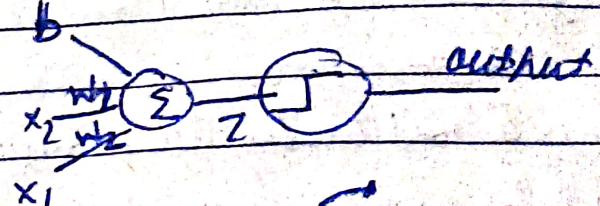


Q. XOR

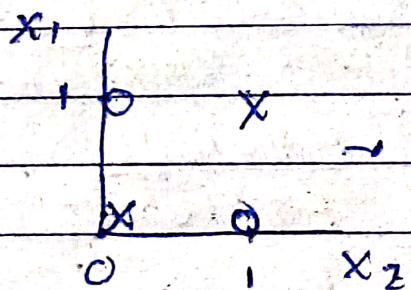
		$x_1$	$x_2$	$y$
0	1	1	0	1
1	0	1	0	0
1	1	0	0	1
0	0	0	0	0



$$z = w_1 x_1 + w_2 x_2 + b$$

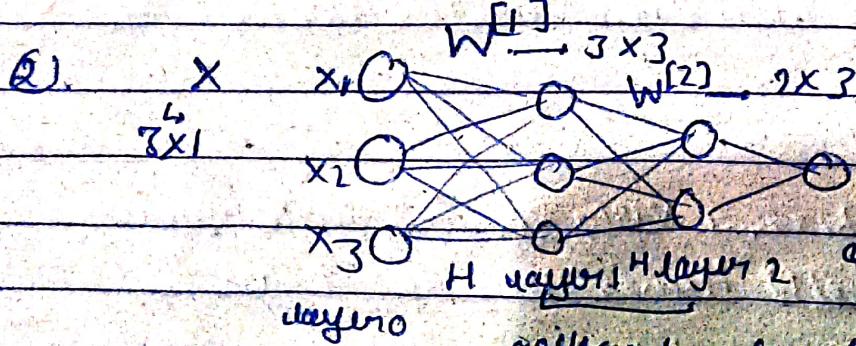
A perceptron can only classify data if it is linearly separable.

$$w_1 x_1 + w_2 x_2 + b = 0 \quad \text{acts as a decision boundary separating two classes}$$



→ cannot be separated by a line

In multilayer perceptrons to introduce non-linearity we apply activation function to the pre-activation w.r.t. sigmoid, tanh, ReLU.



activation function  $f(x) = \text{sigmoid}$

$$z^{[1]} = W^{[1]} x + b^{[1]}$$

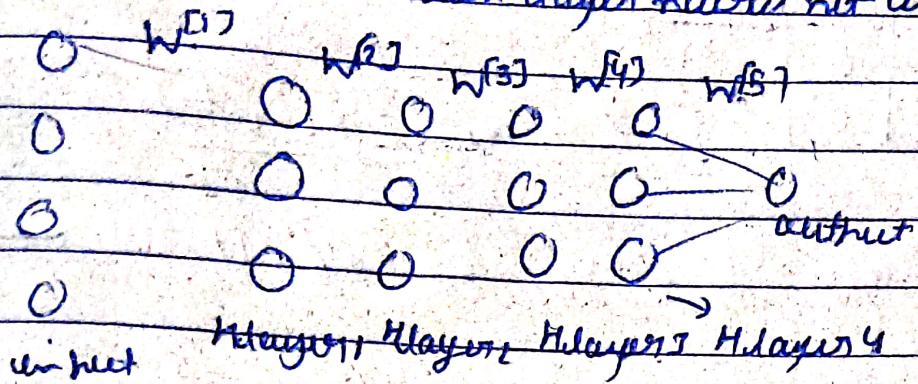
$$z^{[2]} = W^{[2]} (W^{[1]} x + b^{[1]}) + b^{[2]}$$

$$z^{[3]} = W^{[3]} (W^{[2]} x + b^{[2]}) + b^{[3]}$$

$$y = W x + b$$

\* Linear layer stacked on a linear layer is linear;  
stacking linear layers does not increase model capacity.

Consider a non-linear layer neural net work



$$\begin{aligned}\frac{\partial L}{\partial w^{1j}} &= \frac{\partial L}{\partial y} \frac{\partial a(z^1)}{\partial z^1} \frac{\partial z^1}{\partial z^0} \frac{\partial a(z^2)}{\partial z^2} \frac{\partial z^2}{\partial z^1} \frac{\partial a(z^3)}{\partial z^3} \frac{\partial z^3}{\partial z^2} \frac{\partial a(z^4)}{\partial z^4} \frac{\partial z^4}{\partial w^{1j}} \\ &= \frac{\partial L}{\partial y} \prod_{k=2}^4 \left( \frac{\partial a(z^k)}{\partial z^k} \right) \cdot \frac{\partial a(z^1)}{\partial z^1} \frac{\partial z^1}{\partial w^{1j}}\end{aligned}$$

~~weight activation derivatives~~

So, as the no of layers increases the terms that are being multiplication also increases.

causes of ~~the~~ shrinking gradients

~~•~~ - repeated multiplications across many layers

~~•~~ - derivative of activation fn's like (sigmoid, tanh)

& weights  $< 1$

Due to this gradients become very small & early layers stop learning. For ReLU derivatives as 1 when inputs are positive, unlike sigmoid func, solves the vanishing gradient problem.

it helps to capture

Q3

Positional encoding is necessary because tokens give inputs and output to encoder and decoder sequentially.

- \* Positional encoding is necessary because it helps to capture the position of words in the input and output sequence which is given to the encoder & decoder which is important because ~~for~~ self attention in permitting invariant for self attention
- "I am Mayank" = "Mayank am I"

\* Absolute positional encode      ~~and~~      Positional encoding

- \* uses position embedding matrix whose parameters are learned during training
- \* Fixed max length

\* Uses  $PE(pos, 2i) = \sin(\frac{pos}{10000^{2i/10}})$

$PE(pos, 2i+1) = \cos(\frac{pos}{10000^{2i/10}})$

- \* to make position embedding
- \* Encodes relative distance, similarity
- \* No learnable parameters
- \* max length can vary

\* ROPE (Rotary positional embeddings)

- \* It adds positional information by ~~by~~ rotating query & key vectors.

$$\text{Q rotated} = \text{ROPE}(Q)$$

P9105

PAGE NO.

DATE:

- Attention Scans depends on relative distance
- Can be generalised to sequence of unseen lengths
- rotations do not change val to s<sub>i</sub> magnitude

Q4

$$\text{Attention}(\alpha, k, v) = \text{softmax}\left(\frac{\alpha_i k^T}{\sqrt{d_k}}\right) v$$

$$\alpha = \alpha_i w_\alpha, k = x w_k, v = x w_v$$

### Intuition

$\alpha \rightarrow$  which token is looking at?

$k \rightarrow$  what the token offers or how it defines itself as

$v \rightarrow$  the actual information content that is passed on

Attention scores are scaled by  $\sqrt{d_k}$  to ensure variance remains constant as  $V_{dk} \propto d_k$ , so for high dimensions large variance leads to large values due to which softmax saturates and causes the problem of vanishing gradient.

Diagonal values are usually highest because diagonal entries correspond to  $\alpha_{ii} k_i^T$  in a word attending to itself, since they come from the same token & also has same position they have highest value.

Q5

### Multi-head attention

With one single attention head the model focuses on one type of relationship at a time. But with each multiple attention head each head discerns a different relationships - or Head 1 → Find subjects, Head 2 → ~~tracks~~ words

heads : Introducing dimensions of model

$h$  : no of attention heads

$d_{head}$  : Dimension per head

$$d_{model} = d_{head} h$$



Q6 In greedy decoding at each step we choose the token with most probability whereas in Beam Search at each step we keep top k sequences to maximize the global sequence probability. Exploring multiple future possibilities.

Ex at first step: 1

Word	Probability	Word	Probability
"I"	0.6	"I am"	0.3
"You"	0.4	"You are"	0.9

birthday

$$I \text{ am} = 0.6 \times 0.3 = 0.18$$

Beam Search

$$I \text{ am} = 0.6 \times 0.3 = 0.18$$

$$\text{You are} = 0.4 \times 0.9 = 0.36$$

Q1

$$d_{model} = 768$$

$$h = 12$$

$$d_{head} = d_{model} = \frac{768}{12} = 64$$

$$W_a, W_k, W_v \in R^{768 \times 768}$$

$$\text{Total parameters} = 3 \times 768 \times 768$$

$$= 1769472$$

$$\alpha = [2.0, 1.0, 0.0]$$

$$= \left[ \frac{c^2}{d^2 + l'^2 + l^2}, \frac{l'}{d^2 + l'^2 + l^2}, \frac{10^0}{d^2 + l'^2 + l^2} \right]$$

$$= [0.665, 0.243, 0.000]$$