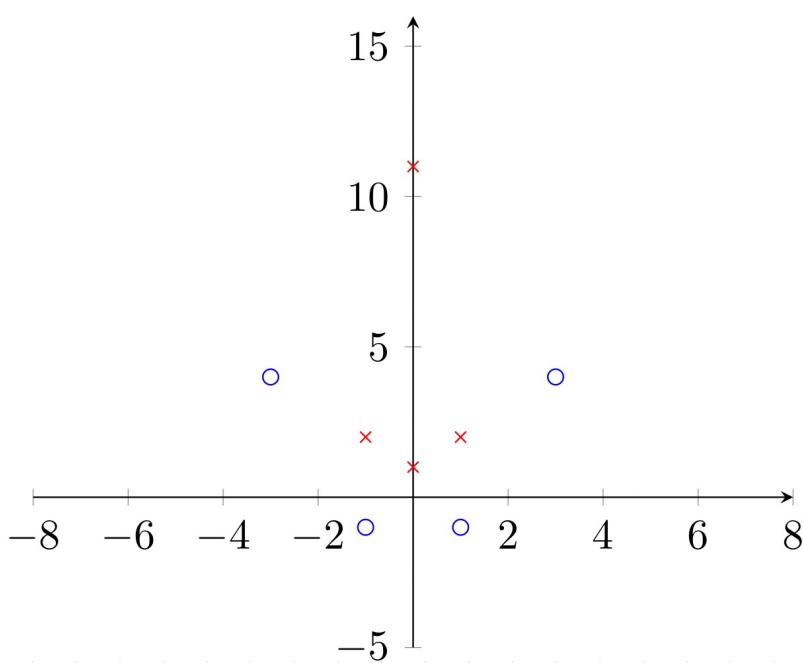


Problem 1. (15 points). SVM's obtain *non-linear* decision boundaries by mapping the feature vectors $\bar{X} \in \mathbb{R}^d$ to a possibly high dimensional space via a function $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^m$, and then finding a linear decision boundary in the new space.

We also saw that to implement SVM, it suffices to know the kernel function $K(\bar{X}_i, \bar{X}_j) = \phi(\bar{X}_i) \cdot \phi(\bar{X}_j)$, without even explicitly specifying the function ϕ .

Recall **Mercer's theorem**. K is a kernel function if and only if for any n vectors, $\bar{X}_1, \dots, \bar{X}_n \in \mathbb{R}^d$, and **any** real numbers c_1, \dots, c_n , $\sum_{i=1}^n \sum_{j=1}^n c_i c_j K(\bar{X}_i, \bar{X}_j) \geq 0$.

1. Prove the following half of Mercer's theorem (which we showed in class). If K is a kernel then $\sum_{i=1}^n \sum_{j=1}^n c_i c_j K(\bar{X}_i, \bar{X}_j) \geq 0$.
2. Let $d = 1$, and $x, y \in \mathbb{R}$. Is the function $K(x, y) = x + y$ a kernel?
3. Let $d = 1$, and $x, y \in \mathbb{R}$. Is $K(x, y) = xy + 1$ a kernel?
4. Suppose $d = 2$, namely the original features are of the form $\bar{X}_i = [\bar{X}^1, \bar{X}^2]$. Show that $K(\bar{X}, \bar{Y}) = (1 + \bar{X} \cdot \bar{Y})^2$ is a kernel function. This is called as **quadratic kernel**.
(Hint: Find a $\phi: \mathbb{R}^2 \rightarrow \mathbb{R}^m$ (for some m) such that $\phi(\bar{X}) \cdot \phi(\bar{Y}) = (1 + \bar{X} \cdot \bar{Y})^2$).
5. Consider the training examples $\langle [0, 1], 1 \rangle, \langle [1, 2], 1 \rangle, \langle [-1, 2], 1 \rangle, \langle [0, 11], 1 \rangle, \langle [3, 4], -1 \rangle, \langle [-3, 4], -1 \rangle, \langle [1, -1], -1 \rangle, \langle [-1, -1], -1 \rangle$. We have plotted the data points below.
 - Is the data **linearly classifiable** in the original 2-d space? If yes, please come up with *any* linear decision boundary that separates the data. If no, please explain why.
 - Is the data linearly classifiable in the feature space corresponding to the quadratic kernel. If yes, please come up with *any* linear decision boundary that separates the data. If no, please explain why.



$$① K(\bar{x}_i, \bar{x}_j) = \phi(\bar{x}_i) \cdot \phi(\bar{x}_j)$$

$$\begin{aligned} \sum_{i=1}^n \sum_{j=1}^n c_i \cdot c_j \cdot K(\bar{x}_i, \bar{x}_j) &= \sum_{i=1}^n \sum_{j=1}^n c_i c_j \cdot \phi(\bar{x}_i) \cdot \phi(\bar{x}_j) \\ &= \left(\sum_{i=1}^n c_i \phi(\bar{x}_i) \right) \left(\sum_{j=1}^n c_j \phi(\bar{x}_j) \right) \\ &\stackrel{?}{=} \left\| \sum_{i=1}^n c_i \phi(\bar{x}_i) \right\|_2^2 \geq 0 \end{aligned}$$

$$\text{So, } \sum_{i=1}^n \sum_{j=1}^n c_i \cdot c_j \cdot K(\bar{x}_i, \bar{x}_j) \geq 0$$

② K is a kernel if $\sum_{i=1}^n \sum_{j=1}^n c_i \cdot c_j \cdot K(\bar{x}_i, \bar{x}_j) \geq 0$

is $K(x, y) = x + y$ a kernel?

Let $n = 1$, $x = \bar{x}_i$ and $y = \bar{x}_j$,

then by Mercer's theorem in ①, K is kernel if:

$$\sum_{i=1}^1 \sum_{j=1}^1 c_i \cdot c_j \cdot K(\bar{x}_i, \bar{x}_j) = c_1 \cdot c_1 \cdot \overbrace{K(\bar{x}_1, \bar{x}_1)}^{x_1 + x_1} \geq 0$$

If $c_1 = -1$ and $\bar{x}_1 = -1$:

$$(-1)(-1)(-1-1) = -2 < 0 \rightarrow \text{does not satisfy Mercer's Theorem}$$

This proof that $K(x, y) = x + y$ is not a kernel.

③ Is $K(x, y) = xy + 1$ a kernel?

If K is kernel, then it can be decomposed into dot product of 2 vectors:

$$\begin{aligned} K(x, y) &= \phi(\bar{x}) \cdot \phi(\bar{y}) \\ &= xy + 1 \\ &= [x \ 1] \cdot \begin{bmatrix} y \\ 1 \end{bmatrix} \\ &= \phi(\bar{x})^\top \cdot \phi(\bar{y}) \end{aligned}$$

$$\phi(\bar{x}) = \begin{bmatrix} x \\ 1 \end{bmatrix} \quad \phi(\bar{y}) = \begin{bmatrix} y \\ 1 \end{bmatrix}$$

Since we have found $\phi(\bar{x})$ and $\phi(\bar{y})$, it satisfies Mercer's theorem proof in ① to further proof that $K(x, y) = xy + 1$ is a kernel.

4. Suppose $d = 2$, namely the original features are of the form $\bar{X}_i = [\bar{X}^1, \bar{X}^2]$. Show that $K(\bar{X}, \bar{Y}) = (1 + \bar{X} \cdot \bar{Y})^2$ is a kernel function. This is called as **quadratic kernel**.

(Hint: Find a $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^m$ (for some m) such that $\phi(\bar{X}) \cdot \phi(\bar{Y}) = (1 + \bar{X} \cdot \bar{Y})^2$).

$$\begin{aligned}
 K(\bar{X}, \bar{Y}) &= (1 + \bar{x} \cdot \bar{y})^2 \quad \rightarrow \bar{X}_i = [x_1, x_2] \text{ & } \bar{Y}_i = [y_1, y_2] \\
 &= (1 + x_1 y_1 + x_2 y_2)^2 \\
 &= (1 + x_1 y_1 + x_2 y_2)(1 + x_1 y_1 + x_2 y_2) \\
 &= 1 + x_1 y_1 + x_2 y_2 + x_1 y_1 + x_1^2 y_1^2 + x_1 x_2 y_1 y_2 \\
 &\quad + x_2 y_2 + x_1 x_2 y_1 y_2 + x_2^2 y_2^2 \\
 &= 1 + 2x_1 y_1 + 2x_2 y_2 + 2x_1 x_2 y_1 y_2 + x_1^2 y_1^2 + x_2^2 y_2^2 \\
 &= \phi(\bar{X})^\top \cdot \phi(\bar{Y}) \quad \leftarrow 1 \times n \cdot n \times 1 \Rightarrow 1 \times 1
 \end{aligned}$$

$$\phi(\bar{X}) = \begin{bmatrix} \sqrt{1} \\ \sqrt{2} x_1 \\ \sqrt{2} x_2 \\ \sqrt{2} x_1 x_2 \\ x_1^2 \\ x_2^2 \end{bmatrix} \quad \phi(\bar{Y}) = \begin{bmatrix} \sqrt{1} \\ \sqrt{2} y_1 \\ \sqrt{2} y_2 \\ \sqrt{2} y_1 y_2 \\ y_1^2 \\ y_2^2 \end{bmatrix}$$

$$\begin{aligned}
 &= (\sqrt{1} \times \sqrt{1}) + (\sqrt{2} x_1 \times \sqrt{2} y_1) + (\sqrt{2} x_2 \times \sqrt{2} y_2) \\
 &\quad + (\sqrt{2} x_1 x_2 \times \sqrt{2} y_1 y_2) + (x_1^2 y_1^2) + (x_2^2 y_2^2) \\
 &= 1 + 2x_1 y_1 + 2x_2 y_2 + 2x_1 x_2 y_1 y_2 + x_1^2 y_1^2 + x_2^2 y_2^2
 \end{aligned}$$

Since we proof to split $K(x, y) = \phi(x) \cdot \phi(y)$,
we can proof Mercer's theorem and
that $K(X, Y)$ is a kernel.

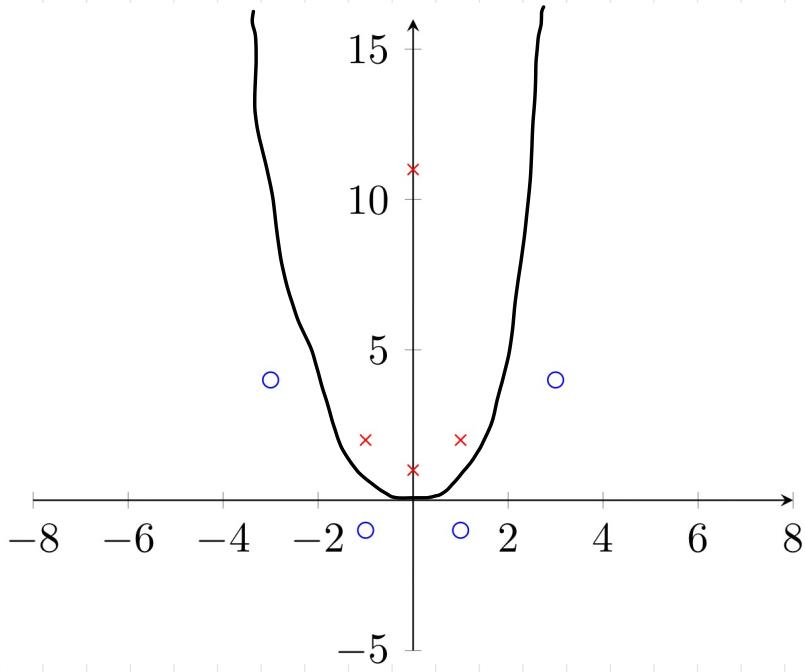
⑤

1. No, it is not linearly classifiable.

There is no single hyperplane w that can correctly classify the data in the original 2-D hyperplane

2. Yes, it is linearly classifiable corresponding to the quadratic kernel.

For example, if the hyperplane looks like this:



This hyperplane has a formula of $y=x^2$ in $x-y$ plane

2 region: $x_2 > (x_1)^2$ or $x_2 < (x_1)^2$

$$(x_1)^2 - x_2 < 0 \quad \text{or} \quad (x_1)^2 - x_2 > 0$$

label = 1

or

X

label = -1

or

O

Problem 2. (10 points). Let $f, h_i, 1 \leq i \leq n$ be real-valued functions and let $\alpha \in \mathbb{R}^n$. Let $L(z, \alpha) = f(z) + \sum_{i=1}^n \alpha_i h_i(z)$. In this problem, we will prove that the following two optimization problems are equivalent.

$$\begin{aligned} & \min_z f(z) \\ & \text{s.t. } h_i(z) \leq 0, i = 1, \dots, n. \end{aligned} \quad (1)$$

$$\min_z \left[\max_{\alpha \geq 0} L(z, \alpha) \right] \quad (2)$$

Let (z^*, α^*) be the solution of (2) and let z_p^* be the solution of (1). Prove that:

$$L(z^*, \alpha^*) = f(z_p^*)$$

Hint: Use the fact that for any $z, \alpha \geq 0$, $L(z^*, \alpha^*) \geq L(z^*, \alpha)$ and $L(z^*, \alpha^*) \leq L(z, \alpha_z)$, where $\alpha_z = \arg \max_{\alpha \geq 0} L(z, \alpha)$.

You may follow the following steps but it is not required as long as your proof is correct.

1. Prove that $L(z^*, \alpha^*) \leq f(z_p^*)$
2. Prove that $L(z^*, \alpha^*) \geq f(z_p^*)$

So if α^* is $\max_{\alpha} L(z, \alpha)$. For any $z, \alpha \geq 0$, $L(z^*, \alpha^*) \geq L(z^*, \alpha)$

Since if we use α^* , it will maximize $L(z^*, \alpha^*)$.

If $z^* \mapsto \min_z \left[\max_{\alpha \geq 0} L(z, \alpha) \right]$. For any $z, \alpha \geq 0$, $L(z^*, \alpha^*) \leq L(z, \alpha_z)$.

Since if we use z^* , it will minimize $L(z, \alpha^*)$

$$L(z^*, \alpha^*) = f(z^*) + \sum_{i=1}^n \alpha_i^* h_i(z^*) \quad \text{s.t. } h_i(z) \leq 0 \quad i = 1, \dots, n$$

$$L(z_p^*, \alpha_{z_p^*}) = f(z_p^*) + \sum_{i=1}^n \alpha_{z_p^* i} h_i(z_p^*)$$

We know that $L(z^*, \alpha^*) \leq L(z, \alpha_z)$, so

$$L(z^*, \alpha^*) \leq L(z_p^*, \alpha_{z_p^*}) \leftarrow \text{because } z_p^* \text{ may not minimize } L(z, \alpha)$$

$$L(z^*, \alpha^*) \leq f(z_p^*) + \sum_{i=1}^n \alpha_{z_p^* i} h_i(z_p^*)$$

$$\text{We also know that } f(z_p^*) + \sum_{i=1}^n \alpha_{z_p^* i} h_i(z_p^*) \leq f(z_p^*)$$

because by definition $\alpha_{z_p^* i} \geq 0$ and $h_i(z_p^*) \leq 0$ for $i = 1, \dots, n$

so $\sum_{i=1}^n \alpha_{z_p^* i} h_i(z_p^*)$ is negative value that will subtract $f(z_p^*)$ term

Thus, we have proven that $L(z^*, \alpha^*) \leq f(z_p^*)$

$$\mathcal{L}(z^*, \alpha^*) = f(z^*) + \sum_{i=1}^n \alpha_i^* h_i(z^*)$$

$$\mathcal{L}(z_p^*, \alpha) = f(z_p^*) + \sum_{i=1}^n \alpha_i h_i(z_p^*)$$

We know that $\mathcal{L}(z^*, \alpha^*) \geq \mathcal{L}(z_p^*, \alpha)$, so

$$\mathcal{L}(z^*, \alpha^*) \geq \mathcal{L}(z_p^*, \alpha) \rightarrow \text{because } \alpha \text{ may not max}$$

$$f(z^*) + \sum_{i=1}^n \alpha_i^* h_i(z^*) \geq f(z_p^*) + \sum_{i=1}^n \alpha_i h_i(z_p^*)$$

$\mathcal{L}(z, \alpha)$

We can also know that z^* satisfies the constraint $h_i(z^*) \leq 0, i=1, \dots, n$

This is because Lagrange multiplier is based on original + linear combination of constraints.

ex: $f(z^*) + \sum_{i=1}^n \alpha_i^* h_i(z^*)$

$\underbrace{\hspace{10em}}$

original function
to max or min
linear combination
of constraint

So, since z^* satisfy the constraint $h_i(z^*) \leq 0$ for $i=1, \dots, n$

we also know that $f(z_p^*) \leq f(z^*)$.

This is because z_p^* will minimize (1), as mentioned in the problem.

So, back to the equation we have: $\rightarrow < 0$, as $\alpha \geq 0$ and $h_i(z) \leq 0$

$$f(z^*) + \sum_{i=1}^n \alpha_i^* h_i(z^*) \geq f(z_p^*) + \sum_{i=1}^n \alpha_i h_i(z_p^*)$$

So, $f(z^*) + \sum_{i=1}^n \alpha_i^* h_i(z^*) \geq f(z_p^*)$, since

$\Rightarrow \mathcal{L}(z^*, \alpha^*) \geq f(z_p^*)$ is proven

$$\begin{aligned} f(z^*) &\geq f(z_p^*) \\ \text{and} \\ \sum_{i=1}^n \alpha_i^* h_i(z^*) &\geq \sum_{i=1}^n \alpha_i h_i(z_p^*) \\ \text{[as } \alpha \text{ will maximize } \mathcal{L}(z, \alpha)] \end{aligned}$$

Since we've proven that $\mathcal{L}(z^*, \alpha^*) \leq f(z_p^*)$

and $\mathcal{L}(z^*, \alpha^*) \geq f(z_p^*)$

This means that $\mathcal{L}(z^*, \alpha^*) = f(z_p^*)$

Problem 3. (15 points). In this problem, we derive the dual formulation of the soft-margin SVM problem with $\bar{\xi} = \xi_1, \dots, \xi_n$.

reference :
Wikipedia on
Lagrange &
KKT.

$$\min_{\bar{w}, \bar{\xi}} \frac{1}{2} \|\bar{w}\|_2^2 + C \cdot \sum_{i=1}^n \xi_i$$

such that

$$1 - \xi_i - y_i(\bar{X}_i \cdot \bar{w} - t) \leq 0, \quad i = 1, \dots, n,$$

$$-\xi_i \leq 0, \quad i = 1, \dots, n.$$

Now we can define $2n$ Lagrangian variables $\alpha = \alpha_1, \dots, \alpha_n$, and $\beta = \beta_1, \dots, \beta_n$ corresponding to these equations and obtain the following Lagrangian.

$$L(\bar{w}, t, \bar{\xi}, \alpha, \beta) = \frac{1}{2} \|\bar{w}\|_2^2 + C \cdot \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - \xi_i - y_i(\bar{X}_i \cdot \bar{w} - t)) - \sum_{i=1}^n \beta_i \xi_i \quad (4)$$

The original problem is now equivalent to

$$\min_{\bar{w}, t, \bar{\xi}} \left[\max_{\alpha \geq 0, \beta \geq 0} L(\bar{w}, t, \bar{\xi}, \alpha, \beta) \right].$$

For this objective, minmax theorem says that the min max problem is equivalent to the max min problem below. You do not have to prove this, but are encouraged to do so (using the argument given in the discussion earlier). This gives the following problem.

$$\max_{\alpha \geq 0, \beta \geq 0} \left[\min_{\bar{w}, t, \bar{\xi}} L(\bar{w}, t, \bar{\xi}, \alpha, \beta) \right].$$

1. For a fixed α, β take the gradient of the the Lagrangian with respect to \bar{w} , and express \bar{w} in terms of the other variables.
2. Differentiate the Lagrangian with respect to t , and equate to zero to obtain another equation.
3. Differentiate the Lagrangian with respect to ξ_i and show that $\alpha_i + \beta_i = C$ at the optimum. This shows that $\alpha_i \leq C$, since $\beta_i \geq 0$.
4. The expressions from 1, 2, 3 define one of the KKT conditions. Show that under these conditions,

$$\min_{\bar{w}, t, \bar{\xi}} L(\bar{w}, t, \bar{\xi}, \alpha, \beta) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (\bar{X}_i \cdot \bar{X}_j).$$

Basically, use the expressions from 1, 2, and 3 to "cancel out" $\bar{w}, t, \bar{\xi}$ in the Lagrangian.

5. Combine the results above to argue that the following optimization is equivalent to the soft margin SVM we started with.

$$\max_{\alpha \geq 0} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (\bar{X}_i \cdot \bar{X}_j) \quad (5)$$

such that

$$0 \leq \alpha_i \leq C, \quad i = 1, \dots, n.$$

(1) $L(\bar{w}, t, \bar{\xi}, \alpha, \beta) = \frac{1}{2} \|\bar{w}\|_2^2 + C \cdot \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - \xi_i - y_i(\bar{X}_i \cdot \bar{w} - t)) - \sum_{i=1}^n \beta_i \xi_i$

$$\begin{aligned} \nabla_{\bar{w}} L(\bar{w}, t, \bar{\xi}, \alpha, \beta) &= \bar{w} + 0 + \frac{d}{d\bar{w}} \sum_{i=1}^n \alpha_i (1 - \xi_i - y_i(\bar{X}_i \cdot \bar{w} - t)) - 0 \\ &= \bar{w} - \sum_{i=1}^n \alpha_i \cdot y_i \cdot \bar{X}_i \Rightarrow 0 \end{aligned}$$

$$\bar{w} = \sum_{i=1}^n \alpha_i \cdot y_i \cdot \bar{X}_i$$

$$\textcircled{2} \quad \nabla_t L(\bar{w}, t, \bar{\xi}, \alpha, \beta) = 0 + 0 + \sum_{i=1}^n \alpha_i y_i - 0 \Rightarrow 0$$

$$\boxed{\sum_{i=1}^n \alpha_i y_i = 0}$$

$$\textcircled{3} \quad \nabla_\xi L(\bar{w}, t, \bar{\xi}, \alpha, \beta) = 0 + \sum_{i=1}^n C - \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \beta_i \Rightarrow 0$$

$$\sum_{i=1}^n C = \sum_{i=1}^n \alpha_i + \sum_{i=1}^n \beta_i \Rightarrow C = \alpha_i + \beta_i$$

$$\textcircled{4} \quad L(\bar{w}, t, \bar{\xi}, \alpha, \beta) = \frac{1}{2} \|\bar{w}\|_2^2 + C \cdot \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - \xi_i - y_i (\bar{X}_i \cdot \bar{w} - t)) - \sum_{i=1}^n \beta_i \xi_i$$

$$= \frac{1}{2} \left\| \sum_{i=1}^n \alpha_i \cdot y_i \cdot \bar{x}_i \right\|_2^2 + (\alpha_i + \beta_i) \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i \xi_i - \sum_{i=1}^n \alpha_i y_i (\bar{x}_i \cdot \bar{w} - t) - \sum_{i=1}^n \beta_i \xi_i$$

$$= \frac{1}{2} \left\| \sum_{i=1}^n \alpha_i \cdot y_i \cdot \bar{x}_i \right\|_2^2 + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i (\bar{x}_i \cdot \bar{w} - t)$$

$$= \frac{1}{2} \left(\sum_{i=1}^n \alpha_i \cdot y_i \cdot \bar{x}_i \right) \left(\sum_{j=1}^n \alpha_j \cdot y_j \cdot \bar{x}_j \right) + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i (\bar{x}_i \cdot \left(\sum_{j=1}^n \alpha_j y_j \bar{x}_j \right) - t)$$

$$= \sum_{i=1}^n \alpha_i + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (\bar{x}_i \cdot \bar{x}_j) - \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (\bar{x}_i \cdot \bar{x}_j) + \sum_{i=1}^n \alpha_i y_i t$$

$$L(\bar{w}, t, \bar{\xi}, \alpha, \beta) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (\bar{x}_i \cdot \bar{x}_j) \quad \underbrace{\sum_{i=1}^n \alpha_i y_i = 0}_{\text{become minimum}}$$

Since from \textcircled{3}, $\alpha_i \leq C$ and $\alpha_i \geq 0$, meet KKT constraint criterion. So, optimized and

$$\min_{\bar{w}, t, \bar{\xi}} L(\bar{w}, t, \bar{\xi}, \alpha, \beta) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (\bar{X}_i \cdot \bar{X}_j).$$

$$\textcircled{5} \quad \text{Since we know that } \min_{\bar{w}, t, \bar{\xi}} L(\bar{w}, t, \bar{\xi}, \alpha, \beta) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (\bar{X}_i \cdot \bar{X}_j).$$

Then the max-min problem is:

$$\max_{\alpha \geq 0, \beta \geq 0} \min L(\bar{w}, t, \bar{\xi}, \alpha, \beta) \Rightarrow \max_{\alpha \geq 0, \beta \geq 0} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (\bar{x}_i \cdot \bar{x}_j)$$

where since $C = \alpha + \beta$ and $\beta \geq 0$, $\alpha_i \leq C$.

So, the equation above follows $0 \leq \alpha_i \leq C$

Equivalent to the soft margin SVM we started with.