

$$= w_1^{[0]} w_2^{[1]} g(z^{[1]}) + b^{[2]} + b^{[3]}.$$

$$= w_1^{[3]} w_2^{[2]} z^{[1]} + b^{[2]} + b^{[3]}.$$

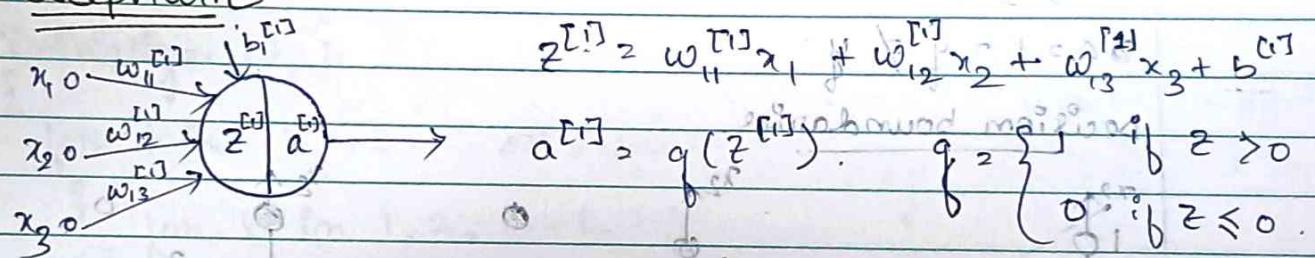
$$= \underbrace{(w_1^{[3]} w_2^{[2]} w_3^{[1]})}_W \underbrace{a^{[0]}}_x + \underbrace{(b^{[1]} + b^{[2]} + b^{[3]})}_b$$

$$\boxed{Wx + b}$$

Representation power of NN.

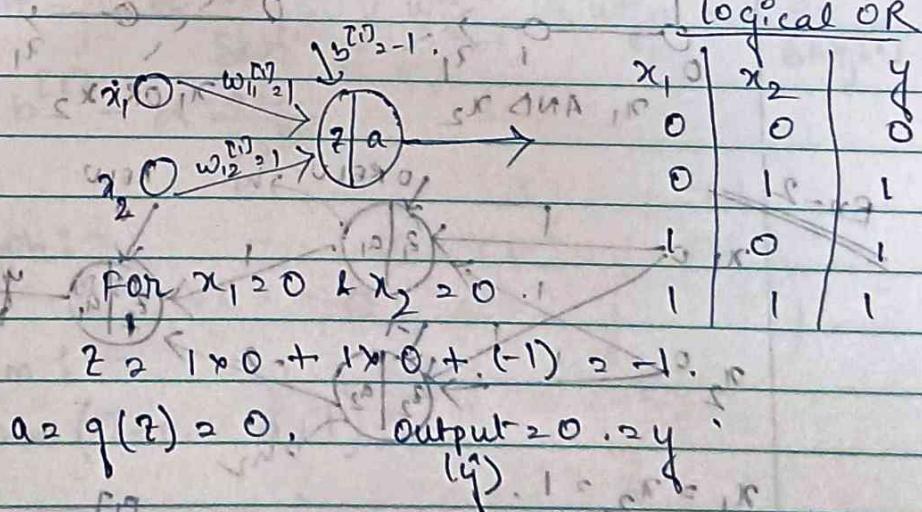
- NN with at least one hidden layer are Universal approximators.
- For any continuous funcⁿ, $f(x)$, and some $\epsilon > 0$, there exists a NN- $g(x)$ with at least 1 hidden layer (with some non-linearity) such that $\forall x |f(x) - g(x)| \leq \epsilon$

Perception



Logical AND

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



For $x_1 = 0$ & $x_2 = 1$.

$$z = 1 \times 0 + 1 \times 1 + (-1) = 0. \quad i.e. (1) \otimes (1) = 1 \times 1 + 1 \times 0 + (-1)$$

$$a = g(z) = 0 \quad \hat{y} = 0^2 + 1^2 + (-1)^2 = 0, \quad g = 0 = y$$

For $x_1 = 1$ & $x_2 = 0$.

$$a = g(z) = 0 \quad \hat{y} = 0^2 + 1^2 + (-1)^2 = 0, \quad g = 0 = y$$

for $x_1 = x_2 = 1$

$$Z = 1 \times 1 + 1 \times 1 + (-1) \\ Z = 1.$$

$$a = g(Z) = y = 1 = y.$$

when $b=0$, logical OR

works for perceptron.

XOR table.

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

for $x_1 = x_2 = 0$.

$$Z = 1 \times 0 + 1 \times 0 + 0.$$

$$a = 0 = y.$$

for $x_1 = x_2 = 1$

$$Z = 1 \times 1 + 1 \times 0 + 0.$$

$$a = 1 = y.$$

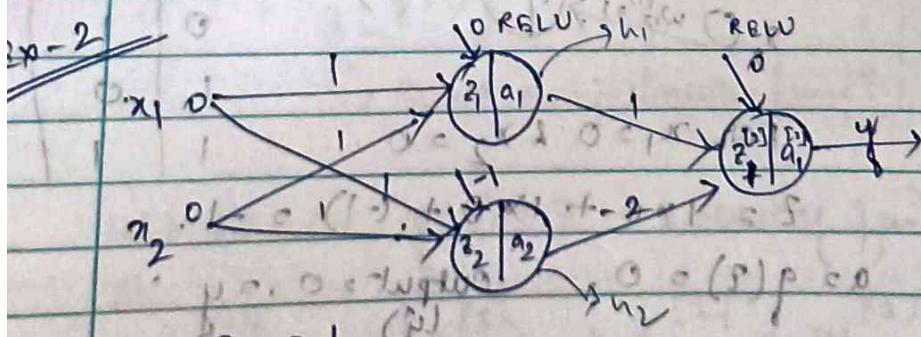
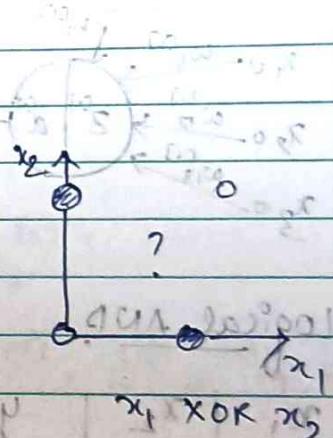
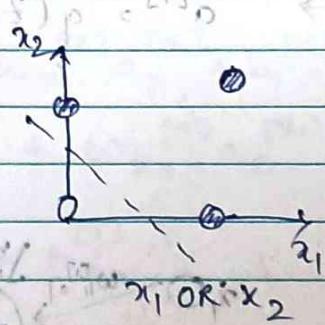
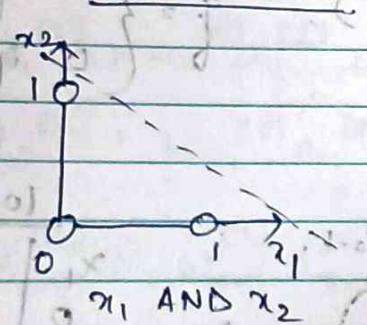
for $x_1 = 1, x_2 = 1$

$$Z = 1 \times 1 + 1 \times 1 + 0$$

$= 2$.

$a = 2 \neq y$.

Decision boundaries.



$$x_1 = x_2 = 1.$$

$$2 \cdot 1 \times 1 + 1 \times 1 + 0 = 2. \quad a_1 = \max(0, 2) = 2.$$

$$2 \cdot 1 \times 1 + 1 \times 1 + (-1) = 1. \quad a_2 = \max(0, 1) = 1.$$

$$2 \cdot 2 \cdot 1 \times 1 + 1 \times (-2) + 0 = 0. \quad a_3 = \max(0, 0) = 0.$$

Steps:-

- ① Identify additional functions.
- ② Draw a computational graph.
- ③ Compute forward pass.
- ④ Perform backward pass starting from end of the circuit.

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

$$w_0 = 2, x_0 = -1, w_1 = -3, x_1 = -2, w_2 = -3$$

$$\textcircled{1} \quad b_3 = w_0x_0 + w_1x_1 + w_2 = b_1 + b_2 + w_2$$

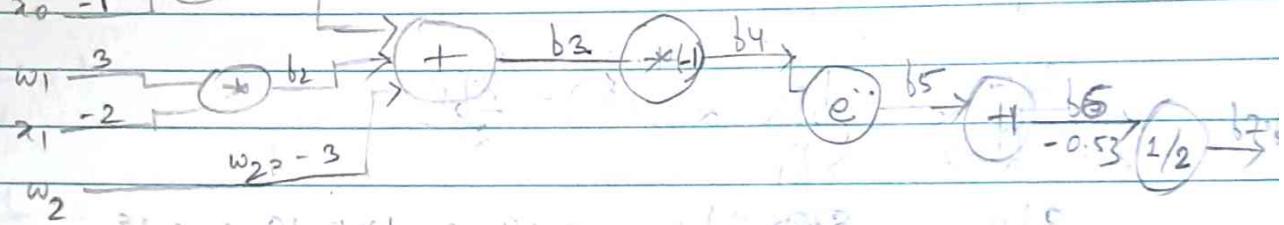
$$f(w, x) = \frac{P_{1.0} \cdot e^{(P_{1.0})}}{1 + e^{-2}} \cdot b_4 = -b_3, b_5 = e^{b_4}$$

$$b_1 = w_0x_0 = 2 \times -1 = -2$$

$$b_2 = w_1x_1 = -3 \times (-2) = +6$$

$$b_3 = b_1 + b_2 + w_2 \\ = -2 + (+6) + (-3) \\ = +1$$

$$\textcircled{2} \quad w_0 = 2, P_{1.0} = P_{1.0} \cdot b_1 = b_7 = 1/b_6$$



$$b_4 = e^1, b_5 = e^2.718 = 14.778$$

$$b_6 = 1 + e^{2.718} = 15.778$$

$$b_7 = \frac{1}{15.778} = 0.064$$

$$\frac{\delta b_7}{\delta b_7} = 1, \frac{\delta b_1}{\delta w_0} = x_0, \frac{\delta b_1}{\delta x_0}, \frac{\delta b_1}{\delta w_0} = w_0$$

local gradients:-

$$\frac{\delta b_1}{\delta w_0} = x_0 = -1$$

$$\frac{\delta b_2}{\delta x_1} = w_1 = -3$$

$$\frac{\delta b_3}{\delta b_1} = 1$$

$$\frac{\delta b_1}{\delta x_0} = w_0 = 2$$

$$\frac{\delta b_3}{\delta b_1} = 1$$

$$\frac{\delta b_5}{\delta b_4} = e^{b_4} = e^{2.718} = 14.778$$

$$\frac{\delta b_2}{\delta w_1} = x_1 = -2$$

$$\frac{\delta b_3}{\delta b_2} = 1$$

$$\frac{\delta b_4}{\delta b_3} = -1$$

Global gradients :-

$$\frac{\delta f_7}{\delta b_6} = \frac{-1}{b_6} \times 1 = \frac{-1}{2.36} \approx -0.54$$

$$\frac{\delta f_7}{\delta b_5} = \frac{\delta b_6}{\delta b_5} \times (-0.54) = 1 \times (-0.54) = -0.54$$

$$\frac{\delta f_7}{\delta b_4} = \frac{\delta b_5}{\delta b_4} \times (-0.54) = e^{\frac{0.19}{2-1}} \times (-0.54) \approx -0.19$$

$$\frac{\delta f_7}{\delta b_3} = \frac{\delta b_4}{\delta b_3} \times (-0.19) = 1 \times (-0.19) = 0.19$$

$$\frac{\delta f_7}{\delta b_2} = \frac{\delta b_3}{\delta b_2} \times 0.19 = 1 \times 0.19 = 0.19$$

~~$$\frac{\delta f_7}{\delta w_2} = \frac{\delta b_3}{\delta w_2} \times 0.19 = 0.19$$~~

$$\frac{\delta f_7}{\delta w_2} = \frac{\delta b_3}{\delta w_2} \times 0.19 = 1 \times 0.19 = 0.19$$

$$\frac{\delta f_7}{\delta b_2} = \frac{\delta b_3}{\delta b_2} \times 0.19 = 1 \times 0.19 = 0.19$$

$$\frac{\delta f_7}{\delta b_1} = \frac{\delta b_3}{\delta b_1} \times 0.19 = 0.19$$

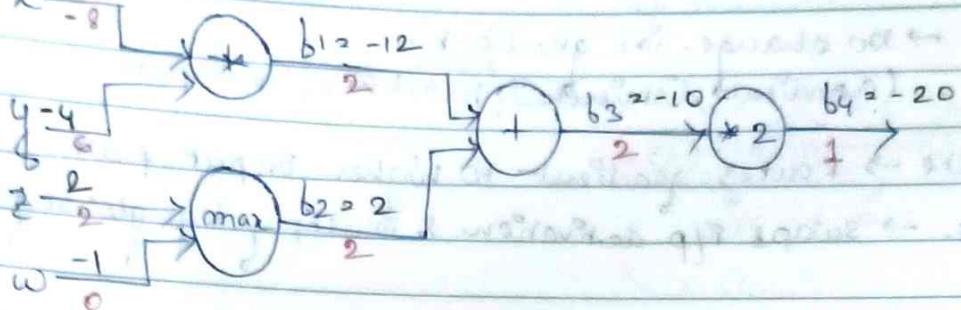
$$\frac{\delta f_7}{\delta w_0} = \frac{\delta b_1}{\delta w_0} \times 0.19 = -1 \times 0.19 = -0.19$$

$$\frac{\delta f_7}{\delta x_0} = \frac{\delta b_1}{\delta x_0} \times 0.19 = 2 \times 0.19 = 0.38$$

$$\frac{\delta f_7}{\delta x_1} = \frac{\delta b_2}{\delta x_1} \times 0.19 = -3 \times 0.19 = -0.57$$

$$\frac{\delta f_7}{\delta w_1} = \frac{\delta b_2}{\delta w_1} \times 0.19 = -2 \times 0.19 = -0.38$$

$$09.3 \quad x = \frac{3}{-8}$$



$$b_1 = x * y$$

$$\frac{\delta b_1}{\delta z} = y = 4$$

$$\frac{\delta b_3}{\delta b_1} = 1 \quad S$$

$$b_2 = 8 \max(z, w)$$

$$\frac{\delta b_1}{\delta y} = x = 3$$

$$\frac{\delta b_3}{\delta b_2} = 1$$

$$b_3 = b_1 + b_2$$

$$\frac{\delta b_2}{\delta w} =$$

$$\frac{\delta b_4}{\delta b_3} = 2$$

$$b_4 = b_3 * 2$$

$$\frac{\delta b_2}{\delta z} =$$

$$\frac{\delta b_4}{\delta b_4} = 1$$

$$\frac{\delta b_4}{\delta b_3} = 2 \times 1 = 2$$

$$\frac{\delta b_4}{\delta b_1} = \frac{\delta b_3}{\delta b_1} \times 2 = 1 \times 2 = 2$$

$$\frac{\delta b_4}{\delta b_2} = \frac{\delta b_3}{\delta b_2} \times 2$$

$$= 1 \times 2 = 2$$

$$\frac{\delta b_4}{\delta z} = \frac{\delta b_1}{\delta z} \times 2 = -4 \times 2 = -8$$

$$\frac{\delta b_2}{\delta z} = 1 \quad \text{if } z > w.$$

$$\frac{\delta b_4}{\delta y} = \frac{\delta b_1}{\delta y} \times 2 = 3 \times 2 = 6$$

$$0 \quad \text{if } z < w.$$

undefined if $z = w$.
arbitrary.

$$\frac{\delta b_2}{\delta w} = 1 \quad \text{if } w > z.$$

undefined if $w = z$.

$\rightarrow 0$.

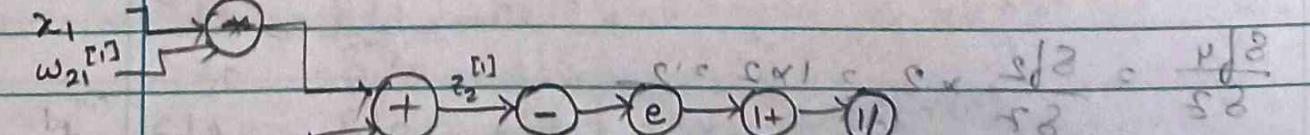
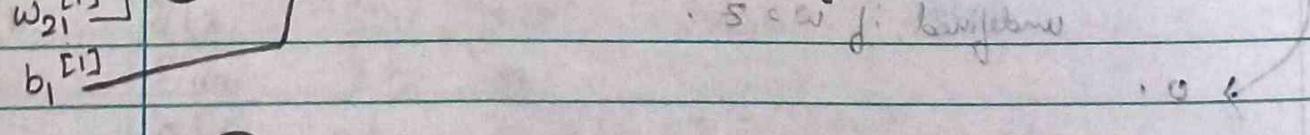
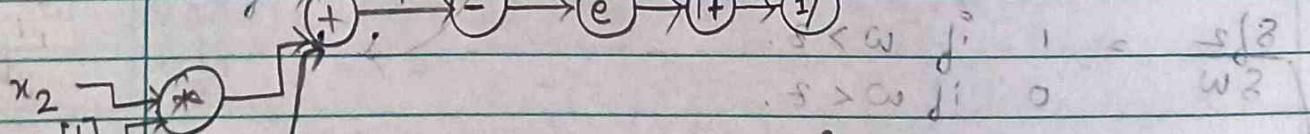
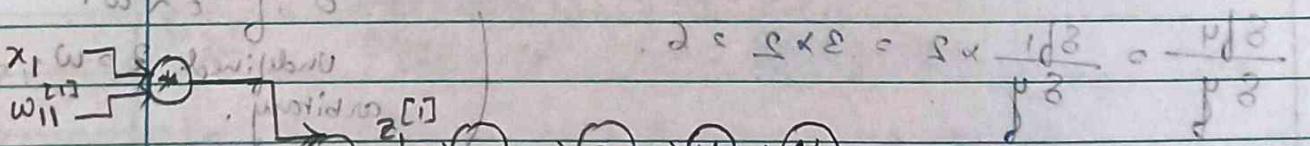
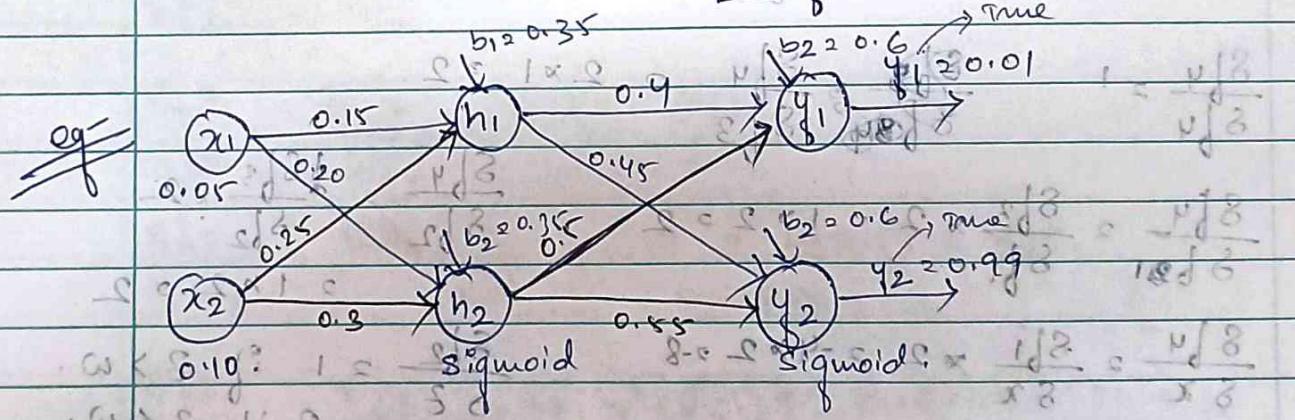
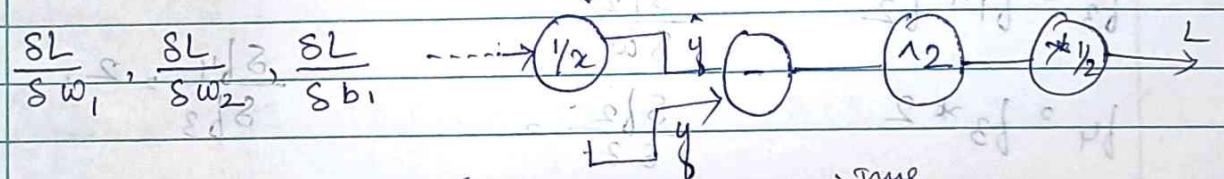
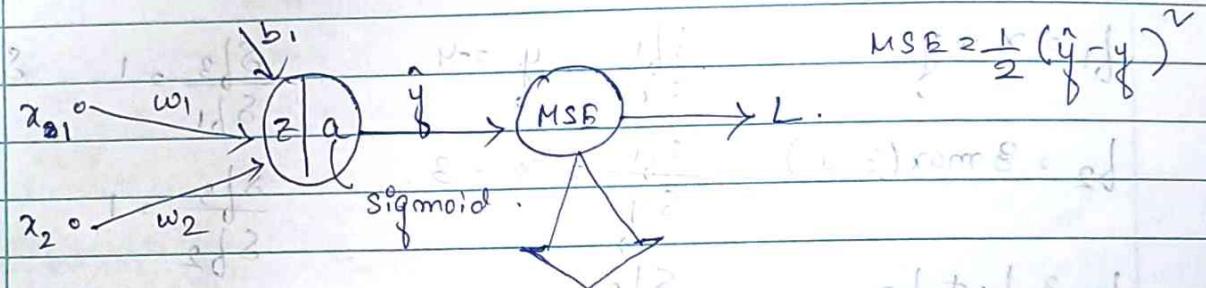
$$\frac{\delta b_4}{\delta z} = \frac{\delta b_2}{\delta z} \times 2 = 1 \times 2 = 2$$

$$\frac{\delta b_4}{\delta w} = \frac{\delta b_2}{\delta w} \times 2 = 0 \times 2 = 0$$

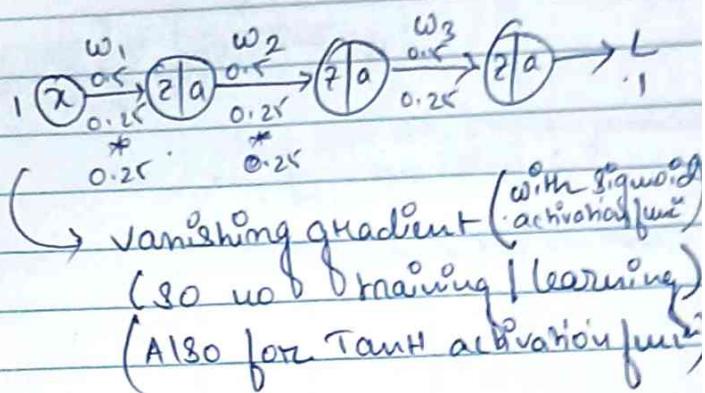
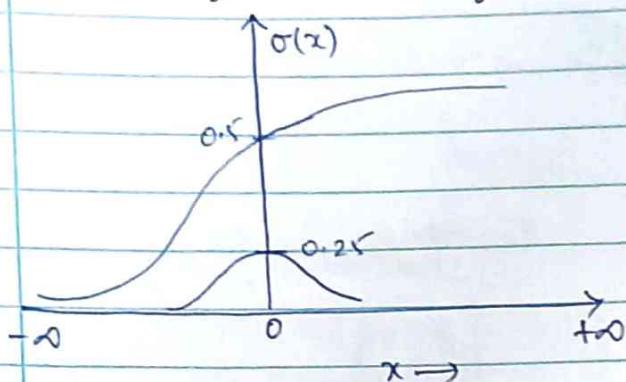
(+) gate \rightarrow no change in gradient; gradient distributed equally to all 2 of its inputs.

(max) gate \rightarrow routes gradient to higher input path.

(*) gate \rightarrow swaps i/p activation & multiply by global gradient.



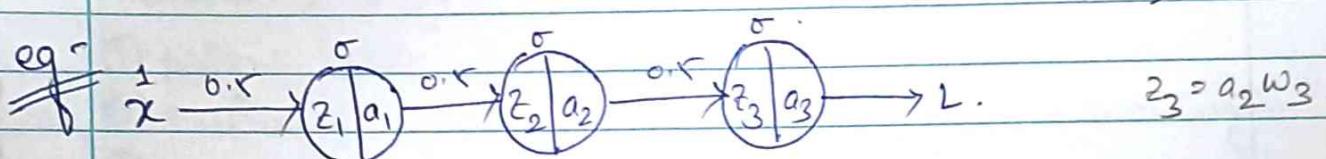
Vanishing & Exploding Gradients Problem



To fix vanishing gradient problem \rightarrow ReLU is used.
but this

) but this leads to

Exploding gradient problem
(clipping the gradient reduces the problem).



$$Z_1 = 1 \times 0.5 > 0.5$$

$$z_2 = 0.622 \times 0.15 \\ \approx 0.311$$

$$z_3 = 0.577 \times 0.5$$

$$= 0.2885$$

$$a_1 = 5(0.5) = 0.622.$$

2 6 2

$$a_2 = \sigma(0.311) = 0.577 \cdot a_3 = \sigma(0.2885) = 0.571 \cdot$$

$$\text{loss}(L) = \frac{1}{2} (0.571 - 1)^2 = 0.092$$

$$\frac{SL}{8a_3} \cdot (a_3 - y) = 0.571 - 1 = 0.428.$$

$$\frac{8a_3}{8z_3} = \frac{0.0009}{0.0001} (1 - 0.288) = 0.288(1 - 0.288) = 0.208$$

$$\frac{8^2 3}{8 w_3} = a_2 = 0.577 \dots$$

Convolution operation :- $3 \times 1 + 0 \times 0 + 1 \times (-1)$

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

Suppose

convolution
operator.

-5	-4	0	8
-10	-2	2	3
0	-2	-4	-7
-3	-2	-3	-16

Input 6×6 image (m)

Filter

Kernel
 3×3 (f).

$$3 \times 1 + 0 \times 0 + 1 \times (-1) + 1 \times 1 + 5 \times 0 + 8 \times (-1) + 2 \times 1 + 7 \times 0 + 2 \times (-1)$$

= -5

Slide filter by 1:

$$1 \times 0 + 0 \times 1 + 2 \times -1 + 5 \times 1 + 8 \times 0 + 9 \times (-1) + 7 \times 1 + 2 \times 0 + 5 \times (-1)$$

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

Vertical filter.

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

Horizontal
filter.

Both for edge detection.

Edge detectors - Sobel edge detection.

Canny

Laplacian

Prewitt

Sobel

$$\text{Output image size} = (m-f+1) \times (m-f+1)$$

Paddingfor every layer \rightarrow input image will shrink. (just like $6 \times 6 \rightarrow 4$)

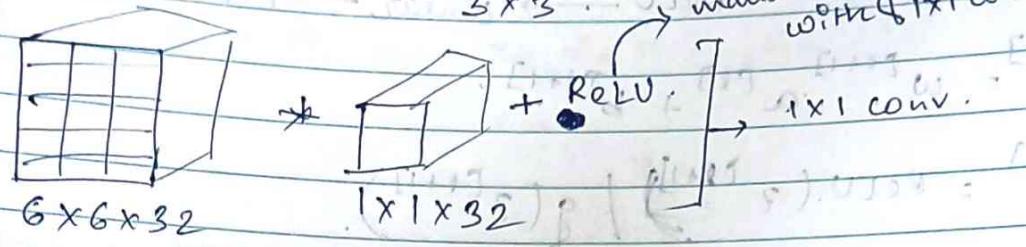
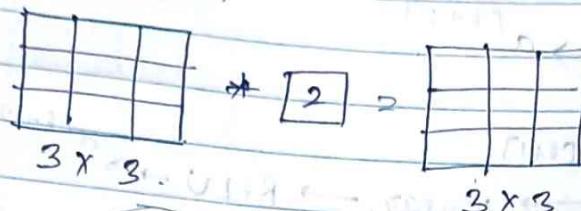
{ filter does not see all the pixels with same no. of times. problem. }

Solve 2 problem —

① Shrinking problem.

② Equal attention to all pixels.

1x1 conv / Network in a Network.



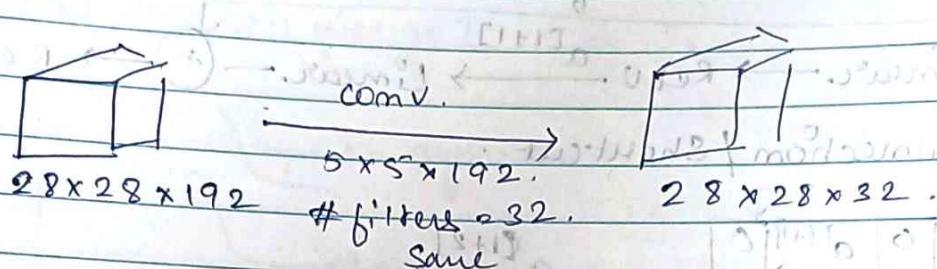
$$\begin{aligned} m &= b+1 \\ &= 28 \times 5+1 \\ &= 241 \end{aligned}$$

$$m-f$$

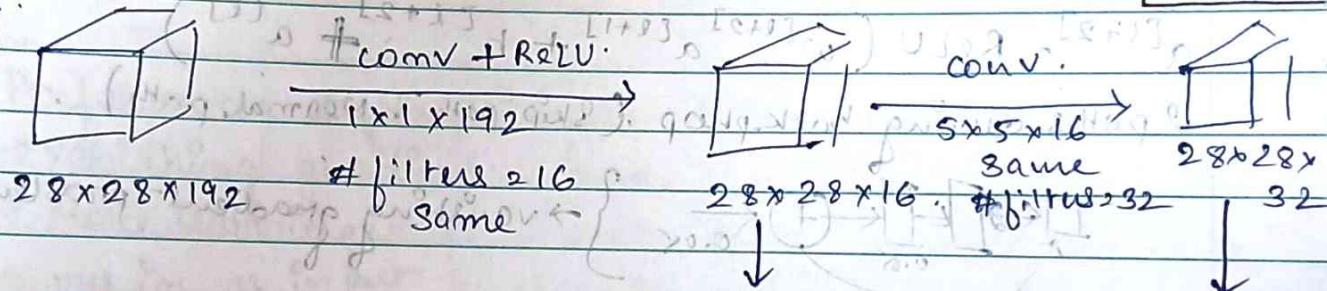
$$\begin{matrix} 3 \times 3 & \times & 2 \times 2 \end{matrix}$$

$$\begin{matrix} 16 \\ 3 \times 3 \times 2 \times 2 \end{matrix}$$

$$\begin{matrix} 9 \times 4 \\ 3c \end{matrix}$$



$$\text{Total operations (multiplications)} = 5 \times 5 \times 192 \times 28 \times 28 \times 32 \approx 120,422,400 \approx [120 \text{ million}]$$



(*) 1x1 conv \rightarrow also called. 2408448. + 10,035,200
bottleneck layer $\approx 2.4 \text{ million}$ $\approx 10 \text{ million}$

(*) From 120M to 12M just by using 1x1 conv. \rightarrow Reduce the no. of channels/increase channels (ReLU, canon)

so,

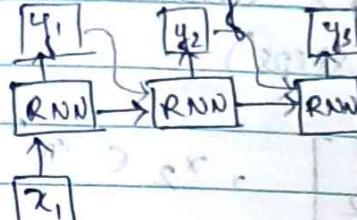
① Modify height, width \rightarrow Padding, striding, pooling

② Modify channels \rightarrow No. of filters, 1x1 conv.

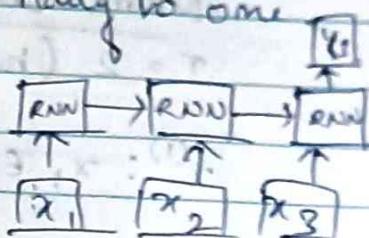
(a) One-to-one



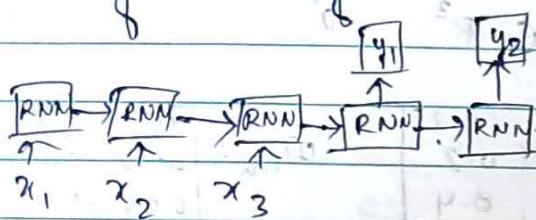
(b) One-to-many



(c) Many-to-one



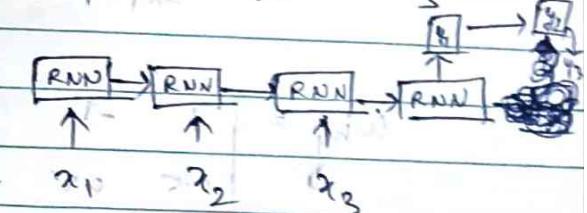
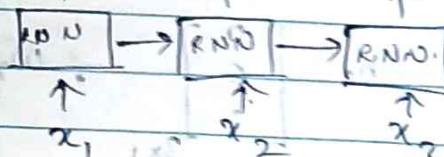
(d) many-to-many



video captioning

y₁y₂y₃y₁y₂y₃y₁y₂y₃

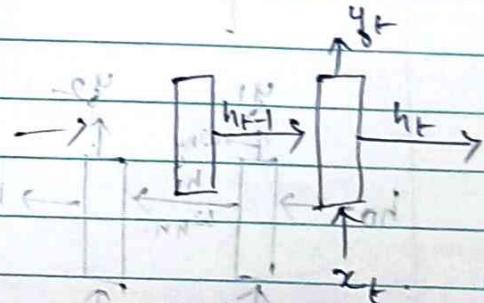
Translation [e.g. English to French]

to see the entire sentence
then give the output.

$$h_t = f_w(h_{t-1}, x_t).$$

$$\Rightarrow \tanh(w_{hn} h_{t-1} + w_{xn} x_t)$$

Note: No bias term



$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}_{h_t} = \tanh \left(\begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}_{w_{hn}} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}_{h_{t-1}} + \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}_{w_{xn}} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}_{x_t} \right)$$

Weight sharing
at every time
stamp.

$$\begin{bmatrix} 0 \\ 0 \end{bmatrix}_{y_t} = w_{ny} \begin{bmatrix} 1 \\ h_t \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ h_t \end{bmatrix}$$

$$h_t = 4 \times 1$$

$$1.0 - 2.0 - 3.0 - 4.0 - 5.0 - 6.0 - 7.0 - 8.0 - 9.0 - 10.0$$

$$2.0 - 3.0 - 4.0 - 5.0 - 6.0 - 7.0 - 8.0 - 9.0 - 10.0$$

$$3.0 - 4.0 - 5.0 - 6.0 - 7.0 - 8.0 - 9.0 - 10.0$$

$$4.0 - 5.0 - 6.0 - 7.0 - 8.0 - 9.0 - 10.0$$

$$5.0 - 6.0 - 7.0 - 8.0 - 9.0 - 10.0$$

$$6.0 - 7.0 - 8.0 - 9.0 - 10.0$$

$$7.0 - 8.0 - 9.0 - 10.0$$

$$8.0 - 9.0 - 10.0$$

$$9.0 - 10.0$$

Numerical example :- (vanilla RNN)

$T = 2$ (two time steps)

$$x^{(1)}: x_1 \in \mathbb{R}^2 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, x_2 \in \mathbb{R}^2 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

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

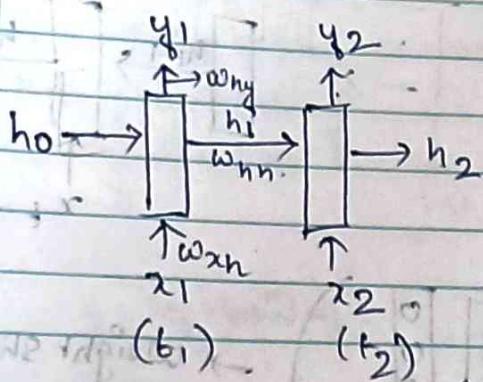
Hidden state dim $= 3 \rightarrow h_t \in \mathbb{R}^3$

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

$$w_{xh} = \begin{bmatrix} 0.5 & -0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.4 \end{bmatrix}_{3 \times 2}$$

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

$$w_{hy} = \begin{bmatrix} 1 & -1 & 0.5 \\ 0.5 & 0.5 & -0.5 \end{bmatrix}_{2 \times 3}$$



$$h_1 = \tanh(w_{hh} \cdot h_0 + w_{xh} \cdot x_1)$$

$$y_1 = w_{hy} \cdot h_1$$

For $t=1$

$$w_{hh} \cdot h_0 + w_{xh} \cdot x_1$$

$$\begin{bmatrix} 0.1 & 0.4 & 0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.2 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0.5 & -0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.4 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$2 \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} -0.1 \\ 1.2 \\ 0.9 \end{bmatrix} = \begin{bmatrix} -0.1 \\ 1.2 \\ 0.9 \end{bmatrix}$$

$$h_1 = \tanh \left(\begin{bmatrix} -0.1 \\ 1.2 \\ 0.9 \end{bmatrix} \right) = \begin{bmatrix} -0.1 \\ 0.83 \\ 0.716 \end{bmatrix}$$

$$\frac{e^{-0.1} - e^{0.1}}{2.05} = 0.32 - 0.301$$

$$\frac{e^{1.2} - e^{-1.2}}{2.05} = 0.406 - 0.301$$

$$\frac{e^{0.9} - e^{-0.9}}{2.05} = 0.904 - 0.105$$

$$\frac{e^{0.9} + e^{-0.9}}{2.05} = 0.904 + 0.105 = 1.009$$

$$\frac{-0.201}{2.05} = -0.098$$

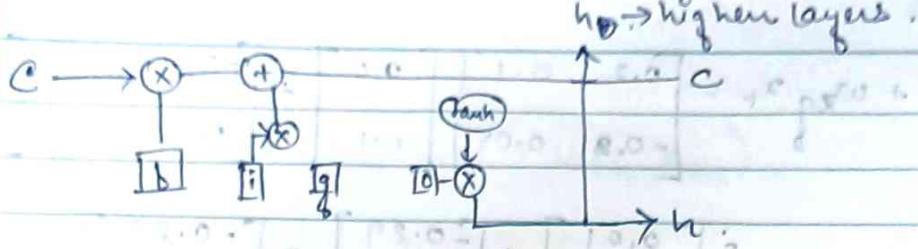
$$q_1 = \begin{bmatrix} 1 & -1 & 0.5 \\ 0.5 & 0.5 & -0.5 \end{bmatrix} \begin{bmatrix} -0.1 \\ 0.83 \\ 0.716 \end{bmatrix} = \begin{bmatrix} -0.572 \\ 0.007 \end{bmatrix}$$

For t = 2

$$\text{Ans } w_{nn} \times h_1 + w_{x_n} \times x_2$$

$$2 \begin{bmatrix} 0.1 & 0.4 & 0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.2 \end{bmatrix} \begin{bmatrix} -0.1 \\ 0.83 \\ 0.716 \end{bmatrix} + \begin{bmatrix} 0.5 & -0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.4 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

$$2 \begin{bmatrix} 0.322 \\ 0.4122 \end{bmatrix} +$$



Numerical example :-

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

$$w_{x0} = \begin{bmatrix} 0.5 & -0.3 \\ 0.4 & 0.1 \end{bmatrix} \quad w_{q0} = \begin{bmatrix} -0.4 & 0.2 \\ 0.3 & 0.3 \end{bmatrix}$$

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

$$w_{h0} = \begin{bmatrix} 0.1 & 0.2 \\ -0.2 & 0.05 \end{bmatrix} \quad w_{h1} = \begin{bmatrix} 0.05 & -0.1 \\ 0.2 & 0.1 \end{bmatrix}$$

$$w_{x0} = \begin{bmatrix} 0.3 & 0.25 \\ -0.2 & 0.2 \end{bmatrix}$$

$$w_{q0} = \begin{bmatrix} 0.15 & 0.05 \\ 0.1 & -0.2 \end{bmatrix}$$

$$w_{h1} = \begin{bmatrix} -0.5 & 0.4 \\ 0.2 & -0.3 \end{bmatrix}$$

$$w_{h1} = \begin{bmatrix} 0.2 & 0.1 \\ -0.1 & 0.05 \end{bmatrix}$$

Do - 1 forward pass for 1 time step.
Find - h_t & c_t values.

$$w_{h1} h_{t-1} + w_{x0} x_t \Rightarrow \begin{bmatrix} 0.1 & 0.2 \\ -0.2 & 0.05 \end{bmatrix} \begin{bmatrix} 0.0 \\ 0.1 \end{bmatrix} + \begin{bmatrix} 0.5 & -0.3 \\ 0.4 & 0.1 \end{bmatrix} \begin{bmatrix} 0.5 \\ -0.1 \end{bmatrix}$$

$$= \begin{bmatrix} 0.02 \\ 0.005 \end{bmatrix} + \begin{bmatrix} 0.28 \\ 0.19 \end{bmatrix} = \begin{bmatrix} 0.3 \\ 0.195 \end{bmatrix}$$

$$i_t = \sigma \begin{bmatrix} 0.3 \\ 0.195 \end{bmatrix} = \begin{bmatrix} 0.574 \\ 0.548 \end{bmatrix}$$

$$w_{h0} h_{t-1} + w_{x0} x_t \Rightarrow \begin{bmatrix} 0.15 & 0.05 \\ 0.1 & -0.2 \end{bmatrix} \begin{bmatrix} 0.0 \\ 0.1 \end{bmatrix} + \begin{bmatrix} 0.3 & 0.25 \\ -0.2 & 0.2 \end{bmatrix} \begin{bmatrix} 0.5 \\ -0.1 \end{bmatrix}$$

$$= \begin{bmatrix} 0.005 \\ -0.02 \end{bmatrix} + \begin{bmatrix} 0.125 \\ -0.12 \end{bmatrix} = \begin{bmatrix} 0.13 \\ -0.14 \end{bmatrix}$$

$$o_t = \sigma \begin{bmatrix} 0.13 \\ -0.14 \end{bmatrix} = \begin{bmatrix} 0.532 \\ 0.468 \end{bmatrix}$$

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

$$1.323 - 0.7557$$

$$1.323 + 0.7552$$

$$0.5673$$

$$2.0787$$

$$w_{hg} h_{t-1} + w_{xg} x_t \Rightarrow \begin{bmatrix} 0.2 & 0.1 \\ -0.21 & 0.05 \end{bmatrix} \begin{bmatrix} 0.0 \\ 0.1 \end{bmatrix} + \begin{bmatrix} -0.5 & 0.4 \\ 0.2 & -0.3 \end{bmatrix} \begin{bmatrix} 0.5 \\ -0.1 \end{bmatrix}$$

$$-0.1201$$

$$2.003 \times$$

$$-1.030$$

$$2 \begin{bmatrix} 0.01 \\ 0.005 \end{bmatrix} + \begin{bmatrix} -0.29 \\ 0.13 \end{bmatrix} = \begin{bmatrix} -0.28 \\ 0.135 \end{bmatrix}$$

$$0.9704$$

$$q_t = \tanh \begin{bmatrix} -0.28 \\ 0.135 \end{bmatrix} = \begin{bmatrix} 0.272 \\ 0.134 \end{bmatrix}$$

$$-0.0506$$

$$2.0004$$

$$w_{hb} h_{t-1} + w_{xb} x_t \Rightarrow \begin{bmatrix} 0.05 & -0.1 \\ 0.2 & 0.1 \end{bmatrix} \begin{bmatrix} 0.0 \\ 0.1 \end{bmatrix} + \begin{bmatrix} -0.4 & 0.2 \\ 0.3 & 0.3 \end{bmatrix} \begin{bmatrix} 0.5 \\ -0.1 \end{bmatrix}$$

$$2 \begin{bmatrix} -0.01 \\ -0.01 \end{bmatrix} + \begin{bmatrix} -0.24 \\ 0.12 \end{bmatrix} = \begin{bmatrix} -0.25 \\ 0.13 \end{bmatrix}$$

$$b_t = 0 \begin{bmatrix} -0.25 \\ 0.13 \end{bmatrix} = \begin{bmatrix} 0.44 \\ 0.53 \end{bmatrix}$$

$$c_t = \begin{bmatrix} 0.44 \\ 0.53 \end{bmatrix} \oplus \begin{bmatrix} 0.2 \\ -0.2 \end{bmatrix} + \begin{bmatrix} 0.57 \\ 0.55 \end{bmatrix} \oplus \begin{bmatrix} -0.27 \\ 0.13 \end{bmatrix}$$

$$2 \begin{bmatrix} 0.088 \\ -0.106 \end{bmatrix} + \begin{bmatrix} -0.1539 \\ 0.0715 \end{bmatrix} = \begin{bmatrix} -0.06 \\ -0.03 \end{bmatrix}$$

$$h_t = \begin{bmatrix} 0.53 \\ 0.46 \end{bmatrix} \oplus \tanh \begin{bmatrix} -0.06 \\ -0.03 \end{bmatrix} = \begin{bmatrix} 0.53 \\ 0.46 \end{bmatrix} \oplus \begin{bmatrix} -0.0599 \\ -0.0297 \end{bmatrix}$$

$$= \begin{bmatrix} -0.03 \\ -0.01 \end{bmatrix}$$

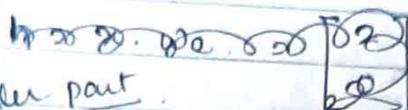
$$x \in \mathbb{R}, h \in \mathbb{R}$$

$x = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$

$$w_e \rightarrow 1 \times 2 = \begin{bmatrix} 0.5 & -1.0 \end{bmatrix}_{1 \times 2}$$

$$w_d \rightarrow 2 \times 1 = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}_{2 \times 1}$$

$$\text{MSE loss} \rightarrow \frac{1}{2} \| \hat{x} - x \|^2$$



Encoder part.



$$h = w_e \cdot x = \begin{bmatrix} 0.5 & -1.0 \end{bmatrix} \begin{bmatrix} 2 \\ 0 \end{bmatrix} \rightarrow 1$$

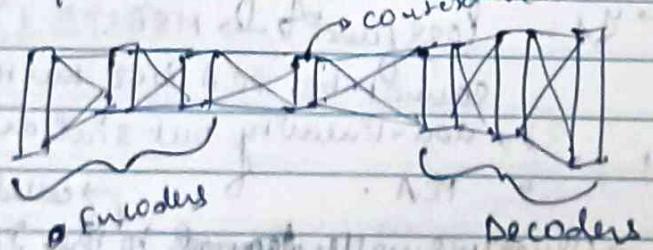
Decoder part.

$$\hat{x} = w_d \cdot h = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix} + 1 = \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$$

$$\text{Loss} = \frac{1}{2} \| (\hat{x} - x) \|^2 \rightarrow \text{Loss} = \frac{1}{2} \sqrt{(1-2)^2 + (0-0.5)^2}$$

$$\begin{aligned} & \frac{1}{2} \sqrt{1 + 0.25} \\ & \rightarrow \frac{1.118}{2} \rightarrow 0.559. \end{aligned}$$

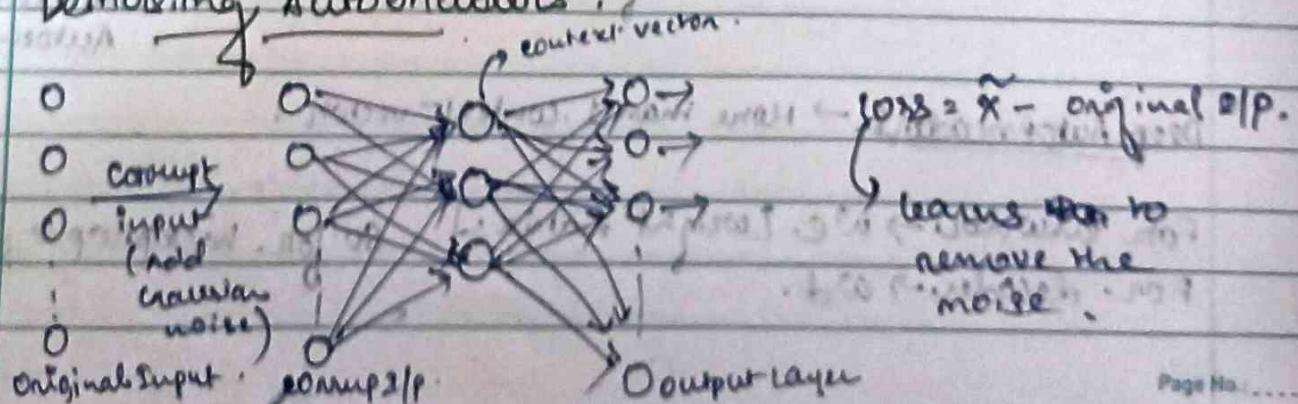
Deep Autoencoders \rightarrow

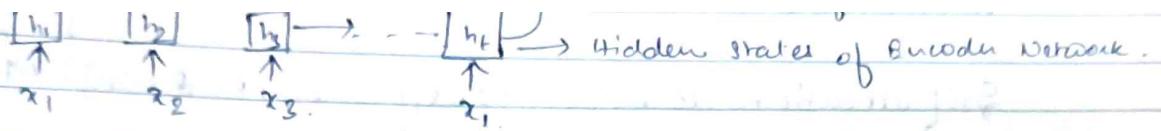


Add Gaussian noise.

Q&A

Demoing Autoencoders.





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

$$\alpha_{t,j} = \frac{\exp(\text{score}(s_{t-1}, h_j))}{\sum_{j=1}^T \exp(\text{score}(s_{t-1}, h_j))} \quad \left| \begin{array}{l} \text{score} \Rightarrow \text{similarity / Alignment} \\ \text{score.} \end{array} \right.$$

~~eg~~

$$h_1 = [1, 0, 1] \quad c_t = \sum_{j=1}^T \alpha_{t,j} h_j^o$$

$$h_2 = [0, 1, 1] \quad c_t = \alpha_{t,1} h_1 + \alpha_{t,2} h_2 + \alpha_{t,3} h_3$$

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

$$s_{t-1} = [1, 0, 1] \quad \alpha_{t,1} = \frac{\exp(\text{score}(s_{t-1}, h_1))}{\exp(\text{score}(s_{t-1}, h_1)) + \exp(\text{score}(s_{t-1}, h_2)) + \exp(\text{score}(s_{t-1}, h_3))}$$

Score function \Rightarrow dot product

$$\alpha_{t,1} = \frac{\exp(2)}{\exp(2) + \exp(1) + \exp(1)} \rightarrow \frac{7.39}{12.826} = 0.57$$

$$\alpha_{t,2} = \frac{\exp(1)}{12.826} = 0.2$$

$$\alpha_{t,3} = \frac{\exp(1)}{12.826} = 0.2$$

$$c_t = 0.57 [1, 0, 1] + 0.2 [0, 1, 1] + 0.2 [1, 1, 0] \\ = [0.57, 0, 0.57] + [0, 0.2, 0.2] + [0.2, 0.2, 0] = [0.77, 0.4, 0.77]$$

0.91.5
3.0629

1.577
0.634

0.943
2.211

apsara

Date 9.10.25

Numerical example :- Compute z_1 and z_2 .

Sentence = "Playing Outside"

For playing :-

$$q_1 = [0.212 \ 0.04 \ 0.63 \ 0.36]$$

$$k_1 = [0.31 \ 0.84 \ 0.963 \ 0.57]$$

$$v_1 = [0.36 \ 0.83 \ 0.1 \ 0.36]$$

For outside :-

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

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

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

Note) Use scaled dot product - attention.

for playing

$$q_1 \cdot k_1 = [0.212 \ 0.04 \ 0.63 \ 0.36] \odot [0.31 \ 0.84 \ 0.963 \ 0.57]$$

$$= 0.06572 + 0.0336 + 0.60669 + 0.2052$$

$$= 0.91151$$

Now divide by $\sqrt{4}$ [dimensions of k_1]

$$\frac{q_1 \cdot k_1}{\sqrt{4}} = \frac{0.91151}{2} = 0.4557$$

$$\text{softmax}(0.4557) = 0.496 \quad q_1 \cdot k_2 = [0.212 \ 0.04 \ 0.63 \ 0.36] \odot [0.45 \ 0.94 \ 0.73 \ 0.58]$$

$$= 0.0954 + 0.0376 + 0.4899 + 0.2088$$

Now divide by $\sqrt{4}$

$$\frac{q_1 \cdot k_2}{\sqrt{4}} = \frac{0.8017}{2} = 0.40085$$

$$\text{softmax}(0.4557, 0.40085) = 0.513, 0.486$$

$$0.513 \times v_1$$

$$0.496 \times [0.31 \ 0.36 \ 0.19 \ 0.72]$$

$$= 0.513 [0.36 \ 0.83 \ 0.1 \ 0.36] = 0.150 \ 0.1746 \ 0.0921$$

$$+ [0.185 \ 0.426 \ 0.0513 \ 0.185] = 0.349$$

$$0.486 \times v_2 = [0.185 \ 0.426 \ 0.0513 \ 0.185] + [0.150 \ 0.1746 \ 0.0921]$$

$$= 0.349$$

$$= [0.33 \ 0.605 \ 0.143 \ 0.534]$$

3.9961

 $v_2 \cdot k_1$

$$\Rightarrow [0.1 \ 0.14 \ 0.86 \ 0.77] @ [0.31 \ 0.884 \ 0.963 \ 0.57]$$

$$\Rightarrow [0.031 + 0.1176 + 0.828 + 0.4389]$$

$$\Rightarrow 1.4155$$

Now,

 $v_2 \cdot k_1$

$$\frac{v_2 \cdot k_1}{\sqrt{4}} \Rightarrow \frac{1.4155}{2} = 0.707$$

 $v_2 \cdot k_2$

$$\Rightarrow [0.1 \ 0.14 \ 0.86 \ 0.77] @ [0.45 \ 0.94 \ 0.73 \ 0.58]$$

$$\Rightarrow [0.045 + 0.1316 + 0.6278 + 0.4466] \Rightarrow 1.251$$

Now,

 $v_2 \cdot k_2$

$$\frac{v_2 \cdot k_2}{\sqrt{4}} \Rightarrow \frac{1.251}{2} = 0.6255$$

$$\text{Softmax}(0.707, 0.6255) \Rightarrow 0.52, 0.48$$

 $0.52 * v_2$

$$\Rightarrow 0.52 [0.31 \ 0.36 \ 0.19 \ 0.72]$$

$$\Rightarrow [0.1612 \ 0.1872 \ 0.0988 \ 0.3744]$$

 $0.48 * v_2$

$$\Rightarrow 0.48 [0.31 \ 0.36 \ 0.19 \ 0.72]$$

~~$$\Rightarrow [0.1612 \ 0.1872 \ 0.0988 \ 0.3744]$$~~

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

$$\Rightarrow [0.1612 \ 0.1872 \ 0.0988 \ 0.3744] + [0.1488 \ 0.1728 \ 0.0912 \ 0.3456]$$

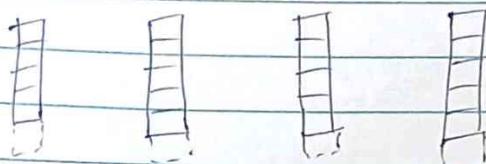
$$\Rightarrow [0.31 \ 0.36 \ 0.19 \ 0.72]$$

Layer Normalization

$$\boxed{1} \boxed{1} \boxed{-1} \rightarrow \underline{x - \mu}$$

Positional encoding :-

Key components of Transformer architecture



Sinusoidal & cosine funcⁿ.

$$PE_{pos, 2i} = \sin\left(\frac{pos \cdot 2i}{10000^{\frac{2i}{d_{emb}}}}\right)$$

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

Positional embedding vector.

Even pos \rightarrow use sine

Odd pos \rightarrow use cosine

				Positional Embedding				
I	am	a	robot.	P ₀₀	P ₀₁	P ₀₂	P ₀₃	0 1 0 1
P _{E₁} , d=4 $\Rightarrow R^4$	am	\rightarrow	1	P ₁₀	P ₁₁	P ₁₂	P ₁₃	0.84 0.54 0.10 1
	a	\rightarrow	2					
	Robot	\rightarrow	3	P ₃₀	P ₃₁	P ₃₂	P ₃₃	

$$P_{00} = \sin\left(\frac{0}{10000^{\frac{0}{4}}}\right) = \sin(0) = 0 \quad P_{10} = \sin\left(\frac{1}{10000^{\frac{0}{4}}}\right) = \sin(1) = 0.8414$$

$$P_{01} = \cos\left(\frac{0}{10000^{\frac{1}{4}}}\right) = \cos(0) = 1 \quad P_{11} = \cos\left(\frac{1}{10000^{\frac{1}{4}}}\right) = \cos(1) = 0.54$$

$$P_{02} = \sin\left(\frac{0}{10000^{\frac{2}{4}}}\right) = \sin(0) = 0 \quad P_{12} = \sin\left(\frac{1}{10000^{\frac{2}{4}}}\right), 0.10$$

$$P_{03} = \cos\left(\frac{0}{10000^{\frac{3}{4}}}\right) = \cos(0) = 1 \quad P_{13} = \cos\left(\frac{1}{10000^{\frac{3}{4}}}\right) = 0.991 = 1$$