

Forward Pass:-

Backward Pass

- * Once we get the output we travel backwards & update the weights

⇒ Derivation of Tan-H Activation function

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

$$\frac{dg(x)}{dx} = (e^x + e^{-x}) \frac{d(e^x - e^{-x})}{dx} - (e^x - e^{-x}) \frac{d(e^x + e^{-x})}{dx}$$
$$(e^x + e^{-x})^2$$

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

$$= \frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2}$$

$$= 1 - \frac{(e^z - e^{-z})^2}{(e^z + e^{-z})^2}$$

$$\boxed{g'(x) = 1 - (g(x))^2}$$

\Rightarrow Derivation of ReLU Activation Function :-

1) for $z < 0$

$$g(z) = 0$$

$$\boxed{\frac{d(g(z))}{dz} = 0}$$

2) for $z > 0$

$$g(z) = x$$

$$\boxed{\frac{d(g(z))}{dz} = \frac{d(x)}{dx} = 1}$$

\Rightarrow Derivation of Leaky ReLU :-

1) for $z < 0$

$$g(z) = \alpha x$$

$$\frac{d(g(z))}{dz} = \frac{d(\alpha x)}{dx}$$

$$= \alpha \frac{d(x)}{dx}$$

$$= \alpha \times 1$$

$$\boxed{g'(x) = \alpha}$$

2) for $z > 0$

$$g(z) = x$$

$$\frac{d(g(z))}{dx} = \frac{d(x)}{dx} \Rightarrow$$

$$\boxed{g'(z) = 1}$$

- * These probabilities are distributed across different classes such that their sum equals 1.
- * It is suitable for classification tasks

~~concept~~ Formula :-

$$S_i = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

Example :-

$$x = \begin{bmatrix} 2.0 \\ 1.0 \\ 0.1 \end{bmatrix}$$

Softmax

$$S = \begin{bmatrix} S_1 \\ S_2 \\ S_3 \end{bmatrix}$$

$$S_1 = \frac{e^{x_1}}{e^{x_1} + e^{x_2} + e^{x_3}} = \frac{e^2}{e^2 + e^1 + e^{0.1}} = \frac{7.389}{11.21} = 0.66$$

$$S_2 = \frac{e^{x_2}}{e^{x_1} + e^{x_2} + e^{x_3}} = \frac{e^1}{e^2 + e^1 + e^{0.1}} = \frac{2.71}{11.21} = 0.24$$

$$S_3 = \frac{e^{x_3}}{e^{x_1} + e^{x_2} + e^{x_3}} = \frac{e^{0.1}}{e^2 + e^1 + e^{0.1}} = \frac{1.05}{11.21} = 0.10$$

$$\therefore S = \begin{bmatrix} S_1 \\ S_2 \\ S_3 \end{bmatrix} = \begin{bmatrix} 0.66 \\ 0.24 \\ 0.10 \end{bmatrix}$$

⇒ Derivative of softmax activation function :-

Step 1 :- Define the softmax function

$$\sigma_i = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

$$\text{Assume } S = \sum_{j=1}^n e^{x_j}$$

$$\therefore \sigma_i = \frac{e^{x_i}}{S}$$

We want to compute $\frac{\partial \sigma_i}{\partial x_j}$

Step 2 :- Apply quotient rule

$$\sigma_i = \frac{e^{x_i}}{S}$$

$$\frac{\partial \sigma_i}{\partial x_j} = \frac{\partial}{\partial x_j} \left(\frac{e^{x_i}}{S} \right)$$

$$= \frac{\partial}{\partial x_j} \left(e^{x_i} \cdot \frac{1}{S} \right) \quad \# \text{ express as a product}$$

$$\textcircled{1} \leftarrow = \left(\frac{\partial e^{x_i}}{\partial x_j} \cdot \frac{1}{S} \right) + \left(\frac{\partial}{\partial x_j} \left(\frac{1}{S} \right) \cdot e^{x_i} \right) \quad \# \text{ use prod rule}$$

$$\textcircled{2} \leftarrow \frac{\partial e^{x_i}}{\partial x_j} = \begin{cases} e^{x_i} & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases}$$

~~$$S = \sum_{j=1}^n e^{x_j} \Rightarrow \frac{\partial S}{\partial x_j} = 0 + 0 + \dots e^{x_j}$$~~
$$= e^{x_j}$$

$$\textcircled{3} \leftarrow \frac{\partial}{\partial x_j} \left(\frac{1}{S} \right) = -\frac{1}{S^2} \cdot \frac{\partial S}{\partial x_j} \quad \# \text{ chain rule of differentiation}$$

$$= -\frac{1}{S^2} \cdot e^{x_j}$$

$$= -\frac{e^{x_j}}{S^2} //$$

step 3 : Plug-in ② and ③ into ①

$$\begin{aligned}\frac{\partial \sigma_i}{\partial x_j} &= \left(\frac{\partial}{\partial x_j} (e^{x_i}) \cdot \frac{1}{S} \right) + \left(\frac{\partial}{\partial x_j} \left(\frac{1}{S} \right) \cdot e^{x_i} \right) \\ &= \left(\frac{\partial}{\partial x_j} (e^{x_i}) \cdot \frac{1}{S} \right) + \left(e^{x_i} \cdot -\frac{e^{x_j}}{S^2} \right)\end{aligned}$$

step 4 : Now take separate cases

case - 1 : $i = j$

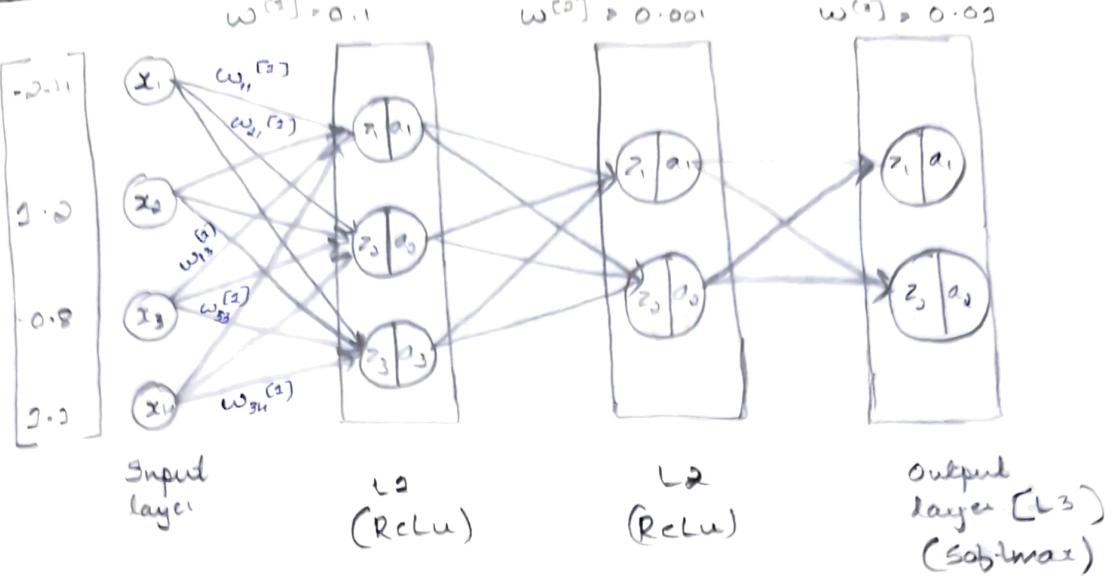
$$\begin{aligned}\frac{\partial (\sigma_i)}{\partial x_j} &= \frac{e^{x_i}}{S} - \frac{e^{x_i} e^{x_i}}{S^2} \\ &= \frac{e^{x_i}}{S} - \frac{e^{x_i} e^{x_i}}{S^2} \\ &= \sigma_i - \sigma_i^2\end{aligned}$$

$$\boxed{\frac{\partial (\sigma_i)}{\partial x_j} = \sigma_i (1 - \sigma_i)}$$

case - 2 : $i \neq j$

$$\begin{aligned}\frac{\partial (\sigma_i)}{\partial x_j} &= \left(0 \cdot \frac{1}{S} \right) + \left(e^{x_i} \cdot -\frac{e^{x_j}}{S^2} \right) \\ &= -\frac{e^{x_i} e^{x_j}}{S^2} \\ &= -\frac{e^{x_i}}{S} \cdot \frac{e^{x_j}}{S} \\ &= -S_i \cdot S_j\end{aligned}$$

$$\boxed{\frac{\partial (\sigma_i)}{\partial x_j} = -\sigma_i \cdot \sigma_j}$$



\Rightarrow Layer 1

$$\begin{aligned}
 z_1 &= (-2.4 \times 0.1) + (1.2 \times 0.1) + (-0.8 \times 0.1) + (1.1 \times 0.1) \\
 &= -0.24 + 0.12 - 0.08 + 0.11 \\
 &> -0.09
 \end{aligned}$$

$$z_2 = \begin{bmatrix} -2.4 \\ 1.2 \\ -0.8 \\ 1.1 \end{bmatrix} \times \begin{bmatrix} 0.1 & 0.1 & 0.1 & 0.1 \end{bmatrix} = -0.09$$

$$z_3 = \begin{bmatrix} -2.4 \\ 1.2 \\ -0.8 \\ 1.1 \end{bmatrix} \times \begin{bmatrix} 0.1 & 0.1 & 0.1 & 0.1 \end{bmatrix} = -0.09$$

after ReLU activation all values in layer 2 will become 0

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

\Rightarrow Layer 2 : L_2

$$z_1 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \times \begin{bmatrix} 0.001 & 0.002 & 0.003 \end{bmatrix} = 0$$

$$z_2 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \times \begin{bmatrix} 0.001 & 0.001 & 0.001 \end{bmatrix} = 0$$

after ReLU activation all values in layer 2 will also become 0

\Rightarrow Layer 3 [Output layer]

$$z_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \times [0.01 \quad 0.02] = 0$$

$$z_2 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \times [0.02 \quad 0.01] = 0$$

$$\alpha_1 = \frac{e^0}{e^0 + e^0} = \frac{0.5}{2} = 0.5$$

$$\alpha_2 = \frac{e^0}{e^0 + e^0} \Rightarrow \frac{e^0}{2e^0} = 0.5$$

After softmax activation all values in the output will become 0.5

$$\text{Output Layer} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$$

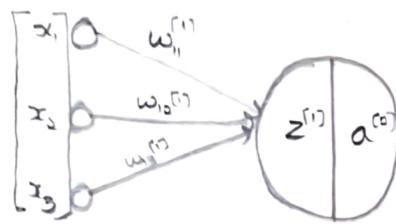
\Rightarrow Sizing :-

- * More the hidden neurons in a layer, more the model will fit to the data
- * Due to the presence of multiple deep layers, what we optimize is not a convex function

\Rightarrow Regularization

- * Overfitting in a neural network can be controlled by applying regularization

\Rightarrow Perception :-



$$z^{(1)} = w_{11}^{(1)} x_1 + w_{12}^{(1)} x_2 + w_{13}^{(1)} x_3 + b^{(1)}$$

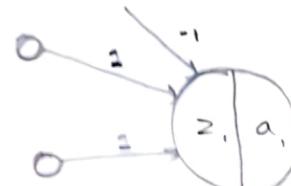
$$a^{(1)} = g(z^{(1)})$$

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

- * Single layer Neural Network with one Neuron is called as a perceptron

\Rightarrow AND

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



$$z_1 = w_{11}^{(1)} x_1 + w_{12}^{(1)} x_2 + b_1$$

$$a_1 = \begin{cases} 1 & \text{if } z_1 > 0 \\ 0 & \text{if } z_1 \leq 0 \end{cases}$$

a) $z_1 = 0 \times 1 + 0 \times 1 + -1 = -1$

$a_1 = 0 (\because z_1 \leq 0)$

b) $z_1 = 0 \times 1 + 1 \times 1 + -1 = 0$

$a_1 = 0 (\because z_1 \leq 0)$

c) $z_1 = 1 \times 1 + 0 \times 1 + -1 = 0$

$a_1 = 0 (\because z_1 \leq 0)$

d) $z_1 = 1 \times 1 + 1 \times 1 + -1 = 1$

$a_1 = 1 (\because z_1 > 0)$

\Rightarrow OR

a) $0 \times 1 + 0 \times 1 + -1 = -1$

$a = 0 (\because z \leq 0)$

b) $= 0$

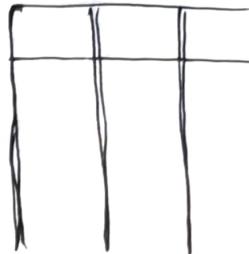
$a = 0 (\because z \leq 0)$

c) $= 0$

$a = 0 (\because z \leq 0)$

d) $= 0 \ 1$

$a = 1 (\because z > 0)$



$$y = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \times , \text{ this is not right for OR gate}$$

change bias to 0

a) $0 \times 1 + 0 \times 1 + 0 = 0$

$a = 0 (\because z \leq 0)$

b) $0 \times 1 + 1 \times 1 + 0 = 1$

$a = 1 (\because z > 0)$

c) $1 \times 1 + 0 \times 1 + 0 = 1$

$a = 1 (\because z > 0)$

d) $1 \times 1 + 1 \times 1 + 0 = 2$

$a = 1 (\because z > 0)$

\Rightarrow XOR

a) $0 \times 1 + 0 \times 1 + 0 = 0$

$a = 0$

b) $0 \times 1 + 1 \times 1 + 0 = 1$

$a = 1$

c) $1 \times 1 + 0 \times 1 + 0 = 1$

$a = 1$

d) $1 \times 1 + 1 \times 1 + 0 = 2$

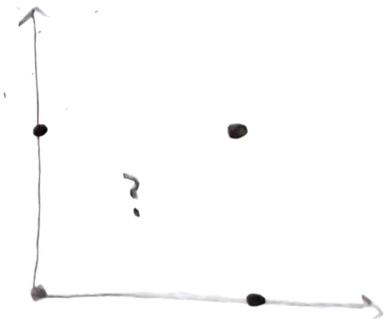
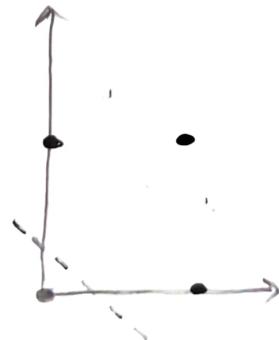
$a = 1$

Truth table

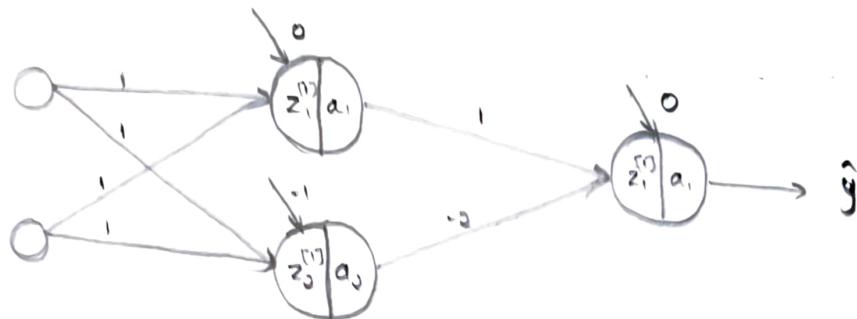
x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0

$$y = \begin{bmatrix} 0 \\ \frac{1}{2} \\ 1 \end{bmatrix} \times$$

* A linear boundary cannot be drawn for XOR because XOR is non linear



* A perceptron can model only linear combinations and therefore by introducing an hidden layer we can model non-linear function



$$z_1^{(1)} = 0 \times 1 + 0 \times 1 + 0 = 0 \quad a = 0 \quad (z \leq 0)$$

~~$$z^{(1)} = \omega^{(1)} \times a^{(0)} + b^{(1)}$$~~

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

$$z^{(1)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 0 \\ -1 \end{bmatrix} \quad a^{(1)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$z^2 = \omega^{(2)} \times a^{(1)} + b^{(2)}$$

~~$$= \begin{bmatrix} 1 & -2 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{bmatrix} 0 \\ -1 \end{bmatrix} + \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$~~

~~$$= \begin{bmatrix} 1 \\ -2 \end{bmatrix} \Rightarrow a = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$~~

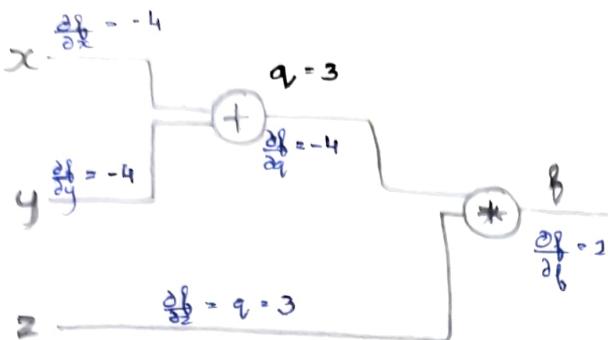
$$z^2 = [1 \ -2] \begin{bmatrix} 0 \\ 0 \end{bmatrix} + 0 = 0$$

$$a^2 = 0 \quad (\because z \leq 0)$$

$$\therefore y = 0 \quad (\because \text{Row 1 is correct})$$

Example

$$f = (x+y)^2$$



$$\text{let } q_1 = (x+y)$$

$$f = q_1^2$$

$$\frac{\partial f}{\partial x} = 1 \quad \frac{\partial f}{\partial y} = 1$$

$$\frac{\partial f}{\partial z} = 2 \quad \frac{\partial f}{\partial v} = 9$$

Computational graph :-

$$x = -2 \quad y = 5 \quad z = 4 \quad (\text{Random initialization})$$

$$f = (x+y)^2 = (-2+5) \cdot 4$$

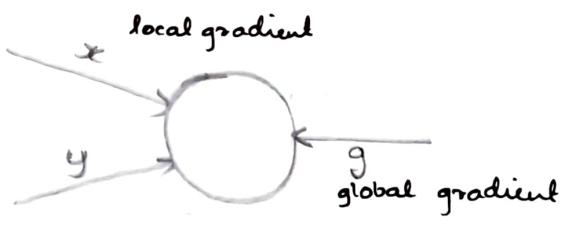
$$\begin{aligned} &= 12 \\ \frac{\partial f}{\partial x} &= 1 \\ \frac{\partial f}{\partial y} &= 1 \\ \frac{\partial f}{\partial z} &= 4 \end{aligned}$$

global gradient local gradient
 $\frac{\partial f}{\partial q_k} = \boxed{\frac{\partial f}{\partial q_k}} \quad \boxed{\frac{\partial q_k}{\partial x}}$
 chain rule

$$= -4 \times 1 = -4$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \cdot \frac{\partial q}{\partial y} = -4 \times 1 = -4$$

$$\frac{\partial f}{\partial z} = q = 3$$



$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} * \frac{\partial q}{\partial x}$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} * \frac{\partial q}{\partial y}$$

Example 2

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

$$x = -2 \quad y = 5 \quad z = -4$$

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

- ① Identify additional functions

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

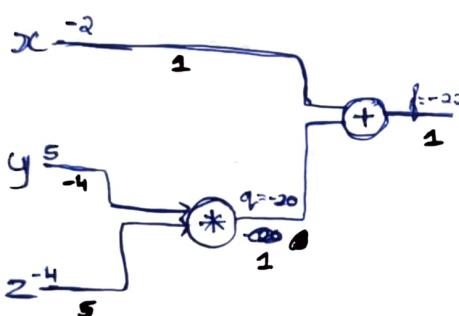
$$q = yz$$

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

$$\frac{\partial f}{\partial x} = 1 \quad \frac{\partial f}{\partial q} = 1$$

$$\frac{\partial q}{\partial y} = z \quad \frac{\partial q}{\partial z} = y$$

- ② Draw a computational graph :-



- ④ Backward pass :-

$$\begin{aligned} \frac{\partial f}{\partial x} &= 1 \\ \frac{\partial f}{\partial y} &= 1 \quad , \quad \frac{\partial f}{\partial z} = 1 \\ \frac{\partial f}{\partial y} &= \left[\frac{\partial f}{\partial q} \right] * \left[\frac{\partial q}{\partial y} \right] = 1 * -4 = -4 \\ \frac{\partial f}{\partial z} &= \left[\frac{\partial f}{\partial q} \right] * \left[\frac{\partial q}{\partial z} \right] = 1 * 5 = 5 \end{aligned}$$

- ③ Forward pass :-

$$q = yz = -20$$

$$f = x + q = -22$$

$$\begin{aligned}w_0 &= 2 \\x_0 &= -1 \\w_1 &= -3 \\x_1 &= -2 \\w_2 &= -3\end{aligned}$$

① Identify additional functions

$$f_2 = w_0 x_0 \quad f_3 = w_1 x_1 \quad f_4 = w_2$$

$$f(\omega, x) = \frac{1}{1 + e^{-(f_1 + f_2 + f_3)}}$$

~~$$f_4 = x_1 + e^{-(f_1 + f_2 + f_3)}$$~~
~~$$f(\omega, x) = \frac{1}{1 + f_4}$$~~

~~$$f_4 = x_1 e^{-(f_1 + f_2 + f_3)}$$~~
~~$$f(\omega, x) = \frac{1}{1 + f_4}$$~~

② Identify additional functions :

$$f_2 = w_0 x_0 \quad f_3 = w_1 x_1 \quad \frac{1}{1 + e^{-(f_1 + f_2 + f_3)}}$$

$$\frac{\partial f_1}{\partial x_0} = w_0, \quad \frac{\partial f_2}{\partial x_1} = w_1$$

$$f_5 = f_2 + f_3 \quad f_4 = f_3 + w_2$$

$$\frac{\partial f_1}{\partial w_2} = x_0, \quad \frac{\partial f_2}{\partial w_1} = x_1$$

$$f_6 = -f_4$$

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

$$f_7 = 1 + f_6$$

$$\frac{\partial f_4}{\partial f_3} = 1, \quad \frac{\partial f_4}{\partial w_2} = 1$$

$$f_8 = \frac{1}{f_7}$$

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

$$f(\omega, x) = f_8$$

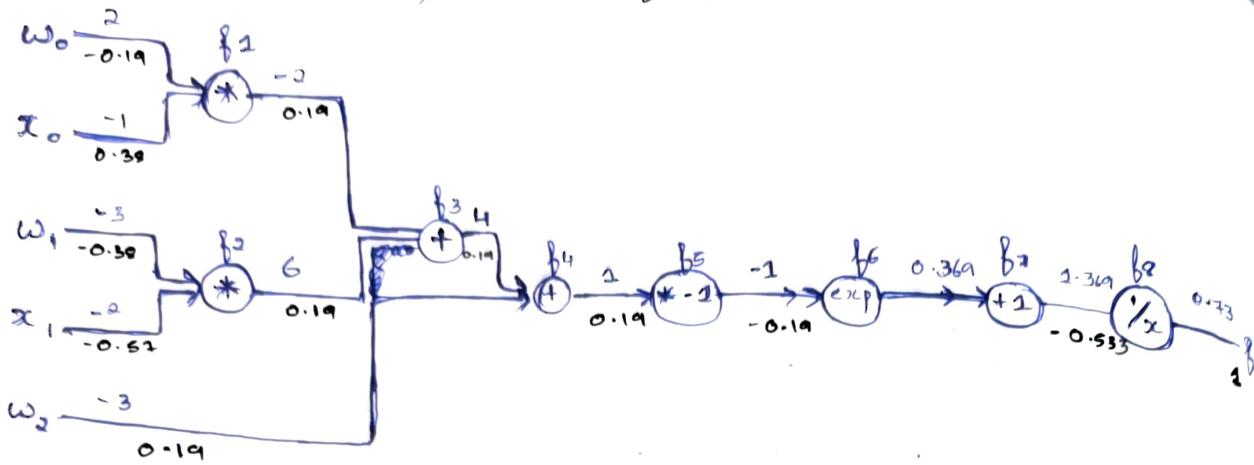
$$\frac{\partial f_6}{\partial f_5} = e^{f_5}$$

$$\frac{\partial f_7}{\partial f_6} = 1$$

$$\frac{\partial f_8}{\partial f_7} = 0 - \frac{1}{f_7^2}$$

$$\frac{\partial f}{\partial f_8} = 1$$

② Draw a computation graph



③ Forward pass :

$$f_1 = w_0 x_0 = -2$$

$$f_2 = w_1 x_1 = +6$$

$$f_3 = -2 + 6 = 4$$

$$f_4 = \cancel{4} - 3 = 1$$

$$f_5 = \cancel{1} + 1 = -1$$

$$f_6 = e^{-2} \cancel{1} = \frac{1}{e} = 0.369$$

$$f_7 = 2 \cdot 0.369$$

$$f_8 = 0.73$$

④ Backward Pass :

$$\frac{\partial f}{\partial f_8} = 1$$

$$\frac{\partial f_8}{\partial f_7} = \frac{1}{2 \cdot 0.369^2} = -0.534 \quad \left(\text{local * global} \right) = \frac{1}{f_7} * 1 = -0.534$$

$$\frac{\partial f_8}{\partial f_6} = \text{local * global} = \frac{\partial f_6}{\partial f_5} * -0.54 = e^{f_5} * -0.54 = e^{-1} * -0.54 = -0.19$$

$$\frac{\partial f_8}{\partial f_5} = \text{local * global} = \frac{\partial f_5}{\partial f_4} * -0.54 = e^{f_4} * -0.54 = e^{-2} * -0.54 = -0.19$$

$$\frac{\partial f_0}{\partial w_4} = \frac{\partial f_5}{\partial f_4} * -0.19 = -1 * -0.19 = 0.19$$

$$\frac{\partial f_0}{\partial w_3} = \frac{\partial f_4}{\partial w_2} * 0.19 = 1 * 0.19 = 0.19$$

$$\frac{\partial f_0}{\partial f_3} = \frac{\partial f_4}{\partial f_3} * 0.19 = 1 * 0.19 = 0.19$$

$$\frac{\partial f_0}{\partial f_2} = \frac{\partial f_3}{\partial f_2} * 0.19 = 1 * 0.19 = 0.19$$

$$\frac{\partial f_0}{\partial f_1} = \frac{\partial f_3}{\partial f_1} * 0.19 = 1 * 0.19 = 0.19$$

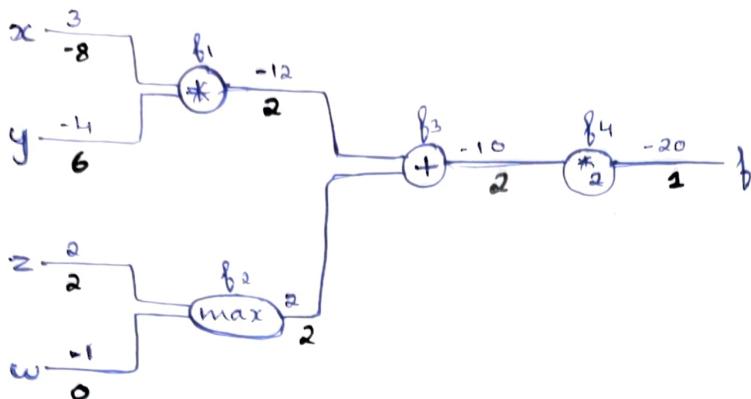
$$\frac{\partial f_0}{\partial x_1} = \frac{\partial f_2}{\partial x_1} * 0.19 = w_1 * 0.19 = -3 * 0.19 = -0.57$$

$$\frac{\partial f_0}{\partial w_1} = \frac{\partial f_2}{\partial w_1} * 0.19 = x_1 * 0.19 = -2 * 0.19 = -0.38$$

$$\frac{\partial f_0}{\partial w_0} = \frac{\partial f_1}{\partial w_0} * \cancel{0.19} = -1 * 0.19 = -0.19$$

$$\frac{\partial f_0}{\partial x_0} = \frac{\partial f_1}{\partial x_0} * 0.19 = w_0 * 0.19 = 2 * 0.19 = 0.38$$

Example - 4



① Identify additional functions :-

$$f_1 = x * y = -12 \quad \frac{\partial f_1}{\partial x} = y \quad \frac{\partial f_1}{\partial y} = x$$

$$f_2 = \max(z, w) = 2 \quad \frac{\partial f_2}{\partial z} = 0$$

$$f_3 = f_1 + f_2 = -10 \quad \frac{\partial f_3}{\partial f_1} = 1 \quad \frac{\partial f_3}{\partial f_2} = 1$$

$$f_4 = 2 * f_3 = -20 \quad \frac{\partial f_4}{\partial f_3} = 2$$

④ Backward Pass

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

$$\frac{\partial f_3}{\partial f_4} = \cancel{\frac{\partial f_3}{\partial f_4}} \quad \frac{\partial f_4}{\partial f_3} * \text{global} = 2 * 1 = \underline{\underline{2}}$$

$$\frac{\partial f_4}{\partial f_1} = \frac{\partial f_3}{\partial f_1} * 2 = 1 * 2 = 2$$

$$\frac{\partial f_4}{x} = \frac{\partial f_1}{x} * 2 = y * 2 = -8$$

$$\frac{\partial f_4}{y} = \frac{\partial f_1}{y} * 2 = x * 2 = 6$$

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

$$\frac{\partial f_2}{\partial z} = \begin{cases} 1 & \text{if } z > w \\ 0 & \text{if } z < w \\ \text{undefined or arbitrary} & \text{if } z = w \end{cases}$$

$$\frac{\partial f_2}{\partial w} = \begin{cases} 1 & \text{if } w > z \\ 0 & \text{if } w < z \\ \text{undefined or arbitrary} & \text{if } z = w \end{cases}$$

\Rightarrow Scale and Shift

- * Mean of the Gaussian will not be 0 and the variance will also be different

\rightarrow How much to shift?

- * Scale = γ and shift = β and γ and β are learnable parameters

- * $\gamma^{(l)}$ $\beta^{(l)}$: Scale & shift will be applied to each layer

Example:

$$z = [4, 8, 6, 10] \Rightarrow \text{No. of samples in a batch}$$

$$\mu = \underline{\underline{7}}$$

$$\sigma^2 = \frac{1}{4} \sum_{i=1}^4 (z_i - \mu)^2 = \frac{(4-7)^2 + (8-7)^2 + (6-7)^2 + (10-7)^2}{4+e} = e = 10^{-5}$$

$$= \frac{(-3)^2 + 1^2 + 1^2 + 3^2}{4}$$

$$\sigma = \underline{\underline{5}}$$

$\hat{z} = \gamma z + \beta$

$$\gamma = 2$$

$$\beta = 1$$

$\hat{z} = \frac{z - \mu}{\sigma}$

~~ẑ~~

$$\hat{z}_1 = 2(4) + 1 = 9.0$$

$$\hat{z}_2 = 2(8) + 1 = 17.0$$

$$\hat{z}_3 = \dots$$

$$\hat{z}_4 = \dots$$

→ Inverted Dropout :-



train

- 1) Sample 1 → output = 2
- 2) Sample 2 → dropout → output = 0
- 3) Sample 3 → dropout → output = 0
- 4) Sample 4 → output = 2

$$\text{Total output} = \frac{4}{4}$$

test

- 1) Sample 1 → output = 2
- 2) Sample 2 → output = 2
- 3) Sample 3 → output = 2
- 4) Sample 4 → output = 2

$$\text{Total output} = \frac{8}{4}$$

* To achieve the test output during training we scale the total output ~~of~~ by p

$$\therefore \text{Total output} = \frac{\text{Total output (train)}}{p}$$

* Since an FNN cannot be applied, we build something different called 'Convolutional Neural Network'.

⇒ Convolution operation :- ②

* It is the heart of a CNN

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

Input 6x6 info

*

Computation operator

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

filter / kernel
3x3

=

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

Super-impose filter on the info and then multiply the value (ie perform convolution operation) of the filter with the info

= 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
$$= -5$$

* Now add this value to ...

$6 \times 6 \times 32$

what it looks like :- for each pixel

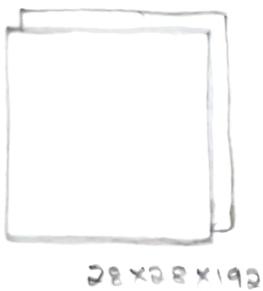


O

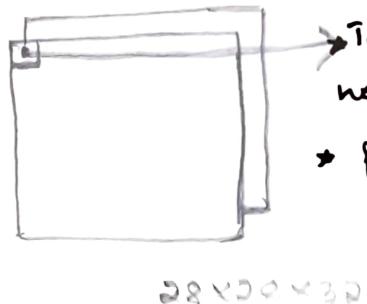
1x32

MLP

- weight sharing :- we tell the filter to slide across the whole image and learn what its assigned to.



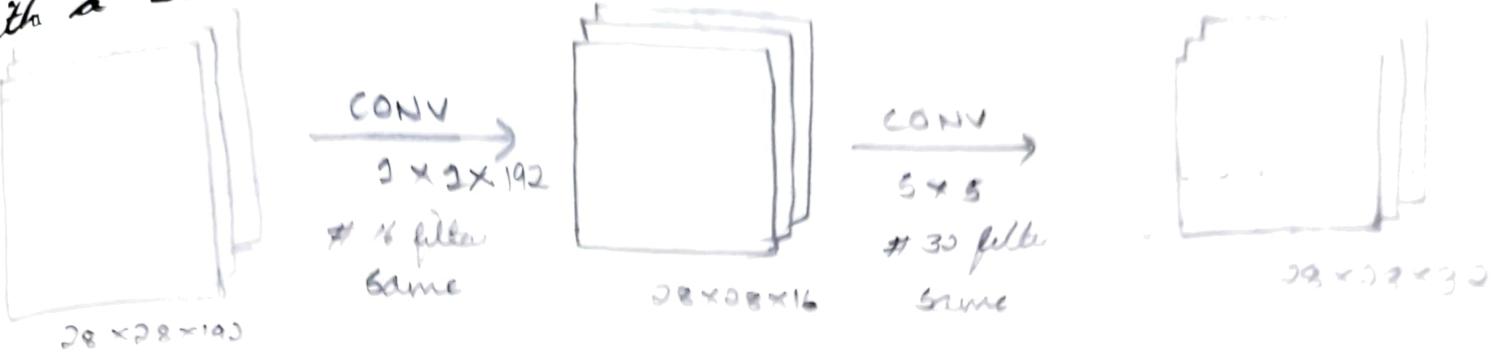
CONV
 $5 \times 5 \times 192$ filter
filters = 32
same convolution



→ to get this one pixel we need to do $5 \times 5 \times 192$ operations
* for whole image then;
 $(5 \times 5 \times 192) \times (28 \times 28 \times 32)$

$$\# \text{Total operations} = \underbrace{5 \times 5 \times 192}_{\text{filter size}} \times \underbrace{28 \times 28}_{\text{No. of stride required to cover the whole picture}} \times \underbrace{32}_{\text{No. of filters}} = 120422400$$

\Rightarrow with a 1×1 conv :-



Total operations

$$\left(\frac{1}{2} \times 1 \times 192 \times 28 \times 28 \right) \times 16 \times 16 = 1,422,016 \approx 2 \text{ million}$$

$$28 = n + 2p - f + 1$$

$$\rightarrow 28 + 2p - 5 + 1$$

$$L = 20$$

$$D = 2$$

$$2) 5 \times 5 \times 16 \times 28 \times 28 \times 32 = 10,035,200 \approx 10, \text{million}$$

$$\text{Total} \approx 12.5 \text{ million}$$

* Almost $\frac{1}{10}$ of the total operations as that of model without 1×1 conv.

\Rightarrow Modifying height and width :-

- 1) Padding .
- 2) Striding
- 3) Pooling

\Rightarrow Modify channels :-

- 1) No of filters
- 2) 1×1 CONV

$T = 2$ (Two time steps)

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

input dimension = 2 $\rightarrow x_t \in \mathbb{R}^2$

hidden state dimensions = 3 $\rightarrow h_{(t)} \in \mathbb{R}^3$

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

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

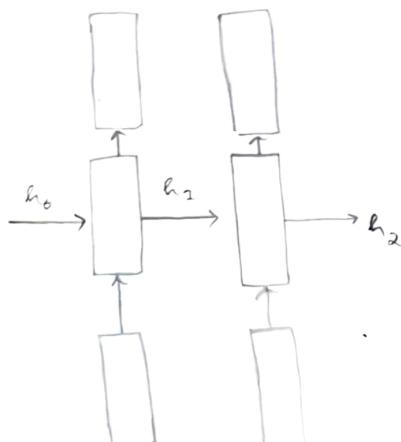
$$w_{2h} = \begin{bmatrix} 0.5 & -0.3 \\ 0.2 & 0.2 \\ 0.2 & 0.4 \end{bmatrix}_{3 \times 2}$$

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

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

$$h_1 = w_{hh} * h_0 + w_{hx} x_1$$

$$\begin{aligned} h_{(1)} &= (3 \times 3) * (3 \times 1) + 3 \times 2 \\ &\text{Both result must be } (3 \times 1) \\ &= (3 \times 2) \cdot (3 \times 1) \\ &= (3 \times 2) \end{aligned}$$



$$\begin{aligned} y^{(1)} &= w_{hy} * h_1 \\ y^{(1)}_{2 \times 1} &= 0.5 * (3 \times 1) \\ &= (2 \times 3) = 2 \end{aligned}$$

$$y^{(1)} = (2 \times 2)$$

\Rightarrow Time-step 2, $t = 1$

$$h_1 = \underbrace{(w_{hh} * h_0) + (w_{hx} x_1)}$$

$$\begin{aligned} w_{hh} * h_0 &= \begin{bmatrix} 0.1 & 0.4 & 0.0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.2 \end{bmatrix} * \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} w_{hx} x_1 &= \begin{bmatrix} 0.5 & -0.3 \\ 0.2 & 0.2 \\ 0.1 & 0.4 \end{bmatrix} * \begin{bmatrix} 1 \\ 2 \end{bmatrix} \\ &= \begin{bmatrix} -0.1 \\ 1.2 \\ 0.9 \end{bmatrix} \end{aligned}$$

$$h_1 = \text{Tanh} \left[\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} -0.1 \\ 1.2 \\ 0.9 \end{bmatrix} \right]$$

$$\Rightarrow \text{Tanh} \left(\begin{bmatrix} -0.1 \\ 1.2 \\ 0.9 \end{bmatrix} \right)$$

$$h_1 = \begin{bmatrix} -0.099 \\ 0.834 \\ 0.71 \end{bmatrix}$$

$$y^{(1)} = \omega_{HY} + h_1$$

$$= \begin{bmatrix} 1 & -2 & 0.5 \\ 0.5 & 0.5 & -0.5 \end{bmatrix} * \begin{bmatrix} -0.099 \\ 0.834 \\ 0.71 \end{bmatrix}$$

$$= \begin{bmatrix} -0.58 \\ 0.01 \end{bmatrix}$$

\Rightarrow Time-step 2, $t=2$

$$h_2 = \text{Tanh} (\omega_{HH} * h_1) + (\omega_{HX} * x_2)$$

$$\omega_{HH} * h_1 = \begin{bmatrix} 0.1 & 0.4 & 0.0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.0 \end{bmatrix} * \begin{bmatrix} -0.099 \\ 0.834 \\ 0.71 \end{bmatrix}$$

$\omega_{HX} * x_2 = \begin{bmatrix} 0.2 & 0.4 & 0.0 \\ 0.0 & 0.3 & 0.2 \\ 0.0 & 0.0 & 0.1 \end{bmatrix} * \begin{bmatrix} -0.8 \\ 1.2 \\ 0.9 \end{bmatrix}$

$$= \begin{bmatrix} 0.5 & -0.3 \\ 0.8 & 0.2 \\ 0.2 & 0.4 \end{bmatrix} * \begin{bmatrix} -1 \\ 2 \end{bmatrix}$$

$$= \begin{bmatrix} 0.8 \\ -0.6 \\ 0.3 \end{bmatrix}$$

$$h_2 \Rightarrow \text{Tanh} \left(\left(\begin{bmatrix} 0.3237 \\ 0.4112 \\ 0.0536 \end{bmatrix} + \begin{bmatrix} -0.8 \\ -0.6 \\ 0.3 \end{bmatrix} \right) \right)$$

$$\Rightarrow \text{Tanh} \left(\begin{bmatrix} -0.47 \\ 0.19 \\ 0.3536 \end{bmatrix} \right)$$

$$h_2 = \begin{bmatrix} -0.66 \\ -0.19 \\ 0.337 \end{bmatrix}$$

Example :-

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

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

$$C_{t-1} = [0.2, -0.2]$$

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

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

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

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

$$\omega_{zo} = \begin{bmatrix} 0.3 & 0.25 \\ -0.2 & 0.1 \end{bmatrix}$$

$$\omega_{eo} = \begin{bmatrix} 0.15 & 0.05 \\ 0.2 & 0.2 \end{bmatrix}$$

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

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

$$h_t = ?$$

$$C_t = ?$$

$$\hat{f}_t = \alpha(\omega_{if} h_{t-1} + \omega_{xf} x_t)$$

$$\begin{aligned} \omega_{if} h_{t-1} &= \begin{bmatrix} 0.05 & -0.1 \\ 0.2 & 0.3 \end{bmatrix} * \begin{bmatrix} 0.0 \\ 0.2 \end{bmatrix} \\ &\Rightarrow \begin{bmatrix} -0.01 \\ 0.02 \end{bmatrix}_{2 \times 1} \end{aligned}$$

$$\begin{aligned} \omega_{xf} x_t &= \begin{bmatrix} -0.4 & 0.2 \\ 0.3 & 0.3 \end{bmatrix} * \begin{bmatrix} 0.5 \\ -0.1 \end{bmatrix} \\ &\Rightarrow \begin{bmatrix} -0.22 \\ 0.12 \end{bmatrix} \end{aligned}$$

$$(\omega_{if} h_{t-1} + \omega_{xf} x_t) = \begin{bmatrix} -0.02 \\ 0.02 \end{bmatrix} + \begin{bmatrix} -0.22 \\ 0.12 \end{bmatrix} = \begin{bmatrix} -0.23 \\ 0.13 \end{bmatrix}$$

$$\alpha(\omega_{if} h_{t-1} + \omega_{xf} x_t) = \alpha \begin{bmatrix} -0.23 \\ 0.13 \end{bmatrix} = \begin{bmatrix} 0.443 \\ 0.532 \end{bmatrix} = \hat{f}_t \quad \text{--- ①}$$

$$\Rightarrow i_t = \sigma(\omega_{h_0} h_{t-1} + \omega_{x_1} x_t)$$

$$\begin{aligned}\omega_{h_0} h_{t-1} &= \begin{bmatrix} 0.2 & 0.2 \\ -0.2 & 0.05 \end{bmatrix} * \begin{bmatrix} 0.0 \\ 0.1 \end{bmatrix} & \omega_{x_1} x_t &= \begin{bmatrix} 0.5 & -0.3 \\ 0.4 & 0.2 \end{bmatrix} * \begin{bmatrix} 0.5 \\ -0.1 \end{bmatrix} \\ &= \begin{bmatrix} 0.02 \\ 0.005 \end{bmatrix} & &= \begin{bmatrix} 0.28 \\ 0.19 \end{bmatrix}\end{aligned}$$

$$\begin{aligned}(\omega_{h_0} h_{t-1} + \omega_{x_1} x_t) &= \begin{bmatrix} 0.02 \\ 0.005 \end{bmatrix} + \begin{bmatrix} 0.28 \\ 0.19 \end{bmatrix} \\ &= \begin{bmatrix} 0.30 \\ 0.195 \end{bmatrix}\end{aligned}$$

$$\sigma(\omega_{h_0} h_{t-1} + \omega_{x_1} x_t) = \sigma \begin{bmatrix} 0.30 \\ 0.195 \end{bmatrix} = \begin{bmatrix} 0.575 \\ 0.549 \end{bmatrix} = i_t \longrightarrow ②$$

$$\Rightarrow o_t = \sigma(\omega_{h_0} h_{t-1} + \omega_{x_0} x_t)$$

$$\begin{aligned}\omega_{h_0} h_{t-1} &= \begin{bmatrix} 0.15 & 0.05 \\ 0.2 & 0.2 \end{bmatrix} * \begin{bmatrix} 0.0 \\ 0.1 \end{bmatrix} & \omega_{x_0} x_t &= \begin{bmatrix} 0.3 & 0.25 \\ -0.2 & 0.2 \end{bmatrix} * \begin{bmatrix} 0.5 \\ -0.1 \end{bmatrix} \\ &= \begin{bmatrix} 0.005 \\ 0.07 \end{bmatrix} & &= \begin{bmatrix} 0.125 \\ -0.12 \end{bmatrix}\end{aligned}$$

$$\begin{aligned}(\omega_{h_0} h_{t-1} + \omega_{x_0} x_t) &= \begin{bmatrix} 0.005 \\ 0.07 \end{bmatrix} + \begin{bmatrix} 0.125 \\ -0.12 \end{bmatrix} & \sigma(\omega_{h_0} h_{t-1} + \omega_{x_0} x_t) &= \sigma \begin{bmatrix} 0.130 \\ -0.10 \end{bmatrix} \\ &= \begin{bmatrix} 0.130 \\ -0.10 \end{bmatrix} & &= \begin{bmatrix} 0.53 \\ 0.47 \end{bmatrix} = 0.7 \longrightarrow ③\end{aligned}$$

$$\Rightarrow g_t = \tanh(\omega_{\text{ag}} h_{t-1} + \omega_{\text{xg}} x_t)$$

$$\left. \begin{aligned} \omega_{\text{ag}} h_{t-1} &= \begin{bmatrix} 0.2 & 0.1 \\ -0.2 & 0.05 \end{bmatrix} * \begin{bmatrix} 0.0 \\ 0.1 \end{bmatrix} & \omega_{\text{xg}} x_t &= \begin{bmatrix} 0.5 & 0.4 \\ 0.2 & -0.3 \end{bmatrix} * \begin{bmatrix} 0.5 \\ -0.1 \end{bmatrix} \\ &= \begin{bmatrix} 0.02 \\ 0.005 \end{bmatrix} & &= \begin{bmatrix} -0.29 \\ 0.13 \end{bmatrix} \end{aligned} \right\}$$

$$\begin{aligned} (\omega_{\text{ag}} h_{t-1} + \omega_{\text{xg}} x_t) &= \begin{bmatrix} 0.02 \\ 0.005 \end{bmatrix} + \begin{bmatrix} -0.29 \\ 0.13 \end{bmatrix} \\ &= \begin{bmatrix} -0.28 \\ 0.135 \end{bmatrix} \end{aligned}$$

$$\sigma(\omega_{\text{ag}} h_{t-1} + \omega_{\text{xg}} x_t) = \tanh \begin{bmatrix} -0.28 \\ 0.135 \end{bmatrix} = \begin{bmatrix} 0.07 \\ 0.134 \end{bmatrix} \quad \textcircled{a}$$

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

$$f_t \odot c_{t-1} = \begin{bmatrix} 0.443 \\ 0.532 \end{bmatrix} \odot \begin{bmatrix} 0.2 \\ -0.1 \end{bmatrix} = \begin{bmatrix} 0.0886 \\ -0.1064 \end{bmatrix} \quad \textcircled{b}$$

$$i_t \odot g_t = \begin{bmatrix} 0.575 \\ 0.549 \end{bmatrix} \odot \begin{bmatrix} -0.28 \\ 0.135 \end{bmatrix} = \begin{bmatrix} 0.161 \\ 0.074 \end{bmatrix} \quad \textcircled{c}$$

$$c_t = \begin{bmatrix} 0.088 \\ -0.106 \end{bmatrix} + \begin{bmatrix} 0.161 \\ 0.074 \end{bmatrix}$$

$$= \begin{bmatrix} -0.02 \\ -0.02 \end{bmatrix} \quad \textcircled{d}$$

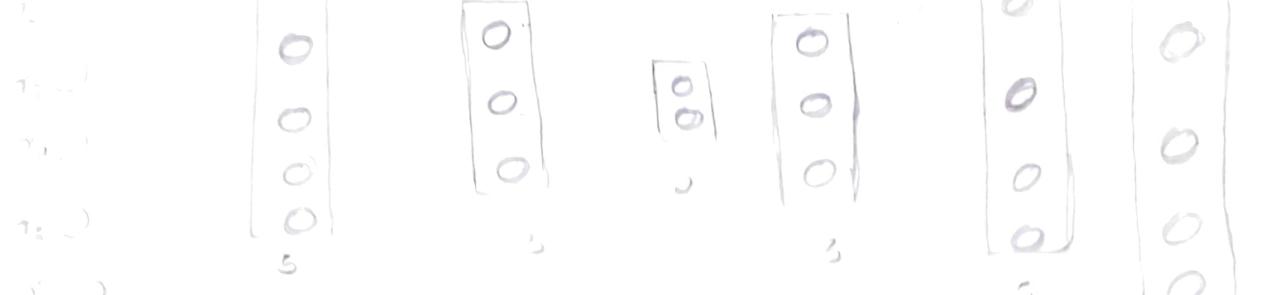
$$\Rightarrow h_2 = o_2 \circ \tanh(c_2)$$

$$\tanh(c_2) = \tanh \begin{bmatrix} -0.07 \\ -0.03 \end{bmatrix}$$

$$= \begin{bmatrix} -0.07 \\ -0.03 \end{bmatrix}$$

$$h_2 = \begin{bmatrix} 0.53 \\ 0.47 \end{bmatrix} * \begin{bmatrix} -0.07 \\ -0.03 \end{bmatrix}$$

$$= \begin{bmatrix} -0.0371 \\ -0.0242 \end{bmatrix}$$



ENCODER DECODER
Must be mirror of the other side of the middle range

→ Example :-

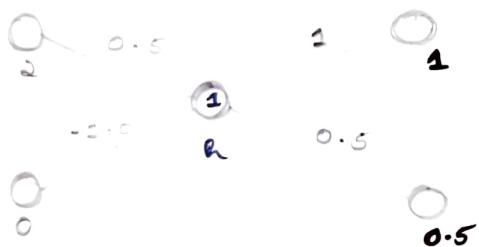
$$x \in \mathbb{R}^2, x = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$$

$$h \in \mathbb{R}^1 \Rightarrow ?$$

$$w_c = ? \Rightarrow 1 \times 2 \Rightarrow [0.5, -1.0]$$

$$w_d = ? \Rightarrow 2 \times 1 \Rightarrow \begin{bmatrix} 2 \\ 0.5 \end{bmatrix}$$

$$MSE = \frac{1}{2} \| \hat{x} - x \|^2$$



$$h = w_c \cdot x$$

$$= 2 \times 0.5 + 0 \times -1.0$$

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

$$\text{decoder} \rightarrow \begin{bmatrix} 1 \\ 0.5 \end{bmatrix} * 2 \quad (w_d * h)$$

$$\rightarrow \begin{bmatrix} 1 \\ 0.5 \end{bmatrix}$$

$$L = \frac{1}{2} (1 - \cancel{0.5})^2 + (0.5 - \cancel{0.5})^2$$

0.625

$$= \frac{1}{2} (1^2) + (0.5)^2$$

$$= \frac{1 + 0.25}{2}$$

$$= 0.625$$

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

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

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

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

$c_2 = ?$

$$c_2 = \sum_{j=1}^3 \alpha_{2,j} h_j$$

$$= \alpha_{2,1} h_1 + \alpha_{2,2} h_2 + \alpha_{2,3} h_3$$

$$\alpha_{2,1} = \frac{\exp(\text{Score}(s_{t-1}, h_1))}{\exp(\text{Score}(s_{t-1}, h_1)) + \exp(\text{Score}(s_{t-1}, h_2)) + \exp(\text{Score}(s_{t-1}, h_3))}$$

~~$$\text{Score}(s_{t-1}, h_1) = [1, 0, 1] \cdot [1, 0, 1]$$~~

$$= 2$$

$$\text{Score}(s_{t-1}, h_2) = [1, 0, 1] \cdot [0, 1, 1]$$

$$= 1$$

$$\text{Score}(s_{t-1}, h_3) = [1, 0, 1] \cdot [1, 1, 0]$$

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

$$\alpha_{2,1} = \frac{e^2}{e^2 + e^1 + e^1} = \frac{7.4}{7.4 + 2.71 + 2.71} = \frac{7.4}{12.82} = 0.598$$

$$\alpha_{2,2} = \frac{e^1}{e^2 + e^1 + e^1} = \frac{2.71}{12.82} = 0.21$$

$$\alpha_{2,3} = \frac{e^1}{e^2 + e^1 + e^1} = \frac{2.71}{12.82} = 0.21$$

$$\begin{aligned}
 C_2 &= \alpha_{2,1} h_1 + \alpha_{2,2} h_2 + \alpha_{2,3} h_3 \\
 &= 0.58 [2, 0, 2] + 0.2 [0, 2, 2] + 0.2 [2, 2, 0] \\
 &= [0.58, 0, 0.58] + [0, 0.2, 0.2] + [0.2, 0.2, 0] \\
 &\Rightarrow [0.78, 0.4, 0.78]
 \end{aligned}$$

→ This is more close to h_2 , it means, it paying more attention to h_2 .

Thinking Machine

Thinking

Machines

① $q_1 \Rightarrow$ Thinking

$k_1 \Rightarrow$ Thinking

① $q_1 \Rightarrow$ Machine

$k_1 \Rightarrow$ Thinking

② $a_1 \Rightarrow$ Thinking

$k_1 \Rightarrow$ Machines

$a_1 \Rightarrow$ Machine

$k_1 \Rightarrow$ Machines

⇒ Example

Input :- Playing outside $\Rightarrow z_1, z_2$

Playing

Outside

$$q_1 = [0.212, 0.04, 0.63, 0.36]^T$$

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

$$k_1 = [0.31, 0.84, 0.963, 0.57]^T$$

$$k_2 = [0.45, 0.94, 0.73, 0.56]^T$$

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

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

Playing

dot prod.

$$\textcircled{1} q_1 * k_1 = \cancel{0.91} 0.91$$

Outside

$$\textcircled{1} q_2 * k_2 = 1.25$$

$$\textcircled{2} q_1 * k_2 = 0.80$$

$$\textcircled{2} q_2 * k_1 = 1.41$$

⇒ divide by $\sqrt{R^F} = \sqrt{4} = 2$

$$\textcircled{1} 0.455$$

$$\textcircled{1} 0.625$$

$$\textcircled{2} 0.40$$

$$\textcircled{2} 0.705$$

⇒ softmax

$$\textcircled{1} 0.51$$

$$\textcircled{1} 0.48$$

$$\textcircled{2} 0.49$$

$$\textcircled{2} 0.52$$

$$\text{softmax}_1 + v_1 + \text{softmax}_2 * v_2$$

$$v_{12} = [0.18, 0.42, 0.05, 0.35]$$

$$v_2 = [0.15, 0.18, 0.09, 0.35]$$

$$0.48 * v_2$$

$$v_{21} = [0.15, 0.17, 0.09, 0.35]$$

$$v_{22} = [0.19, 0.13, 0.05, 0.20]$$

$$0.52 * [v_1]$$

sum

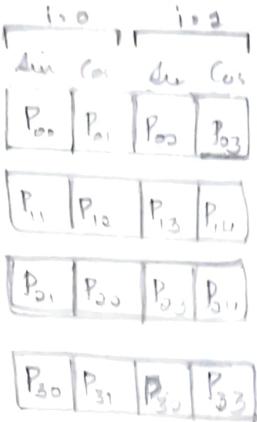
$$z_1 = [0.33, 0.6, 0.14, 0.54]$$

$$z_2 = [0.34, 0.6, 0.14, 0.55]$$

example:

→ 9 am a Robot

$P_e : d = 4 = \mathbb{R}^4$ * dimension of embedding space



$$P_{pos, 2i} = \sin\left(\frac{pos}{100^{(2i/4)}}\right) \quad \# \mathbb{R}^d = 4$$

$$P_{pos, 2i+1} = \cos\left(\frac{pos}{100^{(2i/4)}}\right)$$

$$\Rightarrow i=0$$

$$P_{0,0} = \sin\left(\frac{0}{100^{(0)}}\right) = \sin(0) = 0$$

$$P_{0,1} = \cos\left(\frac{0}{100^{(0)}}\right) = \cos(0) = 1$$

$$\Rightarrow i=1$$

$$P_{0,2} = \sin\left(\frac{0}{100^{(2)}}\right) = \sin(0) = 0$$

$$P_{0,3} = \cos\left(\frac{0}{100^{(2)}}\right) = \cos(0) = 1$$

$$P^{(0)} = [0|1|0|1]$$

$$\Rightarrow P^1$$

$$\Rightarrow i = 0$$

$$\textcircled{2} \quad P_{1,0} = \sin\left(\frac{1}{100^{(o)}}\right) = \sin(1) \Rightarrow 0.84$$

$$P_{1,1} = \cos\left(\frac{1}{100^{(o)}}\right) = \cos(1) = 0.54$$

$$\Rightarrow i = 1$$

$$P_{2,2} = \sin\left(\frac{1}{100^{0.5}}\right) = \sin\left(\frac{1}{10}\right) = \sin(0.1)$$

$$\Rightarrow 0.0998$$

$$\Rightarrow 0.10$$

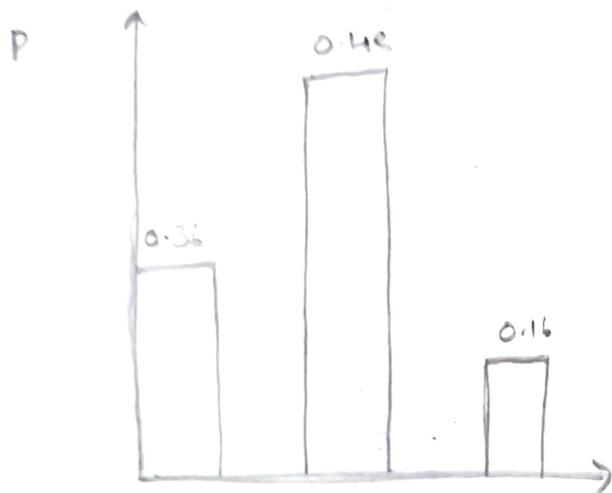
$$P_{2,3} = \cos\left(\frac{1}{10}\right) = \cos(0.1) = 0.99 \Rightarrow 1$$

$$P^1 = \boxed{0.84 \quad 0.54 \quad 0.1 \quad 1}$$

$$\theta^* = \underset{\theta \in M}{\operatorname{arg\! min}} \mathbb{E} \log (\cancel{P_\theta(x)}) - \log P_\theta(x) \quad \# \text{ we are minimizing on } \theta \text{ only}$$

$$\theta^* = \underset{\theta \in M}{\operatorname{arg\! min}} \mathbb{E} \log (\cancel{P_\theta(x)}) \quad \# \text{ Maximal likelihood estimation}$$

\Rightarrow example :-



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

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

$$= \frac{9}{25} \ln \left(\frac{9/25}{1/3} \right) + \frac{12}{25} \ln \left(\frac{12/25}{1/3} \right) + \frac{4}{25} \ln \left(\frac{4/25}{1/3} \right)$$

$$= 0.0852$$

$$\begin{aligned}
 (G||P) &= \sum_{x \in X} P(x=x) \ln \frac{g(x=x)}{P(x=x)} \\
 &= \frac{1}{3} \ln \left(\frac{Y_3}{9/25} \right) + \frac{1}{3} \ln \left(\frac{Y_3}{12/25} \right) + \frac{1}{3} \ln \left(\frac{Y_3}{4/25} \right) \\
 &= \frac{1}{3} \ln (0.926) + \frac{1}{3} \ln (0.69) + \frac{1}{3} \ln (2.08) \\
 &= \frac{1}{3} (-0.07) + \frac{1}{3} (-0.37) + \frac{1}{3} \cancel{\ln} (0.73) \\
 &= -0.02 - 0.12 + 0.244 \\
 &= 0.10
 \end{aligned}$$