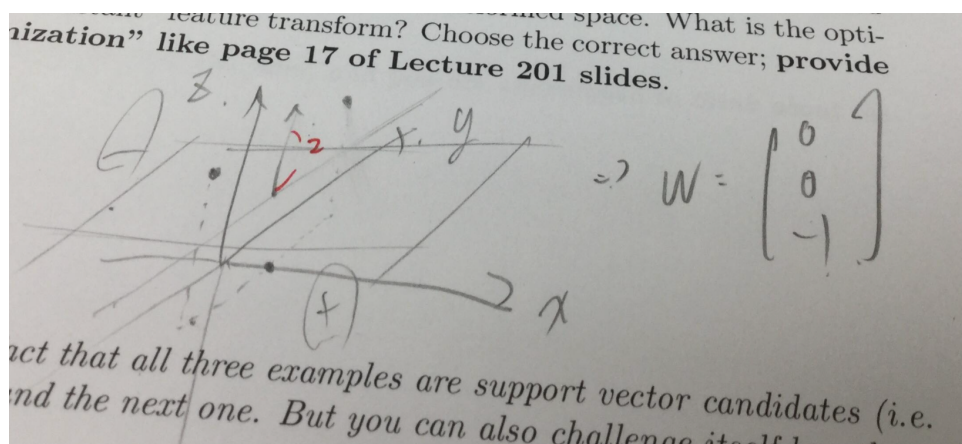


Machine Learning Foundations

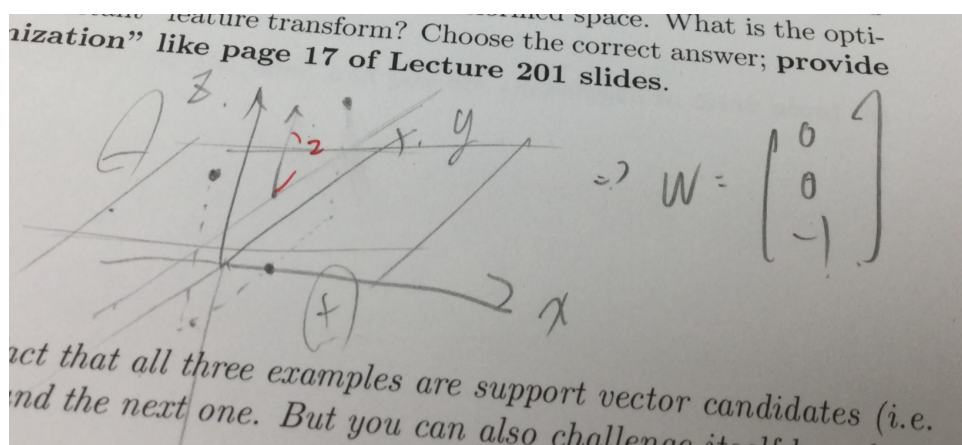
Homework 5

B05505009 謝惟遠

1. [e]



2. [b]



3. [e]

The best separation line will be the center of segment $\overline{x_M x_{M+1}}$

4. [a]

$$\begin{aligned}
P(|x_1 - x_2| \geq 2\rho) &= P(x_1 - x_2 \geq 2\rho) + P(x_2 - x_1 \geq 2\rho) \\
&= 2P(x_2 - x_1 \geq 2\rho) = 2 \int_0^{1-2\rho} \int_{x_1+2\rho}^1 dx_2 dx_1 \\
&= 2 \int_0^{1-2\rho} (1 - 2\rho - x_1) dx_1 = 2 \left((1 - 2\rho)x_1 - \frac{1}{2}x_1^2 \right) \Big|_0^{1-2\rho} \\
&= 2(1 - 2\rho)x_1 - x_1^2 \Big|_0^{1-2\rho} = (1 - 2\rho)^2 \\
2P(|x_1 - x_2| < \rho) + 4P(|x_1 - x_2| \geq \rho) &= 2 - 2(1 - 2\rho)^2 + 4(1 - 2\rho)^2 \\
&= 2 + 2(1 - 2\rho)^2
\end{aligned}$$

5. [c]

Karush-Kuhn-Tucker:

$$\begin{aligned}
\mathcal{L}(b, \mathbf{w}, \boldsymbol{\mu}) &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^n \mu_i (\rho_+ \mathbb{I}[y_i = 1] + \rho_- \mathbb{I}[y_i = -1] - y_i(\mathbf{w}^T \mathbf{x}_i + b)) \\
\nabla_{\mathbf{w}} \mathcal{L}(b, \mathbf{w}, \boldsymbol{\mu}) &= \mathbf{w} - \sum_{i=1}^n \mu_i y_i \mathbf{x}_i \\
\nabla_{\mathbf{w}} \mathcal{L}(b^*, \mathbf{w}^*, \boldsymbol{\mu}) &= \mathbf{w}^* - \sum_{i=1}^n \mu_i y_i \mathbf{x}_i = \mathbf{0} \\
\mathbf{w}^* &= \sum_{i=1}^n \mu_i y_i \mathbf{x}_i \\
\mathcal{L}(b^*, \mathbf{w}^*, \boldsymbol{\mu}^*) &= \frac{1}{2} \mathbf{w}^{*T} \mathbf{w}^* \\
&+ \sum_{i=1}^n \mu_i^* (\rho_+ \mathbb{I}[y_i = 1] + \rho_- \mathbb{I}[y_i = -1] - y_i(\mathbf{w}^{*T} \mathbf{x}_i + b)) \\
&= \max_{\boldsymbol{\mu}} \left(\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \mu_i y_i \mu_j y_j \mathbf{x}_j^T \mathbf{x}_i \right. \\
&+ \left. \sum_{i=1}^n \mu_i (\rho_+ \mathbb{I}[y_i = 1] + \rho_- \mathbb{I}[y_i = -1] - y_i(\mathbf{w}^{*T} \mathbf{x}_i + b)) \right) \\
&= \max_{\boldsymbol{\mu}} \left(\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \mu_i \mu_j y_i y_j \mathbf{x}_j^T \mathbf{x}_i - \sum_{i=1}^n y_i(\mathbf{w}^{*T} \mathbf{x}_i + b) \right. \\
&+ \left. \sum_{i=1}^n \mu_i (\rho_+ \mathbb{I}[y_i = 1] + \rho_- \mathbb{I}[y_i = -1]) \right)
\end{aligned}$$

$$\begin{aligned}
&= \max_{\mu} \left(-\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \mu_i \mu_j y_i y_j \mathbf{x}_j^T \mathbf{x}_i \right. \\
&\quad \left. + \sum_{i=1}^n \mu_i (\rho_+ \mathbb{I}[y_i = 1] + \rho_- \mathbb{I}[y_i = -1]) \right) \\
&= \min_{\mu} \left(\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \mu_i \mu_j y_i y_j \mathbf{x}_j^T \mathbf{x}_i - \sum_{i=1}^n \mu_i (\rho_+ \mathbb{I}[y_i = 1] + \rho_- \mathbb{I}[y_i = -1]) \right)
\end{aligned}$$

6. [e]

$$\begin{aligned}
&\nabla_{\mu} \left(\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \mu_i \mu_j y_i y_j \mathbf{x}_j^T \mathbf{x}_i - \sum_{i=1}^n f(i) \mu_i \right) \\
&= \nabla_{\mu} \left(\frac{1}{2} \mathbf{w}^{*T} \mathbf{w}^* - \sum_{i=1}^n f(i) \mu_i \right) \\
&= \mathbf{w}^{*T} \nabla_{\mu} \left(\sum_{i=1}^n \mu_i y_i \mathbf{x}_i \right) - \nabla_{\mu} \left(\sum_{i=1}^n f(i) \mu_i \right) \\
&= \mathbf{w}^{*T} \begin{bmatrix} | & | & & | \\ y_1 \mathbf{x}_1 & y_2 \mathbf{x}_2 & \dots & y_n \mathbf{x}_n \\ | & | & & | \end{bmatrix} - \begin{bmatrix} f(1) \\ f(2) \\ \vdots \\ f(n) \end{bmatrix}^T = \mathbf{0} \\
&\sum_{i=1}^n \mu_i y_i \mathbf{x}_i^T \begin{bmatrix} | & | & & | \\ y_1 \mathbf{x}_1 & y_2 \mathbf{x}_2 & \dots & y_n \mathbf{x}_n \\ | & | & & | \end{bmatrix} = \begin{bmatrix} f(1) \\ f(2) \\ \vdots \\ f(n) \end{bmatrix}^T \\
&\sum_{i=1}^n \mu_i y_i \mathbf{x}_i^T \begin{bmatrix} | & | & & | \\ y_1 \mathbf{x}_1 & y_2 \mathbf{x}_2 & \dots & y_n \mathbf{x}_n \\ | & | & & | \end{bmatrix} = \begin{bmatrix} f(1) \\ f(2) \\ \vdots \\ f(n) \end{bmatrix}^T
\end{aligned}$$

for original case and new case:

$$\begin{cases} \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^T \begin{bmatrix} | & | & & | \\ y_1 \mathbf{x}_1 & y_2 \mathbf{x}_2 & \dots & y_n \mathbf{x}_n \\ | & | & & | \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}^T \\ \sum_{i=1}^n \mu_i y_i \mathbf{x}_i^T \begin{bmatrix} | & | & & | \\ y_1 \mathbf{x}_1 & y_2 \mathbf{x}_2 & \dots & y_n \mathbf{x}_n \\ | & | & & | \end{bmatrix} = \begin{bmatrix} \rho_+ \mathbb{I}[y_1 = 1] + \rho_- \mathbb{I}[y_1 = -1] \\ \rho_+ \mathbb{I}[y_2 = 1] + \rho_- \mathbb{I}[y_2 = -1] \\ \vdots \\ \rho_+ \mathbb{I}[y_n = 1] + \rho_- \mathbb{I}[y_n = -1] \end{bmatrix}^T \end{cases}$$

$$\begin{cases}
\sum_{i=1}^n (\mu_i - \rho_+ \alpha_i) y_i \mathbf{x}_i^T \begin{bmatrix} | & & | \\ y_1 \mathbf{x}_1 & y_2 \mathbf{x}_2 & \dots & y_n \mathbf{x}_n \\ | & & | \end{bmatrix} = \begin{bmatrix} (\rho_- - \rho_+) \mathbb{I}[y_1 = -1] \\ (\rho_- - \rho_+) \mathbb{I}[y_2 = -1] \\ \vdots \\ (\rho_- - \rho_+) \mathbb{I}[y_n = -1] \end{bmatrix}^T \\
\sum_{i=1}^n (\mu_i - \rho_- \alpha_i) y_i \mathbf{x}_i^T \begin{bmatrix} | & & | \\ y_1 \mathbf{x}_1 & y_2 \mathbf{x}_2 & \dots & y_n \mathbf{x}_n \\ | & & | \end{bmatrix} = \begin{bmatrix} (\rho_+ - \rho_-) \mathbb{I}[y_1 = 1] \\ (\rho_+ - \rho_-) \mathbb{I}[y_2 = 1] \\ \vdots \\ (\rho_+ - \rho_-) \mathbb{I}[y_n = 1] \end{bmatrix}^T
\end{cases}$$

$$\begin{cases}
\sum_{i=1}^n (\mu_i - \rho_+ \alpha_i) y_i \mathbf{x}_i^T \sum_{i=1}^n \alpha_i y_i \mathbf{x}_1 = \begin{bmatrix} (\rho_- - \rho_+) \mathbb{I}[y_1 = -1] \\ (\rho_- - \rho_+) \mathbb{I}[y_2 = -1] \\ \vdots \\ (\rho_- - \rho_+) \mathbb{I}[y_n = -1] \end{bmatrix}^T \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{bmatrix} \\
\sum_{i=1}^n (\mu_i - \rho_- \alpha_i) y_i \mathbf{x}_i^T \sum_{i=1}^n \alpha_i y_i \mathbf{x}_1 = \begin{bmatrix} (\rho_+ - \rho_-) \mathbb{I}[y_1 = 1] \\ (\rho_+ - \rho_-) \mathbb{I}[y_2 = 1] \\ \vdots \\ (\rho_+ - \rho_-) \mathbb{I}[y_n = 1] \end{bmatrix}^T \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{bmatrix}
\end{cases}$$

$$\sum_{i=1}^n (2\mu_i - \rho_+ \alpha_i - \rho_- \alpha_i) y_i \mathbf{x}_i^T \sum_{i=1}^n \alpha_i y_i \mathbf{x}_1 = (\rho_+ - \rho_-) \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}^T \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{bmatrix}$$

given α is optimal:

$$\begin{aligned}
\sum_{j=1}^n \alpha_j y_j &= 0 \\
\sum_{i=1}^n (2\mu_i - \rho_+ \alpha_i - \rho_- \alpha_i) y_i \mathbf{x}_i^T \sum_{i=1}^n \alpha_i y_i \mathbf{x}_1 &= 0 \\
2\mu_i &= (\rho_+ + \rho_-) \alpha_i
\end{aligned}$$

7. [d]

$\log 0$ is not defined

8. [b]

$$\exp(x) \in (0, 1]$$

9. [d]

Consider worst case, where every other point has the opposite label and is at a distance of ϵ :

$$\begin{aligned}
 1 - \sum_{i \neq n} \exp(-\gamma \epsilon^2) &> 0 \\
 1 &> (N-1) \exp(-\gamma \epsilon^2) \\
 \exp(\gamma \epsilon^2) &> (N-1) \\
 \gamma \epsilon^2 &> \log(N-1) \\
 \gamma &> \frac{\log(N-1)}{\epsilon^2}
 \end{aligned}$$

10. [c]

$$\begin{aligned}
 \mathbf{w}_{t+1} &= \mathbf{w}_t + y_{n(t)} \phi(\mathbf{x}_n) \\
 a_{t+1,n} &= y_{n(t)}
 \end{aligned}$$

11. [a]

$$\begin{aligned}
 \mathbf{w}_t &= \sum_{n=1}^N a_{t,n} \phi(\mathbf{x}_n) \\
 \mathbf{w}_t^T \phi(\mathbf{x}) &= \sum_{n=1}^N a_{t,n} \phi(\mathbf{x}_n)^T \phi(\mathbf{x}) = \sum_{n=1}^N a_{t,n} K(\mathbf{x}_n, \mathbf{x})
 \end{aligned}$$

12. [b]

$$b = y_n - y_n \xi_n - w^T z_n$$

If $y_n = 1$,

$$\begin{aligned}
 b &= 1 - \xi_n - w^T z_n \\
 b &\leq 1 - w^T z_n
 \end{aligned}$$

If $y_n = -1$,

$$\begin{aligned}
 b &= -1 + \xi_n - w^T z_n \\
 b &\geq -1 - w^T z_n
 \end{aligned}$$

Upper bound at $y_n = 1$, thus

$$\min_{y_n > 0} (1 - w^T z_n) \geq b$$

13. [e]

$$\begin{aligned}
\mathcal{L}(b, \mathbf{w}, \boldsymbol{\mu}, \boldsymbol{\xi}) &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i^2 + \sum_{i=1}^N \mu_i (1 - \xi_i - y_i (\mathbf{w}^T \mathbf{z}_i + b)) \\
\nabla_{\mathbf{w}} \mathcal{L}(b, \mathbf{w}, \boldsymbol{\mu}, \boldsymbol{\xi}) &= \mathbf{w} - \sum_{n=1}^N y_n \mu_n \nabla_{\mathbf{w}} (\mathbf{w}^T \mathbf{z}_n + b) = \mathbf{w} - \sum_{n=1}^N y_n \mu_n \mathbf{z}_n \\
\mathbf{w}^* &= \sum_{n=1}^N y_n \mu_n \mathbf{z}_n \\
\nabla_{\boldsymbol{\mu}} \mathcal{L}(b, \mathbf{w}, \boldsymbol{\mu}, \boldsymbol{\xi}) &= \begin{bmatrix} 1 - \xi_1 - y_1 (\mathbf{w}^T \mathbf{z}_1 + b) \\ 1 - \xi_2 - y_2 (\mathbf{w}^T \mathbf{z}_2 + b) \\ \vdots \\ 1 - \xi_n - y_n (\mathbf{w}^T \mathbf{z}_n + b) \end{bmatrix} \\
\nabla_{\boldsymbol{\xi}} \mathcal{L}(b, \mathbf{w}, \boldsymbol{\mu}, \boldsymbol{\xi}) &= 2C \boldsymbol{\xi} - \boldsymbol{\mu} \\
\boldsymbol{\xi}^* &= \frac{\boldsymbol{\mu}}{2C} \\
\mathcal{L}(\boldsymbol{\mu}) &= \frac{1}{2} \mathbf{w}^{*T} \mathbf{w}^* + C \sum_{i=1}^N \xi_i^{*2} + \sum_{i=1}^N y_i \mu_i (1 - \xi_i^* - y_i (\mathbf{w}^{*T} \mathbf{z}_i + b)) \\
&= \frac{1}{2} \mathbf{w}^{*T} \mathbf{w}^* + \frac{\boldsymbol{\mu}^T \boldsymbol{\mu}}{4C} + \sum_{i=1}^N \mu_i \left(1 - \frac{\mu_i}{2C} - y_i \left(\sum_{n=1}^N y_n \mu_n \mathbf{z}_n^T \mathbf{z}_i \right) \right) \\
&= \frac{1}{2} \mathbf{w}^{*T} \mathbf{w}^* + \frac{\boldsymbol{\mu}^T \boldsymbol{\mu}}{4C} + \sum_{i=1}^N \mu_i \left(1 - \frac{\mu_i}{2C} - y_i \left(\sum_{n=1}^N y_n \mu_n \mathbf{z}_n^T \mathbf{z}_i \right) \right) \\
&= \frac{1}{2} \mathbf{w}^{*T} \mathbf{w}^* + \frac{\boldsymbol{\mu}^T \boldsymbol{\mu}}{4C} + \sum_{i=1}^N \mu_i - \frac{\boldsymbol{\mu}^T \boldsymbol{\mu}}{2C} - \sum_{i=1}^N \mu_i y_i \left(\sum_{n=1}^N y_n \mu_n \mathbf{z}_n^T \mathbf{z}_i \right) \\
&= \frac{1}{2} \mathbf{w}^{*T} \mathbf{w}^* - \frac{\boldsymbol{\mu}^T \boldsymbol{\mu}}{4C} + \sum_{i=1}^N \mu_i - \sum_{i=1}^N \sum_{n=1}^N \mu_i \mu_n y_i y_n K(\mathbf{x}_n, \mathbf{x}_i) \\
&= \sum_{i=1}^N \mu_i - \frac{1}{2} \sum_{i=1}^N \sum_{n=1}^N \left(\mu_i \mu_n y_i y_n K(\mathbf{x}_n, \mathbf{x}_i) + \frac{\mu_i \mu_n}{2C} \mathbb{I}[i = n] \right)
\end{aligned}$$

14. [e]

$$\begin{aligned}
\mathcal{L}(b, \mathbf{w}, \boldsymbol{\mu}, \boldsymbol{\xi}) &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i^2 + \sum_{i=1}^N \mu_i (1 - \xi_i - y_i (\mathbf{w}^T \mathbf{z}_i + b)) \\
\nabla_{\boldsymbol{\xi}} \mathcal{L}(b, \mathbf{w}, \boldsymbol{\mu}, \boldsymbol{\xi}) &= 2C \boldsymbol{\xi} - \boldsymbol{\mu} \\
\boldsymbol{\xi}^* &= \frac{\boldsymbol{\mu}}{2C}
\end{aligned}$$