

1. **Multi-class Least Squares** In this section, you will determine the parameter matrix $\mathbf{W} \in \mathbb{R}^{m \times p}$ for the Multi-class Least Squares problem.

Given a data matrix $\mathbf{X} \in \mathbb{R}^{n \times m}$ and target matrix $\mathbf{T} \in \mathbb{R}^{n \times p}$, the sum-of-squares error function can be written as

$$J = Er(\mathbf{W}) = \text{Tr}\{(\mathbf{XW} - \mathbf{T})^T(\mathbf{XW} - \mathbf{T})\}$$

where Tr is the trace of a matrix. You can assume that \mathbf{X} has full rank.

We will solve this problem by setting the derivative with respect to \mathbf{W} to be zero and solve for \mathbf{W} . To do this we must first know some matrix derivative properties.

- (a) Let \mathbf{A}, \mathbf{Z} be two matrices. Prove

$$\begin{bmatrix} \frac{\partial \text{Tr}(\mathbf{AZ})}{\partial z_{11}} & \frac{\partial \text{Tr}(\mathbf{AZ})}{\partial z_{12}} & \dots & \dots \\ \frac{\partial \text{Tr}(\mathbf{AZ})}{\partial z_{21}} & \dots & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix} \Leftarrow \frac{d\text{Tr}(\mathbf{AZ})}{d\mathbf{Z}} = \mathbf{A}^T$$

$$\frac{\partial \text{Tr}(\mathbf{AZ})}{\partial z_{ij}} = \frac{\sum_k \sum_l z_{kl} a_{ji}}{\partial z_{ij}} = a_{ji}$$

$$\frac{\partial \text{Tr}(\mathbf{AZ})}{\partial \mathbf{Z}} = \mathbf{A}^T$$

Linear Regression (Least Square Problem)

$$x_1, \dots, x_n \in \mathbb{R}^m$$

$$t_1, \dots, t_n \in \mathbb{R}$$

$$\hat{t}_i = w^T x_i$$

$$\min_w \sum_{i=1}^n (w^T x_i - t_i)^2$$

$$\|Xw - t\|^2$$

Multi-class Least Square Regression $\begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{bmatrix}$

$$x_1, \dots, x_n \in \mathbb{R}^m$$

$$\begin{bmatrix} t_{11} \\ t_{12} \\ \vdots \\ t_{1p} \end{bmatrix} \in \mathbb{R}$$

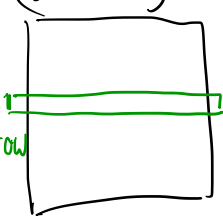
$$\begin{bmatrix} t_{n1} \\ t_{n2} \\ \vdots \\ t_{np} \end{bmatrix} \in \mathbb{R}$$

$$\hat{t}_{ij} = w_j^T x_i$$

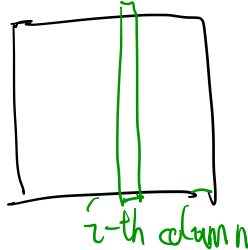
$$\min_{w_1, w_2, \dots, w_p} \sum_{i=1}^n \sum_{j=1}^p (w_j^T x_i - t_{ij})^2$$

matrix-vector form

$$\text{Tr}(A) = \sum_i A_{ii}$$

$$(XW - T)^T$$


i-row

$$XW - T$$


i-th column

$$\text{Tr} \{ (XW - T)^T (XW - T) \}$$

$$\begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix} [w_1, w_2, \dots, w_p] \begin{bmatrix} t_{11} & \dots & t_{1p} \\ t_{21} & \dots & t_{2p} \\ \vdots & & \vdots \\ t_{n1} & \dots & t_{np} \end{bmatrix}$$

$$\begin{bmatrix} x_1^T w_1 & x_1^T w_2 & \dots & x_1^T w_p \\ \vdots & \vdots & & \vdots \\ x_n^T w_1 & x_n^T w_2 & \dots & x_n^T w_p \end{bmatrix}$$

$m \times p$

$$\text{Tr} \{ (XW - T)^T (XW - T) \} = \sum_{i=1}^n \sum_{j=1}^p (w_j^T x_i - t_{ij})^2$$

(b) Let \mathbf{A}, \mathbf{Z} be two matrices. Prove

$$\begin{aligned} \frac{d \text{Tr}(\mathbf{Z} \mathbf{A} \mathbf{Z}^T)}{d \mathbf{Z}} &= \mathbf{Z} \mathbf{A}^T + \mathbf{Z} \mathbf{A} \\ \text{Tr}(\mathbf{Z} \mathbf{A} \mathbf{Z}^T) &= \sum_{i'} \sum_{j'} \sum_{k'} z_{i'j'} a_{j'k'} z_{i'k'} \\ &\quad \textcircled{1} i'=k' \quad \textcircled{2} j' \neq k' \\ &= \sum_{i'} \sum_{j'} z_{i'j'}^2 a_{j'j'} + \sum_{i'} \sum_{j'} \sum_{k' \neq j'} z_{i'j'} a_{j'k'} z_{i'k'} \\ \frac{\partial \text{Tr}(\mathbf{Z} \mathbf{A} \mathbf{Z}^T)}{\partial z_{ij}} &= 2 z_{ij} a_{jj} + \sum_{k' \neq j} a_{jk'} z_{ik'} + \sum_{k' \neq j} z_{ik'} a_{kj} \\ &= \sum_k z_{ik'} (a_{jk'} + a_{kj}) \\ \frac{d \text{Tr}(\mathbf{Z} \mathbf{A} \mathbf{Z}^T)}{d \mathbf{Z}} &= \mathbf{Z} \mathbf{A}^T + \mathbf{Z} \mathbf{A} \quad \frac{d \text{Tr}(\mathbf{Z}^T \mathbf{A} \mathbf{Z})}{d \mathbf{Z}} = \frac{d \text{Tr}(\mathbf{Z} \mathbf{A} \mathbf{Z}^T)}{d \mathbf{Z}^T} = \mathbf{A} \mathbf{Z} + \mathbf{A}^T \mathbf{Z} \end{aligned}$$

(c) Now, we can take the derivative of $\text{Er}(\mathbf{W})$ and set it to zero. Show that this results in

$$\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}$$

$$\begin{aligned} \nabla_{\mathbf{W}} L &= \nabla_{\mathbf{W}} \text{Tr} \{ (\mathbf{X} \mathbf{W} - \mathbf{T})^T (\mathbf{X} \mathbf{W} - \mathbf{T}) \} \\ &= \nabla_{\mathbf{W}} \{ \text{Tr} \{ \mathbf{W}^T \mathbf{X}^T \mathbf{X} \mathbf{W} \} - \text{Tr} \{ \mathbf{W}^T \mathbf{X}^T \mathbf{T} \} - \text{Tr} \{ \mathbf{T}^T \mathbf{X} \mathbf{W} \} + \text{Tr} \{ \mathbf{T}^T \mathbf{T} \} \} \\ &= 2 \mathbf{X}^T \mathbf{X} \mathbf{W} - 2 \mathbf{X}^T \mathbf{T} \end{aligned}$$

$$2 \mathbf{X}^T \mathbf{X} \mathbf{W} - 2 \mathbf{X}^T \mathbf{T} = 0$$

$$\Downarrow \\ \mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}$$

$$[w_1, w_2, \dots, w_p]$$

2. In this problem, we will derive the least square solution for multi-class classification. Consider a general classification problem with K classes, with a 1-of- K binary encoding scheme (defined latter) for the target vector $t, t \in \mathbb{R}^K$. Suppose we are given a training data set $\{x_n, t_n\}, n = 1, \dots, n$ where $x_n \in \mathbb{R}^D$. For the 1-of- K binary encoding scheme, t_n has the k -th element being 1 and all other elements being 0 if the n -th data is in class k . We can use the following linear model to describe each class:

$$y_k(x) = w_k^T x + w_{k0},$$

where $k = 1, \dots, K$. We can conveniently group these together using vector notation so that

$$y(x) = \tilde{\mathbf{W}}^T \tilde{x},$$

where $\tilde{\mathbf{W}}$ is a matrix whose k -th column comprises the $D + 1$ -dimensional vector $\tilde{w} = [w_{k0}, w_k^T]^T$ and \tilde{x} is the corresponding augmented input vector $[1, x^T]^T$. For each new input with feature x , we assign it to the class for which the output $y_k = \tilde{w}_k^T \tilde{x}$ is largest. Define a matrix \mathbf{T} whose n -th row is the vector t_n^T and together a matrix $\tilde{\mathbf{X}}$ whose n -th row is \tilde{x}_n^T , the sum-of-squares error function can be written as

$$J(\tilde{\mathbf{W}}) = \frac{1}{2} \text{Tr} \{ (\tilde{\mathbf{X}} \tilde{\mathbf{W}} - \mathbf{T})^T (\tilde{\mathbf{X}} \tilde{\mathbf{W}} - \mathbf{T}) \}.$$

Find the closed form solution of $\tilde{\mathbf{W}}$ that minimizes the objective function $J(\tilde{\mathbf{W}})$. Hint: You may use the following two matrix derivative about trace, $\frac{\partial}{\partial \mathbf{Z}} \text{Tr}(\mathbf{A} \mathbf{Z}) = \mathbf{A}^T$ and $\frac{\partial}{\partial \mathbf{Z}} \text{Tr}(\mathbf{Z}^T \mathbf{A} \mathbf{Z}) = (\mathbf{A}^T + \mathbf{A}) \mathbf{Z}$.

$$\begin{aligned}
 & x_1, x_2, \dots, x_n \in \mathbb{R}^m & y_i \in \{1, 2, \dots, p\} & \text{p-class classification} \\
 & X_i : [t_{i1} \ t_{i2} \ \dots \ t_{ip}]^T \\
 & t_{ij} = \begin{cases} 1 & \text{if data } i \text{ is in class } j \\ 0 & \text{otherwise} \end{cases} \\
 & W = (X^T X)^{-1} X^T T & p = K \\
 & \hat{T} = X_{\text{test}} W = \begin{bmatrix} \hat{t}_{11} & \hat{t}_{12} & \dots & \hat{t}_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix} & \begin{matrix} n \times p \\ p = K \end{matrix} \\
 & & & X_{i, \text{test}} \text{ is in } \underset{\text{class}}{\arg \max_j} \hat{t}_{ij}
 \end{aligned}$$