

Homework 10

Due Friday, 11/22/19. Please place give your homework to the TA Jiaxuan Duan or place it in his mailbox by Friday afternoon.

1. (Training vs. test error) Suppose $y_i = f(\mathbf{x}_i) + \epsilon_i$ for some unknown function f , where \mathbf{x}_i is non-random, $\mathbb{E}(\epsilon_i) = 0$ and $\mathbb{E}(\epsilon_i^2) = \sigma^2$ for $i = 1, \dots, n$. Suppose you use the training set $\mathcal{T} = \{y_1, \dots, y_n\}$ to obtain \hat{f} , an estimate for f . Define the in-sample **test error** and **training error** to be

$$\text{Err}_{\text{in}} = \mathbb{E} \left[n^{-1} \sum_{i=1}^n \{\tilde{y}_i - \hat{f}(\mathbf{x}_i)\}^2 \mid \mathcal{T} \right], \quad \overline{\text{err}} = n^{-1} \sum_{i=1}^n \{y_i - \hat{f}(\mathbf{x}_i)\}^2$$

where \tilde{y}_i is an independent copy of y_i (i.e. \tilde{y}_i and y_i are independent and have the same distribution). Err_{in} determines how the estimator performs in the prediction of new unobserved data. Note that Err_{in} and $\overline{\text{err}}$ are random quantities, where the randomness is due to the randomness in the training data \mathcal{T} .

- (a) Show that

$$\mathbb{E}(\text{Err}_{\text{in}}) = n^{-1} \sum_{i=1}^n \left(\left[\text{Bias} \{ \hat{f}(\mathbf{x}_i) \} \right]^2 + \text{Var} \{ \hat{f}(\mathbf{x}_i) \} \right) + \sigma^2,$$

where the expectation is taken over the training set \mathcal{T} .

- (b) Show that the expected **optimism** in the training error, $\omega = \text{Err}_{\text{in}} - \overline{\text{err}}$, is

$$\mathbb{E}(\omega) = \frac{2}{n} \sum_{i=1}^n \text{Cov} \{ y_i, \hat{f}(\mathbf{x}_i) \},$$

where the expectation is taken over the training set \mathcal{T} . Why does this show that the training error usually **underestimates** the test (i.e. prediction) error?

2. (Ridge regression) Suppose $\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_n)^T \in \mathbb{R}^{n \times p}$, $\mathbf{Y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n$ and $\lambda \geq 0$.

- (a) Show that $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p$ is invertible for all $\lambda > 0$, regardless of whether or not \mathbf{X} is full rank.

- (b) **Bonus:** Suppose $\text{Var}(y_i) = \sigma^2$. Recall that because ridge regression is linear in \mathbf{Y} , $df_{\lambda} = \frac{1}{\sigma^2} \text{Tr}(\mathbf{H}_{\lambda})$, where $\hat{\mathbf{Y}}_{\lambda}^{(\text{ridge})} = \mathbf{X} \hat{\boldsymbol{\beta}}_{\lambda}^{(\text{ridge})}(\lambda) = \mathbf{H}_{\lambda} \mathbf{Y}$ (you should be able to write down \mathbf{H}_{λ} in terms of \mathbf{X} and λ). Show that for $\lambda_1 < \lambda_2$, $df_{\lambda_1} > df_{\lambda_2}$. That is, increasing λ reduces the effective number of parameters in the model.

3. Consider the data Fat.txt. Remove every tenth observation from the data for use as a test sample. Use the remaining data to fit (i.e. train) the following models where % body fat, `siri`, is the response and all other variables are predictors:

- (i) A simple linear model.
 - (ii) Ridge regression, where the tuning parameter λ is chosen with generalized cross validation (see the lecture slides from 10/22 and 2(b) above). See RidgeExample.R for an example of how to do this in R. Plot the generalized cross validation value as a function of λ so that the minimum value is clearly visible.
- (a) Which model has the smaller training error (see problem 1)? Why is training error a poor judge of how well the model will predict future data?
 - (b) Now use the models you fit in (i) and (ii) to predict the held out data. Which model performs better? Clearly indicate the metric (i.e. loss function) you used to judge model performance. (**Hint:** since you are using squared loss to choose λ , you should be using squared loss to judge prediction...)