

R12945060

Problem 1 (Preliminary) (1 pt)

(a) (0.2 pts)

(i) (0.1 pts) Given $\mathbf{x}, \mathbf{a} \in \mathbb{R}^n$. Show that

$$\frac{\partial \|\mathbf{x} - \mathbf{a}\|_2}{\partial \mathbf{x}} = \frac{\mathbf{x} - \mathbf{a}}{\|\mathbf{x} - \mathbf{a}\|_2}.$$

補: $\|\mathbf{x} - \mathbf{a}\|_2 = \sqrt{(x_1 - a_1)^2 + (x_2 - a_2)^2 + \dots + (x_n - a_n)^2}$
 $= \sqrt{(\mathbf{x} - \mathbf{a})^\top (\mathbf{x} - \mathbf{a})}$

$$y = \|\mathbf{x} - \mathbf{a}\|_2 = \sqrt{(\mathbf{x} - \mathbf{a})^\top (\mathbf{x} - \mathbf{a})}$$

$$\text{令 } S = \mathbf{x} - \mathbf{a} \Rightarrow y = \sqrt{S^\top S}$$

$$\frac{\partial y}{\partial S} = \frac{1}{2} (S^\top S)^{1/2} \cdot 2S = \frac{S}{\|S\|_2}$$

$$\text{鏈式: } \frac{\partial S}{\partial \mathbf{x}} = \frac{\partial (\mathbf{x} - \mathbf{a})}{\partial \mathbf{x}} = I \quad (I \text{為 } n \times n \text{ 幖單位矩阵})$$

$$\frac{\partial y}{\partial \mathbf{x}} = \frac{\partial y}{\partial S} \frac{\partial S}{\partial \mathbf{x}} = \frac{S}{\|S\|_2} \cdot I = \frac{\mathbf{x} - \mathbf{a}}{\|\mathbf{x} - \mathbf{a}\|_2}$$

$$\text{得 } \frac{\partial \|\mathbf{x} - \mathbf{a}\|_2}{\partial \mathbf{x}} = \frac{\mathbf{x} - \mathbf{a}}{\|\mathbf{x} - \mathbf{a}\|_2}$$

(ii) (0.1 pts) Given $\mathbf{a} \in \mathbb{R}^m$, $\mathbf{X} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^n$. Show that

$$\frac{\partial \mathbf{a}^\top \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^\top.$$

$$y = \mathbf{a}^\top \mathbf{X} \mathbf{b} \Rightarrow y = \sum_{i=1}^m \sum_{j=1}^n a_i x_{ij} b_j$$

\mathbf{X} 矩阵中每個元素

$$\frac{\partial}{\partial x_{ij}} \left(\sum_{i=1}^m \sum_{j=1}^n a_i x_{ij} b_j \right) = a_i b_j$$

$$\frac{\partial \mathbf{a}^\top \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^\top$$

(b) (0.2 pts) Let $\mathbf{X} \in \mathbb{R}^{n \times n}$. Show that

$$\frac{\partial \det(\mathbf{X})}{\partial \mathbf{X}} = \det(\mathbf{X}) (\mathbf{X}^{-1})^T.$$

Hint: Recall the cofactor matrix

$$\mathbf{C} = \begin{bmatrix} C_{11} & \cdots & C_{1n} \\ \vdots & \ddots & \vdots \\ C_{n1} & \cdots & C_{nn} \end{bmatrix}$$

where $C_{ij} = (-1)^{i+j} M_{ij}$ and $M_{ij} = \det((x_{mn})_{m \neq i, n \neq j})$. The adjoint matrix is the transpose of the cofactor matrix

$$\text{adj}(\mathbf{X}) = \mathbf{C}^T.$$

We have an identity

$$\mathbf{X} \text{adj}(\mathbf{X}) = \det(\mathbf{X}).$$

You may check Wikipedia for more details.

$$\frac{\partial \det(\mathbf{X})}{\partial X_{ij}} = \text{cofactor}_{ji}(\mathbf{X})$$

$$\mathbf{X}^{-1} = \frac{1}{\det(\mathbf{X})} \text{adj}(\mathbf{X})$$

$$\frac{\partial \det(\mathbf{X})}{\partial \mathbf{X}} = \text{adj}(\mathbf{X}) = \mathbf{C}^T$$

$$\frac{\partial \det(\mathbf{X})}{\partial \mathbf{X}} = \det(\mathbf{X})(\mathbf{X}^T)^T$$

$$\mathbf{X} \cdot \text{adj}(\mathbf{X}) = \det(\mathbf{X}) \mathbf{I}$$

(c) (0.6 pts) Prove that

$$\frac{\partial \log(\det(\mathbf{A}))}{\partial a_{ij}} = \mathbf{e}_j^T \mathbf{A}^{-1} \mathbf{e}_i, \quad (1)$$

where $\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mm} \end{bmatrix} \in \mathbb{R}^{m \times m}$ is a (non-singular) matrix, and \mathbf{e}_j is the unit vector

along the j -th axis (e.g. $\mathbf{e}_3 = [0, 0, 1, 0, \dots, 0]^T$). It is common to write (1) as

$$\frac{\partial \log(\det(\mathbf{A}))}{\partial \mathbf{A}} = (\mathbf{A}^{-1})^T$$

Hint: Same as (b).

$$\frac{\partial \det(\mathbf{A})}{\partial \mathbf{A}} = \det(\mathbf{A})(\mathbf{A}^{-1})^T$$

$$\frac{\partial \log(\det(\mathbf{A}))}{\partial \mathbf{A}} = \frac{1}{\det(\mathbf{A})} \frac{\partial \det(\mathbf{A})}{\partial \mathbf{A}}$$

$$\frac{\partial \log(\det(\mathbf{A}))}{\partial \mathbf{A}} = (\mathbf{A}^{-1})^T$$

$$\frac{\partial \log(\det(\mathbf{A}))}{\partial a_{ij}} = (\mathbf{A}^{-1})_{ji}$$

$$(\mathbf{A}^{-1})_{ji} = \mathbf{e}_j^T \mathbf{A}^{-1} \mathbf{e}_i \Rightarrow \frac{\partial \log(\det(\mathbf{A}))}{\partial a_{ij}} = \mathbf{e}_j^T \mathbf{A}^{-1} \mathbf{e}_i$$

Problem 2 (Classification with Gaussian Mixture Model) (2.4 pts)

In this question, we tackle the binary classification problem through the generative approach, where we assume the data point X (viewed as a \mathbb{R}^d -valued r.v.) and its label Y (viewed as a $\{\mathcal{C}_1, \mathcal{C}_2\}$ -valued r.v.) are generated according to the generative model (parameterized by θ) as follows:

$$\mathbb{P}_\theta[X = \mathbf{x}, Y = \mathcal{C}_k] = \pi_k f_{\mu_k, \Sigma_k}(\mathbf{x}) \quad (k \in \{1, 2\}) \quad (2)$$

where $\theta = (\pi_1, \pi_2, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma_1, \Sigma_2)$ for which

$$f_{\mu_k, \Sigma_k}(\mathbf{x}) = \frac{1}{(2\pi)^{d/2}} \frac{1}{|\Sigma_k|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right).$$

Now suppose we observe data points $\mathbf{x}_1, \dots, \mathbf{x}_N$ and their corresponding labels y_1, \dots, y_N .

(a) (1.2 pt)

- (i) (0.3 pt) Please write down the likelihood function $L(\theta)$ that describes how likely the generative model would generate the observed data $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ in terms of $\theta = (\pi_1, \pi_2, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma_1, \Sigma_2)$.
 - (ii) (0.3 pt) Find the maximum likelihood estimate $\theta^* = (\pi_1^*, \pi_2^*, \boldsymbol{\mu}_1^*, \boldsymbol{\mu}_2^*, \Sigma_1^*, \Sigma_2^*)$ that maximizes the likelihood function $L(\theta)$.
 - (iii) (0.3 pt) Write down $\mathbb{P}_\theta[Y = \mathcal{C}_1 | X = \mathbf{x}]$ and $\mathbb{P}_\theta[X = \mathbf{x} | Y = \mathcal{C}_1]$ in terms of $\theta = (\pi_1, \pi_2, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma_1, \Sigma_2)$. What are the physical meaning of the aforementioned quantities?
 - (iv) (0.3 pt) Express $\mathbb{P}_\theta[Y = \mathcal{C}_1 | X = \mathbf{x}]$ in the form of $\sigma(z)$, where $\sigma(\cdot)$ denotes the sigmoid function, and express z in terms of $\theta = (\pi_1, \pi_2, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma_1, \Sigma_2)$ and x .
- (b) (1.2 pt) Suppose we pose an additional constraint that the covariance matrices of the two Gaussian distributions are identical, namely $\Sigma_1 = \Sigma_2 = \Sigma$, in which the generative model is parameterized by $\vartheta = (\pi_1, \pi_2, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma)$. Redo questions (a) under such setting.

$$(a)(i) \quad L(\theta) = \prod_{i=1}^N P_\theta[X = \mathbf{x}_i, Y = y_i]$$

$$\rightarrow \text{使用 } P_\theta[X = \mathbf{x}, Y = \mathcal{C}_k] \text{ 展開} \quad L(\theta) = \prod_{i=1}^N (\pi_1 f_{\mu_1, \Sigma_1}(\mathbf{x}_i))^{1(y_i=1)} (\pi_2 f_{\mu_2, \Sigma_2}(\mathbf{x}_i))^{1(y_i=2)}$$

$$(ii) \quad \log L(\theta) = \sum_{i=1}^N (1(y_i=1) \log(\pi_1 f_{\mu_1, \Sigma_1}(\mathbf{x}_i)) + 1(y_i=2) \log(\pi_2 f_{\mu_2, \Sigma_2}(\mathbf{x}_i)))$$

$$\Rightarrow \log L(\theta) = \sum_{i=1}^N 1(y_i=1) (\log \pi_1 + \log f_{\mu_1, \Sigma_1}(\mathbf{x}_i)) + \sum_{i=1}^N 1(y_i=2) (\log \pi_2 + \log f_{\mu_2, \Sigma_2}(\mathbf{x}_i))$$

$$\pi_1^* = \frac{1}{N} \sum_{i=1}^N 1(y_i=1), \quad \pi_2^* = \frac{1}{N} \sum_{i=1}^N 1(y_i=2)$$

$$\boldsymbol{\mu}_1^* = \frac{\sum_{i=1}^N 1(y_i=1) \mathbf{x}_i}{\sum_{i=1}^N 1(y_i=1)}, \quad \boldsymbol{\mu}_2^* = \frac{\sum_{i=1}^N 1(y_i=2) \mathbf{x}_i}{\sum_{i=1}^N 1(y_i=2)}$$

$$\boldsymbol{\Sigma}_1^* = \frac{\sum_{i=1}^N 1(y_i=1) (\mathbf{x}_i - \boldsymbol{\mu}_1^*) (\mathbf{x}_i - \boldsymbol{\mu}_1^*)^\top}{\sum_{i=1}^N 1(y_i=1)}, \quad \boldsymbol{\Sigma}_2^* = \frac{\sum_{i=1}^N 1(y_i=2) (\mathbf{x}_i - \boldsymbol{\mu}_2^*) (\mathbf{x}_i - \boldsymbol{\mu}_2^*)^\top}{\sum_{i=1}^N 1(y_i=2)}$$

$$(iii) \text{ (1)} P_0[Y=C_1 | X=x] = \frac{P_0[X=x, Y=C_1]}{P_0[X=x]}$$

$$P_0[X=x, Y=C_1] = \pi_1 f_{\mu_1, \Sigma_1}(x), \quad P_0[X=x] = \pi_1 f_{\mu_1, \Sigma_1}(x) + \pi_2 f_{\mu_2, \Sigma_2}(x)$$

$$P_0[Y=C_1 | X=x] = \frac{\pi_1 f_{\mu_1, \Sigma_1}(x)}{\pi_1 f_{\mu_1, \Sigma_1}(x) + \pi_2 f_{\mu_2, \Sigma_2}(x)}$$

$$(2) P_0[X=x | Y=C_1] = f_{\mu_1, \Sigma_1}(x) = \frac{1}{(2\pi)^{n/2} |\Sigma_1|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1)\right)$$

在類別 C_1 中，數據點 X 的生成過程

$P_0[Y=C_1 | X=x] \Rightarrow$ 在給定數據點 X 的情況下，它屬於類別 C_1 的概率
 \Rightarrow 後驗，在觀察到數據 X 後，對其所屬類別的信念

$P_0[X=x | Y=C_1] \Rightarrow$ 在數據屬於類別 C_1 的條件下，生成數據點 X 的概率
 \Rightarrow 生成式概率，在給定類別的情況下數據的分佈

(iv)

$$P_0[Y=C_1 | X=x] = \sigma(z) = \frac{1}{1+e^{-z}} = \frac{1}{1+\frac{\pi_1 f_{\mu_1, \Sigma_1}(x)}{\pi_2 f_{\mu_2, \Sigma_2}(x)}}$$

取對數：

$$P_0[Y=C_1 | X=x] = \frac{1}{1+e^{-z}} \quad (z = \log \frac{\pi_1 f_{\mu_1, \Sigma_1}(x)}{\pi_2 f_{\mu_2, \Sigma_2}(x)})$$

$$z = \log\left(\frac{\pi_1}{\pi_2}\right) + \log\left(\frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}}\right) - \frac{1}{2} [(x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1) - (x-\mu_2)^T \Sigma_2^{-1} (x-\mu_2)]$$

$$P_0[Y=C_1 | X=x] = \sigma(z)$$

$$z \approx \log\left(\frac{\pi_1}{\pi_2}\right) + \log\left(\frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}}\right) - \frac{1}{2} [(x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1) - (x-\mu_2)^T \Sigma_2^{-1} (x-\mu_2)]$$

$$(b) (i) P_{\theta}[X=x, Y=C_k] = \pi_k f_{M_k, \Sigma}(x), \quad k \in \{1, 2\}$$

$$f_{M_k, \Sigma}(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-M_k)^T \Sigma^{-1}(x-M_k)\right) \Rightarrow \text{統一}\Sigma$$

$$\Rightarrow L(\theta) = \prod_{i=1}^N (\pi_1 f_{M_1, \Sigma}(x_i))^{I(Y_i=1)} (\pi_2 f_{M_2, \Sigma}(x_i))^{I(Y_i=2)}$$

$$(ii) \log L(\theta) = \sum_{i=1}^N I(Y_i=1)(\log \pi_1 + \log f_{M_1, \Sigma}(x_i)) + \sum_{i=1}^N I(Y_i=2)(\log \pi_2 + \log f_{M_2, \Sigma}(x_i))$$

$$\pi_1^* = \frac{1}{N} \sum_{i=1}^N I(Y_i=1), \quad \pi_2^* = \frac{1}{N} \sum_{i=1}^N I(Y_i=2), \quad M_1^* = \frac{\sum_{i=1}^N I(Y_i=1)x_i}{\sum_{i=1}^N I(Y_i=1)}, \quad M_2^* = \frac{\sum_{i=1}^N I(Y_i=2)x_i}{\sum_{i=1}^N I(Y_i=2)}$$

$$\Sigma^* = \frac{\sum_{i=1}^N I(Y_i=1)(x_i - M_1^*)(x_i - M_1^*)^T + \sum_{i=1}^N I(Y_i=2)(x_i - M_2^*)(x_i - M_2^*)^T}{N}$$

$$(iii) P_{\theta}[Y=C_1 | X=x] = \frac{\pi_1 f_{M_1, \Sigma}(x)}{\pi_1 f_{M_1, \Sigma}(x) + \pi_2 f_{M_2, \Sigma}(x)}$$

$$(iv) P_{\theta}[Y=C_1 | X=x] = \sigma(z)$$

$$z = \log\left(\frac{\pi_1}{\pi_2}\right) - \frac{1}{2} \left[(x - M_1)^T \Sigma^{-1} (x - M_1) - (x - M_2)^T \Sigma^{-1} (x - M_2) \right]$$

Problem 3 (Closed-Form Linear Regression Solution) (1 pts + Bonus 1.5 pts)

Consider the linear regression model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\epsilon},$$

where $\mathbf{y} \in \mathbb{R}^n$, $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\boldsymbol{\theta} \in \mathbb{R}^d$ and $\boldsymbol{\epsilon} \in \mathbb{R}^n$. Denote $\mathbf{X}_i \in \mathbb{R}^{1 \times d}$ as the i -th row of \mathbf{X} , with the following interpretations:

- If the linear model has the bias term, then write $\boldsymbol{\theta} = [w_1, \dots, w_m, b]^T$ and $\mathbf{X}_i = [x_{i,1}, x_{i,2}, \dots, x_{i,m}, 1]$, namely $d = m + 1$.
- If the linear model has no bias term, then write $\boldsymbol{\theta} = [w_1, \dots, w_d]^T$ and $\mathbf{X}_i = [x_{i,1}, x_{i,2}, \dots, x_{i,m}]$, namely $d = m$.

- (a) Without the bias term, consider the L^2 -regularized loss function:

$$\sum_i \kappa_i (y_i - \mathbf{X}_i \boldsymbol{\theta})^2 + \lambda \sum_j w_j^2, \quad \lambda > 0.$$

Show that the optimal solution that minimizes the loss function is $\boldsymbol{\theta}^* = (\mathbf{X}^T \mathbf{K} \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{K} \mathbf{y}$, where

$$\mathbf{K} = \begin{bmatrix} \kappa_1 & & 0 \\ & \ddots & \\ 0 & & \kappa_n \end{bmatrix}$$

is a diagonal matrix and \mathbf{I} is the $d \times d$ identical matrix.

$$L(\boldsymbol{\theta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T \mathbf{K} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}^T \boldsymbol{\theta}$$

展开

$$\Rightarrow L(\boldsymbol{\theta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T \mathbf{K} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}^T \boldsymbol{\theta} = \mathbf{y}^T \mathbf{K} \mathbf{y} - 2\mathbf{y}^T \mathbf{K} \mathbf{X} \boldsymbol{\theta} + \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{K} \mathbf{X} \boldsymbol{\theta} + \lambda \boldsymbol{\theta}^T \boldsymbol{\theta}$$

$$\text{对 } \boldsymbol{\theta} \text{ 求导: } \frac{\partial L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = -2\mathbf{X}^T \mathbf{K} \mathbf{y} + 2\mathbf{X}^T \mathbf{K} \mathbf{X} \boldsymbol{\theta} + 2\lambda \boldsymbol{\theta}$$

最小化 loss function, 梯度设为 0:

$$-2\mathbf{X}^T \mathbf{K} \mathbf{y} + 2\mathbf{X}^T \mathbf{K} \mathbf{X} \boldsymbol{\theta} + 2\lambda \boldsymbol{\theta} = 0 \Rightarrow \mathbf{X}^T \mathbf{K} \mathbf{X} \boldsymbol{\theta} + \lambda \boldsymbol{\theta} = \mathbf{X}^T \mathbf{K} \mathbf{y}$$

$$\Rightarrow (\mathbf{X}^T \mathbf{K} \mathbf{X} + \lambda \mathbf{I}) \boldsymbol{\theta} = \mathbf{X}^T \mathbf{K} \mathbf{y} \Rightarrow \boldsymbol{\theta} = (\mathbf{X}^T \mathbf{K} \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{K} \mathbf{y}$$

(b) (Bonus, 1.5 pts) With the bias term, the L^2 -regularized loss function becomes

$$\sum_i \kappa_i (y_i - \mathbf{X}_i \boldsymbol{\theta})^2 + \lambda \sum_j w_j^2, \quad \lambda > 0.$$

Show that the optimal solution that minimizes the loss function is $\boldsymbol{\theta}^* = [\mathbf{w}^{*T}, b^*]^T$, where

$$\begin{aligned} \mathbf{w}^* &= \left(\tilde{\mathbf{X}}^T \mathbf{K} \tilde{\mathbf{X}} + \lambda \mathbf{I} - \frac{1}{\text{Tr}(\mathbf{K})} \tilde{\mathbf{X}}^T \mathbf{K} \mathbf{e} \mathbf{e}^T \mathbf{K} \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T \mathbf{K} \left(\mathbf{y} - \frac{1}{\text{Tr}(\mathbf{K})} \mathbf{e} \mathbf{e}^T \mathbf{K} \mathbf{y} \right), \\ b^* &= \frac{1}{\text{Tr}(\mathbf{K})} (\mathbf{e}^T \mathbf{K} \mathbf{y} - \mathbf{e}^T \mathbf{K} \tilde{\mathbf{X}} \mathbf{w}^*) \end{aligned}$$

for which $\mathbf{e} = [1 \dots 1]^T$ denotes the all one vector, $\mathbf{X} = [\tilde{\mathbf{X}} \mathbf{e}]$, $\text{Tr}(\mathbf{K})$ is the trace of the matrix \mathbf{K} , and that \mathbf{K} and \mathbf{I} are defined as in (a).

$$\begin{aligned} L(\boldsymbol{\theta}) &= \sum_i \kappa_i (y_i - \tilde{\mathbf{X}}_i \mathbf{w} - b)^2 + \lambda \mathbf{w}^T \mathbf{w} \\ \Rightarrow L(\mathbf{w}, b) &= (\mathbf{y} - \tilde{\mathbf{X}} \mathbf{w} - \mathbf{b} \mathbf{e})^T \mathbf{K} (\mathbf{y} - \tilde{\mathbf{X}} \mathbf{w} - \mathbf{b} \mathbf{e}) + \lambda \mathbf{w}^T \mathbf{w} \end{aligned}$$

對 \mathbf{w} 和 b 求梯度：

$$\text{對 } \mathbf{w} \text{ 求導數: } \frac{\partial L(\mathbf{w}, b)}{\partial \mathbf{w}} = -2 \tilde{\mathbf{X}}^T \mathbf{K} (\mathbf{y} - \tilde{\mathbf{X}} \mathbf{w} - \mathbf{b} \mathbf{e}) + 2 \lambda \mathbf{w}$$

$$\text{令等於 } 0 \quad \tilde{\mathbf{X}}^T \mathbf{K} \tilde{\mathbf{X}} \mathbf{w} + \lambda \mathbf{w} = \tilde{\mathbf{X}}^T \mathbf{K} (\mathbf{y} - \mathbf{b} \mathbf{e})$$

$$\text{對 } b \text{ 求導數: } \frac{\partial L(\mathbf{w}, b)}{\partial b} = -2 \mathbf{e}^T \mathbf{K} (\mathbf{y} - \tilde{\mathbf{X}} \mathbf{w} - \mathbf{b} \mathbf{e})$$

$$\text{令等於 } 0 \quad \mathbf{e}^T \mathbf{K} \mathbf{e} \mathbf{b} = \mathbf{e}^T \mathbf{K} (\mathbf{y} - \tilde{\mathbf{X}} \mathbf{w})$$

$$\mathbf{w}^* = (\tilde{\mathbf{X}}^T \mathbf{K} \tilde{\mathbf{X}} + \lambda \mathbf{I} - \frac{1}{\text{Tr}(\mathbf{K})} \tilde{\mathbf{X}}^T \mathbf{K} \mathbf{e} \mathbf{e}^T \mathbf{K} \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{K} (\mathbf{y} - \frac{1}{\text{Tr}(\mathbf{K})} \mathbf{e}^T \mathbf{K} \mathbf{y})$$

$$b^* = \frac{1}{\text{Tr}(\mathbf{K})} (\mathbf{e}^T \mathbf{K} \mathbf{y} - \mathbf{e}^T \mathbf{K} \tilde{\mathbf{X}} \mathbf{w}^*)$$

$$\boldsymbol{\theta}^* = [\mathbf{w}^{*T}, b^*]^T$$

Problem 4 (Noise and Regularization) (1 pts)

Consider the linear model $f_{\mathbf{w}, b} : \mathbb{R}^k \rightarrow \mathbb{R}$, where $\mathbf{w} \in \mathbb{R}^k$ and $b \in \mathbb{R}$, defined as

$$f_{\mathbf{w}, b}(x) = \mathbf{w}^T \mathbf{x} + b$$

Given dataset $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, if the inputs $\mathbf{x}_i \in \mathbb{R}^k$ are contaminated with input noise $\boldsymbol{\eta}_i \in \mathbb{R}^k$, we may consider the expected sum-of-squares loss in the presence of input noise as

$$\tilde{L}_{ss}(\mathbf{w}, b) = \mathbb{E} \left[\frac{1}{2N} \sum_{i=1}^N (f_{\mathbf{w}, b}(\mathbf{x}_i + \boldsymbol{\eta}_i) - y_i)^2 \right]$$

where the expectation is taken over the randomness of input noises $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_N$. Additionally, the inputs (\mathbf{x}_i) and the input noise $(\boldsymbol{\eta}_i)$ are independent.

Now assume the input noises $\boldsymbol{\eta}_i = [\eta_{i,1}, \eta_{i,2}, \dots, \eta_{i,k}]^T$ are random vectors with zero mean $\mathbb{E}[\eta_{i,j}] = 0$, and the covariance between components is given by

$$\mathbb{E}[\eta_{i,j}\eta_{i',j'}] = \delta_{i,i'}\delta_{j,j'}\sigma^2$$

where $\delta_{i,i'} = \begin{cases} 1 & \text{if } i = i' \\ 0 & \text{otherwise.} \end{cases}$ denotes the Kronecker delta.

Please show that

$$\tilde{L}_{ss}(\mathbf{w}, b) = \frac{1}{2N} \sum_{i=1}^N (f_{\mathbf{w}, b}(\mathbf{x}_i) - y_i)^2 + \frac{\sigma^2}{2} \|\mathbf{w}\|^2$$

That is, minimizing the expected sum-of-squares loss in the presence of input noise is equivalent to minimizing noise-free sum-of-squares loss with the addition of a L^2 -regularization term on the weights.
(Hint: $\|\mathbf{x}\|^2 = \mathbf{x}^T \mathbf{x} = \text{tr}(\mathbf{x} \mathbf{x}^T)$ and the square of a vector is dot product with itself)

loss function: $\tilde{L}_{ss}(\mathbf{w}, b) = \mathbb{E} \left[\frac{1}{2N} \sum_{i=1}^N ((f_{\mathbf{w}, b}(\mathbf{x}_i + \boldsymbol{\eta}_i) - y_i)^2) \right]$

展開 $f_{\mathbf{w}, b}(\mathbf{x}_i + \boldsymbol{\eta}_i)$:

$$f_{\mathbf{w}, b}(\mathbf{x}_i + \boldsymbol{\eta}_i) = \mathbf{w}^T (\mathbf{x}_i + \boldsymbol{\eta}_i) + b = \mathbf{w}^T \mathbf{x}_i + \mathbf{w}^T \boldsymbol{\eta}_i + b$$

$$f_{\mathbf{w}, b}(\mathbf{x}_i + \boldsymbol{\eta}_i) - y_i = (\mathbf{w}^T \mathbf{x}_i + b - y_i) + \mathbf{w}^T \boldsymbol{\eta}_i$$

$$\tilde{L}_{ss}(\mathbf{w}, b) = \mathbb{E} \left[\frac{1}{2N} \sum_{i=1}^N ((\mathbf{w}^T \mathbf{x}_i + b - y_i) + \mathbf{w}^T \boldsymbol{\eta}_i)^2 \right]$$

展開平方項:

$$\tilde{L}_{ss}(\mathbf{w}, b) = \frac{1}{2N} \sum_{i=1}^N \mathbb{E} \left[((\mathbf{w}^T \mathbf{x}_i + b - y_i)^2 + 2(\mathbf{w}^T \mathbf{x}_i + b - y_i)(\mathbf{w}^T \boldsymbol{\eta}_i) + (\mathbf{w}^T \boldsymbol{\eta}_i)^2) \right]$$

① 無 noise 的 loss function :

$$\mathbb{E}[(\mathbf{w}^T \mathbf{x}_i + b - y_i)^2] = (\mathbf{w}^T \mathbf{x}_i + b - y_i)^2$$

② 由於 $\boldsymbol{\eta}_i$ 均值為 0, $\therefore \mathbb{E}[2(\mathbf{w}^T \mathbf{x}_i + b - y_i)(\mathbf{w}^T \boldsymbol{\eta}_i)] = 0$

③ $\mathbb{E}[(\mathbf{w}^T \boldsymbol{\eta}_i)^2] = \mathbb{E}[\mathbf{w}^T \boldsymbol{\eta}_i \boldsymbol{\eta}_i^T \mathbf{w}] = \mathbf{w}^T \mathbb{E}[\boldsymbol{\eta}_i \boldsymbol{\eta}_i^T] \mathbf{w} = \mathbf{w}^T (\Sigma^2 I) \mathbf{w} = \Sigma^2 \|\mathbf{w}\|^2$

結合: $\tilde{L}_{ss}(\mathbf{w}, b) = \frac{1}{2N} \sum_{i=1}^N ((\mathbf{w}^T \mathbf{x}_i + b - y_i)^2) + \frac{1}{2N} \sum_{i=1}^N \Sigma^2 \|\mathbf{w}\|^2$

$\Sigma^2 \|\mathbf{w}\|^2$ 和 b 無關: $\tilde{L}_{ss}(\mathbf{w}, b) = \frac{1}{2N} \sum_{i=1}^N ((\mathbf{w}^T \mathbf{x}_i + b - y_i)^2) + \frac{\Sigma^2}{2} \|\mathbf{w}\|^2$

最小化 input noise 期望平方根 = 最小化 noise 的平方根 + 加上 L^2 -regularization $\frac{\Sigma^2}{2} \|\mathbf{w}\|^2$

$$\text{得 } \tilde{L}_{ss}(\mathbf{w}, b) = \frac{1}{2N} \sum_{i=1}^N (f_{\mathbf{w}, b}(\mathbf{x}_i) - y_i)^2 + \frac{\Sigma^2}{2} \|\mathbf{w}\|^2$$

$\frac{\Sigma^2}{2} \|\mathbf{w}\|^2$ 可減少 noise 的影響, 對 model 更穩定

Problem 5 (Gradient descent for Logistic Regression with Vectorized Feature) (0.6 pts)

This problem is related to the appendix of W2_Logistic_Regression.pdf. Consider the following optimization problem

$$\min_{\mathbf{w}} \ell(\mathbf{w}), \quad (3)$$

where

$$\ell(\mathbf{w}) = \frac{1}{d} \sum_{n=1}^d \ell^{(n)}(\mathbf{w}), \quad \ell^{(n)}(\mathbf{w}) = \ln(1 + \exp(-y_n (\mathbf{w}^\top \mathbf{x}_n))).$$

Assume that there are d training data, \mathbf{x}_n is the n -th training data, and the label $y_n = \pm 1$.

- (a) (0.2 pts) Prove that $\frac{1}{\ln 2} \ell^{(n)}(\mathbf{w})$ is an upper bound of $\mathbb{1}\{\text{sign}(\mathbf{w}^\top \mathbf{x}_n) \neq y_n\}$ for any \mathbf{w} , where $\mathbb{1}\{\cdot\}$ is the indicator function. Do not use graph calculator for the arguments.

$$(a) \ell^{(n)}(\mathbf{w}) = \ln(1 + \exp(-y_n (\mathbf{w}^\top \mathbf{x}_n)))$$

\Rightarrow 越 $y_n (\mathbf{w}^\top \mathbf{x}_n)$ 越小 $\ell^{(n)}(\mathbf{w})$ 越大

由 indicator function 知 $\mathbb{1}\{\text{sign}(\mathbf{w}^\top \mathbf{x}_n) \neq y_n\}$ 值是 1 表示錯誤分類, 0 為正確分類

$\Rightarrow y_n (\mathbf{w}^\top \mathbf{x}_n) > 0$, 分類正確, $y_n (\mathbf{w}^\top \mathbf{x}_n) < 0$, 分類錯誤

$$\text{pf. } \frac{1}{\ln 2} \ell^{(n)}(\mathbf{w}) \geq \mathbb{1}\{\text{sign}(\mathbf{w}^\top \mathbf{x}_n) \neq y_n\}$$

$$y_n (\mathbf{w}^\top \mathbf{x}_n) \geq 0 \Rightarrow \ell^{(n)}(\mathbf{w}) = \ln(1 + \exp(-y_n (\mathbf{w}^\top \mathbf{x}_n))) \leq \ln(1+1) = \ln 2$$

$$\frac{1}{\ln 2} \ell^{(n)}(\mathbf{w}) \leq 1, \text{ 分類正確時 loss} \leq 1$$

$$y_n (\mathbf{w}^\top \mathbf{x}_n) < 0 \Rightarrow \ell^{(n)}(\mathbf{w}) = \ln(1 + \exp(-y_n (\mathbf{w}^\top \mathbf{x}_n))) \geq \ln 2$$

$y_n (\mathbf{w}^\top \mathbf{x}_n)$ 越小, loss 越接近 $\ln(1 + \exp(0)) = \ln 2 \Rightarrow \frac{1}{\ln 2} \ell^{(n)}(\mathbf{w}) \geq 1$
分類錯誤, loss ≥ 1

$$\frac{1}{\ln 2} \ell^{(n)}(\mathbf{w}) \geq \mathbb{1}\{\text{sign}(\mathbf{w}^\top \mathbf{x}_n) \neq y_n\}$$

由上知錯誤分類 indicator function 的上界

(b) (0.2 pts) For a given (\mathbf{x}_n, y_n) , derive its gradient $\nabla \ell^{(n)}(\mathbf{w})$.

$$(b) \ell^{(n)}(\mathbf{w}) = \ln(1 + \exp(-y_n(\mathbf{w}^T \mathbf{x}_n)))$$

$$g(\mathbf{w}) = -y_n(\mathbf{w}^T \mathbf{x}_n)$$

$$\Rightarrow \ell^{(n)}(\mathbf{w}) = \ln(1 + \exp(g(\mathbf{w})))$$

$$\nabla \ell^{(n)}(\mathbf{w}) = \frac{d}{d\mathbf{w}} \ln(1 + \exp(g(\mathbf{w})))$$

$$\frac{d}{d\mathbf{w}} \ln(1 + \exp(g(\mathbf{w}))) = \frac{1}{1 + \exp(g(\mathbf{w}))} \cdot \frac{d}{d\mathbf{w}} \exp(g(\mathbf{w}))$$

$$(\because g(\mathbf{w}) = -y_n(\mathbf{w}^T \mathbf{x}_n))$$

$$\frac{d}{d\mathbf{w}} \exp(g(\mathbf{w})) = \exp(g(\mathbf{w})) \cdot \frac{d}{d\mathbf{w}} g(\mathbf{w}) \quad \times \quad \frac{d}{d\mathbf{w}} g(\mathbf{w}) = -y_n \mathbf{x}_n$$

$$\frac{d}{d\mathbf{w}} \exp(g(\mathbf{w})) = \exp(g(\mathbf{w})) \cdot (-y_n \mathbf{x}_n)$$

$$\text{合併 } \nabla \ell^{(n)}(\mathbf{w}) = \frac{1}{1 + \exp(g(\mathbf{w}))} \cdot \exp(g(\mathbf{w})) \cdot (-y_n \mathbf{x}_n)$$

$$\Rightarrow \nabla \ell^{(n)}(\mathbf{w}) = \frac{\exp(g(\mathbf{w}))}{1 + \exp(g(\mathbf{w}))} \cdot (-y_n \mathbf{x}_n)$$

$$\frac{\exp(g(\mathbf{w}))}{1 + \exp(g(\mathbf{w}))} = \frac{1}{1 + \exp(-g(\mathbf{w}))} \Rightarrow \nabla \ell^{(n)}(\mathbf{w}) = \frac{1}{1 + \exp(y_n(\mathbf{w}^T \mathbf{x}_n))} \cdot (-y_n \mathbf{x}_n)$$

(c) (0.2 pts) Prove that the optimization problem 3 is equivalent to minimizing the following objective function

$$\mathcal{L}(\mathbf{w}) = - \sum_{n=1}^d \left(\frac{1+y_n}{2} \ln \frac{1+\tanh(\frac{1}{2}\mathbf{w}^\top \mathbf{x}_n)}{2} + \frac{1-y_n}{2} \ln \frac{1-\tanh(\frac{1}{2}\mathbf{w}^\top \mathbf{x}_n)}{2} \right).$$

$$\tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}, \quad 1 + \tanh(z) = \frac{2\exp(z)}{\exp(z) + \exp(-z)}, \quad 1 - \tanh(z) = \frac{2\exp(-z)}{\exp(z) + \exp(-z)}$$

$$\frac{1+y_n}{2} \ln \frac{1 + \tanh(\frac{1}{2}\mathbf{w}^\top \mathbf{x}_n)}{2}$$

當 $y_n=1$, 分類正確: $\ln(1 + \exp(-\mathbf{w}^\top \mathbf{x}_n))$

$$\frac{1-y_n}{2} \ln \frac{1 - \tanh(\frac{1}{2}\mathbf{w}^\top \mathbf{x}_n)}{2}$$

當 $y_n=-1$, 分類錯誤: $\ln(1 + \exp(\mathbf{w}^\top \mathbf{x}_n))$

得 $\mathcal{L}(\mathbf{w}) = \ell(\mathbf{w})$