

1. (2)

① $f(w) = w^T A w$ given $w \in \mathbb{R}^n$, $A \in \mathbb{R}^{n \times n}$

$$W = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} \quad A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \quad W^T = [w_1, \dots, w_n]$$

$$f(w) = [w_1, \dots, w_n] \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}$$

$$= \left[\sum_{j=1}^n w_j a_{1j}, \dots, \sum_{j=1}^n w_j a_{nj} \right] \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}$$

$$= w_1 \sum_{i=1}^n w_2 a_{i1} + w_2 \sum_{i=1}^n w_2 a_{i2} + \dots + w_n \sum_{i=1}^n w_i a_{in}$$

$$= \sum_{j=1}^n w_j \sum_{i=1}^n w_i a_{ij} = \sum_{i=1}^n \sum_{j=1}^n w_i w_j a_{ji} \quad \text{r.c.}$$

since f is a scalar function, w is a vector

the size of $\frac{df}{dw}$ will be the same as w , i.e., R^n

$$\text{since } \frac{df}{dw_{ir}} = \left(\sum_{j=1}^n w_j a_{ji} + w_i a_{ir} \right) + \sum_{\substack{j=1 \\ j \neq i}}^n w_j a_{ij} = \sum_{j=1}^n w_j a_{ji} + \sum_{j=1}^n w_j a_{ij}$$

Therefore, $\frac{df}{dw} = A^T w + Aw$.

we can examine the size of $\frac{\partial f}{\partial w}$, which is R^n

② $f(w) = Aw$ given $w \in \mathbb{R}^n$, $A \in \mathbb{R}^{m \times n}$.

$$A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix} \quad w = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}$$

Let $f(w) = y$. We can tell that y is a vector of size \mathbb{R}^m .

So $\frac{\partial f}{\partial w}$ will be a matrix of size $\mathbb{R}^{m \times n}$

$$\begin{aligned} f(w) &= \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} \\ &= \begin{bmatrix} \sum_{j=1}^n a_{1j} w_j \\ \sum_{j=1}^n a_{2j} w_j \\ \vdots \\ \sum_{j=1}^n a_{mj} w_j \end{bmatrix} = y \end{aligned}$$

$$\text{since } \frac{\partial y_i}{\partial w_j} = a_{ij}$$

$$\text{therefore, } \frac{\partial f}{\partial w} = A$$

we can examine the correctness of the answer by the size.

1. (3) $f(w) = w^T A w$ given $A \in \mathbb{R}^{n \times n}$ $w \in \mathbb{R}^n$

① $A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix}$ $w = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}$

$$\begin{aligned} df &= d(w^T A w) = d(w^T) A w + w^T (d(A w)) \\ &= d(w^T) A w + w^T (d(A) w + A d(w)) \\ &= d(w^T) A w + w^T A d(w) \end{aligned}$$

since f is a scalar function

$$df = \text{tr}(dy)$$

$$\begin{aligned} df &= \text{tr}(d(w^T) A w + w^T A d(w)) \\ &= \text{tr}(d(w^T) A w) + \text{tr}(w^T A d(w)) \end{aligned}$$

since $\text{tr}(X^T) = \text{tr}(X)$

$$\text{tr}(XY) = \text{tr}(YX)$$

$$\begin{aligned} \text{therefore, } \text{tr}(d(w^T) A w) &= \text{tr}((dw)^T A w) \\ &= \text{tr}(A w)^T dw \end{aligned}$$

$$\text{tr}(w^T A dw) = \text{tr}((A^T w)^T dw)$$

$$df = \text{tr}((A w + A^T w)^T dw)$$

$$\frac{df}{dw} = A w + A^T w$$

$$\textcircled{2} \quad f(W) = \text{tr}(W^T A W) \quad \text{given } W \in \mathbb{R}^{m \times n}, A \in \mathbb{R}^{m \times m}$$

$$A = \begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & & \vdots \\ a_{m1} & \dots & a_{mm} \end{bmatrix} \quad W = \begin{bmatrix} w_1 & \dots & w_n \\ \vdots & & \vdots \\ w_{m1} & \dots & w_{mn} \end{bmatrix}$$

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

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

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

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

$$= \text{tr}((dW)^T A W + W^T A dW)$$

$$= \text{tr}((dW)^T A W) + \text{tr}(W^T A dW)$$

$$= \text{tr}((AW)^T dW) + \text{tr}(W^T A dW)$$

$$= \text{tr}((AW + A^T W)^T dW)$$

$$\text{therefore, } \frac{df}{dW} = AW + A^T W.$$

$$(4) \quad \text{Since } \frac{dl}{dw} = \left(\frac{dz}{dw^T} \right)^T \frac{dl}{dz}$$

a.

$$z = Xw - y.$$

$$\frac{dz}{dw} = X^T \quad \frac{dz}{dw^T} = X$$

$$l = z^T z$$

$$\frac{dl}{dz} = 2z$$

$$\begin{aligned} \text{Therefore: } \frac{dl}{dw} &= X^T \cdot 2(Xw - y) \\ &= 2X^T \cdot Xw - 2X^T y. \end{aligned}$$

(4). b. since $\frac{\partial Y_{kl}}{\partial X_{ij}} = \frac{\partial \sum_s A_{ks} X_{sl}}{\partial X_{ij}} = \frac{\partial A_{ki} X_{il}}{\partial X_{ij}} = A_{ki} \delta_{lj}$

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

$$\frac{\partial L}{\partial X_{ij}} = \sum_{kl} \frac{\partial L}{\partial Y_{kl}} A_{ki} \delta_{lj} = \sum_k \frac{\partial L}{\partial Y_{kj}} A_{ki} = A^T_{:,i} \left(\frac{\partial L}{\partial Y} \right)_{:,j}$$

from $\frac{\partial L}{\partial X_{ij}}$, we can derive $\frac{\partial L}{\partial X}$

$$\frac{\partial L}{\partial X} = A \frac{\partial L}{\partial Y}$$

1.2.

$$(1). f(x) = \text{Relu}(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases} \quad x \in \mathbb{R}. \quad \nabla f(x) = \begin{cases} 1 & x > 0 \\ 0 & x < 0 \end{cases}$$

f is not differentiable.

By the definition of convexity,

$$f(\theta x + (1-\theta)y) = \begin{cases} \theta x + (1-\theta)y & \text{if } \theta x + (1-\theta)y \geq 0 \\ 0 & \theta x + (1-\theta)y < 0 \end{cases}$$
$$\theta \in [0, 1]$$

$$\text{since } f(x) \geq x \quad f(y) \geq y$$

Therefore, if $\theta x + (1-\theta)y \geq 0$,

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

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

$$\text{if } \theta x + (1-\theta)y < 0$$

$$\text{since } f(x) \geq 0, f(y) \geq 0$$

$$\theta f(x) + (1-\theta)f(y) \geq 0 = f(\theta x + (1-\theta)y)$$

Therefore, $\text{Relu}(x)$ is convex.

$$(2) f(x) = |x| = \begin{cases} x & x \geq 0 \\ -x & x < 0 \end{cases}$$

f is not differentiable at $x=0$.

By the definition of convexity,

$$f(\theta x + (1-\theta)y) = \begin{cases} \theta x + (1-\theta)y & \text{if } \theta x + (1-\theta)y \geq 0 \\ (\theta-1)y - \theta x & \text{if } \theta x + (1-\theta)y < 0 \end{cases}$$

$$\theta \in [0, 1]$$

$$\text{since } f(x) \geq x \quad f(y) \geq y$$

Therefore, if $\theta x + (1-\theta)y \geq 0$.

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

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

$$\text{if } \theta x + (1-\theta)y < 0$$

$$\textcircled{1} \text{ if } x < 0, y > 0.$$

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

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

$$\textcircled{2} \text{ if } x < 0, y < 0,$$

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

③ if $x > 0, y < 0$,

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

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

Therefore, $f(x) = |x|$ is convex for $x \in \mathbb{R}$.

(3). $f(x) = \|Ax - b\|_2^2$ dom $f \in \mathbb{R}$. given $A \in \mathbb{R}^{m \times n}$ $x \in \mathbb{R}^n$

f is twice differentiable

By 2nd-order conditions, f is convex iff $\nabla^2 f(x) \succeq 0$.

$$\nabla f(x) = 2A^T(Ax - b)$$

$$\nabla^2 f(x) = 2A^T A \quad (\text{from the note})$$

notice that $A^T A$ is non-negative

therefore, $\nabla^2 f(x) \succeq 0$

$\Rightarrow f(x) = \|Ax - b\|_2^2$ is convex.

$$1.3, \quad \sum_{i=1}^N \alpha_i \|y_i - Wx_i - b\|_2^2$$

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

$$Y = (y_1, y_2, \dots, y_N)^T \in \mathbb{R}^{N \times k}$$

$$X = [(x_1^T, 1), (x_2^T, 1), \dots, (x_N^T, 1)^T]^T \in \mathbb{R}^{N \times (d+1)}$$

$$W = (w, b)^T \in \mathbb{R}^{(d+1) \times k}$$

$$A = \text{diag}(\alpha_1, \alpha_2, \dots, \alpha_N).$$

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

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

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

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

$$W^T (X^T A X) W$$

$$\text{let } \frac{d}{dW} J(W) = 0$$

$$-(X^T A Y)^T - (Y^T A X)^T + (X^T A X) W + (X^T A^T X) W = 0.$$

$$(X^T A X + X^T A^T X) W = Y^T A^T X + X^T A^T Y$$

$$W = (X^T A X + X^T A^T X)^{-1} (Y^T A^T X + X^T A^T Y)$$

(v) given that

$$\frac{\partial J(w)}{\partial w} = -(X^T A Y)^T - (Y^T A X)^T + (X^T A X)w + (X^T A^T X)w$$

where

$$Y = (y_1, y_2, \dots, y_N)^T \in \mathbb{R}^{N \times k}$$

$$X = [(x_1^T, 1), (x_2^T, 1), \dots, (x_N^T, 1)]^T \in \mathbb{R}^{N \times (d+1)}$$

$$w = (w, b)^T \in \mathbb{R}^{(d+1) \times k}$$

$$A = \text{diag}(\alpha_1, \alpha_2, \dots, \alpha_N).$$

step 1: initialize w and b , α (learning rate)

step 2: iterate:

$$J(w) = wX.$$

$$w \leftarrow w - \alpha \frac{\partial J(w)}{\partial w}$$

if breaking criteria satisfy

$$(\text{e.g. } |\alpha \frac{\partial J(w)}{\partial w}| < \epsilon_0)$$

then

save w, b .

break.

1.4.

since $X_i \sim N(\mu, \sigma^2)$

$$\text{set } J(\omega) = \prod_{i=1}^N \frac{1}{\sqrt{\pi\sigma^2}} \exp\left(-\frac{(X_i - \mu)^2}{2\sigma^2}\right)$$

when $\frac{dJ(\omega)}{d\sigma^2} = 0$, we can get the optimal σ_{MLE}^2

when $\frac{dJ(\omega)}{d\mu} = 0$, we can get the optimal μ_{MLE} .

for simplicity, we make $z = \sigma^2$

$$\begin{aligned} J(\omega) &= \prod_{i=1}^N (2\pi z)^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(X_i - \mu)^2 z^{-1}\right) \\ &= (2\pi z)^{-\frac{N}{2}} \prod_{i=1}^N \exp\left(-\frac{1}{2}(X_i - \mu)^2 z^{-1}\right) \end{aligned}$$

$$\ln J(\omega) = -\frac{N}{2} \ln 2\pi z - z^{-1} \sum_{i=1}^N \frac{1}{2}(X_i - \mu)^2$$

$$\frac{d \ln J(\omega)}{d\mu} = -2z^{-1} \sum_{i=1}^N \frac{1}{2}(X_i - \mu) = 0$$

$$\Rightarrow \mu_{MLE} = \frac{1}{N} \sum_{i=1}^N X_i$$

$$\frac{\partial \ln J(\omega)}{\partial z} = -\frac{N}{2} \cdot z^{-1} + z^{-2} \sum_{i=1}^N \frac{1}{2} (X_i - \mu)^2 = 0.$$

$$\frac{N}{2} \cdot z^{-1} = z^{-2} \sum_{i=1}^N \frac{1}{2} (X_i - \mu)^2$$

$$z = \frac{1}{N} \sum_{i=1}^N (X_i - \mu)^2$$

When $\mu = \mu_{MLE}$

$$\sigma_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (X_i - \mu_{MLE})^2$$