

IE643: Deep Learning - Theory and Practice

Quiz 1 Solution

July-Dec 2024

1. Consider the following data set: $D = \{(x^1, y^1), (x^2, y^2), (x^3, y^3), (x^4, y^4)\}$, where $(x^1, y^1) = ((0, 1), -1)$, $(x^2, y^2) = ((1, 0), -1)$, $(x^3, y^3) = ((0, 3), +1)$ and $(x^4, y^4) = ((3, 0), +1)$.

- (a) [3 marks] Find a linear separator for dataset D . Clearly write the expression for the linear separator and justify why your construction is in fact a linear separator for D .

let, take a linear separator $x_1 w_1 + x_2 w_2 - b = 0$

x_1	x_2	y		
0	1	$\text{sign}(w_2 - b)$	-1	$w_2 - b < 0 \quad \text{---(1)}$
1	0	$\text{sign}(w_1 - b)$	-1	$w_1 - b < 0 \quad \text{---(2)}$
0	3	$\text{sign}(3w_2 - b)$	+1	$3w_2 - b > 0 \quad \text{---(3)}$
3	0	$\text{sign}(3w_1 - b)$	+1	$3w_1 - b > 0 \quad \text{---(4)}$

from eq (3) and (4)

$$3w_1 + 3w_2 - 2b > 0$$

$$\Rightarrow \frac{3}{2} (w_1 + w_2) > b \quad \text{---(5)}$$

from eq (5) and (6)

$$\boxed{\frac{1}{2} (w_1 + w_2) < b < \frac{3}{2} (w_1 + w_2)}$$

from eq (1) and (2)

$$w_1 + w_2 - 2b < 0$$

$$\Rightarrow \frac{1}{2} (w_1 + w_2) < b \quad \text{---(6)}$$

from eq (1) & (3)

$$\begin{aligned} 3w_2 - b &> 0 \\ -w_2 + b &> 0 \\ \hline 2w_2 &> 0 \end{aligned}$$

$$\boxed{w_2 > 0}$$

from (2) and (4)

$$3w_1 - b > 0$$

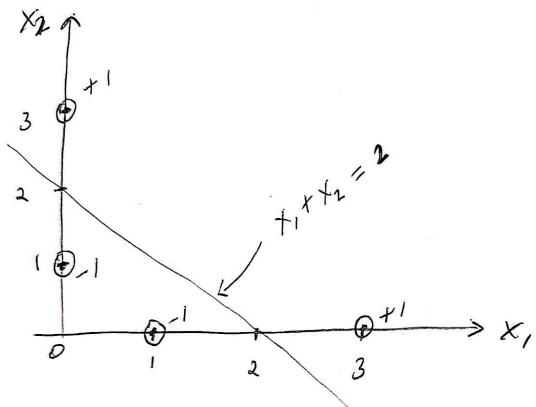
$$\begin{array}{r} -w_1 + b > 0 \\ \hline 2w_1 > 0 \end{array}$$

$$\boxed{w_1 > 0}$$

let take $w_1 = w_2 = 1$, then $\boxed{1 < b < 3}$

\therefore linear separator will be, taking $b = 2$

$$x_1 + x_2 - 2 = 0$$



1. Consider the following data set: $D = \{(x^1, y^1), (x^2, y^2), (x^3, y^3), (x^4, y^4)\}$, where $(x^1, y^1) = ((0, 1), -1)$, $(x^2, y^2) = ((1, 0), -1)$, $(x^3, y^3) = ((0, 3), +1)$ and $(x^4, y^4) = ((3, 0), +1)$.

(a) [3 marks] Find a linear separator for dataset D . Clearly write the expression for the linear separator and justify why your construction is in fact a linear separator for D .

(b) [6 marks] Suppose the linear separator you found in part 1a is represented as (w, b) such that for any $x \in \mathbb{R}^2$ on the linear separator, $\langle w, x \rangle - b = 0$ holds. Transform D into a new dataset $D' = \{(q^1, y^1), (q^2, y^2), (q^3, y^3), (q^4, y^4)\}$ where the features of i -th sample x^i in D are transformed into q^i so that the feature $q_1^i = \alpha_1 x_1^i + \beta_1$ and the feature $q_2^i = \alpha_2 x_2^i + \beta_2$ for suitable real numbers $\alpha_1, \beta_1, \alpha_2, \beta_2$, such that the linear separator of the new dataset D' is of the form $\langle v, 0 \rangle$ (that is for any point $x \in \mathbb{R}^2$ on the linear separator, $\langle v, x \rangle = 0$). Give details of the resultant new dataset D' and the corresponding new linear separator $\langle v, 0 \rangle$. How do v and w relate to each other?

In this part we have to transform the point such that the linear separator will pass through origin.

for some $x \in \mathbb{R}^2$, and $x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ the transformation
is applied such that

$$\begin{aligned} q_1 &= \alpha_1 x_1 + \beta_1 \\ q_2 &= \alpha_2 x_2 + \beta_2 \end{aligned} \quad \left. \right\} \quad \text{--- (1)}$$

let $q = \begin{pmatrix} q_1 \\ q_2 \end{pmatrix}$

Now, as per question we have to make sure that the linear separator for transformed dataset is of the form $\langle v, q \rangle = 0$

$$\Rightarrow v_1 q_1 + v_2 q_2 = 0 \quad \text{--- (2)}$$

we know from part (a) that separator
is $x_1 + x_2 - 2 = 0$

i.e. of the form $\omega_1 x_1 + \omega_2 x_2 - b = 0$

replacing the values of x_1 and x_2 based upon
eq ①.

$$\Rightarrow \omega_1 \left(\frac{q_1 - \beta_1}{\alpha_1} \right) + \omega_2 \left(\frac{q_2 - \beta_2}{\alpha_2} \right) = b$$

Simplifying,

$$\frac{\omega_1}{\alpha_1} q_1 + \frac{\omega_2}{\alpha_2} q_2 - \left(\frac{\omega_1 \beta_1}{\alpha_1} + \frac{\omega_2 \beta_2}{\alpha_2} \right) = b \quad \text{--- (3)}$$

In new separator we don't need bias term

$$\text{so, } \frac{\omega_1 \beta_1}{\alpha_1} + \frac{\omega_2 \beta_2}{\alpha_2} + b = 0$$

$$\frac{\beta_1}{\alpha_1} + \frac{\beta_2}{\alpha_2} = -2 \quad \text{--- (4)}$$

[take values of ω_1, ω_2 and b
from part (a)
i.e. $\omega_1 = \omega_2 = 1$ and $b = 2$]

eq (3) becomes

$$\frac{\omega_1}{\alpha_1} q_1 + \frac{\omega_2}{\alpha_2} q_2 = 0 \rightarrow \frac{1}{2} q_1 + \frac{1}{2} q_2 = 0$$

comparing with $v_1 q_1 + v_2 q_2 = 0$

$$v_1 = \frac{\omega_1}{\alpha_1} \quad \text{and} \quad -v_2 = \frac{\omega_2}{\alpha_2}$$

Now for new dataset D'

let $\alpha_1 = -\beta_1$ and $\alpha_2 = -\beta_2$ satisfying eq (4)
applying transformation to each point in D

x_1	q_1	x_2	q_2	y
0	$0 + \beta_1 = \beta_1$	1	$-\beta_2 + \beta_2 = 0$	-1
1	$-\beta_1 + \beta_1 = 0$	0	$0 + \beta_2 = \beta_2$	-1
0	$0 + \beta_1 = \beta_1$	3	$-3\beta_2 + \beta_2 = -2\beta_2$	1
3	$-3\beta_1 + \beta_1 = -2\beta_1$	0	$0 + \beta_2 = \beta_2$	1

Verifying linear separator

$$v_1 q_1 + v_2 q_2 = 0$$

$$\frac{1}{\alpha_1} q_1 + \frac{1}{\alpha_2} q_2 = 0$$

$$\text{Sign}\left(\frac{\beta_1}{-\beta_1} + \frac{0}{-\beta_2}\right) = -1 = \hat{y}'$$

$$\text{Sign}\left(0 + \frac{\beta_2}{-\beta_2}\right) = -1$$

$$\text{Sign}\left(\frac{\beta_1}{-\beta_1} + \frac{(-2\beta_2)}{-\beta_2}\right) = 1$$

$$\text{Sign}\left(\frac{(-2\beta_1)}{-\beta_1} + \frac{\beta_2}{-\beta_2}\right) = 1$$

let take $\alpha_1 = \alpha_2 = 1$ and $\beta_1 = \beta_2 = -1$

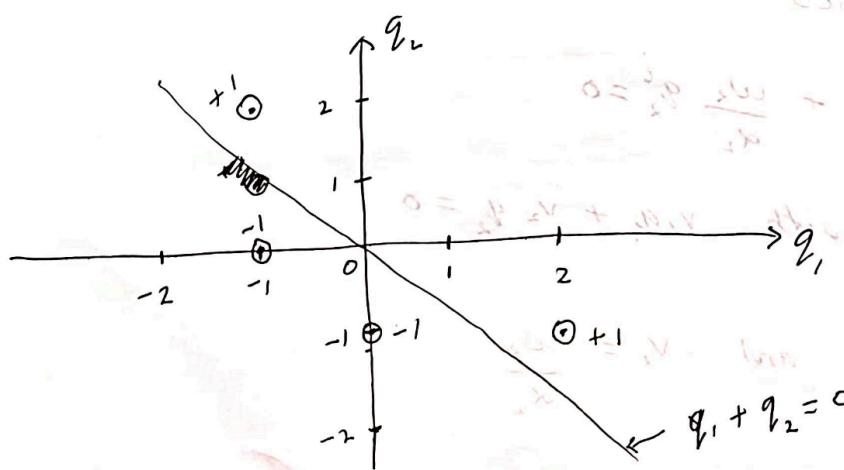
applying transformation to each point in D

x_1	q_1	x_2	q_2	y
0	$0 - 1 = -1$	1	$1 - 1 = 0$	-1
1	$1 - 1 = 0$	0	$0 - 1 = -1$	-1
0	$0 - 1 = -1$	3	$3 - 1 = 2$	1
3	$3 - 1 = 2$	0	$0 - 1 = -1$	1

$$\therefore D' = \{((-1, 0), -1), ((0, -1), -1), ((-1, 2), 1), ((2, -1), 1)\}$$

$$v_1 = 1 \text{ and } v_2 = 1$$

\therefore linear separator $\rightarrow q_1 + q_2 = 0$



1. Consider the following data set: $D = \{(x^1, y^1), (x^2, y^2), (x^3, y^3), (x^4, y^4)\}$, where $(x^1, y^1) = ((0, 1), -1)$, $(x^2, y^2) = ((1, 0), -1)$, $(x^3, y^3) = ((0, 3), +1)$ and $(x^4, y^4) = ((3, 0), +1)$.
- [3 marks] Find a linear separator for dataset D . Clearly write the expression for the linear separator and justify why your construction is in fact a linear separator for D .
 - [6 marks] Suppose the linear separator you found in part 1a is represented as (w, b) such that for any $x \in \mathbb{R}^2$ on the linear separator, $\langle w, x \rangle - b = 0$ holds. Transform D into a new dataset $D' = \{(q^1, y^1), (q^2, y^2), (q^3, y^3), (q^4, y^4)\}$ where the features of i -th sample x^i in D are transformed into q^i so that the feature $q_1^i = \alpha_1 x_1^i + \beta_1$ and the feature $q_2^i = \alpha_2 x_2^i + \beta_2$ for suitable real numbers $\alpha_1, \beta_1, \alpha_2, \beta_2$, such that the linear separator of the new dataset D' is of the form $\langle v, 0 \rangle$ (that is for any point $x \in \mathbb{R}^2$ on the linear separator, $\langle v, x \rangle = 0$). Give details of the resultant new dataset D' and the corresponding new linear separator $\langle v, 0 \rangle$. How do v and w relate to each other?
 - [4 marks] Consider the dataset $A = \{(x^t, y^t)\}_{t=1,2,3,\dots}$ where $x^t \in \mathbb{R}^2$, and $y^t \in \{+1, -1\}$. Recall that this dataset A is linearly separable if there exist a non-zero $w^* \in \mathbb{R}^2$ and $\gamma > 0$ such that $\langle w^*, x^t \rangle > \gamma$ when $y^t = 1$ and $\langle w^*, x^t \rangle < -\gamma$ when $y^t = -1$. Justify if this definition directly fits for the dataset D . Also justify if this definition fits directly for the dataset D' . Using your analysis in part 1b, explain why such a general definition of linear separability for A is sufficient.

$$\left. \begin{array}{l} \langle w^*, x^t \rangle > \gamma \text{ when } y^t = 1 \\ \langle w^*, x^t \rangle < -\gamma \text{ when } y^t = -1 \end{array} \right\} \quad \text{--- (1)}$$

where $w^* \in \mathbb{R}^2$, $x^t \in \mathbb{R}^2$ and $\gamma > 0$

Checking condition (1) for dataset D

x_1	x_2	$\langle w^*, x^t \rangle$	y	
0	1	w_2^*	-1	$\rightarrow w_2^* < -\gamma \quad \text{--- (2)}$
1	0	w_1^*	-1	$\rightarrow w_1^* < -\gamma \quad \text{--- (3)}$
0	3	$3w_2^*$	+1	$\rightarrow 3w_2^* > \gamma \quad \text{--- (4)}$
3	0	$3w_1^*$	+1	$\rightarrow 3w_1^* > \gamma \quad \text{--- (5)}$

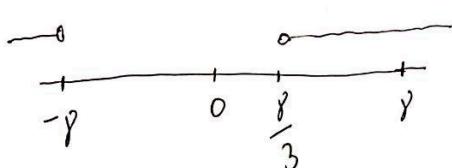
from eq (1) and (3)

$$+w_2^* < -\gamma$$

$$3w_2^* > \gamma \rightarrow w_2^* > \frac{\gamma}{3}$$

} contradiction

Similarly from eq (2) and (4)
it can be shown the these
inequalities contradicts.



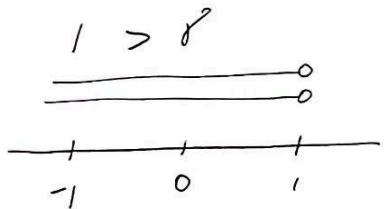
Therefore condition ① does not fits on dataset D
 it is due to the fact that, for separating dataset
 D we need a bias term in the linear separator.

Checking condition ① for dataset D'

q_1	q_2	y	$\langle v^*, q^t \rangle$
β_1	0	-1	$v_1 \beta_1 = \frac{\beta_1}{\alpha_1} = -1 \rightarrow -\gamma \text{ } \textcircled{1}$
0	β_2	-1	$v_2 \beta_2 = \frac{\beta_2}{\alpha_2} = -1 \rightarrow -\gamma \text{ } \textcircled{2}$
β_1	$-2\beta_2$	+1	$v_1 \beta_1 - 2v_2 \beta_2 = \frac{\beta_1}{\alpha_1} - 2 \frac{\beta_2}{\alpha_2} = 1 \rightarrow >\gamma \text{ } \textcircled{3}$
$-2\beta_1$	β_2	+1	$-2v_1 \beta_1 + v_2 \beta_2 = -2 \frac{\beta_1}{\alpha_1} + \frac{\beta_2}{\alpha_2} = 1 \rightarrow >\gamma \text{ } \textcircled{4}$

from eq ① and ③

$$-1 < -\gamma \rightarrow 1 > \gamma$$



from eq ② and ④

$$-1 < -\gamma \rightarrow 1 > \gamma$$

$$1 > \gamma$$

Same as previous

It can be seen that condition ① fits directly to dataset D', it is due to the fact that linear separator $\langle v^*, q^t \rangle = 0$ passes from origin therefore not need of bias term.

But dataset D' is formed by transforming the dataset D, maintaining the relative position of points,

So any dataset can be transform into the dataset like D' , on which the linear separability condition ① i.e. $y^t \langle w^*, x^t \rangle > \gamma$ can fit.

Therefore such general definition of linear separability is sufficient.

2. [5 marks] Consider a multi-layer perceptron (MLP) with L layers and the layers are numbered as $\ell = 1, 2, \dots, L$. Let W^ℓ denote the weights connecting to the ℓ -th layer of the MLP where $\ell \geq 2$. Let e denote the error with respect to a sample (x, y) . Derive a suitable expression to compute the error gradients $\nabla_{W^\ell} e$ and show that it is of the form $\nabla_{W^\ell} e = \text{Diag}(\phi')^\ell \delta^\ell (a^{\ell-1})^\top$ where the j -th diagonal element of $\text{Diag}(\phi')^\ell$ is $\phi'(z_j^\ell)$, δ^ℓ denotes the error gradients with respect to the activations at ℓ -th layer and $a^{\ell-1}$ denotes the activations at layer $\ell - 1$. (Note: You can reuse the expressions for δ^ℓ , derived in class.)

2nd Answer:

JANUARY Backward Prop Derivation 2022

$x^i \rightarrow \begin{array}{|c|c|c|c|} \hline & 1 & 2 & \cdots & L \\ \hline \end{array} \rightarrow \hat{y}^i$

$$e^i = e(x^i, y^i, \hat{y}^i, \omega)$$

$$E(\omega) = \sum_{i=1}^I e^i(\omega)$$

$$\nabla E(\omega) = \sum_{i=1}^I \nabla e^i(\omega)$$

$$I \in \{1, 2, \dots, I\}$$

$$\nabla_\omega e(\omega) = \nabla_\omega e(x, y, \hat{y}, \omega) = \text{Error Gradient of Sample } (x, y)$$

$$= \nabla_\omega (y - \hat{y})^2$$

$$= \nabla_\omega (y - \phi(w^T \phi(w^T x)))^2$$

Forward Pass

08 Saturday

$$z_i^l = \sum_{j=1}^{N_{l-1}} w_{ij}^l a_j^{l-1}$$

$$a_i^l = \phi(z_i^l)$$

$$\frac{\partial a_i^l}{\partial z_i^l} = \phi'(z_i^l) \quad (1)$$

09 Sunday

$$\frac{\partial e}{\partial w_{ij}^l} = \frac{\partial e}{\partial z_i^l} \cdot \frac{\partial z_i^l}{\partial w_{ij}^l} = \frac{\partial e}{\partial a_i^l} \cdot \frac{\partial a_i^l}{\partial z_i^l} \cdot \frac{\partial z_i^l}{\partial w_{ij}^l}$$

JANUARY

10 Monday

$$\frac{\partial e}{\partial w_{ij}} = \boxed{\frac{\partial e}{\partial a_i^l} \cdot \phi'(z_i^{l+1}) \cdot a_j^{l+1}} \quad \textcircled{2}$$

$$\frac{\partial e}{\partial a_i^l} = \frac{\partial e}{\partial a_i^l} \cdot \frac{\partial z_1^{l+1}}{\partial a_i^l} + \frac{\partial e}{\partial z_2^{l+1}} \cdot \frac{\partial z_2^{l+1}}{\partial a_i^l} + \\ \frac{\partial e}{\partial z_3^{l+1}} \cdot \frac{\partial z_3^{l+1}}{\partial a_i^l} + \dots + \frac{\partial e}{\partial z_{N^{l+1}}} \cdot \frac{\partial z_{N^{l+1}}}{\partial a_i^l}$$

$$= \sum_{m=1}^{N^{l+1}} \frac{\partial e}{\partial z_m^{l+1}} \cdot \frac{\partial z_m^{l+1}}{\partial a_i^l}$$

$$= \sum_{m=1}^{N^{l+1}} \frac{\partial e}{\partial z_m^{l+1}} \cdot w_{mj}^{l+1}$$

$$= \sum_{m=1}^{N^{l+1}} \frac{\partial e}{\partial a_m^{l+1}} \cdot \frac{\partial a_m^{l+1}}{\partial z_m^{l+1}} \cdot w_{mj}^{l+1}$$

$$= \sum_{m=1}^{N^{l+1}} \frac{\partial e}{\partial a_m^{l+1}} \cdot \phi'(z_m^{l+1}) \cdot w_{mj}^{l+1}$$

$$z_m^{l+1} = \sum_{m=1}^{N^l} w_{mj}^{l+1} a_j^l$$

11 Tuesday

$$\frac{\partial z_m^{l+1}}{\partial w_{mj}^{l+1}} = a_j^l$$

$$\frac{\partial z_m^{l+1}}{\partial a_j^l} = w_{mj}^{l+1}$$

$$\frac{\partial e}{\partial a_i^l} = \left[w_{1i}^{l+1} \phi'(z_1^{l+1}) - w_{N^{l+1}i}^{l+1} \phi'(z_{N^{l+1}}^{l+1}) \right] \begin{bmatrix} \frac{\partial e}{\partial a_1^{l+1}} \\ \vdots \\ \frac{\partial e}{\partial a_{N^{l+1}}^{l+1}} \end{bmatrix}$$

$(1 \times N^{l+1}) \quad (N^{l+1} \times 1)$

JANUARY 2022						
M	T	W	T	F	S	S
31					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30

JANUARY
12 Wednesday

2022

$$\begin{bmatrix} \frac{\partial e}{\partial w_1^l} \\ \vdots \\ \frac{\partial e}{\partial w_{N^l}^l} \end{bmatrix}_{(N^l \times 1)} = \begin{bmatrix} w_{1,1}^{l+1} \phi'(z_1^{l+1}) & w_{1,N^l}^{l+1} \phi'(z_{N^l}^{l+1}) \\ \vdots & \vdots \\ w_{1,1}^{l+1} \phi'(z_1^{l+1}) & w_{N^l,N^l}^{l+1} \phi'(z_{N^l}^{l+1}) \end{bmatrix}_{(N^l \times N^{l+1})} \begin{bmatrix} \frac{\partial e}{\partial a_1^{l+1}} \\ \vdots \\ \frac{\partial e}{\partial a_{N^{l+1}}^{l+1}} \end{bmatrix}_{(N^{l+1} \times 1)}$$

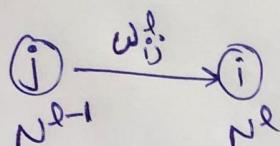
$$= \begin{bmatrix} w_{1,1}^{l+1} & w_{1,2}^{l+1} & \cdots & w_{1,N^l}^{l+1} \\ w_{1,2}^{l+1} & w_{2,2}^{l+1} & \cdots & w_{2,N^l}^{l+1} \\ \vdots & \vdots & \ddots & \vdots \\ w_{1,N^l}^{l+1} & w_{2,N^l}^{l+1} & \cdots & w_{N^l,N^l}^{l+1} \end{bmatrix}_{(N^l \times N^{l+1})} \begin{bmatrix} \phi'(z_1^{l+1}) \\ \phi'(z_2^{l+1}) \\ \vdots \\ \phi'(z_{N^l}^{l+1}) \end{bmatrix}_{(N^{l+1} \times 1)}$$

$$\boxed{\delta^l = (w^{l+1})^T \text{Diag}((\phi'(z))^{l+1}) \sigma^{l+1}} - \textcircled{2}$$

$$\begin{array}{ccc} \textcircled{j} \xrightarrow[N^{l-1}]{w_{ij}^l} \textcircled{i} & \Rightarrow [w^l] \text{ have shape } & (N^l \times N^{l-1}) \\ & \Rightarrow \left[\frac{\partial e}{\partial w^l} \right] \text{ have shape } & (N^l \times N^{l-1}) \end{array}$$

$$\begin{aligned} i &\in \{1, 2, \dots, N^l\}, & N^l &= \text{Number of Neuron in } l^{\text{th}} \text{ layer.} \\ j &\in \{1, 2, \dots, N^{l-1}\}, & N^{l-1} &= \text{Number of Neuron in } (l-1)^{\text{th}} \text{ layer.} \end{aligned}$$

$$\frac{\partial e}{\partial w_{ij}^l} = \frac{\partial e}{\partial a_i^l} \cdot \phi'(z_i^l) \cdot a_j^{l-1}$$

 $\Rightarrow [w^l]$ have shape $(N^l \times N^{l-1})$

$$\Rightarrow \left[\frac{\partial e}{\partial w_{ij}^l} \right] \text{ have shape } (N^l \times N^{l-1})$$

$$i \in \{1, 2, \dots, N^{l-1}\}$$

$$j \in \{1, 2, \dots, N^l\}$$

$$\left[\frac{\partial e}{\partial w_{1j}^l} \quad \frac{\partial e}{\partial w_{2j}^l} \quad \frac{\partial e}{\partial w_{3j}^l} \quad \dots \quad \frac{\partial e}{\partial w_{N^l j}^l} \right] = \begin{array}{|c|c|} \hline \frac{\partial e}{\partial a_1^l} & \frac{\partial e}{\partial a_2^l} \\ \hline \end{array} \dots$$

$$= \left[\frac{\partial e}{\partial a_1^l} \phi'(z_1^l) \quad \frac{\partial e}{\partial a_2^l} \phi'(z_2^l) \quad \dots \quad \frac{\partial e}{\partial a_{N^l}^l} \phi'(z_{N^l}^l) \right] a_j^{l-1}$$

$$\begin{bmatrix} \frac{\partial e}{\partial w_{11}} & \frac{\partial e}{\partial w_{12}} & \cdots & \frac{\partial e}{\partial w_{1N^l}} \\ \frac{\partial e}{\partial w_{21}} & \frac{\partial e}{\partial w_{22}} & \cdots & \frac{\partial e}{\partial w_{2N^l}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial e}{\partial w_{N^l 1}} & \frac{\partial e}{\partial w_{N^l 2}} & \cdots & \frac{\partial e}{\partial w_{N^l N^l}} \end{bmatrix}_{(N^{l-1} \times N^l)} = \begin{bmatrix} a_1^{l-1} \\ a_2^{l-1} \\ \vdots \\ a_{N^l}^{l-1} \end{bmatrix}_{N^{l-1} \times 1} \begin{bmatrix} \frac{\partial e}{\partial a_1^l} \phi'(z_1^l) & \frac{\partial e}{\partial a_2^l} \phi'(z_2^l) & \cdots & \frac{\partial e}{\partial a_{N^l}^l} \phi'(z_{N^l}^l) \end{bmatrix}_{1 \times N^l}$$

$$(\nabla_{w^l})^T e \quad \cancel{\text{Diagram}} \quad = \begin{bmatrix} a_1^{l-1} \\ a_2^{l-1} \\ \vdots \\ a_{N^l}^{l-1} \end{bmatrix}_{N^{l-1} \times 1} \begin{bmatrix} \frac{\partial e}{\partial a_1^l} & \frac{\partial e}{\partial a_2^l} & \cdots & \frac{\partial e}{\partial a_{N^l}^l} \end{bmatrix}_{1 \times N^l} \begin{bmatrix} \phi'(z_1^l) & \cdots & 0 & \cdots & 0 \\ 0 & \phi'(z_2^l) & \cdots & 0 \\ \vdots & 0 & \cdots & 0 & \cdots & \phi'(z_{N^l}^l) \end{bmatrix}_{N^l \times N^l}$$

$$\nabla_{w^l} e = \cancel{\text{Diagram}} = a^{l-1} (s^l)^T \text{Diag}(\phi'(z^l))$$

Take Transpose on Both Sides,

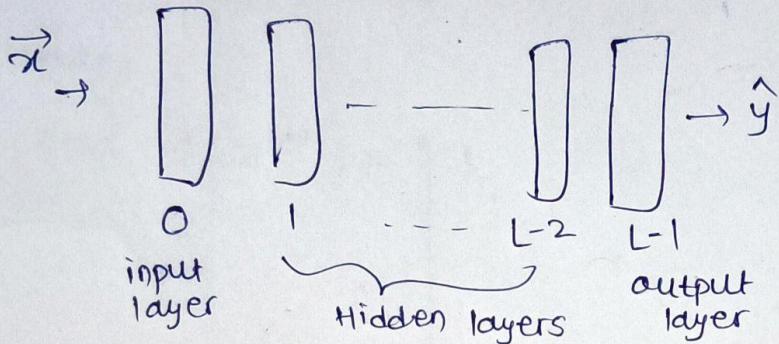
$$(\nabla_{w^l})^T e = \cancel{\text{Diagram}} = (a^{l-1} (s^l)^T \text{Diag}(\phi'(z^l)))^T = \text{Diag}(\phi'(z^l)) s^l (a^{l-1})^T$$

$$\boxed{\nabla_{w^l} e = \text{Diag}(\phi'(z^l)) s^l (a^{l-1})^T}$$

3. Consider a MLP (assume an input layer, an output layer and $L - 2$ hidden layers amounting to a total of L layers in the MLP) such that at the hidden layers and at the output layer, the linear activation function given by $\phi(z) = z$ is used. Consider that a sample (x, y) is used during training.

- [2 marks] Given the sample (x, y) , write the expression to denote a prediction \hat{y} from the MLP.
- [3 marks] Let e be the error for the sample (x, y) . Write the expression to compute the error gradients with respect to the activations at the ℓ -th layer. Compare and contrast this expression with that derived in class.
- [4 marks] Discuss based on the expression obtained in part 3b, when the network can experience vanishing and exploding gradient problems.

There are total L layers indexing from 0 to $L-1$.



$$a) \hat{y} = \phi^{L-1}(w^{L-1}\phi^{L-2}(w^{L-2}\phi^{L-3}(w^{L-3}\dots\phi^1(w^1\vec{x})\dots)))$$

since $\phi^1(z^1) = z^1$ for all hidden layers & output layer,

$$\hat{y} = w^{L-1} \cdot w^{L-2} \cdot w^{L-3} \dots w^1 \vec{x} \quad (\text{all the weight matrices are pre-multiplied to their next term})$$

$$b) \delta^L = v^{L+1} v^{L+2} \dots v^{L-1} \delta^{L-1}$$

where
 $L-1 = \cancel{\text{out}} \text{ index of output layer for given question.}$

$$\text{where } v^{L+1} = (w^{L+1})^T \text{ Diag}((\phi'(z))^{L+1})$$

$$\text{where, } \text{Diag}(\phi'(z)^{L+1}) = \begin{bmatrix} \phi'(z_1^{L+1}) \\ \phi'(z_2^{L+1}) \\ \phi'(z_3^{L+1}) \\ \vdots \\ \phi'(z_{N_{LH}}^{L+1}) \end{bmatrix}$$

Expression derived in the class

since $\phi(z^l) = z^l$ (linear activation)

$$\therefore \phi'(z^l) = 1$$

Hence,

$$\text{Diag}(\phi'(z)^{lH}) = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & & \ddots \\ & & & 1 \end{bmatrix}_{N_{lH} \times N_{lH}} \quad (\text{Identity matrix})$$

$$\therefore \text{Diag}(\phi'(z)^{lH}) = I_{N_{lH}} \quad (\text{identity matrix } N_{lH} \times N_{lH})$$

Hence,

$$V^{lH} = (W^{lH})^T \text{Diag}(\phi'(z)^{lH}) = (W^{lH})^T I_{N_{lH}}$$
$$\therefore V^{lH} = (W^{lH})^T$$

putting in the expression of δ^l .

$$\therefore \delta^l = V^{lH} V^{l+2} \dots V^{L-1} \delta^{L-1}$$

$$\delta^l = (W^{lH})^T (W^{l+2})^T \dots (W^{L-1})^T \delta^{L-1}$$

where $\delta^{L-1} = \frac{\partial e}{\partial a^{L-1}}$
error gradient wrt
activation of
output layer

c) error gradient wrt weight =

$$\begin{aligned}\frac{\partial e}{\partial w^l} &= \nabla_{w^l} e = \underbrace{\text{Diag}(\phi'(z^l))}_{\because \phi'(z^l)=1} \delta^l (\vec{a}^{l-1})^T \\ &= I_{N_l} \delta^l (\vec{a}^{l-1})^T\end{aligned}$$

$$\nabla_{w^l} e = \delta^l (\vec{a}^{l-1})^T$$

putting δ^l from part b)

$$\nabla_{w^l} e = (w^{l+1})^T (w^{l+2})^T \dots (w^{L-1})^T \delta^{L-1} (\vec{a}^{L-1})^T$$

Exploding gradient problem will occur \rightarrow

① when $(\vec{a}^{L-1})^T$ takes very high values.
since we are using linear activation ($\phi(z^l) = z^l$),

$$\vec{a}^{L-1} = \vec{z}^{L-1}$$

Hence, there is no limit for the values which \vec{a}^{L-1} can take. ② when the term $(w^{L+1})^T (w^{L+2})^T \dots (w^{L-1})^T$ takes very high values.

Vanishing gradient problem will occur \rightarrow

when the term $(w^{L+1})^T (w^{L+2})^T \dots (w^{L-1})^T$ takes very small values.

4. Consider a data set $D = \{(x^t, y^t)\}$, $t = 1, 2, 3, \dots$, where $x^t \in \mathbb{R}^d$ and $y^t \in \{+1, -1\}$, $\forall t = 1, 2, 3, \dots$. Note that D is linearly separable if there exists a hyperplane $H = (w^*, b^*)$ (with $w^* \neq 0$) and $\gamma > 0$ such that $y^t(\langle w^*, x^t \rangle - b^*) \geq \gamma$ for every $t = 1, 2, 3, \dots$
- [4 marks] Consider a point x^p in D with $y^p = +1$. Find a point u on H closest to x^p and derive an expression for $\text{dist}(x^p, u)$, the distance between x^p and u . Explain how you constructed u .
 - [3 marks] Consider a point x^q in D with $y^q = -1$. Find a point v on H closest to x^q and derive an expression for $\text{dist}(x^q, v)$, the distance between x^q and v . Explain how you constructed v .
 - [6 marks] Let \hat{x} denote the point with label $+1$ which is closest to H and let \tilde{x} denote the point with label -1 which is closest to H . Suppose that the distance of \hat{x} to the linear separator H is α and the distance of \tilde{x} to the linear separator H is also α . Justify with proper reasons if and when the relations $\alpha = \gamma$, $\alpha > \gamma$ and $\alpha < \gamma$ can hold. (Note: The distance of a point $x \in \mathbb{R}^d$ to H is given by the distance $\text{dist}(x, u)$ where u is a point on H and is closest to x .)

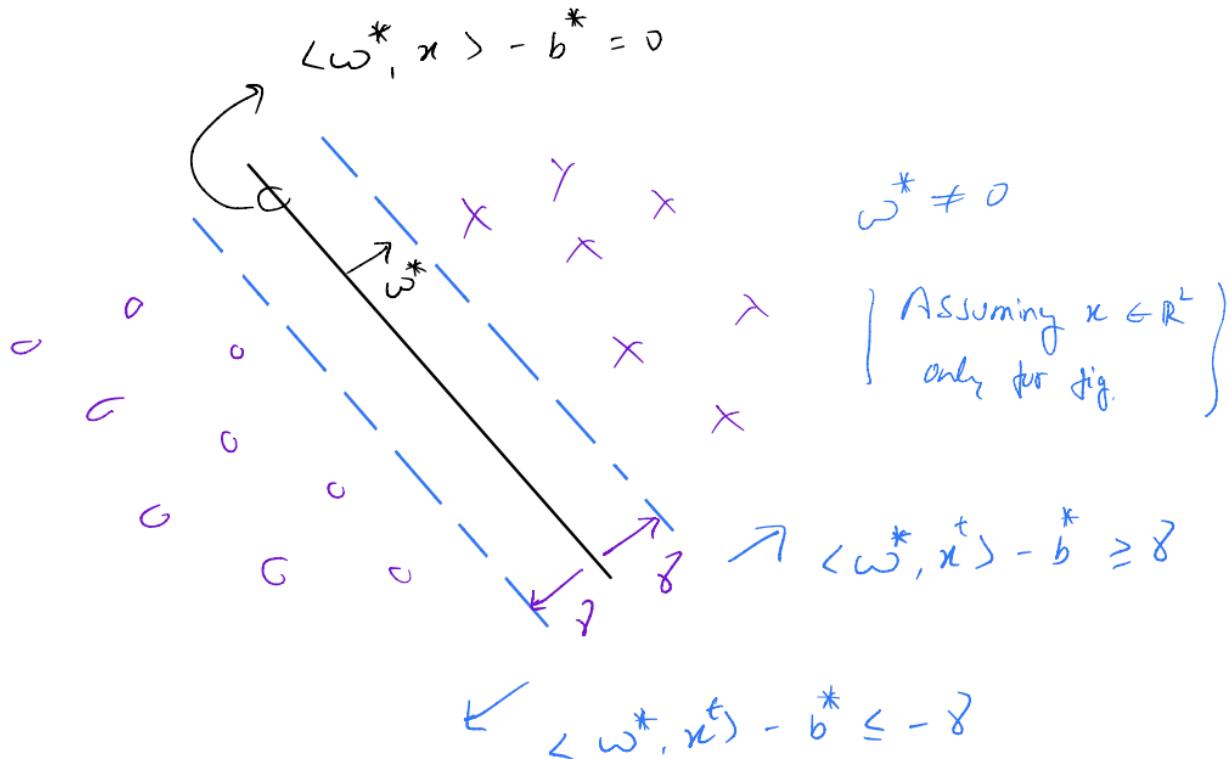
Q4/ Given:

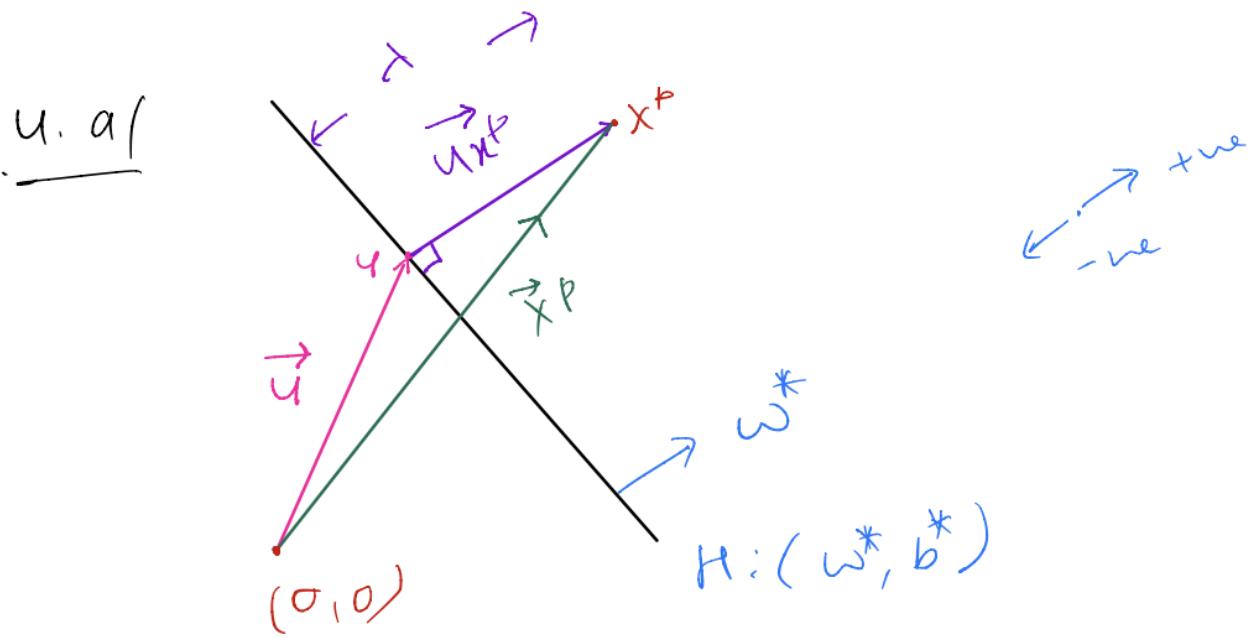
Dataset : $D = \{(x^t, y^t)\}$

feature vectors : $x^t \in \mathbb{R}^d$

class labels : $y^t \in \{+1, -1\}$

D is linearly separable if





- Assuming $(0,0)$ not lie on H
- \vec{u} : Closest point to \vec{x}^p on H
- \vec{u} is chosen as the point intersecting H when a \perp^r line is dropped on H from \vec{x}^p
- $\vec{u}x^p$ is vector joining point u & point x^p with magnitude λ (λ to find)

Since, u is in the hyperplane H ,

$$\langle \omega^*, u \rangle - b^* = 0 \quad - (1)$$

By triangle law of vector addition ($\vec{x}^p = \vec{u} + \vec{u}x^p$)

$$\vec{x}^P = \vec{u} + \lambda \frac{\vec{\omega}^*}{\|\vec{\omega}^*\|} \quad - \textcircled{2}$$

$$\text{eq-}\textcircled{2} \times (\vec{\omega}^*)^T$$

$$\vec{\omega}^{*T} \vec{x}^P = \vec{\omega}^{*T} \cdot \vec{u} + \lambda \frac{\vec{\omega}^{*T} \cdot \vec{\omega}^*}{\|\vec{\omega}^*\|}$$

$$\langle \vec{\omega}^*, \vec{x}^P \rangle = \langle \vec{\omega}^*, \vec{u} \rangle + \frac{\lambda \langle \vec{\omega}^*, \vec{\omega}^* \rangle}{\|\vec{\omega}^*\|}$$

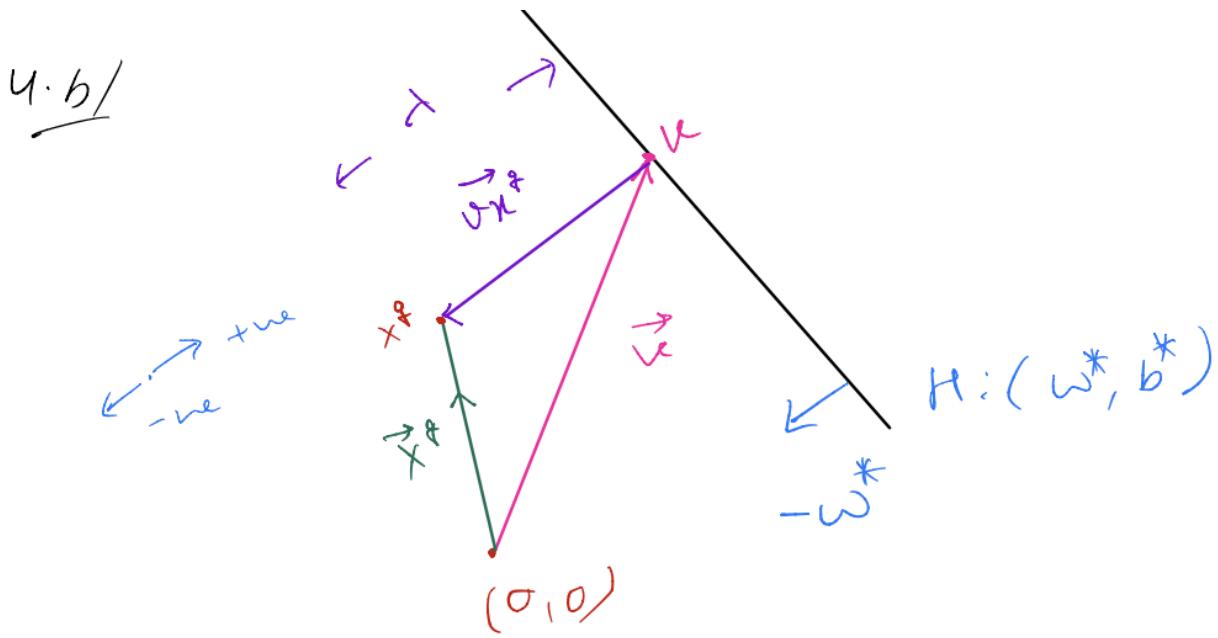
$$\langle \vec{\omega}^*, \vec{x}^P \rangle = b^* + \frac{\lambda \|\vec{\omega}^*\|^2}{\|\vec{\omega}^*\|} \quad (\text{from } \textcircled{1})$$

$$\lambda = \frac{\langle \vec{\omega}^*, \vec{x}^P \rangle - b^*}{\|\vec{\omega}^*\|}$$

$$\Rightarrow \lambda = \left| \frac{\langle \vec{\omega}^*, \vec{x}^P \rangle - b^*}{\|\vec{\omega}^*\|} \right| \quad \begin{array}{l} \text{Since } \lambda \text{ is distance} \\ \text{it cannot be negative.} \end{array}$$

Now substitute value of λ in eq- $\textcircled{2}$ we get \vec{u} as follows :

$$\vec{u} = \vec{x}^P - \left| \frac{\langle \vec{\omega}^*, \vec{x}^P \rangle - b^*}{\|\vec{\omega}^*\|} \right| * \frac{\vec{\omega}^*}{\|\vec{\omega}^*\|}$$



- Assuming $(0,0)$ not lie on H
- j : closest point to \vec{x}^* on H
- j is chosen as the point intersecting H when a \perp^r line is dropped on H from \vec{x}^*
- $\vec{v}x^*$ is vector joining point v & point x^* with magnitude λ (λ to find)

Since, j is in the hyperplane H ,

$$\langle w^*, j \rangle - b^* = 0 \quad \text{--- (1)}$$

By triangle law of vector addition ($\vec{x}^* = \vec{j} + \vec{v}x^*$)

$$\vec{x}^* = \vec{v} + \lambda * \left(-\frac{\vec{\omega}^*}{\|\vec{\omega}^*\|} \right) \quad - \textcircled{2}$$

$$\text{eq-}\textcircled{2} \times (\vec{\omega}^*)^T$$

$$\vec{\omega}^{*T} \vec{x}^* = \vec{\omega}^{*T} \cdot \vec{v} - \lambda * \left(\frac{\vec{\omega}^{*T} \cdot \vec{\omega}^*}{\|\vec{\omega}^*\|} \right)$$

$$\langle \vec{\omega}^*, \vec{x}^* \rangle = \langle \vec{\omega}^*, \vec{v} \rangle - \frac{\lambda \langle \vec{\omega}^*, \vec{\omega}^* \rangle}{\|\vec{\omega}^*\|}$$

$$\langle \vec{\omega}^*, \vec{x}^* \rangle = b^* - \frac{\lambda \|\vec{\omega}^*\|^2}{\|\vec{\omega}^*\|} \quad (\text{from } \textcircled{1})$$

$$\lambda = \frac{-\langle \vec{\omega}^*, \vec{x}^* \rangle + b^*}{\|\vec{\omega}^*\|}$$

$$\Rightarrow \lambda = \left| \frac{-\langle \vec{\omega}^*, \vec{x}^* \rangle + b^*}{\|\vec{\omega}^*\|} \right| \quad \begin{array}{l} \text{since } \lambda \text{ is distance} \\ \text{it cannot be negative.} \end{array}$$

Now substitute value of λ in eq- $\textcircled{2}$ we get \vec{v} as follows :

$$\vec{v} = \vec{x}^* + \left| \frac{-\langle \vec{\omega}^*, \vec{x}^* \rangle + b^*}{\|\vec{\omega}^*\|} \right| * \frac{\vec{\omega}^*}{\|\vec{\omega}^*\|}$$

4. c)

using eq⁻ⁿ of 2 from 4.a & 4.b

$$\alpha = \frac{|\langle \omega^*, \hat{x} \rangle - b|}{\|\omega^*\|} = \frac{|\langle \omega^*, \tilde{x} \rangle - b|}{\|\omega^*\|} \quad \text{--- (1)}$$

From given linear separability condition:

$$y^t (\langle \omega^*, x^t \rangle - b^*) \geq 8$$

If x^t has label $y^t = +1$,

$$\text{then, } (\langle \omega^*, x^t \rangle - b^*) \geq 8 \quad \text{--- (2)}$$

If x^t has label $y^t = -1$,

$$\text{then, } -(\langle \omega^*, x^t \rangle - b^*) \geq 8$$

$$\Rightarrow (\langle \omega^*, x^t \rangle - b^*) \leq -8 \quad \text{--- (3)}$$

from (2) & (3) we can write,

$$|\langle \omega^*, x^t \rangle - b^*| \geq 8 \quad \forall x^t \in D$$

$\therefore \hat{x} \& \tilde{x} \in D$

$$\text{thus } |\langle \omega^*, \hat{x} \rangle - b^*| \geq 8 \quad \text{--- (4)}$$

(At place of \hat{x} , \tilde{x} can also be used)

In eq-(4) Dividing $\|\omega^*\|$ on both side

$$\frac{|\langle \omega^*, \hat{x} \rangle - b^*|}{\|\omega^*\|} \geq \frac{\gamma}{\|\omega^*\|}$$
$$\underbrace{\left[\alpha \geq \frac{\gamma}{\|\omega^*\|} \right]}_{\text{From eq-(1)}} - \textcircled{5}$$

Now we will analyze the case for
 $\|\omega^*\|$ in 3 intervals $\|\omega^*\|=1$, $\|\omega^*\| < 1$ &
 $\|\omega^*\| > 1$

Case 1: $\|\omega^*\|=1$: From eq-(5)

$$\alpha \geq \frac{\gamma}{\|\omega^*\|} \Rightarrow \alpha \geq \gamma$$

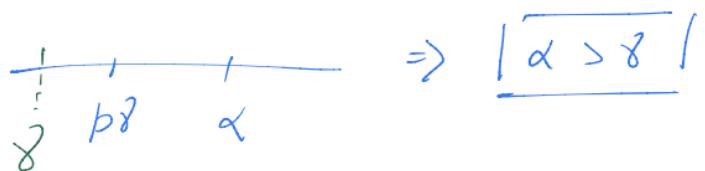
{when } $\alpha = \gamma$: This holds when the point lies exactly in the margin defined by γ , meaning the distance from the point to the hyperplane is exactly γ .

{when } $\alpha > \gamma$: This happens for closest points (\hat{x} & \tilde{x}) that lie farther from the hyperplane than the margin γ but still satisfy the linear separability condition.

Case 2: $0 < \|\omega^*\| < 1$

from eqⁿ ②: $\alpha \geq \frac{\gamma}{\|\omega^*\|} \Rightarrow$ let $\|\omega^*\| = 1/p$ ($p > 1$)

so, $\alpha \geq \frac{\gamma}{1/p} \Rightarrow \alpha \geq p\gamma$



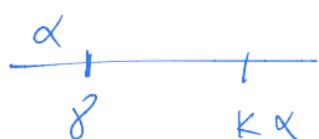
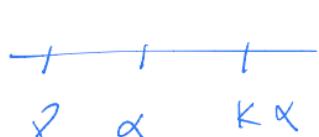
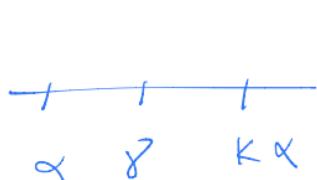
$\left\{ \begin{array}{l} \text{when } \|\omega^*\| > 1 \\ \alpha > \gamma \end{array} \right\}$: This happens for closest points ($\hat{x} \neq \tilde{x}$) that lie farther from the hyperplane than the margin γ but still satisfy the linear separability condition.

Case 3: $\|\omega^*\| > 1$

from eqⁿ ⑤: $\alpha \geq \frac{\gamma}{\|\omega^*\|} \rightarrow$ let $\|\omega^*\| = K$ ($K > 1$)

$$K\alpha \geq \gamma \Rightarrow \gamma \leq K\alpha$$

Now 3 cases can be possible



$$\alpha < \gamma$$

$$\alpha > \gamma$$

$$\alpha = \gamma$$

Here $\alpha < \gamma$, $\alpha > \gamma$ & $\alpha = \gamma$ all 3 are possible

$\left\{ \begin{array}{l} \text{when } \alpha = \gamma \\ \text{when } \alpha > \gamma \end{array} \right\}$: This holds when the point lies exactly on the margin defined by γ , meaning the distance from the point to the hyperplane is exactly γ .

$\left\{ \begin{array}{l} \text{when } \alpha < \gamma \\ \text{when } \alpha > \gamma \end{array} \right\}$: This happens for closest points (\tilde{x} & \tilde{y}): they lie farther from the hyperplane than the margin γ but still satisfy the linear separability condition.

$\left\{ \begin{array}{l} \text{when } \alpha < \gamma \\ \text{when } \alpha > \gamma \end{array} \right\}$: Happens when a point lies inside the margin, meaning it's closer to the hyperplane than the margin threshold. This is possible for points that are still correctly classified but do not respect the margin condition.

Conclusion:

- when $\|\omega^*\| = 1$, Possible case: $\alpha = \gamma$, $\alpha > \gamma$
- when $\|\omega^*\| > 1$, Possible case: $\alpha > \gamma$
- when $\|\omega^*\| < 1$, Possible case: $\alpha = \gamma$, $\alpha > \gamma$, $\alpha < \gamma$

