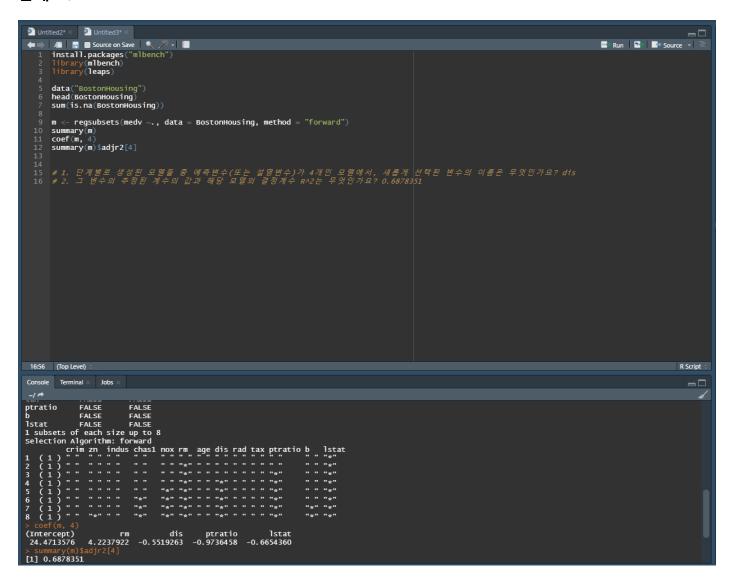
ESC 21-SUMMER week 2_Homework

W.J.Park

문제 1.



Ex. 7.4 Consider the in-sample prediction error (7.18) and the training error $\overline{\text{err}}$ in the case of squared-error loss:

$$\operatorname{Err}_{\text{in}} = \frac{1}{N} \sum_{i=1}^{N} \operatorname{E}_{Y^{0}} (Y_{i}^{0} - \hat{f}(x_{i}))^{2}$$

$$\overline{\operatorname{err}} = \frac{1}{N} \sum_{i=1}^{N} (y_{i} - \hat{f}(x_{i}))^{2}.$$

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

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

as given in (7.21).

Sol)
$$\int_{\mathbb{R}^{2}} \int_{\mathbb{R}^{2}} \left[E_{3}(0) - E_{3}(0) - E_{3}(0) \right] = \frac{1}{2} \sum_{i=1}^{N} C_{i}(\hat{\beta}_{i}, \hat{\beta}_{i})$$

$$= E_{3}(0) = E_{3}(0) = \frac{1}{2} \sum_{i=1}^{N} C_{i}(\hat{\beta}_{i}, \hat{\beta}_{i})$$

$$= E_{3}\left(\frac{1}{N}\sum_{i=1}^{N} E_{3} E_{3}(\hat{\gamma}_{i}) - \hat{f}(\kappa_{i})\right)^{2} - \frac{1}{N}\sum_{i=1}^{N} (y_{i} - \hat{f}(\kappa_{i}))^{2}\right)$$

$$= \frac{1}{N}\sum_{i=1}^{N} E_{3} E_{3} E_{3}(\hat{\gamma}_{i}) - E_{3}(y_{i} - \hat{f}(\kappa_{i}))^{2} - \frac{1}{N}\sum_{i=1}^{N} (y_{i} - \hat{f}(\kappa_{i}))^{2}$$

$$= \frac{1}{N}\sum_{i=1}^{N} \left[E_{3}E_{5}(\hat{\gamma}_{i}) - E_{3}(y_{i} - \hat{\gamma}_{i}) \right]$$

$$= \frac{1}{N}\sum_{i=1}^{N} \left[E_{3}(x_{i} - \hat{\gamma}_{i}) - E_{3}(y_{i} - \hat{\gamma}_{i}) + E_{3}(y_{i} - \hat{\gamma}_{i}) - E_{3}(y_{i} - \hat{\gamma}_{i}) + E_{3}(y_{i} - \hat{\gamma}_{i}) - E_{3}(y_{i} - \hat{\gamma}_{i}) \right]$$

$$= \frac{1}{N}\sum_{i=1}^{N} \left[E_{3}(y_{i} - \hat{\gamma}_{i}) - E_{3}(y_{i} - \hat{\gamma}_{i}) + E_{3}(y_{i} - \hat{\gamma}_{i}) + E_{3}(y_{i} - \hat{\gamma}_{i}) + E_{3}(y_{i} - \hat{\gamma}_{i}) \right]$$

$$= \frac{1}{N}\sum_{i=1}^{N} \left[E_{3}(y_{i} - \hat{\gamma}_{i}) - E_{3}(y_{i} - \hat{\gamma}_{i}) + E_{3}(y_{i} - \hat{\gamma}_{i}) + E_{3}(y_{i} - \hat{\gamma}_{i}) \right]$$

$$= \frac{1}{N}\sum_{i=1}^{N} \left[E_{3}(y_{i} - \hat{\gamma}_{i}) - E_{3}(y_{i} - \hat{\gamma}_{i}) + E_{3}(y_{i} - \hat{\gamma}_{i}) \right]$$

$$= \frac{1}{N}\sum_{i=1}^{N} \left[E_{3}(y_{i} - \hat{\gamma}_{i}) - E_{3}(y_{i} - \hat{\gamma}_{i}) + E_{3}(y_{i} - \hat{\gamma}_{i}) \right]$$

$$= \frac{1}{N}\sum_{i=1}^{N} \left[E_{3}(y_{i} - \hat{\gamma}_{i}) - E_{3}(y_{i} - \hat{\gamma}_{i}) + E_{3}(y_{i} - \hat{\gamma}_{i}) \right]$$

$$= \frac{1}{N}\sum_{i=1}^{N} \left[E_{3}(y_{i} - \hat{\gamma}_{i}) - E_{3}(y_{i} - \hat{\gamma}_{i}) + E_{3}(y_{i} - \hat{\gamma}_{i}) \right]$$

$$= \frac{1}{N}\sum_{i=1}^{N} \left[E_{3}(y_{i} - \hat{\gamma}_{i}) - E_{3}(y_{i} - \hat{\gamma}_{i}) + E_{3}(y_{i} - \hat{\gamma}_{i}) \right]$$