



K. N. Toosi University of Technology

Faculty of Physics
Educational Group of
Atomic-Molecular and Astronomy

Answer to Question 2 (Minimizer in Linear Regression)

by
Ali Bagheri

Teacher
Dr. Mohammad Hossein Zhoolideh

Academic Year 1401-1402
(First Semester)

Question

Show that the minimizer for least-squares linear regression with L_2 regularization is $\mathbf{w} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$.

Note: To solve this question, you can use chapter 3 of the Pattern Recognition and Machine Learning (Bishop) book.

Answer

In Section 1.1 of the Pattern Recognition and Machine Learning (Bishop) book, the idea of adding an regularization term to an error function in order to control the overfitting is introduced, so that the total error function to be minimized takes the form

$$J(\mathbf{w}) = E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$$

where λ is the regularization coefficient that controls the relative importance of the data-dependent error ($E_D(\mathbf{w})$) and the regularization term ($E_W(\mathbf{w})$). One of the simplest forms of regularizer is given by the sum-of-squares of the weight vector elements

$$E_W(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

If we also consider the sum-of-squares error function given by

$$E_D(\mathbf{w}) = \frac{1}{2} (\mathbf{t} - \mathbf{w} \Phi)^2,$$

then the total error function becomes

$$\begin{aligned} J(\mathbf{w}) &= E_D(\mathbf{w}) + \lambda E_W(\mathbf{w}) \\ &= \frac{1}{2} (\mathbf{t} - \mathbf{w} \Phi)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \end{aligned}$$

Now we take the derivative of this equation in terms of \mathbf{w} and set it equal to zero

$$\begin{aligned} \frac{\partial}{\partial \mathbf{w}} J(\mathbf{w}) &= -(\mathbf{t} - \mathbf{w} \Phi) \Phi + \lambda \mathbf{w} \\ &= 0 \end{aligned}$$

As a result, we have

$$\begin{aligned} \mathbf{w} &= (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \\ &= (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \end{aligned}$$

A Proof of some important equations

$$A = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{pmatrix}_{n \times n}$$

$$B = \begin{pmatrix} b_{11} & \cdots & b_{1n} \\ \vdots & \ddots & \vdots \\ b_{n1} & \cdots & b_{nn} \end{pmatrix}_{n \times n}$$

$$AB = \begin{pmatrix} a_{11}b_{11} + \cdots + a_{1n}b_{n1} & \cdots & a_{11}b_{1n} + \cdots + a_{1n}b_{nn} \\ \vdots & \ddots & \vdots \\ a_{n1}b_{11} + \cdots + a_{nn}b_{n1} & \cdots & a_{n1}b_{1n} + \cdots + a_{nn}b_{nn} \end{pmatrix}_{n \times n}$$

$$= \begin{pmatrix} \sum_{j=1}^n a_{1j}b_{j1} & \cdots & \sum_{j=1}^n a_{1j}b_{jn} \\ \vdots & \ddots & \vdots \\ \sum_{j=1}^n a_{nj}b_{j1} & \cdots & \sum_{j=1}^n a_{nj}b_{jn} \end{pmatrix}_{n \times n}$$

$$\text{tr}(AB) = \sum_{i=1}^n \sum_{j=1}^n a_{ij}b_{ji} \equiv f(AB)$$

$$\nabla_A f(AB) = \begin{pmatrix} \frac{\partial f(AB)}{\partial a_{11}} & \cdots & \frac{\partial f(AB)}{\partial a_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f(AB)}{\partial a_{n1}} & \cdots & \frac{\partial f(AB)}{\partial a_{nn}} \end{pmatrix}$$

$$= \begin{pmatrix} b_{11} & \cdots & b_{n1} \\ \vdots & \ddots & \vdots \\ b_{1n} & \cdots & b_{nn} \end{pmatrix} = B^T$$

$$\rightarrow \nabla_A f(AB) = B^T \quad (1)$$

$$A = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{pmatrix}_{n \times n}$$

$$A^T = \begin{pmatrix} a_{11} & \dots & a_{n1} \\ \vdots & \ddots & \vdots \\ a_{1n} & \dots & a_{nn} \end{pmatrix}_{n \times n}$$

$$\nabla_{A^T} f(A) = \begin{pmatrix} \frac{\partial f(A)}{\partial a_{11}} & \dots & \frac{\partial f(A)}{\partial a_{n1}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f(A)}{\partial a_{1n}} & \dots & \frac{\partial f(A)}{\partial a_{nn}} \end{pmatrix}_{n \times n}$$

$$\nabla_A f(A) = \begin{pmatrix} \frac{\partial f(A)}{\partial a_{11}} & \dots & \frac{\partial f(A)}{\partial a_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f(A)}{\partial a_{n1}} & \dots & \frac{\partial f(A)}{\partial a_{nn}} \end{pmatrix}_{n \times n}$$

$$\rightarrow \nabla_{A^T} f(A) = (\nabla_A f(A))^T \quad (2)$$

In the same way, we can prove the following equation:

$$\nabla_A \text{tr}(ABA^T C) = CAB + C^T AB^T \quad (3)$$

Combining Equations (2) and (3), we find that

$$\begin{aligned} \nabla_{A^T} \text{tr}(ABA^T C) &= (\nabla_A \text{tr}(ABA^T C))^T \\ &= (CAB + C^T AB^T)^T \\ &= B^T A^T C^T + BA^T C \end{aligned} \quad (4)$$

$$\begin{aligned}
\nabla_{\theta} J(\theta) &= \nabla_{\theta} \frac{1}{2} (X_{(n \times m)} \theta_{(m \times 1)} - \vec{y}_{(n \times 1)})^T (X_{(n \times m)} \theta_{(m \times 1)} - \vec{y}_{(n \times 1)}) \\
&= \frac{1}{2} \nabla_{\theta} \left(\theta_{(1 \times m)}^T X_{(m \times n)}^T X_{(n \times m)} \theta_{(m \times 1)} \right. \\
&\quad - \theta_{(1 \times m)}^T X_{(m \times n)}^T \vec{y}_{(n \times 1)} \\
&\quad - \vec{y}_{(1 \times n)}^T X_{(n \times m)} \theta_{(m \times 1)} \\
&\quad \left. + \vec{y}_{(1 \times n)}^T \vec{y}_{(n \times 1)} \right)
\end{aligned} \tag{5}$$

The above equation is in the form of the following equation:

$$\frac{1}{2} \nabla_{\theta} \left(A_{(1 \times 1)} - B_{(1 \times 1)} - B_{(1 \times 1)}^T + C_{(1 \times 1)} \right)$$

In fact, the result of a 1×1 matrix is a polynomial expression, and since each polynomial expression is equal to its transpose and also the trace of each polynomial expression is equal to that expression itself, the equation can be written as follows:

$$\begin{aligned}
\nabla_{\theta} J(\theta) &= \frac{1}{2} \nabla_{\theta} \text{tr} (\theta^T X^T X \theta - 2 \vec{y}^T X \theta + \vec{y}^T \vec{y}) \\
&= \frac{1}{2} \nabla_{\theta} (\text{tr} (\theta^T X^T X \theta) - 2 \text{tr} (\vec{y}^T X \theta) + \text{tr} (\vec{y}^T \vec{y})) \\
&= \frac{1}{2} (\nabla_{\theta} \text{tr} (\theta^T X^T X \theta) - 2 \nabla_{\theta} \text{tr} (\vec{y}^T X \theta) + \nabla_{\theta} \text{tr} (\vec{y}^T \vec{y})) \\
&= \frac{1}{2} (\nabla_{\theta} \text{tr} (\theta^T X^T X \theta) - 2 \nabla_{\theta} \text{tr} (\vec{y}^T X \theta))
\end{aligned}$$

From equations (1) and (4), we know:

$$\begin{aligned}
\nabla_{A^T} \text{tr} (A B A^T C) &= B^T A^T C^T + B A^T C \\
\frac{A^T = \theta, B = B^T = X^T X, C = I}{\longrightarrow} \nabla_{\theta} \text{tr} (\theta^T X^T X \theta) &= X^T X \theta I + X^T X \theta I \\
&= X^T X \theta + X^T X \theta
\end{aligned}$$

$$\begin{aligned}
\nabla_A \text{tr} (A B) &= \nabla_A \text{tr} (B A) = B^T \\
\frac{A = \theta, B = \vec{y}^T X}{\longrightarrow} \nabla_{\theta} \text{tr} (\vec{y}^T X \theta) &= X^T \vec{y}
\end{aligned}$$

As a result, we have:

$$\begin{aligned}
\nabla_{\theta} J(\theta) &= \frac{1}{2} (X^T X \theta + X^T X \theta - 2 X^T \vec{y}) \\
&= X^T X \theta - X^T \vec{y}
\end{aligned} \tag{6}$$

By comparing equation (5) with equation (6), we find that in the matrix derivative, when we want to take a derivative with respect to θ , θ^T is the same as θ itself.