

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

$$z_1^{(1)} = 1 \times 0 + 0 \times 1 + 0 = 0.$$

$$z_2^{(1)} = 1 \times 0 + 0 \times 1 + (-1) = -1$$

$$a_1(z_1^{(1)}) = \text{ReLU}(z_1^{(1)}) = 0$$

$$a_2(z_2^{(1)}) = \text{ReLU}(z_2^{(1)}) = 0.$$

$$z_1^{(2)} = w_{11}^{(2)} \times a_1^{(1)} + w_{12}^{(2)} \times a_2^{(1)} + b_1^{(2)}$$

$$z_1^{(2)} = 1 \times 0 + (-2) \times 0 + 0$$

$$a_1^{(2)} = \text{ReLU}(z_1^{(2)}) = 0.$$

$$y = 0.$$

x_1	x_2	y	\hat{y}
0	0	0	0
0	1	1	1
1	0	1	1
1	1	0	0



→ handcrafting the initialization on the basis of domain knowledge.

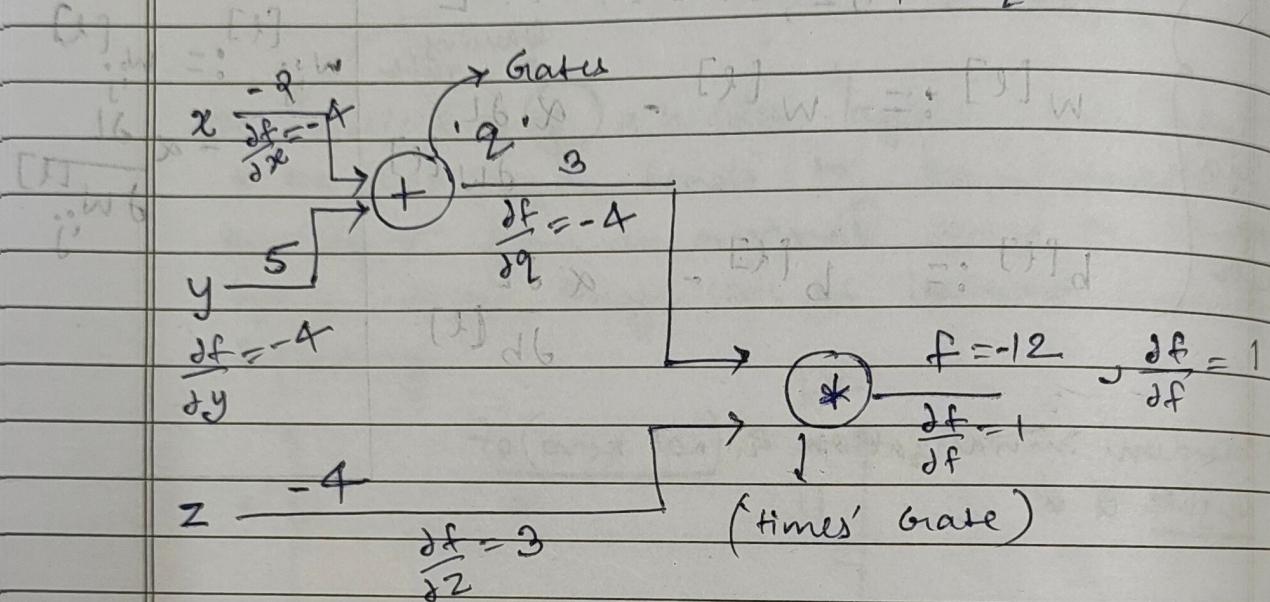
papergrid
Date: 11/11/2023

Computational Graph

* Computation

$$f = (x + y)z \quad ; \quad \text{let } q = x + y$$

$$f = qz$$



• How changes in the input leads to change in op

• $\frac{\partial f}{\partial x} \rightarrow$ how x is influencing f ?

chain rule.

$$\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$

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

$$\frac{\partial q}{\partial x} = \frac{\partial(x+y)}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$= \frac{\partial f}{\partial q} \cdot 1$$

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

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

$$= -4 \times 1$$

$$= -4$$

define additional functions $\rightarrow \{ q = yz; f = x + q \}$

papergrid

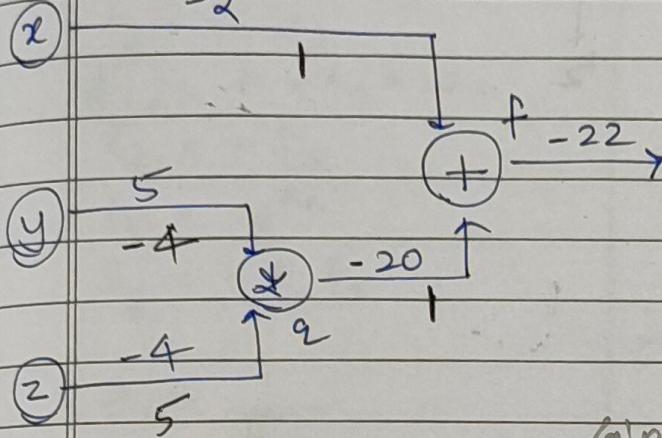
* begin computational graph

Date: / /

-2

($x, w + y, w + z$)

define forward pass.



$$q = 5x - 4 = 5 - 4 = 1$$

$$f = x + q = -2 + 1 = -3$$

$$s = -1$$

$$x = 5$$

$$y = 5$$

Global (gradient flowing from the final output).

Perform Backward pass

Global.

$$\frac{\partial f}{\partial q} = \frac{\partial(x+q)}{\partial q} = 1, \quad \frac{\partial f}{\partial x} = 1 \quad \text{# local}$$

$$\begin{cases} \frac{\partial q}{\partial y} = \frac{\partial(yz)}{\partial y} = z & \frac{\partial q}{\partial z} = \frac{\partial(yz)}{\partial z} = y \\ & = -4 & = 5 \end{cases} \quad \text{# local}$$

$$\text{Now, } \frac{\partial f}{\partial y} = \frac{\partial q}{\partial y} \times \frac{\partial f}{\partial q} = z \times 1 = -4$$

$$\frac{\partial f}{\partial z} = \frac{\partial q}{\partial z} \times \frac{\partial f}{\partial q} = y \times 1 = 5$$

∴ Overall influence of x on f is 1, y on f is -4 & z on f is 5.

Give only 2 inputs to each gate

papergrid

Date: / /

$$f_1 = w_0 x_0$$

$$f_2 = w_1 x_1$$

$$f_3 = f_1 + f_2 + w_2$$

$$f_4 = -f_3$$

$$f_5 = e^{f_4}$$

$$f_6 = 1 + f_5$$

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

$$f = f_7$$

$$f_1 = w_0 x_0$$

$$f_2 = w_1 x_1$$

$$f_3 = f_1 + f_2$$

$$f_4 = f_3 + w_2$$

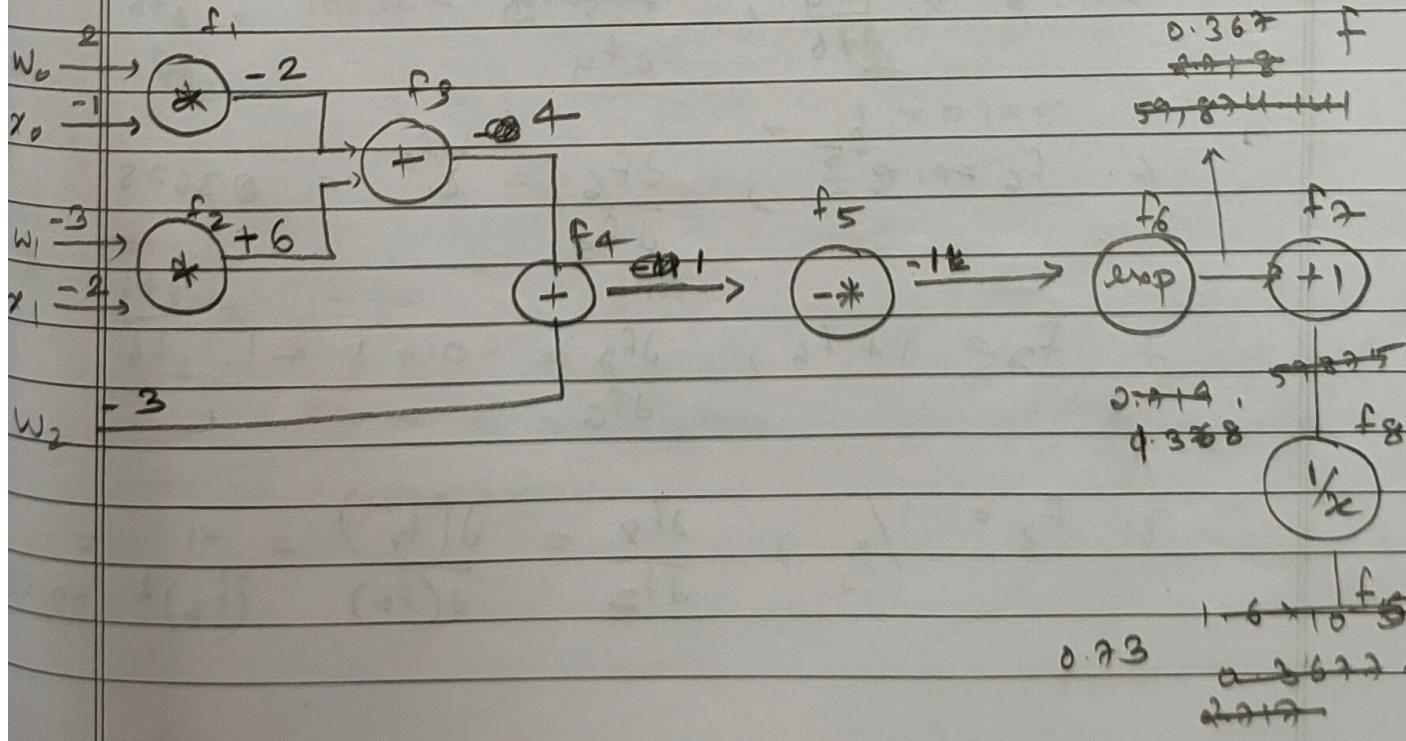
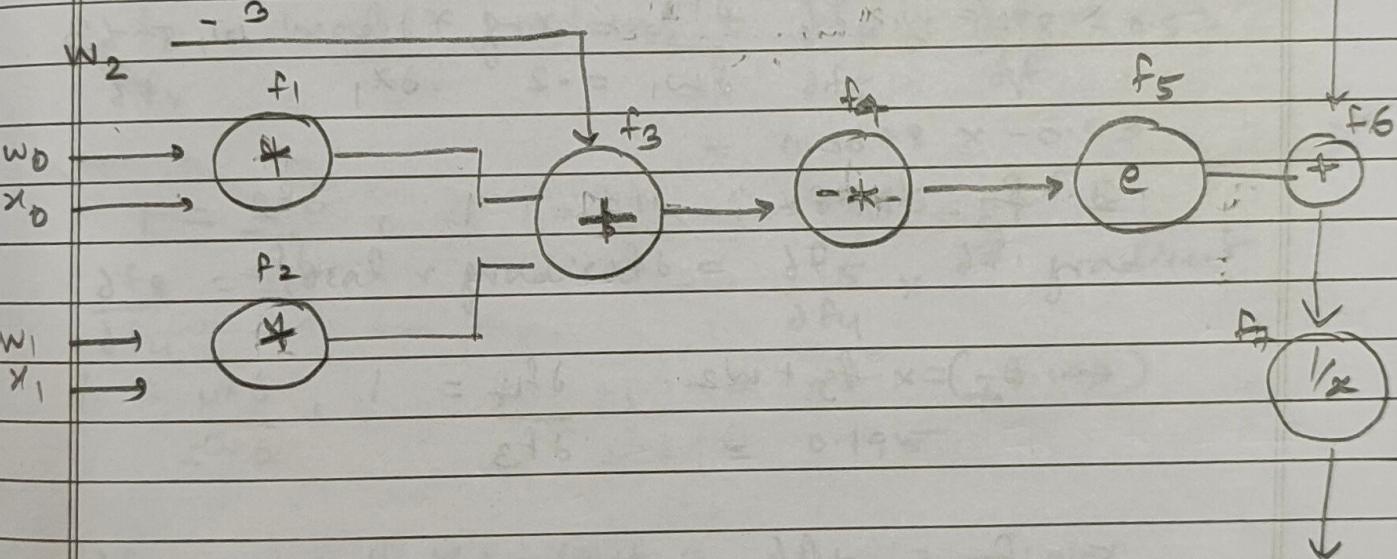
$$f_5 = -f_4$$

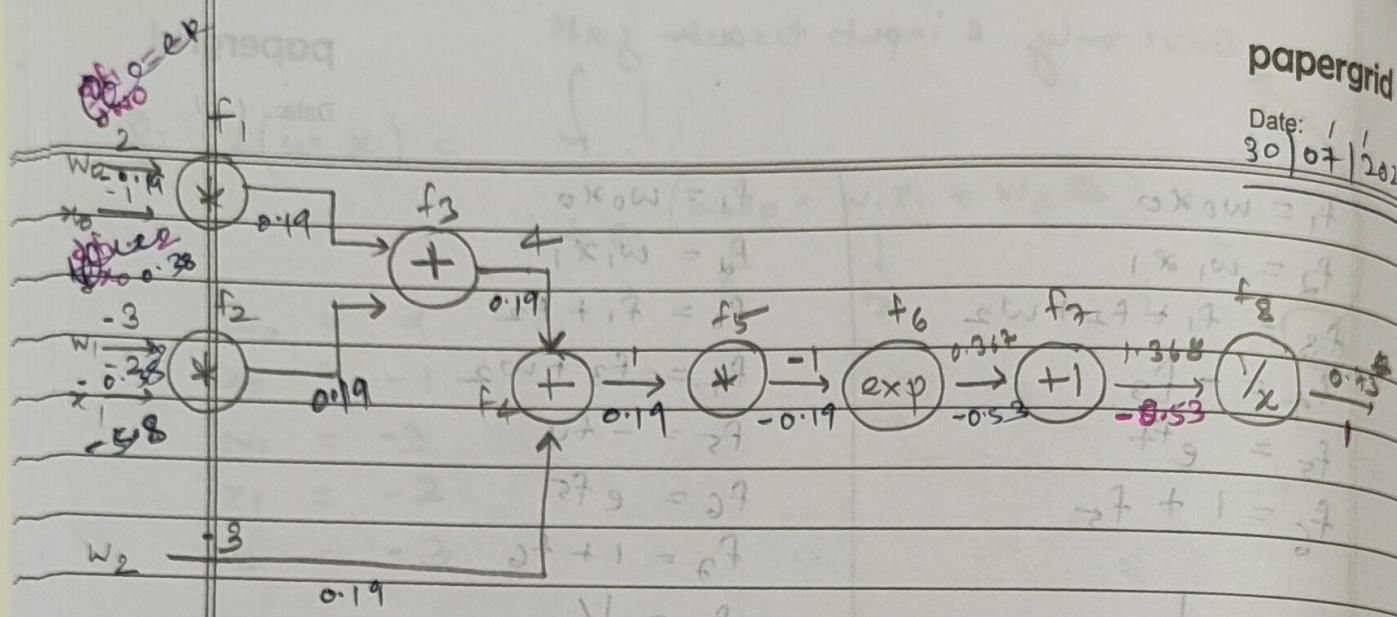
$$f_6 = e^{f_5}$$

$$f_7 = 1 + f_6$$

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

$$f = f_8$$





$$1. \quad f_1 = w_0 x_0 ; \frac{\partial f_1}{\partial w_0} = x_0, \frac{\partial f_1}{\partial x_0} = w_0 = 2$$

$$2. \quad f_2 = w_1 x_1 ; \frac{\partial f_2}{\partial w_1} = x_1, \frac{\partial f_2}{\partial x_1} = w_1 = -3$$

$$3. \quad f_3 = f_1 + f_2 ; \frac{\partial f_3}{\partial f_1} = 1, \frac{\partial f_3}{\partial f_2} = 1$$

$$4. \quad f_4 = -f_3 + w_2 ; \frac{\partial f_4}{\partial f_3} = 1, \frac{\partial f_4}{\partial w_2} = 1$$

$$5. \quad f_5 = -f_4, \frac{\partial f_5}{\partial f_4} = -1$$

$$6. \quad f_6 = e^{f_5}, \frac{\partial f_6}{\partial f_5} = e^{f_5} = 0.3688$$

$$7. \quad f_7 = 1 + f_6, \frac{\partial f_7}{\partial f_6} = 0 + 1 = 1$$

$$8. \quad f_8 = \frac{1}{f_7}, \frac{\partial f_8}{\partial f_7} = \frac{\partial(\frac{1}{f_7})}{\partial f_7} = \frac{-1}{(f_7)^2} = -0.33$$

* Bulk prop.

$$\frac{\partial f_8}{\partial f_2} = \text{local } \times \text{gradient}$$

$$\frac{\partial f_8}{\partial f_2} = \frac{1}{\frac{\partial f_8}{\partial f_6}} \times 1 = -0.53$$

$$\begin{aligned} \frac{\partial f_8}{\partial f_6} &= \text{local } \times \text{gradient} = \frac{\partial f_2}{\partial f_6} \times \frac{\partial f_8}{\partial f_2} \\ &= 1 \times (-0.53) \\ &= -0.53 \end{aligned}$$

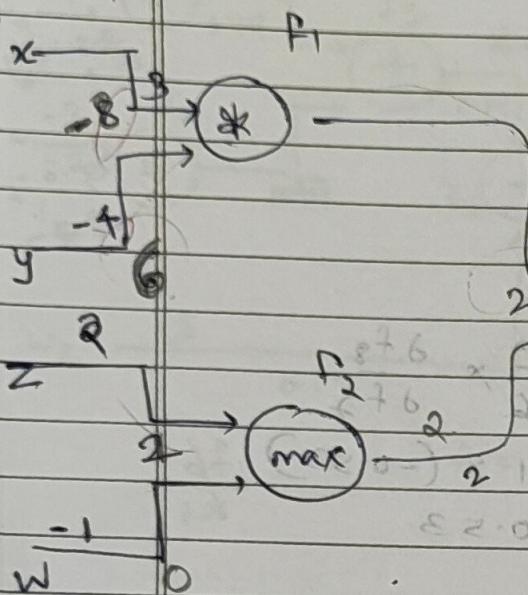
$$\begin{aligned} \frac{\partial f_8}{\partial f_5} &= \text{local } \times \text{gradient} = \frac{\partial f_6}{\partial f_5} \times \frac{\partial f_8}{\partial f_6} = -0.53 \\ &= 0.3678 \times -0.53 \end{aligned}$$

$$\begin{aligned} \frac{\partial f_8}{\partial f_4} &= \text{local } \times \text{gradient} = \frac{\partial f_5}{\partial f_4} \times \text{gradient} \\ &= -1 \times (-0.195) \end{aligned}$$

$$\begin{aligned} \frac{\partial f_8}{\partial f_2} &= \text{local } \times \text{gradient} = \frac{\partial f_4}{\partial f_2} \times 0.195 \\ &= 1 \times 0.195 \end{aligned}$$

$$\begin{aligned} \frac{\partial f_8}{\partial f_2} &= \text{local } \times \text{gradient} \\ &\leq w \text{ if } l = \overline{w} \\ &\geq w \text{ if } l = \overline{w} \\ &= 0.195 \text{ if } l = \overline{w} \end{aligned}$$

#



$$f_2 = \max(z, w) \quad f_3 = \begin{cases} 2 & \text{if } z > w \\ -10 & \text{if } z \leq w \end{cases}$$

$$f_1 = x + 8y \quad \frac{\partial f_1}{\partial x} = 1, \quad \frac{\partial f_1}{\partial y} = 8$$

$$f_2 = \max(z, w) \quad \frac{\partial f_2}{\partial z} = 1, \quad \frac{\partial f_2}{\partial w} = 0$$

$$f_3 = f_1 + f_2 \quad \frac{\partial f_3}{\partial z} = 1, \quad \frac{\partial f_3}{\partial w} = 0$$

$$f_4 = 2 \times f_3 \quad \frac{\partial f_4}{\partial z} = 1 \text{ if } z > w$$

$$\frac{\partial f_4}{\partial z} = 0 \text{ if } z \leq w$$

undefined if $z = w$

read:

$$\frac{\partial f_2}{\partial z} = \frac{\partial (\max(z, w))}{\partial z}; \quad z > w \quad \frac{\partial f_3}{\partial z} = 1$$

then due to
partial derivative
 $w \rightarrow \text{const}$

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

$$= \frac{\partial (z)}{\partial z} = 1$$

$$\frac{\partial f_4}{\partial z} = 2$$

$$\frac{\partial f_2}{\partial z} = \frac{\partial (\max(z, w))}{\partial z}; \quad z < w \quad \frac{\partial f_3}{\partial z} = 0$$

$$\frac{\partial f_4}{\partial z} \neq \frac{\partial f_3}{\partial z} \times$$

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

undefined if $w = z$

$$\frac{\delta f_4}{\delta f_3} = 2, \quad \frac{\delta f_4}{\delta f_2} = \frac{\delta f_3}{\delta f_2} \times \frac{\delta f_4}{\delta f_3} = 1 \times 2 = 2$$

$$6 \xrightarrow{16} 16$$

$$\frac{\delta f_4}{\delta f_1} = \frac{\delta f_3}{\delta f_1} \times 2 = 2$$

$$\frac{\delta f_4}{\delta f_2} = 2$$

$$(1-2) \rightarrow 1 = 32M$$

$$\frac{\delta f_4}{\delta 2} = \frac{\delta f_2}{\delta 2} \times 2 = 2 \times 1 \text{ (FR)}$$

~~remember~~
~~Note - 2~~

31/07/2025.

1. \oplus gate acts as a gradient distributor., the gradient gets distributed equally to all the inputs. ✓

2. \max gate , in this gate the gradient flows towards the path where the input values are higher . ✓

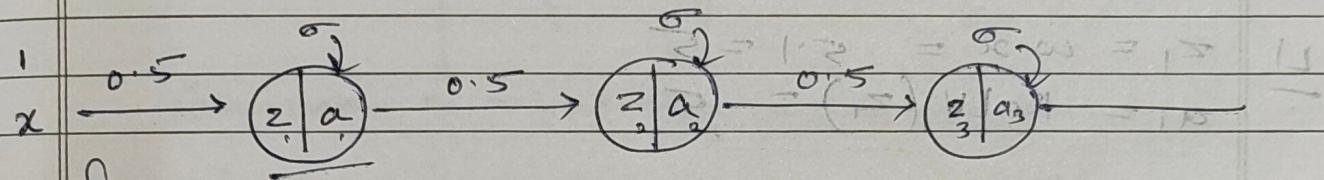
In previous eg: , f_2 will have no impact or change of own input . ✓

→ routes the gradient to the hi

3. \star gate \rightarrow swaps the input activations \mathbf{x}_1 (times or multiply) multiply by global gradient . (applicable with two inputs only) *

- as it prevents saturation of the input grid
- ReLU, used to fix the vanishing gradient problem.
 - Opposite of vanishing gradients → exploding gradients caused by ReLU.
 - To fix exploding gradients, caused by ReLU, the 'clipping' technique is used.
 - fixing gradients

Example



~~forward pass~~

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

$$a_1 = \sigma(z_1) = \frac{1}{1 + e^{-0.5}} = \frac{1}{1 + 0.606} = \underline{\underline{0.622}}$$

$$z_2 = 0.5 \times 0.622 = 0.311$$

$$a_2 = \sigma(z_2) = \frac{1}{1 + e^{-0.311}} = \frac{1}{1 + 0.732} = \underline{\underline{0.577}}$$

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

$$a_3 = \frac{1}{1 + e^{-0.288}} = \frac{1}{1.749} = \underline{\underline{0.572}}$$

$$L = \frac{1}{2} (a_3 - y)^2 \quad ; \quad \frac{\partial L}{\partial a_3} = a_3 - y = \underline{\underline{0.572 - 1}} = -0.428$$

$$\frac{\partial a_1}{\partial x}$$

During the test time, the entire architecture is kept as it is, the entire data is shown to the neurons. To account for the keep probability at the training time, a scaling factor is applied at each neuron during testing.

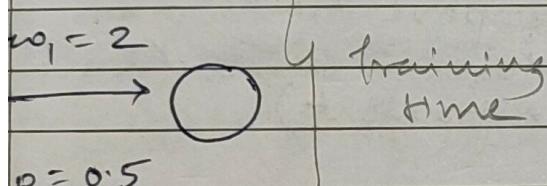
* scaling factor is keep prob. multiplied to the activation

→ Standard Dropout. → scaling factor is keep prob. * to the activation / weights.

noscaling at test-time.

to maintain consistent activation magnitudes

* Truncated Dropout.



training
time

testing time.

no dropout.

Sample 1 → o/p → 2

Sample 2 → o/p → 2

Sample 3 → o/p → 2

Sample 4 → o/p → 2

Total o/p = 8.

$$\text{Total output} = 2+2=4$$

multiply by $\frac{1}{p}$

$$\therefore 4 \times \frac{1}{0.5} = 8$$

scale the activation during training by multiplying with $(\frac{1}{0.5}) = \frac{1}{p}$

high pixel \rightarrow feature
paper

Date: / /

Convolution operation.

values represent how
the feature is detected
in the original image.

$$\begin{array}{|c|c|c|c|c|c|} \hline
 3 & 0 & 1 & 2 & 7 & 4 \\ \hline
 1 & 5 & 8 & 9 & 3 & 1 \\ \hline
 2 & 7 & 2 & 5 & 1 & 3 \\ \hline
 0 & 1 & 3 & 1 & 7 & 8 \\ \hline
 4 & 2 & 1 & 6 & 2 & 8 \\ \hline
 2 & 4 & 5 & 2 & 3 & 9 \\ \hline
 \end{array}
 \begin{array}{|c|c|c|} \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 \end{array}
 =
 \begin{array}{|c|c|c|} \hline
 -5 & -4 & 0 & 8 \\ \hline
 -10 & -2 & 2 & 3 \\ \hline
 0 & -2 & -4 & -7 \\ \hline
 -3 & -2 & -3 & 6 \\ \hline
 \end{array}$$

Input 6×6 image

convolution
operator

feature map
= Activation map

filter or kernel (3×3) \rightarrow learn this parameter?

Element wise multiplication
of summation

\rightarrow Take the filter & superimpose on the input.

$$\begin{aligned}
 & 3 \times 1 + 0 \times 0 + (-1) \times (-1) + 1 \times 1 + 5 \times 0 + 8 \times -1 \\
 & + 2 \times 1 + 7 \times 0 + 2 \times -1
 \end{aligned}$$

$$= -5$$

\rightarrow slide the filter by one

$$\begin{aligned}
 & 0 \times 1 + (1 \times 0) + 2 \times -1 + 5 \times 1 + 8 \times 0 + 9 \times -1 + 7 \times 1 + 2 \times 0 \\
 & + 2 \times -1 \\
 & = -2 - 4 + 5 \\
 & = -2 + 1
 \end{aligned}$$

$$= -4$$

3x3 - 14/08/2025 - 9

papergrid

Example.

Date: 1-11

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

$$\begin{array}{|c|c|} \hline 1 & -1 \\ \hline 0 & -1 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline -9 & -11 & -7 \\ \hline \end{array}$$

$$\begin{array}{|c|c|c|} \hline -11 & 4 & -3 \\ \hline \end{array}$$

3x3

$$n = 5 \times 5$$

$$7 \times 7$$

$$p = 1$$

$$5 + 2 \times 1 - 2 + 1$$

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

$$s = 2$$

$$p + 1$$

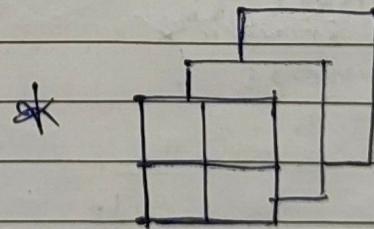
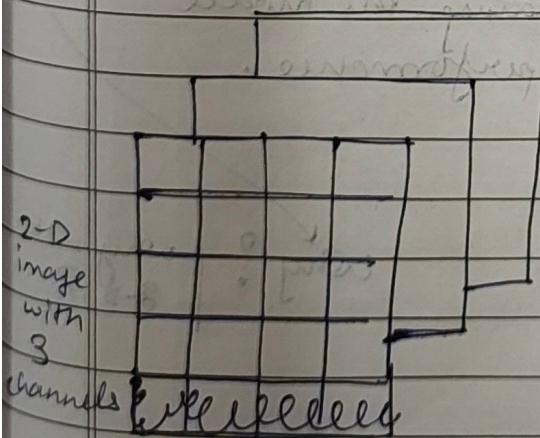
$$\left[\frac{n + 2p - f + 1}{s} \right]$$

3x3

* General formula :- $\left[\frac{n + 2p - f + 1}{s} \right] \frac{f + 4 - 2}{2} = 1$

$$\left[\frac{n + 2p - f + 1}{s} \right] \frac{f + 4 - 2}{2} = 1$$

Convolutions on 3D Volume. → eg: RGB,

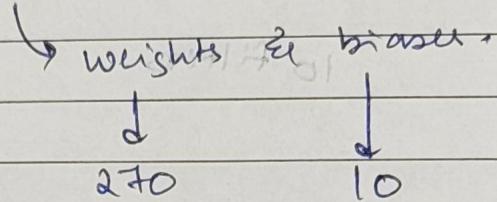


4 x 4 x (3)
height
width
no. of channels

- Filters for each channel
stacked behind each other.

a) 10 filters } in one layer of Conv: Net
 $3 \times 3 \times 3$

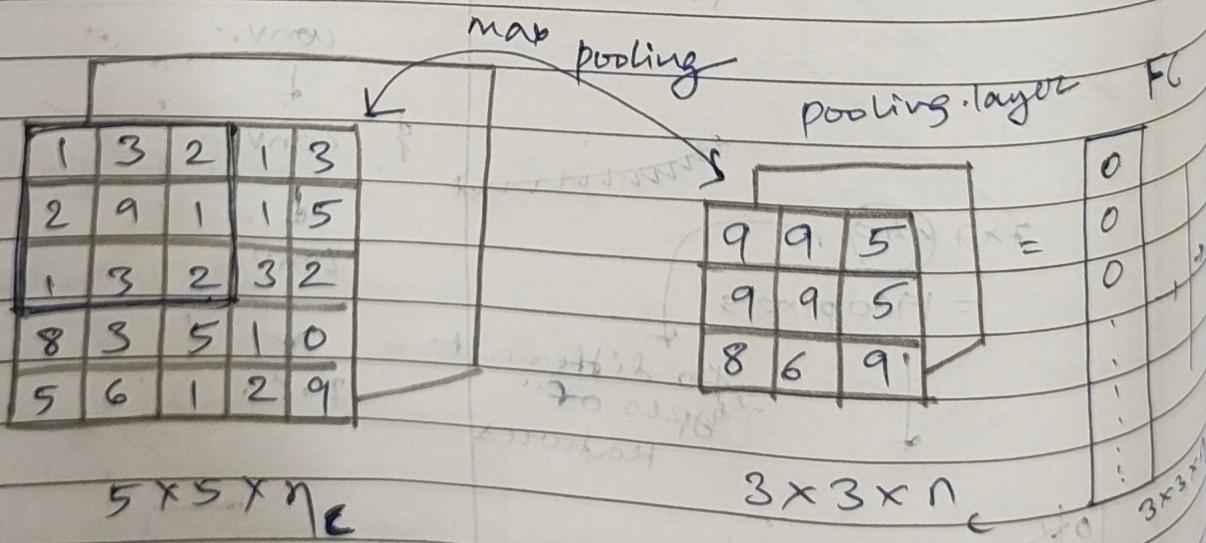
How many parameters does this layers have?



1 filter \rightarrow 27 weights + 1 bias

10 filter \rightarrow $270 \cdot 10 + 10 \cdot 1 = 280$ ✓

- * 1. CONV layer
- 2. Pooling layer \rightarrow Max Pooling
- 3. Fully Connected layer (FC) Avg. Pooling



conv layer \rightarrow pooling layer \rightarrow conv layer

Character level language model.

papergrid

Date: / /

$T = 2$ (two time steps) \rightarrow "Hello World"

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

$$\text{input dim} = 2 \Rightarrow x_t \in \mathbb{R}^2$$

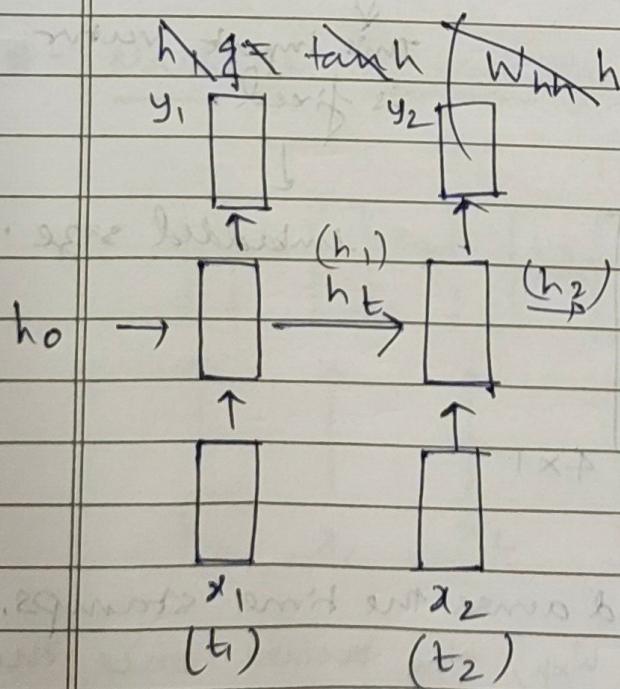
$$\text{hidden state dim} = 3 \Rightarrow h_t \in \mathbb{R}^3$$

$$h_0 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, W_{xh} = \begin{bmatrix} 0.5 & -0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.4 \end{bmatrix}, W_{hh} = \begin{bmatrix} 0.1 & 0.4 & 0.5 \\ 0.2 & 0.3 & 0.6 \\ 0.05 & 0.1 & 0.1 \end{bmatrix}$$

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

$$W_{hy} = \begin{bmatrix} 1.0 & -1.0 & 0.5 \\ 0.5 & 0.5 & -0.5 \end{bmatrix}$$

2×3



$$h_1 = (W_{hh} h_0 + W_{xh} x_1)$$

$$= \begin{bmatrix} 0.1 & 0.4 & 0.0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.2 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\textcircled{1} \quad h_1 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} -0.1 \\ 1.2 \\ 0.9 \end{bmatrix} = \begin{bmatrix} -0.1 \\ 1.2 \\ 0.9 \end{bmatrix}$$

$$y_1 = w_{hy} \times h_1 = \begin{bmatrix} 1.0 & -1.0 & 0.5 \\ 0.5 & 0.5 & -0.5 \end{bmatrix} \begin{bmatrix} -0.1 \\ 1.2 \\ 0.9 \end{bmatrix}$$

2×3 3×1

$$\begin{bmatrix} (-28.0) + (8.0) \\ 28.0 - 4.0 \\ 28.0 + 8.0 \end{bmatrix} = \begin{bmatrix} -0.85 \\ 1.0 \\ 0.1 \end{bmatrix}$$

2×1

$$h_2 = \begin{bmatrix} 0.1 & 0.4 & 0.0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.2 \end{bmatrix} \begin{bmatrix} -0.1 \\ 1.2 \\ 0.9 \end{bmatrix} + \begin{bmatrix} 0.5 & -0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.4 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

$$= -0.01 + 0.48 + 0 + 0.02$$

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

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

$$g(-0.1) = \frac{1.105 - 0.904}{1.105 + 0.904}$$

$$h_1 = \begin{bmatrix} -0.1 \\ 0.8 \\ 0.7 \end{bmatrix} + \frac{1.105 + 0.904}{2} = \frac{0.209}{2} = 0.1045$$

$$g(1.2) = \frac{3.32 - 0.301}{3.32 + 0.301}$$

$$g(0.9) = \frac{2.5 - 0.406}{2.5 + 0.406} = \frac{3.019}{3.621} = 0.834$$

$$= 0.72$$

$$y_1 = W_{hy} \times h_1 = \begin{bmatrix} 1.0 & -1.0 & 0.5 \\ 0.5 & 0.5 & -0.5 \end{bmatrix} \begin{bmatrix} -0.1 \\ -0.8 \\ 0.35 \end{bmatrix}$$

$$= \cancel{\left[-0.1 + (-0.8) + (0.35) \right]} + \cancel{0.05 + 0.4 - 0.35}$$

$$= \begin{bmatrix} -0.35 \\ 0.1 \end{bmatrix} \begin{bmatrix} -0.55 \\ 0.008 \end{bmatrix}$$

$$= \begin{bmatrix} -0.18 - 0.8 + 0.35 \\ -0.05 + 0.40 - 0.35 \end{bmatrix} = \begin{bmatrix} -0.55 \\ 0.008 \end{bmatrix}$$

$$h_2 = \begin{bmatrix} 0.1 & 0.4 & 0.0 \\ -0.25 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.2 \end{bmatrix} \begin{bmatrix} -0.099 \\ 0.83 \\ 0.716 \end{bmatrix} + \begin{bmatrix} 0.5 & -0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.4 \end{bmatrix} \begin{bmatrix} -0.5 \\ 0.008 \end{bmatrix}$$

$$= \begin{bmatrix} -0.01 + 0.32 + 0 \\ -0.02 + 0.24 + 0.14 \\ -0.005 - 0.08 + 0.14 \end{bmatrix} + \begin{bmatrix} -0.5 - 0.3 \\ -0.8 + 0.2 \\ -0.1 + 0.4 \end{bmatrix}$$

$$= \begin{bmatrix} 0.31 \\ 0.36 \\ 0.056 \end{bmatrix} + \begin{bmatrix} -0.8 \\ -0.6 \\ 0.3 \end{bmatrix} = \begin{bmatrix} -0.5 \\ -0.24 \\ 0.356 \end{bmatrix}$$

$$- 0.0099 + 0.332 + 0$$

$$0.3221$$

$$0.00198 + 0.249 + 0.1432$$

$$0.3941$$

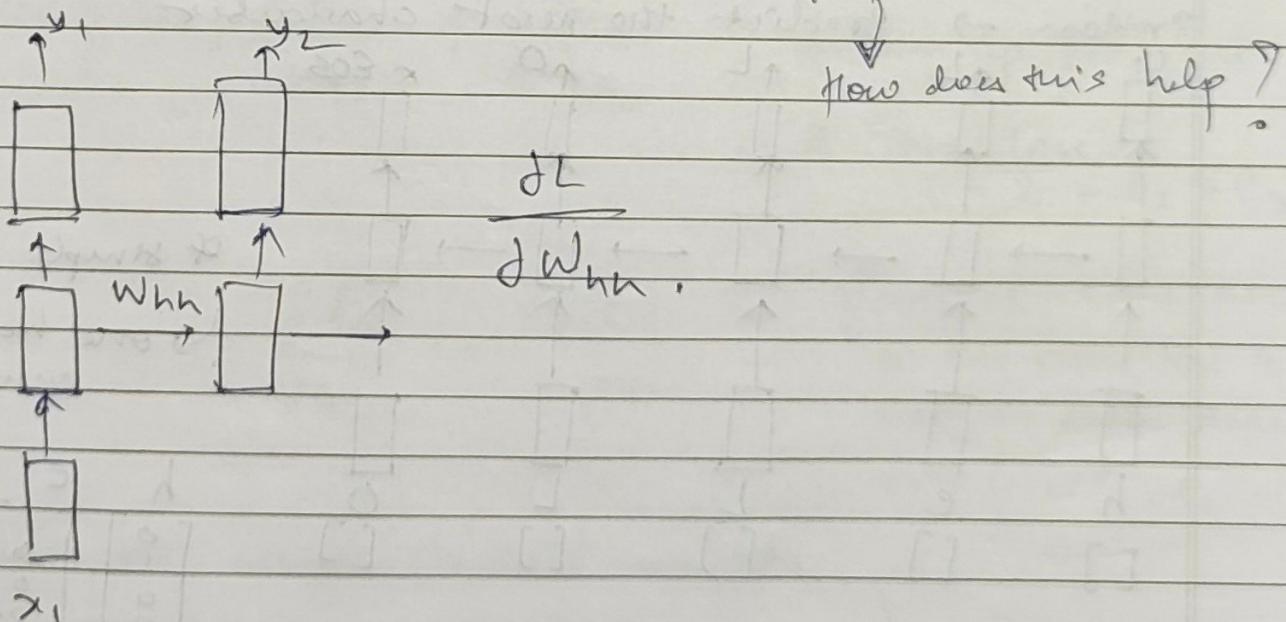
$$- 0.0049 - 0.083 + 0.1432$$

$$0.0553$$

$$h_2 = \begin{bmatrix} -0.4778 \\ -0.2059 \\ 0.3553 \end{bmatrix}$$

$$y_2 = \begin{bmatrix} -0.0886 \\ -0.4844 \end{bmatrix}$$

Weights and are the same, therefore parameters are shared across the time stamps in RNN.



Example -

$$\begin{aligned}x_t &= [0.5, -0.1]^T \\h_{t-1} &= [0.0, 0.1]^T\end{aligned}$$

$$c_t = [0.2, -0.2]^T$$

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

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

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

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

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

$$W_{h0} = \begin{bmatrix} 0.15 & 0.05 \\ 0.1 & -0.2 \end{bmatrix}$$

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

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

$$\left\{ h_t = ?, c_t = ? \right\}$$

$$h_t = o_t \odot \tanh(c_t) \rightarrow \text{element wise multiplication}$$

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

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

$i_t = s_t \rightarrow \text{if } f_t \rightarrow o_t \rightarrow h_t$

0.002

0.02

0.01

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

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

$$\sigma \left(\begin{bmatrix} -0.02 \\ 0.005 \end{bmatrix} + \begin{bmatrix} 0.28 \\ 0.49 \end{bmatrix} \right) = \sigma \begin{pmatrix} 0.30 \\ 0.95 \end{pmatrix}$$

$$i_t = \begin{pmatrix} 0.57 \\ 0.55 \end{pmatrix}$$

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

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

$$\tanh \left(\begin{bmatrix} 0.01 \\ 0.005 \end{bmatrix} + \begin{bmatrix} -0.29 \\ 0.43 \end{bmatrix} \right) = \begin{pmatrix} -0.28 \\ 0.43 \end{pmatrix}$$

$$= \begin{pmatrix} -0.272 \\ 0.13 \end{pmatrix}$$

$$g_t = f_t \odot c_t + i_t \odot g_t$$

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

$$= \sigma \left(\begin{bmatrix} 0.05 & -0.1 \\ -0.2 & 0.1 \end{bmatrix} \begin{bmatrix} 0.6 \\ 0.1 \end{bmatrix} + \begin{bmatrix} -0.4 & 0.2 \\ 0.3 & 0.3 \end{bmatrix} \begin{bmatrix} 0.8 \\ -0.1 \end{bmatrix} \right)$$

$$= \sigma \left(\begin{bmatrix} -0.01 \\ 0.01 \end{bmatrix} + \begin{bmatrix} -0.22 \\ 0.12 \end{bmatrix} \right) = \sigma \begin{bmatrix} -0.23 \\ 0.73 \end{bmatrix}$$

$$= \begin{bmatrix} 0.44 \\ 0.53 \end{bmatrix}$$

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

$$= \tanh \left(\begin{bmatrix} 0.2 & 0.1 \\ -0.1 & 0.05 \end{bmatrix} \begin{bmatrix} 0.6 \\ 0.1 \end{bmatrix} + \begin{bmatrix} -0.5 & 0.4 \\ 0.2 & -0.3 \end{bmatrix} \begin{bmatrix} 0.8 \\ -0.1 \end{bmatrix} \right)$$

$$= \begin{bmatrix} -0.23 \\ 0.13 \end{bmatrix}$$

$$o_t = \begin{bmatrix} 0.53 \\ 0.46 \end{bmatrix}$$

output of gate is sigmoid \rightarrow therefore b/w 0-1
Date: 1/1/2023
mostpapergrid

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

$$= \begin{bmatrix} 0.44 \\ 0.53 \end{bmatrix} \odot \begin{bmatrix} -0.2 \\ -0.2 \end{bmatrix} + \begin{bmatrix} 0.57 \\ 0.55 \end{bmatrix} \odot \begin{bmatrix} -0.22 \\ 0.13 \end{bmatrix}$$

$$c_t = \begin{bmatrix} 0.088 \\ -0.106 \end{bmatrix} + \begin{bmatrix} -0.1539 \\ 0.0215 \end{bmatrix} = \begin{bmatrix} -0.0659 \\ -0.0345 \end{bmatrix}$$

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

$$= \begin{bmatrix} 0.53 \\ 0.46 \end{bmatrix} \odot \tanh \begin{bmatrix} -0.0659 \\ -0.0345 \end{bmatrix}$$

$$= \begin{bmatrix} 0.53 \\ 0.46 \end{bmatrix} \odot \begin{bmatrix} -0.0658 \\ -0.0344 \end{bmatrix}$$

$$h_t = \begin{bmatrix} -0.0348 \\ -0.0158 \end{bmatrix} = \begin{bmatrix} -0.03 \\ -0.01 \end{bmatrix}$$

Interview tips.

- 1) Understand the Dataset properly
- 2) seurat Bioinfo tools. papergrid
 - Generative models (Data Simulation)

Date: / /

30/09/2025

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

$$h \in \mathbb{R}^1; \text{ what should be size of } w_e$$

context
vector

$$\therefore w_d$$

- denoising
- Autoencoders
- Attention

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

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

$$\text{MSE Loss} = \frac{1}{2} \left\| \hat{x} - x \right\|_2^2$$

~~Autoencoder - Encoder~~ $\cancel{\text{Encoder}}$ $h = x w_e = \begin{bmatrix} 2 \\ 0 \end{bmatrix} \begin{bmatrix} 0.5 & -1.0 \end{bmatrix}_{1 \times 2}$

$$h = 1$$

$$x = \cancel{w_d h} = h \cdot w_d = 1 \times \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$$

$$\Delta = \sqrt{(-1)^2 + 0.25} = \sqrt{0.25} = \frac{\sqrt{3}}{2} = 0.8$$

$$L = \sqrt{(1-2)^2 + (0.5-0)^2} = \sqrt{1.25}$$

$$L = \frac{1}{2} \left[(1-2)^2 + (0.5-0)^2 \right] = \underline{\underline{0.625}}$$

11.6.17

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

grid

25

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

swing function = 0.0

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

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

$$c_t = ?$$

$$s_{t+1} = [1, 0, 1] \times 6.0 +$$

$$d_{t,j} = \frac{\exp(\text{sure}(s_{t-1}, h_j))}{\sum_{j=1}^3 \exp(\text{sure}(s_{t-1}, h_j))}$$

$$\sum_{j=1}^3 \exp(\text{sure}(s_{t-1}, h_j))$$

$$c_t = d_{t,1}h_1 + d_{t,2}h_2 + d_{t,3}h_3$$

$$d_{t,1} = \frac{\exp([1, 0, 1] \cdot [1, 0, 1])}{\exp([1, 0, 1] \cdot [1, 0, 1]) + \exp([1, 0, 1] \cdot [0, 1, 1]) + \exp([1, 0, 1] \cdot [1, 1, 0])}$$

$$\exp([1, 0, 1] \cdot [1, 0, 1]) + \exp([1, 0, 1] \cdot [0, 1, 1])$$

$$+ \exp([1, 0, 1] \cdot [1, 1, 0])$$

$$= \frac{\exp(2)}{\exp(2) + \exp(1) + \exp(1)}$$

=

~~5.43~~

~~5.43 + 2.718 + 2.718~~

$$= \frac{\exp(2)}{\exp(2) + \exp(1) + \exp(1)}$$

$$d_{t,1} = \frac{5.43}{10.86} = 0.5$$

$$= 2.389$$

$$2.389 + 2.71 + 2.71$$

$$= 0.576$$

$$\alpha_{t,2} = 0.2$$

$$\alpha_{t,3} = 0.2$$

$$c_t = 0.526 \times [1, 0, 1] + 0.2 \times [0, 1, 1] \\ + 0.2 \times [1, 1, 0]$$

$$= [0.526, 0, 0.526] + [0, 0.2, 0.2] \\ + [0.2, 0.2, 0]$$

$$c_t = [0.776, 0.4, 0.776]^T$$

* c_t closer to h_1 → gets aligned to one of the hidden states.

(focus on this part of input to generate next word (total)).

$$([1, 0, 0] \cdot [1, 0, 1]) \text{ you} + ([1, 0, 1] \cdot [1, 0, 1]) \text{ you}$$

$$([0, 1, 1] \cdot [1, 0, 1]) \text{ you} +$$

$$216.6 + 216.6 + 216.6$$

$$(1) \text{ you} + (1) \text{ you} + (1) \text{ you}$$

$$216.6 + 216.6 + 216.6$$

$$+ 216.6 =$$

$$182.6 =$$

$$16.6 + 16.6 + 182.6$$

Thinking Machines

Thinking

Machines

$$\textcircled{1} \quad q_1 = \text{Thinking} \quad \checkmark$$

$$k_1 = \text{Thinking} \quad \checkmark$$

$$v_1 = \text{Thinking}$$

$$\textcircled{1} \quad q_1 = \text{Machines}$$

$$k_1 = \text{Machines}$$

Thinking

$$\textcircled{2} \quad q_1 = \text{Thinking}$$

$$k_2 = \text{Machines}$$

$$\textcircled{2} \quad q_1 = \text{Machines}$$

$$k_2 = \text{Machines}$$

I/P

Playing outside

O/P

$z_1 \downarrow$ & z_2

Playing,

$$q = [0.212 \quad 0.04 \quad 0.63 \quad 0.36]^T$$

$$k_1 = [0.31 \quad 0.84 \quad 0.903 \quad 0.57]^T$$

$$v_1 = [0.36 \quad 0.83 \quad 0.1 \quad 0.35]^T$$

outside

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

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

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

playing .

papergrid

Date: / /

Score

$$\text{score} \quad q_1 \times k_1$$

$$q_1 \times k_2$$



$$\begin{bmatrix} 0.212 & 0.04 & 0.63 & 0.36 \end{bmatrix}^T$$

$$0.91121$$

$$0.8017$$

score

divide
by
Jd_k

$$\frac{0.91121}{\sqrt{4}} =$$

$$= 0.456$$

$$\frac{0.8012}{\sqrt{4}} =$$

$$= 0.401$$

softmax

$$\frac{e^{-0.456}}{e^{-0.456} + e^{-0.401}}$$

$$\frac{e^{-0.401}}{e^{-0.456} + e^{-0.401}}$$

$$= \frac{0.633}{0.633 + 0.669}$$

$$= \frac{0.669}{0.633 + 0.669}$$

$$= \frac{0.633}{1.302}$$

$$= \frac{0.669}{1.302}$$

$$= 0.486$$

$$= 0.5138$$

$$z_1 = \text{softmax} \times v_1 + \text{softmax} \times v_2$$

$$\Rightarrow 0.486 \times [0.36 \ 0.83 \ 0.1 \ 0.35]^T$$

$$+ 0.5138 \times [0.31 \ 0.36 \ 0.19 \ 0.22]^T$$

$$= [0.13096 \ 0.40338 \ 0.0486 \ 0.1701]$$

$$[0.1592 \ 0.1849 \ 0.097 \ 0.3699]$$

Example

I am a robot.

$$P_E: d = 4 \rightarrow 1R$$

I.
am
a
robot

0
1
2
3

P_{00}
P_{01}
P_{02}
P_{03}

$i=0$	$i=1$		
P_{00}	P_{01}	P_{02}	P_{03}

P_{10}	P_{11}	P_{12}	P_{13}
----------	----------	----------	----------

P_{30}	P_{31}	P_{32}	P_{33}
----------	----------	----------	----------

$$P_E_{pos, 2i} = \sin\left(\frac{pos}{100(2i/d_{emb})}\right)$$

For $i=0$, (position)

$$P_E_{0,0} = \sin\left(\frac{0}{100(200/4)}\right) = \sin(0^\circ) = 0.$$

dimension index.

$$P_E_{0,1} = \cos(0) = \cos(0) = 1$$

$$P_E_{0,2} = \sin\left(\frac{0}{100(200/4)}\right) = \sin(0) = 0$$

$$P_E_{0,3} = \cos\left(\frac{0}{100(200/4)}\right) = \cos(0) = 1$$

For $i = 1$

$$P_{10} = \sin\left(\frac{1}{100(2/\sqrt{2})}\right) = \sin\left(\frac{1}{100\sqrt{2}}\right)$$

$$= \sin(0.1)$$

$$= 0.0017.$$

$$P_{11} = \cos\left(\frac{1}{100(2/\sqrt{2})}\right) = \cos(0.1) = 0.99$$

$$P_{12} = \sin\left(\frac{1}{100(2/\sqrt{4})}\right) = \sin(0.1) = 0.0017.$$

$$P_{13} = \cos(0.1) = 0.99.$$

For $i = 3$.

$$P_{30} = \sin\left(\frac{3}{100(2/\sqrt{4})}\right) = \sin\left(\frac{3}{100\sqrt{2}}\right)$$

$$= \sin\left(\frac{3}{2\sqrt{100}}\right)$$

$$= \sin\left(\frac{3}{1000}\right)$$

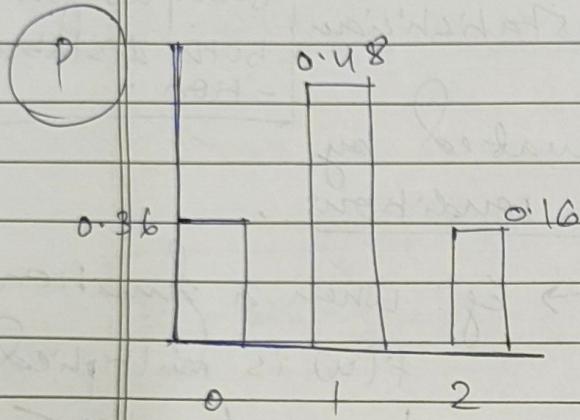
$$= \sin(0.003)$$

$$= 5.23 \times 10^{-5}$$

$$\text{argmin}_{\theta \in M} E[-\log(p_\theta(x))]$$

{ Maximum likelihood estimation. }

Example -



Binomial distribution

distribution X	0	1	2
p(x)	9/25	12/25	4/25
q(x)	1/3	1/3	1/3

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

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

≈ 0.0852.

$$\begin{aligned}
 D_{KL}(Q || P) &= \frac{1}{25} \frac{1}{3} \ln \left(\frac{\frac{1}{3}}{\frac{9}{25}} \right) + \frac{1}{3} \ln \left(\frac{\frac{1}{3}}{\frac{12}{25}} \right) \\
 &\quad + \frac{1}{3} \ln \left(\frac{\frac{1}{3}}{\frac{4}{25}} \right) \\
 &= \cancel{\frac{1}{3} \ln \left(\frac{1}{3} \times \frac{25}{9} \right)}^{0.076} \\
 &= \frac{1}{3} x - \cancel{0.09} + \frac{1}{3} x - \cancel{0.36} \\
 &\quad + \frac{1}{3} x \cancel{0.733}
 \end{aligned}$$

$$D_{KL}(Q || P) = 0.09$$

* Two kinds of Generative Models -

- i) Fully visible models. \rightarrow considering each pixel value as an observation.
- ii) Latent variable models.