

# HW5 Handwritten Assignment Solution

Lecturer: Pei-Yuan Wu

TAs: Yuan-Chia Chang(Problem 1, 2), Chun-Lin Huang(Problem 3, 4, 5)

December 2023, First Edition

## Problem 1 (Kernel)(0.5%)

Consider the following data points:

- $c_1 = \{(3, 3), (3, -3), (-3, 3), (-3, -3)\}$
- $c_2 = \{(6, 6), (6, -6), (-6, 6), (-6, -6)\}$

The data are not linearly separable in this case. Write down a feature map and kernel function to transform the data into a new space, in which the data are linearly separable. Note that you do not just give me a feature map; please explain why.

*Solution:* Feature map  $\Phi : (x, y) \rightarrow (\sqrt{2}xy, x^2, y^2)$ . Kernel function:  $\Phi(\mathbf{x})^T \Phi(\mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2$ . These points can be separated by a plane  $z = 20$ .

## Problem 2 (SVM with Gaussian kernel)(1.5%)

Consider the task of training a support vector machine using the Gaussian kernel  $K(x, z) = \exp(-\frac{\|x-z\|^2}{\tau^2})$ . We will show that as long as there are no two identical points in the training set, we can always find a value for the bandwidth parameter  $\tau$  such that the SVM achieves zero training error.

Recall from class that the decision function learned by the support vector machine can be written as

$$f(x) = \sum_{i=1}^N \alpha_i y_i k(x_i, x) + b$$

Assume that the training data  $\{(x_1, y_1), \dots, (x_N, y_N)\}$  consists of points which are separated by at least distance of  $\epsilon$ ; that is,  $\|x_j - x_i\| \geq \epsilon$ , for any  $i \neq j$ . For simplicity, we assume  $\alpha_i = 1$  for all  $i = 1, \dots, N$  and  $b = 0$ . Find values for the Gaussian kernel width  $\tau$  such that  $x_i$  is correctly classified, for all  $i = 1, \dots, N$ , e.g.,  $f(x_i)y_i > 0$  for all  $i = 1, \dots, N$ .

Hint: Notice that for  $y \in \{-1, +1\}$  the prediction on  $x_i$  will be correct if  $|f(x_i) - y_i| < 1$ , so find a value of  $\tau$  that satisfies this inequality for all  $i$ .

*Solution:* For a training example  $(x_i, y_i)$ ,  $f(x) = \sum_{j=1, j \neq i}^N \alpha_i y_i k(x_i, x_j) = \sum_{j=1, j \neq i}^N \alpha_i y_i k(x_i, x_j) + y_i$ . Hence, if  $|\sum_{j=1, j \neq i}^N \alpha_i y_i k(x_i, x_j)| < |y_i| = 1$ , the training data will not be misclassified. We get

$$\begin{aligned} \left| \sum_{j \neq i} y_j k(x_i, x_j) \right| &= \left| \sum_{j \neq i} y_j \exp(-\|x_j - x_i\|^2 / \tau^2) \right| \\ &\leq \sum_{j \neq i} \left| y_j \exp(-\|x_j - x_i\|^2 / \tau^2) \right| = \sum_{j \neq i} |y_j| \cdot \exp(-\|x_j - x_i\|^2 / \tau^2) \\ &= \sum_{j \neq i} \exp(-\|x_j - x_i\|^2 / \tau^2) \leq \sum_{j \neq i} \exp(-\epsilon^2 / \tau^2) = (N-1) \exp(-\epsilon^2 / \tau^2). \end{aligned}$$

Thus, we need to choose a  $\tau$  such that

$$\tau < \frac{\epsilon}{\sqrt{\log(N-1)}}$$

By choosing, for example,  $\tau = \epsilon / \sqrt{\log N}$ , we are done.

### Problem 3 (Support Vector Regression)(2%)

Suppose we are given a training set  $\{(x_1, y_1), \dots, (x_m, y_m)\}$ , where  $x_i \in \mathbb{R}^{(n+1)}$  and  $y_i \in \mathbb{R}$ . We would like to find a hypothesis of the form  $f(x) = w^T x + b$ . It is possible that no such function  $f(x)$  exists to satisfy these constraints for all points. To deal with otherwise infeasible constraints, we introduce slack variables  $\xi_i$  for each point. The (convex) optimization problem is

$$\min_{w, b, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i \quad (1)$$

$$\text{s.t. } y_i - w^T x_i - b \leq \epsilon + \xi_i \quad i = 1, \dots, m \quad (2)$$

$$w^T x_i + b - y_i \leq \epsilon + \xi_i \quad i = 1, \dots, m \quad (3)$$

$$\xi_i \geq 0 \quad i = 1, \dots, m \quad (4)$$

where  $\epsilon > 0$  is a given, fixed value and  $C > 0$ . Denote that  $\xi = (\xi_1, \dots, \xi_m)$ .

- Write down the Lagrangian for the optimization problem above. Consider the sets of Lagrange multiplier  $\alpha_i, \alpha_i^*, \beta_i$  corresponding to the (2), (3), and (4), so that the Lagrangian would be written as  $\mathcal{L}(w, b, \xi, \alpha, \alpha^*, \beta)$ , where  $\alpha = (\alpha_1, \dots, \alpha_m)$ ,  $\alpha^* = (\alpha_1^*, \dots, \alpha_m^*)$ , and  $\beta = (\beta_1, \dots, \beta_m)$ .
- Derive the dual optimization problem. You will have to take derivatives of the Lagrangian with respect to  $w, b$ , and  $\xi$ .
- Suppose that  $(\bar{w}, \bar{b}, \bar{\xi})$  and  $(\bar{\alpha}, \bar{\alpha}^*, \bar{\beta})$  are the optimal solutions to a primal and dual optimization problem, respectively.

$$\text{Denote } \bar{w} = \sum_{i=1}^m (\bar{\alpha}_i - \bar{\alpha}_i^*) x_i$$

- Prove that

$$\bar{b} = \arg \min_{b \in \mathbb{R}} C \sum_{i=1}^m \max(|y_i - (\bar{w}^T x_i + b)| - \epsilon, 0) \quad (5)$$

- Define  $e = y_i - (\bar{w}^T x_i + \bar{b})$  Prove that

$$\begin{cases} \bar{\alpha}_i = \bar{\alpha}_i^* = 0, & \bar{\xi}_i = 0, & \text{if } |e| < \epsilon \\ 0 \leq \bar{\alpha}_i \leq C, & \bar{\xi}_i = 0, & \text{if } e = \epsilon \\ 0 \leq \bar{\alpha}_i^* \leq C, & \bar{\xi}_i = 0, & \text{if } e = -\epsilon \\ \bar{\alpha}_i = C, & \bar{\xi}_i = e - \epsilon & \text{if } e > \epsilon \\ \bar{\alpha}_i^* = C, & \bar{\xi}_i = -(e + \epsilon) & \text{if } e < -\epsilon \end{cases} \quad (6)$$

- Show that the algorithm can be kernelized and write down the kernel form of the decision function. For this, you have to show that

- The dual optimization objective can be written in terms of inner products or training examples
- At test time, given a new  $x$  the hypothesis  $f(x)$  can also be computed in terms of inner products.

*Solution:*

- Let  $\alpha_i, \alpha_i^*, \beta \geq 0 (i = 1, \dots, m)$  be the Lagrange multiplier for the primal problem. Then the Lagrangian can be written as:

$$\begin{aligned} L(w, b, \xi, \alpha, \alpha^*, \beta, ) &= \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \beta_i \xi_i \\ &\quad - \sum_{i=1}^m \alpha_i (\epsilon + \xi_i - y_i + w^T x_i + b) \\ &\quad - \sum_{i=1}^m \alpha_i^* (\epsilon + \xi_i + y_i - w^T x_i - b) \end{aligned} \quad (7)$$

2. Note that by  $\alpha_i^{(*)}$ , we refer to  $\alpha_i$  and  $\alpha_i^*$ . First, the dual function can be written as:

$$\theta(\alpha, \alpha^*, \beta) = \inf_{w, b, \xi} L(w, b, \xi, \alpha, \alpha^*, \beta) \quad (8)$$

Now, taking the derivatives of Lagrangian w.r.t. all primal variables, we get

$$\frac{\partial}{\partial w} L = w - \sum_{i=1}^m (\alpha_i - \alpha_i^*) x_i = 0 \Rightarrow w = \sum_{i=1}^m (\alpha_i - \alpha_i^*) x_i \quad (9)$$

$$\frac{\partial}{\partial b} L = \sum_{i=1}^m (\alpha_i^* - \alpha_i) = 0 \quad (10)$$

$$\frac{\partial}{\partial \xi} L = C - \alpha_i^{(*)} - \beta_i = 0 \quad (11)$$

Note that

$$\begin{aligned} \theta_D(\alpha, \alpha^*, \beta) &= \frac{1}{2} \|w\|^2 - \epsilon \sum_{i=1}^m (\alpha_i + \alpha_i^*) + \sum_{i=1}^m y_i (\alpha_i - \alpha_i^*) + b \sum_{i=1}^m (\alpha_i^* - \alpha_i) \\ &\quad + \sum_{i=1}^m (\alpha_i^* - \alpha_i) w^T x_i + \sum_{i=1}^m (C - \beta_i - \alpha_i - \alpha_i^*) \xi_i \end{aligned} \quad (12)$$

By the above equation(9)(10) and (11), we get

$$\begin{aligned} \theta_D(\alpha, \alpha^*) &= \frac{1}{2} \|w\|^2 - \epsilon \sum_{i=1}^m (\alpha_i + \alpha_i^*) + \sum_{i=1}^m y_i (\alpha_i - \alpha_i^*) + \sum_{i=1}^m (\alpha_i^* - \alpha_i) w^T x_i \\ &= \frac{1}{2} \left\| \sum_{i=1}^m (\alpha_i - \alpha_i^*) x_i \right\|^2 - \sum_{i=1}^m (\alpha_i - \alpha_i^*) \left( \sum_{j=1}^m (\alpha_j - \alpha_j^*) x_j^T x_i \right) \\ &\quad - \epsilon \sum_{i=1}^m (\alpha_i + \alpha_i^*) + \sum_{i=1}^m y_i (\alpha_i - \alpha_i^*) \\ &= -\frac{1}{2} \sum_{i=1, j=1}^m (\alpha_i - \alpha_i^*) (\alpha_j - \alpha_j^*) x_i^T x_j - \epsilon \sum_{i=1}^m (\alpha_i + \alpha_i^*) + \sum_{i=1}^m y_i (\alpha_i - \alpha_i^*) \end{aligned} \quad (13)$$

Now, the dual problem can be formulated as:

$$\begin{aligned} \max_{\alpha_i, \alpha_i^*} & -\frac{1}{2} \sum_{i=1, j=1}^m (\alpha_i - \alpha_i^*) (\alpha_j - \alpha_j^*) x_i^T x_j - \epsilon \sum_{i=1}^m (\alpha_i + \alpha_i^*) + \sum_{i=1}^m y_i (\alpha_i - \alpha_i^*) \\ \text{s.t.} & \sum_{i=1}^m (\alpha_i^* - \alpha_i) = 0 \\ & 0 \leq \alpha_i, \alpha_i^* \leq C \end{aligned} \quad (14)$$

3. (a) Write the primal problem in the form:

$$\min \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \max(|y_i - (w^T x_i + b)| - \epsilon, 0)$$

Since  $\bar{w}$  is optimal, the optimal bias  $\bar{b}$  is

$$\operatorname{argmin}_b \sum_{i=1}^m \max(|y_i - \bar{w}^T x_i + b| - \epsilon, 0)$$

- (b) Since  $(\bar{w}, \bar{b}, \bar{\xi})$  and  $(\bar{\alpha}, \bar{\alpha}^*, \bar{\beta})$  satisfies the KKT conditions that the following satisfies for all  $i = 1, \dots, N$ :

$$\begin{aligned}
(S_1) \sum_{i=1}^m (\bar{\alpha}_i - \bar{\alpha}_i^*) &= 0 & (P_1) y_i - (\bar{w}^T x_i + \bar{b}) - \epsilon - \bar{\xi}_i &\leq 0 \\
(S_2) C &= \bar{\alpha}_i + \bar{\alpha}_i^* + \bar{\beta}_i & (P_2) (\bar{w}^T x_i + \bar{b}) - y_i - \epsilon - \bar{\xi}_i &\leq 0 \\
(S_3) \bar{w} &= \sum_{i=1}^m (\bar{\alpha}_i - \bar{\alpha}_i^*) x_i & (P_3) -\bar{\xi}_i &\leq 0 \\
(D_1) \bar{\alpha}_i, \bar{\alpha}_i^*, \bar{\beta}_i &\geq 0 & (C_1) \bar{\alpha}_i (y_i - (\bar{w}^T x_i + \bar{b}) - \epsilon - \bar{\xi}_i) &= 0 \\
& & (C_2) \bar{\alpha}_i^* ((\bar{w}^T x_i + \bar{b}) - y_i - \epsilon - \bar{\xi}_i) &= 0 \\
& & (C_3) \bar{\beta}_i (-\bar{\xi}_i) &= 0
\end{aligned}$$

$(C_3)$  is rewritten as  $(C - \bar{\alpha}_i - \bar{\alpha}_i^*)\bar{\xi}_i = 0$  by  $(S_2)$

Define  $e = y_i - (\bar{w}^T x_i + \bar{b})$

- If  $|e| < \epsilon$ , then  $\bar{\alpha}_i = \bar{\alpha}_i^* = 0$  by  $(C_1)(C_2)$ ,  $\bar{\xi}_i = 0$  by  $(C_3)$
- If  $e = \epsilon$ , then  $\bar{\alpha}_i^* = 0$  by  $(C_2)$ ,  $\bar{\xi}_i = 0$ ,  $0 \leq \bar{\alpha}_i \leq C$  by  $(C_1)(C_3)(D_1)$
- If  $e = -\epsilon$ , then  $\bar{\alpha}_i = 0$  by  $(C_1)$ ,  $\bar{\xi}_i = 0$ ,  $0 \leq \bar{\alpha}_i^* \leq C$  by  $(C_2)(C_3)(D_1)$
- If  $e > \epsilon$ , then  $\bar{\alpha}_i^* = 0$  by  $(C_2)$ ,  $\bar{\xi}_i \neq 0$  by  $(P_1)$ ,  $\bar{\alpha}_i = C$  by  $(C_3)$ ,  $\bar{\xi}_i = e - \epsilon$  by  $(C_1)$
- If  $e < -\epsilon$ , then  $\bar{\alpha}_i = 0$  by  $(C_1)$ ,  $\bar{\xi}_i \neq 0$  by  $(P_2)$ ,  $\bar{\alpha}_i^* = C$  by  $(C_3)$ ,  $\bar{\xi}_i = -e - \epsilon$  by  $(C_2)$

In fact,  $\bar{\alpha}_i \bar{\alpha}_i^* = 0$  (its easily to prove by contradiction)

4. By equation (10) in (b), we have  $w = \sum_{i=1}^m (\alpha_i - \alpha_i^*) x_i$ , then

$$f(w, x) = w^T x + b = \sum_{i=1}^m (\alpha_i - \alpha_i^*) x_i^T x + b = \sum_{i=1}^m (\alpha_i - \alpha_i^*) k(x_i, x) + b \quad (15)$$

This shows that the decision function can be written as a kernel form.

## 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:

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

which can be rewritten in the following primal problem:

$$\begin{aligned}
&\text{minimize} && f(\alpha, b, \xi) = \frac{1}{2} \sum_{i=1}^N \alpha_i^2 + \sum_{i=1}^N C_i \xi_i \\
&\text{subject to} && \left. \begin{aligned} g_{1,i}(\alpha, b, \xi) &= 1 - \xi_i - y_i \left( \sum_{j=1}^N \alpha_j y_j \mathbf{x}_i \cdot \mathbf{x}_j + b \right) \leq 0 \\ g_{2,i}(\alpha, b, \xi) &= -\alpha_i \leq 0 \\ g_{3,i}(\alpha, b, \xi) &= -\xi_i \leq 0 \end{aligned} \right\} i \in 1, N \\
&\text{variables} && \alpha = (\alpha_1, \dots, \alpha_N) \in \mathbb{R}^N, b \in \mathbb{R}, \xi = (\xi_1, \dots, \xi_N) \in \mathbb{R}^N
\end{aligned} \quad (16)$$

as well as its Lagrangian dual problem:

$$\begin{aligned}
&\text{maximize} && \theta(\omega, \beta, \gamma) = \inf_{\alpha \in \mathbb{R}^N, b \in \mathbb{R}, \xi \in \mathbb{R}^N} L(\alpha, b, \xi, \omega, \beta, \gamma) \\
&\text{subject to} && \omega_i \geq 0, \beta_i \geq 0, \gamma_i \geq 0, \quad i \in 1, N \\
&\text{variables} && \omega = (\omega_1, \dots, \omega_N) \in \mathbb{R}^N, \beta = (\beta_1, \dots, \beta_N) \in \mathbb{R}^N, \gamma = (\gamma_1, \dots, \gamma_N) \in \mathbb{R}^N
\end{aligned} \quad (17)$$

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 (16) and (17) is zero.

3. Derive  $\theta(\omega, \beta, \gamma)$  in explicit form of dual variables  $\omega, \beta, \gamma$ .
4. Show that the dual problem can be simplified as

$$\begin{aligned} & \text{maximize} && \sum_{i=1}^N \omega_i - \frac{1}{2} \sum_{i=1}^N \left( \sum_{j=1}^N \omega_j y_j y_i \mathbf{x}_j \cdot \mathbf{x}_i \right)_+^2 \\ & \text{subject to} && \sum_{i=1}^N \omega_i y_i = 0 \\ & \text{variables} && 0 \leq \omega_i \leq C_i, \quad i = 1, \dots, N \end{aligned} \quad (18)$$

5. Suppose  $(\bar{\alpha}, \bar{b}, \bar{\xi})$  and  $(\bar{\omega}, \bar{\beta}, \bar{\gamma})$  are the optimal solutions to problems (16) and (17) respectively. Denote  $\bar{\mathbf{w}} = \sum_{j=1}^N \bar{\alpha}_j y_j \mathbf{x}_j$ .

- (a) Prove that

$$\bar{\alpha}_i = \max \left( \sum_{j=1}^N \bar{\omega}_j y_j y_i \mathbf{x}_j \cdot \mathbf{x}_i, 0 \right) \quad \forall i = 1, \dots, N \quad (19)$$

- (b) Prove that

$$\bar{b} = \min_{b \in \mathbb{R}} \sum_{i=1}^N C_i \max(1 - y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + b), 0), \quad (20)$$

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

- (d) Prove that

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

*Solution:*

1. The Lagrangian function can be explicitly written as

$$\begin{aligned} L(\alpha, b, \xi, \omega, \beta, \gamma) &= f(\alpha, b, \xi) + \sum_{i=1}^N \omega_i g_{1,i}(\alpha, b, \xi) + \sum_{i=1}^N \beta_i g_{2,i}(\alpha, b, \xi) + \sum_{i=1}^N \gamma_i g_{3,i}(\alpha, b, \xi) \\ &= \frac{1}{2} \sum_{i=1}^N \alpha_i^2 + \sum_{i=1}^N C_i \xi_i + \sum_{i=1}^N \omega_i \left( 1 - \xi_i - y_i \left( \sum_{j=1}^N \alpha_j y_j \mathbf{x}_i \cdot \mathbf{x}_j + b \right) \right) + \sum_{i=1}^N \beta_i (-\alpha_i) + \sum_{i=1}^N \gamma_i (-\xi_i) \end{aligned}$$

2. We can verify that the problem satisfy the condition of strong duality theorem. Hence, the duality gap is zero.
3. Take partial deriviates of the Lagranian function  $L$  over  $\alpha, b, \xi$  yields

$$\frac{\partial}{\partial \alpha_i} L = \alpha_i - \sum_{j=1}^N \omega_j y_j y_i \mathbf{x}_j \cdot \mathbf{x}_i - \beta_i, \quad \frac{\partial}{\partial b} L = - \sum_{i=1}^N \omega_i y_i, \quad \frac{\partial}{\partial \xi_i} L = C_i - \omega_i - \gamma_i.$$

- If the following conditions hold:

$$\sum_{i=1}^N \omega_i y_i = 0, \quad \omega_i + \gamma_i = C_i \quad \forall i = 1, \dots, N \quad (21)$$

Then  $\theta(\omega, \beta, \gamma) = L(\alpha, b, \xi, \omega, \beta, \gamma)$  iff

$$\alpha_i = \sum_{j=1}^N \omega_j y_j y_i \mathbf{x}_j \cdot \mathbf{x}_i + \beta_i \quad \forall i \in 1, N, \quad (22)$$

at which

$$\theta(\omega, \beta, \gamma) = \sum_{i=1}^N \omega_i - \frac{1}{2} \sum_{i=1}^N \left( \sum_{j=1}^N \omega_j y_j y_i \mathbf{x}_j \cdot \mathbf{x}_i + \beta_i \right)^2. \quad (23)$$

- Otherwise, since  $\frac{\partial}{\partial b}L, \frac{\partial}{\partial \xi_1}L, \dots, \frac{\partial}{\partial \xi_N}L$  are constants not identically zero, one has  $\theta(\omega, \beta, \gamma) = -\infty$ .

4. Because  $\beta_i \geq 0$ ,

$$\left( \sum_{j=1}^N \omega_j y_j y_i \mathbf{x}_j \cdot \mathbf{x}_i + \beta_i \right)^2 = \left( \max \left( \sum_{j=1}^N \omega_j y_j y_i \mathbf{x}_j \cdot \mathbf{x}_i, 0 \right) \right)^2 = \left( \sum_{j=1}^N \omega_j y_j y_i \mathbf{x}_j \cdot \mathbf{x}_i \right)_+^2 \quad (24)$$

By (c), one may rewrite (17) as

$$\begin{aligned} & \text{maximize} && \sum_{i=1}^N \omega_i - \frac{1}{2} \sum_{i=1}^N \left( \sum_{j=1}^N \omega_j y_j y_i \mathbf{x}_j \cdot \mathbf{x}_i \right)_+^2 \\ & \text{subject to} && \sum_{i=1}^N \omega_i y_i = 0, \quad \omega_i + \gamma_i = C_i, \quad \forall i = 1, \dots, N \\ & \text{variables} && \omega_i \geq 0, \beta_i \geq 0, \gamma_i \geq 0, \quad i = 1, \dots, N \end{aligned} \quad (25)$$

which can be further simplified as (18).

5. By (b), the duality gap is zero, so the primal and dual optimal solutions satisfy the KKT conditions. Note that in (c) we have shown that the stationary condition  $\theta(\omega, \beta, \gamma) = L(\alpha, b, \xi, \omega, \beta, \gamma)$  holds iff (21) and (22) are both satisfied. As such, we may write down the KKT conditions, including **Stationary**, **Primal feasible**, **Dual feasible**, and **Complimentary** conditions, as follows (for all  $i = 1, \dots, N$ )

$$\begin{array}{ll} \text{(S1)} & \bar{\omega}_i + \bar{\gamma}_i = C_i \\ \text{(S2)} & \sum_{i=1}^N \bar{\omega}_i y_i = 0 \\ \text{(S3)} & \bar{\alpha}_i = \sum_{j=1}^N \bar{\omega}_j y_j y_i \mathbf{x}_j \cdot \mathbf{x}_i + \bar{\beta}_i \\ \text{(P1)} & y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + \bar{b}) \geq 1 - \bar{\xi}_i \\ \text{(P2)} & \bar{\alpha}_i \geq 0 \\ \text{(P3)} & \bar{\xi}_i \geq 0 \\ \text{(D1)} & \bar{\omega}_i \geq 0 \\ \text{(D2)} & \bar{\beta}_i \geq 0 \\ \text{(D3)} & \bar{\gamma}_i \geq 0 \\ \text{(C1)} & \bar{\omega}_i (1 - \bar{\xi}_i - y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + \bar{b})) = 0 \\ \text{(C2)} & \bar{\beta}_i \bar{\alpha}_i = 0 \\ \text{(C3)} & \bar{\gamma}_i \bar{\xi}_i = 0 \end{array}$$

Given the optimal  $\bar{\omega}_i$ , it is clear by (25) and (S3) that  $\bar{\beta}_i = \left( \sum_{j=1}^N \bar{\omega}_j y_j y_i \mathbf{x}_j \cdot \mathbf{x}_i \right)_-$  and (19) holds.

Given the optimal coefficient vector  $\bar{\alpha}$ , observe that (16) can be rewritten as the following optimization problem

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

hence the optimal bias is given by (20).

- If  $y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + \bar{b}) < 1$ , then  $\bar{\xi}_i > 0$  by (P1),  $\bar{\gamma}_i = 0$  by (C3),  $\bar{\omega}_i = C_i$  by (S1),  $\bar{\xi}_i = 1 - y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + \bar{b})$  by (C1).
- If  $y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + \bar{b}) > 1$ , then  $\bar{\omega}_i = 0$  by (P3, C1), so  $\bar{\gamma}_i = C_i$  by (S1), so  $\bar{\xi}_i = 0$  by (C3).
- If  $y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + \bar{b}) = 1$ , then  $\bar{\xi}_i = 0$  by (S1, C1, C3), and  $0 \leq \bar{\omega}_i \leq C_i$  by (D1, D3, S1).

## Problem 5 (Spherical one class SVM)(2%)(Bonus)

Suppose we aim to fit a hypersphere which encompasses a majority of data points  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^M$  by considering the following optimization problem: (here  $\boldsymbol{\mu}$  and each  $\mathbf{x}_i$  are considered as column vectors)

$$\begin{aligned} & \text{minimize} && R^2 + \frac{1}{\nu} \sum_{i=1}^N C_i \xi_i \\ & \text{subject to} && \left. \begin{array}{l} \|\mathbf{x}_i - \boldsymbol{\mu}\|^2 \leq R^2 + \xi_i \\ \xi_i \geq 0 \end{array} \right\} \quad \forall i \in 1, N \\ & && R \geq 0 \\ & \text{variables} && R \in \mathbb{R}, \boldsymbol{\mu} \in \mathbb{R}^M, \boldsymbol{\xi} = (\xi_1, \dots, \xi_N) \in \mathbb{R}^N \end{aligned} \quad (26)$$

where  $C_i > 0$  for each  $i \in 1, N$ , and  $0 < \nu < \sum_{i=1}^N C_i$ . Let  $\rho = R^2$  and rewrite (26) in the form of primal problem:

$$\begin{aligned} & \text{minimize} && f(\rho, \boldsymbol{\mu}, \boldsymbol{\xi}) = \rho + \frac{1}{\nu} \sum_{i=1}^N C_i \xi_i \\ & \text{subject to} && \left. \begin{aligned} g_{1,i}(\rho, \boldsymbol{\mu}, \boldsymbol{\xi}) &= \|\mathbf{x}_i - \boldsymbol{\mu}\|^2 - \rho - \xi_i \leq 0 \\ g_{2,i}(\rho, \boldsymbol{\mu}, \boldsymbol{\xi}) &= -\xi_i \leq 0 \\ g_3(\rho, \boldsymbol{\mu}, \boldsymbol{\xi}) &= -\rho \leq 0 \end{aligned} \right\} \forall i \in 1, N \\ & \text{variables} && \rho \in \mathbb{R}, \boldsymbol{\mu} \in \mathbb{R}^M, \boldsymbol{\xi} \in \mathbb{R}^N \end{aligned} \quad (27)$$

as well its Lagrangian dual problem:

$$\begin{aligned} & \text{maximize} && \theta(\alpha, \beta, \gamma) = \inf_{\rho \in \mathbb{R}, \boldsymbol{\mu} \in \mathbb{R}^M, \boldsymbol{\xi} \in \mathbb{R}^N} L(\rho, \boldsymbol{\mu}, \boldsymbol{\xi}, \alpha, \beta, \gamma) \\ & \text{subject to} && \alpha_i \geq 0, \beta_i \geq 0 \quad \forall i \in 1, N \\ & && \gamma \geq 0 \\ & \text{variables} && \alpha = (\alpha_1, \dots, \alpha_N) \in \mathbb{R}^N, \beta = (\beta_1, \dots, \beta_N) \in \mathbb{R}^N, \gamma \in \mathbb{R} \end{aligned} \quad (28)$$

1. Write down the Lagrangian function  $L(\rho, \boldsymbol{\mu}, \boldsymbol{\xi}, \alpha, \beta, \gamma)$  in explicit form of  $\rho, \boldsymbol{\mu}, \boldsymbol{\xi}, \alpha, \beta, \gamma$ .
2. Show that the duality gap between (27) and (28) is zero.
3. Derive  $\theta(\alpha, \beta, \gamma)$  in explicit form of dual variables  $\alpha, \beta, \gamma$ .
4. Show that the dual problem can be simplified as

$$\begin{aligned} & \text{maximize} && \|\alpha\|_1 \left( \sum_{i=1}^N \hat{\alpha}_i \|\mathbf{x}_i\|^2 - \sum_{1 \leq i, j \leq N} \hat{\alpha}_i \hat{\alpha}_j \mathbf{x}_i^T \mathbf{x}_j \right) \\ & \text{subject to} && \sum_{i=1}^N \alpha_i \leq 1 \\ & \text{variables} && 0 \leq \alpha_i \leq \frac{C_i}{\nu}, i \in 1, N \end{aligned} \quad (29)$$

where  $\|\alpha\|_1 = \sum_{i=1}^N \alpha_i$  and  $\alpha_i = \|\alpha\|_1 \hat{\alpha}_i$ .

5. Suppose  $(\bar{\rho}, \bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\xi}})$  and  $(\bar{\alpha}, \bar{\beta}, \bar{\gamma})$  are optimal solutions to problems (27) and (28), respectively.

- (a) Show that  $\|\bar{\alpha}\|_1 \bar{\boldsymbol{\mu}} = \sum_{i=1}^N \bar{\alpha}_i \mathbf{x}_i$ .
- (b) Show that

$$\bar{\rho} \in_{\rho \geq 0} \left( \rho + \frac{1}{\nu} \sum_{i=1}^N C_i \max(\|\mathbf{x}_i - \bar{\boldsymbol{\mu}}\|^2 - \rho, 0) \right).$$

- (c) Show that

$$\min \left\{ \rho \geq 0 : \sum_{i: \|\mathbf{x}_i - \bar{\boldsymbol{\mu}}\|^2 > \rho} C_i \leq \nu \right\} \leq \bar{\rho} \leq \min \left\{ \rho \geq 0 : \sum_{i: \|\mathbf{x}_i - \bar{\boldsymbol{\mu}}\|^2 > \rho} C_i < \nu \right\}. \quad (30)$$

- (d) Prove that  $\bar{\xi}_i = \max(\|\mathbf{x}_i - \bar{\boldsymbol{\mu}}\|^2 - \bar{\rho}, 0)$  for each  $i \in 1, N$ .
- (e) Prove that

$$\begin{cases} \bar{\alpha}_i = C_i/\nu & , \text{ if } \|\mathbf{x}_i - \bar{\boldsymbol{\mu}}\|^2 > \bar{\rho} \\ \bar{\alpha}_i = 0 & , \text{ if } \|\mathbf{x}_i - \bar{\boldsymbol{\mu}}\|^2 < \bar{\rho} \\ 0 \leq \bar{\alpha}_i \leq C_i/\nu & , \text{ if } \|\mathbf{x}_i - \bar{\boldsymbol{\mu}}\|^2 = \bar{\rho} \end{cases}.$$

6. Suppose  $C_i = 1/n$  for each  $i \in 1, n$ . What is the physical meaning of  $\nu$ ?

*Solution:*

1. The Lagrangian function can be explicitly written as

$$L(\rho, \boldsymbol{\mu}, \boldsymbol{\xi}, \alpha, \beta, \gamma) = \rho + \frac{1}{\nu} \sum_{i=1}^N C_i \xi_i + \sum_{i=1}^N \alpha_i (\|\mathbf{x}_i - \boldsymbol{\mu}\|^2 - \rho - \xi_i) + \sum_{i=1}^N \beta_i (-\xi_i) + \gamma(-\rho).$$

2. We can verify that the problem satisfy the condition of strong duality theorem. Hence, the duality gap is zero.

3. Take partial derivatives of the Lagrangian function  $L$  over  $\rho, \boldsymbol{\mu}, \boldsymbol{\xi}$ , one has

$$\frac{\partial}{\partial \rho} L = 1 - \sum_{i=1}^N \alpha_i - \gamma, \quad \nabla_{\boldsymbol{\mu}} L = 2 \sum_{i=1}^N \alpha_i (\boldsymbol{\mu} - \mathbf{x}_i), \quad \frac{\partial}{\partial \xi_i} L = \frac{C_i}{\nu} - \alpha_i - \beta_i.$$

• If the following conditions hold:

$$\sum_{i=1}^N \alpha_i + \gamma = 1, \quad \alpha_i + \beta_i = \frac{C_i}{\nu} \quad \forall i \in 1, N, \quad (32)$$

then  $\theta(\alpha, \beta, \gamma) = L(\rho, \boldsymbol{\mu}, \boldsymbol{\xi}, \alpha, \beta, \gamma)$  iff

$$\|\alpha\|_1 \boldsymbol{\mu} = \sum_{i=1}^N \alpha_i \mathbf{x}_i, \quad (33)$$

at which

$$\theta(\alpha, \beta, \gamma) = \|\alpha\|_1 \left( \sum_{i=1}^N \hat{\alpha}_i \|\mathbf{x}_i\|^2 - \sum_{1 \leq i, j \leq N} \hat{\alpha}_i \hat{\alpha}_j \mathbf{x}_i^T \mathbf{x}_j \right) \quad (34)$$

• Otherwise, as  $\frac{\partial}{\partial \rho} L, \frac{\partial}{\partial \xi_1} L, \dots, \frac{\partial}{\partial \xi_N} L$  are constants not identically zero, one has  $\theta(\alpha, \beta, \gamma) = -\infty$ .

4. By (c), one may rewrite (28) as

$$\begin{aligned} & \text{maximize} && \|\alpha\|_1 \left( \sum_{i=1}^N \hat{\alpha}_i \|\mathbf{x}_i\|^2 - \sum_{1 \leq i, j \leq N} \hat{\alpha}_i \hat{\alpha}_j \mathbf{x}_i^T \mathbf{x}_j \right) \\ & \text{subject to} && \left. \begin{aligned} & \alpha_i \geq 0, \beta_i \geq 0 \\ & \alpha_i + \beta_i = \frac{C_i}{\nu} \end{aligned} \right\} \quad \forall i \in 1, N \\ & && \sum_{i=1}^N \alpha_i + \gamma = 1 \\ & && \gamma \geq 0 \\ & \text{variables} && \alpha \in \mathbb{R}^N, \beta \in \mathbb{R}^N, \gamma \in \mathbb{R} \end{aligned}$$

which can be further simplified as (29).

5. By (b), the duality gap is zero, so the primal and dual optimal solutions satisfy the KKT conditions. Note that in (c) we have shown that the stationary condition  $\theta(\alpha, \beta, \gamma) = L(\rho, \boldsymbol{\mu}, \boldsymbol{\xi}, \alpha, \beta, \gamma)$  holds iff (32) and (33) are both satisfied. As such, we may write down the KKT conditions, including **S**tationary, **P**rimal feasible, **D**ual feasible, and **C**omplimentary conditions, as follows (for all  $i \in 1, N$ )

$$\begin{array}{ll} \text{(S1)} & \bar{\alpha}_i + \bar{\beta}_i = C_i/\nu & \text{(D1)} & \bar{\alpha}_i \geq 0 \\ \text{(S2)} & \sum_{i=1}^N \bar{\alpha}_i + \bar{\gamma} = 1 & \text{(D2)} & \bar{\beta}_i \geq 0 \\ \text{(S3)} & \|\bar{\alpha}\|_1 \bar{\boldsymbol{\mu}} = \sum_{i=1}^N \bar{\alpha}_i \mathbf{x}_i & \text{(D3)} & \bar{\gamma} \geq 0 \\ \text{(P1)} & \|\mathbf{x}_i - \bar{\boldsymbol{\mu}}\|^2 \leq \bar{\rho} + \bar{\xi}_i & \text{(C1)} & \bar{\alpha}_i (\|\mathbf{x}_i - \bar{\boldsymbol{\mu}}\|^2 - \bar{\rho} - \bar{\xi}_i) = 0 \\ \text{(P2)} & \bar{\xi}_i \geq 0 & \text{(C2)} & \bar{\beta}_i \bar{\xi}_i = 0 \\ \text{(P3)} & \bar{\rho} \geq 0 & \text{(C3)} & \bar{\gamma} \bar{\rho} = 0 \end{array}$$

Note that (27) can be rewritten as the following optimization problem

$$\begin{aligned} & \text{minimize} && \rho + \frac{1}{\nu} \sum_{i=1}^N C_i \max(\|\mathbf{x}_i - \boldsymbol{\mu}\|^2 - \rho, 0) \\ & \text{subject to} && \rho \geq 0 \\ & \text{variables} && \rho \in \mathbb{R}, \boldsymbol{\mu} \in \mathbb{R}^M \end{aligned}$$

Given the optimal centroid  $\bar{\boldsymbol{\mu}}$ , the optimal squared radius  $\bar{\rho}$  is therefore given by

$$\bar{\rho} \in_{\rho \geq 0} \underbrace{\left( \rho + \frac{1}{\nu} \sum_{i=1}^N C_i \max(\|\mathbf{x}_i - \bar{\boldsymbol{\mu}}\|^2 - \rho, 0) \right)}_{G(\rho)}.$$



Note that  $G$  is a convex and continuous function with right derivative

$$\partial_+ G(\rho) := \lim_{\Delta\rho \downarrow 0} \frac{G(\rho + \Delta\rho) - G(\rho)}{\Delta\rho} = 1 - \frac{F(\rho)}{\nu},$$

where  $F : \rho \in [0, \infty) \mapsto \sum_{i: \|\mathbf{x}_i - \bar{\boldsymbol{\mu}}\|^2 > \rho} C_i$  is a non-increasing function. Let  $\rho_1 = \min \{\rho \geq 0 : F(\rho) \leq \nu\}$  and  $\rho_2 = \min \{\rho \geq 0 : F(\rho) < \nu\}$ .

- If  $\rho < \rho_1$ , then  $F(\rho) > \nu$  and  $\partial_+ G(\rho) < 0$ . Thus  $G$  is strictly decreasing on  $[0, \rho_1]$ .
- If  $\rho \geq \rho_2$ , then  $F(\rho) < \nu$  and  $\partial_+ G(\rho) > 0$ . Thus  $G$  is strictly increasing on  $[\rho_2, \infty)$ .
- If  $\rho \in [\rho_1, \rho_2)$ , then  $F(\rho) = \nu$  and  $\partial_+ G(\rho) = 0$ . Thus  $G$  remains constant on  $[\rho_1, \rho_2]$ .
- If  $\|\mathbf{x}_i - \bar{\boldsymbol{\mu}}\|^2 > \bar{\rho}$ , then  $\bar{\xi}_i > 0$  by (P1),  $\bar{\beta}_i = 0$  by (C2),  $\bar{\alpha}_i = C_i/\nu$  by (S1),  $\bar{\xi}_i = \|\mathbf{x}_i - \bar{\boldsymbol{\mu}}\|^2 - \bar{\rho}$  by (C1).
- If  $\|\mathbf{x}_i - \bar{\boldsymbol{\mu}}\|^2 < \bar{\rho}$ , then  $\bar{\alpha}_i = 0$  by (P2,C1),  $\bar{\beta}_i = C_i/\nu$  by (S1),  $\bar{\xi}_i = 0$  by (C2).
- If  $\|\mathbf{x}_i - \bar{\boldsymbol{\mu}}\|^2 = \bar{\rho}$ , then  $\bar{\xi}_i = 0$  by (S1,C1,C2), and  $0 \leq \bar{\alpha}_i \leq C_i/\nu$  by (D1,D2,S1).

6. If  $C_i = \frac{1}{n}$ , by 5(c),  $\nu$  is the acceptable violation rate of the data points. For example, if  $\nu = 0.1$ ,  $\bar{\rho}$  is chosen to be the value that 10% of the data points are outside of the ball.

## Version Description

1. First Edition: Finish Problem 1 to 5