

1. In this problem, we will derive the least square solution for multi-class classification. Consider a general classification problem with  $K$  classes, with a 1-of- $K$  binary encoding scheme (defined latter) for the target vector  $t, t \in \mathbb{R}^K$ . Suppose we are given a training data set  $\{x_n, t_n\}, n = 1, \dots, n$  where  $x_n \in \mathbb{R}^D$ . For the 1-of- $K$  binary encoding scheme,  $t_n$  has the  $k$ -th element being 1 and all other elements being 0 if the  $n$ -th data is in class  $k$ . We can use the following linear model to describe each class:

$$y_k(x) = w_k^T x + w_{k0},$$

where  $k = 1, \dots, K$ . We can conveniently group these together using vector notation so that

$$y(x) = \tilde{W}^T \tilde{x},$$

where  $\tilde{W}$  is a matrix whose  $k$ -th column comprises the  $D + 1$ -dimensional vector  $\tilde{w} = [w_{k0}, w_k^T]^T$  and  $\tilde{x}$  is the corresponding augmented input vector  $[1, x^T]^T$ . For each new input with feature  $x$ , we assign it to the class for which the output  $y_k = \tilde{w}_k^T \tilde{x}$  is largest. Define a matrix  $T$  whose  $n$ -th row is the vector  $t_n^T$  and together a matrix  $\tilde{X}$  whose  $n$ -th row is  $\tilde{x}_n^T$ , the sum-of-squares error function can be written as

$$J(\tilde{W}) = \frac{1}{2} \text{Tr} \left\{ (\tilde{X}\tilde{W} - T)^T (\tilde{X}\tilde{W} - T) \right\}.$$

- (a) Find the closed form solution of  $\tilde{W}$  that minimizes the objective function  $J(\tilde{W})$ . Hint: You may use the following two matrix derivative about trace,  $\frac{\partial}{\partial Z} \text{Tr}(AZ) = A^T$  and  $\frac{\partial}{\partial Z} \text{Tr}(Z^T AZ) = (A^T + A)Z$ .

**Solution:** We first expand the term in the trace operator and get

$$J(\tilde{W}) = \frac{1}{2} \text{Tr} \left\{ \tilde{W}^T \tilde{X}^T \tilde{X} \tilde{W} - 2T^T \tilde{X} \tilde{W} + T^T T \right\}.$$

Setting the derivative with respect to  $\tilde{W}$  to 0, we get

$$0 = \tilde{X}^T \tilde{X} \tilde{W} - \tilde{X}^T T.$$

This shows that

$$\tilde{W} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T T.$$

Note: This is the same solution for multi-dimensional least square if we treat  $t_n$  as the targeted value for  $x_n$ .

- (b) Show that  $J(\tilde{W})$  has a unique minimum. Hint: show that the double derivative of  $J(\tilde{W})$  with respect to  $\tilde{W}$  is positive semi-definite.

**Solution:** We find the double derivative

$$\frac{\partial}{\partial Z} \left[ \tilde{X}^T \tilde{X} \tilde{W} - \tilde{X}^T T \right] = \tilde{X}^T \tilde{X}.$$

$\tilde{X}^T \tilde{X}$  is positive semidefinite because for any vector  $v$ ,  $v^T \tilde{X}^T \tilde{X} v = \|\tilde{X} v\|^2 \geq 0$ .

2. Show that a kernel function  $K(x_1, x_2)$  satisfies the following generalization of the Cauchy-Schwartz inequality:

$$K(x_1, x_2)^2 \leq K(x_1, x_1)K(x_2, x_2).$$

Hint: The Cauchy-Schwartz inequality states that: for two vectors  $u$  and  $v$ ,  $|u^T v|^2 \leq \|u\|^2 \|v\|^2$ .

**Solution 1:** From the definition of kernel, we have

$$\begin{aligned} K(x_1, x_2)^2 &= (\phi(x_1)^T \phi(x_2))^2 \\ &\leq (\phi(x_1)^T \phi(x_1))(\phi(x_2)^T \phi(x_2)) \\ &= K(x_1, x_1)K(x_2, x_2). \end{aligned}$$

The inequality comes from the Cauchy-Schwartz inequality.

**Solution 2:** For an alternative solution, we consider the  $2 \times 2$  Gram matrix

$$\mathbf{K} = \begin{bmatrix} K(x_1, x_1) & K(x_1, x_2) \\ K(x_2, x_1) & K(x_2, x_2) \end{bmatrix}$$

Since  $K(x_1, x_2)$  is a valid kernel,  $\mathbf{K}$  is positive definite with  $|\mathbf{K}| \geq 0$ . This shows that  $K(x_1, x_2)^2 \leq K(x_1, x_1)K(x_2, x_2)$ .

3. Given valid kernels  $K_1(x, x')$  and  $K_2(x, x')$ , show that the following kernels are also valid:

(a)  $K(x, x') = K_1(x, x') + K_2(x, x')$ .

**Solution:** Suppose  $K_1(x, x')$  has positive semi-definite Kernel matrix  $\mathbf{K}_1$  and  $K_2(x, x')$  has positive semi-definite Kernel matrix  $\mathbf{K}_2$  with same dimension. Then it is easy to show that  $K(x, x')$  has Kernel matrix  $\mathbf{K} = \mathbf{K}_1 + \mathbf{K}_2$  which is also positive semi-definite.

(b)  $K(x, x') = K_1(x, x')K_2(x, x')$ .

**Solution:** We assume the mapping function for  $K_1(x, x')$  is  $\phi^{(1)}(x)$  and similarly  $\phi^{(2)}(x)$  for  $K_2(x, x')$ . Moreover, we further assume the dimension of  $\phi^{(1)}(x)$  is  $M$  and the dimension of  $\phi^{(2)}(x)$  is  $N$ . We can then expand  $K(x, x')$ .

$$\begin{aligned}
 K(x, x') &= K_1(x, x')K_2(x, x') \\
 &= \phi^{(1)}(x)^T \phi^{(1)}(x') \phi^{(2)}(x)^T \phi^{(2)}(x') \\
 &= \sum_{i=1}^M \phi_i^{(1)}(x) \phi_i^{(1)}(x') \sum_{j=1}^N \phi_j^{(2)}(x) \phi_j^{(2)}(x') \\
 &= \sum_{i=1}^M \sum_{j=1}^N \left[ \phi_i^{(1)}(x) \phi_j^{(2)}(x) \right] \left[ \phi_i^{(1)}(x') \phi_j^{(2)}(x') \right] \\
 &= \sum_{k=1}^{MN} \phi_k(x) \phi_k(x') = \phi(x)^T \phi(x').
 \end{aligned}$$

In above equation,  $\phi(x)$  is a  $MN \times 1$  column vector with the  $k$ -th element given by  $\phi_i^{(1)}(x) \times \phi_j^{(2)}(x)$ . For given  $k$ , the corresponding  $i$  and  $j$  is calculated as follows:  $i = \lfloor (k-1)/N \rfloor + 1$ , and  $j = (k-1) \bmod N + 1$ .

(c)  $K(x, x') = \exp(K_1(x, x'))$ . Hint: use your results in (a) and (b).

**Solution:** Consider the Taylor series expansion for the exponential function:

$$K(x, x') = \sum_{n=0}^{\infty} \frac{K_1(x, x')^n}{n!}.$$

Using results from (a) and (b) repetitively shows that  $K(x, x')$  is a valid kernel.

4. In class, we learned that the soft margin SVM have the primal problem:

$$\begin{aligned} \min_{\xi, w, b} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq 1 - \xi_i, \quad i = 1, \dots, m \\ & \xi_i \geq 0, \quad i = 1, \dots, m \end{aligned}$$

and the dual problem:

$$\begin{aligned} \max_{\alpha} \quad & W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} a_i a_j \langle x^{(i)}, x^{(j)} \rangle \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, i = 1, \dots, m \\ & \sum_{i=1}^m \alpha_i y^{(i)} = 0. \end{aligned}$$

Now suppose we have solved the dual problem and have the optimal  $\alpha$ . Show that the parameter  $b$  can be determined using the following equation:

$$b = \frac{1}{N_{\mathcal{M}}} \sum_{n \in \mathcal{M}} \left( y^{(n)} - \sum_{m \in \mathcal{S}} \alpha_m y^{(m)} \langle x^{(n)}, x^{(m)} \rangle \right). \quad (1)$$

In (1),  $\mathcal{M}$  denotes the set of indexes of data points having  $0 < \alpha_n < C$  and  $\mathcal{S}$  denotes the set of indexes of data points having  $\alpha_n \neq 0$ .

**Solution:** From the KKT condition (complementary slackness), we find that for each data points with  $0 < \alpha_n < C$ , i.e.,  $n \in \mathcal{M}$ , we have

$$y^{(n)}(w^T x^{(n)} + b) = 1.$$

Multiply  $y^n$  on both sides and then summing over  $\mathcal{M}$ , we have:

$$b = \frac{1}{N_{\mathcal{M}}} \sum_{n \in \mathcal{M}} (y^{(n)} - w^T x^{(n)}).$$

Rewrite  $w$  in terms of  $\alpha$  by using  $w = \sum_{m \in \mathcal{S}} \alpha_m y^{(m)} x^{(m)}$ . We find

$$b = \frac{1}{N_{\mathcal{M}}} \sum_{n \in \mathcal{M}} \left( y^{(n)} - \sum_{m \in \mathcal{S}} \alpha_m y^{(m)} \langle x^{(n)}, x^{(m)} \rangle \right).$$

5. Suppose you are given 6 one-dimensional points: 3 with negative labels  $x_1 = -1, x_2 = 0, x_3 = 1$  and 3 with positive labels  $x_4 = -3, x_5 = -2, x_6 = 3$ . In this question, we first compare the performance of linear classifier with or without kernel. Then we solve for the maximum margin classifier using SVM.

- (a) Consider a linear classifier of form  $f(x) = \text{sign}(w_1x + w_0)$ . Write down the optimal value of  $w$  and its classification accuracy on the above 6 points. There might be more than one optimal solution, writing down one of them is enough.

**Solution:** One optimal solution is  $w = [-1, -3/2]^T$  which gives the accuracy of  $5/6$ .

- (b) Given two samples  $x$  and  $z$  in  $\mathbb{R}$ , define the kernel  $K : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  as

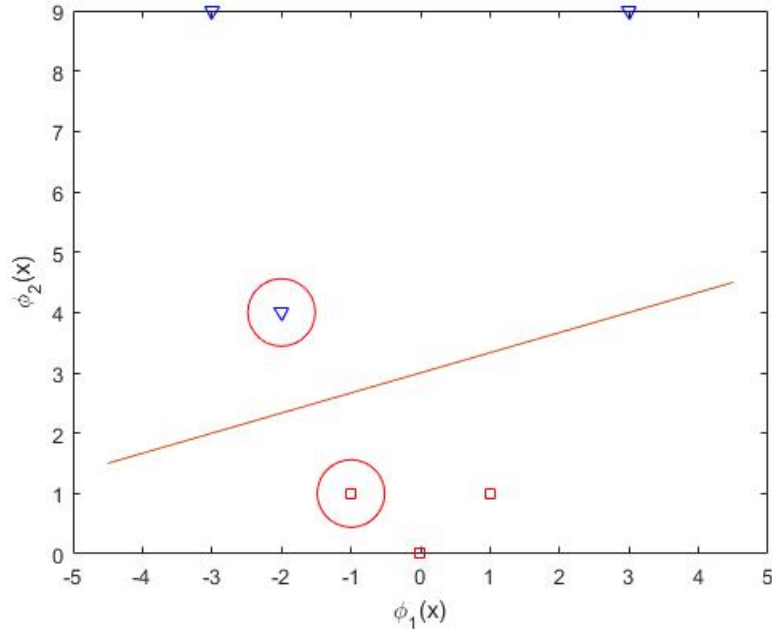
$$K(x, z) = xz(1 + xz).$$

Find the corresponding feature map  $\phi(x)$ .

**Solution:**  $\phi(x) = [x, x^2]^T$ .

- (c) Apply  $\phi(x)$  to the data and plot the points in the induced feature space  $\mathbb{R}^2$ . Are these points linearly separable now?

**Solution:** Yes.



- (d) Look at the points in the induced feature space. Construct a maximum margin separating hyperplane in  $\mathbb{R}^2$  which can be parameterized by  $w_1\phi_1(x) + w_2\phi_2(x) + w_0 = 0$ . Draw this hyperplane on your plot and circle the support vectors.

**Solution:** In the induced feature space, the two points that are closest are  $(-1, 1)$  with negative label and  $(-2, 4)$  with positive label. They are the support vectors and the maximum margin separating hyperplane is given by  $-\phi_1(x) + 3\phi_2(x) - 9 = 0$  by finding a line that has normal vector  $[-1, 3]^T$  and also passes through the mid-point of the support vectors, i.e.,  $(-\frac{3}{2}, \frac{5}{2})$ . The line is drawn on the above figure.

- (e) Draw the decision boundary of the separating hyperplane you found in (d) in the original  $\mathbb{R}$  feature space.

**Solution:** To find the decision boundary, we solve for  $-x + 3x^2 - 9 = 0$  ( $x \approx (-1.6, 1.9)$ ). The decision boundary are shown in the following figure:



- (f) Find the  $\alpha_i$ ,  $w$  and  $b$  in

$$h(x) = \text{sign} \left( \sum_{n \in \mathcal{S}} \alpha_n y_n K(x_n, x) + b \right) = \text{sign} (w^T \phi(x) + b).$$

Do this by solving the dual form of the quadratic program. How is  $w$  and  $b$  related to your solution in part (d)?

**Solution:** Since we only have two support vectors, only the Lagrange multiplier corresponding to the support vectors are non-zero. Let  $\alpha_1$  denote the Lagrange multiplier for  $x_1 = -1$  and similarly  $\alpha_5$  for  $x_5 = -2$ . From the condition  $\sum_{i=1}^6 \alpha_i y_i = 0$ , we get  $\alpha_1 = \alpha_5 = \alpha_0$ . Write down the objective of the dual problem of SVM

$$\begin{aligned} W(\alpha) &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y_i y_j \alpha_i \alpha_j K(x_i, x_j) \\ &= 2\alpha_0 - \frac{1}{2} \alpha_0^2 K(-1, -1) + \alpha_0^2 K(-1, -2) - \frac{1}{2} \alpha_0^2 K(-2, -2) \\ &= 2\alpha_0 - 5\alpha_0^2. \end{aligned}$$

Maximizing  $W(\alpha)$  over  $\alpha_0$ , we get  $\alpha_1 = \alpha_5 = \alpha_0 = \frac{1}{5}$ . Using  $w = \sum_{m \in \mathcal{S}} \alpha_m y_m \phi(x_m)$ , we get  $w = [-\frac{1}{5}, \frac{3}{5}]^T$ . To find  $b$ , recall that

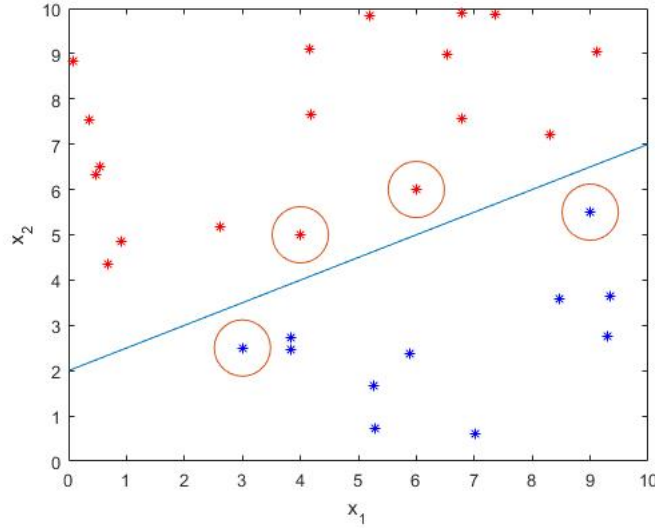
$$y_i (w^T \phi(x_i) + b) = 1$$

for any support vectors  $\phi(x_i)$ . Use any support vector, we can get  $b = -\frac{9}{5}$ . The  $w$  and  $b$  we find by solving the dual problem is a scaled version of  $[w_1, w_2]^T$  and  $w_0$  in part (d). These solutions therefore give the same separating hyperplane.

6. In this exercise, we will use MATLAB to solve both the primal and the dual problem of SVM. In *Data.csv*, the first two columns contain feature vectors  $x^{(i)} \in \mathbb{R}^2$  and the last column contains the label  $y^{(i)} \in \{-1, 1\}$ . We will use CVX as the optimization solver in this problem. For help with CVX, refer to the CVX Users' Guide. Attach your code for submission. For Python user, feel free to use the following libraries: math, csv, numpy, matplotlib and cvxpy.

- (a) **Visualization** Use different color to plot data with different labels in the 2-D feature space. Is the data linearly separable?

**Solution:** Yes, the data is linearly separable.



- (b) **The Primal Problem** Use CVX to solve the primal problem of this form:

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq 1, \quad i = 1, \dots, m \end{aligned}$$

Report  $w$  and  $b$ . Plot the hyperplane defined by  $w$  and  $b$ .

**Solution:** We get  $w = [-0.5, 1]^T$  and  $b = -2$ . The hyperplane is shown in the figure above.

- (c) **The Dual Problem** Use CVX to solve the dual problem of this form:

$$\begin{aligned} \max_{\alpha} \quad & W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} a_i a_j \langle x^{(i)}, x^{(j)} \rangle \\ \text{s.t.} \quad & 0 \leq \alpha_i, i = 1, \dots, m \\ & \sum_{i=1}^m \alpha_i y^{(i)} = 0. \end{aligned}$$

Use the resulting  $a$  to identify the support vectors on the plot. Report your non-zero  $a_i$ 's. How many support vectors do you have? Circle those support vectors.

Note: The latter part of  $W(a)$  is in quadratic form, i.e.,  $a^T P a$ . To use CVX, first find  $P$  and then use `quad_form(a,P)`. For Python user, you will need to add a small number to the diagonal of  $P$  matrix to make cvxpy work. i.e. Run the following code before using cvxpy: “`P += 1e-13 * numpy.eye(29)`”, where 29 is the total number of data. Also, assume it is 0 if a number is less than 1e-9.

**Solution:** There are 4 support vectors.

The corresponding  $\alpha = [0.3347, 0.2903, 0.3165, 0.3085]^T ([0.3289, 0.2961, 0.2993, 0.3257]^T$  (Python)). The support vectors are circled in the above figure. In order to use CVX to solve this problem, first find a matrix  $P$  where  $P_{ij} = y^{(i)}y^{(j)}\langle x^{(i)}, x^{(j)} \rangle$ . Then let CVX to maximize  $\sum_{i=1}^m \alpha_i - \frac{1}{2}\alpha^T P \alpha$ .