

Linear regression

Machine Learning II
2021-2022 - UMONS
Souhaib Ben Taieb

1

Do Exercise 3.3 in LFD.

Exercise 3.3

Consider the hat matrix $H = X(X^T X)^{-1} X^T$, where X is an N by $d + 1$ matrix, and $X^T X$ is invertible.

- (a) Show that H is symmetric.
- (b) Show that $H^K = H$ for any positive integer K .
- (c) If I is the identity matrix of size N , show that $(I - H)^K = I - H$ for any positive integer K .
- (d) Show that $\text{trace}(H) = d + 1$, where the trace is the sum of diagonal elements. *[Hint: $\text{trace}(AB) = \text{trace}(BA)$.]*

Figure 1: Source: Abu-Mostafa et al. Learning from data. AMLbook.

Do Exercise 3.4 in LFD.

Exercise 3.4

Consider a noisy target $y = \mathbf{w}^{*\top} \mathbf{x} + \epsilon$ for generating the data, where ϵ is a noise term with zero mean and σ^2 variance, independently generated for every example (\mathbf{x}, y) . The expected error of the best possible linear fit to this target is thus σ^2 .

For the data $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$, denote the noise in y_n as ϵ_n and let $\boldsymbol{\epsilon} = [\epsilon_1, \epsilon_2, \dots, \epsilon_N]^\top$; assume that $\mathbf{X}^\top \mathbf{X}$ is invertible. By following

(continued on next page)

the steps below, show that the expected in-sample error of linear regression with respect to \mathcal{D} is given by

$$\mathbb{E}_{\mathcal{D}}[E_{\text{in}}(\mathbf{w}_{\text{lin}})] = \sigma^2 \left(1 - \frac{d+1}{N}\right).$$

- Show that the in-sample estimate of \mathbf{y} is given by $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}^* + \mathbf{H}\boldsymbol{\epsilon}$.
- Show that the in-sample error vector $\hat{\mathbf{y}} - \mathbf{y}$ can be expressed by a matrix times $\boldsymbol{\epsilon}$. What is the matrix?
- Express $E_{\text{in}}(\mathbf{w}_{\text{lin}})$ in terms of $\boldsymbol{\epsilon}$ using (b), and simplify the expression using Exercise 3.3(c).
- Prove that $\mathbb{E}_{\mathcal{D}}[E_{\text{in}}(\mathbf{w}_{\text{lin}})] = \sigma^2 \left(1 - \frac{d+1}{N}\right)$ using (c) and the independence of $\epsilon_1, \dots, \epsilon_N$. [Hint: The sum of the diagonal elements of a matrix (the trace) will play a role. See Exercise 3.3(d).]

For the expected out-of-sample error, we take a special case which is easy to analyze. Consider a test data set $\mathcal{D}_{\text{test}} = \{(\mathbf{x}_1, y'_1), \dots, (\mathbf{x}_N, y'_N)\}$, which shares the same input vectors \mathbf{x}_n with \mathcal{D} but with a different realization of the noise terms. Denote the noise in y'_n as ϵ'_n and let $\boldsymbol{\epsilon}' = [\epsilon'_1, \epsilon'_2, \dots, \epsilon'_N]^\top$. Define $E_{\text{test}}(\mathbf{w}_{\text{lin}})$ to be the average squared error on $\mathcal{D}_{\text{test}}$.

- Prove that $\mathbb{E}_{\mathcal{D}, \boldsymbol{\epsilon}'}[E_{\text{test}}(\mathbf{w}_{\text{lin}})] = \sigma^2 \left(1 + \frac{d+1}{N}\right)$.

The special test error E_{test} is a very restricted case of the general out-of-sample error. Some detailed analysis shows that similar results can be obtained for the general case, as shown in Problem 3.11.

Figure 2: Source: Abu-Mostafa et al. Learning from data. AMLbook.

3

Solve Problem 3.11 in LFD.

Problem 3.11 Consider the linear regression problem setup in Exercise 3.4, where the data comes from a genuine linear relationship with added noise. The noise for the different data points is assumed to be iid with zero mean and variance σ^2 . Assume that the 2nd moment matrix $\Sigma = \mathbb{E}_{\mathbf{x}}[\mathbf{x}\mathbf{x}^T]$ is non-singular. Follow the steps below to show that, with high probability, the out-of-sample error on average is

$$E_{\text{out}}(\mathbf{w}_{\text{lin}}) = \sigma^2 \left(1 + \frac{d+1}{N} + o\left(\frac{1}{N}\right) \right).$$

- (a) For a test point \mathbf{x} , show that the error $y - g(\mathbf{x})$ is

$$\epsilon - \mathbf{x}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon},$$

where ϵ is the noise realization for the test point and $\boldsymbol{\epsilon}$ is the vector of noise realizations on the data.

- (b) Take the expectation with respect to the test point, i.e., \mathbf{x} and ϵ , to obtain an expression for E_{out} . Show that

$$E_{\text{out}} = \sigma^2 + \text{trace} \left(\Sigma (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \right).$$

[Hints: $a = \text{trace}(a)$ for any scalar a ; $\text{trace}(\mathbf{AB}) = \text{trace}(\mathbf{BA})$; expectation and trace commute.]

- (c) What is $\mathbb{E}_{\boldsymbol{\epsilon}}[\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T]$?

- (d) Take the expectation with respect to $\boldsymbol{\epsilon}$ to show that, on average,

$$E_{\text{out}} = \sigma^2 + \frac{\sigma^2}{N} \text{trace} \left(\Sigma \left(\frac{1}{N} \mathbf{X}^T \mathbf{X} \right)^{-1} \right).$$

Note that $\frac{1}{N} \mathbf{X}^T \mathbf{X} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T$ is an N sample estimate of Σ . So $\frac{1}{N} \mathbf{X}^T \mathbf{X} \approx \Sigma$. If $\frac{1}{N} \mathbf{X}^T \mathbf{X} = \Sigma$, then what is E_{out} on average?

- (e) Show that (after taking the expectation over the data noise) with high probability,

$$E_{\text{out}} = \sigma^2 \left(1 + \frac{d+1}{N} + o\left(\frac{1}{N}\right) \right).$$

[Hint: By the law of large numbers $\frac{1}{N} \mathbf{X}^T \mathbf{X}$ converges in probability to Σ , and so by continuity of the inverse at Σ , $\left(\frac{1}{N} \mathbf{X}^T \mathbf{X} \right)^{-1}$ converges in probability to Σ^{-1} .]

Figure 3: Source: Abu-Mostafa et al. Learning from data. AMLbook.

4

Solve Problem 3.14 in LFD.

Problem 3.14 In a regression setting, assume the target function is linear, so $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}^*$, and $\mathbf{y} = \mathbf{Z}\mathbf{w}^* + \boldsymbol{\epsilon}$, where the entries in $\boldsymbol{\epsilon}$ are zero mean, iid with variance σ^2 . In this problem derive the bias and variance as follows.

- (a) Show that the average function is $\bar{g}(\mathbf{x}) = f(\mathbf{x})$, no matter what the size of the data set. What is the bias?
- (b) What is the variance? [*Hint: Problem 3.11*]

Figure 4: Source: Abu-Mostafa et al. Learning from data. AMLbook.

5

Solve Problem 3.15 in LFD.

Problem 3.15 In the text we derived that the linear regression solution weights must satisfy $X^T X \mathbf{w} = X^T \mathbf{y}$. If $X^T X$ is not invertible, the solution $\mathbf{w}_{\text{lin}} = (X^T X)^{-1} X^T \mathbf{y}$ won't work. In this event, there will be many solutions for \mathbf{w} that minimize E_{in} . Here, you will derive one such solution. Let ρ be the rank of X . Assume that the singular value decomposition (SVD) of X is $X = U \Gamma V^T$, where $U \in \mathbb{R}^{N \times \rho}$ satisfies $U^T U = I_\rho$, $V \in \mathbb{R}^{(d+1) \times \rho}$ satisfies $V^T V = I_\rho$, and $\Gamma \in \mathbb{R}^{\rho \times \rho}$ is a positive diagonal matrix.

- (a) Show that $\rho < d + 1$.
- (b) Show that $\mathbf{w}_{\text{lin}} = V \Gamma^{-1} U^T \mathbf{y}$ satisfies $X^T X \mathbf{w}_{\text{lin}} = X^T \mathbf{y}$, and hence is a solution.
- (c) Show that for any other solution that satisfies $X^T X \mathbf{w} = X^T \mathbf{y}$, $\|\mathbf{w}_{\text{lin}}\| < \|\mathbf{w}\|$. That is, the solution we have constructed is the minimum norm set of weights that minimizes E_{in} .

Figure 5: Source: Abu-Mostafa et al. Learning from data. AMLbook.