

Problem 1 (Trace Optimization) (1%)

1. Let $\Sigma \in R^{m \times m}$ be a symmetric positive semi-definite matrix, $\mu \in R^m$. Please construct a set of points $x_1, \dots, x_n \in R^m$ such that

$$\frac{1}{n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T = \Sigma, \quad \frac{1}{n} \sum_{i=1}^n x_i = \mu$$

Hint: n is not given by the problem. WLOG, you may assume $\mu = 0 \in \mathbb{R}^d$.

2. Let $1 \leq k \leq m$, solve the following optimization problem and justify with proof:

$$\begin{array}{ll} \text{minimize} & \text{Trace}(\Phi^T \Sigma \Phi) \\ \text{subject to} & \Phi^T \Phi = I_k \\ \text{variables} & \Phi \in R^{m \times k} \end{array}$$

In other words, you need to find Φ and verify that your Φ minimize the trace.

1. 當 $\mu=0$ 時, $x_1, x_2, \dots, x_n \in R^m$, 使得

$$\frac{1}{n} \sum_{i=1}^n x_i = 0, \quad \frac{1}{n} \sum_{i=1}^n x_i x_i^T = \Sigma$$

$$\Sigma = U \Lambda U^T \quad \left\{ \begin{array}{l} U \in R^{m \times m} \quad (U^T U = I_m) \\ \Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m) \text{ 是包含非負特徵值的對角矩阵} \end{array} \right.$$

希望產生一組 x_1, \dots, x_n 使 sample 共變異數為 Σ

① 為變量常態分佈 $N(\mu, \Sigma)$, 從中取無限多組獨立樣本平均必接近 μ 且 Σ

② 獨立高斯分佈 $z_i \sim N(0, I_m)$ 產生 n 組樣本 z_1, z_2, \dots, z_n

$x_i = U \sqrt{\lambda} z_i$, 因 z_i 的共變異數為 I_m , 則 x_i 的共變異數為 $U \Lambda U^T = \Sigma$

令 $\mu=0$, $x_i = U \sqrt{\lambda} z_i$ 的樣本平均將隨著 n 越來越大而接近 0 .

當樣本數 n 足夠大, 樣本均值和樣本共變異數會近似期望值

2. 令 $E = U^T \Sigma$, 由於 U 是正交矩陣, $E^T E = I_k$ 等價於 $E^T E = I_k$

$$\text{Trace}(E^T \Sigma E) = \text{Trace}(E^T \Lambda E) \Rightarrow \text{Trace}(E^T \Lambda E) = \sum_{j=1}^k \lambda_j^2 = \sum_{j=1}^k \lambda_j (\lambda_j)$$

λ 是權重, 目標是將 E 的列向量 ϕ_i 對應到 λ 方向

選擇 E 的列向量應對屬於 Σ 最小的大個特徵值 $\lambda_1, \lambda_2, \dots, \lambda_k$ 時特徵向量

$$\text{將 } E = [u_1, u_2, \dots, u_k] \Rightarrow \text{Trace}(E^T \Sigma E) = \sum_{i=1}^k \lambda_i \text{ 為最小值}$$

Problem 2 (Gradient Boosting)(1%) 二元分類

Consider the binary classification problem, where we are given training data set $\{(x_i, y_i)\}_{i=1}^N$ with $x_i \in \mathbb{R}^d$ and $y_i \in \{1, -1\}$. Let $F = \{f \mid f : \mathbb{R}^d \rightarrow \{1, -1\}\}$ be the collection of classifiers. Given number of epochs $T \in \mathbb{N}$. Suppose that we want to find the function

$$g(x) = \sum_{t=1}^N \alpha_t f_t(x)$$

加權 跟分類器

where $f_t \in F$ and $\alpha_t \in \mathbb{R}$ for all $t = 1, \dots, T$, by which the aggregated classifier is given by

$$h(x) = \begin{cases} 1, & \text{if } g(x) > 0 \\ -1, & \text{if } g(x) \leq 0. \end{cases}$$

Please apply gradient boosting to show how the functions f_t and the coefficients α_t are computed with an aim to minimize the following loss function

$$L(g) = \sum_{i=1}^N \log \left(1 + e^{-y_i g(x_i)} \right).$$

$$g_0(x) = \operatorname{argmin}_{\ell=1}^N \ell \lg (1 + e^{y_i x})$$

Iterative Steps:

For $t = 1, 2, \dots, T$:

compute residuals: (計算損失函數對模型預測的真捕獲，作為模型的改進方向)

$$r_i^{(t)} = -\frac{\partial L(g)}{\partial g(x_i)} = \frac{y_i}{1 + e^{y_i g_t(x_i)}}$$

Fit a weak classifier $f_t(x)$:

$$\text{最小平方損失: } \sum_{i=1}^n (r_i^{(t)} - f_t(x_i))^2$$

$$f_t(x) = \operatorname{argmin}_{f \in F} \sum_{i=1}^n (r_i^{(t)} - f(x_i))^2$$

weight α_t :

舊的模型: \quad 新的學習基底 $f_t(x)$

$$\alpha_t = \operatorname{argmin}_{\alpha \in \mathbb{R}} \sum_{i=1}^n \log (1 + e^{-y_i (f_t(x_i) + \alpha f_t(x_i))})$$

★ 找一個最佳的 α_t ，使得在當前模型 $g^{(t-1)}(x)$ 的基礎上，加上新的學習基底

$f_t(x)$ 的修正後，損失函數 $L(g)$ 最小

Update the model

$$g_t(x) = g_{t-1}(x) + \alpha_t f_t(x)$$

迭代 T 次，或 another convergence criterion is satisfied.

$$g_T(x) = \sum_{t=1}^T \alpha_t f_t(x)$$

Final Aggregated Classifier:

$$h(x) = \begin{cases} 1 & \text{if } g_T(x) > 0, \\ -1 & \text{if } g_T(x) \leq 0. \end{cases}$$

Problem 3 (EM algorithm for mixture of exponential model) (1%)

Given N samples $x_1, \dots, x_N \in [0, \infty)$, we would like to cluster them into K clusters. Assume the samples are generated according to Exponential mixture models

$$X \sim \sum_{j=1}^K \pi_j \text{Exp}(\tau_j)$$

where $\pi_1 + \dots + \pi_K = 1$, and $\text{Exp}(\tau)$ denotes the exponential distribution with probability density function

$$f_\tau(x) = \begin{cases} (1/\tau)e^{-x/\tau}, & x \geq 0 \\ 0, & x < 0. \end{cases}$$

We would like to apply Expectation Maximization algorithm to find the maximum likelihood estimation of parameters $\theta = \{(\pi_k, \tau_k)\}_{k=1}^K$.

- (a) Please write down the E-step and M-step and show that the parameters are updated from $\theta^{(t)} = \{(\pi_k^{(t)}, \tau_k^{(t)})\}_{k=1}^K$ to $\theta^{(t+1)} = \{(\pi_k^{(t+1)}, \tau_k^{(t+1)})\}_{k=1}^K$ in the following form:

$$\tau_k^{(t+1)} = \frac{\sum_{i=1}^N \delta_{ik}^{(t)} x_i}{\sum_{i=1}^N \delta_{ik}^{(t)}}, \quad \pi_k^{(t+1)} = \frac{1}{N} \sum_{i=1}^N \delta_{ik}^{(t)}$$

- (b) What is the closed form expression of $\delta_{ik}^{(t)}$?

(a) E-step

$$\delta_{ik}^{(t)} = P(Z_i = k | X_i; \theta^{(t)})$$

$$\text{Bayes' rule: } \rightarrow \delta_{ik}^{(t)} = \frac{P(Z_i = k, X_i | \theta^{(t)})}{P(X_i | \theta^{(t)})}$$

$$\because Z_i = k \text{ 表示 } X_i \text{ 由第 } k \text{ 个成分生成} \Rightarrow P(Z_i = k, X_i | \theta^{(t)}) = \pi_k^{(t)} f_{T_k}^{(t)}(X_i)$$

$$P(X_i | \theta^{(t)}) = \sum_{j=1}^K \pi_j^{(t)} f_{T_j}^{(t)}(X_i)$$

$$f_{T_i}(x) = \frac{1}{T_i} e^{-x/T_i} \text{ 带入:}$$

$$\delta_{ik}^{(t)} = \frac{\pi_k^{(t)} \frac{1}{T_k} e^{-x_i/T_k}}{\sum_{j=1}^K \pi_j^{(t)} \frac{1}{T_j} e^{-x_i/T_j}}$$

$$Q(\theta | \theta^{(t)}) = \mathbb{E}_{Z|X, \theta^{(t)}} [\log p(X, Z | \theta)] = \sum_{k=1}^K \left[\sum_{i=1}^N \delta_{ik}^{(t)} \log \pi_k^{(t)} - \sum_{i=1}^N \delta_{ik}^{(t)} \log T_k - \sum_{i=1}^N \delta_{ik}^{(t)} x_i \right]$$

M-step: 最大化 $Q(\theta | \theta^{(t)})$

$$\theta^{(t+1)} = \arg \max_{\theta} Q(\theta | \theta^{(t)})$$

对 T_k 和 π_k 求偏导数 并设为 0, 用时考虑 $\sum_{k=1}^K \pi_k = 1$

1. updating T_k :
令 $L(\pi_k) = \sum_{i=1}^N \delta_{ik}^{(t)} \cdot \ln \pi_k$. 在 $Q(\theta | \theta^{(t)})$ 中和几个有用的式子

$$\sum_{i=1}^N \delta_{ik}^{(t)} \log \pi_k = N_k^{(t)} \log \pi_k$$

最大化 $L(\pi_k)$ 使 $L(\pi_k) = 1$:

$$L(\pi_1, \dots, \pi_K, \lambda) = \sum_{i=1}^N N_k^{(t)} \log \pi_k + \lambda (1 - \sum_{k=1}^K \pi_k)$$

$$\frac{\partial L}{\partial \pi_k} = \frac{N_k^{(t)}}{\pi_k} - \lambda = 0 \Rightarrow \pi_k = \frac{N_k^{(t)}}{\lambda}$$

$$\sum_{k=1}^K \frac{N_k^{(t)}}{\lambda} = 1 \Rightarrow \lambda = \sum_{k=1}^K N_k^{(t)} = N$$

$$\therefore \sum_{k=1}^K \frac{N_k^{(t)}}{\lambda} = \sum_{i=1}^N \sum_{k=1}^K \delta_{ik}^{(t)} = N$$

$$\therefore \pi_K^{(t+1)} = \frac{N_k^{(t)}}{N} = \frac{1}{N} \sum_{i=1}^N \delta_{ik}^{(t)}$$

2. updating T_k :

含有 T_k 的部分 $\bar{T}_k = \frac{N}{N} \sum_{i=1}^N \delta_{ik}^{(t)} \log T_k - \sum_{i=1}^N \delta_{ik}^{(t)} \frac{x_i}{T_k} = Q_k(T_k)$

$$\frac{dQ_k}{dT_k} = -\frac{\sum_{i=1}^N \delta_{ik}^{(t)}}{T_k} + \sum_{i=1}^N \delta_{ik}^{(t)} \frac{x_i}{T_k^2} = 0 \Rightarrow \frac{1}{T_k} \left[-\bar{T}_k \sum_{i=1}^N \delta_{ik}^{(t)} + \sum_{i=1}^N \delta_{ik}^{(t)} x_i \right] = 0$$

$$T_k^{(t+1)} = \frac{\sum_{i=1}^N \delta_{ik}^{(t)} x_i}{\sum_{i=1}^N \delta_{ik}^{(t)}}$$

(b) from E-step: $\delta_{ik}^{(t)} = \frac{\pi_k^{(t)} f_{T_k}^{(t)}(X_i)}{\sum_{j=1}^K \pi_j^{(t)} f_{T_j}^{(t)}(X_i)}$, $f_{T_k}(x) = \frac{1}{T_k} e^{-x/T_k}$, $x=0$ 带入

$$\delta_{ik}^{(t)} = \frac{\pi_k^{(t)} \frac{1}{T_k} e^{-x_i/T_k}}{\sum_{j=1}^K \pi_j^{(t)} \frac{1}{T_j} e^{-x_i/T_j}}$$

Problem 4 (Sparse SVM)(2%)

Given training data of N input-output pairs $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$, where $x_i \in \mathcal{X}$ and $y_i \in \{\pm 1\}$. One can give two types of arguments in favor of the SVM algorithm: one based on the sparsity of the support vectors, another based on the notion of margin. Suppose instead of maximizing the margin, we choose instead to maximize sparsity by minimizing the p -norm of the vector $\alpha = (\alpha_1, \dots, \alpha_N)$ that defines the weight vector \mathbf{w} , for some $p \geq 1$. In this question we consider the case $p = 2$, which leads to the following optimization problem:

$$\begin{aligned} & \text{minimize} \quad f(\alpha, b, \xi) = \frac{1}{2} \sum_{i=1}^N \alpha_i^2 + \sum_{i=1}^N C_i \xi_i \\ & \text{subject to} \quad y_i \left(\sum_{j=1}^N \alpha_j y_j \mathbf{x}_i \cdot \mathbf{x}_j + b \right) \geq 1 - \xi_i, \quad i \in \{1, \dots, N\} \\ & \text{variables} \quad b \in \mathbb{R}, \alpha_i \geq 0, \xi_i \geq 0, \quad i \in \{1, \dots, N\} \end{aligned}$$

which can be rewritten in the following primal problem:

$$\begin{aligned} & \text{minimize} \quad f(\alpha, b, \xi) = \frac{1}{2} \sum_{i=1}^N \alpha_i^2 + \sum_{i=1}^N C_i \xi_i \\ & \text{subject to} \quad \begin{cases} g_{1,i}(\alpha, b, \xi) = 1 - \xi_i - y_i \left(\sum_{j=1}^N \alpha_j y_j \mathbf{x}_i \cdot \mathbf{x}_j + b \right) \leq 0 \\ g_{2,i}(\alpha, b, \xi) = -\alpha_i \leq 0 \\ g_{3,i}(\alpha, b, \xi) = -\xi_i \leq 0 \end{cases} \quad i \in \{1, \dots, N\} \\ & \text{variables} \quad \alpha = (\alpha_1, \dots, \alpha_N) \in \mathbb{R}^N, b \in \mathbb{R}, \xi = (\xi_1, \dots, \xi_N) \in \mathbb{R}^N \end{aligned} \quad (1)$$

as well as its Lagrangian dual problem:

$$\begin{aligned} & \text{maximize} \quad \theta(\omega, \beta, \gamma) = \inf_{\alpha \in \mathbb{R}^N, b \in \mathbb{R}, \xi \in \mathbb{R}^N} L(\alpha, b, \xi, \omega, \beta, \gamma) \\ & \text{subject to} \quad \omega_i \geq 0, \beta_i \geq 0, \gamma_i \geq 0, \quad i \in \{1, \dots, N\} \\ & \text{variables} \quad \omega = (\omega_1, \dots, \omega_N) \in \mathbb{R}^N, \beta = (\beta_1, \dots, \beta_N) \in \mathbb{R}^N, \gamma = (\gamma_1, \dots, \gamma_N) \in \mathbb{R}^N \end{aligned} \quad (2)$$

1. Write down the Lagrangian function $L(\alpha, b, \xi, \omega, \beta, \gamma)$ in explicit form of $\alpha, b, \xi, \omega, \beta, \gamma$.
2. Show that the duality gap between (1) and (2) is zero.
3. Derive $\theta(\omega, \beta, \gamma)$ in explicit form of dual variables ω, β, γ .

$$1. L(\alpha, b, \xi, \omega, \beta, \gamma) = f(\alpha, b, \xi) + \sum_{i=1}^N \omega_i g_{1,i}(\alpha, b, \xi) + \sum_{i=1}^N \beta_i g_{2,i}(\alpha, b, \xi) + \sum_{i=1}^N \gamma_i g_{3,i}(\alpha, b, \xi)$$

$$L(\alpha, b, \xi, \omega, \beta, \gamma) = \frac{1}{2} \sum_{i=1}^N \alpha_i^2 + \sum_{i=1}^N C_i \xi_i + \sum_{i=1}^N \omega_i [1 - \xi_i - y_i \left(\sum_{j=1}^N \alpha_j y_j \mathbf{x}_i \cdot \mathbf{x}_j + b \right)] - \sum_{i=1}^N \beta_i \alpha_i - \sum_{i=1}^N \gamma_i \xi_i;$$

$$2. f(\alpha, b, \xi) = \frac{1}{2} \sum_{i=1}^N \alpha_i^2 + \sum_{i=1}^N C_i \xi_i$$

是 α 的二次可微分凸式加上 ξ 的线性凸式：所有限制式

$$1 - \xi_i - y_i \left(\sum_j \alpha_j y_j \mathbf{x}_i \cdot \mathbf{x}_j + b \right) \leq 0, \quad -\alpha_i \leq 0, \quad -\xi_i \leq 0 \Rightarrow \text{线性不等式, convex optimization problem}$$

Slater:

對線性約束的凸問題，只要存在某組 (α, b, ξ) 使所有不等式嚴格成立，即保證 strong duality

例： $\alpha_i = 0, b = 0, \xi_i = k \Rightarrow 1 - k - y_i(0) = 1 - k < 0, \alpha_i = 0 \Rightarrow -\alpha_i = 0, \xi_i = k > 0 \Rightarrow -\xi_i < 0$

$k > 1$ 即可嚴格滿足全部原始限制式

∴ 滿足 slater, 此凸問題的對偶問題必然不存在濁溝 \Rightarrow Duality Gap = 0

$$3. \theta(\omega, \beta, \gamma) = \inf_{\alpha, b, \xi} L(\alpha, b, \xi, \omega, \beta, \gamma)$$

① 對 b 做極小化

$$\text{Lagrangian } b \text{ 部分: } - \sum_{i=1}^N \omega_i y_i b \quad \frac{\partial L}{\partial b} = 0 \quad - \sum_{i=1}^N \omega_i y_i = 0 \Rightarrow \sum_{i=1}^N \omega_i y_i = 0$$

② 對 ξ 做極小化

$$\text{Lagrangian } \xi \text{ 部分: } \sum_{i=1}^N C_i \xi_i - \sum_{i=1}^N \omega_i \xi_i - \sum_{i=1}^N \beta_i \xi_i \quad \frac{\partial L}{\partial \xi_i} = 0 \quad C_i - \omega_i - \beta_i = 0 \Rightarrow C_i = \omega_i + \beta_i$$

③ 對 α 做極小化

$$\text{Lagrangian } \alpha \text{ 部分: } \frac{1}{2} \sum_{i=1}^N \alpha_i^2 - \sum_{i=1}^N \beta_i \alpha_i - \sum_{i=1}^N \omega_i y_i \left(\sum_j \alpha_j y_j \mathbf{x}_i \cdot \mathbf{x}_j \right) \quad \frac{\partial L}{\partial \alpha_i} = 0 \quad \alpha_i - \sum_{j=1}^N \omega_j y_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \beta_i = 0 \Rightarrow \alpha_i = \sum_{j=1}^N \omega_j y_j (\mathbf{x}_i \cdot \mathbf{x}_j) + \beta_i$$

④ 得 $\theta(\omega, \beta, \gamma)$ 式解

$$\theta(\omega, \beta, \gamma) = \sum_{i=1}^N \omega_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \omega_i \omega_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) \quad \text{且同時滿足: } \sum_{i=1}^N \omega_i y_i = 0, \quad 0 \leq \omega_i \leq C_i, \quad \beta_i \geq 0, \quad \gamma_i = C_i - \omega_i \geq 0$$

$$\Rightarrow \theta(\omega, \beta, \gamma) = \sum_{i=1}^N \omega_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \omega_i \omega_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) \quad \text{subject to: } \sum_{i=1}^N \omega_i y_i = 0, \quad 0 \leq \omega_i \leq C_i$$

4. Show that the dual problem can be simplified as

$$\begin{array}{ll} \text{maximize}_{\omega} & \sum_{i=1}^N \omega_i - \frac{1}{2} \sum_{i=1}^N \left(\sum_{j=1}^N \omega_j y_j y_i \mathbf{x}_j \cdot \mathbf{x}_i \right)_+^2 \\ \text{subject to} & \sum_{i=1}^N \omega_i y_i = 0 \\ \text{variables} & 0 \leq \omega_i \leq C_i, \quad i = 1, \dots, N \end{array} \quad (3)$$

5. Suppose $(\bar{\alpha}, \bar{b}, \bar{\xi})$ and $(\bar{\omega}, \bar{\beta}, \bar{\gamma})$ are the optimal solutions to problems (1) and (2) respectively. Denote $\bar{\mathbf{w}} = \sum_{j=1}^N \bar{\alpha}_j y_j \mathbf{x}_j$.

(a) Prove that

$$\bar{\alpha}_i = \max \left(\sum_{j=1}^N \bar{\omega}_j y_j y_i \mathbf{x}_j \cdot \mathbf{x}_i, 0 \right) \quad \forall i = 1, \dots, N \quad (4)$$

(b) Prove that

$$\bar{b} = \arg \min_{b \in \mathbb{R}} \sum_{i=1}^N C_i \max (1 - y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + b), 0), \quad (5)$$

(c) Prove that $\bar{\xi}_i = \max (1 - y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + \bar{b}), 0)$ for all $i = 1, \dots, N$.

(d) Prove that

$$\left. \begin{array}{ll} \bar{\omega}_i = C_i, & \text{if } y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + \bar{b}) < 1 \\ \bar{\omega}_i = 0, & \text{if } y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + \bar{b}) > 1 \\ 0 \leq \bar{\omega}_i \leq C_i, & \text{if } y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + \bar{b}) = 1 \end{array} \right\} \quad \forall i = 1, \dots, N$$

$$4. L(\alpha, b, \beta, \omega, \beta, \gamma) = \frac{1}{2} \sum_{i=1}^N \alpha_i + \frac{1}{2} \sum_{i=1}^N C_i \xi_i + \sum_{i=1}^N \omega_i \left(1 - \xi_i - y_i \left(\sum_{j=1}^N \alpha_j y_j \mathbf{x}_j \cdot \mathbf{x}_i + b \right) \right) - \sum_{i=1}^N \beta_i \alpha_i - \sum_{i=1}^N \beta_i \xi_i$$

Dual Function $\Theta(\omega) = \inf_{\alpha, b, \beta} L(\alpha, b, \beta, \omega, \beta, \gamma) \rightarrow \text{minimize } L(\alpha, b, \beta, \omega, \beta, \gamma) \text{ with respect to } \alpha, b, \text{ and } \beta$

$$\textcircled{1} \text{ 和 } \beta \text{ 相關: } \sum_{i=1}^N \alpha_i \xi_i - \sum_{i=1}^N (\omega_i + \beta_i) \xi_i \stackrel{\text{極小化}}{\rightarrow} \sum_{i=1}^N \alpha_i = C_i - \omega_i - \beta_i \Rightarrow \dots \omega_i + \beta_i \leq C_i \rightarrow \xi_i = 0 \Rightarrow \beta_i = C_i - \omega_i \text{ 且 } 0 \leq \omega_i \leq C_i$$

$$\textcircled{2} \text{ 和 } \beta \text{ 相關: } -\frac{\partial}{\partial \beta} \sum_{i=1}^N \beta_i \alpha_i \stackrel{\text{極小化}}{\rightarrow} \sum_{i=1}^N \beta_i = \sum_{i=1}^N \omega_i \text{ 約束條件 } \sum_{i=1}^N \beta_i = 0$$

$$\textcircled{3} \text{ 和 } \gamma \text{ 相關: } \frac{1}{2} \sum_{i=1}^N \alpha_i^2 - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j \stackrel{\text{極小化}}{\rightarrow} \alpha_i = \sum_{j=1}^N \omega_j y_j \mathbf{x}_i \cdot \mathbf{x}_j \Rightarrow \alpha_i = \sum_{j=1}^N \omega_j y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

$$\Theta(\omega) = \sum_{i=1}^N \omega_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \omega_i \omega_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j \rightarrow \max_{\omega} \sum_{i=1}^N \omega_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \omega_i \omega_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j \text{ 約束條件: } \left[\begin{array}{l} \sum_{i=1}^N \omega_i = 0 \\ 0 \leq \omega_i \leq C_i, \quad \forall i \end{array} \right]$$

5 (a) KKT:

$$\begin{aligned} \alpha_i &= \sum_{j=1}^N \omega_j y_j \mathbf{x}_j \cdot \mathbf{x}_i + \beta_i \quad \text{where } \beta_i \geq 0 \text{ ensures } \alpha_i \geq 0 \\ \Rightarrow \alpha_i &= \max \left(\sum_{j=1}^N \omega_j y_j \mathbf{x}_j \cdot \mathbf{x}_i, 0 \right) \end{aligned}$$

(b) 找 b 使得 $\sum_{i=1}^N C_i \max (1 - y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + b), 0)$ 最小化

\bar{b} 控制决策邊界 $\bar{\mathbf{w}} \cdot \mathbf{x}_i + b = 0$. 最佳 b 是誤失指向量 (類別) $-y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + b) \Rightarrow$ 正負位於分界線上 的偏置

$$\bar{b} = \arg \min_{b \in \mathbb{R}} \sum_{i=1}^N C_i \max (1 - y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + b), 0)$$

(c) 原始問題中約束 $1 - \xi_i - y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + b) \leq 0 \Rightarrow \xi_i \geq 1 - y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + b)$ 同時 $\xi_i \geq 0$

$$\bar{\xi}_i = \max (1 - y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + b), 0)$$

(d) KKT:

$\bar{\omega}_i = C_i \cdot \text{if sample } i \text{ 位於圓盤內部 } (y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + b) < 1) \rightarrow \text{需滿足 } \bar{\xi}_i > 0 \Rightarrow \bar{\omega}_i = C_i$

$\bar{\omega}_i = 0 \cdot \text{if sample } i \text{ 位於分界線外部 } (y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + b) > 1) \text{ 則 } \bar{\xi}_i = 0 \text{ 且 } \bar{\omega}_i = 0$

$0 \leq \bar{\omega}_i \leq C_i \cdot \text{if sample } i \text{ 位於分界線上 } (y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + b) = 1) \text{ 則 } \bar{\omega}_i \text{ 取值於區間 } [0, C_i]$

Problem 5 (Invited Talk)(1%)

Please share your thoughts about this talk.

本次演講探討隱私保護已成為大數據與AI領域不可忽視的議題。隨著機器學習應用的普及，如何平衡技術進步與隱私風險將是未來研究的重要方向。其風險可能忘從模型推測出訓練數據，判斷特定數據是否參與模型訓練。因此，若能透過加入雜訊來隱藏單一數據樣本的貢獻，或將模型訓練過程分佈在本地端設備上，減少數據集中化的風險，也可更進一步做加密技術，保護數據處理過程中的機密性。在未來，我想可以結合隱私保護技術與模型壓縮技術，在降低隱私風險的同時提升計算效率，這對於隱私保護機器學習的應用有重大突破。

Problem 6 (Bellman Optimality Equations)(2%)

In this problem, we aim to help students review and understand the proof of the Bellman optimality equations.

Theorem. Let Π be the set of nonstationary and randomized policies. Define

$$V^*(s) = \sup_{\pi \in \Pi} V^\pi(s)$$
$$Q^*(s, a) = \sup_{\pi \in \Pi} Q^\pi(s, a).$$

Then there exists a stationary and deterministic policy π such that for all $s \in \mathcal{S}$ and $a \in \mathcal{A}$,

$$V^\pi(s) = V^*(s)$$
$$Q^\pi(s, a) = Q^*(s, a).$$

Remark. The notations are consistent with those in the [lecture notes](#).

1. Verify that V^* and Q^* is bounded between 0 and $\frac{1}{1-\gamma}$. Hence, V^* and Q^* must be finite.
2. Show that given $(s_0, a_0, r_0, s_1) = (s, a, r, s')$, the optimal discounted value γV^* , from $t = 1$ onwards, does not depend on the initial conditions s , a , and r :

$$\sup_{\pi \in \Pi} \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, (s_0, a_0, r_0, s_1) = (s, a, r, s') \right] = \gamma V^*(s').$$

3. Let π^* be a policy such that

$$\forall s \in \mathcal{S}, \quad \pi^*(s) \in \arg \max_{a \in \mathcal{A}} r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V^*(s')].$$

- (a) Explain that π^* is deterministic.
- (b) Now suppose that the transition of states and actions is deterministic. In order to show that π^* is an optimal policy, i.e. $V^*(s) = V^{\pi^*}(s)$, we have to show two inequalities: $V^* \geq V^{\pi^*}$ and $V^* \leq V^{\pi^*} < \infty$. The first one is trivial since $\pi^* \in \Pi$. Now, please show the other inequality $V^* \leq V^{\pi^*} < \infty$.
- (c) Similarly, show that $Q^{\pi^*} = Q^*$ under the assumption that the transition is deterministic.

(a) 依據 Bellman, π^* 是通過對於每個狀態 $s \in S$ 選擇使得 $Q^*(s, a)$ 最大的行動 a 定義的:

$$\pi^*(s) \in \arg \max_{a \in A} r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot|s, a)} [V^*(s')]$$

對於每個狀態 s , π^* 都選擇能使上式最大化的行動 a .

∴ 最大化操作 $\arg \max$ 是針對一個確定的行動值 $a \in A$, 而每個 a 的期望值 $r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot|s, a)} [V^*(s')]$ 是確定的數值。

∴ 每個 S , 最多只有一個 a 能達到該最大值

⇒ π^* 是不隨機的, 每個狀態對應一個唯一的行動選擇

(b) 證明 $V^*(s) \leq V^{\pi^*}(s)$

針對任何狀態 s 和動作 a , 下一個狀態是函數 $s' = f(s, a)$

Bellman optimality equation for V^* :

$$V^*(s) = \max_{a \in A} [r(s, a) + \gamma V^*(f(s, a))] \quad \text{--- ①}$$

∵ $\pi^*(s)$ achieves that maximum: $V^*(s) = r(s, \pi^*(s)) + \gamma V^*(f(s, \pi^*(s)))$

$$V^*(s) = r(s, \pi^*(s)) + \gamma V^{\pi^*}(f(s, \pi^*(s))) \quad \text{--- ②}$$

Comparing ① and ②, the difference is just $V^*(f(s, \pi^*(s))) - V^{\pi^*}(f(s, \pi^*(s)))$

by forward iteration that for the next state $s_1 = f(s, \pi^*(s))$, we also have $V^*(s_1) \leq V^{\pi^*}(s_1)$

Repeating this forward along the entire state sequence generated by π^* yields $V^*(s) \leq V^{\pi^*}(s)$ for all $s \in S$

證明 $V^{\pi^*}(s) < \infty$

Assume rewards are bounded: $0 \leq r(s, a) \leq R_{\max} < \infty$

$$V^*(s) = \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \leq \sum_{t=0}^{\infty} \gamma^t R_{\max} = \frac{R_{\max}}{1-\gamma} < \infty$$

$V^*(s) \leq V^{\pi^*}(s) < \infty$ (Bellman argument + bounded rewards)

∴ $V^{\pi^*}(s) = V^*(s)$, π^* is optimal in the deterministic transition MDP

(c) $Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot|s, a)} [V^*(s')]$

In the deterministic case, if next state is $s' = f(s, a) \Rightarrow Q^*(s, a) = r(s, a) + \gamma V^*(f(s, a))$

For π^* , $Q^{\pi^*}(s, a) = r(s, a) + \gamma V^{\pi^*}(f(s, a))$

∴ $V^*(s') = V^{\pi^*}(s')$ for all states s' , it follows: $Q^{\pi^*}(s, a) = Q^*(s, a)$ for all s, a

∴ $Q^{\pi^*} = Q^*$, completing the argument for action-value optimality as well.