

Sigmoid

ReLU

Tanh

Leaky ReLU

Softmax

Sigmoid

$$\text{derivative} \rightarrow f'(x) = f(x) \cdot (1 - f(x))$$

$$f(x) = \frac{1}{1 + \exp(-x)}$$

Tanh

$$f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

derivative

of  $\tanh(x)$

$$2(\text{sign}) - 1 = \tanh'$$

$$\text{ReLU} \rightarrow f(x) = \max(0, x) \Rightarrow \text{output range} = [0, \infty)$$

$$\text{Leaky ReLU} \rightarrow f(x) = \begin{cases} x, & x \geq 0 \\ \alpha x, & x \leq 0 \end{cases}$$

$$\frac{d}{dx} \left( \frac{e^x - e^{-x}}{e^x + e^{-x}} \right)$$

derivative of  $\tan(u)$ ?

$$\frac{uv' + u'v}{v^2}$$

$$\Rightarrow \frac{(e^x - e^{-x})(e^x + e^{-x}) - (e^x + e^{-x})(e^x - e^{-x})}{(e^x + e^{-x})^2}$$

$$\Rightarrow \frac{(e^x - e^{-x})^2 - (e^x + e^{-x})^2}{(e^x + e^{-x})^2}$$

$$\Rightarrow \frac{(e^x + e^{-x})(e^x + e^{-x}) - (e^x - e^{-x})(e^x - e^{-x})}{(e^x + e^{-x})^2}$$

$$\Rightarrow \frac{(e^x + e^{-x})^2 - (e^x - e^{-x})^2}{(e^x + e^{-x})^2}$$

$$\Rightarrow \frac{(e^x + e^{-x})^2 - (e^x - e^{-x})^2}{(e^x + e^{-x})^2}$$

$$\Rightarrow \frac{1 - \tan(-h)^2}{\cos^2(-h)} \cdot \text{with } \tan(-h) = \text{opp. / adj.}$$

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \xrightarrow{\text{softmax}} \begin{bmatrix} s_1 \\ s_2 \\ s_3 \end{bmatrix}$$

$$f(x) = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

Jacobian matrix =  $\begin{bmatrix} \frac{\partial s_1}{\partial x_1} & \frac{\partial s_1}{\partial x_2} & \frac{\partial s_1}{\partial x_3} \\ \frac{\partial s_2}{\partial x_1} & \frac{\partial s_2}{\partial x_2} & \frac{\partial s_2}{\partial x_3} \\ \frac{\partial s_3}{\partial x_1} & \frac{\partial s_3}{\partial x_2} & \frac{\partial s_3}{\partial x_3} \end{bmatrix} \rightarrow mxn \text{ matrix}$

Q)  $x = \begin{bmatrix} 2.0 \\ 1.0 \\ 0.1 \end{bmatrix}$  (a) Compute softmax o/p  
 (b) Compute the Jacobian

$$s_1 = \frac{e^2}{e^2 + e^1 + e^{0.1}} = \frac{7.38}{7.38 + 2.718 + 1.106} = 0.684$$

$$s_2 \Rightarrow \frac{e^1}{e^2 + e^1 + e^{0.1}} = \frac{2.718}{11.38} = 0.23$$

$$s_3 \Rightarrow \frac{e^{0.1}}{e^2 + e^1 + e^{0.1}} = \frac{1.106}{11.38} = 0.096$$

(1)

$$\begin{pmatrix} 0.68 \\ 0.23 \\ 0.096 \end{pmatrix} \quad \begin{pmatrix} 0.66 \\ 0.24 \\ 0.04 \end{pmatrix}$$

software uses sophisticated sigmoid.

$$\frac{\partial S_1}{\partial x_1} = S_1(1-S_1) \Rightarrow 0.66(0.24) \\ \Rightarrow 0.2244$$

$$S_{ij} = \frac{\partial S_i}{\partial x_j}$$

case 1 :  $i=j$  (diagonal element).

$$\frac{\partial S_1}{\partial x_2} = -S_1 S_2 \Rightarrow 0.66 \times 0.24 \\ \Rightarrow -0.1584$$

$$\frac{\partial S_i}{\partial x_i} = S_i(1-S_i)$$

case 2 = off diagonal element.

$$\frac{\partial S_1}{\partial x_3} = -S_1 S_3 \Rightarrow 0.66 \times 0.1 \\ \Rightarrow -0.06$$

$$\frac{\partial S_i}{\partial x_j} = -S_i S_j$$

$$\frac{\partial S_2}{\partial x_1} = -S_2 S_1 \Rightarrow -0.1584$$

$$\frac{\partial S_2}{\partial x_2} = S_2(1-S_2) \Rightarrow 0.1824$$

$$\frac{\partial S_2}{\partial x_3} \Rightarrow 0.24 \times 0.1 = -0.024$$

$$\frac{\partial S_3}{\partial x_1} = 0.1 \times 0.66 \Rightarrow -0.06$$

$$\frac{\partial S_3}{\partial x_2} \Rightarrow 0.24 \times 0.1 \Rightarrow -0.024 + 0.024$$

$$\frac{\partial S_3}{\partial x_3} \Rightarrow -0.0024 + \frac{1.0}{1.0+1.0+1.0}$$

$$\begin{bmatrix} 0.66 & -0.1584 & -0.06 \\ -0.1584 & 0.1824 & -0.024 \\ -0.06 & -0.024 & 0.0024 \end{bmatrix}$$

$$\begin{bmatrix} 0.2244 & -0.1584 & -0.06 \\ -0.1584 & 0.1824 & -0.024 \\ -0.06 & -0.024 & 0.09 \end{bmatrix}$$

$3 \times 3$

(Matrix) of weights for  $\phi_1$  input (3)

at step (3)

at step (3)

1) Datasets

2) DataLoaders -

Tensor is an n-dimensional array; we can plug in the datastructures.

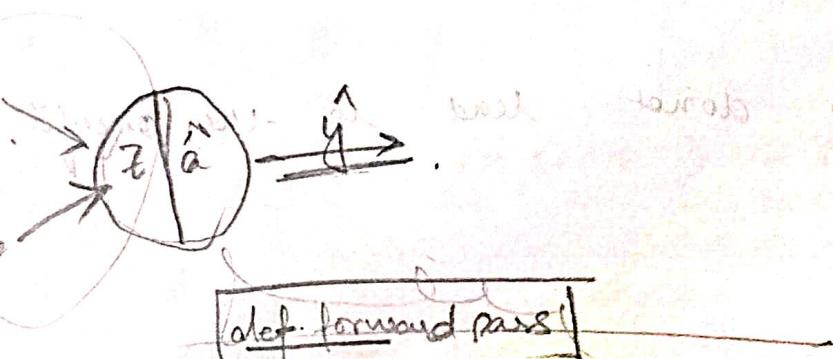
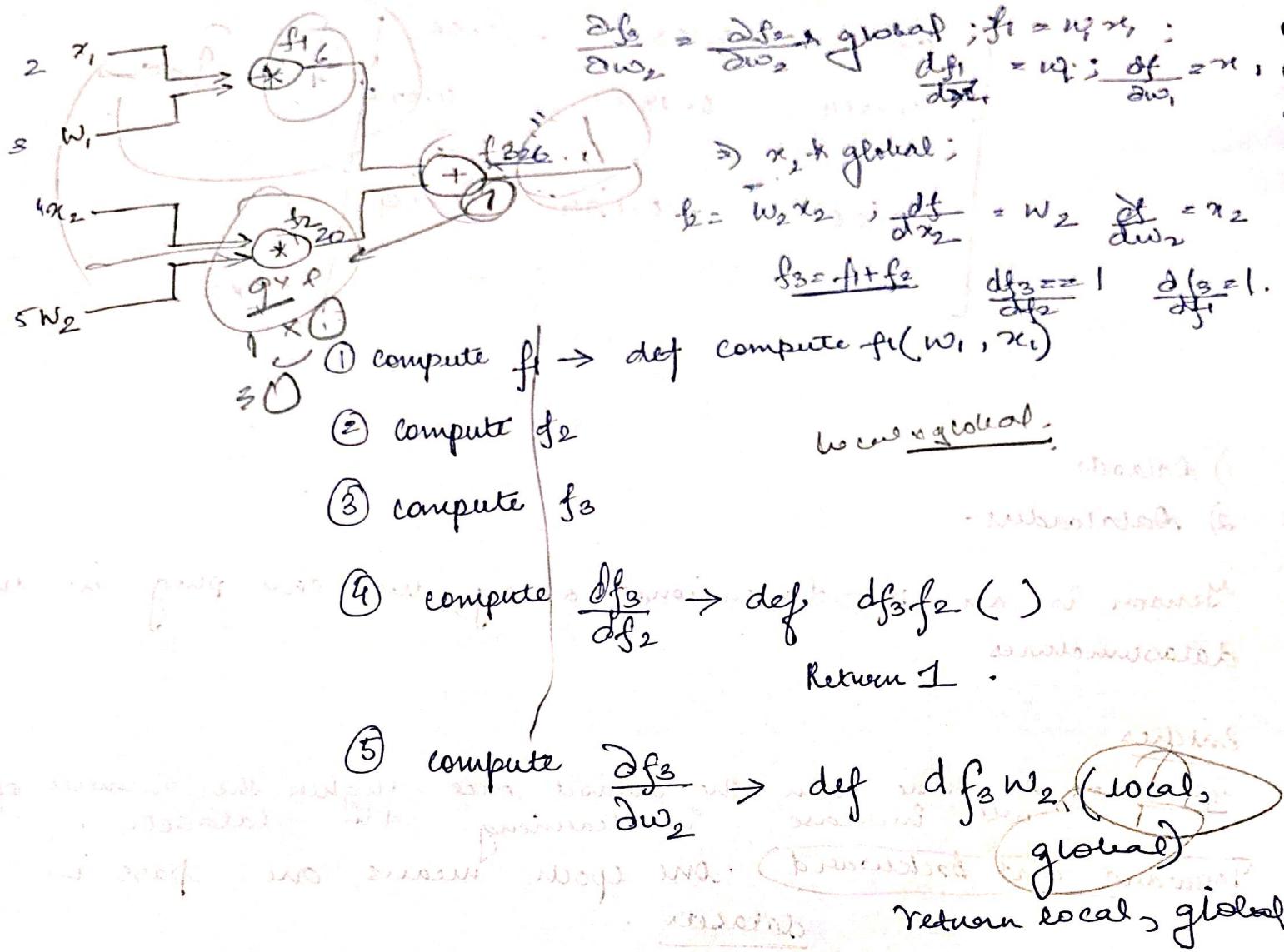
Batches

Epochs → 1 means seen the dataset once. Higher the number of epochs well. increase in learning the datasets.

Forward and backward → one epoch means one pass in the forward and backward dataset.

Batch size → hyperparameter.

- Running multiple times do not lead to the overfit.



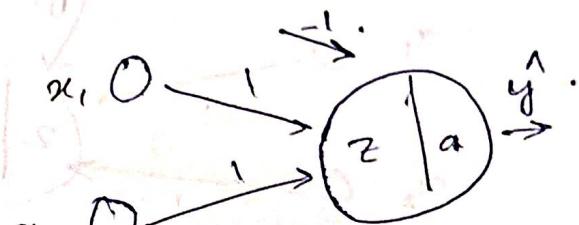
$$a^{(i)} = g(z^{(i)})$$

$$g = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z \leq 0 \end{cases}$$

Put this model into logistic regression and minimize the cost entropy loss.

AND

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



Case 1

$$\text{first case} = 0 \times 0 + (-1) \\ \Rightarrow -1$$

$$g = 0$$

$$\hat{y} = 0$$

This perceptron works perfectly as AND.

$$\text{Second} = 0 \times 1 - 1 \\ \Rightarrow -1$$

$$\Rightarrow \hat{y} = 0$$

$$\text{Third} = 0 - 1 \\ \hat{y} \Rightarrow 0$$

$$\text{Fourth} = \frac{1+1-1}{2} \\ \Rightarrow 0.5$$

OR

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

Case 1:-  $0+1=1 \Rightarrow 0+0=0$   
 $\Rightarrow 2.$   
 $\Rightarrow 0 \Rightarrow 0.$

Case 2:-  $0+1=0$   
 $\Rightarrow 1$

case 3:-

$$1+0+0$$

$$\Rightarrow 1.$$
  
$$y = 1.$$

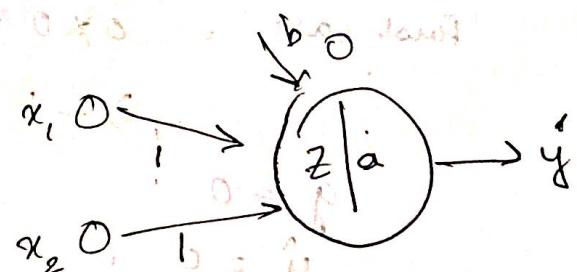
case 4:-

$$1+1+0$$

$$\Rightarrow 2$$
  
$$y = 1.$$

Now; it will work as OR function.

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

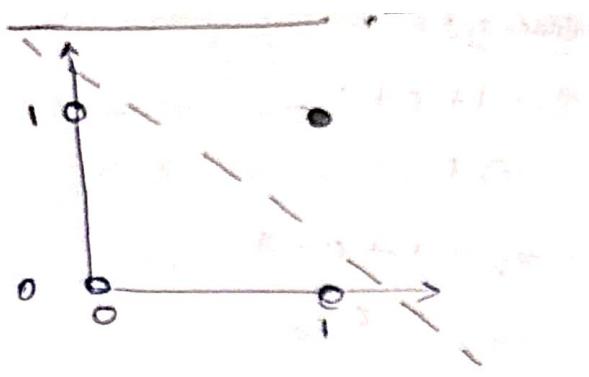


case 4:-  $1+1+0$ .

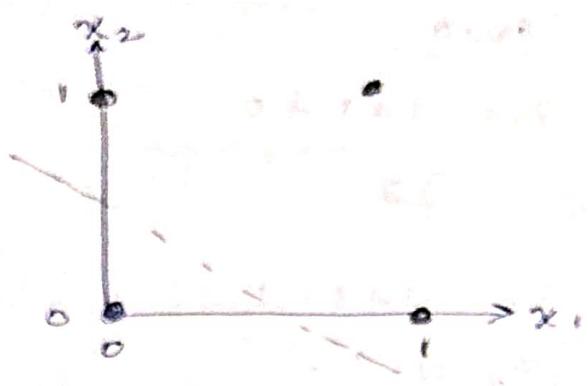
$$\Rightarrow 2-2$$

$$y = 1.$$

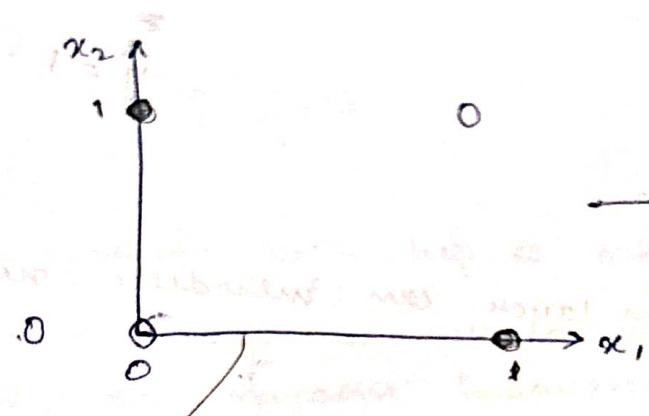
even if we change last term it's difficult to satisfy the four conditions to get 1.



a)  $x_1$  and  $x_2$



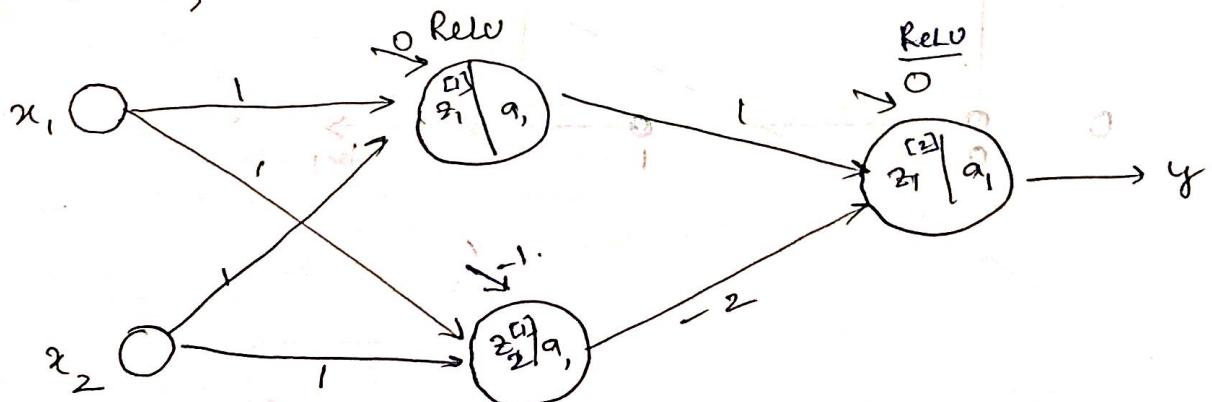
b)  $x_1$  OR  $x_2$



a) original space

c)  $x_1$  XOR  $x_2$

It is hard to draw boundary over here as the conditions are not satisfied and perceptron can't act as a XOR; so we will introduce a hidden layer.



$$\begin{aligned} z_1 &= 0+0 \xrightarrow{\text{ReLU}} 0+0 \\ \Rightarrow y &= 0 \end{aligned}$$

$$\begin{aligned} z_2 &= 0+0-1 \xrightarrow{\text{ReLU}} 0-1 \\ \Rightarrow y &= 2 \end{aligned}$$

$$\begin{aligned} &0+2+0 \\ \Rightarrow y &= 2 \end{aligned}$$

Case 4

$$z_1 = 1 + 1 + 0$$

$$\Rightarrow 2$$

$$z_2 = 1 + 1 - 1$$

$$z_2 \Rightarrow 1$$

$$\Rightarrow 1 \times 2 - 2 \times 1 + 0$$

$$\Rightarrow 2 - 2 + 0$$

$$\hat{y} = 0$$

Case 2 :-

$$z_1 = 1 + 0 + 0$$

$$\Rightarrow 1$$

$$\begin{aligned} z_2 &= 1 + 0 - 1 \\ &= 0 \end{aligned}$$

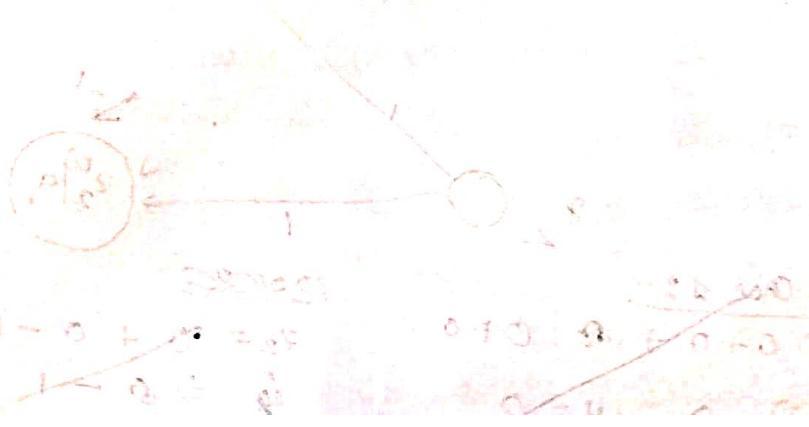
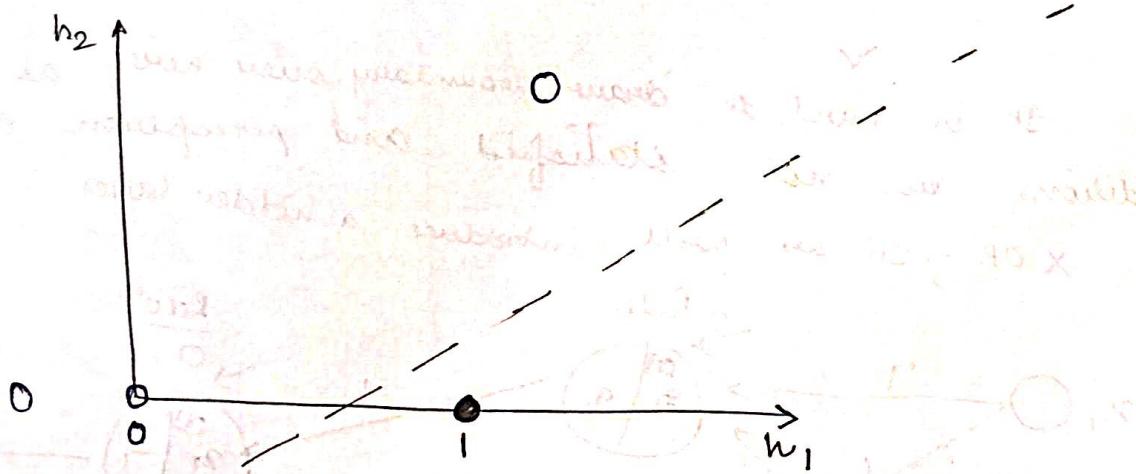
$$\Rightarrow 1 \times 1 - 2 \times 0 + 0$$

$$\Rightarrow 1$$

$$\hat{y} = 1$$

⇒ Introducing a addition layer can introduce non-linearity functions.

ii) The new (linearly separable) hspace



$$f = (x+y)z$$

$$\frac{df}{dx} = z$$

$$x \xrightarrow{-2} -1 \xrightarrow{+y} q \xrightarrow{\frac{df}{dq} = 2} -4$$

$$y \xrightarrow{-5} -4 \xrightarrow{+q} z$$

$$z \xrightarrow{-4} -4 \xrightarrow{\frac{df}{dz} = 3} f = -12$$

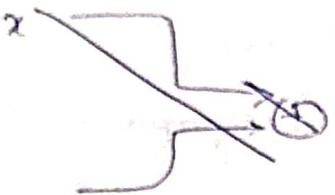
Black one  $\rightarrow$  last number that we get from backward pass

$$\text{Let } q = x+y \Rightarrow \frac{dq}{dx} = 1; \frac{dq}{dy} = 1.$$

$$f = \underline{qz} = \frac{df}{dq} \Rightarrow -z \quad \frac{df}{dz} = q$$

$\frac{df}{dx} \rightarrow$  how much  $x$  is influencing  $f \Rightarrow \underbrace{\frac{dq}{dx} \times \frac{df}{dq}}_{\text{local gradient}} = z \Rightarrow -4$

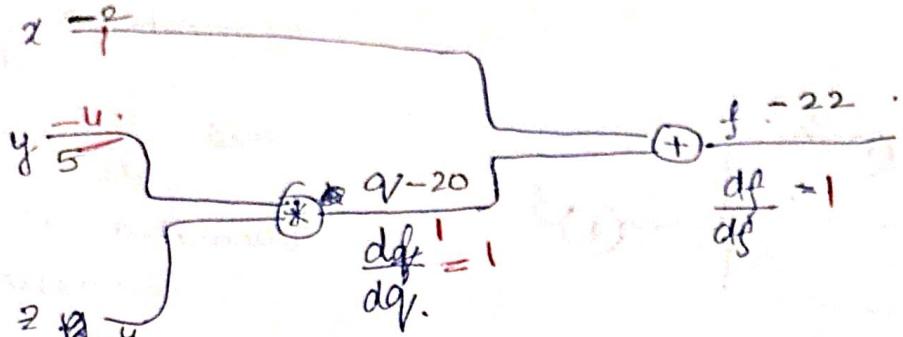
$\frac{df}{du} \rightarrow \frac{df}{dq} \times \frac{dq}{du} \Rightarrow -4 \times 1 \Rightarrow -4$   $\downarrow$  global gradient



$$Q = \underline{x + yz} .$$

$$\underline{q_1 = yz} .$$

b)



~~for~~ 5 → for a unit change in  $x$  there is an output change in  $f$  which is 5 times.

$$\frac{df}{dx} = 1 .$$

$$q = x + yz .$$

$$\frac{dq}{dx} = 1$$


---

$$\frac{dq}{dy} \Rightarrow z .$$

$$\frac{dq}{dz} = 2 \quad \frac{dq}{dy} = y \quad \frac{dq}{dx} = x$$

$$f = x + q .$$

$$\frac{df}{dx} = 1 \quad \frac{df}{dq} = 1 .$$


---

~~local gradient x global gradient~~

$$\frac{df}{dz} \Rightarrow \frac{dq}{dz} \times \frac{df}{dq} \quad \frac{df}{dy} \Rightarrow \frac{dq}{dy} \times \frac{df}{dq}$$

$$\Rightarrow q \times 1 \quad \text{since } q \text{ is constant} \Rightarrow 2 \times 1 \text{ if } q = 2$$

$$\Rightarrow \underline{\underline{5}} \quad \Rightarrow -4$$

$$\frac{df}{dz} = \frac{dq}{dx} \times \frac{df}{dq}$$

$$\Rightarrow 1 \times 1$$

Q)  $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$

$$w_0 = 2$$

$$x_0 = -1$$

$$w_1 = 3$$

$$x_1 = -2 \quad w_2 = -3$$

$$f_1 = w_0 x_0$$

$$f_2 = w_1 x_1$$

$$\Rightarrow \frac{1}{1} + \frac{1}{e^{-1}} + \frac{1}{e^{-6}} = \frac{1}{e^{-8}}$$

$$f_3 = f_1 + f_2 + w_2$$

$$f_4 = f_3 + w_2$$

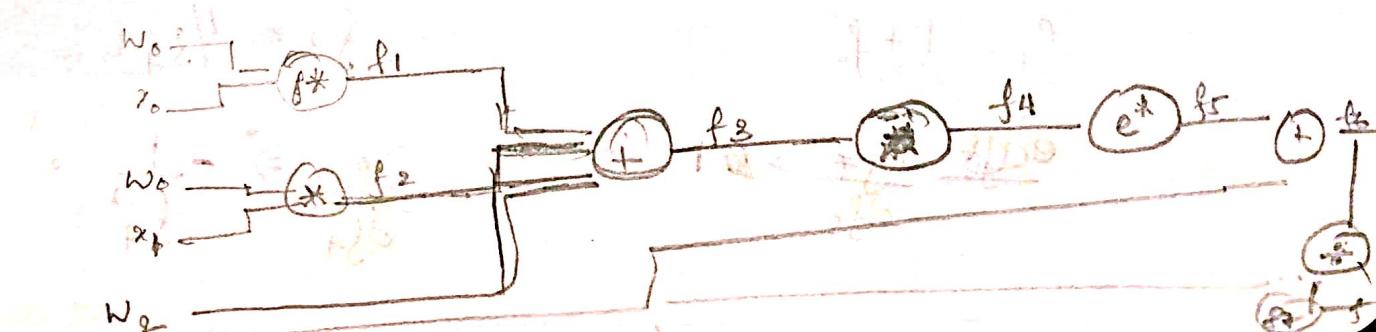
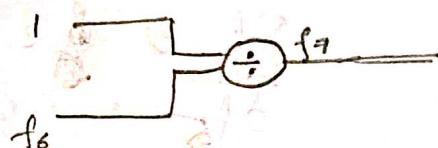
$$f_5 = -f_4$$

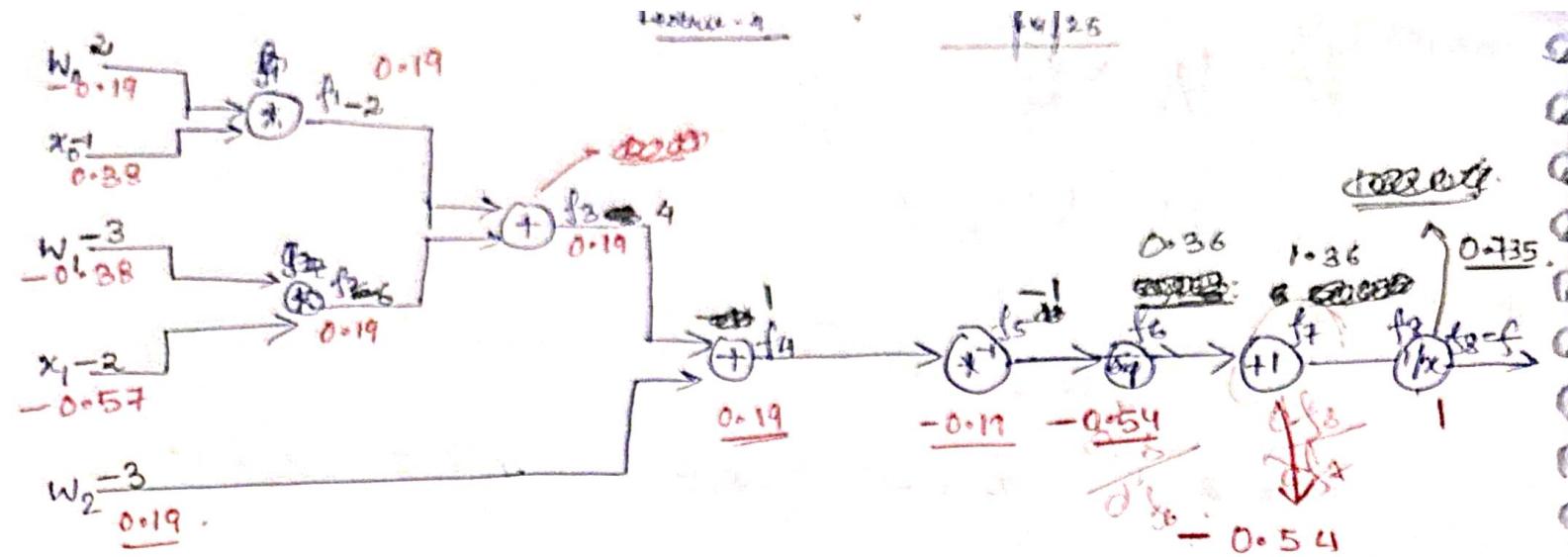
$$f_6 = e^{f_5}$$

$$f_7 = 1 + f_6$$

$$f_8 = 1/f_7$$

$$\boxed{f = f_8}$$





$$\frac{\delta f_8}{\delta f_8} = 1$$

$$f_1 = \underline{w_0 x_0}$$

$$\frac{\delta f_1}{\delta w_0} = x_0 \quad \frac{\delta f_1}{\delta x_0} = w_0$$

$$f_2 = w_1 x_1 \quad \frac{\delta f_2}{\delta w_1} \Rightarrow x_1 \quad \frac{\delta f_2}{\delta x_1} = w_1$$

$$f_3 = f_1 + f_2$$

$$\frac{\delta f_3}{\delta f_1} = 1 \quad \frac{\delta f_3}{\delta f_2} = 1$$

$$f_4 = f_3 + w_2$$

$$\frac{\delta f_4}{\delta f_3} = 1 \quad \frac{\delta f_4}{\delta w_2} = 1$$

$$f_5 = -f_4$$

$$\frac{\delta f_5}{\delta f_4} = -1$$

$$f_6 = e^{f_5}$$

$$\frac{\delta f_6}{\delta f_5} \Rightarrow e^{f_5}$$

$$f_7 = 1 + f_6$$

$$\cancel{\frac{\delta f_8}{\delta f_7}} \quad \frac{\delta f_7}{\delta f_6} = 1$$

$$f_8 = 1/f_7$$

$$\frac{\delta f_8}{\delta f_7} \Rightarrow -\frac{1}{f_7^2}$$

$$\cancel{\frac{df_7}{df_8}} \Rightarrow \frac{df_7}{df_6} \times$$

$$\frac{df_6}{df_7} \Rightarrow \frac{df_6}{df_5} \times \cancel{\frac{df_5}{df_7}}$$

$$\Rightarrow 1 \times f_5 e^{f_5}$$

$$\Rightarrow 1 \times (-1) e^{-1}$$

$$\Rightarrow -0.36$$

$$\cancel{\frac{df_6}{df_8}} \Rightarrow \cancel{\frac{df_5}{df_4}} \times \cancel{\frac{df_4}{df_7}}$$

$$\cancel{\frac{df_5}{df_8}} \Rightarrow \frac{df_6}{df_5} \Rightarrow \frac{df_6}{df_5} \times d$$

$$\frac{df_8}{df_7} = \text{local} * \text{global}$$

$$\frac{df_8}{df_7} \Rightarrow -\frac{1}{f_7^2} * 1$$

$$\Rightarrow -\frac{1}{(1-36)^2} * 1$$

$$\Rightarrow -0.54.$$

$$\frac{df_8}{df_5} = \text{local} * \text{global}$$

$$\frac{df_8}{df_5} \Rightarrow \frac{df_8}{df_6} \times \cancel{\frac{df_6}{df_5}} 0.54$$

$$\Rightarrow 0.34 \times 0.54$$

$$\Rightarrow -0.19.$$

$$\frac{df_8}{df_6} = \text{local} * \text{global}$$

$$\frac{df_8}{df_6} = \frac{df_7}{df_5} \times \frac{df_8}{df_7}$$

$$P1 \cdot 0 \Rightarrow 1 \times -0.54$$

$$P1 \cdot 0 \Rightarrow -0.54$$

$$88.3 \%$$

$$\frac{df_8}{df_4} \Rightarrow \text{local} * \text{global}$$

$$\frac{df_8}{df_4} \Rightarrow \frac{df_5}{df_4} \times \cancel{\frac{df_6}{df_5}} 0.19$$

$$P1 \cdot 3 \times 88.3 \% \Rightarrow 1 \times (-0.19)$$

$$K2 \cdot 3 \cdot 88.3 \% \Rightarrow 0.19$$

$$\frac{df_8}{df_2} = \text{local} \times \text{global}$$

$$1 \times 0.19$$

$$\Rightarrow 0.19$$

$$\frac{df_8}{df_3} = 1 \times 0.19$$

$$\Rightarrow 0.19$$

$$\frac{df_8}{df_2} = 1 \times 0.19$$

$$\Rightarrow 0.19$$

$$\frac{df_8}{df_1} = 1 \times 0.19$$

$$\Rightarrow 0.19$$

$$\frac{df_8}{dw_0} \times \text{local} = 0.19$$

$$\Rightarrow x_0 \times 0.19$$

$$\Rightarrow -0.19$$

$$\frac{df_8}{dx_0} \Rightarrow w_0 \times 0.19$$

$$\Rightarrow \alpha \times 0.19$$

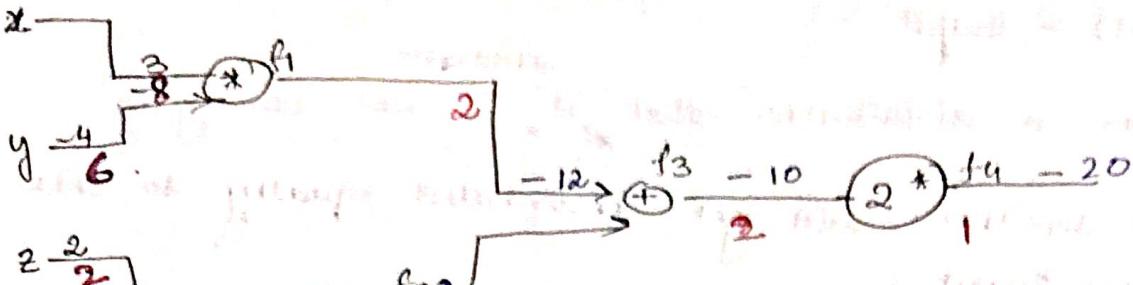
$$\Rightarrow 0.38$$

$$\frac{df_8}{dw_1} \Rightarrow -2 \times 0.19$$

$$\Rightarrow -0.38$$

$$\frac{df_8}{d\alpha} \Rightarrow -3 \times 0.19$$

$$\Rightarrow -0.57$$



$$f_1 = x * y$$

$$\frac{\partial f_1}{\partial x} \Rightarrow y \quad \frac{\partial f_1}{\partial y} = x$$

$$f_1 + f_2$$

$$f_3 \Rightarrow \text{min}(z, w) / \text{max}(z, w)$$

$$\frac{\partial f_3}{\partial f_1} \Rightarrow 1$$

$$\frac{\partial f_3}{\partial f_2} = 1$$

$$f_4 \Rightarrow 2 * f_3$$

$$f_2 = \min(z, w) / \max(z, w)$$

$$\frac{\partial f_2}{\partial z} = 1 \quad \frac{\partial f_2}{\partial w} = 1$$

$$\frac{\partial f_4}{\partial f_3} \Rightarrow 2$$

$$\frac{\partial f_4}{\partial f_2} = 1$$

$$\begin{aligned} \frac{\partial f_2}{\partial z} &= 1 \text{ if } z > w \\ &= 0 \text{ if } z < w \\ &\Rightarrow \text{Undef. if } z = w \end{aligned}$$

$$\frac{\partial f_4}{\partial f_2} \Rightarrow 2 \times 1$$

$$\frac{\partial f_4}{\partial f_3} \Rightarrow 2$$

$$\begin{aligned} \frac{\partial f_1}{\partial f_4} &\Rightarrow \text{local} \times \text{global} \\ &\Rightarrow x^1 \end{aligned}$$

$$\begin{aligned} \frac{\partial f_4}{\partial f_n} &\Rightarrow y^2 \\ &\Rightarrow -4 \times 2 \end{aligned}$$

$$\frac{\partial f_2}{\partial w} = 1 \text{ if } w > z$$

$$\Rightarrow 0 \text{ if } w < z$$

$$\begin{aligned} &\Rightarrow \text{Undef} \\ &\text{or} \\ &\text{arbitrary} \end{aligned}$$

$$\begin{aligned} \frac{\partial f_4}{\partial f_y} &\Rightarrow 0^3 \times 2 \\ &\Rightarrow 0^6 \end{aligned}$$



$$z_1 = 1 \times 0.5 \Rightarrow 0.5$$

$$a_1 = g(z_1) = 0.222$$

$\approx 0.222$

$\approx$

$$z_2 \Rightarrow 0.622 \times 0.5$$

$$\Rightarrow 0.3110$$

$$a_2 \Rightarrow g(z_2) \approx 0.572$$

$$\Rightarrow 0.572$$

$$z_3 \Rightarrow 0.572 \times 0.5$$

$$\approx 0.2885$$

$$a_3 \Rightarrow g(z_3)$$

$$\Rightarrow 0.571$$

$$\frac{\partial L}{\partial a_3} \Rightarrow a_3 - y = 0.572 - 1 \\ \Rightarrow -0.428$$

$$\Rightarrow \frac{1}{2} (0.572 - 1)^2$$

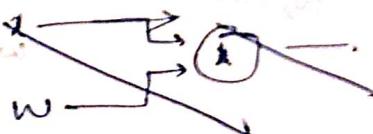
$$\approx 0.0915$$

$$\approx 0.0912$$

$$\frac{\partial L}{\partial a_3} \wedge \frac{\partial a}{\partial z_3}$$

$$\frac{\partial a}{\partial z_3} \Rightarrow g(z) \left(1 - g(z)\right)$$

$$\approx 0.082$$



(\*)

Layer: L<sub>3</sub>; Gradient = -0.105

Layer: L<sub>2</sub>; Gradient  $\Rightarrow -0.013$

Layer: L<sub>1</sub>; Gradient  $\Rightarrow -0.0015$

on 1024 pixels we need to train the network.

$$W = 3 \times 1000$$

$\Rightarrow 8^{\text{th}}$  million

## Convolution Operation

3	0	1	2	7	4
1	5	-1	9	3	1
2	4	2	5	1	3
0	1	3	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

$\downarrow$   
Convolution  
operated.

1	0	-1
1	0	-1
1	0	-1

~~Input~~  $\Rightarrow 6 \times 6$  image filter  $\Rightarrow 3 \times 3$

-5	-4	0	8
-10	-2	2	3
0	-2	-4	-1
-3	-2	-3	-6

$$\begin{aligned}
 & 2 + 0 - 4 + 9 + 0 - 1 + 5 + 0 \\
 & - 3 = 8
 \end{aligned}$$

Openings into next layer

$$8 \times 1 + 0 \times 0 + 1 \times -1 + 1 \times 1 + 5 \times 0 + 8 \times -1 + 2 \times 1 + 7 \times 0 + 2 \times -1 \\
 \Rightarrow -5$$

Output

$$0 + 0 + -2 + 8 + 7 + 0 + 0 = 15$$

$$\Rightarrow -4$$

$$\Rightarrow 1 + 0 + -7 + 8 + 0 - 3 + 2 - 0 = 1$$

$$\Rightarrow 1$$

→ output will be skewed and less visibility of the pixels.

→ If we want attention to put all some middle pixels or want to miss some pixels then we can skip the layer and use the strided convolution.

→ we can use padding; if we want do more attention at the edges.

$\Rightarrow \text{if } f=2$

2	3	7	4	6
6	6	9	8	7
3	4	8	3	8
7	8	3	6	6
4	2	1	8	3

$$m = 5 \times 5$$

$$p = 1; s = 2$$

0	0	0	0	0	0	0
0	2	3	7	4	6	0
0	6	6	9	8	7	0
0	3	4	8	3	8	0
0	7	8	3	6	6	0
0	4	2	1	8	3	0
0	0	0	0	0	0	0

$$= 0 + 0 + 0 - 2$$

$$= 0 + 0 + 0$$

$$= 0 - 6 - 3$$

1	1
0	1

$$= \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$$

$$f = 2 \times 2$$

$$\Rightarrow n+2p-f+1$$

$$\Rightarrow n+2-2+1$$

$$\Rightarrow \frac{6}{2} \Rightarrow 3$$

$$\begin{array}{c} \cancel{n+2p-f+1} \\ \cancel{5+2-2+1-2} \\ \Rightarrow 3 \end{array}$$

1	-1
0	-1

$$= \begin{bmatrix} 1 & -1 \\ 0 & -1 \end{bmatrix}$$

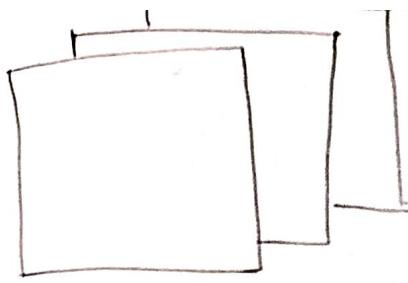
-2	-7	-6
-9	-11	-7
-11	8	0

$$3 \times 3$$

$$\Rightarrow 6 - 9 + 0 - 8$$

$$\Rightarrow 6 - 15$$

$$\Rightarrow 8 - 7 + 0 - 8$$



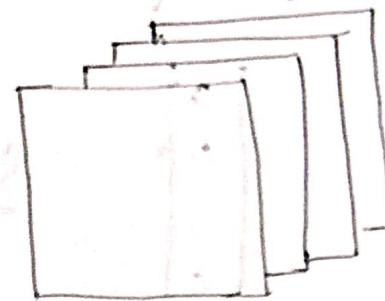
$39 \times 39 \times 3$

$$f[1] = 3$$

$$s[1] = 1$$

$$p[1] = 0$$

10 filters



$34 \times 37 \times 10$

The size will decrease as padding is zero.

so the new size will become  $37 \times 37 \times 3$ .

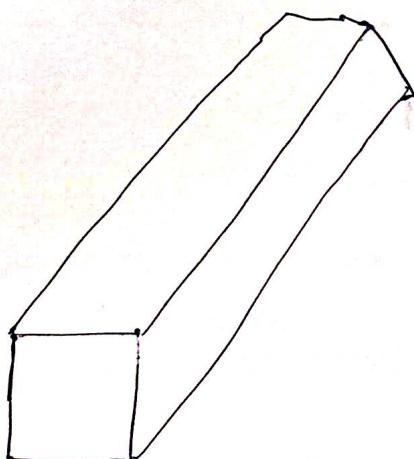
$$f[2] = 5$$

$$s[2] = 2$$

$$p[1] = 0$$

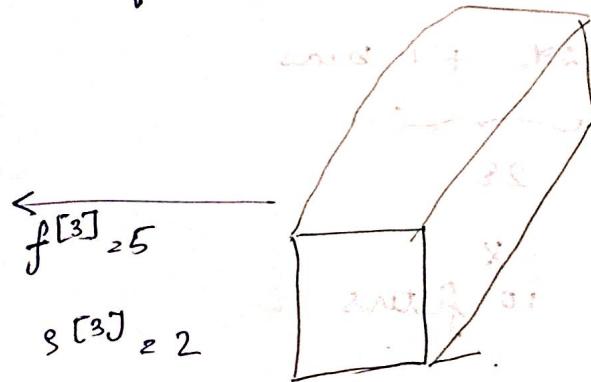
20 filters

- After applying the convet the size get decreases  
No. of channels depend on filters



$$7 \times 7 \times 40 = 1960 \text{ pixels}$$

PTO



$$f[3] = 5$$

$$s[3] = 2$$

$$p[1] = 0$$

$$14 \times 17 \times 20$$

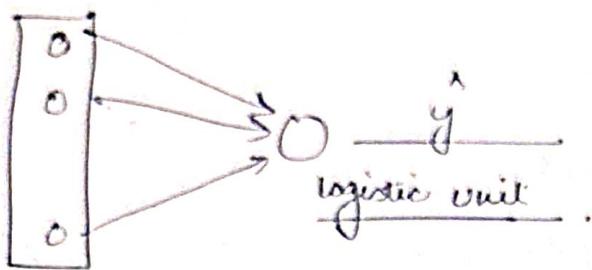
40 pixels.

formulae -

$$\Rightarrow \left[ \frac{n+2p+f}{s} + 1 \right]$$

$$\Rightarrow \left[ \frac{37-5}{2} + 1 \right]$$

→ Apply a linear layer



I/P  $\rightarrow$  CONV - CONV - CONV - FC

discriminatory features

19/8/25

We have 10 filters each of size  $3 \times 3 \times 3$  in one layer of convnet.

How many parameters does this layer have?

$$9 \times 3 \Rightarrow 27$$

$$\underbrace{3 \times 3 \times 3}_{27} \rightarrow 27 + 1 \text{ bias}$$

28

$$X \\ 10 \text{ filters}$$

$$\Rightarrow 280 \text{ parameters}$$

If we know that  $a[g(z)]$  is better than the output will be  $a[z]$ .

$$a^{[l+2]} = g\left(w^{[l+2]_{l+1}} + b^{[l+2]} + a^{[l]}\right).$$

$$\approx g\left(z^{[l+2]} + \underline{Na^{[l]}}\right).$$

$$g(256 + 128)$$

$$a^{[l+2]} = \text{ReLU}(z^{[l+2]} + a^{[l]})$$

$$32 \times 1 = 32 \times 1$$

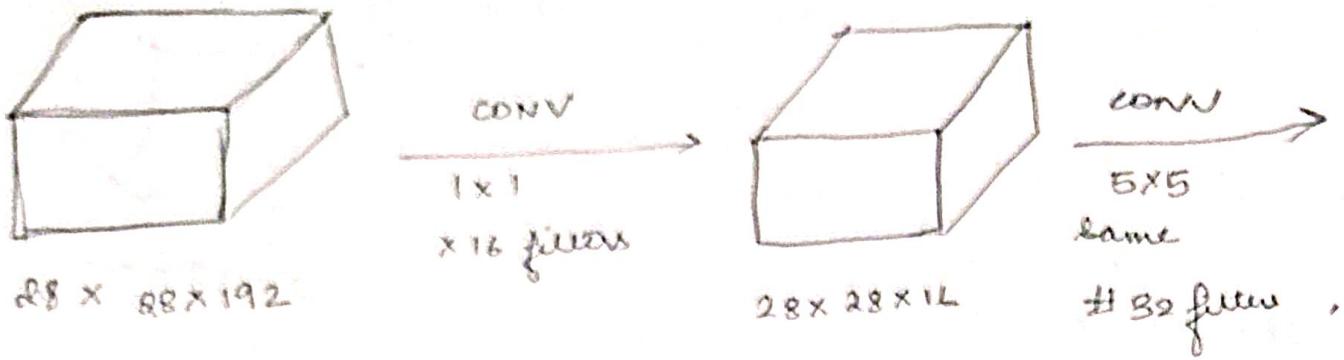
$$32 \times 1 = 32 \times 1 + \boxed{W^1} \rightarrow \text{This } 64 \times 1$$

$\downarrow$   
This should be  $\frac{32}{64}$

• ResNet can do deeper.

$$\begin{aligned} \text{ResNet} &= 16 \\ n &= 32 \end{aligned}$$

- First use the existing resnet and see and the transform the images.



1st layer

$$\text{For } 1 \times 1 \times 192 = 192.$$

$$192 \times 16 = 3072.$$

$$28 \times 28 = 784.$$

8,808

24,08,448-

for 2nd layer  $\Rightarrow$

$$\Rightarrow 5 \times 5 \times 16 = 400.$$

$$400 \times 32 = \underline{\underline{12,800}}$$

4080

$$\Rightarrow 784 = 1,00,35,200$$

$$\therefore \text{Total} = 1,00,35,200 + 24,08,448.$$

$$\Rightarrow 12,44,36,48 \approx 12 \text{ million}.$$

If we need to ① modify height, width.

$\rightarrow$  Padding

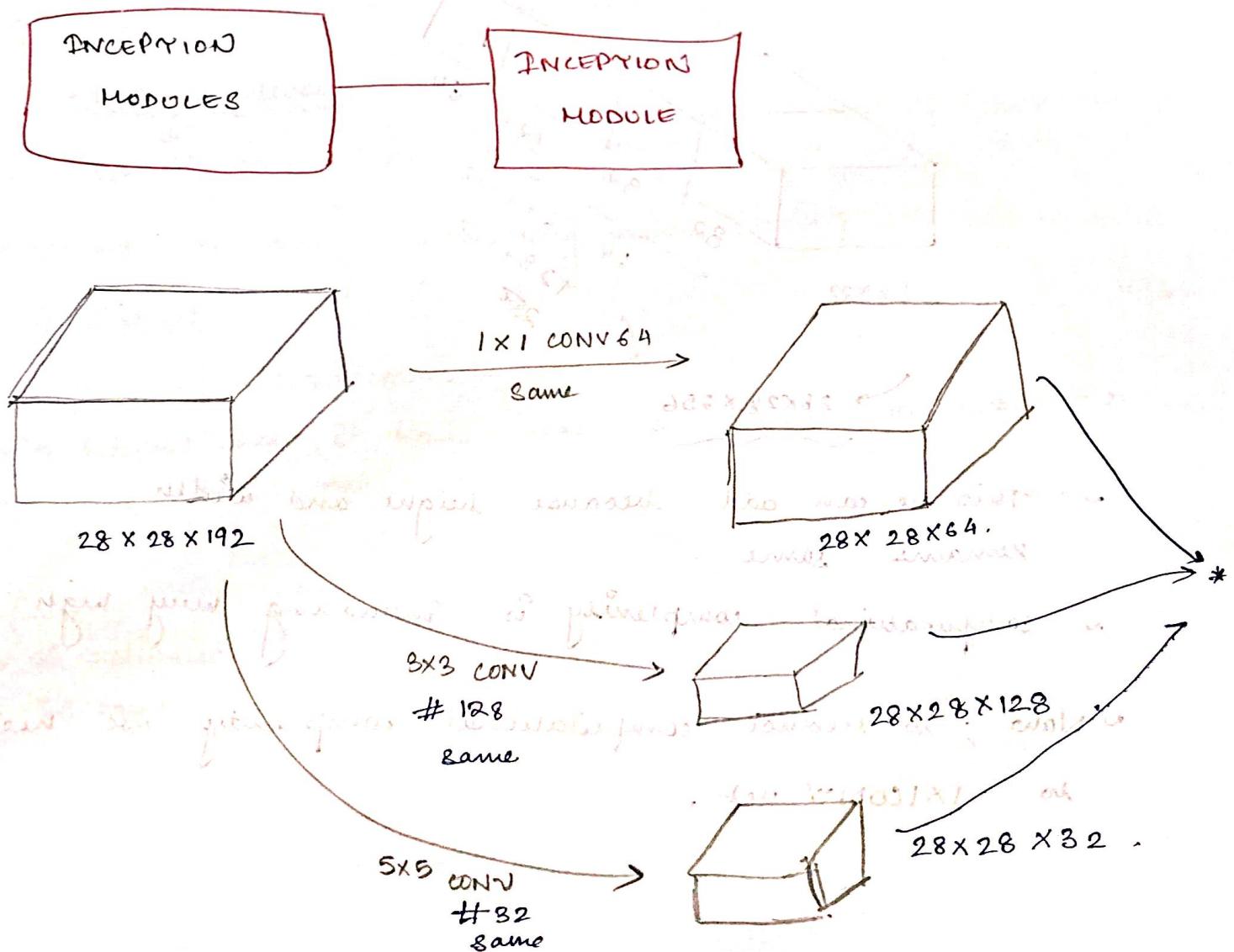
$\rightarrow$  Striding, Pooling.

② modify channels  $\rightarrow$  filters,  $1 \times 1$  CONV net.

The layers are known as bottleneck layers

- ReLU is included in the conv layer by default but if we needed we can remove them.
- ReLU is important in the conv net, Resnet or feed neural network it is very important so we therefore the non-linearity to be put on them.

### INCEPTION NETWORK



- It will be shared across the timesteps
- Input is of varying length, then the input size is different.
- Output is same that input's dimension that is size sharing of weights.

Vanilla RNN → to handle variable length

817/9/25

$$x = \mathbb{R} \text{ (two time steps).}$$

$$x^{(1)}: x_1 \in \mathbb{R}^2 \\ = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$x_2 \in \mathbb{R}^2 \\ = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

$$\text{input dimensions} = 2.$$

$$\Rightarrow x_t \in \mathbb{R}^2.$$

$$\text{hidden state dim} = 3 \Rightarrow h_t \in \mathbb{R}^3$$

$$h_0 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad 3 \times 1$$

$$w_{xh} = \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} \quad 3 \times 2$$

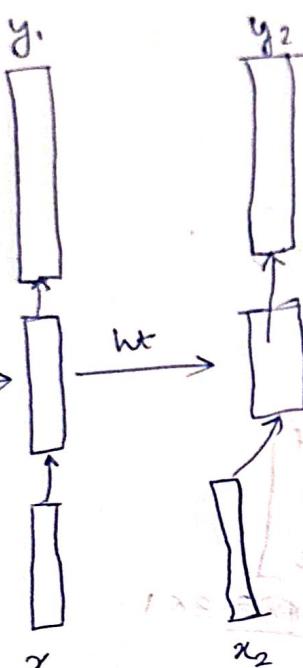
$$w_{hh} = \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} \quad 3 \times 3$$

$$3 \times 1 \times 1 \times 2$$

$$w_{hy} = \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} \quad 3 \times 3$$

$$3 \times 3$$

$$w_{xh} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \quad w_{hh} = \begin{bmatrix} 0.5 & 0.3 \\ 0.2 & 0.2 \\ 0.1 & 0.4 \end{bmatrix}_{3 \times 2}, \quad w_{hy} = \begin{bmatrix} 0.1 & 0.4 & 0.7 \\ 0.2 & 0.3 & 0.2 \\ 0.05 & 0.8 & 0.2 \end{bmatrix}_{3 \times 3}, \quad b_{hy} = \begin{bmatrix} -0.1 \\ 0.5 \\ 0.5 \end{bmatrix}_{3 \times 1}$$



$$h_k = \tanh(w_{hk} h_{k-1} + w_{xh} x_k)$$

$$y_k = W_{hy} h_k$$

At time step = 1; ( $k=1$ )

$$\begin{aligned} & \text{Input: } x_1, x_2 \text{ at } t=0 \\ & \text{Hidden state: } h_0 \text{ at } t=0 \\ & \text{Weights: } w_{xh} = \begin{bmatrix} 0.1 & 0.4 & 0.0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.2 \end{bmatrix}_{3 \times 3}, \quad w_{hh} = \begin{bmatrix} 0.5 & 0.3 \\ 0.2 & 0.2 \\ 0.1 & 0.4 \end{bmatrix}_{3 \times 2}, \quad w_{hy} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}_{3 \times 1} \\ & \text{Bias: } b_{hy} = \begin{bmatrix} -0.1 \\ 0.5 \\ 0.5 \end{bmatrix}_{3 \times 1} \end{aligned}$$

$$\begin{aligned} & h_1 = \tanh(0) + \begin{bmatrix} -0.1 \\ 0.1 \\ 0.9 \end{bmatrix} = \begin{bmatrix} 0.5 - 0.6 \\ 0.8 + 0.4 \\ 0.1 + 0.8 \end{bmatrix} = \begin{bmatrix} -0.1 \\ 1.2 \\ 0.9 \end{bmatrix} \\ & y_1 = \begin{bmatrix} 0.5 & -0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.4 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \end{aligned}$$

$$\tanh(0.1) \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} -0.1 \\ 1.2 \\ 0.9 \end{bmatrix}$$

$$\tanh(-0.1) = -0.0997$$

$$\tanh(1.2) = 0.8337$$

$$\tanh(0.9) \approx 0.7163$$

$$= \tanh \begin{bmatrix} -0.1 \\ 1.2 \\ 0.9 \end{bmatrix}$$

$$\Rightarrow y = \begin{bmatrix} -0.0997 \\ 0.8337 \\ 0.7163 \end{bmatrix} \begin{bmatrix} 1.0 & -1.0 & 0.5 \\ 0.5 & 0.5 & -0.5 \end{bmatrix}$$

At timestep ( $t=2$ )

$$\text{eq} =$$

$$h_2 = \tanh(\omega_{hh} h_{t-1} + w_{xh} x_t)$$

$$h_2 = \tanh(\omega_{hh} h_1 + w_{xh} x_2)$$

$$h_2 = \tanh \left( \begin{bmatrix} 0.1 & 0.4 & 0.07 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.01 & 0.02 \end{bmatrix} \begin{bmatrix} -0.0997 \\ 0.8337 \\ 0.7163 \end{bmatrix} + \begin{bmatrix} 0.5 & -0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.4 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} \right)$$

$$h_2 = \tanh \left( \begin{bmatrix} -0.0997 \\ 0.8337 \\ 0.7163 \end{bmatrix} + \begin{bmatrix} 0.5 & -0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.4 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} \right)$$

$$\Rightarrow \tan^{-1} \left( \begin{bmatrix} -0.0997 + 0.33 + 0.4163 \\ 0.0997 + 0.25011 + 0.14326 \\ -0.0049 - 0.8337 + 0.14326 \end{bmatrix} + \begin{bmatrix} -0.5 - 0.3 \\ -0.8 + 0.2 \\ -0.1 + 0.4 \end{bmatrix} \right)$$

$$\tan^{-1} \left( \begin{bmatrix} 0.2383 \\ 0.4129 \\ -0.69534 \end{bmatrix} + \begin{bmatrix} -0.8 \\ -0.6 \\ 0.3 \end{bmatrix} \right)$$

$$\tan^{-1} \begin{pmatrix} -0.5647 \\ -0.1871 \\ -0.395 \end{pmatrix}$$

$$= \begin{pmatrix} -0.5090 \\ -0.1849 \\ -0.3753 \end{pmatrix}.$$

~~Eq 2 is identity -~~

$$y_1 = \text{Wday} \times h_{(1)} \dots$$

$$\rightarrow y_1 \begin{bmatrix} -0.57 \\ 0.008 \end{bmatrix},$$

$$\Rightarrow \begin{bmatrix} 1.0 & -1.0 & 0.5 \\ 0.5 & 0.5 & -0.5 \end{bmatrix} \begin{bmatrix} 0.0997 \\ 0.8337 \\ 0.14326 \end{bmatrix}.$$

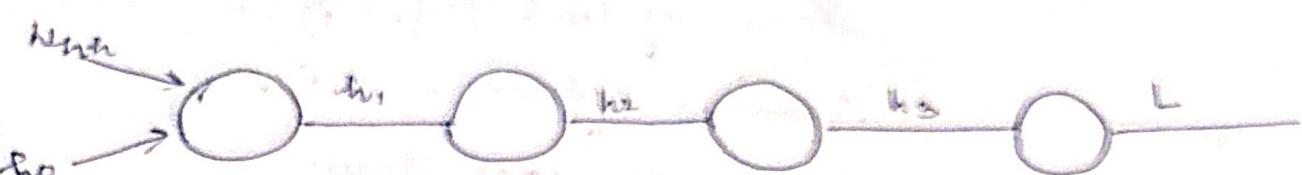
$$\begin{bmatrix} 0.0997 \\ 0.8337 \\ 0.14326 \end{bmatrix}$$

$$\begin{bmatrix} 0.0997 - 0.8337 + 0.35815 \\ 0.0498 + 0.4168 - 0.35815 \end{bmatrix}$$

$$\approx \begin{bmatrix} -0.37 \\ 0.1085 \end{bmatrix}.$$

## Vanishing Gradient Problem

Vanilla RNN or multilayer RNN will lead to Vanishing Gradient Problem.



$$\frac{\partial L}{\partial w_{hi}} = \frac{\partial L}{\partial h_3} \cdot \frac{\partial h_3}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial w_{hi}}$$

$$\Rightarrow \frac{\partial L}{\partial h_3} \left( \frac{\partial h_3}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} \right) \frac{\partial h_1}{\partial w_{hi}}$$

$$= \frac{\partial L}{\partial h_3} \left( \prod_{t=2}^{T=3} \frac{\partial h_t}{\partial h_{t-1}} \right) \frac{\partial h_1}{\partial w_{hi}}$$

for T time steps.

$$\frac{\partial L}{\partial w_{hi}} = \frac{\partial L}{\partial h_T} \left( \prod_{t=2}^T \frac{\partial h_t}{\partial h_{t-1}} \right) \frac{\partial h_1}{\partial w_{hi}}$$

$$h_t = \tanh(w_{hi} h_{t-1} + w_{xi} x_t)$$

this is gradient

because it deals with multivariate class or vector  
derivative  $\times$  single layer

$$\frac{\partial h_t}{\partial h_{t-1}} = \tanh' (w_{hi} h_{t-1} + w_{xi} x_t) w_{hi}$$

we need to calculate

derivative  
through

$$c_t = i_t \odot c_{t-1} + f_t \odot g_t$$

$$h_t = Q \odot \tanh(c_t)$$

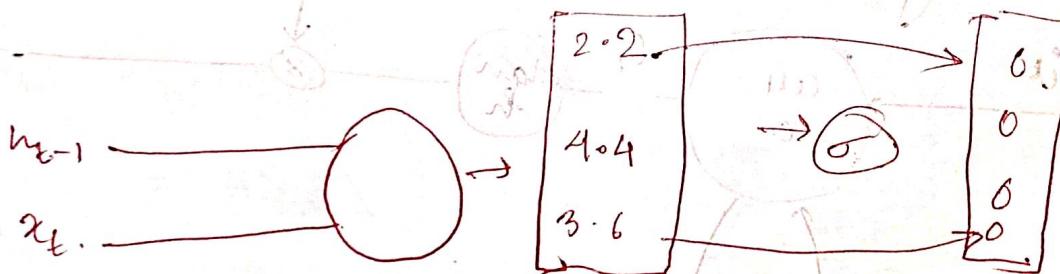
$$f_t = \sigma(w_{hf} h_{t-1} + w_{xf} x_t)$$

$$i_t = \sigma(w_{hi} h_{t-1} + w_{xi} x_t)$$

$$o_t = \sigma(w_{ho} h_{t-1} + w_{xo} x_t)$$

$$g_t = \tanh(w_{hg} h_{t-1} + w_{xg} x_t)$$

These gates are sigmoid because they have a range between (0, 1). That's why they are known as smooth gate.



$$i_t = \sigma(w_{hi} h_{t-1} + w_{xi} x_t)$$

because  
h is input  
and i depends  
on it.

Precious  
if  
Present  
output.

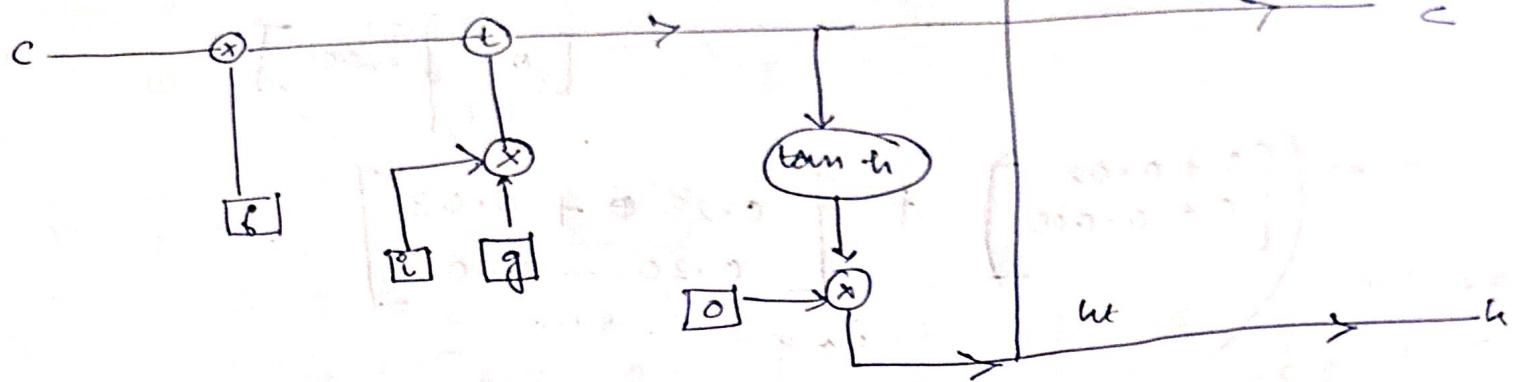
→ if it is a vector whose values (0, 1).

$g_t$  is another gate.

$$g_t = \tanh(w_{hg} h_{t-1} + w_{hg} x_t).$$

$i_t$  and  $g_t$  will be element wise operation.

$h = \text{higher analysis}$ .



We use LSTMs ① to avoid vanishing gradient. (not 100% but minimizes the problem)

Example:

$$x_t = [0.5, -0.1]^T$$

$$h_{t-1} = [0.0, 0.1]^T$$

$$c_{t-1} = [0.2, -0.2]^T$$

$$Wx_i = \begin{bmatrix} 0.5 & -0.3 \\ 0.4 & 0.1 \end{bmatrix}$$

$$Wh_i = \begin{bmatrix} 0.1 & 0.2 \\ -0.2 & 0.05 \end{bmatrix}$$

$$Wx_f = \begin{bmatrix} -0.4 & 0.2 \\ 0.3 & 0.3 \end{bmatrix}$$

$$Wh_f = \begin{bmatrix} 0.05 & -0.1 \\ 0.2 & 0.1 \end{bmatrix}$$

$$Wx_o = \begin{bmatrix} 0.3 & 0.25 \\ -0.2 & 0.2 \end{bmatrix}$$

$$Who = \begin{bmatrix} 0.15 & 0.05 \\ 0.1 & -0.2 \end{bmatrix}$$

$$Wng = \begin{bmatrix} -0.5 & 0.4 \\ 0.2 & -0.3 \end{bmatrix}$$

$$Wng = \begin{bmatrix} 0.2 & 0.1 \\ -0.1 & 0.05 \end{bmatrix}$$

$$\Rightarrow i_t = - (w_{12} i_{t-1} + w_{22} o_2)$$

$$\Rightarrow - \left( \begin{bmatrix} 0.1 & 0.2 \\ -0.2 & 0.05 \end{bmatrix} \begin{bmatrix} 0.0 \\ 0.1 \end{bmatrix} \right) + \begin{bmatrix} 0.5 & -0.3 \\ 0.0 & 0.1 \end{bmatrix}$$

$$= - \left( \begin{bmatrix} 0 + 0.02 \\ 0 + 0.005 \end{bmatrix} \right) + \begin{bmatrix} 0.25 + 0.03 \\ 0.20 - 0.01 \end{bmatrix}$$

$$= - \left( \begin{bmatrix} 0.02 \\ 0.005 \end{bmatrix} \right) + \begin{bmatrix} 0.28 \\ 0.19 \end{bmatrix}$$

$$= - \left( \begin{bmatrix} 0.30 \\ 0.024 \\ 0.224 \\ 0.5460 \end{bmatrix} \right)$$

$$\therefore \Rightarrow \begin{bmatrix} 0.5744 \\ 0.055344 \end{bmatrix}$$

$$0.5744 = - \left( \begin{bmatrix} 0.15 & 0.05 \\ 0.1 & -0.2 \end{bmatrix} \begin{bmatrix} 0.0 \\ 0.1 \end{bmatrix} \right) + \begin{bmatrix} 0.3 & 0.25 \\ -0.2 & 0.2 \end{bmatrix} \begin{bmatrix} 0.5 \\ -0.1 \end{bmatrix}$$

$$= - \left[ \begin{bmatrix} 0.0 + 0.05 \\ 0.0 - 0.02 \end{bmatrix} \right] + \begin{bmatrix} 0.15 + 0.025 \\ -0.10 - 0.00 \end{bmatrix}$$

$$= \begin{bmatrix} 0.05 \\ -0.02 \end{bmatrix} + \begin{bmatrix} 0.12[5] \\ -0.12 \end{bmatrix}$$

$$\Rightarrow \begin{pmatrix} 0.175 \\ -0.14 \end{pmatrix}$$

$$q_t = \begin{bmatrix} 0.5336 \\ -0.46 \end{bmatrix} \dots$$

$$q_t = \tanh \left( \begin{bmatrix} 0.2 & 0.1 \\ -0.1 & 0.05 \end{bmatrix} \begin{bmatrix} 0.0 \\ 0.1 \end{bmatrix} + \begin{bmatrix} -0.5 & 0.4 \\ 0.2 & -0.3 \end{bmatrix} \begin{bmatrix} 0.5 \\ -0.1 \end{bmatrix} \right)$$

$$q_t = \tanh \left( \begin{bmatrix} 0.0 + 0.01 \\ -0.0 + 0.005 \end{bmatrix} + \begin{bmatrix} -0.25 & 0.04 \\ 0.10 & 0.03 \end{bmatrix} \right)$$

$$\tanh \left( \begin{bmatrix} 0.01 \\ 0.005 \end{bmatrix} + \begin{bmatrix} -0.29 \\ 0.13 \end{bmatrix} \right)$$

$$q_t = \tanh \begin{bmatrix} -0.28 \\ 0.135 \end{bmatrix}$$

$$q_t = [-0.273, 0.134]$$

$$h_t = [0.44 \quad 0.53]^T$$

$$c_t = [i_t \odot c_{t-1} + o_t \odot g_t]$$

$$= \begin{bmatrix} 0.44 & 0.53 \end{bmatrix} (0.2, -0.2) + \begin{bmatrix} 0.0 & 0.0 \end{bmatrix} \begin{bmatrix} -0.273 \\ 0.137 \end{bmatrix}$$

$$c_t = \left( \begin{bmatrix} 0.088 & -0.106 \end{bmatrix} + \begin{bmatrix} 0 & 0.013 \end{bmatrix} \right).$$

$$c_t = \begin{bmatrix} 0.088 & -0.106 \end{bmatrix} + \begin{bmatrix} 0.1539 & 0.073 \end{bmatrix}$$

$$c_t = \begin{bmatrix} 0.88 & -0.093 \end{bmatrix} \cdot \begin{bmatrix} -0.06 & -0.03 \end{bmatrix}$$

$$h_t = o_t \times \tanh \begin{bmatrix} 0.88 & -0.093 \end{bmatrix}$$

$$h_t = \begin{bmatrix} 0.53 & 0.46 \end{bmatrix} \times \begin{bmatrix} 0.707 & -0.0916 \end{bmatrix}$$

$$h_t = \begin{bmatrix} 0.53 & 0.46 \end{bmatrix} \times \begin{bmatrix} -0.059 & -0.029 \\ 0.0014 & 0.0003 \end{bmatrix}$$

$$h_t = \begin{bmatrix} -0.031 & -0.013 \end{bmatrix}$$

### RNN

$$\frac{\partial h}{\partial h_{t-1}} = \tanh'(w_{hh} h_{t-1} + w_{hx} x_{t-1})$$

### LSTM

$$\frac{\partial h_t}{\partial h_{t-1}} = O_t \tanh'(c) \boxed{\frac{\partial c_t}{\partial h_{t-1}}} + \frac{\partial O_t}{\partial h_{t-1}} \tanh(c)$$

$$\frac{\partial c_t}{\partial h_{t-1}} = f_t \frac{\partial c_{t-1}}{\partial h_{t-1}} + c_{t-1} \frac{\partial f_t}{\partial h_{t-1}} + g_t \frac{\partial i_t}{\partial h_{t-1}} + \boxed{u_t \frac{\partial g_t}{\partial h_{t-1}}}$$

It is the one which LSTM prevents vanishing gradient problem is minimized but it can go to zero so rather than neglecting that sometimes it can lead to zero.

Disadvantages: It is computationally very expensive.

$$O_t \tanh'(c) \frac{\partial c_t}{\partial h_{t-1}}$$

$$O_t \tanh'(c) \frac{\partial c_t}{\partial h_{t-1}} + \frac{\partial O_t}{\partial h_{t-1}} \tanh(c)$$

So we need to calculate

$$= h_1 \cdot \text{score}[[1, 0, 1], t_{t+1}] \text{ and } a = [x_{t+1}, h_1] + \alpha_{t+2} h_2 + \alpha_{t+3} h_3$$

$$h_2 = [0, 1, 1]$$

(dot product)

$$h_3 = [1, 1, 0]$$

(dot product)

$$s_{t+1} = [1, 0, 1] [F, F+1, F+2, F+3]$$

$$x_{t+1} \geq \frac{\exp(\text{score}[1, 0, 1], [1, 0, 1])}{\exp(\text{score}[0, 1, 0], [1, 0, 1]) + \exp(\text{score}[1, 0, 1], [0, 1, 0]) + \exp(\text{score}[1, 0, 1], [1, 1, 0])}$$

$$x_{t+1} \geq \frac{\exp[1+0+1]}{\exp[1+0+1] + \exp[0+0+1] + \exp[1+0+0]}$$

$$x_{t+1} = \frac{\exp[2]}{\exp[2] + \exp[1] + \exp[0]}$$

$$x_{t+2} \geq \frac{2 \cdot 7}{12 \cdot 8}$$

$$\alpha_{t+1} = \frac{7 \cdot 4}{7 \cdot 4 + 2 \cdot 7 + 2 \cdot 7}$$

$$\geq 0.21$$

$$\alpha_{t+1} \geq \frac{7 \cdot 4}{12 \cdot 8}$$

$$x_{t+3} \geq 0.21$$

$$\Rightarrow 0.578$$

$$C_t = 0.57 [1, 0, 1] + 0.2 [0, 1, 1] + 0.2 [1, 1, 0]$$

$$\Rightarrow [0.57 \quad 0 \quad 0.57] + [0.2 \quad 0.2 \quad 0.2] + [0.2 \quad 0.2 \quad 0]$$

$$\Rightarrow [0.57 \quad 0.2 \quad 0.77] + [0.2 \quad 0.2 \quad 0]$$

$$q \Rightarrow [0.77 \quad 0.2 \quad 0.2]$$

As  $C_t$  is closer to  $h_1$ , it is paying more attention to  $h_1$ .

$$E_t = \frac{[1 \ 0 \ 1]}{100} + \frac{[1 \ 1 \ 1]}{100} + \frac{[1 \ 1 \ 0]}{100}$$

$$= \frac{[1 \ 0 \ 1]}{100} + \frac{[1 \ 1 \ 1]}{100} + \frac{[1 \ 1 \ 0]}{100}$$

$$= \frac{[1 \ 0 \ 1]}{100} + \frac{[1 \ 1 \ 1]}{100} + \frac{[1 \ 1 \ 0]}{100}$$

$\Delta E_t$

$E_t$  ~~100~~ 100

$\Delta E_t$

$E_t$  ~~100~~ 100

$\Delta E_t$

$E_t$  ~~100~~ 100

vector database	Key	Value	
	$[1.0, 0.5, 0.2]$	"Action movie Die hard"	$\rightarrow \textcircled{2}$
	$[0.3, 0.0, 0.1]$	"Romantic Titanic"	$\rightarrow \textcircled{3}$
	$[0.9, 0.4, 0.8]$	"Action Movie" John Wick	$\rightarrow \textcircled{1}$

$\alpha$  = looking for action movie

$\Downarrow$  word2vec

$$[1.0, 0.3, 0.5] \cdot [1 + 0.15 + 0.01]$$

Dot product

$$\Downarrow [1 + 0.15 + 0.01]$$

$$\Rightarrow 1.25$$

$$\Rightarrow 0.3 + 0.3 + 0.05$$

$$\Rightarrow 0.65$$

$$\Rightarrow 0.9 + 0.12 + 0.40$$

$$\Rightarrow 1.42$$

# Transformers

9/10/25

It has encoder and decoder



Self attention is a part of Multi Head Attention.

Project the embedding vector into some space.

Compute the  $z_1 \ 3 \ z_2$  for the input sentence "playing outside".

Embedding

First step: - Word 2 Vec.

Multiply that by  $W_q \ W_K \ W_V$

and  $q_1 = \begin{bmatrix} 0.04 \\ 0.212 \\ 0.089 \\ 0.63 \\ 0.36 \end{bmatrix}^T$ .

$$k_1 = [0.31 \ 0.89 \ 0.963 \ 0.57]^T$$

$$v_1 = [0.36 \ 0.83 \ 0.1 \ 0.38]^T$$

outside :-

$$q_2 = [0.1 \ 0.14 \ 0.86 \ 0.77]^T$$

$$k_2 = [0.45 \ 0.94 \ 0.73 \ 0.58]^T$$

$$v_2 = [0.31 \ 0.36 \ 0.19 \ 0.72]^T$$

$$q_1 \cdot K_1 = [0.212 \times 0.31 + 0.04 \times 0.84 + 0.63 \times 0.963 + 0.36 \times 0.57] \quad \text{PH-0}$$

$$+ 0.36 \times 0.57 \quad \text{PH-0}$$

$$q_1 \cdot K_1 = 0.06572 + 0.0336 + 0.60669 + 0.2052$$

$$q_1 \cdot K_1 = 10.9112 \rightarrow 0.91$$

Divide by  $\sqrt{dk}$

dimension of key

$$\sqrt{4}$$

$$\frac{2}{\cancel{2}}$$

$$\frac{10.9112}{\cancel{2}} = 0.4556$$

$$\Rightarrow 0.46$$

$\rightarrow$  score

Softman.

$$q_1 \cdot K_2 = [0.212 \times 0.45 + 0.04 \times 0.94 + 0.63 \times 0.73 + 0.36 \times 0.5] \quad \text{PH-0}$$

$$\Rightarrow 0.0954 + 0.0376 + 0.4599 + 0.2088$$

$$\Rightarrow 0.8017$$

$$\Rightarrow 0.4008 \rightarrow \text{score}.$$

Divide by  
2

softman

$\Rightarrow$

$$0.5138$$

$$0.4862$$

~~0.52~~

$$0.52$$

$$0.49$$

$$0.52x \begin{bmatrix} 0.36 & 0.83 & 0.1 & 0.38 \\ 0.1872 & 0.4816 & 0.052 & 0.1976 \end{bmatrix}$$

~~200~~

$$0.72x \begin{bmatrix} 0.31 & 0.36 & 0.19 & 0.72 \\ 0.0372 & 0.0432 & 0.0228 & 0.086 \end{bmatrix}$$

$$\cancel{\begin{bmatrix} 0. \\ 0. \end{bmatrix}}$$

$$0.49 \begin{bmatrix} 0.31 & 0.36 & 0.19 & 0.72 \\ 0.1519 & 0.1764 & 0.0931 & 0.3528 \end{bmatrix}$$

$\Sigma^2$ )

$$\begin{bmatrix} 0.3391 & 0.608 & 0.146 & 0.5504 \end{bmatrix}$$

$$q_2 = \underline{\underline{q_2}} \\ q_2 \cdot k_1 = [0.1 \quad 0.14 \quad 0.86 \quad 0.77] \begin{bmatrix} 0.31 & 0.84 & 0.963 \\ 0.57 \end{bmatrix}$$

$$q_2 \cdot k_1 = 0.031 + 0.1176 + 0.828 + 0.4389$$

$$q_2 \cdot K = 1.42$$

divide  
by  
2

$$\therefore 0.71$$

$$q_2 \cdot k_2 = [0.1 \quad 0.14 \quad 0.86 \quad 0.77] \begin{bmatrix} 0.45 & 0.94 & 0.73 \end{bmatrix}$$

$$q_2 \cdot k_2 = [0.045 + 0.13 + 0.63 + 0.45]$$

$$q_2 \cdot k_2 = 1.25$$

$$q_2 \cdot k_2 \Rightarrow 0.625$$

$$\text{softmax} = 0.5203$$

$$\Rightarrow 0.52$$

$$0.4797$$

$$\Rightarrow 0.48$$

$$0.52 \times [0.36, \quad 0.83, \quad \cancel{0.22}, \quad 0.38]$$

$$= \begin{bmatrix} 0.1872 & \cancel{0.4316} & 0.052 & 0.1976 \end{bmatrix}$$

$$0.48 \times [0.31 \quad 0.86 \quad 0.19 \quad 0.72]$$

$$\Rightarrow [0.1488 \quad 0.1728 \quad 0.0912 \quad 0.3456]$$

$\text{Enc}_2 = [0.336 \quad 0.60 \quad 0.14 \quad 0.54]$ .  
so  $\text{Enc}_2 = [0.336 \quad 0.60 \quad 0.14 \quad 0.54]$  is step

$$z_2 = [0.30 \quad 0.59 \quad 0.14 \quad 0.54].$$

Whenever we want to concatenate bring MLP,

If we transform classification we will take out p  
the output of the encoder and then we will

feed to  $xg_{\text{boost}}$  or SVM.

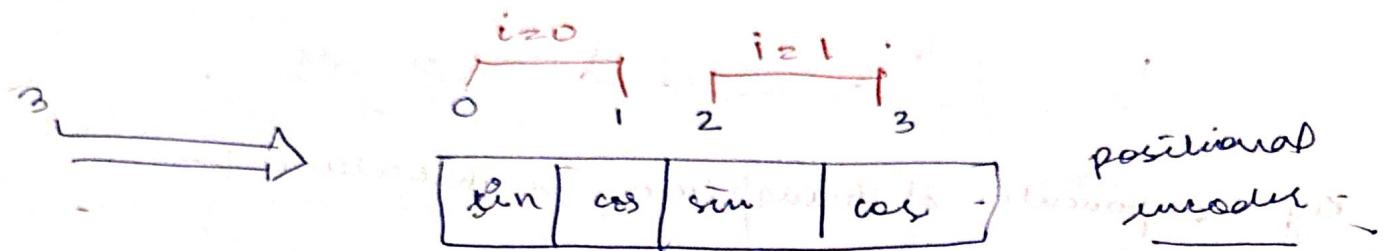
If we want to do transformer translation we will  
feed the output of the encoder to the decoder -

$\Rightarrow t=4$  Positional embedding vector  $P_E \in \mathbb{R}^{512}$

$\Rightarrow$   $P_E = [P_{00}, P_{01}, P_{02}, P_{03}, P_{10}, P_{11}, P_{12}, P_{13}, P_{20}, P_{21}, P_{22}, P_{23}, P_{30}, P_{31}, P_{32}, P_{33}]$

For  $i$  all even function use sin function

For  $i$  all odd function use cos function



$$PE_{3, 2i} =$$

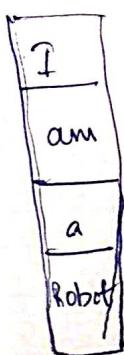
$$PE_{\text{pos}, 2i} = \sin \left( \frac{\text{pos}}{10000^{\frac{2i}{d}}} \right)$$

$$PE_{\text{pos}, 2i+1} = \cos \left( \frac{\text{pos}}{10000^{\frac{2i+1}{d}}} \right)$$

If we have d vector  
 $0 \leq i < d/2$ .

I am a robot.

$$P.E : d=4 \Rightarrow \mathbb{R}^4$$



- $\rightarrow 0 \rightarrow [P_{00} | P_{01} | P_{02} | P_{03}]$
- $\rightarrow 1 \rightarrow [P_{10} | P_{11} | P_{12} | P_{13}]$
- $\rightarrow 2 \rightarrow [P_{20} | P_{21} | P_{22} | P_{23}]$
- $\rightarrow 3 \rightarrow [P_{30} | P_{31} | P_{32} | P_{33}]$

$$P_{\text{pos}, 2i} = \sin\left(\frac{\text{Pos}}{100} \times 2i/\pi\right) \rightarrow \sin\left(\frac{2i}{100}\right)$$

only for this case.

$$\Rightarrow \sin\left(\frac{0}{100} \times 2 \times 0/\pi\right) = \sin 0.$$

$$\Rightarrow \cos 0.$$

$$= [\dots] + [\dots] + [\dots] + [\dots]$$

$\sin(0)$	$\cos(0)$	$\sin(0)$	$\cos(0)$
0	1	0	1

$$P_{1,20} = \sin\left(\frac{1}{0}\right) = \sin 0.$$

<del>0.998</del>	0.998	<del>0.001</del>	0.99
0.01	1	-	1

$$P_{1,4} = \cos 0.$$

$$P_{1,2} = \sin\left(\frac{1}{100} \times 2/4\pi^2\right).$$

$$P_{1,3} = \sin\frac{1}{10}$$

$$\sin 0.01$$

$$\Rightarrow 0.099.$$

$$\Rightarrow 0.01.$$

$$P_{1,4} = \cos 0.$$

$$\Rightarrow 0.99 \Rightarrow 1.$$