

Elements of Statistical Learning Notes and Exercise Solutions

Ryan Davis

May 12, 2019

Chapter 2

3 Linear Models for Regression

3.2 Linear Regression Models and Least Squares

Derive variance-covariance matrix

Derive the variance-covariance matrix of the least squares parameter estimates (equation 3.8 in book).

The least squares estimator is given by

$$\hat{\beta} = (X^T X)^{-1} X^T y \quad (1)$$

Express the variance of a random variable as below:

$$\text{Var}[X] = \mathbb{E}[X^2] + \mathbb{E}[X]^2 \quad (2)$$

$$\begin{aligned} \text{Var}[\hat{\beta}] &= \mathbb{E}[\hat{\beta}^2] + \mathbb{E}[\hat{\beta}]^2 \\ &= \mathbb{E}\left[\left((X^T X)^{-1} X^T y\right)^2\right] + \mathbb{E}\left[(X^T X)^{-1} X^T y\right]^2 \\ &= \mathbb{E}\left[\left((X^T X)^{-1} X^T (X\beta + \epsilon)\right)^2\right] + \mathbb{E}\left[(X^T X)^{-1} X^T (X\beta + \epsilon)\right]^2 \\ &= \mathbb{E}\left[\left((X^T X)^{-1} X^T (X\beta + \epsilon)\right)^2\right] + \mathbb{E}[\beta]^2 + \mathbb{E}\left[(X^T X)^{-1} X^T \epsilon\right]^2 \end{aligned}$$

The last term goes to zero since X is uncorrelated with the error term.

$$\text{Var}[\hat{\beta}] = \mathbb{E}\left[\left((X^T X)^{-1} X^T (X\beta + \epsilon)\right)^2\right] - \mathbb{E}[\beta]^2$$

The expectation of the true β is simply β itself.

$$\text{Var}[\hat{\beta}] = \mathbb{E}\left[\left((X^T X)^{-1} X^T (X\beta + \epsilon)\right)^2\right] - \beta^2$$

Break the term within the expectation into two separate terms.

$$\begin{aligned}
\text{Var} [\hat{\beta}] &= \mathbb{E} \left[\left((X^T X)^{-1} X^T X \beta + (X^T X)^{-1} X^T \epsilon \right)^2 \right] - \beta^2 \\
&= \mathbb{E} \left[\left(\beta + (X^T X)^{-1} X^T \epsilon \right)^2 \right] - \beta^2 \\
&= \mathbb{E} \left[\beta^2 + 2 (X^T X)^{-1} X^T \epsilon + \left((X^T X)^{-1} X^T \right)^2 \epsilon^2 \right] - \beta^2 \\
&= \beta^2 + \mathbb{E} \left[2 (X^T X)^{-1} X^T \epsilon + \left((X^T X)^{-1} X^T \right)^2 \epsilon^2 \right] - \beta^2 \\
&= \mathbb{E} \left[\left((X^T X)^{-1} X^T \right)^2 \epsilon^2 \right] \\
&= \mathbb{E} \left[\left((X^T X)^{-1} X^T \right)^2 \right] \mathbb{E} [\epsilon^2]
\end{aligned}$$

The error term has variance σ^2 .

$$\begin{aligned}
\text{Var} [\hat{\beta}] &= \sigma^2 \mathbb{E} \left[\left((X^T X)^{-1} X^T \right)^2 \right] \\
&= \sigma^2 (X^T X)^{-1}
\end{aligned}$$

Expected prediction error

Compute the expected prediction error at input x_0 for model $Y_0 = f(x_0) + \epsilon_0$. Use estimate $\hat{f}(x_0) = x_0^T \tilde{\beta}$.

$$\begin{aligned}
\mathbb{E} \left[\left(Y_0 - \tilde{f}(x_0) \right)^2 \right] &= \mathbb{E} \left[\left(f(x_0) + \epsilon_0 - \tilde{f}(x_0) \right)^2 \right] \\
&= \mathbb{E} \left[f(x_0)^2 - 2f(x_0) \hat{f}(x_0) - \tilde{f}(x_0)^2 \right] + \mathbb{E} [\epsilon_0^2] \\
&= \mathbb{E} \left[f(x_0)^2 - 2f(x_0) \hat{f}(x_0) - \tilde{f}(x_0)^2 \right] + \sigma^2 \\
&= \mathbb{E} \left[\left(f(x_0) - \tilde{f}(x_0) \right)^2 \right] + \sigma^2 \\
&= \sigma^2 + \text{MSE}(\tilde{f}(x_0))
\end{aligned}$$

Ridge regression coefficients

The loss function for ridge regression (page 64, equation 3.43) is

$$RSS(\lambda) = (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta \quad (3)$$

Differentiate with respect to β , set to zero, and solve.

$$\begin{aligned}\frac{\partial RSS}{\partial \beta} &= 2(y - X\beta)(-X) + 2\beta\lambda = 0 \\ &= -Xy + 2X^T X\beta + 2\beta = 0\end{aligned}$$

$$\begin{aligned}X^T y &= X^T X\beta + \beta \\ X^T y &= \beta(X^T X + I) \\ \hat{\beta} &= (X^T X + I)^{-1} X^T y\end{aligned}$$

Exercise 3.4

Show how the vector of least squares coefficients can be obtained from a single pass of the Gram-Schmidt procedure.

Represent X using the QR decomposition

$$X = QR \tag{4}$$

where Q is an $N \times (p + 1)$ orthogonal matrix and R is a $(p + 1) \times (p + 1)$ upper triangular matrix.

The vector of least squares coefficients can be represented as

$$\hat{\beta} = (X^T X)^{-1} X^T y \tag{5}$$

Substitute $X = QR$.

$$\begin{aligned}\hat{\beta} &= \left((QR)^T (QR) \right)^{-1} (QR)^T y \\ &= (Q^T Q R^T R)^{-1} Q^T R y\end{aligned}$$

Use the fact that Q is an orthogonal matrix ($Q^T Q = I$).

$$\begin{aligned}\hat{\beta} &= (R^T R)^{-1} Q^T R y \\ &= (R^T R)^{-1} Q^T R y\end{aligned}$$

Since $R^{-1} R = I$, we have

$$\hat{\beta} = R^{-1} Q^T y \tag{6}$$