1.1.(2>0

$w^T A w = A_{11} w_1 w_1 + A_{12} w_1 w_2 + \cdots + A_{1n} w_1 w_n$

$\qquad + \cdots$

$\qquad + A_{n1} w_n w_1 + A_{n2} w_n w_2 + \cdots + A_{mn} w_n w_n.$

$$\frac{df}{dw} = \begin{bmatrix} \frac{w^T A w}{\partial w_1} \\ \frac{w^T A w}{\partial w_2} \\ \vdots \\ \frac{w^T A w}{\partial w_n} \end{bmatrix} = \begin{bmatrix} (A_{11} w_1 + A_{12} w_2 + \cdots + A_{1n} w_n) + (A_{11} w_1 + A_{21} w_2 + \cdots + A_{n1} w_n) \\ \vdots \\ (A_{n1} w_n + A_{n2} w_1 + \cdots + A_{nn} w_n) + (A_{1n} w_1 + A_{2n} w_2 + \cdots + A_{nn} w_n) \end{bmatrix}$$

$$= \begin{bmatrix} A_{11} + A_{12} + \cdots + A_{1n} \\ \vdots \quad \vdots \quad \vdots \quad \vdots \\ A_{n1} + A_{n2} + \cdots + A_{nn} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} + \begin{bmatrix} A_{11} + A_{21} + \cdots + A_{n1} \\ \vdots \quad \vdots \quad \vdots \quad \vdots \\ A_{1n} + A_{2n} + \cdots + A_{nn} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}$$

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

②. $Aw = \begin{bmatrix} A_{11} w_1 & A_{12} w_2 & \cdots & A_{1n} w_n \\ \vdots & \vdots & \vdots & \vdots \\ A_{m1} w_1 & A_{mn} w_1 & \cdots & A_{mn} w_n \end{bmatrix}$

$$\frac{df}{dw} = \left[ \frac{\partial A_{ij} w_j}{\partial w_j} \right] = [A_{ij}] = A.$$

1.1.(3)

① $\dfrac{df}{dw} = \dfrac{d(w^T A)}{dw} \cdot w + w^T \dfrac{d(Aw)}{dw}$

$\qquad = Aw + w^T A$

$\qquad = (A + A^T) w$

② $\dfrac{df}{dw} = tr(d(w^T A w))$

$\qquad = tr[(A + A^T) w]$

1.1.(4) (a) $\dfrac{\partial l}{\partial w} = \left( \dfrac{\partial z}{\partial w^T} \right)^T \dfrac{\partial l}{\partial z}$

$\qquad = (x)^T \cdot 2z$

$\qquad = 2 x^T z$

1.2 ⑴ $f(\lambda x_1 + (1-\lambda) x_2):$ $\qquad (\lambda \in (0,1))$

1) $x_1, x_2$ share same Positivity/Negativity

As Relu is linear at $x \leq 0$ or $x \geq 0$, it is easy to find that:

$\qquad f(\lambda x_1 + (1-\lambda) x_2) = \lambda f(x_1) + (1-\lambda) f(x_2)$

2) Take $x_1 < 0$, $x_2 < 0$ for example.

$\qquad f(\lambda x_1 + (1-\lambda) x_2) \leq f((1-\lambda) x_2)$

$\qquad\qquad\qquad = (1-\lambda) f(x_2) + \lambda f(x_1)$   which is 0.

$\qquad$ As $\lambda x_1 + (1-\lambda) x_2 < (1-\lambda) x_2$

⟹ As a conclusion, $f(\lambda x_1 + (1-\lambda) x_2) \leq \lambda f(x_1) + (1-\lambda) x_2.$

$\qquad\qquad\qquad$ it is convex

2)

1) For $x_1, x_2$ with same Positivity/Negativity

Apparently $f(\lambda x_1 + (1-\lambda)x_2) = \lambda f(x_1) + (1-\lambda)x_2$

2) take $x_1 < 0, \quad x_2 \geq 0$

$$f(\lambda x_1 + (1-\lambda)x_2) = |\lambda x_1 + (1-\lambda)x_2| \leq |\lambda x_1| + |(1-\lambda)x_2|$$
$$= \lambda |x_1| + (1-\lambda)|x_2|$$
$$= \lambda f(x_1) + (1-\lambda)f(x_2)$$

As a conclusion, it is convex

3) $f(x) = (Ax-b)^T(Ax-b)$

$$= (x^TA^T - b^T)(Ax-b)$$
$$= x^TA^TAx - x^TA^Tb - b^TAx + b^Tb$$

$\frac{\partial f(x)}{\partial x} = 2A^TAx - 2A^Tb$

$\frac{\partial^2 f(x)}{\partial^2 x} = 2A^TA$ , $A^TA$ is a square matrix $\Rightarrow$ symmetric.

Let's prove $A^TA$ is semidefinite $\Leftrightarrow$ Prove $\forall$ vector $V$, $V^T(A^TA)V \geq 0$

$V^T(A^TA)V = 2(AV)^T(AV)$

$\Rightarrow \|AV\|_2^2 \geq 0 \Rightarrow$ Proved.

$\Rightarrow \frac{\partial^2 f(x)}{\partial^2 x} = 2A^TA$ is semidefinite $\Rightarrow f(x)$ is convex

1.3 (1) $tr[(Y-xw)^TA(Y-xw)]$

$= tr[Y^TAY - Y^TAxw - w^Tx^TAY + w^Tx^TAxw]$

$= tr[Y^TAY] - 2tr[w^Tx^TAY] + tr[w^Tx^TAxw]$

Taking derivative W.R.T $w$:

$\cdots = \frac{\partial}{\partial w}[-2tr[w^Tx^TAY] + tr[w^Tx^TAxw]]$

$= -2x^TAY + 2x^TAxw$

let $\frac{\partial}{\partial w} = 0 \Rightarrow x^TAxw = x^TAY \Rightarrow w = (x^TAx)^{-1}x^TAY$.

(2) $\frac{\partial}{\partial w} = -2x^TAY + 2x^TAxw$

$\Rightarrow w^{(t+1)} = w^{(t)} - \alpha(-2x^TAY + 2x^TAxw^{(t)})$ , $\alpha$ is the learn rate.

1.4

Likelihood func for given data:

$L(\mu, \sigma^2) = \prod_{n=1}^{N} \frac{1}{\sqrt{2\pi\sigma^2}} exp(-\frac{(x_n-\mu)^2}{2\sigma^2})$

$\ln L(\mu, \sigma^2) = -\frac{N}{2}\ln(2\pi) - N\ln\sigma - \sum_{n=1}^{N}\frac{(x_n-\mu)^2}{2\sigma^2}$

$\frac{\partial}{\partial\sigma^2}\ln L(\mu,\sigma^2) = -\frac{N}{\sigma^2} + \sum_{n=1}^{N}\frac{(x_n-\mu)^2}{2(\sigma^2)^2} = 0$

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