# CSCI 567: Discussion 3
## Linear Algebra II

- Practice Problems 3&4
- Intuition For Least Squares Solution
- Quadratic Form
- Eigenvalues & Eigenvectors
- Practice Problems 1&2

**Q3:** Given $f(\omega) = \|X\omega - y\|_2^2 + \omega^T M \omega$, $\quad (X \in \mathbb{R}^{n \times d}, y \in \mathbb{R}^n,$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad M \in \mathbb{R}^{d \times d}, P.D.)$

$\qquad$ Find $\omega^* = \arg\min_{\omega \in \mathbb{R}^d} f(\omega)$

$f(\omega) = \|X\omega - y\|_2^2 + \omega^T M \omega$

$\quad = (X\omega - y)^T (X\omega - y) + \omega^T M \omega$

$\quad = (X\omega)^T X\omega + y^T y - y^T X\omega - (X\omega)^T y + \omega^T M \omega$

$\quad = \omega^T X^T X\omega + y^T y - 2y^T X\omega + \omega^T M \omega$

$\quad = \omega^T \underbrace{(X^T X + M)}_{A} \omega - 2\underbrace{y^T X}_{b^T}\omega + y^T y$

$\qquad\qquad\qquad\qquad\qquad\qquad A = X^T X + M,$
$\qquad\qquad\qquad\qquad\qquad\qquad b = X^T y \quad —①$

$\quad = f_1(\omega) - 2 f_2(\omega) + y^T y$

$f_1(\omega) = \omega^T A \omega. \qquad f_2(\omega) = b^T \omega.$

$\dfrac{\partial f(\omega)}{\partial \omega} = \boxed{\dfrac{\partial f_1(\omega)}{\partial \omega}} - 2 \boxed{\dfrac{\partial f_2(\omega)}{\partial \omega}} = 0. \qquad\qquad —②.$

$$b^T \omega = \sum_{i=1}^{d} \omega_i b_i$$

**1.** $\dfrac{\partial f_2(\omega)}{\partial \omega}$ :

$$\frac{\partial f_2(\omega)}{\partial \omega_k} = b_k$$

$$\Rightarrow \quad \frac{\partial f_2(\omega)}{\partial \omega} = \begin{bmatrix} \dfrac{\partial f_2(\omega)}{\partial \omega_1} \\ \vdots \\ \dfrac{\partial f_2(\omega)}{\partial \omega_d} \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_d \end{bmatrix} = \boxed{b} \qquad ---\text{(iii)}$$

**2.** $\dfrac{\partial f_1(\omega)}{\partial \omega}$ : <u>TWO WAYS:</u>

$$v = A\omega$$
$$\omega^T A\omega = \sum_i \omega_i v_i$$
$$v_i = a_i^T \omega = \sum_j A_{ij}\omega_j$$

a) $f_1(\omega) = \displaystyle\sum_{i=1}^{d} \sum_{j=1}^{d} \omega_i A_{ij} \omega_j$

$$= \sum_{i=1}^{d} \omega_i^2 A_{ii} + \sum_{i=1, i\neq j}^{d} \sum_{j=1}^{d} \omega_i A_{ij} \omega_j$$

$$\frac{\partial f_1(\omega)}{\partial \omega_k} = 2\omega_k A_{kk} + \sum_{i=1, i\neq k}^{d} \omega_i A_{ik} + \sum_{j=1, j\neq k}^{d} A_{kj}\omega_j$$

$$= \sum_{i=1}^{d} A_{ik}\omega_i + \sum_{j=1}^{d} A_{kj}\omega_j$$

$$= \underset{\underset{\text{column vector}}{\downarrow}}{(a^k)^T \omega} + \underset{\underset{\text{row vector}}{\downarrow}}{a_k^T \omega}$$

$$\Rightarrow \quad \boxed{\frac{\partial f_1(\omega)}{\partial \omega}} = \underbrace{\begin{bmatrix} - (a^1)^T - \\ \vdots \\ - (a^d)^T - \end{bmatrix}}_{A^T} \omega + \underbrace{\begin{bmatrix} - a_1^T - \\ \vdots \\ - a_d^T - \end{bmatrix}}_{A} \omega$$

$$= \boxed{(A^T + A)\, \omega} \qquad\qquad ---\text{(iv)}$$

b) $f_1(\omega) = g^T(\omega)\,\omega, \quad g(\omega) = A^T\omega.$

$$\boxed{\frac{\partial f_1(\omega)}{\partial \omega}} = \left(\boxed{\frac{\partial g(\omega)}{\partial \omega}}\right)^T \omega + \left(\boxed{\frac{\partial \omega}{\partial \omega}}\right)^T g(\omega).$$

$$\boxed{\frac{\partial g(\omega)}{\partial \omega}} = \left[\frac{\partial g(\omega)}{\partial \omega_1} \cdots \frac{\partial g(\omega)}{\partial \omega_d}\right] = \begin{bmatrix} -\left(\frac{\partial g_1(\omega)}{\partial \omega}\right)^T - \\ \vdots \\ -\left(\frac{\partial g_d(\omega)}{\partial \omega}\right)^T - \end{bmatrix}$$

$$g_K(\omega) = (a^k)^T\omega \quad \Rightarrow \quad \frac{\partial g_K(\omega)}{\partial \omega} = a^k$$

$$\boxed{\frac{\partial g(\omega)}{\partial \omega}} = \begin{bmatrix} -a^{1T}- \\ \vdots \\ -a^{dT}- \end{bmatrix} = \boxed{A^T}$$

$$\boxed{\frac{\partial \omega}{\partial \omega}} = \begin{bmatrix} \frac{\partial \omega_1}{\partial \omega_1} & \frac{\partial \omega_1}{\partial \omega_2} & \cdots & \frac{\partial \omega_1}{\partial \omega_d} \\ \frac{\partial \omega_2}{\partial \omega_1} & \frac{\partial \omega_2}{\partial \omega_2} & \cdots & \vdots \\ \vdots & \vdots & \ddots & \\ \frac{\partial \omega_d}{\partial \omega_1} & \cdots & & \frac{\partial \omega_d}{\partial \omega_d} \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & \vdots \\ \vdots & \vdots & \ddots & \\ 0 & \cdots & \cdots & 1 \end{bmatrix} = \boxed{I_{d\times d}}$$

$$\boxed{\frac{\partial f_1(\omega)}{\partial \omega}} = A\omega + I^T A^T \omega = \boxed{(A + A^T)\omega} \qquad -\text{\textcircled{\scriptsize iv}}$$

$$\frac{\partial f(w)}{\partial w} = \boxed{\frac{\partial f_1(w)}{\partial w}} - 2 \boxed{\frac{\partial f_2(w)}{\partial w}} = 0.$$

Using ⟨III⟩, ⟨IV⟩ in ⟨I⟩:

$$\frac{\partial f(w)}{\partial w} = \boxed{(A + A^T) w} - 2\boxed{b} = 0.$$

$$\Rightarrow \quad w^* = \left(\frac{A + A^T}{2}\right)^{-1} b \qquad \text{(when } A + A^T \text{ is invertible)}$$

Using ⟨II⟩:  $A = X^T X + M$,  $b = X^T y$

$$\Rightarrow \quad A + A^T = 2 X^T X + M + M^T$$

$$\Rightarrow \quad \boxed{w^* = \left[ X^T X + \frac{1}{2}(M + M^T) \right]^{-1} (X^T y)}. \qquad \text{(Ans.)}$$

Special case:  $M = \lambda I_{d \times d}$

$$\Rightarrow \quad w^* = (X^T X + \lambda I)^{-1} (X^T y).$$

$$\lambda = 0$$

$$\Rightarrow \quad w^* = (X^T X)^{-1} X^T y.$$

**SHORT-ANSWER QUESTION.** *The following questions use linear algebra and calculus in ML formulations. They particularly test your knowledge of gradients of multivariate functions.*

**Q3** Consider the following optimization problem:

$$w_* = \arg\min_{w \in \mathbb{R}^d} \|Xw - y\|_2^2 + w^T M w$$

Here, $X \in \mathbb{R}^{n \times d}$, $y \in \mathbb{R}^n$, $M \in \mathbb{R}^{d \times d}$ is a positive definite matrix and $\|\cdot\|_2$ stands for the $\ell_2$ norm. Find the closed form solution for $w_*$. Proceed in a similar way as how we derived the general least-squares solution in class. (This optimization problem is a generalization of $\ell_2$ *regularization*, which we will see in class.)

*Answer:* Setting the gradient $2X^T(Xw - y) + (M + M^T)w$ to be $\mathbf{0}$ and using the fact that $M$ is invertible gives

$$w'_* = \left(X^T X + \frac{M + M^T}{2}\right)^{-1} X^T y.$$

**Q4** Assume we have a training set $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^d \times \mathbb{R}$, where each outcome $y_i$ is generated by a probabilistic model $w_*^T x_i + \epsilon_i$ with $\epsilon_i$ being an independent Gaussian noise with zero-mean and variance $\sigma^2$ for some $\sigma > 0$. In other words, the probability of seeing any outcome $y \in \mathbb{R}$ given $x_i \in \mathbb{R}^d$ is

$$\varepsilon_i = y_i - w_*^T x_i \qquad \Pr(y \mid x_i; w_*, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(\frac{-(y - w_*^T x_i)^2}{2\sigma^2}\right).$$

Assume $\sigma$ is fixed and given, find the maximum likelihood estimation for $w_*$. In other words, first write down the probability of seeing the outcomes $y_1, \dots, y_n$ given $x_1, \dots, x_n$ as a function of the value of $w_*$; then find the value of $w_*$ that maximizes this probability. You can assume $X^T X$ is invertible, where $X$ is the data matrix with each row corresponding to the features of an example. You may find it helpful to review the steps we took in Lecture 2 to find the maximum likelihood solution for the logistic model.

*Answer:* The probability of seeing the outcomes $y_1, \dots, y_n$ given $x_1, \dots, x_n$ for a linear model $w$ is

$$\mathcal{P}(w) = \prod_{i=1}^n \Pr(y_i \mid x_i; w, \sigma) = \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} \exp\left(\frac{-(y_i - w^T x_i)^2}{2\sigma^2}\right).$$

Taking the negative log, this becomes

$$F(w) = n \ln\sqrt{2\pi} + n\ln\sigma + \frac{1}{2\sigma^2}\sum_{i=1}^n (y_i - w^T x_i)^2 = n\ln\sqrt{2\pi} + n\ln\sigma + \frac{1}{2\sigma^2}\|Xw - y\|_2^2.$$

$$(Xw - y)_i$$
$$= x_i^T w - y_i$$

Maximizing $\mathcal{P}$ is the same as minimizing $F$, which is clearly the same as just minimizing $\|Xw - y\|_2^2$, the same objective as for least square regression. Therefore the MLE for $w_*$ is exactly the same as the least square solution:

$$w_* = (X^T X)^{-1} X^T y.$$