

## Writing Problems:

## 1.1 Matrix Derivation

For the matrix derivation, denominator-layout is adopted

## (2) Derivation by definition

① (5 pts)

$$f(w) = w^T A w$$

$$= [w_1, w_2, \dots, w_n] \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & & \ddots & \\ \vdots & & & \\ a_{n1} & & & a_{nn} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

$$= [a_{11}w_1 + a_{21}w_2 + \dots + a_{n1}w_n, \dots, a_{1n}w_1 + a_{2n}w_2 + \dots + a_{nn}w_n] \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

$$= a_{11}w_1^2 + a_{21}w_2w_1 + \dots + a_{n1}w_nw_1 + \dots + a_{1n}w_1w_n + a_{2n}w_2w_n + \dots + a_{nn}w_n^2$$

$$\frac{df}{dw} = \begin{bmatrix} a_{11}w_1 + a_{21}w_2 + \dots + a_{n1}w_n & a_{11}w_1 + a_{12}w_2 + \dots + a_{1n}w_n \\ a_{12}w_1 + a_{22}w_2 + \dots + a_{n2}w_n & a_{21}w_1 + a_{22}w_2 + \dots + a_{2n}w_n \\ \vdots & \vdots \\ a_{n1}w_1 + a_{n2}w_2 + \dots + a_{nn}w_n & a_{1n}w_1 + a_{2n}w_2 + \dots + a_{nn}w_n \end{bmatrix}_{n \times 1}$$

$$= \begin{bmatrix} a_{11}w_1 + a_{21}w_2 + \dots + a_{n1}w_n \\ a_{12}w_1 + a_{22}w_2 + \dots + a_{n2}w_n \\ \vdots \\ a_{n1}w_1 + a_{n2}w_2 + \dots + a_{nn}w_n \end{bmatrix} + \begin{bmatrix} a_{11}w_1 + a_{12}w_2 + \dots + a_{1n}w_n \\ a_{21}w_1 + a_{22}w_2 + \dots + a_{2n}w_n \\ \vdots \\ a_{1n}w_1 + a_{2n}w_2 + \dots + a_{nn}w_n \end{bmatrix}$$

$$= A^T w + A w$$

② (5 pts)

$$f(w) = A w$$

$$= \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & & \ddots & \\ \vdots & & & \\ a_{n1} & & & a_{nn} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

$$= \begin{bmatrix} a_{11}w_1 + a_{12}w_2 + \dots + a_{1n}w_n \\ a_{21}w_1 + a_{22}w_2 + \dots + a_{2n}w_n \\ \vdots \\ a_{n1}w_1 + a_{n2}w_2 + \dots + a_{nn}w_n \end{bmatrix}_{n \times 1}$$

$$\frac{df}{dw} = \begin{bmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \dots & a_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & a_{n3} & \dots & a_{nn} \end{bmatrix}$$

$$= A^T$$

(3) (Bonus track: 5 pts) Derivation by differentiation

① (2 pts)  $f(w) = w^T A w$  is a scalar

Let  $g(w) = w^T$ ,  $h(w) = A w$

Since  $d(xy) = (dx)y + x(dy)$ ,  $\text{tr}(A+B) = \text{tr}(A) + \text{tr}(B)$

$$\begin{aligned} d(w^T A w) &= (dw^T) A w + w^T (dA w) \\ &= (dw^T) A w + w^T \cdot dA \cdot w + w^T A (dw) \\ &= (dw^T) A w + w^T A (dw) \\ &= \text{tr}((dw^T) A w + w^T A (dw)) \\ &= \text{tr}((dw^T) A w) + \text{tr}(w^T A (dw)) \\ &= \text{tr}(w^T A^T (dw)) + \text{tr}(w^T A (dw)) \\ &= \text{tr}[w^T A^T (dw) + w^T A (dw)] \\ &= \text{tr}([(A^T + A)w]^T (dw)) \\ &= [(A^T + A)w]^T (dw) \end{aligned}$$

$$\text{Therefore } \frac{df}{dw} = \frac{d}{dw}(w^T A w) = (A^T + A)w$$

② (3 pts)  $f(W) = \text{tr}(W^T A W)$

Since  $d(\text{tr}(X)) = \text{tr}(dX)$ ,

$$df = d(\text{tr}(W^T A W)) = \text{tr}(d(W^T A W))$$

$$\begin{aligned} \text{tr}(d(W^T A W)) &= \text{tr}(d(W^T) \cdot A W + W^T d(A) W + W^T A (dW)) \\ &= \text{tr}(d(W^T) \cdot A W + W^T A (dW)) \\ &= \text{tr}(d(W^T) \cdot A W) + \text{tr}(W^T A (dW)) \\ &= \text{tr}(W^T A^T (dW)) + \text{tr}(W^T A (dW)) \\ &= \text{tr}([(A^T + A)W]^T (dW)) \end{aligned}$$

For the formula  $df = \text{tr}((\frac{df}{dW})^T dW)$ ,

$$\text{we get: } \frac{df}{dW} = (A^T + A)W.$$

(4) (a) (10 pts)

Given that  $\ell = \|Xw - y\|_2^2$ ,  $z = Xw - y$ .

$$\text{Since } \ell = z^T z, \frac{d\ell}{dz} = 2z.$$

$$\frac{\partial z}{\partial w^T} = \frac{\partial}{\partial w^T} (Xw - y)$$

$$= X$$

$$\begin{aligned}\Rightarrow \frac{dL}{dW} &= \left( \frac{\partial Z}{\partial W^T} \right)^T \cdot \frac{\partial L}{\partial Z} \\ &= X^T \cdot 2Z \\ &= 2X^T(XW - y)\end{aligned}$$

(b) (Bonus Task, 5 pts)

Suppose  $A \in \mathbb{R}^{m \times n}$ ,  $X \in \mathbb{R}^{n \times p}$ ,  $\Rightarrow B \in \mathbb{R}^{m \times p}$

$$A = \begin{bmatrix} A_{11} & \dots & A_{1n} \\ \vdots & \ddots & \vdots \\ A_{m1} & \dots & A_{mn} \end{bmatrix}, X = \begin{bmatrix} x_{11} & \dots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{np} \end{bmatrix}, B = \begin{bmatrix} B_{11} & \dots & B_{1p} \\ \vdots & \ddots & \vdots \\ B_{m1} & \dots & B_{mp} \end{bmatrix}$$

For each entry of  $Y$ ,

$$\begin{aligned}Y_{kl} &= \sum_s A_{ks} X_{sl} \\ \Rightarrow \frac{\partial Y_{kl}}{\partial X_{ij}} &= \frac{\partial \sum_s A_{ks} X_{sl}}{\partial X_{ij}} = \frac{\partial A_{ki} X_{il}}{\partial X_{ij}} = A_{ki} \delta_{lj}\end{aligned}$$

where  $\delta_{lj} = 1$  when  $l=j$ , otherwise  $\delta_{lj} = 0$ .

$$\text{Therefore } \frac{\partial L}{\partial X_{ij}} = \sum_k \frac{\partial L}{\partial Y_{kl}} (A_{ki} \delta_{lj}) = \sum_k \frac{\partial L}{\partial Y_{kj}} A_{ki}$$

$$\Rightarrow \frac{\partial L}{\partial X} = \begin{bmatrix} \sum_k \left( \frac{\partial L}{\partial Y_{k1}} \right) \cdot A_{k1} & \sum_k \left( \frac{\partial L}{\partial Y_{k1}} \right) \cdot A_{k2} & \dots & \sum_k \left( \frac{\partial L}{\partial Y_{k1}} \right) \cdot A_{kn} \\ \sum_k \left( \frac{\partial L}{\partial Y_{k2}} \right) \cdot A_{k1} & \sum_k \left( \frac{\partial L}{\partial Y_{k2}} \right) \cdot A_{k2} & \dots & \sum_k \left( \frac{\partial L}{\partial Y_{k2}} \right) \cdot A_{kn} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_k \left( \frac{\partial L}{\partial Y_{kp}} \right) \cdot A_{k1} & \sum_k \left( \frac{\partial L}{\partial Y_{kp}} \right) \cdot A_{k2} & \dots & \sum_k \left( \frac{\partial L}{\partial Y_{kp}} \right) \cdot A_{kn} \end{bmatrix}$$

## 1.2 Convexity of functions, 10 pts

Proof: (1) Proved by definition.

Without loss of generality, for  $x, y \in \mathbb{R}$ , assume  $x \leq y$ .

We have the following 3 situations:

$$\textcircled{1} x \leq y \leq 0$$

$$f(x) = f(y) = 0.$$

$$\forall \theta \in [0, 1], \theta x + (1-\theta)y \leq 0.$$

$$\text{Therefore } f(\theta x + (1-\theta)y) = 0 = \theta f(x) + (1-\theta)f(y) \quad \forall \theta \in [0, 1], x \leq y \leq 0.$$

$$\textcircled{2} x \leq 0 < y$$

$$f(x) = 0, f(y) = y > 0.$$

$$(i) \text{ If } \theta x + (1-\theta)y \leq 0,$$

$$\text{Then } f(\theta x + (1-\theta)y) = 0 < \theta f(x) + (1-\theta)f(y) = (1-\theta)y$$

$$(ii) \text{ If } \theta x + (1-\theta)y > 0,$$

$$f(\theta x + (1-\theta)y) = \theta x + (1-\theta)y < (1-\theta)y = \theta f(x) + (1-\theta)f(y)$$

Hence we conclude that: for  $x \leq 0 < y$ ,  $\forall \theta \in [0, 1]$ ,  
 $f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y)$

(3)  $0 < x \leq y$

$$f(x) = x. \quad f(y) = y.$$

$$\theta x + (1-\theta)y > 0.$$

$$f(\theta x + (1-\theta)y) = \theta x + (1-\theta)y = \theta f(x) + (1-\theta)f(y)$$

In conclusion,  $\forall x, y \in \mathbb{R}, \theta \in [0, 1]$ ,

$$f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y).$$

So  $f(x) = \max(0, x)$  is convex for  $x \in \mathbb{R}$ .

(2) Proved by second-order conditions.

$$f(x) = |x|. \Rightarrow f(x) = \begin{cases} x & \text{for } x \geq 0 \\ -x & \text{for } x < 0. \end{cases}$$

$$\nabla f(x) = \begin{cases} 1 & \text{for } x \geq 0 \\ -1 & \text{for } x < 0. \end{cases}$$

$$\nabla^2 f(x) = 0 \quad \text{for all } x \in \mathbb{R}.$$

By the second-order conditions,  $f(x) = |x|$  is convex.

(3) Proved by second-order conditions.

$$f(x) = \|Ax - b\|_2^2$$

$$\text{By 1.1(4)-(a), } \nabla f(x) = 2A^T(Ax - b)$$

$$\nabla^2 f(x) = 2A^T A \geq 0 \quad \text{for } \forall x \in \mathbb{R}^n$$

By the second-order conditions,  $f(x) = \|Ax - b\|_2^2$  is convex.

### 1.3 Gradient Descent, 10 pts

(1) Define  $f(W) = \text{tr}[(Y - XW)^T A (Y - XW)]$

where  $Y = (y_1, y_2, \dots, y_N)^T \in \mathbb{R}^{N \times k}$ ,  $X = [(x_1^T, 1)^T, (x_2^T, 1)^T, \dots, (x_N^T, 1)^T] \in \mathbb{R}^{N \times (d+1)}$   
 $W = (w, b)^T \in \mathbb{R}^{(d+1) \times k}$ ,  $A = \text{diag}(a_1, a_2, \dots, a_N)$ .

To get the closed-form solution, let  $\frac{df}{dW} = 0$

$$df = d(\text{tr}[(Y - XW)^T A (Y - XW)])$$

$$= \text{tr}[d((Y - XW)^T A (Y - XW))]$$

$$d((Y - XW)^T A (Y - XW))$$

$$= d((Y - XW)^T) A (Y - XW) + (Y - XW)^T d(A (Y - XW))$$

$$= (d(Y - XW))^T A (Y - XW) + (Y - XW)^T d(A (Y - XW))$$

$$\begin{aligned}
&= (d(-xw))^T A(Y-xw) + (Y-xw)^T A(d(Y-xw)) \\
&= [(Y-xw)^T A^T (d(Y-xw))]^T + [(Y-xw)^T A (d(Y-xw))]
\end{aligned}$$

$$\begin{aligned}
\Rightarrow df &= \text{tr}[d((Y-xw)^T A(Y-xw))] \\
&= \text{tr}([(Y-xw)^T A^T (d(Y-xw))]^T + [(Y-xw)^T A (d(Y-xw))]) \\
&= \text{tr}([(Y-xw)^T A^T + (Y-xw)^T A](d(Y-xw)))
\end{aligned}$$

$\Rightarrow$  Use differentiation method.

$$\begin{aligned}
\text{we get: } \frac{df}{d(Y-xw)} &= [(Y-xw)^T A^T + (Y-xw)^T A]^T \\
&= A(Y-xw) + A^T(Y-xw) \\
&= 2A(Y-xw) \quad \text{since } A^T = A.
\end{aligned}$$

Since  $\frac{d(-xw)}{dW^T} = -x$ , we get:

$$\begin{aligned}
\frac{df}{dW} &= \left(\frac{d(-xw)}{dW^T}\right)^T \cdot \frac{df}{d(Y-xw)} = 2(-x^T)A(Y-xw) = 0 \\
&\Rightarrow x^T A Y = x^T A x W \\
&\Rightarrow W = (x^T A x)^{-1} x^T A Y
\end{aligned}$$

(2) The goal is to minimize  $f(W) = \text{tr}[(Y-xw)^T A(Y-xw)]$

where  $Y = (y_1, y_2, \dots, y_N)^T \in \mathbb{R}^{N \times k}$ ,  $X = [(x_1^T, 1)^T, (x_2^T, 1)^T, \dots, (x_N^T, 1)^T] \in \mathbb{R}^{N \times (d+1)}$   
 $W = (w, b)^T \in \mathbb{R}^{(d+1) \times k}$   $A = \text{diag}(a_1, a_2, \dots, a_N)$ .

The step of gradient descent is as follows:

- ① Set an initial  $W \in \mathbb{R}^{(d+1) \times k}$
- ② Find the gradient matrix of  $f(W)$ :  $\frac{df}{dW}$
- ③ Choose  $\alpha$  s.t.  $\alpha = \arg\min_{\alpha} f(W + \alpha \cdot \frac{df}{dW})$
- ④ Set  $W = W + \alpha \cdot \frac{df}{dW}$ , repeat the above steps until  $\nabla f(W)$  is small enough.

## 1.4 Maximum Likelihood Estimation, 10 pts

Proof: The likelihood function is:

$$\begin{aligned}
L(\mu, \sigma^2) &= \prod_{n=1}^N (2\pi\sigma^2)^{-\frac{1}{2}} \cdot \exp\left(-\frac{(x_n - \mu)^2}{2\sigma^2}\right) \\
&= (2\pi\sigma^2)^{-\frac{N}{2}} \cdot \exp\left(-\frac{1}{2\sigma^2} \cdot \sum_{n=1}^N (x_n - \mu)^2\right)
\end{aligned}$$

$$\ell(\mu, \sigma^2) = \log L(\mu, \sigma^2)$$

$$= -\frac{N}{2} \cdot \log(2\pi\sigma^2) + \left(-\frac{1}{2\sigma^2}\right) \cdot \sum_{n=1}^N (X_n - \mu)^2.$$

$$\hat{\mu}_{MLE} := \arg \max_{\mu} \ell(\mu, \sigma^2)$$

$$\hat{\sigma}_{MLE}^2 := \arg \max_{\sigma^2} \ell(\mu, \sigma^2)$$

$$\Rightarrow \left. \frac{\partial f}{\partial \mu} \right|_{\mu = \hat{\mu}_{MLE}} = -\frac{1}{2\sigma^2} \cdot \sum_{n=1}^N (X_n - \mu) \cdot (-2) \Big|_{\mu = \hat{\mu}_{MLE}} = 0$$

$$\text{Therefore } \sum_{n=1}^N (X_n - \hat{\mu}_{MLE}) = 0 \Rightarrow \hat{\mu}_{MLE} = \frac{1}{N} \sum_{n=1}^N X_n$$

$$\left. \frac{\partial f}{\partial \sigma^2} \right|_{\sigma^2 = \hat{\sigma}_{MLE}^2, \mu = \hat{\mu}_{MLE}} = \left(-\frac{N}{2}\right) \cdot \frac{2\pi}{2\pi\sigma^2} + \left(-\frac{1}{2} \sum_{n=1}^N (X_n - \mu)^2\right) \cdot \frac{1}{\sigma^4} (-1)$$

$$\Rightarrow N = \frac{\sum_{n=1}^N (X_n - \hat{\mu}_{MLE})^2}{\hat{\sigma}_{MLE}^2}$$

$$\text{Therefore } \hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{n=1}^N (X_n - \hat{\mu}_{MLE})^2.$$