

$$\frac{\partial f_1}{\partial w_0} = x_0$$

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

$$\frac{\partial f_1}{\partial x_0} = w_{o_1} \quad (2)$$

$$\frac{\partial f_4}{\partial w_2} = 1$$

$$\frac{\partial f_2}{\partial w_1} = x_1 \quad (2)$$

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

$$\frac{\partial f_2}{\partial x_1} = \omega_1$$

$$\frac{\partial f_6}{\partial f_5} = 0.37$$

$$\frac{\partial f_3}{\partial f_1} = v$$

$$\frac{\partial f}{\partial f_6} = -1$$

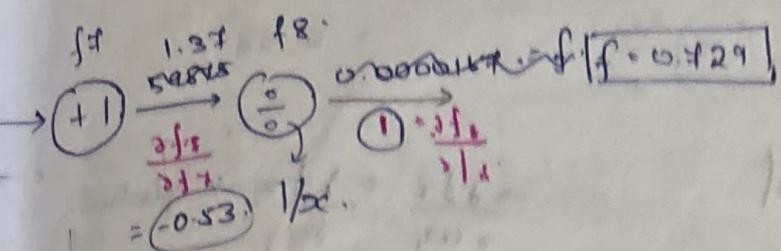
$$\frac{\partial f_3}{\partial f_2} = 1$$

$$\frac{25}{27} \times \left(\frac{-1}{1+T}\right)^2 = \left(\frac{-1}{1+3T}\right)^2 = -0.5^{33}$$

$$\frac{tfs}{2fs} = 1$$

$$\rightarrow \frac{\partial f_8}{\partial f_7} = 1.37$$

$$\rightarrow \frac{\partial f_8}{\partial f_7} + \partial \frac{1}{f_7}$$



equations

$$w_0 x_0 = f_1$$

$$w_1 x_1 = f_2$$

$$f_1 + f_2 = f_3$$

$$f_3 + w_2 \cdot f_4$$

$$1 + f_4 \cdot f_5 = (-1) + f_4 = f_5$$

$$(-1) + f_5 \cdot f_6 = e^{f_5} = f_6$$

$$e^{f_6} \cdot 1 = 1 + f_6 = f_7$$

$$1 + f_7 = f_8$$

$f_1 = -2$	$f_5 = -1$
$f_2 = +6$	$f_6 = 0.37$
$f_3 = -8.4$	$f_7 = 1.37$
$f_4 = -0.429$	$f_8 = 0.429$

Backward propagation

$$\rightarrow \frac{\partial f_8}{\partial f_7} = \text{local} * \text{global}$$

$$= \frac{-1}{(f_7)^2} * 1$$

$$= -0.533$$

$$\rightarrow \frac{\partial f_8}{\partial f_6} = \text{local} * \text{global}$$

$$= \frac{\partial f_7}{\partial f_6} * (-0.53)$$

$$= 1 * (-0.53)$$

$$= -0.53$$

$$\rightarrow \frac{\partial f_8}{\partial f_5} = \frac{\partial f_6}{\partial f_5} \times \text{global}$$

$$= 0.37 \times (-0.53)$$

$$= -0.19$$

$$\rightarrow \frac{\partial f_8}{\partial x_0} = \frac{\partial f_6}{\partial x_0} \times \text{global}$$

$$= w \cdot 2 \times 0.19$$

$$= 0.38$$

$$\rightarrow \frac{\partial f_8}{\partial f_4} = \frac{\partial f_5}{\partial f_4} \times \text{global}$$

$$= (-) \times -0.19$$

$$= 0.19$$

$$\rightarrow \frac{\partial f_8}{\partial w_1} = \frac{\partial f_2}{\partial w_1} \times \text{global}$$

$$= -2 \times 0.19$$

$$= -0.38$$

$$\rightarrow \frac{\partial f_8}{\partial w_2} = \frac{\partial f_4}{\partial w_2} \times \text{global}$$

$$= 0.19$$

$$\rightarrow \frac{\partial f_8}{\partial x_1} = \frac{\partial f_2}{\partial x_1} + \frac{\partial f_2}{\partial x_1} \times \text{global}$$

$$= -3 \times 0.19$$

$$= -0.57$$

$$\rightarrow \frac{\partial f_8}{\partial f_3} = \frac{\partial f_4}{\partial f_3} \times \text{global}$$

$$= 1 \times 0.19$$

$$= 0.19$$

$$\rightarrow \frac{\partial f_8}{\partial f_2} = \frac{\partial f_3}{\partial f_2} \times \text{global}$$

$$= 1 \times 0.19$$

$$= 0.19$$

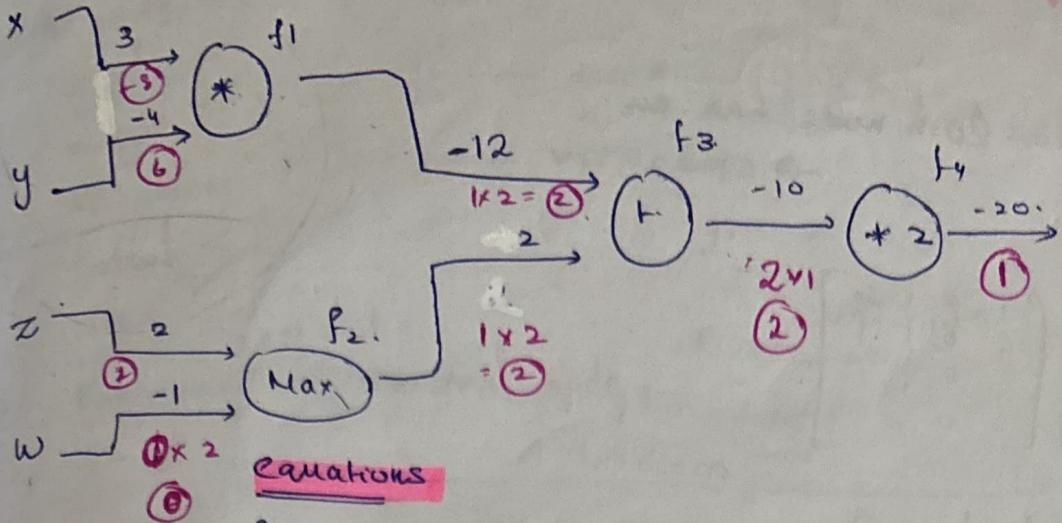
$$\rightarrow \frac{\partial f_8}{\partial f_1} = \frac{\partial f_3}{\partial f_1} \times \text{global}$$

$$\text{global} = \frac{\partial f_6}{\partial f_1} \leftarrow 0.19$$

$$\rightarrow \frac{\partial f_8}{\partial w_0} = \frac{\partial f_1}{\partial w_0} \times \text{global}$$

$$= (-) \times 0.19$$

$$= -0.19$$



$$f_1 = x * y$$

$$f_2 = z + \omega \cdot \max(z, w)$$

$$f_3 = f_1 + f_2.$$

$$f_4 = f_3 * 2.$$

$$x = 3$$

$$y = -4$$

$$z = 2$$

$$\omega = -1.$$

$$f_2 = \begin{cases} z & z > \omega \\ \omega & z < \omega \\ \frac{\partial f_2}{\partial z} & z = \omega \end{cases}$$

$$\frac{\partial f_1}{\partial x} = y ; \quad \frac{\partial f_1}{\partial y} = x = 3.$$

$$\frac{\partial f_1}{\partial z} = -4 \quad \frac{\partial f_1}{\partial \omega} = 0$$

$$\frac{\partial f_2}{\partial z} = 1 ; \quad \frac{\partial f_2}{\partial \omega} = 0 \quad f_2 = \begin{cases} z & z > \omega \\ \omega & z < \omega \end{cases}$$

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

$$\frac{\partial f_4}{\partial f_3} = 2.$$

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

$$\frac{\partial f_2}{\partial z}$$

$$f_2' = \begin{cases} 1 & z > \omega \\ 0 & z < \omega \end{cases}$$

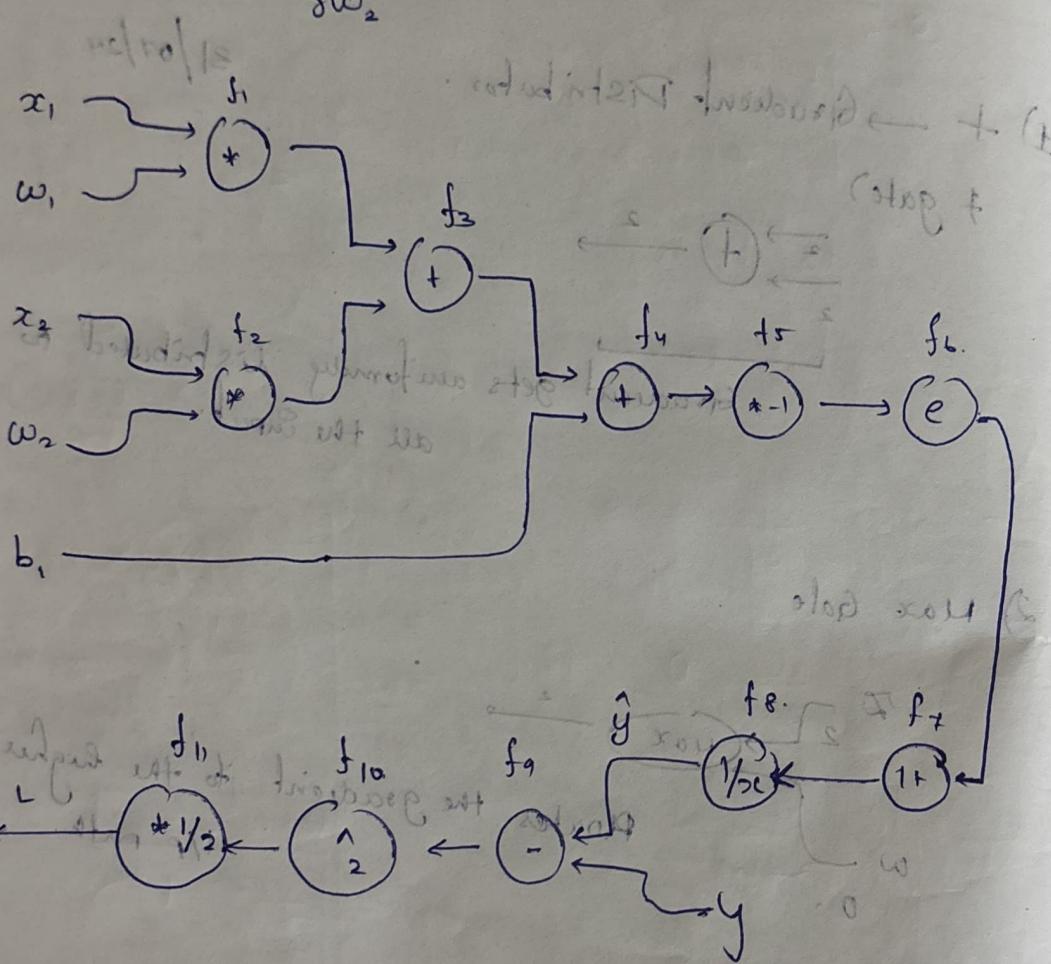
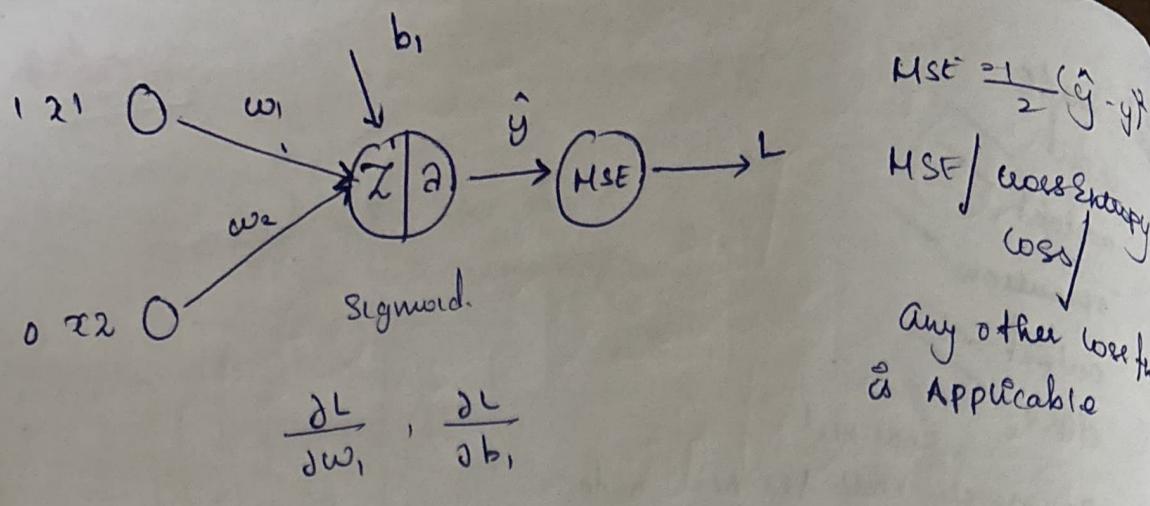
arbitrarily choose
Since $z > \omega$ ①

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

$$\frac{\partial f_2}{\partial \omega} = 0.$$

$$\frac{\partial f_2}{\partial \omega}$$

$$f_2' = \begin{cases} 0 & z > \omega \\ 1 & z < \omega \end{cases}$$

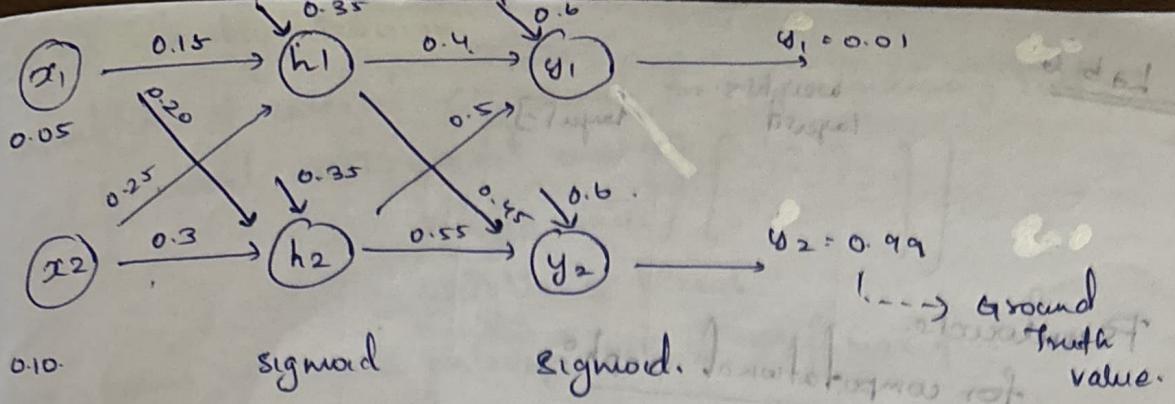


Note:- In future we can replace the whole

represented in sigmoid eqn with one sigmoid gate

cause we know the derivative of sigmoid

function



Two output

$$\text{Loss} = \sum \frac{1}{2} (\text{target} - \text{output})^2$$

for all neurons

Shows losses w.r.t. target values

$$h_1 = x_1 w_{11} + x_2 w_{12} + b_1$$

$$h_2 = x_1 w_{21} + x_2 w_{22} + b_2$$

$$y_1 = h_1 w_{11}^{(1)} + h_2 w_{12}^{(1)} + b_1^{(1)}$$

$$y_2 = h_1 w_{21}^{(1)} + h_2 w_{22}^{(1)} + b_2^{(1)}$$

$$\begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} = \mathbf{w}^{(1)}$$

$$\mathbf{w}^{(1)} = \begin{bmatrix} 0.15 & 0.25 \\ 0.2 & 0.3 \end{bmatrix}$$

$$\mathbf{w}^{(2)} = \begin{bmatrix} 0.4 & 0.5 \\ 0.45 & 0.55 \end{bmatrix}$$

$$s_1 = \text{Sigmoid}(h_1)$$

$$s_2 = \text{Sigmoid}(h_2)$$

$$b^{(1)} = \begin{bmatrix} 0.35 \\ 0.35 \end{bmatrix}$$

$$b^{(2)} = \begin{bmatrix} 0.6 \\ 0.6 \end{bmatrix}$$

$$h_1 = 0.05 \times 0.15 + 0.10 \times 0.25 + 0.35$$

$$= 0.3825$$

$$h_2 = 0.05 \times 0.20 + 0.10 \times 0.3 + 0.35$$

$$= 0.39$$

$$s_1 = \frac{1}{1 + e^{-0.3825}} = 0.595$$

$$s_2 = \frac{1}{1 + e^{-0.39}} = 0.596$$

$$y_1 = 0.595 \times 0.4 + 0.596 \times 0.5 + 0.6$$

$$y_1 = 1.136$$

$$y_2 = 0.595 \times 0.45 + 0.596 \times 0.55 + 0.6$$

$$y_2 = 1.195$$

$$s_1 = 0.757$$

$$s_2 = 0.767$$

Example :-

$$x \rightarrow z_1(a_1) \xrightarrow{0.5} z_2(a_2) \xrightarrow{0.5} z_3(a_3) \xrightarrow{0.5} z_4(a_4)$$

-0.0125 -0.105 -0.572 -0.428

$$z_1 = 0.5x = 0.5$$

$$a_1 = \sigma(z_1)$$

$$= \sigma(0.5x)$$

$$= \sigma(0.5)$$

$$= \frac{1}{1 + e^{-0.5}} = 0.622$$

$$z_2 = 0.622 \times 0.5$$

$$= 0.311$$

$$a_2 = \sigma(z_2)$$

$$= \frac{1}{1 + e^{-0.311}}$$

$$(a_2)$$

$$z_3 = 0.577 \times 0.5 = 0.288$$

$$a_3 = \frac{1}{1 + e^{-0.288}} = 0.572$$

$$= \frac{1556}{1 + e^{-0.288}}$$

need to do

with

$$\frac{\partial L}{\partial a_3} \times \frac{\partial a_3}{\partial z_3}$$

Exploding Gradients



by DRC

solved by clipping the gradient

[basically thresholding]

ground truth
 $y = 1$

$$\frac{\partial a_3}{\partial a_2} = 1$$

$$\frac{\partial a_3}{\partial z_3} = \left(\frac{1}{1 + e^{-z_3}} \right) \left(1 - \frac{1}{1 + e^{-z_3}} \right)$$

$$= 2.45 \left(1 - 0.572 \right)$$

$$L = \frac{1}{2} (a_3 - y)^2$$

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

$$\textcircled{2} \quad \frac{\partial L}{\partial z_3} = \frac{\partial}{\partial z_3} \left(\frac{1}{1 + e^{-z_3}} - y \right)^2 \\ = \left(\frac{1}{1 + e^{-z_3}} - y \right) 2 \cdot \frac{1}{1 + e^{-z_3}} \cdot (1 - \frac{1}{1 + e^{-z_3}})$$

$$= \left(\frac{1}{1 + e^{-0.288}} - 1 \right) 2 \cdot \frac{1}{1 + e^{-0.288}} \cdot (1 - \frac{1}{1 + e^{-0.288}})$$

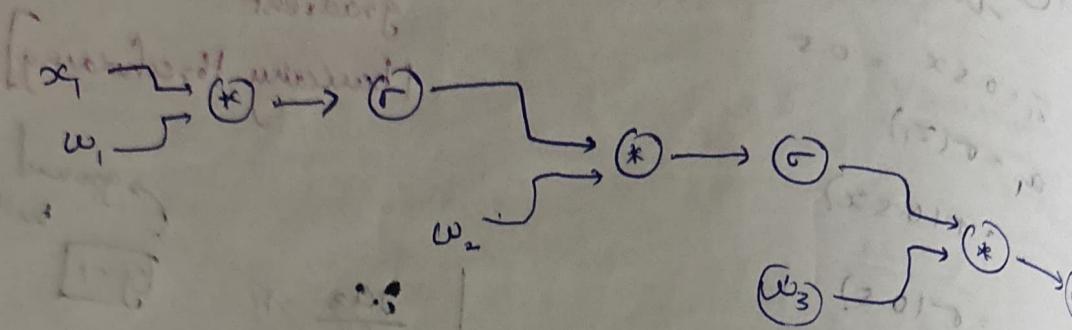
$$= \frac{1}{1 + e^{-0.288}} - 1$$

$$= (0.572 - 1)(0.572 + 1)(1 - 0.572)$$

$$= -0.105$$

$$③ \frac{\partial L}{\partial a_2} = \frac{\partial L}{\partial a_3} \frac{\partial a_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial a_2} = -0.105 \times \frac{\partial 0.5a_2}{\partial a_2}$$

→ getting stuck, we computational graph



$$\frac{\partial L}{\partial z_2} = \frac{\partial L}{\partial a_3} \frac{\partial a_3}{\partial z_3} \frac{\partial z_3}{\partial a_2} \frac{\partial a_2}{\partial z_2}$$

$$= -0.0525 \times \frac{1}{1+e^{-z_2}}$$

$$= -0.0525 \times g(z_2)(1-g(z_2))$$

$$= -0.0525 \times 0.577 \times 0.423$$

$$= -0.0128 = \underline{-0.013}$$

$$\frac{\partial L}{\partial a_1} = \frac{\partial L}{\partial a_3} \frac{\partial a_3}{\partial z_3} \frac{\partial z_3}{\partial a_2} \frac{\partial a_2}{\partial z_2} \frac{\partial z_2}{\partial a_1}$$

$$= -0.013 \times 0.5 = \underline{-0.0064}$$

$$\frac{\partial L}{\partial z_1} = -0.0064 \times \frac{\partial a_1}{\partial z_1}$$

$$= -0.0064 \times g(z_1)(1-g(z_1))$$

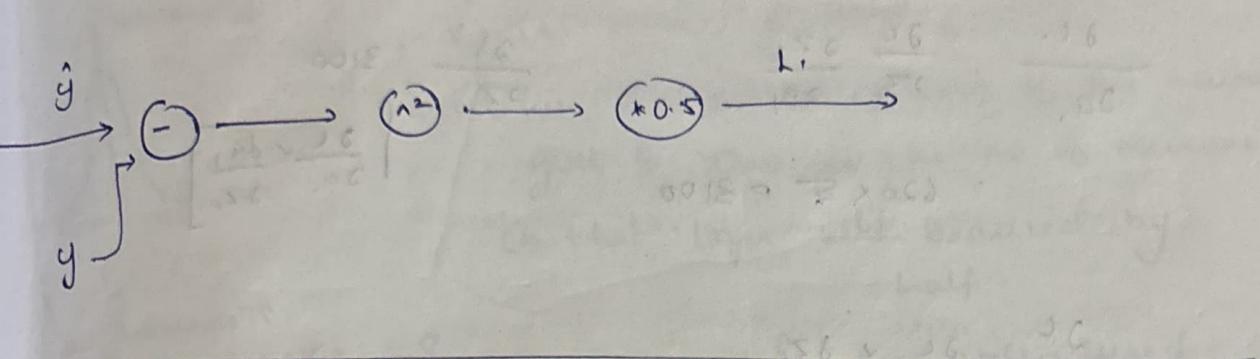
$$= -0.0064 \times 0.622 \times 0.378$$

$$= \underline{-0.001}$$

$$\frac{\partial L}{\partial x} = -0.0015 \times \frac{\partial z_1}{\partial x}$$

$$= -0.0015 \times 0.5$$

$$= \underline{-0.00075} \rightarrow \text{Vernachl. / gradients.}$$



Exploding gradients

$$\frac{\partial L}{\partial z_3} = 125$$

$$\frac{\partial L}{\partial z_2} = 625$$

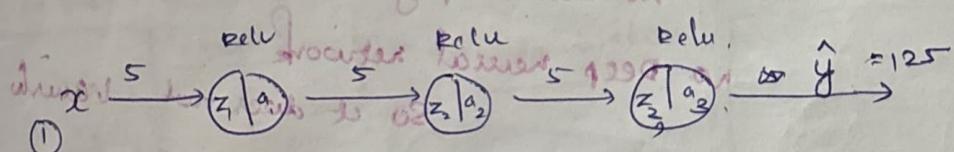
$$\frac{\partial L}{\partial z_1} = 3125$$

$$z_1 = 5x = 5$$

$$a_1 = \text{ReLU}(z_1) = 5$$

$$z_2 = 5 \cdot 5 = 25$$

$$a_2 = \text{ReLU}(z_2) = 25$$



$$z_1 = 5x = 5$$

$$a_1 = \text{ReLU}(5) = 5$$

$$z_2 = 5 \cdot 5 = 25$$

$$a_2 = \text{ReLU}(25) = 25$$

$$z_3 = 25 \cdot 5 = 125$$

$$a_3 = \text{ReLU}(125) = 125$$

$$= 25^4$$

$$\boxed{a_3 = y = 125}$$

$$L = \frac{1}{2} (\hat{y} - y)^2 = \frac{1}{2} (125 - 1)^2 = 7688$$

$$\frac{\partial L}{\partial a_3} \approx (a_3 - y) = 124$$

$$\frac{\partial L}{\partial z_3} = \frac{\partial L}{\partial a_3} \frac{\partial a_3}{\partial z_3} = 124 \cdot 0.25(1-0.25) \cdot 124 \frac{\partial \text{max}(0, z_3)}{\partial z_3} = 124 \cdot 0.25(1-0.25) \cdot 124 = 124$$

$$\frac{\partial L}{\partial a_2} = \frac{\partial L}{\partial z_3} \times \frac{\partial z_3}{\partial a_2}$$

$$= 124 \times 5$$

$$\therefore 620 \times \frac{\max(a_{z_2})}{\partial z_2}$$

Explanation: $\frac{\partial L}{\partial z_2} = 620$

$$\frac{\partial L}{\partial z_2} = 620 \times 1$$

$$\frac{\partial L}{\partial a_1} = \frac{\partial L}{\partial z_2} \frac{\partial z_2}{\partial a_1}$$

$$\therefore 620 \times 5 = 3100$$

$$\frac{\partial L}{\partial z_1} = 3100$$

$$\left| \frac{\partial L}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \right|$$

~~$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z_1} \times \frac{\partial z_1}{\partial x}$$~~

$$= 3100 \times 5 = \underline{15500}$$

Explanation: exploding gradient

Even if exploding gradient we prefer ReLU
 because it's easier to handle.

Why not found in ML
 no backpropagation

no Deep neural network

so it does not vanish

$$\|g\| = \sqrt{\sum g_i^2} \quad \text{if } \|g\| > c$$

$$g_{\text{clipped}} = g \times \frac{c}{\|g\|}$$

14/08/25 → Theory

[Theory]

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

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

$$P = 1$$

$$S = 2$$

$$n = 5 \times 5$$

①.

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

$$(-2)$$

②

$$(-7)$$

③

$$(-6)$$

④

$$-6 - 3$$

$$(-9)$$

⑤

$$6 - 9 - 8$$

$$-3 - 8$$

$$(-11)$$

⑥

$$8 - 7 - 8$$

$$(-7)$$

$$(F)$$

$$-7 - 4$$

$$(-15)$$

$$(8)$$

$$8 - 3$$

$$-1$$

$$(4)$$

$$(9) - 6 - 6 - 3$$

$$(-3)$$

(11)

outside

14/08/25 [Theory]

*

1	-1
0	1

$$f = 2\sqrt{2}$$

⇒

-2	-7	-6
-9	-11	-7
-11	4	-3

3x3

floor

$$5+2-2+1 = 6/2$$

but I need to get 3x3

3x3

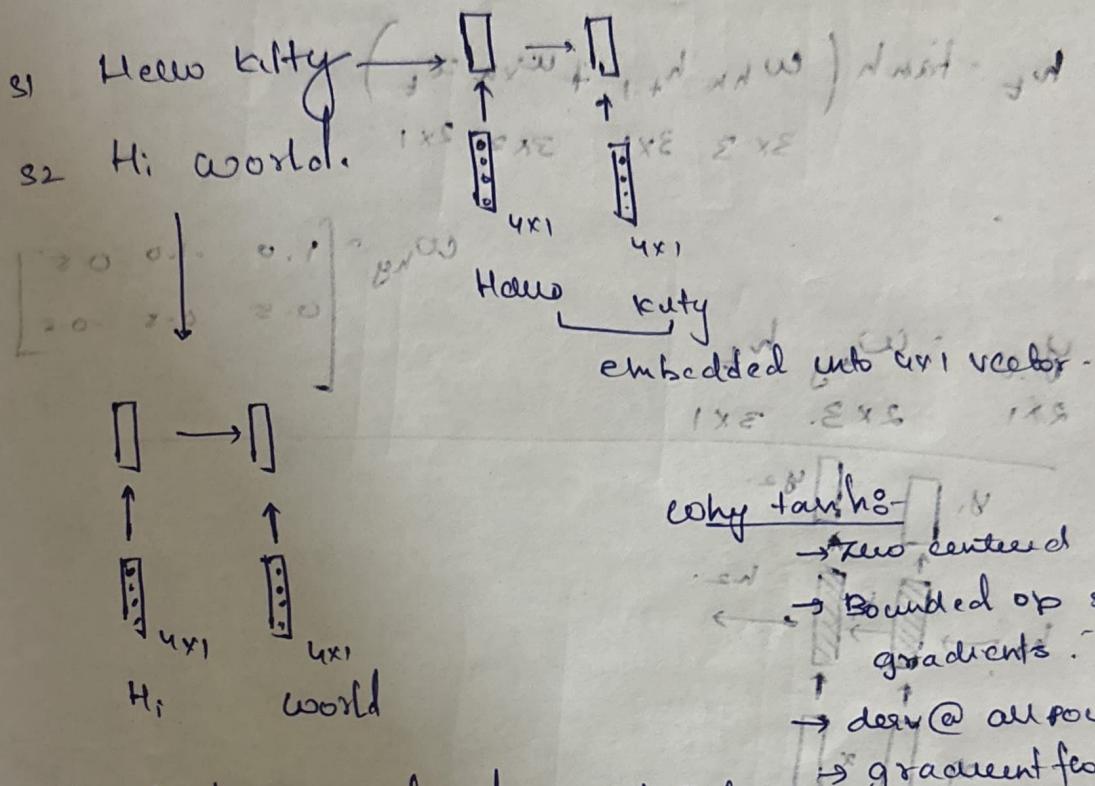
↓

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

b/p.

Note :- $w_{kh}, w_{ch} \rightarrow$ Shared across the timesteps.

Note :- I/P change of length (no of words changing in each sample) \rightarrow handled by Recurrence by RNN block.



Why tanh as activation not ReLU?

I/P to RNN always a vector, use a ResNet to get vector \vec{s}_t then do RNN.

Character level Language Model
Predict the next character.

Vanilla RNN.

$T=2$ (two time steps)

$$X^{(1)}: x_1 \in \mathbb{R}^2 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}; x_2 \in \mathbb{R}^2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

Input dim = 2. $\rightarrow x_t \in \mathbb{R}^2$

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

$$h_0 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

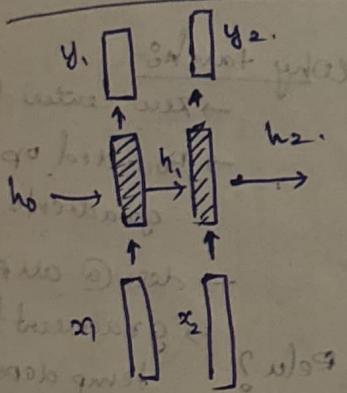
$$y_6 \in \mathbb{R}^2$$

$$w_{2h} = \begin{bmatrix} 0.5 & -0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.4 \end{bmatrix}_{3 \times 2}, w_{hh} = \begin{bmatrix} 0.1 & 0.4 & 0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.2 \end{bmatrix}_{3 \times 3}$$

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

$$y_t = w_{hy} h_t$$

$$w_{hy} = \begin{bmatrix} 1.0 & -1.0 & 0.5 \\ 0.5 & 0.5 & -0.5 \end{bmatrix}_{2 \times 3}$$



① Forward Pass (output is episode 417 of 91)

$$h_t = \tanh \left(\begin{bmatrix} 0.1 & 0.4 & 0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.2 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right)$$

$$\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0.1 & 0.4 & 0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.2 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$w_t = \tanh \left([0 \ 0 \ 0] + \begin{bmatrix} 0.5 - 0.6 & 0.8 + 0.4 & 0.1 - 0.8 \\ 0.1 + 0.8 & 0.1 + 0.8 & 0.1 + 0.8 \end{bmatrix} \right)$$

$$h_+ = \tanh \left([-0.1 \ 0.2 \ 0.2] \right)$$

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

$$= \frac{1 - e^{-2x}}{1 + e^{-2x}}$$

$$\textcircled{1} = \frac{-0.221}{2.221} = -0.099$$

$$\textcircled{2} = \frac{0.551}{1.449} = 0.380 \quad 0.83$$

$$\textcircled{3} = -0.604 \quad 0.716$$

$$h_{\pm} = \begin{bmatrix} -0.099 \\ 0.380 \\ 0.83 \\ -0.604 \\ 0.716 \end{bmatrix}_{3 \times 1}$$

$$y_2 = \begin{bmatrix} 1 & -1 & 0.5 \\ 0.5 & 0.5 & -0.5 \end{bmatrix} \begin{bmatrix} -0.099 \\ 0.380 \\ 0.83 \\ -0.604 \\ 0.716 \end{bmatrix}$$

$$= \begin{bmatrix} -0.099 - 0.380 - 0.302 \\ -0.045 - 0.19 + 0.302 \end{bmatrix}$$

$$= \begin{bmatrix} -0.571 \\ -0.481 \\ 0.064 \end{bmatrix}$$

$$y_1 = \begin{bmatrix} -0.571 \\ -0.481 \\ 0.064 \\ 0.0045 \end{bmatrix}_{2 \times 1}$$

$$h_t = \tanh \left(w_{nh} h_{t-1} + w_{xh} y_{t-1} \right)$$

$$h_2 = \tanh \left(w_{nh} h_1 + w_{xh} y_1 \right)$$

$$\Rightarrow \tanh \left(\begin{bmatrix} 0.1 & 0.4 & 0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.2 \end{bmatrix} \begin{bmatrix} -0.699 \\ 0.83 \\ 0.715 \end{bmatrix} \right)$$

$$+ \begin{bmatrix} 0.5 & 0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.4 \end{bmatrix} \begin{bmatrix} 186.571 \\ 188.8 \\ 192.0 \end{bmatrix}$$

$$= \tanh \left(\begin{bmatrix} 40.3221 \\ 0.412 \\ 0.05525 \end{bmatrix} + \begin{bmatrix} -0.28775 \\ -0.4553 \\ -0.0541 \end{bmatrix} \right)$$

$$\tanh \left(\begin{bmatrix} 0.06435 \\ 0.080 \\ -0.0433 \\ 0.00015 \end{bmatrix} \right) = \begin{bmatrix} 1 & 1 & 1 \\ 2.0 & 2.0 & 2.0 \end{bmatrix} \quad \checkmark$$

$$h_2 = \begin{bmatrix} 0.0643 \\ 0.080 \\ -0.043 \\ 0.00015 \end{bmatrix} \cdot b_2 \cdot \begin{bmatrix} -0.4434 \\ 0.1849 \\ 0.00340 \end{bmatrix}$$

$$186.571 \begin{bmatrix} 8.0 \\ 8.0 \\ 8.0 \end{bmatrix} = 18$$

$$y_2 = \begin{bmatrix} -0.0888 \\ -0.9844 \end{bmatrix}$$

$$w_t = \tanh \left(\begin{bmatrix} 0.1 & 0.4 & 0 \\ -0.2 & 0.3 & 0.2 \\ 0.65 & 0.05 & 0.7 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0.5 & 0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.4 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} \right)$$

add rows good bad rows of inputs

$$= \tanh \left(\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 6.1 \\ 1.2 \\ 0.9 \end{bmatrix} \right)$$

$$h_1 = \begin{bmatrix} -0.099 \\ 0.83 \\ 0.716 \end{bmatrix}$$

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

$$g_1 = \begin{bmatrix} 1 & -1 & 0.5 \\ 0.5 & 0.5 & -0.5 \end{bmatrix} \begin{bmatrix} -0.099 \\ 0.83 \\ 0.716 \end{bmatrix}$$

$$y_1 = \begin{bmatrix} -0.571 \\ 0.0075 \end{bmatrix}_{2 \times 1}$$

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

$$= \tanh \left(\begin{bmatrix} 0.1 & 0.4 & 0 \\ -0.2 & 0.3 & 0.2 \\ 0.65 & 0.05 & 0.7 \end{bmatrix} \begin{bmatrix} -0.099 \\ 0.83 \\ 0.716 \end{bmatrix} + \begin{bmatrix} 0.5 & 0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.4 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} \right)$$

feature space no bias in $\begin{bmatrix} -1 \\ 1 \\ 3 \end{bmatrix}$

$$= \tanh \left(\begin{bmatrix} 0.3221 \\ 0.42 \\ 0.0107 \end{bmatrix} + \begin{bmatrix} -0.8 \\ -0.6 \\ 0.3 \end{bmatrix} \right)$$

$$= \tanh \left(\begin{bmatrix} -0.4479 \\ 0.1821 \\ -0.22 \\ 0.8107 \end{bmatrix} \right) = \begin{bmatrix} 0.2149 & -0.4434 \\ -0.1847 & 0.801 \end{bmatrix}$$

24/09/2025

Note :-

LSTM minimises vanishing gradient but cannot completely eliminate vanishing gradient.

Why LSTM minimises vanishing gradient?

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

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

$$c_{t-1} = [0.2, -0.2] \quad \text{of norm = 0.2}$$

$$w_{xg} = \begin{bmatrix} 0.5 & -0.3 \\ 0.4 & 0.1 \end{bmatrix} \quad \text{of norm = 0.7}$$

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

$$w_{xf} = \begin{bmatrix} -0.4 & 0.2 \\ 0.3 & 0.3 \end{bmatrix} \quad \text{(forget gate bias)}$$

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

$$w_{xo} = \begin{bmatrix} 0.3 & 0.2 \\ -0.2 & 0.2 \end{bmatrix}$$

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

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

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

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

$$\begin{pmatrix} i_0 \\ i_1 \end{pmatrix} = \sigma \left(\begin{pmatrix} 0.1 & 0.2 \\ -0.2 & 0.05 \end{pmatrix} \begin{pmatrix} 0 & 0.1 \\ 0 & 0.2 \end{pmatrix}^T + \begin{pmatrix} 0.5 & -0.3 \\ 0.4 & 0.1 \end{pmatrix} \begin{pmatrix} 0.5 & -0.1 \\ 1.0 & 1.0 \end{pmatrix} \right)$$

$$= \sigma \left(\begin{pmatrix} 0.1 + 0.02 \\ 0 + 0.005 \end{pmatrix}, \begin{pmatrix} 0.25 + 0.03 \\ 0.2 - 0.01 \end{pmatrix} \right) =$$

$$= \sigma \left(\begin{pmatrix} 0.02 \\ 0.005 \end{pmatrix}, \begin{pmatrix} 0.28 \\ 0.19 \end{pmatrix} \right) + \begin{pmatrix} 0.01 \\ 0.01 \end{pmatrix} =$$

$$= \sigma \left(\begin{pmatrix} 0.207 \\ 0.195 \end{pmatrix}, \begin{pmatrix} 0.57 \\ 0.55 \end{pmatrix} \right) =$$

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

$$= \tanh \left(\begin{pmatrix} 0.2 & 0.1 \\ -0.1 & 0.05 \end{pmatrix} \begin{pmatrix} 0 \\ 0.1 \end{pmatrix} + \begin{pmatrix} 0.5 & -0.3 \\ 0.2 & -0.3 \end{pmatrix} \begin{pmatrix} 0.5 \\ -0.1 \end{pmatrix} \right)$$

$$= \tanh \left(\begin{pmatrix} 0.01 \\ 0.005 \end{pmatrix} + \begin{pmatrix} -0.25 - 0.04 \\ 0.1 + 0.03 \end{pmatrix} \right) =$$

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

$$= \tanh \left(\begin{pmatrix} -0.28 \\ 0.135 \end{pmatrix} \right) = \begin{pmatrix} 0.27 \\ -0.14 \\ 0.087 \\ 0.13 \end{pmatrix}$$

$$\begin{aligned}
 f(t) &= \Gamma (\omega_{hf} h_{t-1} + \omega_{zf} z_t^0) + \epsilon_t + w(t) \\
 &= \Gamma \left(\begin{bmatrix} 0.05 \\ -0.1 \\ 0.2 \end{bmatrix} + \begin{bmatrix} 0 \\ 0.1 \end{bmatrix} \right) + \begin{bmatrix} -0.45 \\ 0.2 \\ 0.3 \end{bmatrix} \begin{bmatrix} 0.2 \\ 0.3 \end{bmatrix} \\
 &= \Gamma \left(\begin{bmatrix} 0.05 + 0.0 \\ -0.1 + 0.1 \\ 0.2 + 0.3 \end{bmatrix} \right) + \begin{bmatrix} 0.01 \\ -0.2 \\ 0.15 \end{bmatrix} \begin{bmatrix} 0.02 \\ 0.02 \end{bmatrix} \\
 &= \Gamma \left(\begin{bmatrix} 0.05 \\ 0 \\ 0.5 \end{bmatrix} + \begin{bmatrix} 0.02 \\ 0.22 \\ 0.12 \end{bmatrix} \right) + \begin{bmatrix} 0.02 \\ 0.02 \\ 0.08 \end{bmatrix} \\
 &= \Gamma \left(\begin{bmatrix} 0.07 \\ 0.23 \\ 0.13 \end{bmatrix} \right) + \epsilon_t + w(t) + v(t)
 \end{aligned}$$

$$\begin{aligned}
 \left(\begin{bmatrix} 0.443 \\ 0.53 \end{bmatrix} \begin{bmatrix} 10.0 \\ 20.0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} \begin{bmatrix} 1.0 \\ 0.0 \end{bmatrix} \right) \text{element wise multiplication} \\
 C(t) &= \Gamma \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix} C_{t-1} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} g_t + \right. \\
 &\quad \left. \begin{bmatrix} 10.0 \\ 20.0 \end{bmatrix} + \begin{bmatrix} 1.0 \\ 0.0 \end{bmatrix} \begin{bmatrix} 0.57 \\ 0.58 \end{bmatrix} \begin{bmatrix} -0.27 \\ 0.13 \end{bmatrix} \right) \\
 &= \begin{bmatrix} 0.443 \\ 0.53 \end{bmatrix} \begin{bmatrix} 10.0 \\ 20.0 \end{bmatrix} + \begin{bmatrix} 0.2 \\ -0.2 \end{bmatrix} + \begin{bmatrix} 10.0 \\ 20.0 \end{bmatrix} \begin{bmatrix} 0.57 \\ 0.58 \end{bmatrix} \begin{bmatrix} -0.27 \\ 0.13 \end{bmatrix} \\
 &= \begin{bmatrix} 0.0886 \\ 0.106 \end{bmatrix} + \begin{bmatrix} 10.0 \\ 20.0 \end{bmatrix} \begin{bmatrix} -0.152 \\ 0.075 \end{bmatrix} \\
 &= \begin{bmatrix} 0.0886 \\ 0.106 \\ 10.0 \\ 20.0 \end{bmatrix} + \begin{bmatrix} -0.0654 \\ -0.031 \end{bmatrix}
 \end{aligned}$$

$$o_t = \sigma(\omega_0 h_{t-1} + \omega_{x0} x_t)$$

$$o_t = \sigma \left(\begin{bmatrix} 0.15 & 0.05 \\ 0.1 & -0.2 \end{bmatrix} \begin{bmatrix} 0 \\ 0.1 \end{bmatrix} + \omega \begin{bmatrix} 0.3 & 0.2 \\ -0.2 & 0.2 \end{bmatrix} \begin{bmatrix} 0.5 \\ -0.1 \end{bmatrix} \right)$$

$$= \sigma \left(\begin{bmatrix} 0.05 \\ -0.02 \end{bmatrix} + \begin{bmatrix} 0.15 - 0.02 \\ 0.1 - 0.02 \end{bmatrix} \right)$$

$$= \sigma \left(\begin{bmatrix} 0.05 \\ -0.02 \end{bmatrix} + \begin{bmatrix} 0.13 \\ -0.12 \end{bmatrix} \right)$$

$$= \sigma \left(\begin{bmatrix} 0.18 \\ -0.14 \end{bmatrix} \right) = \begin{bmatrix} 0.55 \\ 0.46 \end{bmatrix}$$

$$h_t = o_t \tanh(o_t)$$

$$\text{available update} \\ o_t = \begin{bmatrix} 0.55 \\ 0.46 \end{bmatrix} \cdot \tanh \begin{bmatrix} -0.0654 \\ -0.0316 \end{bmatrix} = \begin{bmatrix} 0.555 \\ 0.449 \end{bmatrix} \begin{bmatrix} -0.065 \\ -0.031 \end{bmatrix}$$

① old value

$$h_t = \begin{bmatrix} -0.03575 \\ -0.014 \end{bmatrix}$$

$$\begin{bmatrix} 0.5 \\ 0.46 \end{bmatrix} +$$

$$(1-\beta) \begin{bmatrix} 0.5 \\ 0.46 \end{bmatrix} + \beta \begin{bmatrix} 0.555 \\ 0.449 \end{bmatrix} =$$

$$\beta w$$

$$\alpha w (1-\beta) + \beta p +$$

Sum

$$h_1 = [1, 0, 0]$$

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

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

t=3

$$S_{t-1} = [1, 0, 1]$$

Score fn = dot product $(1, 1, 0) \cdot (1, 1, 0)$

Cf = ?

$$C_t = \sum_{j=1}^T \alpha_{t,j} h_j$$

$$C_t = \frac{\alpha_{t,1} h_1 + \alpha_{t,2} h_2 + \alpha_{t,3} h_3}{\exp(score(S_{t-1}, h_1)) + \exp(score(S_{t-1}, h_2)) + \exp(score(S_{t-1}, h_3))}$$

$$\text{dot product} \sum_{j=1}^T \exp(score(S_{t-1}, h_j))$$

$$= \frac{\exp(score(S_{t-1}, h_3))}{\exp(score(S_{t-1}, h_1)) + \exp(score(S_{t-1}, h_2)) + \exp(score(S_{t-1}, h_3))}$$

$$+ \exp(score(S_{t-1}, h_2))$$

$$S_{t-1} \cdot h_2 = [1 \ 0 \ 1] \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}^T \cdot (([1 \ 0] \cdot w_2) \cdot v_2)$$

$$\cdot ([1 \ 0] \cdot [1 \ 0] \cdot w_2) \cdot v_2$$

$$S_{t-1} \cdot h_2 = [1 \ 0 \ 1] \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$$

Since the second row contains lots more values
softmax will focus mostly on the first two

to return with weights around 0.5

$$S_{t-1} \cdot h_2 = [1 \ 0 \ 1] \begin{bmatrix} 0 \\ 0.5 \\ 0.5 \end{bmatrix}$$

~~(2.3)~~ 2

$$[1 \ 0 \ 1] = w$$

$$[1 \ 1 \ 0] = v$$

$$d_{t+1,1} = \frac{\exp(2)}{\exp(1) + \exp(1) + \exp(2)} = \frac{0}{0.576} = 0.576$$

$$d_{t+1,2} = \frac{\exp(1)}{\exp(1) + \exp(1) + \exp(2)} = \frac{0}{0.576} = 0.212$$

$$d_{t+1,3} = 0.212$$

$C_t \rightarrow$ vector

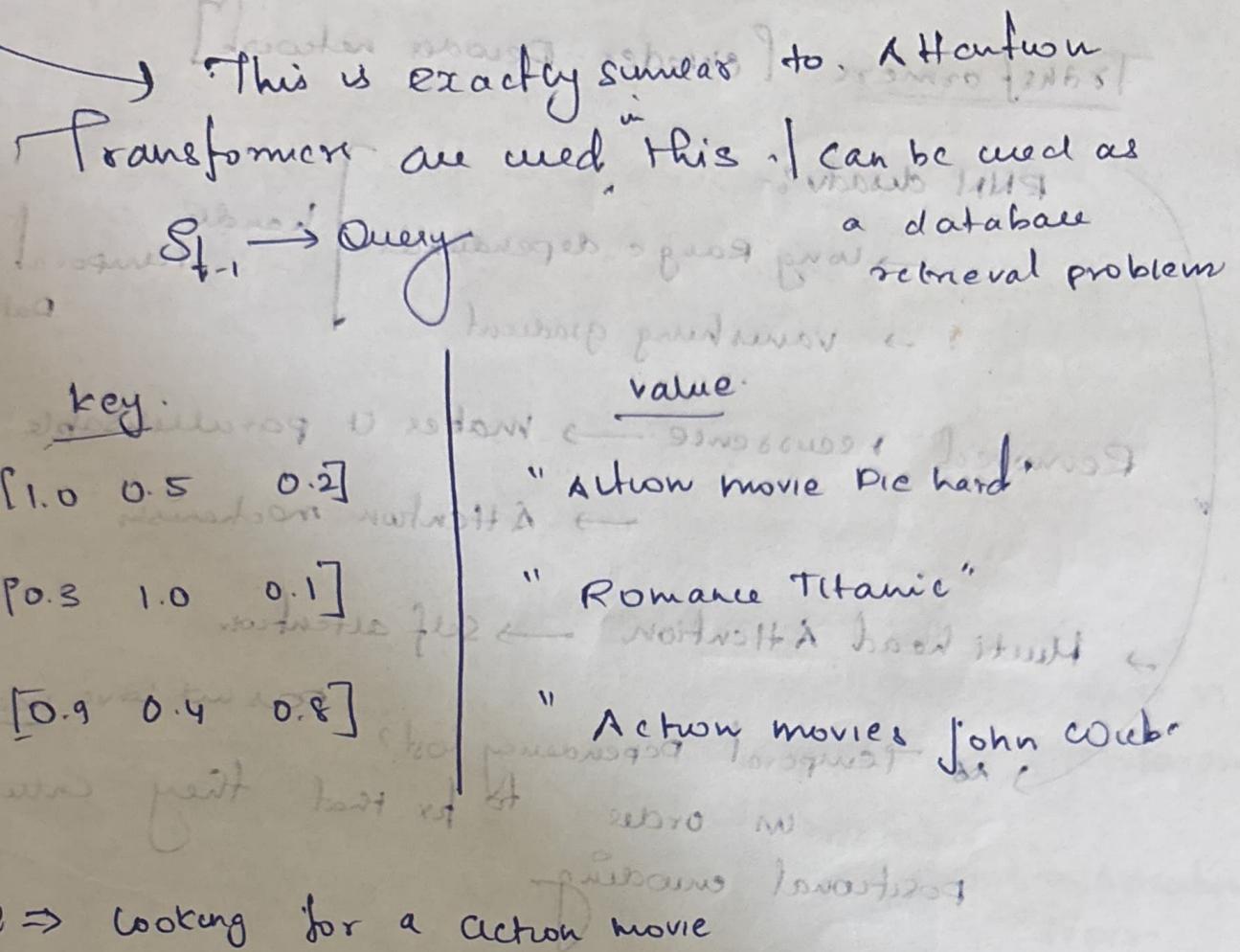
$$= \begin{pmatrix} 0.576 \\ 0 \\ 0.576 \end{pmatrix} \quad \begin{pmatrix} 0 \\ 0.212 \\ 0.212 \end{pmatrix} \quad \begin{pmatrix} 0.212 \\ 0.212 \\ 0 \end{pmatrix}$$

$$= \begin{pmatrix} 0.576 \\ 0.424 \\ 0.424 \end{pmatrix} \xrightarrow{\text{more closer to } h_1}$$

that's how the attention mechanism works.

How does chatGPT work?
How does it help?

How does it work?



$$\begin{bmatrix} 1.0 & 0.5 & 0.2 \end{bmatrix} \begin{bmatrix} 1 \\ 0.3 \\ 0.5 \end{bmatrix} = [1 + 0.15 + 0.10] = [1.25]$$

$$\begin{bmatrix} 0.3 & 1 & 0.1 \end{bmatrix} \begin{bmatrix} 1 \\ 0.3 \\ 0.5 \end{bmatrix} = [0.3 + 1 + 0.05] = [0.65]$$

$$\begin{bmatrix} 0.9 & 0.4 & 0.8 \end{bmatrix} \begin{bmatrix} 1 \\ 0.2 \\ 0.5 \end{bmatrix} = [0.9 + 0.12 + 0.4] = [1.42]$$

Highest similarity

Transformers of Pehudor - Dewudor network

RNN disadvantages

→ long range dependencies

→ vanishing gradient

handle

temporal

Data

Removed recurrence → makes it parallelizable
 → Attention mechanism

→ Multi head Attention → self attention

as temporal dependency lost?

in order to fix that they introduced positional encoding

Diagram

[20, 80, 0.1]

[0.1, 0.01]

[1, 0.2, 0.1]

[0.1]

[0.1]

Example 8

I/P # Playing outside
embedding [] []

(W) not good
O/P $\cdot T_1 T_2$?

(?) not good

C/W return value

Follow show how transitions go between two states

initial matrix $W^{(k)}$ $\begin{pmatrix} 0.212 & 0.04 & 0.63 & 0.36 \end{pmatrix}$

success $\times N^V$ success

failure $\times N^V$ failure

$$q_1 = [0.212 \ 0.04 \ 0.63 \ 0.36] \quad q_2 = [0.1 \ 0.14 \ 0.86]$$

$$k_1 = [0.31 \ 0.84 \ 0.963 \ 0.57] \quad k_2 = [0.45 \ 0.94 \ 0.73]$$

$$v_1 = [0.36 \ 0.83 \ 0.1 \ 0.38] \quad v_2 = [0.31 \ 0.36 \ 0.19]$$

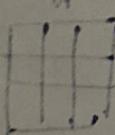
Playing

return

$$\Rightarrow q_1^T \cdot k_1 = p_{11}$$

return

$$p_{11} = [0.212 \ 0.04 \ 0.63 \ 0.36] \begin{bmatrix} 0.31 \\ 0.84 \\ 0.963 \\ 0.57 \end{bmatrix}$$



success

0.72

$$= 0.91121 / \Gamma_4 = 0.4556$$

success probability = 0.4556

$$\Rightarrow q_1^T \cdot k_2 = p_{12}$$

$$p_{12} = [0.212 \ 0.04 \ 0.63 \ 0.36] \begin{bmatrix} 0.45 \\ 0.94 \\ 0.73 \\ 0.58 \end{bmatrix}$$

success after fail in 1 iteration = 0.40085

$$= 0.8017 / \Gamma_4 = 0.40085$$

$$\text{Softmax} = \frac{e^{-P_{11}}}{e^{-P_{11}} + e^{-P_{12}}} \quad \frac{e^{-P_{12}}}{e^{-P_{11}} + e^{-P_{12}}}$$

$$= \frac{\begin{matrix} 0.486 \\ 0.38 \\ 0.1 \\ 0.31 \end{matrix}}{1.3038} \quad \frac{0.513}{1.3038}$$

$$0.486 \quad 0.513$$

$$\text{softmax}^x = \frac{0.486}{1.3038} \begin{pmatrix} 0.26 \\ 0.36 \\ 0.1 \\ 0.38 \end{pmatrix} + \frac{0.513}{1.3038} \begin{pmatrix} 0.31 \\ 0.36 \\ 0.79 \\ 0.712 \end{pmatrix}$$

$$T_1 = \frac{0.486}{1.3038} \begin{pmatrix} 0.334 \\ 0.588 \\ 0.146 \\ 0.554 \end{pmatrix} + \frac{0.513}{1.3038} \begin{pmatrix} 0.31 \\ 0.36 \\ 0.79 \\ 0.712 \end{pmatrix}$$

$$\Rightarrow \alpha_2^T k_1 = P_{21}$$

$$= [0.1 \ 0.14 \ 0.86 \ 0.77] \begin{pmatrix} 0.31 \\ 0.84 \\ 0.963 \\ 0.57 \end{pmatrix}$$

$$= 1041568 / 220 = 0.470784$$

$$\Rightarrow q_2^T k_2$$

$$\begin{bmatrix} 0.1 & 0.14 & 0.86 & 0.37 \end{bmatrix} \begin{bmatrix} 0.31 \\ 0.36 \\ 0.19 \\ 0.72 \end{bmatrix}$$

$$\therefore 0.7992/2 = 0.3996$$

$$\text{Softmax: } \frac{e^{-0.70784}}{e^{-0.70784} + e^{-0.3996}} \quad \frac{e^{-0.3996}}{e^{-0.70784} + e^{-0.3996}}$$

$$= \frac{0.493}{1.1633} \quad \frac{0.6706}{1.1633}$$

$$= 0.4238 \quad 0.5765$$

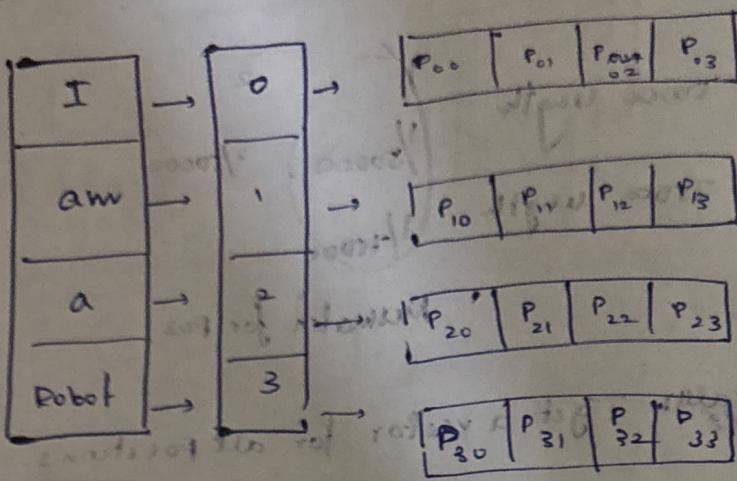
$$= 0.4238 \begin{bmatrix} 0.36 \\ 0.83 \\ 0.1 \\ 0.38 \end{bmatrix} \quad 0.5765 \begin{bmatrix} 0.31 \\ 0.36 \\ 0.19 \\ 0.72 \end{bmatrix}$$

$$X_2 = \begin{bmatrix} 10.0 \\ 8.0 \\ 2.0 \\ 22.0 \end{bmatrix} \begin{bmatrix} 0.3812 \\ 0.5593 \\ 0.151915 \\ 0.5612 \end{bmatrix}$$

$0 < i < d_2$

I am a Robot

$P_G : cl = H \rightarrow \mathbb{R}^H$



$$\theta = 0^\circ \text{ - measured}$$

$$P_{00}, P_{01}$$

$$P_{00} = \cos\left(\frac{\theta}{10000}\right)$$

$$P_{POS, 2i} = \sin\left(\frac{\theta}{10000}\right)$$

$$P_{POS, 2i+1} = \cos\left(\frac{\theta}{10000}\right)$$

$$P_{01} = \cos\left(\frac{\theta}{10000}\right)$$

for presentation

$$P_{02} = \cos\left(\frac{\theta}{10000}\right)$$

$$P_{03} = \cos\left(\frac{\theta}{10000}\right)$$

$$I \rightarrow [0 \mid 1 \mid 0 \mid 1]$$

$i = 0$

P_{10}] \uparrow \leftarrow $\left\{ \begin{array}{l} \text{1st year fail} \\ \text{1st year pass} \end{array} \right.$

P_{11}] \uparrow \leftarrow $\left\{ \begin{array}{l} \text{2nd year fail} \\ \text{2nd year pass} \end{array} \right.$

$$P_{10} = 1 - \text{Sum} \left(\frac{1}{10000} \cdot 2x0/4 \right)$$

$$= \text{Sum} \left(\frac{1}{10000} \right)$$

$$= \text{Sum } 1 = 0.0175$$

$$P_{11} = 0.98 \cdot \left(\frac{1}{10000} \cdot 2x0/4 \right)$$

$$= 0.981 = 0.9998$$

$$i = 1$$

$$\left[\begin{array}{l} P_{12} \\ P_{13} \end{array} \right]$$

Instructional areas

$$P_{12} = \text{Sum} \left(\frac{1}{10000} \cdot 2x1/4^2 \right)$$

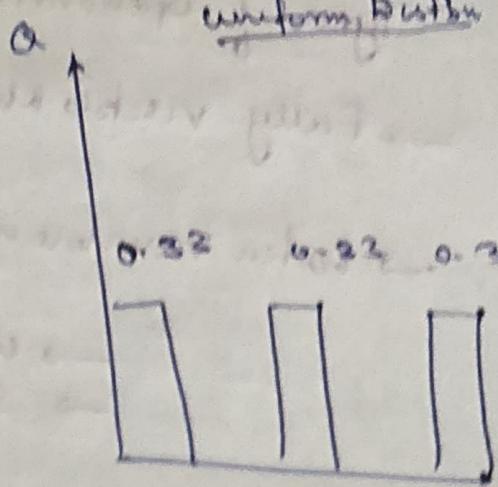
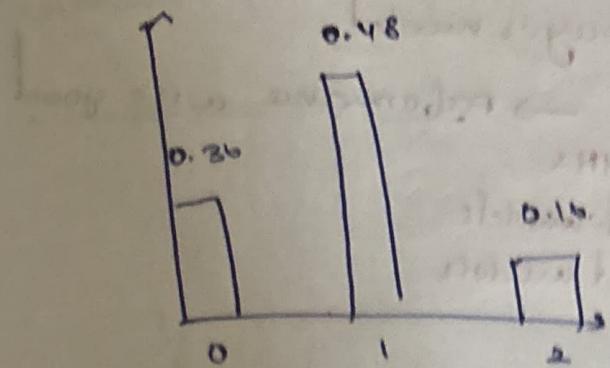
(Instructional first failure)

$$= \text{Sum} \left(\frac{1}{100} \right) = 0.000175$$

$$P_{13} = 0.98 \left(\frac{1}{100} \right) = 0.9999$$

KL example

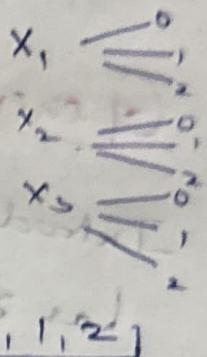
P. Binomial distibn



distribution

	0	1	2
$P(X)$	$\frac{9}{25}$	$\frac{12}{25}$	$\frac{4}{25}$
$Q(X)$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

for bunch
of samples.



$$D_kl(Q||P) = \sum_{x \in X} P(X=x) \ln \frac{Q(x=x)}{P(x=x)}$$

scalars in interval

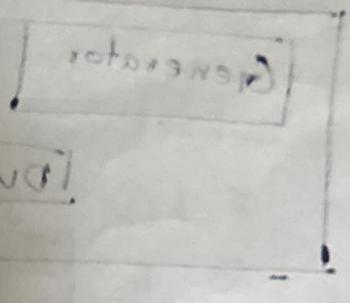
bunch of R.V

$$= \frac{9}{25} \ln \left(\frac{9}{25} \times 3 \right)$$

$$+ \frac{12}{25} \ln \left(\frac{12}{25} \times 3 \right)$$

$$+ \frac{4}{25} \ln \left(\frac{4}{25} \times 3 \right)$$

HAD



bunch of x b. 0.852996

$$D_{kl}(Q||P) = \sum_{x \in X} Q(x=x) \ln \frac{Q(x=x)}{P(x=x)}$$

$$= \frac{1}{3} \ln \left(\frac{1}{3} \times \frac{25}{1} \right) + \frac{1}{3} \ln \left(\frac{1}{3} \times \frac{25}{12} \right) \\ + \frac{1}{3} \ln \left(\frac{1}{3} \times \frac{25}{4} \right)$$

$$= -0.6256 + (-0.1215) + 0.24465 = 0.0978$$