

$$x_1, x_2, \dots, x_M, \dots, x_N$$

Let x be $y = -1$, o be $y = 1$

We know margin = $\min_{n=1, \dots, N} \text{distance}(x_n, w)$

Hence, in order to separate $y = -1, y = 1$, and also achieve max margin, we pick w be the middle

of x_{M+1}, x_M , $w = \frac{x_{M+1} + x_M}{2}$, the margin will be $\frac{|x_{M+1} - x_M|}{2}$



2.

$$\min_{w,b} \frac{1}{2} w^T w$$

subject to $(w^T x_n + b) \geq 1$ for $y_n = 1$ $0 \geq 1 - (w^T x_n + b)$ for $y_n = -1$

$- (w^T x_n + b) \geq -1$ for $y_n = -1$ $0 \geq 1 + (w^T x_n + b)$ for $y_n = -1$

First, we find the lagrange function, we suppose $y_1 \sim y_M = 1$, $y_{M+1} \sim y_N = -1$

$$L(b, w, \alpha) = \frac{1}{2} w^T w + \sum_{n=1}^M \alpha_n (1 - y_n (w^T x_n + b)) + \sum_{n=M+1}^N \alpha_n (-1 - y_n (w^T x_n + b))$$

Then, we solve the Lagrange dual.

$$\max_{\text{all } \alpha > 0} \left[\min_{b,w} \frac{1}{2} w^T w + \sum_{n=1}^M \alpha_n (1 - y_n (w^T x_n + b)) + \sum_{n=M+1}^N \alpha_n (-1 - y_n (w^T x_n + b)) \right]$$

$$\text{Inner problem "unconstrained" at optimal: } \frac{\partial L(b, w, \alpha)}{\partial b} = 0 = -\sum_{n=1}^N y_n \alpha_n$$

$$\text{then it becomes } \max_{\text{all } \alpha > 0, \sum_{n=1}^N y_n \alpha_n = 0} \left[\min_{b,w} \frac{1}{2} w^T w + \sum_{n=1}^M \alpha_n (1 - y_n w^T x_n) + \sum_{n=M+1}^N \alpha_n (-1 - y_n w^T x_n) \right]$$

$$\text{Inner problem "unconstrained" at optimal: } \frac{\partial L(b, w, \alpha)}{\partial w_i} = 0 = w_i - \sum_{n=1}^N \alpha_n y_n x_{ni} \Rightarrow w = \sum_{n=1}^N \alpha_n y_n x_n$$

$$\text{then it becomes } \max_{\text{all } \alpha > 0, \sum y_n \alpha_n = 0, w = \sum \alpha_n y_n x_n} -\frac{1}{2} \left| \left| \sum_{n=1}^N \alpha_n y_n x_n \right| \right|^2 + \sum_{n=1}^M \alpha_n + \sum_{n=M+1}^N \alpha_n$$



3.

constituted by (w_f, b_f)

We know that the points on the hyperplane will have the following property

$$\begin{cases} w_f^T x_n + b_f = 1 & \text{if } y_n = 1 \\ w_f^T x_n + b_f = -1 & \text{if } y_n = -1 \end{cases}$$

We can rewrite into the following equation.

$$w_f^T x_n + b_f = \frac{1+y_n}{2} - \frac{1-y_n}{2} \varphi = \frac{1}{2}(1-\varphi) + \frac{1}{2}(1+\varphi)y_n \dots \textcircled{1}$$

Moreover, we know the hyperplane obtained by $\varphi = 1$ has the property that

$$w_i^T x_n + b_1 = y_n \dots \textcircled{2}$$

We insert \textcircled{2} into \textcircled{1}

$$\underline{w_f^T x_n + b_f} = \frac{1}{2}(1-\varphi) + \frac{1}{2}(1+\varphi)(w_i^T x_n + b_1) = \underline{\frac{1}{2}(1+\varphi)(w_i^T x_n)} + \underline{\frac{1}{2}(1-\varphi) + \frac{1}{2}(1+\varphi)b_1}$$

Hence we have $w_f^T x_n = \frac{1}{2}(1+\varphi)(w_i^T x_n)$

$$b_f = \frac{1}{2}(1-\varphi) + \frac{1}{2}(1+\varphi)b_1 = \frac{1+b_1}{2} - \frac{1-b_1}{2}\varphi$$

since $w = \sum y_n x_n$, w and α are linear relationship.

$\alpha_\varphi = \alpha_1 (\frac{1}{2} + \frac{1}{2}\varphi)$, hence w_f is obtained correctly.

for $b_f = \frac{1+y_n}{2} - \frac{1-y_n}{2}\varphi - w_f^T x_n$ with any SV (z_n, y_n) , which the same

SV as $\varphi = 1$, since we only change the weight.

From \textcircled{1} & \textcircled{2}, b_f can be obtained correctly.

$$\text{Hence } w_f = \left(\frac{1}{2} + \frac{1}{2}\varphi\right) w_i$$

$$b_f = \frac{1+b_1}{2} - \frac{1-b_1}{2}\varphi$$

$$\text{We know } w_{11>b} = \left(\frac{1}{2} + 5b_3\right) w_i$$

$$b_{11>b} = \frac{1+b_1}{2} - \frac{1-b_1}{2} \times 112b \quad \text{※}$$

4. From Q3, we derive the value of $\alpha_f = \alpha_1 (\frac{1}{2} + \frac{1}{2} f)$

Hence, we can see that α_i and α_f have different optimal lagrange multiplier.

Therefore, meaning that α_f is the only optimal solution for the uneven margin SVM.

α^* is not an optimal solution of the uneven SVM -

$$5. K(x, x') = K_1(x, x') K_2(x, x')$$

$$= \phi_1(x)^T \phi_1(x') \phi_2(x)^T \phi_2(x')$$

$$= \left(\sum_{i=1}^d \phi_{1i}(x) \phi_{1i}(x') \right) \left(\sum_{j=1}^d \phi_{2j}(x) \phi_{2j}(x') \right)$$

$$= \sum_{i=1}^d \sum_{j=1}^d (\phi_{1i}(x) \phi_{2j}(x)) (\phi_{1i}(x') \phi_{2j}(x'))$$

$$= \sum_{i=1}^d \sum_{j=1}^d \phi_{ij}(x) \phi_{ij}(x')$$

$$= \phi(x) \cdot \phi(x')$$

where $\phi(x)$ is a $d \times d$ vector, s.t. $\phi_{ij}(x) = \phi_{1i}(x) \phi_{2j}(x)$

Now, we see that we can write $K(x, x')$ as an inner product

using $\phi(x)$, we know $K(x, x')$ is a valid kernel.

$$\phi_1(x) = \begin{pmatrix} \phi_{11}(x) \\ \phi_{12}(x) \\ \vdots \\ \phi_{1d}(x) \end{pmatrix}$$

$$\phi_2(x) = \begin{pmatrix} \phi_{21}(x) \\ \phi_{22}(x) \\ \vdots \\ \phi_{2d}(x) \end{pmatrix}$$

$$b. \quad \|\phi(x)\|^2 = \phi(x)^\top \phi(x) = K(x, x)$$

$$\|\phi(x)\| = \sqrt{K(x, x)}$$

$$\|\phi(x) - \phi(x')\|^2 = \|\phi(x)\|^2 + \|\phi(x')\|^2 - 2\phi(x)^\top \phi(x')$$

$$= K(x, x) + K(x', x') - 2K(x, x')$$

$$\|\phi(x) - \phi(x')\| = \sqrt{K(x, x) + K(x', x') - 2K(x, x')} \quad \Rightarrow \text{Since } x \text{ & } x' \text{ are unit vectors, } K(x, x) \text{ & } K(x', x') \text{ are both constant.}$$

yields when

largest possible distance $\sqrt{-2K(x, x')}$ is as large as possible:

$$\Rightarrow K(x, x') = -2(1 + x^\top x') \text{ achieves its largest value when } x^\top x' = -1$$

which is x and x' point toward the opposite direction,

Suppose $x = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$, $x' = \begin{bmatrix} -1 \\ 0 \end{bmatrix}$, then we have

$$\|\phi(x) - \phi(x')\| = \sqrt{4 + 4 - 0} = 2\sqrt{2}$$

Smallest distance yields when $-2K(x, x')$ is as small as possible

$$-2K(x, x') = -2(1 + x^\top x') \text{ achieves its smallest value when } x^\top x' = 1$$

which means x and x' are parallel, pointing toward the same

direction. Suppose $x = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$, $x' = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$, then we have

$$\|\phi(x) - \phi(x')\| = \sqrt{4 + 4 - 8} = 0$$

7.

$$\|\tilde{\phi}(x)\| = \sqrt{\sum_{i=0}^{\infty} \frac{x^i}{i!}} = \sqrt{\sum_{i=0}^{\infty} \frac{1}{i!} x^{2i}} = \sqrt{\sum_{i=0}^{\infty} \frac{(2x)^i}{i!}} = \sqrt{e^{2x}} = e^x = \exp(x)$$

$\Rightarrow e^x = \frac{1}{\|\tilde{\phi}(x)\|}$

↑ Taylor's series of e^{2x}

8.

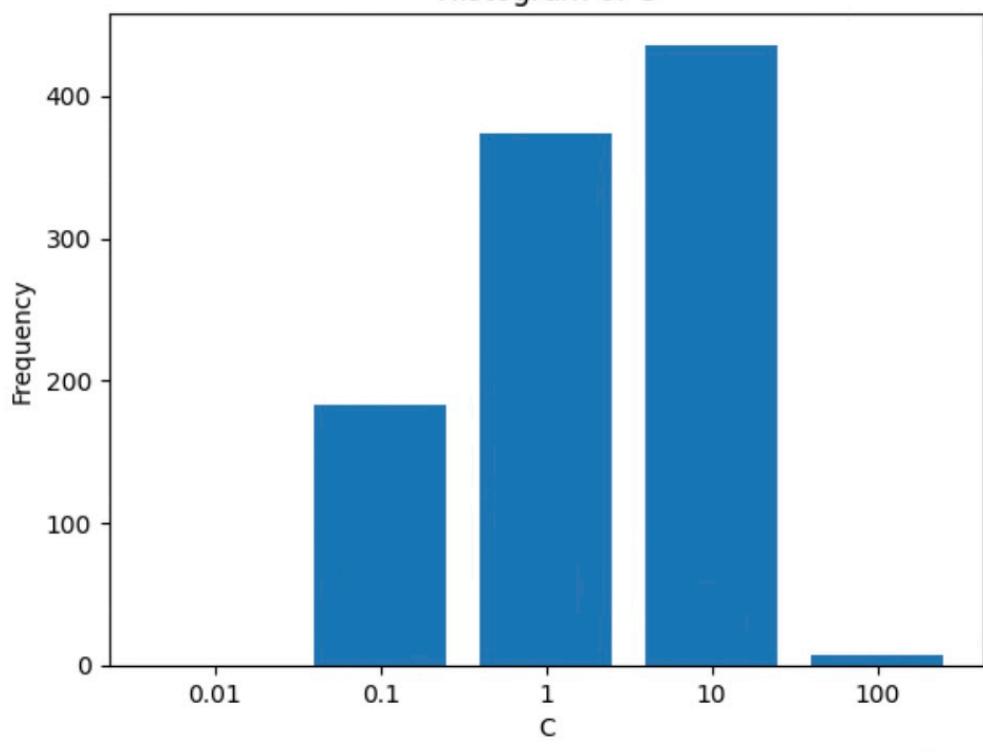
$$k(x, x') = \frac{x^T x'}{\|x\| \|x'\|} = \frac{1}{\|x^T\|} x^T \cdot \frac{1}{\|x'\|} x' = \phi(x)^T \phi(x')$$

$\phi(x) = \frac{1}{\|x\|} x$, we can write $k(x, x')$ as inner product of $\phi(x)$, hence $k(x, x')$ is a valid kernel.

9 Minimum number of support vectors = 629, C = 10, Q = 4

10 Minimum Eout = 0.005000000000000044, C = 1

Histogram of C



We can see that $C = 10$ has the highest frequency to have the best Eval.

Secondly, $C = 1$, then $C = 0.1$, lastly $C = 100$.

We know that C limits the value of a_n , when C is low it has higher margin, when C is high it will have lower margin.

This is the trade-off between low training error and low testing error.

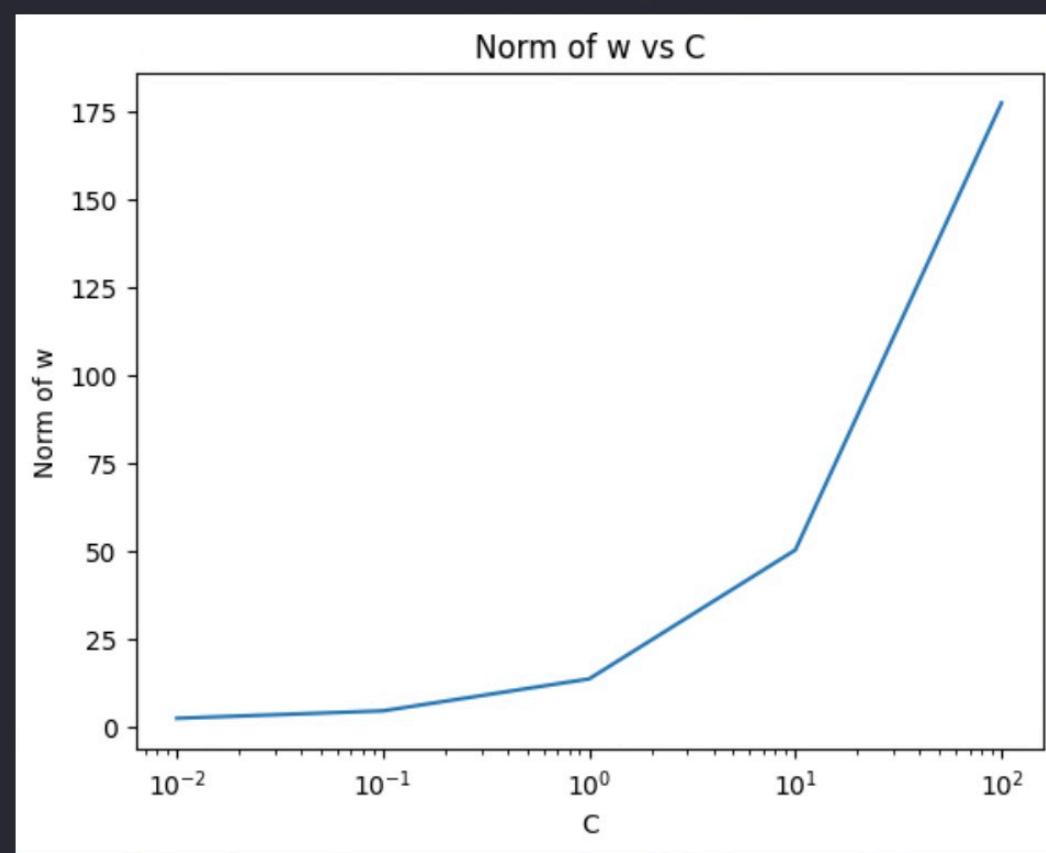
We want C to be big enough that it achieves good Ein but not overfitting, which results in bad Eval.

In this case that we see when C rises from 0.1, the frequency of each value increases.

Eventually, we reach $C = 10$, which is the best value of C , having highest frequency of best Eval.

Meaning when we train with $C = 10$, we have high chance to get both good Ein and Eval.

When we increase C to 100, the Ein might be better, but the model overfitted which result in bad Eval, we can see the frequency dropped significantly.



From the plot we can see that the norm of w have a bit of positive correlation to C.

With small $|w|$ results in small C, large $|w|$ results in large C.

Moreover, we know that small C will have larger margin, large C will have small margin.

Hence, we can conclude that large norm $|w|$ will results in small margin, which reduces the E_{in} , however might cause E_{val} to be worse if we made it too big.



13.

$$\text{Primal SVM: } \min_{w, b} \frac{1}{2} w^T w \quad \text{s.t. } y_n(w^T x_n + b) \geq 1 \quad \text{for } n = 1, 2, \dots, N$$

$$L(b, w, \alpha) = \underbrace{\frac{1}{2} w^T w}_{\text{obj.}} + \sum_{n=1}^N \alpha_n (1 - \underbrace{y_n(w^T x_n + b)}_{\text{constraint}})$$

$$\text{Lagrange problem for primal SVM: } \min_{b, w} \left(\max_{\alpha \geq 0} L(b, w, \alpha) \right)$$

$$\text{Lagrange dual problem: } \max_{\alpha \geq 0} \left(\min_{b, w} L(b, w, \alpha) \right)$$

$$\text{Dual SVM: } \min_{\alpha} \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y_n y_m z_n^T z_m - \sum_{n=1}^N \alpha_n$$

$$\text{subject to } \sum_{n=1}^N y_n \alpha_n = 0, \quad \alpha_n \geq 0, \quad \text{for } n = 1, 2, \dots, N$$

$$L(\alpha, \lambda) = \underbrace{\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y_n y_m z_n^T z_m - \sum_{n=1}^N \alpha_n}_{\text{obj.}} - \underbrace{\sum_{n=1}^N \lambda_n \alpha_n}_{\text{constraint.}} = \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y_n y_m z_n^T z_m - \sum_{n=1}^N (1 + \lambda_n) \alpha_n$$

$$\text{Lagrange problem for Dual SVM: } \min_{\alpha} \left(\max_{\lambda \geq 0} (L(\alpha, \lambda)) \right)$$

$$= \min_{\alpha} \left(\infty \text{ if violate, } \underbrace{\text{obj.}}_{\downarrow} \text{ if feasible} \right)$$

$$\alpha \downarrow \quad \alpha_n < 0 \quad \alpha \geq 0$$

$$\alpha_n \text{ goes to } \infty \quad \alpha = 0$$

$$\| - (1 + \lambda_n) \alpha_n \| \rightarrow + \quad \| - (1 + \lambda_n) \alpha_n \| \rightarrow -$$

$$\min_{\alpha} \left(\max_{\lambda \geq 0} (L(\alpha, \lambda)) \right) \geq \min_{\lambda} L(\alpha, \lambda) \quad \text{for any fix } \lambda' \text{ with all } \lambda'_n \geq 0 \\ \text{if } \max_{\lambda} \geq \text{any}$$

$$\min_{\alpha} \left(\max_{\lambda \geq 0} (L(\alpha, \lambda)) \right) \geq \max_{\lambda} \left(\min_{\alpha} L(\alpha, \lambda') \right) \quad \text{for best } \lambda' \text{ on RHS.} \\ \text{if best is one of any}$$

$$\text{Lagrange dual for dual SVM: } \max_{\lambda \geq 0} \left(\min_{\alpha} L(\alpha, \lambda) \right)$$

$$\text{Lagrange for primal SVM: } \min_{b, w} \left(\max_{\alpha \geq 0} L(b, w, \alpha) \right)$$

As we can see, the max and min are opposite, and they do optimization on different variables. From chat bPT, the dual of dual is not always the same as the primal problem can be confirm.