϵ_i are iid, $\rightarrow \text{Cov}(Y_i, Y_j)$

$$\text{Var}(Y_i) = \text{Var}(\epsilon_i) = \sigma^2, \text{ for } i=1, \dots, n.$$

Linear smoother: $\hat{Y} = SY$

$$\underline{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} \quad \begin{matrix} \nearrow \\ \nearrow \\ \nearrow \end{matrix} \begin{matrix} n \times n \\ \text{Smother matrix} \\ \text{Predicted values} \end{matrix}$$

$$S = \begin{bmatrix} s_{11} & \dots & s_{1n} \\ \vdots & \ddots & \vdots \\ s_{n1} & \dots & s_{nn} \end{bmatrix} = \begin{bmatrix} s_1 \\ \vdots \\ s_n \end{bmatrix}$$

Then we can write

(1) $\hat{Y}_i = s_i Y$

We want to show that

$$\sum_{i=1}^n \text{Cov}(\hat{Y}_i, Y_i) = \text{tr}(S)\sigma^2$$

Let's look at $\text{Cov}(\hat{Y}_i, Y_i)$

$$\text{Cov}(\hat{Y}_i, Y_i) \stackrel{(1)}{=} \text{Cov}(s_i Y, Y_i)$$

$$\begin{aligned} &= \text{Cov}(s_{i1}Y_1 + \dots + s_{in}Y_n, Y_i) \\ &\stackrel{\text{Rules of covariance}}{=} \text{Cov}(s_{i1}Y_1, Y_i) + \dots + \text{Cov}(s_{in}Y_n, Y_i) \\ &= s_{i1}\text{Cov}(Y_1, Y_i) + \dots + s_{ii}\text{Cov}(Y_i, Y_i) \\ &\quad + \dots + s_{in}\underbrace{\text{Cov}(Y_n, Y_i)}_{=0 \text{ due to independence}} \\ &\stackrel{(2)}{=} s_{ii}\text{Cov}(Y_i, Y_i) \\ &= s_{ii}\text{Var}(Y_i) \\ &= s_{ii}\sigma^2 \end{aligned}$$

We can now look at the sum

$$\begin{aligned} \sum_{i=1}^n \text{Cov}(\hat{Y}_i, Y_i) &= \sum_{i=1}^n s_{ii}\sigma^2 \\ &= \sigma^2 \sum_{i=1}^n s_{ii} \\ &= \sigma^2 \cdot \text{tr}(S) \end{aligned}$$

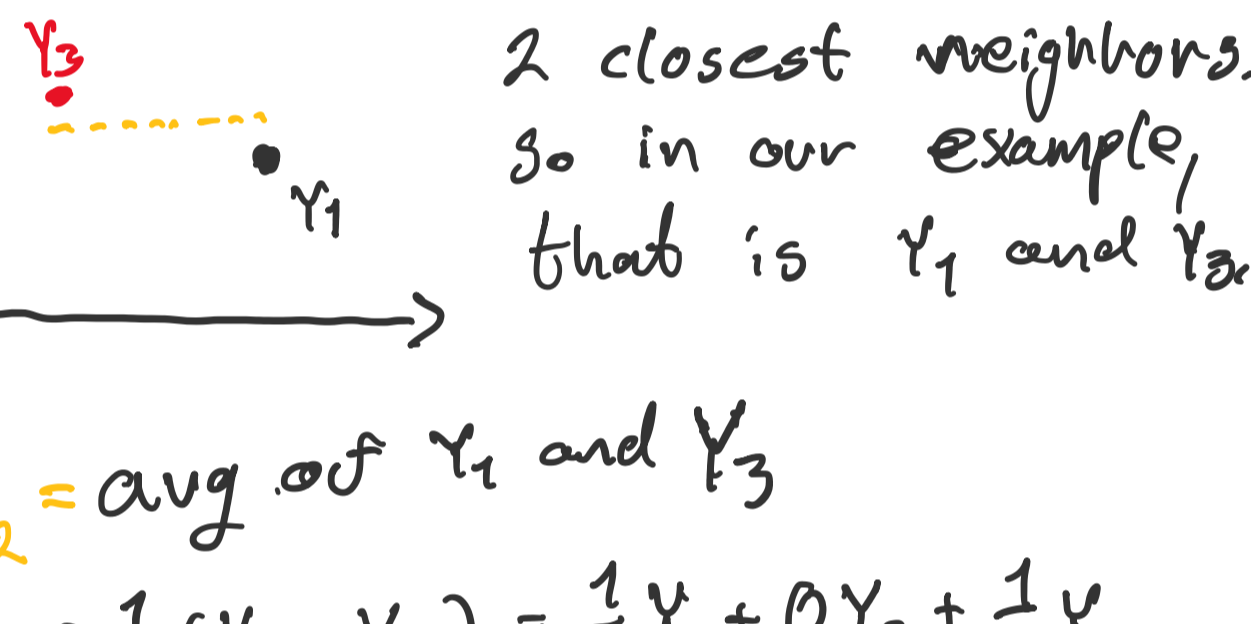
$$\text{tr}(X) = \text{tr} \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{pmatrix} = \sum_{i=1}^n x_{ii}$$

Sum of diagonal elements

$$\hat{Y} = HY, \quad H = X(X^T X)^{-1} X^T$$

symmetric

KNN - we look at $k=2$.



$$\hat{Y}_2 = \text{avg. of } Y_1 \text{ and } Y_3$$

$$= \frac{1}{2}(Y_1 + Y_3) = \frac{1}{2}Y_1 + 0Y_2 + \frac{1}{2}Y_3$$

$$\text{So if } \underline{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \end{bmatrix} \quad \hat{Y}_2 = \begin{bmatrix} \frac{1}{2} & 0 & \frac{1}{2} \end{bmatrix} \underline{Y} = \underline{s}_2 \underline{Y}$$

original response

So $s_{21} = 1/2$, that means that Y_1 contributes half of the value of \hat{Y}_2 .

4.8)

Tree growth algorithm on page 99-103.

Show that $D_j - D_j^* > 0$.

1D-example

(Indicates that training RSS always decreases in each new iteration of the tree growth algorithm.)



So in iteration k of the algorithm we split R_j into two subregions R_j' and R_j'' .

$$D_j = \sum_{i \in R_j} (y_i - \hat{c}_j)^2, \text{ where } \hat{c}_j \text{ is the average response in region } R_j.$$

This means that any other value for \hat{c}_j will yield a larger D_j .

$$D_j^* = \sum_{i \in R_j'} (y_i - \hat{c}_j')^2 + \sum_{i \in R_j''} (y_i - \hat{c}_j'')^2$$

$$\begin{aligned} &< \sum_{i \in R_j'} (y_i - \hat{c}_j)^2 + \sum_{i \in R_j''} (y_i - \hat{c}_j)^2 \\ &= \sum_{i \in R_j} (y_i - \hat{c}_j)^2 \\ &= D_j \end{aligned}$$

We have that $D_j^* < D_j$, this means that $0 < D_j - D_j^*$.

Degenerate cases.

1. If there is only 1 obs. in R_j .
2. If all the responses in R_j have the same response value.

Then $\hat{c}_j = \hat{c}_j' = \hat{c}_j''$.



(Maybe enough with the means of R_j' and R_j'' being the same).