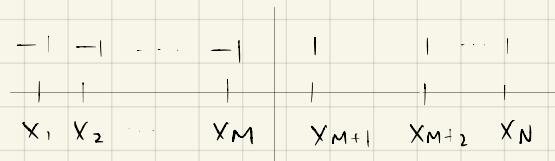


1. (20 points) Consider  $N$  “linearly separable” 1D examples  $\{(x_n, y_n)\}_{n=1}^N$ . That is,  $x_n \in \mathbb{R}$ . Without loss of generality, assume that  $x_1 \leq x_2 \leq \dots \leq x_M < x_{M+1} \leq x_{M+2} \leq \dots \leq x_N$ ,  $y_n = -1$  for  $n = 1, 2, \dots, M$ , and  $y_n = +1$  for  $n = M+1, M+2, \dots, N$ . Apply hard-margin SVM without transform on this data set. What is the largest margin achieved in terms of  $\{(x_n, y_n)\}_{n=1}^N$  and  $M$ ? Prove your answer.

(Note: This is a “hard-margin” decision stump. Have we mentioned that a decision stump is just a 1D perceptron, and the hard-margin SVM is an extension of the perceptron model? :-))



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

$$\frac{x_{M+1} - x_M}{2}$$

\*

2. (20 points) In some situations, we expect to achieve a smaller margin for the positive examples and a larger margin for the negative examples. For instance, when there are very few negative examples and a lot more positive examples, giving the negative examples a bigger margin could be more robust. Consider an *uneven-margin* SVM that solves

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{subject to} \quad & (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \text{ for } y_n = +1 \\ & -(\mathbf{w}^T \mathbf{x}_n + b) \geq \rho \text{ for } y_n = -1. \end{aligned}$$

Derive the (Lagrange) dual problem for the uneven margin SVM above. List your derivation steps.

(Note: This is a very useful extension of SVM, and more importantly, "forces" you to go through the derivation steps of the class slides to gain a deeper understanding of how to derive a dual problem.)

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \quad \text{subject to} \quad \begin{cases} (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 & y_n = +1 \\ -(\mathbf{w}^T \mathbf{x}_n + b) \geq \rho & y_n = -1 \end{cases} \\ \Rightarrow \quad & y_n (\mathbf{w}^T \mathbf{x}_n + b + [\![y_n = -1]\!] (\rho - 1)) \geq 1 \end{aligned}$$

$$\mathcal{L}(\mathbf{w}, b, \alpha_n) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{n=1}^N \alpha_n (1 - y_n (\mathbf{w}^T \mathbf{x}_n + b + [\![y_n = -1]\!] (\rho - 1)))$$

$$\begin{aligned} \text{SVM} = \min_{\mathbf{w}, b} \left( \max_{\alpha_n \geq 0} \mathcal{L}(\mathbf{w}, b, \alpha_n) \right) & \geq \max_{\alpha_n \geq 0} \left( \min_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b, \alpha_n) \right) \\ & = \max_{\alpha_n \geq 0} \left( \min_{\mathbf{w}, b} \left( \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{n=1}^N \alpha_n (1 - y_n (\mathbf{w}^T \mathbf{x}_n + b + [\![y_n = -1]\!] (\rho - 1))) \right) \right) \\ \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w}_i &= \sum_{n=1}^N \alpha_n y_n z_{n,i} \\ \frac{\partial \mathcal{L}}{\partial b} = 0 \Rightarrow \sum_{n=1}^N \alpha_n y_n &= 0 \\ = \max_{\alpha_n \geq 0, \sum \alpha_n y_n = 0, \mathbf{w} = \sum_{n=1}^N \alpha_n y_n z_n} & - \frac{1}{2} \left\| \sum_{n=1}^N \alpha_n y_n z_n \right\|^2 + \sum_{n=1}^N \alpha_n \end{aligned}$$

$\Rightarrow$  Lagrange dual problem

3. (20 points) Let's call the primal solution to the uneven margin SVM above  $(b_1^*, \mathbf{w}_1^*)$ . Assume that Dr. Even has solved the original even-margin SVM

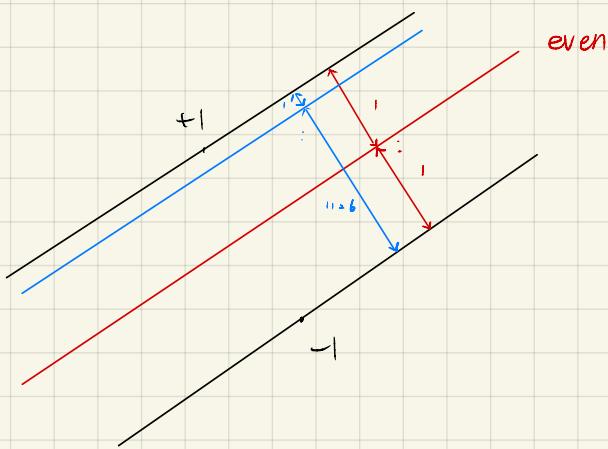
$$\begin{array}{ll} \min_{\mathbf{w}, b} & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{subject to} & (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \text{ for } y_n = +1 \\ & -(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \text{ for } y_n = -1. \end{array}$$

That is, Dr. Even obtained  $(b_1^*, \mathbf{w}_1^*)$ . Somehow zir boss likes the magic number of 1126 better and wants zir to return  $(b_{1126}^*, \mathbf{w}_{1126}^*)$  instead. Describe a systematic process to convert  $(b_1^*, \mathbf{w}_1^*)$  to  $(b_{1126}^*, \mathbf{w}_{1126}^*)$  without solving the QP problem again. Prove the correctness of your process.

(Note: The result should reinforce your geometric sense of what SVM is doing.)

even margin  $\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w}$  subject to  $y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1$

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w} \quad \text{subject to} \quad \begin{cases} (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 & y = +1 \\ -(\mathbf{w}^T \mathbf{x}_n + b) \geq -1 & y = -1 \end{cases}$$



$$\text{even : } (\mathbf{x}^+ - \mathbf{x}^-) \cdot \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

$$= \frac{1}{\|\mathbf{w}\|} (\mathbf{x}^+ \mathbf{w} - \mathbf{x}^- \mathbf{w})$$

$$= \frac{1}{\|\mathbf{w}\|} ((1-b) - (-1-b))$$

$$= \frac{2}{\|\mathbf{w}\|}$$

$$\frac{2}{\|\mathbf{w}\|} = \frac{1127}{\|W_{1126}\|} \quad \|W_{1126}\| = \frac{1127}{2} \|\mathbf{w}\|$$

$$\text{uneven : } (\mathbf{x}^+ - \mathbf{x}^-) \cdot \frac{\mathbf{w}_{1126}}{\|W_{1126}\|}$$

$$= \frac{1}{\|W_{1126}\|} (W_{1126} \mathbf{x}^+ - W_{1126} \mathbf{x}^-)$$

$$= \frac{1}{\|W_{1126}\|} ((1-b) - (-1126-b))$$

$$= \frac{1127}{\|W_{1126}\|}$$

$$\mathbf{w} \cdot \mathbf{x}^+ + b_1 = W_{1126} \mathbf{x}^+ + b_{1126} = 1$$

$$\Rightarrow b_{1126} = b_1 + (\mathbf{w} - W_{1126}) \mathbf{x}^+$$

$$= b_1 + -\frac{1125}{2} \mathbf{w} \cdot \mathbf{x}^+ = b_1 - \frac{1125}{2} (1-b_1) = \frac{1127}{2} b_1 - \frac{1125}{2}$$

#

4. (20 points) Let's call the dual solution to the uneven margin SVM above  $\alpha_\rho^*$ . Assume that Dr. Even has solved the dual problem of the even-margin SVM and obtained  $\alpha_1^*$ . Prove or disprove that  $\alpha_1^*$  is also an optimal solution of the uneven margin SVM above for any  $\rho > 0$ .

(Note: Think about how your geometric sense from the previous problem can be extended to algebraic results.)

$$\mathcal{L}(w, b, \alpha_n) = \frac{1}{2} w^\top w + \sum_{n=1}^N \alpha_n (1 - y_n (w^\top x_n + b + [\mathbb{y} = -1](\rho - 1)))$$

$$\text{even} \Rightarrow \rho = 1 \quad \text{uneven} \Rightarrow \rho = \ell$$

for  $\rho = 1$  or uneven, support vector remain the same  $\Rightarrow \alpha_n y_n$  are the same

$$\frac{\partial \mathcal{L}}{\partial b} = -\sum_{n=1}^N \alpha_n y_n = 0 \Rightarrow \sum_{n=1}^N \alpha_n y_n = 0 \Rightarrow \text{independent of } \rho$$

$$\frac{\partial \mathcal{L}}{\partial w_i} = w_i - \sum_{n=1}^N \alpha_n y_n x_n = 0 \Rightarrow w_i = \sum_{n=1}^N \alpha_n y_n x_n \Rightarrow \text{independent of } \rho$$

the optimal solution  $\alpha_\rho^*$  is also the optimal solution for  $\alpha_1^*$

5. (20 points) For two valid kernels  $K_1(\mathbf{x}, \mathbf{x}') = \phi_1(\mathbf{x})^T \phi_1(\mathbf{x}')$  and  $K_2(\mathbf{x}, \mathbf{x}') = \phi_2(\mathbf{x})^T \phi_2(\mathbf{x}')$ , consider a kernel function  $K(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}') \cdot K_2(\mathbf{x}, \mathbf{x}')$ . Prove that  $K$  is a valid kernel by deriving a transform function  $\phi(\mathbf{x})$  such that

$$K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}').$$

(Note: The result shows that a multiplication of valid kernels is still a valid kernel. So, for instance, a multi-variate Gaussian kernel is a valid kernel; a polynomial kernel—repeated multiplications of the identity  $\phi$ —is a valid kernel.)

$$\begin{aligned} K(\mathbf{x}, \mathbf{x}') &= \phi_1(\mathbf{x})^T \phi_1(\mathbf{x}') + \phi_2(\mathbf{x})^T \phi_2(\mathbf{x}') \\ &= \sum_{i=1}^n \phi_{1,i}(\mathbf{x}) \cdot \phi_{1,i}(\mathbf{x}') + \sum_{j=1}^m \phi_{2,j}(\mathbf{x}) \cdot \phi_{2,j}(\mathbf{x}') \\ &= \sum_{i=1}^n \sum_{j=1}^m (\phi_{1,i}(\mathbf{x}) \cdot \phi_{2,j}(\mathbf{x})) \cdot (\phi_{1,i}(\mathbf{x}') \cdot \phi_{2,j}(\mathbf{x}')) \\ &= \sum_{k=1}^n \sum_{l=1}^m \phi_{3,k}(\mathbf{x}) \cdot \phi_{3,l}(\mathbf{x}') \quad \phi_{3,k} = \phi_{1,i} \cdot \phi_{2,j} \\ &= \phi_3(\mathbf{x})^T \phi_3(\mathbf{x}') \end{aligned}$$

$\Rightarrow$  still a valid kernel.

6. (20 points) For any feature transform  $\phi$  from  $\mathcal{X}$  to  $\mathcal{Z}$ , the  $\mathcal{Z}$ -space distance between two examples  $\mathbf{x}$  and  $\mathbf{x}'$  is  $\|\phi(\mathbf{x}) - \phi(\mathbf{x}')\|$ , which can be computed with the kernel trick. Consider the degree-2 quadratic kernel  $K_2(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^T \mathbf{x}')^2$ . For two unit vectors  $\mathbf{x}$  and  $\mathbf{x}'$ , what are their largest and smallest possible distances in the  $\mathcal{Z}$  space? Prove your claim.

(Note: This shall sharpen your understanding of the kernel trick.)

$$K_2(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^T \mathbf{x}')^2$$

$$\phi(\mathbf{x}) = (1, \sqrt{2}\mathbf{x}_1, \sqrt{2}\mathbf{x}_2, \dots, \mathbf{x}_1^2, \mathbf{x}_2^2, \mathbf{x}_3^2, \dots)$$

$$\|\mathbf{x}\| = \|\mathbf{x}'\| = 1$$

$$\|\phi(\mathbf{x})\| = \sqrt{\phi(\mathbf{x})^T \phi(\mathbf{x})} = \sqrt{K_2(\mathbf{x}, \mathbf{x})} = \sqrt{(1 + \mathbf{x}^T \mathbf{x})^2} = 2$$

$$\max K_2(\mathbf{x}, \mathbf{x}') = (1 + 1)^2 = 4 \quad \text{occurs when } \mathbf{x}, \mathbf{x}' \text{ are same direction}$$

$$\min K_2(\mathbf{x}, \mathbf{x}') = (1 - 1)^2 = 0 \quad \text{occurs when } \mathbf{x}, \mathbf{x}' \text{ are opposite direction}$$

$$\|\phi(\mathbf{x}) - \phi(\mathbf{x}')\| = \sqrt{\phi(\mathbf{x})^T \phi(\mathbf{x}) + \phi(\mathbf{x}')^T \phi(\mathbf{x}') - 2\phi(\mathbf{x})^T \phi(\mathbf{x}')}}$$

$$= \sqrt{\|\phi(\mathbf{x})\|^2 + \|\phi(\mathbf{x}')\|^2 - 2\phi(\mathbf{x})^T \phi(\mathbf{x}')}}$$

$$= \sqrt{8 - 2\phi(\mathbf{x})^T \phi(\mathbf{x}')}$$

$$\max \|\phi(\mathbf{x}) - \phi(\mathbf{x}')\| = \sqrt{8 - 2 \times 0} = 2\sqrt{2}$$

$$\min \|\phi(\mathbf{x}) - \phi(\mathbf{x}')\| = \sqrt{8 - 2 \times 4} = 0$$

7. (20 points) Recall that in the derivation of (one-dimensional) Gaussian kernel, we derived its associated transform

$$\phi(x) = \exp(-x^2) \cdot \left(1, \sqrt{\frac{2}{1!}}x, \sqrt{\frac{2^2}{2!}}x^2, \dots\right).$$

Let

$$\tilde{\phi}(x) = \left(1, \sqrt{\frac{2}{1!}}x, \sqrt{\frac{2^2}{2!}}x^2, \dots\right).$$

Prove that  $\exp(-x^2) = \frac{1}{\|\tilde{\phi}(x)\|}$ .

(Note: The result means  $\phi(x)$  can be viewed as a normalized version of  $\tilde{\phi}(x)$ .)

Taylor series expansion of  $e^{x^2} = \sum_{n=0}^{\infty} \frac{(x^2)^n}{n!}$

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

$$\|\tilde{\phi}(x)\| = \sum_{n=0}^{\infty} \frac{(x^2)^n}{n!} = e^{x^2}$$

$$\exp(-x^2) = \frac{1}{\|\tilde{\phi}(x)\|}$$

\*

8. (20 points) Let  $\cos(\mathbf{x}, \mathbf{x}')$  measure the cosine of the angle between two non-zero vectors  $\mathbf{x}$  and  $\mathbf{x}'$  in  $\mathbb{R}^d$ . The function  $\cos$  is typically called the cosine similarity between two vectors. Prove or disprove that  $\cos(\mathbf{x}, \mathbf{x}')$  is a valid kernel on  $\mathcal{X} = \mathbb{R}^d - \{\mathbf{0}\}$ .

(Hint: to prove, you'd better construct its associated transform; to disprove, you may use Mercer's condition to construct a counter-example on positive semi-definiteness.)

$$K_{ij} = \cos(x_i, x_j) = \frac{x_i \cdot x_j}{\|x_i\| \|x_j\|} = \frac{x_i}{\|x_i\|} \cdot \frac{x_j}{\|x_j\|} = p_i(x) \cdot p_j(x)$$

$$\cos(x, x') = \frac{x^T}{\|x\|} \cdot \frac{x'}{\|x'\|} = \frac{x'^T}{\|x'\|} \cdot \frac{x}{\|x\|} = \cos(x', x) \Rightarrow \text{symmetric}$$

To prove that  $K$  is positive semi-definite, we consider a vector  $c = \begin{pmatrix} c_1 \\ c_2 \\ c_3 \\ \vdots \end{pmatrix}$

$$\begin{aligned} c^T K c &= \sum_{i=0}^n \sum_{j=0}^n c_i c_j \cos(x_i, x_j) = \sum_{i=0}^n \sum_{j=0}^n \frac{c_i c_j}{\|x_i\| \|x_j\|} x_i \cdot x_j \\ &= \sum_{i=0}^n \sum_{j=0}^n \left( \frac{c_i}{\|x_i\|} x_i \cdot \frac{c_j}{\|x_j\|} x_j \right) \\ &= \left| \sum_{i=0}^n \frac{c_i}{\|x_i\|} x_i \right|^2 \geq 0 \end{aligned}$$

$K$  is positive semi-definite

$\Rightarrow \cos(x, x')$  is valid kernel.

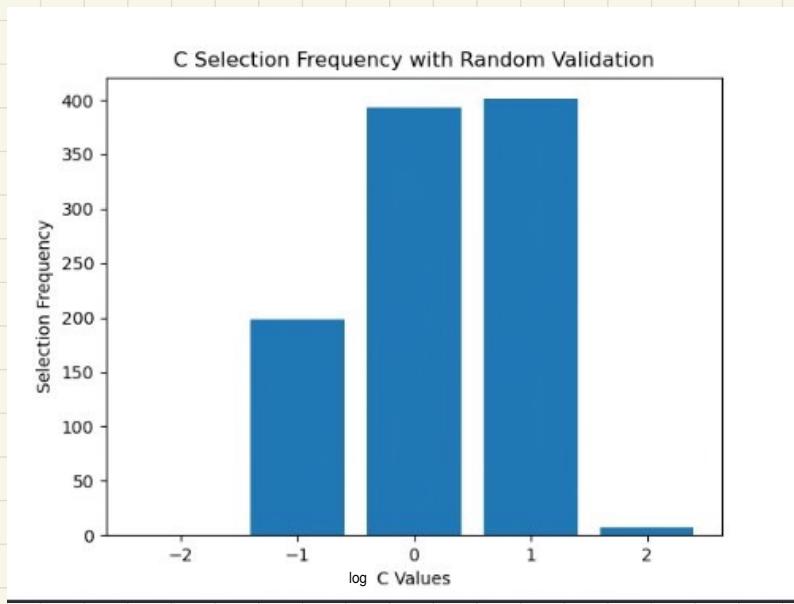
9. (20 points, \*) Consider the polynomial kernel  $K(\mathbf{x}_n, \mathbf{x}_m) = (1 + \mathbf{x}_n^T \mathbf{x}_m)^Q$ , where  $Q$  is the degree of the polynomial. Consider a binary classifier of “4” versus “not 4”. Among  $C \in \{0.1, 1, 10\}$  and  $Q \in \{2, 3, 4\}$ , count the number of support vectors of the resulting SVM classifier. Which of the  $(C, Q)$  combination(s) result in the smallest number of support vectors? Provide your command/code.

Best combination: (C=10, Q=4), Number of support vectors = 629

10. (20 points, \*) Consider the Gaussian kernel  $K(\mathbf{x}_n, \mathbf{x}_m) = \exp(-\gamma \|\mathbf{x}_n - \mathbf{x}_m\|^2)$ . For the binary classification problem of “1” versus “not 1”, when fixing  $\gamma = 1$ , among  $C \in \{0.01, 0.1, 1, 10, 100\}$ , which value(s) of  $C$  result in the lowest  $E_{\text{out}}$ ? Provide your command/code.

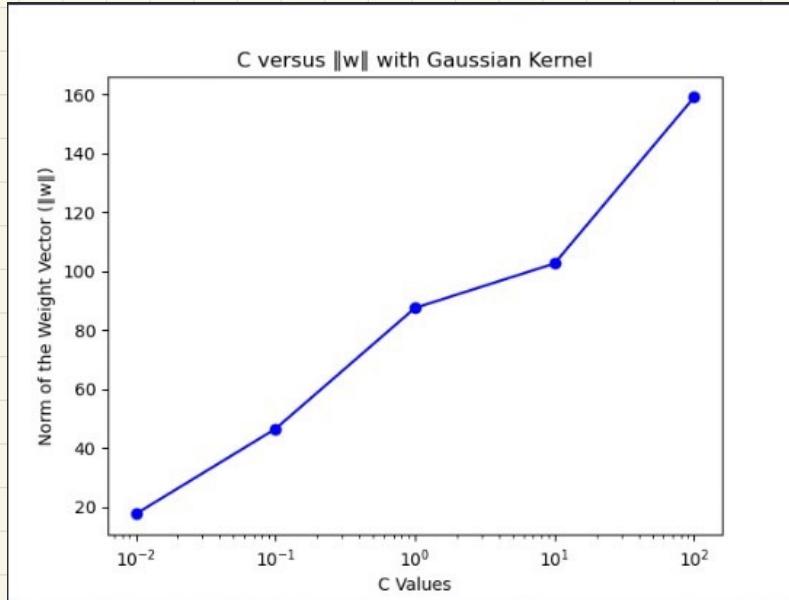
Best value of C: 1, Minimum Eout = 0.5%

11. (\*) Following Problem 10 and consider a validation procedure that randomly samples 200 examples from the training set for validation and leaves the other examples for training  $g_{\text{SVM}}^-$ . Fix  $\gamma = 1$  and use the validation procedure to choose the best  $C$  among  $\{0.01, 0.1, 1, 10, 100\}$  according to  $E_{\text{val}}$ . If there is a tie of  $E_{\text{val}}$ , choose the smallest  $C$ . Repeat the procedure 1000 times and count the number of times that each  $C$  is selected. Plot a bar chart of  $C$  versus its selection frequency. Compare the result to that of Problem 10. Describe your findings. Provide your command/code.



$C = 10$  has the lowest Eval for most of the times.  $C = 1$  has a little bit lower selection frequency than  $C=1$ .  $C = 0.01$  or  $100$  are rarely be the lowest Eval for most of the time.

12. (20 points, \*) Consider the Gaussian kernel  $K(\mathbf{x}_n, \mathbf{x}_m) = \exp(-\gamma \|\mathbf{x}_n - \mathbf{x}_m\|^2)$ . For the binary classification problem of “3” versus “not 3”, when fixing  $\gamma = 1$ , vary  $C \in \{0.01, 0.1, 1, 10, 100\}$  and calculate  $\|\mathbf{w}\|$  of the resulting SVM classifier with the kernel trick. Plot a line chart of  $C$  versus  $\|\mathbf{w}\|$ . Describe your findings. Provide your command/code.



with the  $C$  being larger .  $\|\mathbf{w}\|$  is getting larger . The value and  $|\mathbf{w}|$  are positive correlation .

13. (Bonus 20 points) Derive a simplified Langrange dual problem of the hard-margin SVM dual. Is your Lagrange dual problem of the hard-margin SVM dual the same as the hard-margin SVM primal? Are they “similar” in any sense? Describe your findings.

For your information, this is what chatGPT says.



For quadratic programming, is the dual problem of the dual problem the same as the primal problem?



In quadratic programming, the relationship between the primal and dual problems is more intricate than in linear programming. The dual problem of the dual problem may not necessarily be the same as the primal problem.

In quadratic programming, the primal problem involves minimizing a quadratic objective function subject to linear constraints. The dual problem is derived from this primal problem, and it involves maximizing a concave quadratic function subject to certain constraints.

When you take the dual of the dual problem in quadratic programming, it may not directly lead you back to the primal problem. The relationship is more complex due to the quadratic nature of the objective function.

In summary, for quadratic programming, the dual of the dual is not always the same as the primal problem.

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

$$\begin{aligned} \text{SVM} &= \min_{w, b} \left( \max_{\alpha_n > 0} \mathcal{L}(w, b, \alpha_n) \right) \geq \max_{\alpha_n > 0} \left( \min_{w, b} \mathcal{L}(w, b, \alpha_n) \right) \\ &= \max_{\alpha_n > 0} \left( \min_{w, b} \left( \frac{1}{2} w^T w + \sum_{n=1}^N \alpha_n (1 - y_n (w^T x_n + b)) \right) \right) \end{aligned}$$

$\min(\max \mathcal{L}(w, b, \alpha_n)) = \max(\min \mathcal{L}(w, b, \alpha_n))$  is called strong duality

if convex primal . feasible primal , linear constraints . the dual of the dual problem can be reversed back to the primal .