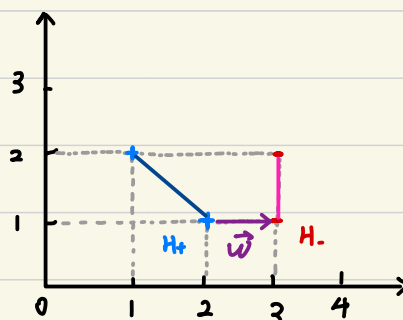


Q1:

$$\text{margin} = \min_{\substack{u \in \Delta, \\ v \in \Delta}} \frac{1}{2} \left\| \sum_{i \in H_+} u_i X_i - \sum_{i \in H_-} v_i X_i \right\|_2^2$$



By inspection, the minimum distance between these two hyperplanes H_+ and H_- is

the distance between $(2, 1)$ and $(3, 1)$

$$\begin{aligned} \text{Let } r \text{ be a scale such that } \vec{w} &= r \left(\begin{bmatrix} 2 \\ 1 \end{bmatrix} - \begin{bmatrix} 3 \\ 1 \end{bmatrix} \right) \\ &= r \begin{bmatrix} -1 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} -r \\ 0 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \text{Since } \vec{w}^T x + b &= y, \quad \text{and} \quad x_1 = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \quad y_1 = 1 \\ x_2 &= \begin{bmatrix} 3 \\ 1 \end{bmatrix}, \quad y_2 = -1, \end{aligned}$$

$$\begin{cases} \vec{w}^T x_1 + b = y_1 \\ \vec{w}^T x_2 + b = y_2 \end{cases}$$

$$\begin{aligned} \Rightarrow \begin{cases} [-r \ 0] \begin{bmatrix} 2 \\ 1 \end{bmatrix} + b = 1 \\ [-r \ 0] \begin{bmatrix} 3 \\ 1 \end{bmatrix} + b = -1 \end{cases} &\Rightarrow \begin{cases} -2r + b = 1 & \textcircled{1} \\ -3r + b = -1 & \textcircled{2} \end{cases} \end{aligned}$$

$$\Rightarrow \textcircled{1} - \textcircled{2} : r = 2.$$

$$\text{Sub. } r = 2 \text{ into } \textcircled{1} : b = 1 + 2r = 5$$

$$\Rightarrow \begin{cases} r = 2 \\ b = 5 \end{cases}$$

Therefore, $\underline{w^* = \begin{bmatrix} -r \\ 0 \end{bmatrix} = \begin{bmatrix} -2 \\ 0 \end{bmatrix}}$ and $\underline{b^* = 5}$ is the hard-margin SVM optimum solution.

Q2:

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \max \{1 - y_i(w^T x_i + b), 0\}^2$$

This is equivalent to

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}, \epsilon \in \mathbb{R}^n} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \epsilon_i^2$$

$$\text{s.t. } \forall i, (1 - y_i \hat{y}_i)^2 = \max \{1 - y_i \hat{y}_i, 0\}^2 \leq \epsilon_i^2$$

$$\Rightarrow \begin{cases} (1 - y_i \hat{y}_i)^2 \leq \epsilon_i^2 \\ 0 \leq \epsilon_i^2 \end{cases}$$

$$\hat{y}_i = w^T x_i + b$$

Since $\epsilon_i^2 \neq 0$ for every ϵ_i , we don't need the second condition

Since $(1 - y_i \hat{y}_i)^2 \leq \epsilon_i^2$, $1 - y_i \hat{y}_i \leq \epsilon_i$

The problem is converted to:

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}, \epsilon \in \mathbb{R}^n} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \epsilon_i^2$$

$$\text{s.t. } \forall i, 1 - y_i \hat{y}_i \leq \epsilon_i$$

Lagrangian is: $L(w, b, \epsilon, \alpha) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \epsilon_i^2 + \sum_{i=1}^n \alpha_i (1 - y_i (w^T x_i + b) - \epsilon_i)$

where $\alpha \geq 0$

Lagrangian Dual: $\min_{w, b, \epsilon} \max_{\alpha \geq 0, \beta \leq 0} L(w, b, \epsilon, \alpha, \beta)$

$$\Rightarrow \max_{\alpha \geq 0} \min_{w, b, \epsilon} L(w, b, \epsilon, \alpha)$$

$$\Rightarrow \max_{\alpha \geq 0} \min_{w, b, \epsilon} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \epsilon_i^2 + \sum_{i=1}^n \alpha_i (1 - y_i (w^T x_i + b) - \epsilon_i)$$

Gradients: $\frac{\partial L}{\partial w} = \frac{1}{2} \frac{\partial (\|w\|_2^2)}{\partial w} + \sum_{i=1}^n \frac{\partial (-\alpha_i y_i w^T x_i)}{\partial w}$

$$= \frac{1}{2} \cdot 2w + \sum_{i=1}^n (-\alpha_i y_i x_i)$$

$$= w - \sum_{i=1}^n \alpha_i y_i x_i = 0$$

$$\Rightarrow w^* = \sum_{i=1}^n \alpha_i y_i x_i$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^n \frac{\partial (\alpha_i y_i b)}{\partial b} = \sum_{i=1}^n \alpha_i y_i = 0$$

$$\begin{aligned} \frac{\partial L}{\partial \varepsilon_i} &= C \frac{\partial (\varepsilon_i^2)}{\partial \varepsilon_i} + \frac{\partial (-\alpha_i \varepsilon_i)}{\partial \varepsilon_i} \\ &= 2C \varepsilon_i - \alpha_i = 0 \end{aligned}$$

$$\Rightarrow \varepsilon_i^* = \frac{\alpha_i}{2C} \geq 0 \quad \text{as } \alpha_i \geq 0$$

Sub. into L :

$$\begin{aligned} & \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \varepsilon_i^2 + \sum_{i=1}^n \alpha_i (1 - y_i (w^T x_i + b) - \varepsilon_i) \\ &= \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \varepsilon_i^2 + \sum_{i=1}^n \alpha_i - \underbrace{\sum_{i=1}^n \alpha_i y_i x_i}_{w} w^T - \underbrace{\sum_{i=1}^n \alpha_i y_i}_0 b \\ & \quad - \sum_{i=1}^n \alpha_i \varepsilon_i \\ &= \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \varepsilon_i^2 + \sum_{i=1}^n \alpha_i - w w^T - \sum_{i=1}^n \alpha_i \varepsilon_i \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \varepsilon_i^2 - \sum_{i=1}^n \alpha_i \varepsilon_i \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \left(\frac{\alpha_i}{2C} \right)^2 - \sum_{i=1}^n \alpha_i \frac{\alpha_i}{2C} \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \|w\|_2^2 + C \cdot \frac{1}{4C^2} \sum_{i=1}^n \alpha_i^2 - \frac{1}{2C} \sum_{i=1}^n \alpha_i^2 \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \|w\|_2^2 + \frac{1}{4C} \sum_{i=1}^n \alpha_i^2 - \frac{1}{2C} \sum_{i=1}^n \alpha_i^2 \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \|w\|_2^2 - \frac{1}{4C} \sum_{i=1}^n \alpha_i^2 \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \left\| \sum_{i=1}^n \alpha_i y_i x_i \right\|_2^2 - \frac{1}{4C} \sum_{i=1}^n \alpha_i^2 \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j - \frac{1}{4C} \sum_{i=1}^n \alpha_i^2 \end{aligned}$$

The dual is : $\max_{d \geq 0} \sum_{i=1}^n d_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n d_i d_j y_i y_j x_i^T x_j - \frac{1}{4C} \sum_{i=1}^n d_i^2$

such that $\sum_{i=1}^n d_i y_i = 0,$

which is equivalent to :

$\min_{d \geq 0} - \sum_{i=1}^n d_i + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n d_i d_j y_i y_j x_i^T x_j + \frac{1}{4C} \sum_{i=1}^n d_i^2$

such that $\sum_{i=1}^n d_i y_i = 0.$

Q3: $C (1 - y_i \hat{y}_i)_+^2 \quad \hat{y}_i = w^T x_i + b \quad (2)$

Case 1: If $1 - y_i \hat{y}_i \leq 0$, $(1 - y_i \hat{y}_i)_+ = 0$

$$C (1 - y_i \hat{y}_i)_+^2 = 0$$

$$\frac{\partial}{\partial w} = 0$$

$$\frac{\partial}{\partial b} = 0$$

Case 2: If $1 - y_i \hat{y}_i > 0$, $(1 - y_i \hat{y}_i)_+ = 1 - y_i \hat{y}_i$

$$\begin{aligned} C (1 - y_i \hat{y}_i)_+^2 &= C (1 - y_i \hat{y}_i)^2 \\ &= C (1 - y_i (w^T x_i - b))^2 \\ &= C (1 - y_i x_i w^T - y_i b)^2 \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial w} &= C \cdot 2 (1 - y_i x_i w^T - y_i b) (-y_i x_i) \\ &= -2C y_i x_i (1 - y_i \hat{y}_i) \quad \text{as } 1 - y_i x_i w^T - y_i b = 1 - y_i \hat{y}_i \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial b} &= C \cdot 2 (1 - y_i x_i w^T - y_i b) (-y_i) \\ &= -2C y_i (1 - y_i x_i w^T - y_i b) \\ &= -2C y_i (1 - y_i \hat{y}_i) \quad \text{as } 1 - y_i x_i w^T - y_i b = 1 - y_i \hat{y}_i \end{aligned}$$

Therefore, by combining the two cases,

the gradient of equation (2) with respect to w is $-2C y_i x_i (1 - y_i \hat{y}_i)_+$,
and the gradient of equation (2) with respect to b is $-2C y_i (1 - y_i \hat{y}_i)_+$

Q4:

$$P_{w, \eta}^n = \arg \min_z \frac{1}{2\eta} \|z - w\|_2^2 + \frac{1}{2} \|z\|_2^2$$

$$\frac{\partial}{\partial z} = \frac{1}{2\eta} \frac{\partial}{\partial z} (\|z - w\|_2^2) + \frac{1}{2} \frac{\partial}{\partial z} (\|z\|_2^2)$$

$$= \frac{1}{2\eta} \cdot 2(z - w) + \frac{1}{2} \cdot 2z$$

$$= \frac{1}{\eta} (z - w) + z$$

$$\text{Let } \frac{\partial}{\partial z} = 0.$$

$$\frac{1}{\eta} (z - w) + z = 0$$

$$\frac{1}{\eta} z - \frac{1}{\eta} w + z = 0$$

$$\frac{1+\eta}{\eta} z = \frac{1}{\eta} w$$

$$z = \frac{1}{\eta} w \cdot \frac{\eta}{1+\eta}$$

$$\underline{z = \frac{1}{1+\eta} w}$$

Q5:

$$\text{From Q3, } \frac{\partial C(1-y_i\hat{y}_i)_+}{\partial w} = -2C y_i x_i (1-y_i\hat{y}_i)_+$$

$$\frac{\partial C(1-y_i\hat{y}_i)_+}{\partial b} = -2C y_i (1-y_i\hat{y}_i)_+$$

Equation (1) can be converted to:

$$\min_{w, b} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n (1-y_i\hat{y}_i)_+^2 \quad \hat{y}_i = w^T x_i + b$$

$$\frac{\partial}{\partial w} = \frac{1}{2} \cdot 2w + (-2C) \sum_{i=1}^n (1-y_i\hat{y}_i)_+$$

$$= w - 2C \sum_{i=1}^n (1-y_i\hat{y}_i)_+$$

$$\frac{\partial}{\partial b} = -2C \sum_{i=1}^n y_i (1-y_i\hat{y}_i)_+$$

$$\text{From Q4, } z = \frac{1}{1+\eta} w$$

Algorithm 1: SGD for SVM.

Input: $X \in \mathbb{R}^{n \times d}$, $y \in \{-1, 1\}^n$, $w = \mathbf{0}_d$, $b = 0$, $\text{max_pass} \in \mathbb{N}$, **step size** η

Output: w, b

1 for $t = 1, 2, \dots, \text{max_pass}$ do

2 for $i = 1, 2, \dots, n$ do

3 choose step size η

4 if $y_i(\langle x_i, w \rangle + b) \leq 1$ then

5 $w_t \leftarrow w_{t-1} - \eta \nabla w_{t-1}$

6 $b_t \leftarrow b_{t-1} - \eta \nabla b_{t-1}$

7 $w_t \leftarrow \frac{1}{1+\eta} w_t$

// x_i is the i -th row of X

// the proximal step to maximize margin

The implementation is in "code-part. ipynb".

Qb:

After running the dataset in "code-part.ipynb", the result is :

C

10000

W

array([-1.99805291e+00, -7.34489503e-04])

b

4.995740369998062

$$w = \begin{bmatrix} -1.998 \\ -7.344 \times 10^{-4} \end{bmatrix} \text{ is very closed to } w^* = \begin{bmatrix} -2 \\ 0 \end{bmatrix}$$

$$b = 4.99574 \text{ is very closed to } b^* = 5.$$

Hence, by using a large C, which is 10000, I recover the hard-margin SVM solution in Q1 be