

# HW4 Handwritten Assignment

Lecturer: Pei-Yuan Wu

TAs: Po-Yang Hsieh, Le-Rong Hsu

November 2024

## Problem 1 (True or false)(0.6%)

You don't need to give any explanation for this problem. Just determine the correctness of the following statement and answer "T" for true and "F" for false. The notation here please refer to the course video.

1. (0.2 %) The strong duality holds only when the primal problem is convex and satisfies the Slater's conditions.
2. (0.2%) The complementary slackness condition,  $u_i g_i(x) = 0, \forall i = 1, \dots, m$ , in a minimization problem, implies that "whenever  $g(\bar{x}) = 0$ , then  $u_i > 0$ ".
3. (0.2%) The dual function  $\theta(\mathbf{u}, \mathbf{v})$  gives a lower bound of the optimal value of the primal problem (as a convex minimization problem in standard form) when  $\theta(\mathbf{u}, \mathbf{v}) > -\infty$  and  $\mathbf{u} \geq 0$ .

## Problem 2 (SVM with Gaussian kernel)(0.9%)

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  $i = 1, \dots, m$  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  $i$ .

## Problem 3 (Support Vector Regression)(1.5%)

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

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

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

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

$$\xi_i \geq 0 \quad i = 1, \dots, m \tag{4}$$

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

- (a) 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  $\beta = (\beta_1, \dots, \beta_m)$ .
- (b) Derive the dual optimization problem. You will have to take derivatives of the Lagrangian with respect to  $w, b$ , and  $\xi$
- (c) 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.

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

- (1) Prove that

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

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

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

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

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

## Problem 4 (Hinge loss with $L^1$ regularization)(1.5%)

Given data points  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^m$  as well as their labels  $y_1, \dots, y_N \in \{\pm 1\}$  and penalty coefficients  $C_1, \dots, C_N > 0$ , where each  $\mathbf{x}_i = [x_{i,1}, \dots, x_{i,m}]^T$  is a column vector, consider the following optimization problem:

$$\begin{aligned} & \text{minimize} && \|\mathbf{w}\|_1 + \sum_{i=1}^N C_i \xi_i \\ & \text{subject to} && y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \\ & && \xi_i \geq 0 \\ & \text{variables} && \mathbf{w} \in \mathbb{R}^m, b \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^N \end{aligned} \quad \left. \vphantom{\begin{aligned} & \text{minimize} \\ & \text{subject to} \\ & \text{variables} \end{aligned}} \right\} i = 1, \dots, N \quad (7)$$

Note that in this formulation, we replace the  $L^2$ -regularization term  $\frac{1}{2} \|\mathbf{w}\|^2$  by the  $L^1$ -regularization term  $\|\mathbf{w}\|_1 = \sum_{j=1}^m |w_j|$ .

- (a) (0.5%) Show that  $(\bar{\mathbf{w}}, \bar{b}, \bar{\boldsymbol{\xi}})$  is an optimal solution of (7) if and only if  $\bar{\mathbf{w}} = \bar{\mathbf{u}} - \bar{\mathbf{v}}$ , where  $(\bar{\mathbf{u}}, \bar{\mathbf{v}}, \bar{b}, \bar{\boldsymbol{\xi}})$  is an optimal solution of the following problem:

$$\begin{aligned} & \text{minimize} && f(\mathbf{u}, \mathbf{v}, b, \boldsymbol{\xi}) = \sum_{j=1}^m (u_j + v_j) + \sum_{i=1}^N C_i \xi_i \\ & \text{subject to} && y_i((\mathbf{u} - \mathbf{v})^T \mathbf{x}_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, N \\ & && \xi_i \geq 0, \quad i = 1, \dots, N \\ & && u_j \geq 0, \quad v_j \geq 0, \quad j = 1, \dots, m \\ & \text{variables} && \mathbf{u} \in \mathbb{R}^m, \mathbf{v} \in \mathbb{R}^m, b \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^N \end{aligned} \quad (8)$$

Following (a), we can now rewrite (8) as the following primal problem:

$$\begin{aligned}
& \text{minimize} && f(\mathbf{u}, \mathbf{v}, b, \boldsymbol{\xi}) = \sum_{j=1}^m (u_j + v_j) + \sum_{i=1}^N C_i \xi_i \\
& \text{subject to} && g_i^1(\mathbf{u}, \mathbf{v}, b, \boldsymbol{\xi}) = 1 - \xi_i - y_i((\mathbf{u} - \mathbf{v})^T \mathbf{x}_i + b) \leq 0, \quad i = 1, \dots, N \\
& && g_i^2(\mathbf{u}, \mathbf{v}, b, \boldsymbol{\xi}) = -\xi_i \leq 0, \quad i = 1, \dots, N \\
& && g_j^3(\mathbf{u}, \mathbf{v}, b, \boldsymbol{\xi}) = -u_j \leq 0, \quad j = 1, \dots, m \\
& && g_j^4(\mathbf{u}, \mathbf{v}, b, \boldsymbol{\xi}) = -v_j \leq 0, \quad j = 1, \dots, m \\
& \text{variables} && \mathbf{u} \in \mathbb{R}^m, \mathbf{v} \in \mathbb{R}^m, b \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^N
\end{aligned} \tag{9}$$

as well as its Lagrangian dual problem:

$$\begin{aligned}
& \text{maximize} && \theta(\alpha, \beta, \boldsymbol{\mu}, \boldsymbol{\nu}) = \inf \{ L(\mathbf{u}, \mathbf{v}, b, \boldsymbol{\xi}, \alpha, \beta, \boldsymbol{\mu}, \boldsymbol{\nu}) : \mathbf{u} \in \mathbb{R}^m, \mathbf{v} \in \mathbb{R}^m, b \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}^N \} \\
& \text{subject to} && \alpha_i \geq 0, \beta_i \geq 0, \quad i = 1, \dots, N \\
& && \mu_j \geq 0, \nu_j \geq 0, \quad j = 1, \dots, m \\
& \text{variables} && \alpha \in \mathbb{R}^N, \beta \in \mathbb{R}^N, \boldsymbol{\mu} \in \mathbb{R}^m, \boldsymbol{\nu} \in \mathbb{R}^m
\end{aligned} \tag{10}$$

where  $L$  denotes the Lagrangian function.

- (b) (0.2%) Associate dual variables  $\alpha_i, \beta_i, \mu_j, \nu_j$  to constraints  $g_i^1, g_i^2, g_j^3, g_j^4$ , respectively. Show that  $L$  can be written in the following explicit form:

$$\begin{aligned}
& L(\mathbf{u}, \mathbf{v}, b, \boldsymbol{\xi}, \alpha, \beta, \boldsymbol{\mu}, \boldsymbol{\nu}) \\
& = \mathbf{1}^T(\mathbf{u} + \mathbf{v}) + \sum_{i=1}^N C_i \xi_i + \sum_{i=1}^N \alpha_i (1 - \xi_i - y_i((\mathbf{u} - \mathbf{v})^T \mathbf{x}_i + b)) + \sum_{i=1}^N \beta_i (-\xi_i) - \boldsymbol{\mu}^T \mathbf{u} - \boldsymbol{\nu}^T \mathbf{v},
\end{aligned} \tag{11}$$

where  $\mathbf{1}$  denotes the all-one vector.

- (c) (0.1%) Show that (9) satisfies Slater's condition.

- (d) (0.2%) Show that:

- (i) (0.1%)  $\theta(\alpha, \beta, \boldsymbol{\mu}, \boldsymbol{\nu}) = -\infty$  unless the following conditions hold:

$$\sum_{i=1}^N \alpha_i y_i = 0, \quad \boldsymbol{\mu} = \mathbf{1} - \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i, \quad \boldsymbol{\nu} = \mathbf{1} + \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i, \tag{12}$$

and

$$\alpha_i + \beta_i = C_i, \quad \forall i = 1, \dots, N, \tag{13}$$

at which case  $\theta(\alpha, \beta, \boldsymbol{\mu}, \boldsymbol{\nu}) = \sum_{i=1}^N \alpha_i$ .

- (ii) (0.1%) The stationary condition holds if and only if (12) and (13) are satisfied.

- (e) (0.2%) Show that the dual problem (10) can be simplified as:

$$\begin{aligned}
& \text{maximize} && \sum_{i=1}^N \alpha_i \\
& \text{subject to} && \sum_{i=1}^N \alpha_i y_i = 0 \\
& && -\mathbf{1} \leq \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \leq \mathbf{1} \\
& \text{variables} && 0 \leq \alpha_i \leq C_i, \quad i = 1, \dots, N
\end{aligned} \tag{14}$$

- (f) (0.2%) Write down the KKT conditions for the primal and dual problems (9)(10).

- (g) (0.1%) Is it true that  $\bar{\mathbf{w}}$  must be a linear combination of  $\mathbf{x}_1, \dots, \mathbf{x}_N$ ? Justify your answer.

## Problem 5 (Spherical one class SVM)(1.5%)

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)

$$\begin{aligned} & \text{minimize} && R^2 + \frac{1}{\nu} \sum_{i=1}^N C_i \xi_i \\ & \text{subject to} && \left. \begin{aligned} & \|\mathbf{x}_i - \boldsymbol{\mu}\|^2 \leq R^2 + \xi_i \\ & \xi_i \geq 0 \end{aligned} \right\} \forall i \in [1, N] \\ & && R \geq 0 \\ & \text{variables} && R \in \mathbb{R}, \boldsymbol{\mu} \in \mathbb{R}^M, \boldsymbol{\xi} = (\xi_1, \dots, \xi_N) \in \mathbb{R}^N \end{aligned} \quad (15)$$

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 (15) in the form of primal problem:

$$\begin{aligned} & \text{minimize} && f(\rho, \boldsymbol{\mu}, \boldsymbol{\xi}) = \rho + \frac{1}{\nu} \sum_{i=1}^N C_i \xi_i \\ & \text{subject to} && \left. \begin{aligned} & g_{1,i}(\rho, \boldsymbol{\mu}, \boldsymbol{\xi}) = \|\mathbf{x}_i - \boldsymbol{\mu}\|^2 - \rho - \xi_i \leq 0 \\ & g_{2,i}(\rho, \boldsymbol{\mu}, \boldsymbol{\xi}) = -\xi_i \leq 0 \\ & g_3(\rho, \boldsymbol{\mu}, \boldsymbol{\xi}) = -\rho \leq 0 \end{aligned} \right\} \forall i \in [1, N] \\ & \text{variables} && \rho \in \mathbb{R}, \boldsymbol{\mu} \in \mathbb{R}^M, \boldsymbol{\xi} \in \mathbb{R}^N \end{aligned} \quad (16)$$

as well its Lagrangian dual problem:

$$\begin{aligned} & \text{maximize} && \theta(\alpha, \beta, \gamma) = \inf_{\rho \in \mathbb{R}, \boldsymbol{\mu} \in \mathbb{R}^M, \boldsymbol{\xi} \in \mathbb{R}^N} L(\rho, \boldsymbol{\mu}, \boldsymbol{\xi}, \alpha, \beta, \gamma) \\ & \text{subject to} && \left. \begin{aligned} & \alpha_i \geq 0, \beta_i \geq 0 \forall i \in [1, N] \\ & \gamma \geq 0 \end{aligned} \right\} \\ & \text{variables} && \alpha = (\alpha_1, \dots, \alpha_N) \in \mathbb{R}^N, \beta = (\beta_1, \dots, \beta_N) \in \mathbb{R}^N, \gamma \in \mathbb{R} \end{aligned} \quad (17)$$

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 (16) and (17) 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

$$\begin{aligned} & \text{maximize} && \|\alpha\|_1 \left( \sum_{i=1}^N \hat{\alpha}_i \|\mathbf{x}_i\|^2 - \sum_{1 \leq i, j \leq N} \hat{\alpha}_i \hat{\alpha}_j \mathbf{x}_i^T \mathbf{x}_j \right) \\ & \text{subject to} && \sum_{i=1}^N \alpha_i \leq 1 \\ & \text{variables} && 0 \leq \alpha_i \leq \frac{C_i}{\nu}, i \in [1, N] \end{aligned} \quad (18)$$

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 (16) and (17), respectively.

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

(b) Show that

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

(c) Show that

$$\min \left\{ \rho \geq 0 : \sum_{i: \|\mathbf{x}_i - \bar{\boldsymbol{\mu}}\|^2 > \rho} C_i \leq \nu \right\} \leq \bar{\rho} \leq \min \left\{ \rho \geq 0 : \sum_{i: \|\mathbf{x}_i - \bar{\boldsymbol{\mu}}\|^2 > \rho} C_i < \nu \right\}. \quad (19)$$

(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  $\nu$ ?