

## Seq. to Seq. Models

→ Transcription, chatgpt, translation etc  
 LLMs comes under this category, etc.

## Reinforcement learning

→ very hard to train, very active  
 → games are classical use case.

Then came the craziness  
 language models, speech recognition,  
 machine translation, conversation modeling,  
 object detection, recognition, etc.  
 visual tracking, image/ video  
 generating authentic faces / new audio etc.

## LAR

Forward pass: input  $\rightarrow$  layers  $\rightarrow$  output

Backward pass: final  $\leftarrow$  layers  $\leftarrow$  update weight

## Tanh:

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

$$f'(z) = (e^z + e^{-z}) \frac{d}{dz} (e^z - e^{-z}) - (e^z - e^{-z}) \frac{d}{dz} (e^z + e^{-z})$$

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

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

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

(2) ReLU

$$\pi(z) = \max(0, z)$$

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

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

(3) leaky ReLU

$$\pi(z) = \begin{cases} a \cdot 0.01(z) & \text{if } z \leq 0 \\ z & \text{if } z > 0 \end{cases}$$

$$\pi'(z) = \begin{cases} 0.01 & \text{if } z \leq 0 \\ 1 & \text{if } z > 0 \end{cases}$$

(4) softmax:

$$s(z)_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \quad \begin{array}{l} \rightarrow \text{specific element in } z \\ \rightarrow \text{sum of exp. of } z \end{array}$$

$$\text{Assume, } \sum_{j=1}^K e^{z_j} = S, \quad s(z)_i = \frac{e^{z_i}}{S}$$

$\rightarrow$  Jacobian matrix

$$J = \left[ \begin{array}{ccc} \frac{\partial s_1}{\partial x_1}, & \frac{\partial s_1}{\partial x_2}, & \frac{\partial s_1}{\partial x_3} \\ \vdots & \ddots & \vdots \\ \frac{\partial s_K}{\partial x_1}, & \frac{\partial s_K}{\partial x_2}, & \frac{\partial s_K}{\partial x_3} \end{array} \right]_{m \times n}$$

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

$$x = \begin{bmatrix} 2.0 \\ 1.0 \\ 0.1 \end{bmatrix}_{3 \times 1}$$

1) compute

softmax output  
jacobian output

2) compute

jacobian

output

softmax output

output

softmax

output

$$\frac{\partial s_i}{\partial x_j} \text{ (case 1)} \Rightarrow s_{ij} = \frac{\partial s_i}{\partial x_j} \quad (\text{diag no 1})$$

case 1: when  $i = j$

$$\frac{\partial s_i}{\partial x_i} = s_{ii}(1 - s_{ii})$$

case 2:  $i \neq j$  (off diag no 1)

$$\text{as } \frac{\partial s_i}{\partial x_j} = -s_{ij}s_{ji}$$

$$J = \begin{bmatrix} 0.2244 & -0.1584 & -0.066 \\ -0.1584 & 0.1824 & -0.027 \\ -0.066 & -0.027 & 0.01 \end{bmatrix}$$

Now how did the derivatives come?

①  $s(x_i) = \sum e^{x_j}$ ,  $\frac{\partial s_i}{\partial x_i} = \frac{\partial}{\partial x_i} \sum e^{x_j} = 1$

But  $s(x_i) = \frac{e^{x_i}}{\sum e^{x_j}} = \frac{e^{x_i}}{\text{Sum}}$  ( $\text{Sum} = x_1 + x_2 + x_3$ )

$x = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \left( \frac{\partial s_i}{\partial x_i} \right)$

$s_1 = \frac{e^{x_1}}{e^{x_1} + e^{x_2} + e^{x_3}}$

$s_2 = \frac{e^{x_2}}{e^{x_1} + e^{x_2}}$

$s_3 = \frac{e^{x_3}}{e^{x_1} + e^{x_2} + e^{x_3}}$

$s_1 = \frac{2.418}{30.192}, s_2 = \frac{7.389}{30.192}, s_3 = \frac{0.20085}{30.192}$

Partial derivatives: ( $i \neq j$ )

$\frac{\partial s_i}{\partial x_j} = \frac{\partial}{\partial x_j} \left( \frac{e^{x_i}}{\sum} \right) = \frac{d}{dx_j} (e^{x_i}) \cdot \text{Sum} - e^{x_i} \left( \frac{d}{dx_j} (\text{Sum}) \right)$

$= \frac{e^{x_i} \cdot \text{Sum} - e^{x_i} \cdot e^{x_i}}{(\text{Sum})^2} = \frac{e^{x_i} (\text{Sum} - e^{x_i})}{(\text{Sum})^2}$

$$\frac{\partial s_1}{\partial x_1} = s(x_i) (1 - s(x_i))$$

$$\frac{\partial s_1}{\partial x_2} = \frac{\partial s_1}{\partial x_2}$$

$$\frac{\partial s_1}{\partial x_2} = \left( \frac{d}{dx_2} (e^{x_1}) \cdot \text{sum} \right) - e^{x_1} \cdot \frac{d}{dx_2} (\text{sum})$$

$$= \frac{0 \cdot \text{sum} - e^{x_1} \cdot e^{x_2}}{(\text{sum})^2} = - \left( \frac{e^{x_1}}{\text{sum}} \right) \left( \frac{e^{x_2}}{\text{sum}} \right)$$

=  $-s_1 \cdot s_2.$

$$\textcircled{1} \quad i=j \quad s = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}, \quad \frac{e^{x_i}}{\text{sum}}$$

$$\frac{\partial s_1}{\partial x_1} = \frac{\partial}{\partial x_1} \left( \frac{e^{x_1}}{\text{sum}} \right)$$

$$= \frac{\frac{d}{dx_1} (e^{x_1}) \cdot \text{sum} - e^{x_1} \frac{d}{dx_1} (\text{sum})}{(\text{sum})^2}$$

$$= \frac{e^{x_1} \cdot \text{sum} - e^{x_1} \cdot e^{x_1}}{(\text{sum})^2} = \frac{e^{x_1} (\text{sum} - e^{x_1})}{(\text{sum})^2}$$

$$= e^{x_1} s (1 - s)$$

$$\textcircled{2} \quad i \neq j \quad \frac{\partial s_1}{\partial x_2} = \frac{\partial}{\partial x_2} \left( \frac{e^{x_1}}{\text{sum}} \right)$$

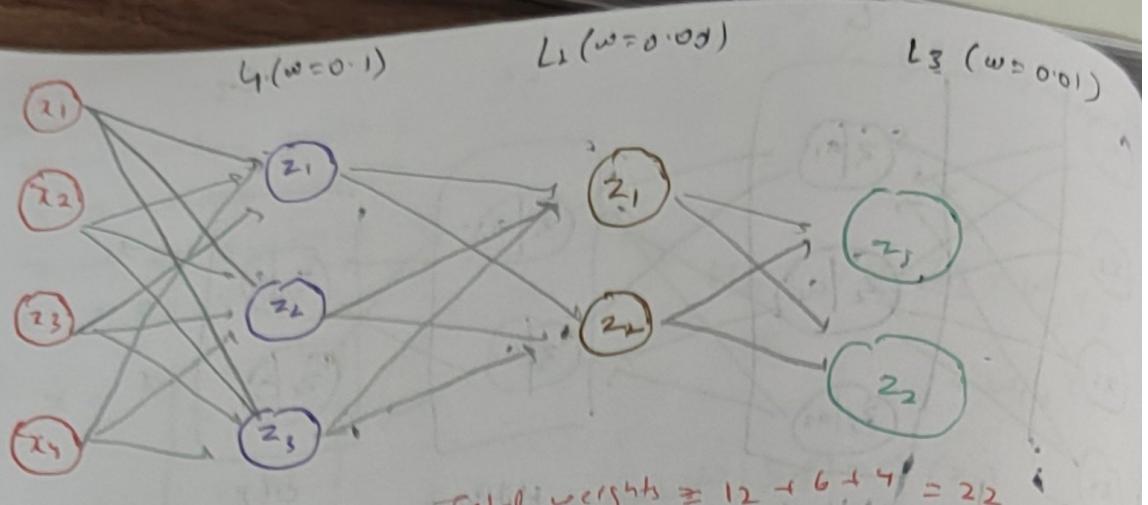
$$= \frac{\frac{d}{dx_2} (e^{x_1}) \cdot \text{sum} - e^{x_1} \cdot \frac{d}{dx_2} (\text{sum})}{(\text{sum})^2} = \frac{0 \cdot \text{sum} - e^{x_1} \cdot e^{x_2}}{(\text{sum})^2}$$

$$= -s \cdot s$$

& & so derivative for softmax:

if  $i \neq j$ ,  $s(i) - s \cdot s$

if  $i = j$ ,  $s(1 - s)$



$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} -2.4 \\ 1.2 \\ -0.8 \\ 1.1 \end{bmatrix}$$

Total weights  $\geq 12 + 6 + 4 = 22$

$$L1: z_1 = \begin{bmatrix} -2.4 \\ 1.2 \\ -0.8 \\ 1.1 \end{bmatrix} \begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \\ 0.1 \end{bmatrix} = \begin{bmatrix} -0.24 \\ 0.12 \\ -0.08 \\ 0.11 \end{bmatrix} = -0.09$$

instead of each time, do matrix multiplication

$$L2: z_2 = w \cdot x = \begin{bmatrix} 0.1 & 0.1 & 0.1 & 0.1 \\ 0.1 & 0.1 & 0.1 & 0.1 \\ 0.1 & 0.1 & 0.1 & 0.1 \\ 0.1 & 0.1 & 0.1 & 0.1 \end{bmatrix} \begin{bmatrix} -0.24 \\ 0.12 \\ -0.08 \\ 0.11 \end{bmatrix} = \begin{bmatrix} -0.09 \\ -0.09 \\ -0.09 \\ -0.09 \end{bmatrix}$$

$$L3: z_3 = w \cdot z_2 = \begin{bmatrix} 0.01 & 0.01 & 0.01 & 0.01 \\ 0.01 & 0.01 & 0.01 & 0.01 \\ 0.01 & 0.01 & 0.01 & 0.01 \\ 0.01 & 0.01 & 0.01 & 0.01 \end{bmatrix} \begin{bmatrix} -0.09 \\ -0.09 \\ -0.09 \\ -0.09 \end{bmatrix} = \begin{bmatrix} -0.09 \\ -0.09 \\ -0.09 \\ -0.09 \end{bmatrix}$$

$$= 1$$

$$\rightarrow \text{ReLU } \begin{bmatrix} -0.09 \\ -0.09 \\ -0.09 \end{bmatrix} = 0.$$

$$z_2 = w \cdot z_1 = \begin{bmatrix} 0.01 & 0.01 & 0.01 & 0.01 \\ 0.01 & 0.01 & 0.01 & 0.01 \\ 0.01 & 0.01 & 0.01 & 0.01 \\ 0.01 & 0.01 & 0.01 & 0.01 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}. \quad \text{ReLU } \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} = 0.$$

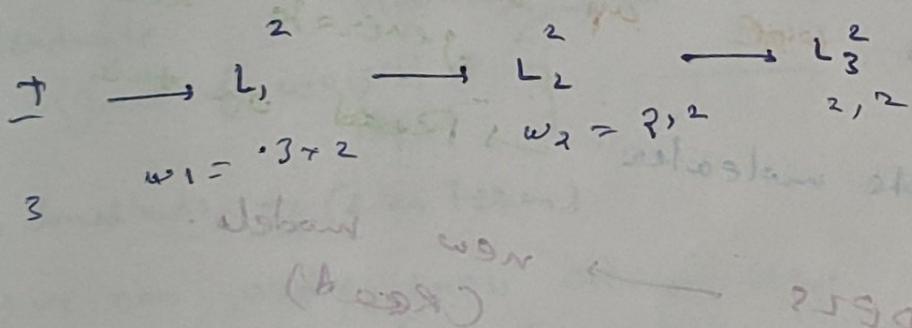
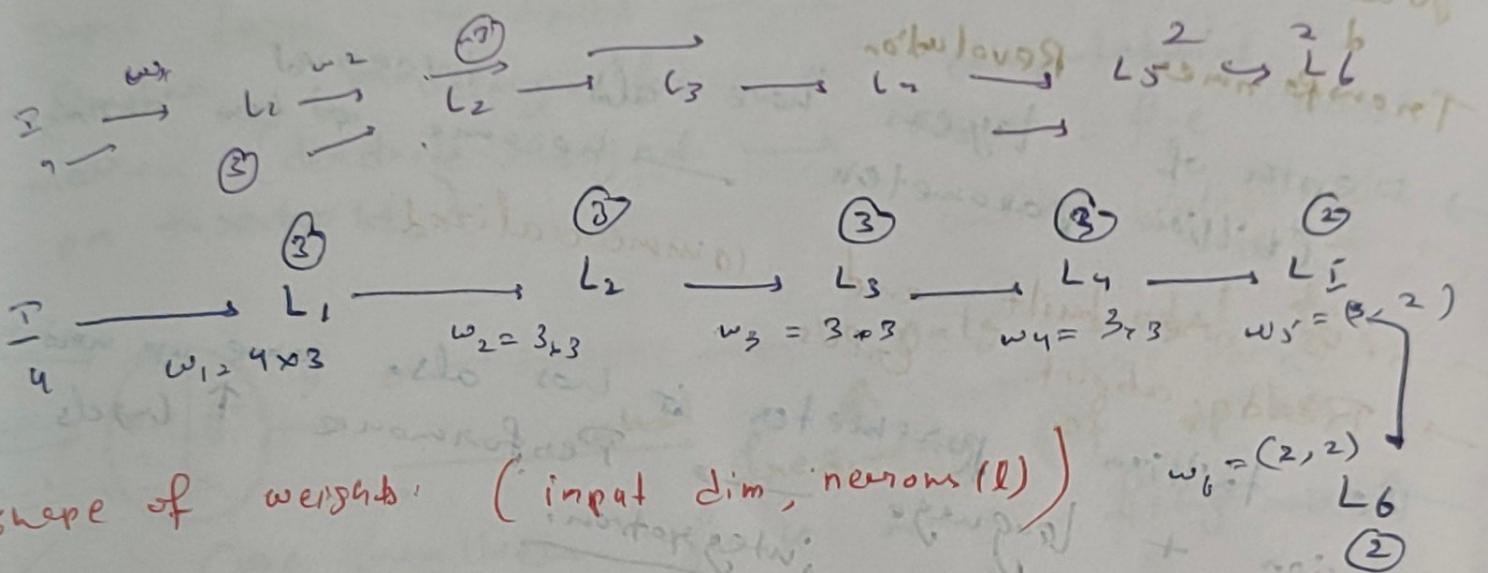
$$z_3 = \begin{bmatrix} 0.01 & 0.01 & 0.01 \\ 0.01 & 0.01 & 0.01 \\ 0.01 & 0.01 & 0.01 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = \text{softmax } \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = ?$$

$$f = \frac{e^x}{e^0 + e^x} = \frac{1}{1+1} = \frac{1}{2} = 0.5$$

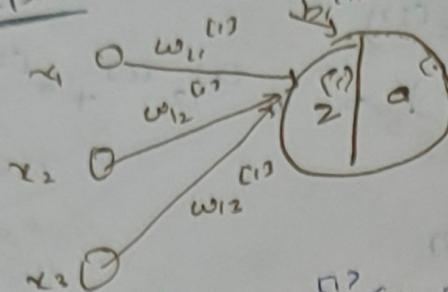
$$\underline{y} = 0.5$$

→ Bias terms are also introduced to the forward pass

→ if u input, 3 neuron → matrix  $(4, 3)$   
 if 3 input, 2 neuron  $(3, 2)$



Perception: single layer straight neural network with no hidden layers here.



A binary classifier

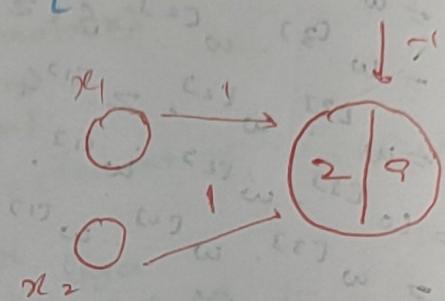
$$z^{(1)} = w_{11}x_1 + w_{12}x_2 + w_{13}x_3 - b_1 \quad a^{(1)} = g(z^{(1)})$$

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

example:

AND

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



①

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

$$= (0)(0) + (0)(1) + (-1)$$

$$= -1$$

②

$$z^{(1)} = (0)(0) + (0)(1) + (-1)$$

$$= 0 + 1 - 1 = 0 \quad g = 0$$

③

$$z^{(1)} = (1)(0) + (0)(1) + (-1) = 0 - 1 = -1 \quad g = 1 \quad (z > 0)$$

④

$$z^{(1)} = (1)(0) + (0)(1) + (-1) = 0 - 1 = -1 \quad g = 1 \quad (z > 0)$$

OR

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

$$b_1 = -1$$

$$\textcircled{1} \quad z^{(1)} = (0)(0) + (1)(0) + (-1)$$

$$= 0 + 0 - 1 = -1 \quad g(-1) = 0$$

$$\textcircled{2} \quad z^{(1)} = (1)(0) + (1)(1) + (-1)$$

$$= 0 + 1 - 1 = 0 \quad g(0) = 0 \quad X$$

$$\textcircled{3} \quad z^{(1)} = (0)(0) + (0)(1) + (1)(1) + (-1)$$

$$= 0 + 0 + 1 - 1 = 0 \quad g(0) = 1 \quad X$$

$$\begin{aligned} b_1 a_1 &= 0 \\ \text{now, } \quad ① \quad z &= 0 \\ ② \quad z &= 0 \\ ③ \quad z &= 0 \\ ④ \quad z &= 0. \end{aligned}$$

$$\begin{aligned} ⑤ \quad z &= (1)(0) + (1)(0) + (-1) = -1, g(-1) = 0. \\ ⑥ \quad z &= (1)(0) + (1)(1) + (-1) = 0, g(0) = 0 \quad \times \\ ⑦ \quad z &= (0)(1) + (1)(0) + (-1) = 0, g(0) = 0 \quad \times \\ ⑧ \quad z &= (1)(1) + (1)(1) + (-1) = 1, g(1) = 1 \end{aligned}$$

$$z = (0)(1) + (0)(1)$$

$$z = 0 \rightarrow$$

↙ shows why it is pdw

→ we can change the bias

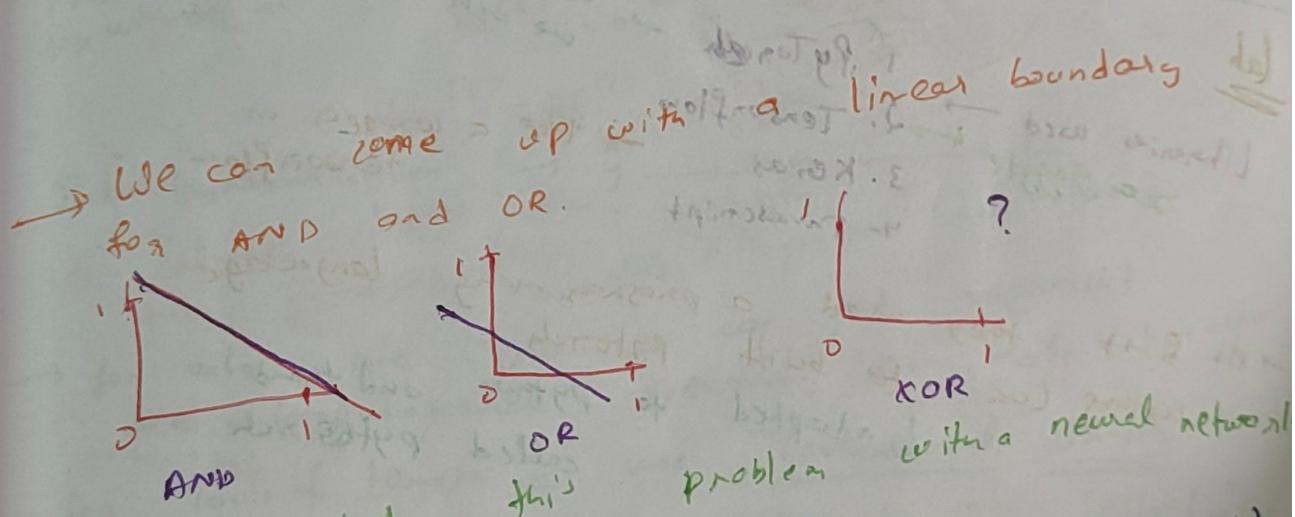
XOR

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

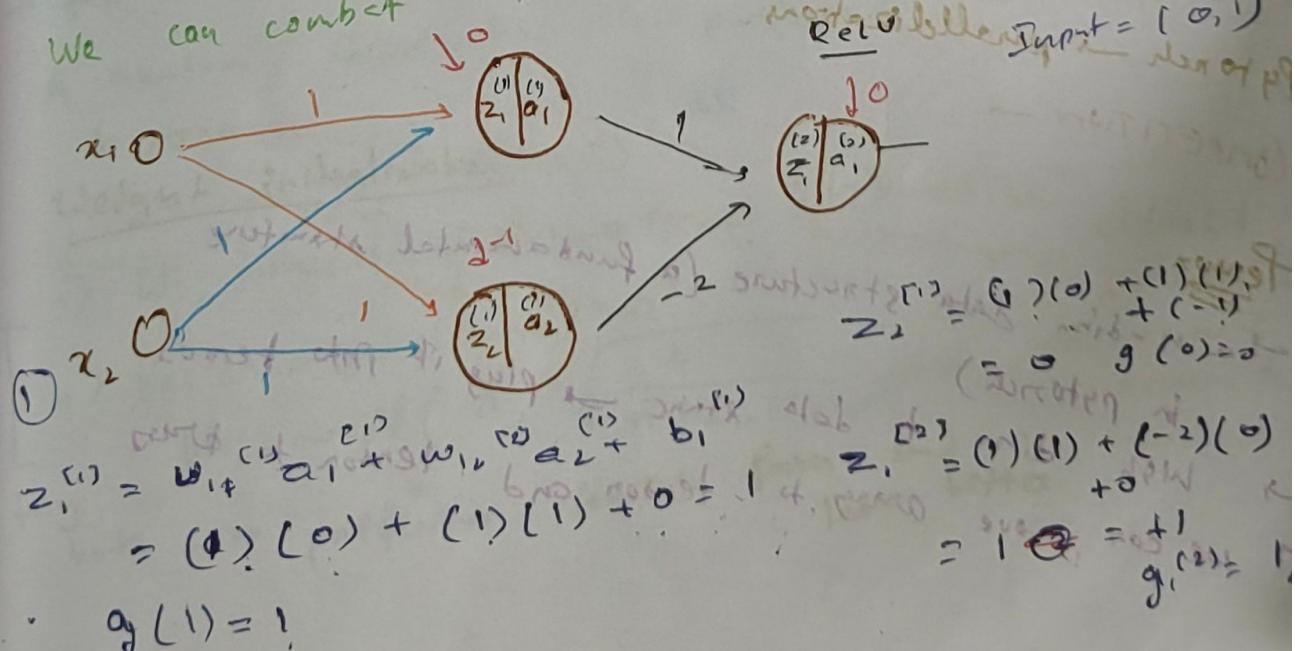
$$\begin{aligned} z &= (1)(1) + (1)(1) + (-2) \\ &= 2 - 2 = 0 \end{aligned}$$

$$\begin{aligned} \text{bias} &= -1 \\ ⑨ \quad z^{(1)} &= (1)(1) + (1)(1) \\ &\quad + (-2) \\ &= 2 - 1 = 1 \\ g(1) &= 1 \end{aligned}$$

$$\begin{aligned} z^{(1)} &= w_{11}x_1 + w_{12}x_2 + b_1 \\ z &= (1)(1) + (1)(1) + (-2) \\ &= 1 + 1 - 2 = 0 \end{aligned}$$



We can combat



$$\begin{aligned} z_1^{(1)} &= w_{11}^{(1)}a_1^{(1)} + w_{12}^{(1)}a_2^{(1)} + b_1 \\ &= (1)(0) + (1)(1) + 0 = 1 \end{aligned}$$

$$\begin{aligned} z_2^{(1)} &= (1)(0) + (-2)(1) + 0 \\ &= 1 - 2 = -1 \end{aligned}$$

$$\begin{aligned} g_1(1) &= 1 \\ g_2(-1) &= 0 \end{aligned}$$

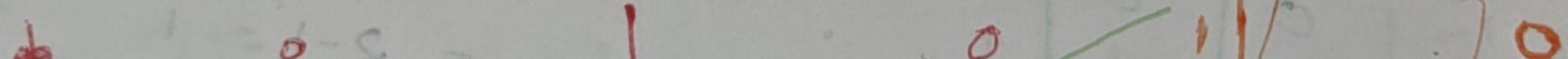
$$g_2(-1) = 0$$

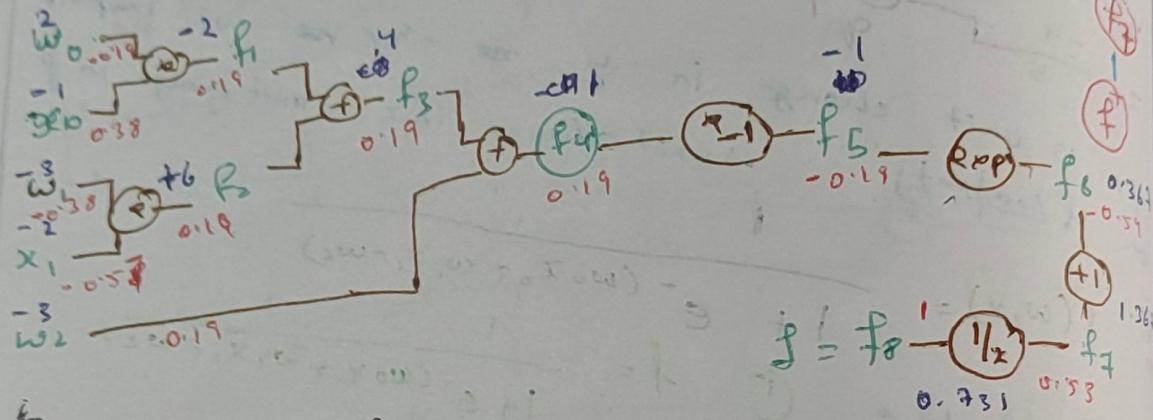
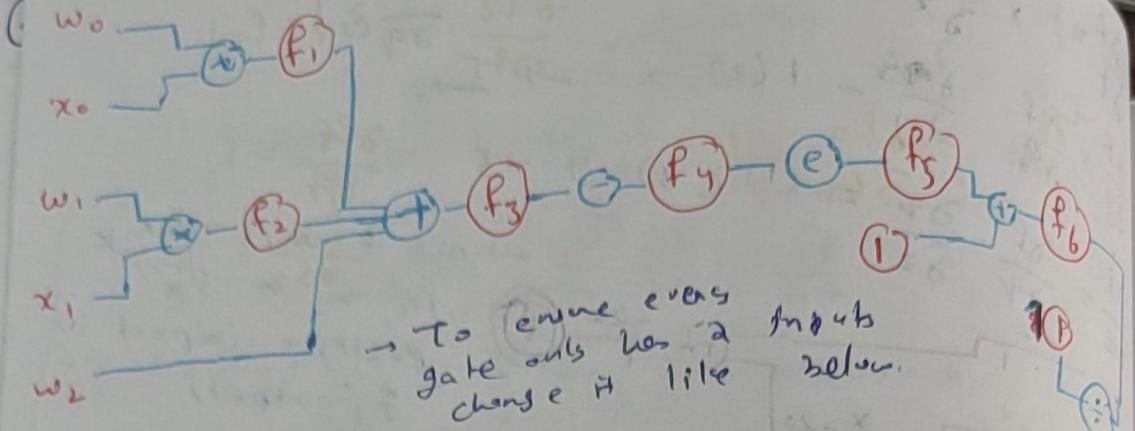
$$\textcircled{2} \quad \text{Input } (1, 1) \quad z_1^{(1)} = 0(1) + 1(1) + 0 = 2 \quad z_2^{(1)} = 1(1) + 0(1) + (-1) \\ g(4) = 1 \quad = +, \text{ SC } \hookrightarrow 2, \\ z_1^{(2)} = 1(1) + 1(-2) = 1 - 2 = -1, \quad g(-1) = 0 \\ +0 \quad g = 0 \text{ if } .$$

Now it give a correct answer with a hidden layer.

why is this working?

we transform our input data (from original space to new space)  $\rightarrow$  in this space there are linearly separable.





Computing gradients: (3) forward pass:

$$f_1 = 2(-1) = -2, \quad f_2 = 6, \quad f_3 = 4, \quad f_4 = 1, \quad f_5 = 1, \quad f_6 = 0.36, \quad f_7 = 0.39, \quad f_8 = 0.73$$

(1) Compute gradients

~~$$\frac{\partial f}{\partial f_8} = 1, \quad \frac{\partial f}{\partial f_7} = \frac{1}{f_7}$$~~

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

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

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

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

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

$$f_6 = \frac{\partial f_6}{\partial F_5} = \exp(F_5)$$

$$f_6 = \frac{\partial f_6}{\partial F_5} \quad f_6 = \exp(F_5) \quad \frac{\partial f_6}{\partial F_5} = \exp(F_5)$$

$$f_7 = f_6 + 1, \quad \frac{\partial f_7}{\partial F_6} = 1 \quad f_8 = f_8 \quad \frac{\partial f_8}{\partial F_7} = 1$$

now calc gradient

$$f_8 = 1, \quad f_7 = \frac{1}{(f_7)^2} = \frac{1}{(0.36)^2} = -0.56$$

$$f_6 = \frac{\partial f_6}{\partial F_7} = \text{local } \times \text{global} = -\frac{1}{f_7^2} \times 1 = -0.56$$

$$\frac{\partial f_8}{\partial f_6} = \frac{\partial f_8}{\partial F_7} \cdot \frac{\partial F_7}{\partial f_6} = (-0.56) (1) = -0.56$$

$$\frac{\partial f_8}{\partial f_5} = \frac{\partial f_8}{\partial F_7} \text{ local } \times \text{global} = \frac{\partial f_6}{\partial f_5} \cdot \frac{\partial f_7}{\partial F_6} = \exp(-1) \cdot (-0.56) \\ = (0.36)(-0.56) \\ = -0.1944 \\ = -0.19$$

$$\frac{\partial f_8}{\partial f_4} = \frac{\partial f_5}{\partial f_4} \cdot \frac{\partial f_8}{\partial f_5} = (-1)(-0.19) = 0.19$$

$$\frac{\partial f_8}{\partial f_3} = \frac{\partial f_4}{\partial f_3} \cdot \frac{\partial f_8}{\partial f_4} = (1)(0.19) = 0.19$$

$$\frac{\partial f_8}{\partial w_2} = \frac{\partial f_3}{\partial w_2} \cdot \frac{\partial f_8}{\partial f_3} = (1)(0.19) = 0.19$$

$$\frac{\partial f_8}{\partial f_1} = \frac{\partial f_3}{\partial f_1} \cdot \frac{\partial f_8}{\partial f_3} = (1)(0.19) = 0.19$$

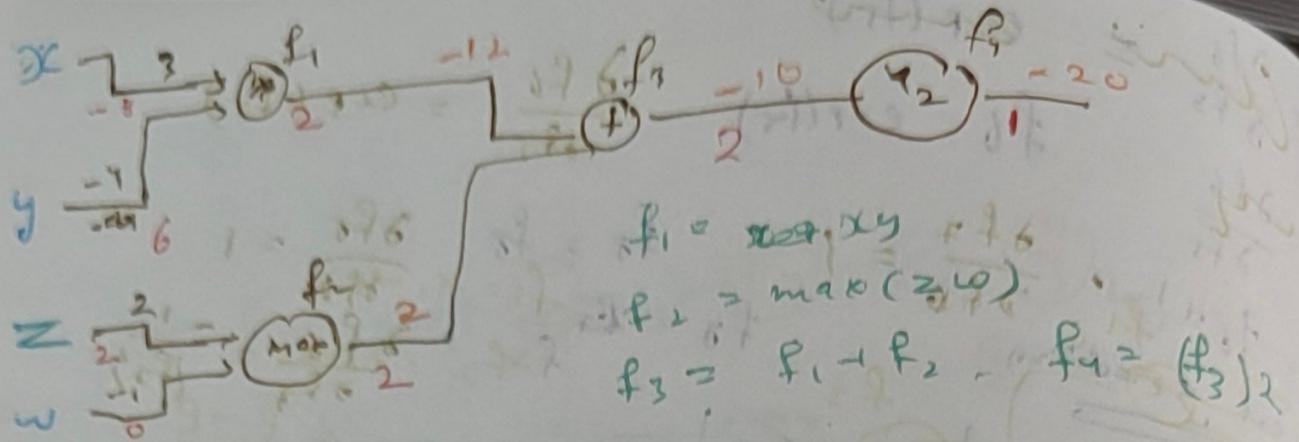
$$\frac{\partial f_8}{\partial f_2} = \frac{\partial f_3}{\partial f_2} \cdot \frac{\partial f_8}{\partial f_3} = (1)(0.19) = 0.19$$

$$\frac{\partial f_8}{\partial w_0} = \frac{\partial f_1}{\partial w_0} \cdot \frac{\partial f_8}{\partial f_1} = (-1)(0.19) = -0.19$$

$$\frac{\partial f_8}{\partial x_0} = \frac{\partial f_1}{\partial x_0} + \frac{\partial f_8}{\partial f_1} = (2)(0.19) = 0.38$$

$$\frac{\partial f_8}{\partial w_0} = \frac{\partial f_2}{\partial w_0} \cdot \frac{\partial f_8}{\partial f_2} = (x_1)(0.19) = -0.38$$

$$\frac{\partial f_8}{\partial x_1} = \frac{\partial f_2}{\partial x_1} \cdot \frac{\partial f_8}{\partial f_2} = (w_1)(0.19) = -0.57$$



local!  $f_p = x \cdot y$

$$\frac{\partial f_i}{\partial x} = y, \quad \frac{\partial f_i}{y}$$

$$f_3 = \text{max}(\omega)$$

$$f_2 \rightarrow \max(z_{\text{out}})$$

$$P = P_1 + P_2$$

$$f_3 = f_1 - f_2 \quad f_4 = (f_3)_2$$

$$f_3 = f_1 + f_2$$

$$\frac{\partial f_2}{\partial h_1} = 1 \quad \text{and} \quad \frac{\partial f_2}{\partial h_2}$$

$$f_4 = (\mathbb{P}_3)^2, \frac{\partial f_4}{\partial f_2} = 2$$

$$\underline{\partial f_2} = 1 \text{ if } z > 6$$

$$\frac{\partial z}{\partial z} = 0 \text{ if } z < w \\ = \text{undef if } z = w \\ = -\infty \text{ if } z > w$$

$$\frac{z}{w} = \begin{cases} 1 & \text{if } w \neq 0 \\ 0 & \text{if } w < 0 \\ \text{undefined} & \text{if } w = 0 \end{cases}$$

(arbitrary)

Labeled :

$$\frac{\partial f_4}{\partial f_1} = 1, \quad \frac{\partial f_4}{\partial f_3} = 2, \quad \frac{\partial f_4}{\partial f_1} = \frac{\partial f_3}{\partial f_1} = \frac{\partial f_3}{\partial f_1} = \frac{\partial f_3}{\partial f_1} \cdot \frac{\partial f_4}{\partial f_3}$$

$$\frac{\partial f_4}{\partial x} = \frac{\partial f_4}{\partial t} \cdot \frac{\partial t}{\partial x} = (y)(2) - e^t(2)$$

$$= (1)(2) = 2$$

$$\frac{\partial f_1}{\partial x} = \frac{\partial f_1}{\partial y} \cdot \frac{\partial f_2}{\partial h} = (2)(2) = (3)(2) > 6$$

$$\text{Now } / \quad \frac{\partial f_4}{\partial f_2} = \frac{\partial f_3}{\partial f_2} \cdot \frac{\partial f_4}{\partial f_3} = (1)(2) = 2$$

$$\frac{\partial f_1}{\partial w} = \frac{\partial f_2}{\partial w} \cdot \frac{\partial f_4}{\partial f_2} = \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}(2) = 0 \text{ if } z \geq w \\ = 2 \text{ if } z < w$$

Histogram of oriented gradients (HOG)  
 Histogram of oriented flows (HOF)  
 Optical flow features

Features w/  
at first  
Clustering

Most natural images  $\rightarrow$  more?

↳ gradient pixel values

↳ adjacent gradients are generally:  $43, -24, 1$

↳ homogeneous before  $\times$  (not edge)

↳ only if strong edge comes in

This becomes saturated so

CNN

$\rightarrow$  scaling is a problem if we wanna use a CNN

Suppose,  $1000 \times 1000 \times 3$  image  $\rightarrow 3M \times 600$

= 3 billion parameters

↳ so CNN is better, weights depends on size heavily here which is a problem, overfitting

→ Convert

① Convolution Operation

3x3 kernel			1x1 stride		
3	0	1	2	3	4
1	5	8	4	3	1
2	7	2	5	1	3
0	1	3	1	4	8
4	2	1	6	2	8
2	9	5	2	3	9

input  $6 \times 6$  image

Convolved operator

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

Filter w/ kernel  $(3 \times 3)$  and stride 1 instead of 2  
neg values and constant instead

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

$$= 3 + 0 - 1 + 1 + 0 - 8 + 2 + 0 - 2$$

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

slide left by one position.

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

$$\begin{aligned}
 &= (0 \times 1) + (0 \times 1) + (2 \times -1) \\
 &+ (5 \times 1) + (8 \times 0) + (9 \times -1) + \\
 &(7 \times 1) + (2 \times 0) + 5 \times 1 \\
 &= 0 + 0 - 2 + 5 + 0 - 9 + 7 - 7 \\
 &= -14
 \end{aligned}$$

→ now slide again

1	2	7
8	9	3
2	3	1

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

2	7	4
9	3	1
5	1	3

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

1	0	1
1	0	1
1	0	1

-10

1	5	8
2	3	2
0	1	5

$$5 - 9 + 7 - 5 + 1 - 1$$

8	9	3
2	5	1
3	1	7

$$\begin{aligned}
 &8 - 3 + 2 - 1 \\
 &+ 3 - 7 \\
 &= 2
 \end{aligned}$$

5	8	9
7	2	5
1	3	1

$$= -2$$

2	7	2
0	1	3
4	2	1

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

1	3	1
5	1	3
1	2	8

$$\begin{aligned}
 &9 + -1 + 5 - 3 + 1 \\
 &- 8
 \end{aligned}$$

2	5	1
3	1	2
1	6	2

$$\begin{aligned}
 &2 - 1 + 3 - 7 \\
 &+ 1 - 2 \\
 &= -4
 \end{aligned}$$

9	2	5
1	3	1
2	1	6

$$\begin{aligned}
 &7 - 5 + 1 - 1 \\
 &+ 2 - 6
 \end{aligned}$$

0	1	3
4	2	1
2	7	5

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

5	1	3
1	7	8
6	2	3

$$\begin{aligned}
 &5 - 3 + 1 \\
 &- 8 + 6 - 8
 \end{aligned}$$

3	7	2
1	6	2
5	2	3

$$\begin{aligned}
 &3 - 2 + 1 - 2 \\
 &+ 5 - 3 \\
 &= -3
 \end{aligned}$$

1	3	1
2	1	6
7	5	2

$$\begin{aligned}
 &1 - 1 + 2 \\
 &- 6 + 4 - 2
 \end{aligned}$$

1	2	8
6	2	8
2	3	9

$$\begin{aligned}
 &1 + 8 + 6 - 8 \\
 &+ 2 - 9 \\
 &= -10
 \end{aligned}$$

we got  $6 \times 6 \rightarrow$  image back.

Formula to compute,

$$\text{O/p size} = (n + 2p - f + 1)$$

choice is a hyperparameter

"Valid" convolution  $\Rightarrow$  O/P size  $\neq$  O/P size

Same convolution  $\Rightarrow$  O/P size = O/P size

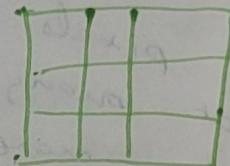
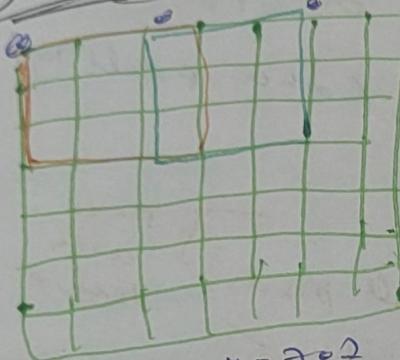
what should

value of p to get same convolution

$$p = \frac{f-1}{2}$$

③

### Strided Convolution



stride,  $s=2$  (move by 2 positions, not 1)

even if we lose info (cannot over all)  
is it okay?

→ output size  $\downarrow$  → add padding?

Stride is based on how much attention you want to give to the pixels (which pixels?)

Example on padding and striding

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

$n = 5 \times 5$

$p = 1, s = 2$

+1	-1
0	-1

$f = 2 \times 2$

$2 \times 2$

$$\begin{array}{l}
 \textcircled{1} \quad \begin{matrix} 0 & 0 \\ 0 & 2 \end{matrix} + \begin{matrix} 1 & -1 \\ 0 & -1 \end{matrix} = \begin{matrix} 0 & 0 \\ 0 & -2 \end{matrix} = \begin{matrix} -2 \\ \cancel{-2} \end{matrix} \\
 \textcircled{2} \quad \begin{matrix} 0 & 0 & 0 \\ 0 & 2 & 0 \end{matrix} + \begin{matrix} 1 & -1 \\ 0 & -1 \end{matrix} = \begin{matrix} 0 & 0 \\ 0 & -2 \end{matrix} = -7 \\
 \textcircled{3} \quad \begin{matrix} 0 & 0 \\ 4 & 6 \end{matrix} + \begin{matrix} 1 & -1 \\ 0 & -1 \end{matrix} = \begin{matrix} 0 & -6 \\ 0 & -6 \end{matrix} = -6 \\
 \textcircled{4} \quad \begin{matrix} 0 & 6 \\ 0 & 3 \end{matrix} + \begin{matrix} 1 & -1 \\ 0 & -1 \end{matrix} = \begin{matrix} 0 & -3 \\ 0 & -3 \end{matrix} = -9 \\
 \textcircled{5} \quad \begin{matrix} 0 & 9 \\ 8 & 7 \\ 3 & 8 \end{matrix} + \begin{matrix} 1 & -1 \\ 0 & -1 \end{matrix} = \begin{matrix} 6 & -9 \\ 7 & -8 \\ 0 & -7 \end{matrix} = -11 \\
 \textcircled{6} \quad \begin{matrix} 0 & 7 \\ 0 & 4 \\ 8 & 2 \\ 2 & 1 \end{matrix} + \begin{matrix} 1 & -1 \\ 0 & -1 \end{matrix} = \begin{matrix} 8 & -3 \\ 7 & -1 \\ 6 & -6 \end{matrix} = 4 \\
 \textcircled{7} \quad \begin{matrix} 0 & 7 \\ 0 & 4 \\ 8 & 2 \\ 2 & 1 \end{matrix} + \begin{matrix} 1 & -1 \\ 0 & -1 \end{matrix} = \begin{matrix} 8 & -3 \\ 7 & -1 \\ 6 & -6 \end{matrix} = -3
 \end{array}$$

$$\begin{matrix}
 -2 & -7 & -6 \\
 -9 & -11 & -7 \\
 -11 & 4 & -3
 \end{matrix}$$

3x3

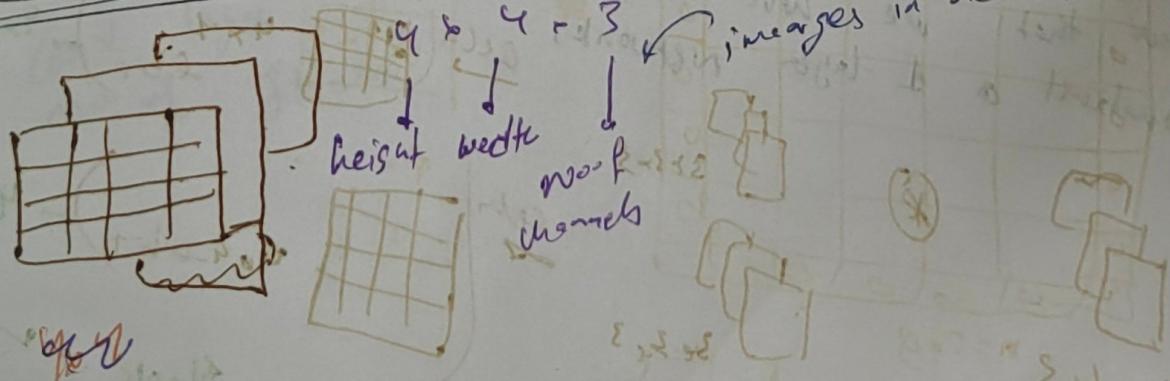
so with the padding and  
striding, the dim formula of output  
becomes,

$$\text{of dim} = \left[ \frac{n+2p-f+1+(s+1)}{s} + 1 \right]$$

take floor value (round up  
decimals)

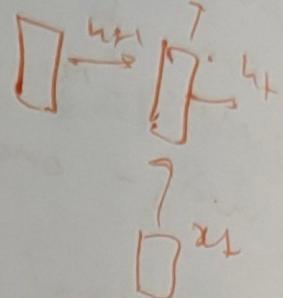
Convolution on a 3D Volume

we use 3D (RGB)



Temporal dependency → handled  
 we need min input by h to something.  
 something generated by h  
 → the choice of RNN block is optional coming out of denoted by DSK  
 on problem network

$$\hat{y}_t = f_w(h_{t-1}, x_t)$$



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

$$h_t = \tanh \left( \underbrace{w_{hh} h_{t-1}}_{4 \times 4} + \underbrace{w_{xh} x_t}_{3 \times 1} \right)$$

$$w_{hh} \in \mathbb{R}^{4 \times 1}$$

→ weight matrix is shared across timesteps.

$$② y_t = w_{hy} h_t$$

$$y_t = \underbrace{\begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}}_{w_{hy} \in \mathbb{R}^{2 \times 4}} \underbrace{\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}}_{h_t \in \mathbb{R}^{4 \times 1}}$$

→ input is fixed size.

## Character-level Language Model

T=2 (two time steps)

$$x_1 \in \mathbb{R}^2, x_2 \in \mathbb{R}^2$$

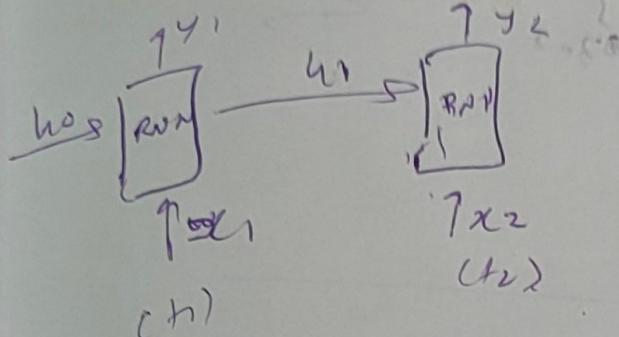
$$\text{input dim} = 2 \Rightarrow x_t \in \mathbb{R}^2$$

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

$$h_0 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, w_{xh} = \begin{bmatrix} 0.5 & 0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.4 \end{bmatrix}, w_{hh} = \begin{bmatrix} 0.1 & 0.2 & 0.4 \\ 0.2 & 0.3 & 0.1 \\ 0.0 & 0.1 & 0.2 \end{bmatrix}$$

$$y^{(t)} \in \mathbb{R}^2 \Rightarrow w_{hy} = \begin{bmatrix} 0.0 & -0.5 \\ 0.5 & 0.5 & -0.5 \end{bmatrix}_{3 \times 3}$$

$$w_{nn} = \begin{bmatrix} 0.1 & 0.4 & 0.0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.2 \end{bmatrix}_{3 \times 3}$$



compute,  $h_1, y_1$ ,  
 $h_2, y_2$

$$h_t = \tanh(w_{nh} h_{t-1} + w_{nx} x_t + b)$$

$$y_t = w_{hy} h_t$$

$$h_t = \tanh \left[ \begin{bmatrix} 0.1 & 0.4 & 0.0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.2 \end{bmatrix}_{3 \times 3} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}_{3 \times 1} + \begin{bmatrix} 0.5 & -0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.7 \end{bmatrix}_{3 \times 2} \begin{bmatrix} 1 \\ 2 \end{bmatrix}_{2 \times 1} \right]$$

$$= \tanh \left( \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}_{3 \times 1} + \begin{bmatrix} -0.1 \\ 1.2 \\ 0.9 \end{bmatrix}_{3 \times 1} \right)$$

$$= \tanh \begin{bmatrix} -0.1 \\ 1.2 \\ 0.9 \end{bmatrix}_{3 \times 1} = \begin{bmatrix} -0.099 \\ 0.833 \\ 0.416 \end{bmatrix}_{3 \times 1}$$

$$y_t = \begin{bmatrix} 0.4 & -0.2 & 0.1 \\ 0.5 & 0.15 & 0.15 \end{bmatrix}_{3 \times 3} \begin{bmatrix} 1.0 & -1.0 & 0.5 \\ 0.5 & 0.15 & 0.15 \end{bmatrix}_{3 \times 3} \begin{bmatrix} -0.099 \\ 0.833 \\ 0.416 \end{bmatrix}_{3 \times 1} \begin{bmatrix} -0.01 \\ 0.83 \\ 0.42 \end{bmatrix}_{3 \times 1}$$

$$\begin{bmatrix} -0.48 \\ 0.05 \end{bmatrix} = \begin{bmatrix} -0.58 \\ 0.05 \end{bmatrix}_{2 \times 1} = \begin{bmatrix} -0.57 \\ 0.008 \end{bmatrix}_{2 \times 1}$$

$$\textcircled{2} \quad h_t = \tanh \left[ \begin{bmatrix} 0.1 & 0.4 & 0.0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.2 \end{bmatrix}_{3 \times 3} \begin{bmatrix} 0.01 \\ 0.83 \\ 0.42 \end{bmatrix}_{3 \times 1} + \begin{bmatrix} 0.5 & -0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.9 \end{bmatrix}_{3 \times 2} \begin{bmatrix} -0.099 \\ 0.833 \\ 0.416 \end{bmatrix}_{3 \times 1} \right] = -1$$

$$= \tanh \begin{bmatrix} 0.1831 \\ 0.395 \end{bmatrix} + \begin{bmatrix} -0.18 \\ 0.1612 \end{bmatrix}_{2 \times 1} = \tanh \begin{bmatrix} -0.409 \\ 0.995 \\ 0.3605 \end{bmatrix}_{3 \times 1}$$

$$u_t = \begin{pmatrix} -0.437 \\ 0.759 \\ 0.345 \end{pmatrix}$$

$$y_t = \begin{pmatrix} 1 & -1 & 0.5 \\ 0.5 & 0.5 & -0.5 \\ 0.5 & 0.5 & 0.5 \end{pmatrix} \begin{pmatrix} -0.437 \\ 0.759 \\ 0.345 \end{pmatrix} = 0.302$$

$$= 8^{\circ}$$

Time step ② ( $x_{\text{new}} + h_t \Delta \text{new}$ )  $\rightarrow$  new AN end

$$h_t = \tanh(w_{ht-1} + w_{ht} x_t)$$

$$= \tanh \left( \begin{pmatrix} 0.1 & 0.4 & 0.0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.1 \end{pmatrix} \begin{pmatrix} -0.029 \\ 0.83 \\ 0.716 \end{pmatrix} + \begin{pmatrix} 0.5 & 0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.6 \end{pmatrix} \begin{pmatrix} -1 \\ 1 \end{pmatrix} \right)$$

$$= \tanh \left[ \begin{pmatrix} 0.3221 \\ 0.912 \\ 0.05525 \end{pmatrix} + \begin{pmatrix} -0.8 \\ -0.6 \\ 0.3 \end{pmatrix} \right] = \tanh \begin{pmatrix} -0.443 \\ 0.188 \\ 0.370 \end{pmatrix}$$

$$h_t = \begin{pmatrix} -0.443 \\ 0.188 \\ 0.370 \end{pmatrix}$$

$$w_t = \begin{pmatrix} -0.4434 \\ -0.1889 \\ 0.390 \end{pmatrix}_{D=1}$$

$$y_t = w_{hs} h_t$$

$$\begin{pmatrix} 1 & -1 & 0.5 \\ 0.5 & 0.5 & -0.5 \\ 0.5 & 0.5 & 0.5 \end{pmatrix} \begin{pmatrix} -0.4434 \\ -0.1889 \\ 0.390 \end{pmatrix}_{D=1} = \begin{pmatrix} 18.8 & 0 \\ 18.8 & 0 \\ 18.8 & 0 \end{pmatrix}_{D=1}$$

$$\begin{pmatrix} 0.05 & 0.1 \\ 0.05 & 0 \\ 0.05 & 0 \end{pmatrix} \begin{pmatrix} 18.8 & 0 \\ 18.8 & 0 \\ 18.8 & 0 \end{pmatrix}_{D=1} = \begin{pmatrix} 18.8 & 0 \\ 18.8 & 0 \\ 18.8 & 0 \end{pmatrix}_{D=1}$$

$$\begin{aligned} \text{Extrinsic} \\ x_t &= [0.5, -0.1]^T \\ h_{t-1} &= [0, 0.1]^T \\ c_{t-1} &= [0.2, -0.2]^T \\ w_{xi} &= \begin{bmatrix} 0.5 & -0.3 \\ 0.1 & 0.1 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} w_{hi} &= \begin{bmatrix} 0.1 & 0.2 \\ -0.2 & 0.5 \end{bmatrix} \\ w_{xf} &= \begin{bmatrix} -0.4 & 0.2 \\ 0.5 & 0.3 \end{bmatrix} \\ w_{hf} &= \begin{bmatrix} 0.05 & -0.1 \\ 0.12 & 0.1 \end{bmatrix} \end{aligned}$$

$$w_{x_0} = \begin{bmatrix} 0.3 & 0.25 \\ -0.2 & 0.2 \end{bmatrix}$$

$$w_{h_0} = \begin{bmatrix} 0.15 & 0.05 \\ 0.1 & -0.2 \end{bmatrix}$$

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

$$w_{hs} = \begin{bmatrix} 0.2 & 0.1 \\ -0.1 & 0.5 \end{bmatrix}$$

(not top)  
w

h<sub>t</sub> = ?

c<sub>t</sub> = ?

① → element  
wire multipl.  
cation.

$$\rightarrow \text{lik} \\ \text{Extrinsic} : O_t \odot \tanh(c_t)$$

$$h_t = O_t \odot (c_{t-1} + s_t \odot s_t)$$

$$c_t = f_t \odot (c_{t-1} + s_t \odot s_t)$$

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

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

$$O_t = \sigma (w_{h_0} h_{t-1} + w_{x_0} x_t)$$

$$\approx \tanh(\omega_{hs} h_{t-1} + \omega_{xs} x_t)$$

$$y_t = \sigma(\omega_{hs} h_{t-1} + \omega_{xs} x_t)$$

$$= \sigma \left[ \begin{bmatrix} 0.05 & -0.1 \\ 0.3 & 0.3 \end{bmatrix} \begin{bmatrix} 0 \\ 0.1 \end{bmatrix} + \begin{bmatrix} 0.7 & 0.2 \\ 0.2 & 0.5 \end{bmatrix} \begin{bmatrix} 0 \\ 0.1 \end{bmatrix} \right]$$

$$= \sigma \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} -0.1 \\ 0.03 \end{bmatrix} + \begin{bmatrix} -0.9 \\ 0.1 \end{bmatrix} \right) \rightarrow 3)$$

$$= \sigma \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0.22 \\ -0.15 \end{bmatrix} \right)$$

$$= \sigma \left( \begin{bmatrix} 0.22 \\ -0.15 \end{bmatrix} \right) \approx \begin{bmatrix} 0.49 \\ 0.42 \end{bmatrix}$$

$$= \begin{bmatrix} 0.49 \\ 0.42 \end{bmatrix} / 0.53$$

$$\overline{y_t} = \sigma(\omega_{hs} h_{t-1} + \omega_{xs} x_t)$$

$$= \sigma \left[ \begin{bmatrix} 0.5 & 0.2 \\ -0.2 & 0.5 \end{bmatrix} \begin{bmatrix} 0 \\ 0.1 \end{bmatrix} + \begin{bmatrix} 0.5 & -0.3 \\ 0.4 & 0.1 \end{bmatrix} \begin{bmatrix} 0 \\ 0.1 \end{bmatrix} \right]$$

$$= \sigma \left[ \begin{bmatrix} 0.5 & 0.2 \\ -0.2 & 0.5 \end{bmatrix} \begin{bmatrix} 0 \\ 0.1 \end{bmatrix} + \begin{bmatrix} 0.28 \\ 0.18 \end{bmatrix} \right] = \sigma \begin{bmatrix} 0.53 \\ 0.23 \end{bmatrix}$$

$$= \begin{bmatrix} 0.53 \\ 0.23 \end{bmatrix}$$

$$= \begin{bmatrix} 0.53 \\ 0.23 \end{bmatrix} \approx \sigma(\omega_{hs} h_{t-1} + \omega_{xs} x_t)$$

$$= \sigma \left( \omega_{hs} h_{t-1} + \omega_{xs} x_t \right) = \sigma \left[ \begin{bmatrix} 0.3 & 0.25 \\ 0.2 & 0.2 \end{bmatrix} \begin{bmatrix} 0 \\ 0.1 \end{bmatrix} \right]$$

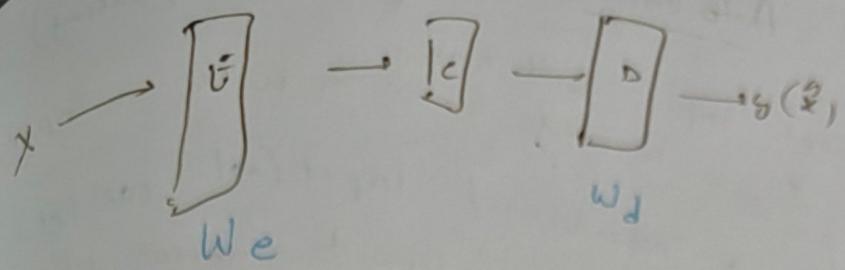
$$= \sigma \left[ \begin{bmatrix} 0.15 & 0.05 \\ 0.1 & -0.2 \end{bmatrix} \begin{bmatrix} 0 \\ 0.1 \end{bmatrix} + \begin{bmatrix} 0.125 \\ 0.08 \end{bmatrix} \right] = \sigma \begin{bmatrix} 0.13 \\ 0.06 \end{bmatrix}$$

$$= \sigma \left( \begin{bmatrix} 0.15 & 0.05 \\ 0.1 & -0.2 \end{bmatrix} \begin{bmatrix} 0 \\ 0.1 \end{bmatrix} + \begin{bmatrix} 0.13 \\ 0.06 \end{bmatrix} \right) = \sigma \begin{bmatrix} 0.13 \\ 0.06 \end{bmatrix}$$

$$= \sigma \begin{bmatrix} 0.13 \\ 0.06 \end{bmatrix} = 0.46$$

$$\begin{aligned}
 g_t &= \tanh(w_{hg} h_{t-1} + w_{xg} x_t) \\
 &= \tanh \left( \begin{bmatrix} 0.12 & 0.11 \\ -0.11 & 0.5 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} + \begin{bmatrix} -0.5 & 0.7 \\ 0.2 & 0.7 \end{bmatrix} \begin{bmatrix} 0.5 \\ -0.1 \end{bmatrix} \right) \\
 &= \tanh \left( \begin{bmatrix} 0.12 \\ 0.05 \end{bmatrix} + \begin{bmatrix} 0.29 \\ 0.07 \end{bmatrix} \right) = \tanh \begin{bmatrix} 0.28 \\ 0.12 \end{bmatrix} \\
 &= \begin{bmatrix} 0.22 \\ 0.12 \end{bmatrix} \quad \text{stretches} \\
 &\quad \text{Rounds down} \\
 c_t &= f_t \odot c_{t-1} + i_t \circ g_t \\
 &= \begin{bmatrix} 0.49 \\ 0.55 \end{bmatrix} \odot \begin{bmatrix} 0.2 \\ -0.2 \end{bmatrix} + \begin{bmatrix} 0.57 \\ 0.59 \end{bmatrix} \odot \begin{bmatrix} -0.28 \\ 0.13 \end{bmatrix} \\
 &= \begin{bmatrix} 0.083 \\ -0.106 \end{bmatrix} + \begin{bmatrix} -0.1539 \\ 0.0702 \end{bmatrix} \\
 &= \begin{bmatrix} 0.083 \\ -0.106 \end{bmatrix} = \begin{bmatrix} -0.06 \\ 0.03 \end{bmatrix} + \tanh(c_t) = \begin{bmatrix} -0.059 \\ -0.021 \end{bmatrix}
 \end{aligned}$$

$$\begin{aligned}
 h_t &= o_t \odot \tanh(c_t) \\
 &= \begin{bmatrix} 0.43 \\ 0.46 \end{bmatrix} \odot \begin{bmatrix} -0.059 \\ -0.021 \end{bmatrix} = \begin{bmatrix} 0.05 \\ 0.05 \end{bmatrix} \\
 &= \begin{bmatrix} -0.03 \\ -0.01 \end{bmatrix} \odot \begin{bmatrix} 0.05 \\ 0.05 \end{bmatrix} = \begin{bmatrix} -0.03 \\ -0.01 \end{bmatrix} \\
 &= \begin{bmatrix} 0.05 \\ 0.05 \end{bmatrix} + \begin{bmatrix} 0.05 \\ 0.05 \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix}
 \end{aligned}$$



perform backpropagation  
 to update weights  $\rightarrow$  from loss  $L = \frac{1}{2} \|x - s\|^2$   
 which helps us know  
 $\frac{\partial L}{\partial w_e}, \frac{\partial L}{\partial w_d}$

$x \in \mathbb{R}^2 \Rightarrow x = \begin{bmatrix} e_1 \\ e_2 \end{bmatrix}$   
 $w \in \mathbb{R}^1$   
 (context vector)  
 $w_e = 1 \times 2, w_d = 2 \times 2$   
 $= \begin{bmatrix} 0.5 & -1.0 \end{bmatrix} \quad \begin{bmatrix} 1.0 \\ 0.5 \end{bmatrix}$

~~loss = MSE~~  
 $\equiv \frac{1}{2} \|x - s\|^2$   
forward pass  
encoder:  $\hat{x} := w_e x$   
 $\begin{bmatrix} 2 \\ 2 \end{bmatrix} \begin{bmatrix} 0.5 & -1.0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$   
 $\hat{x} = \begin{bmatrix} 0.5 & -1.0 \end{bmatrix} \begin{bmatrix} 2 \\ 0 \end{bmatrix} = [1]$

decoder:  $\hat{x} = w_d h$   
 $\hat{x} = \begin{bmatrix} 1 \\ 0.5 \end{bmatrix} \begin{bmatrix} 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0.5 \end{bmatrix}$

$l = \frac{1}{2} \|x - \hat{x}\|^2$   
 $\Rightarrow l = \frac{1}{2} [(e_1 - 1)^2 + (0.5)^2] = \frac{1}{2} [1 + 0.25] = \frac{1.25}{2} = 0.625$   
 (Note:  $(e_1 - 1)^2 = 1 - 2e_1 + e_1^2$ )

$$\alpha_{t,i} = \frac{\exp(score\_fn(s_{t-1}, h_i))}{\sum_{j=1}^T \exp(score\_fn(s_{t-1}, h_j))}$$

we can come up with diff scoring fn  
(not just dot product like above)

~~expt~~

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

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

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

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

8th fn = dot product

$$\alpha_{t,i} = \frac{\exp(score\_fn(s_{t-1}, h_i))}{\sum_{j=1}^T \exp(score\_fn(s_{t-1}, h_j))}$$

$$Q = \sum_{j=1}^T \alpha_{t,i} h_j \quad (T=3)$$

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

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

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

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

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

~~$$= \exp$$~~

$$= \frac{7.389}{7.389 + 2.718 + 2.98}$$

$$= 0.587$$

$$\alpha_{t,2} = \frac{\exp(1) \left( \left( \frac{1}{2} h_1 + \frac{1}{2} \right) \sin \left( \frac{\pi}{2} \left( \frac{1}{2} h_1 + \frac{1}{2} \right) \right) \right)}{\exp(1) + \exp(1) + \exp(1)} \xrightarrow{\frac{2.718}{12.8} = 0.213}$$

$$\alpha_{t,3} = 0.12$$

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

$$= 0.5 \cdot (0,0,1) + 0.2 (0,1,1) + 0.2 (1,1,0)$$

$$= 0.5 (0.57, 0.57) + (0.2, 0.2, 0.2)$$

$$= (0.77, 0.77, 0.22)$$

*referred*

Self attention:  
How is it different from cross attention?

Steps  
① Create  $\rightarrow$  query vector ( $k_i$ )  
key vector  
value vector ( $v_i$ ) } From encoder's input vector ( $e_i$ )  
by initializing weight matrices

② calculate

dotproduct ( $e_i, k_i$ ) } similarity  
with all?

dotproduct ( $e_i, k_2$ )

→ use scaled dot product  
③ softmax on dot product

example

compute  $z_1, z_2$

~~without vec~~  $q_1 = [0.212 \quad 0.04 \quad 0.62 \quad 0.36]^T$

planning  $k_1 = [0.31 \quad 0.84 \quad 0.463 \quad 0.507]^T$

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

abide  $\Rightarrow q_2 = [0.1 \quad 0.14 \quad 0.86 \quad 0.86 \quad 0.44]^T$

$k_2 = [0.45 \quad 0.9 \quad 0.28 \quad 0.583]^T$

$v_2 = [0.31 \quad 0.36 \quad 0.19 \quad 0.32]^T$

$$\begin{aligned}
 & \text{Ansatz: } \left( \frac{Q_1 + P_1}{52} \right) \rightarrow \text{value} \rightarrow 0.01 \rightarrow z_1 \\
 & \Rightarrow \left[ \begin{array}{cccc} 0.212 & 0.04 & 0.63 & 0.26 \end{array} \right] \left[ \begin{array}{cccc} 0.3 & 0.84 & 0.83 & 0.23 \end{array} \right] \\
 & \text{①} \quad \left[ \begin{array}{cccc} 0.212 & 0.04 & 0.63 & 0.26 \end{array} \right] \left[ \begin{array}{cccc} 0.1368 & 0.1368 & 0.1368 & 0.1368 \end{array} \right] \\
 & \Rightarrow \left[ \begin{array}{cccc} 0.06572 & 0.0336 & 0.60669 & 0.26527 \end{array} \right] \\
 & \text{Ansatz: } \left( \frac{Q_1 + P_1}{52} \right) \rightarrow \text{value} \rightarrow 0.1026 \rightarrow z_2 \\
 & \Rightarrow \left[ \begin{array}{cccc} 0.286187 & 0.02286 & 0.168 & 0.1026 \end{array} \right] \left[ \begin{array}{cccc} 0.303345 & 0.222323 & 0.222323 & 0.222323 \end{array} \right] \\
 & \Rightarrow \left[ \begin{array}{cccc} 0.1001644 & 0.09557 & 0.265 & 0.265 \end{array} \right] \\
 & \text{②} \quad \left[ \begin{array}{cccc} 0.212 & 0.07 & 0.63 & 0.186 \end{array} \right] \left[ \begin{array}{cccc} 0.95 & 0.94 & 0.73 & 0.50 \end{array} \right] \\
 & \Rightarrow \left[ \begin{array}{cccc} 0.09554 & 0.0376 & 0.4599 & 0.2088 \end{array} \right] \\
 & \Rightarrow \left[ \begin{array}{cccc} 0.09557 & 0.0396 & 0.4599 & 0.2088 \end{array} \right] \\
 & \Rightarrow \left[ \begin{array}{cccc} 0.0472 & 0.0188 & 0.22995 & 0.1044 \end{array} \right] \\
 & \Rightarrow \left[ \begin{array}{cccc} 0.322325 & 0.203345 & 0.251708 & 0.22062 \end{array} \right] \\
 & \Rightarrow \left[ \begin{array}{cccc} -0.31 & 0.36 & 0.19 & 0.72 \end{array} \right] \\
 & \Rightarrow \left[ \begin{array}{cccc} 0.09892075 & 0.0733842 & 0.0478592 & 0.1588462 \end{array} \right]
 \end{aligned}$$

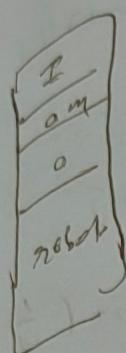
$$\begin{aligned}
 & \text{③} + \text{④} \\
 & \Rightarrow \left[ \begin{array}{cccc} 0.20008 & 0.24654876 & 0.03562142 & 0.24523008 \end{array} \right] \\
 & \text{X wrong}
 \end{aligned}$$

$$\textcircled{1} \text{ Xesta softwari } \begin{bmatrix} 0.225986 & 0.225347 & 0.300122 & 0.245521 \\ 0.225986 & 0.83 & 0.1 & 0.887 \\ [0.036] & [0.036] & [0.036] & [0.036] \\ = [0.08243856] & [0.18703801] & [0.09329908] & [0.09329908] \end{bmatrix}$$

$$\textcircled{2} \text{ ofex softwari } \begin{bmatrix} 0.225685 & 0.203681 & 0.250215 \\ 0.225685 & 0.90 & 0.90 & 0.90 \\ [0.036] & [0.036] & [0.036] & [0.036] \\ = [0.0732504] & [0.0826866] & [0.18015264] & [0.18015264] \end{bmatrix}$$

$$Z_1 = \begin{bmatrix} 0.316 & 0.6046 & 0.1424 & 0.5415 \\ 0.316 & 0.547870 & 0.4483 & 0.5415 \\ [0.036] & [0.036] & [0.036] & [0.036] \\ = [0.09010] & [0.09220] & [0.09220] & [0.09220] \end{bmatrix}$$

I am a robot  
 $P_E: d=4 \rightarrow \mathbb{R}^4$  (using base 1000)



$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

$$\begin{aligned} &\rightarrow \begin{bmatrix} P_{00} & P_{01} & P_{02} & P_{03} \end{bmatrix} \\ &\rightarrow \begin{bmatrix} P_{10} & P_{11} & P_{12} & P_{13} \end{bmatrix} \\ &\rightarrow \begin{bmatrix} P_{20} & P_{21} & P_{22} & P_{23} \end{bmatrix} \\ &\rightarrow \begin{bmatrix} P_{30} & P_{31} & P_{32} & P_{33} \end{bmatrix} \end{aligned}$$

$$\text{pos} = \sin\left(\frac{\text{pos}}{100 \text{ deg}}\right)$$

forward eva!

$$\text{pos} = \cos\left(\frac{\text{pos}}{100 \text{ deg}}\right)$$

$$\frac{d}{2} = \frac{9}{2} = \frac{2}{2}$$

$$= \sin\left(\frac{\text{pos}}{100 \text{ deg}}\right)$$

$$P_{\text{pos}, 210} = \sin(\text{pos}) \Rightarrow \sin(0) = 0$$

$$P_{\text{pos}, 210+1} = \cos\left(\frac{\text{pos}}{100 \text{ deg}}\right) = \cos(0) = 1$$

$$P_{00} = 0, \quad P_{01} = 1$$

$$\left(\frac{0}{100 \frac{2}{4}}\right) = \sin(0)$$

$$P_{\text{pos}, 211} = \sin\left(\frac{0}{100 \frac{2}{4}}\right) = \sin(0) = 0$$

$$\left(\frac{0}{100 \frac{2}{4}}\right) = \cos(0) = 1$$

$$P_{\text{pos}, 211+1} = \cos(0) = 1$$

$$P_{02} = 1$$

$$P_{02} = 0, \quad P_{03} = 1$$

$$\left(\frac{0}{100 \frac{2}{4}}\right) = \sin(0)$$

$$P_{\text{pos}, 212} = \sin\left(\frac{0}{100 \frac{2}{4}}\right) = \sin(0)$$

Ans = 1

P = 0

P<sub>0000</sub>

$$P_{1,2(0)} = 8^{20} \times \left( \frac{1}{100} \times \left( \frac{2(0)}{100} \right) \right) = \frac{1}{100}$$

$$\begin{array}{|c|c|c|c|} \hline & 2 & 0 & 1 \\ \hline 2 & 0 & 9 & 1 \\ \hline 0 & 9 & 1 & 0 & 9 \\ \hline \end{array} = \sin(0.01) = 0.0001745$$

$$\begin{array}{|c|c|c|c|} \hline & 2 & 1 & 0 \\ \hline 2 & 1 & 9 & 1 \\ \hline 0 & 9 & 1 & 0 & 9 \\ \hline \end{array} = 103(0.01)$$

$$\begin{array}{|c|c|c|c|} \hline & 2 & 1 & 1 \\ \hline 2 & 1 & 9 & 1 \\ \hline 0 & 9 & 1 & 0 & 9 \\ \hline \end{array} = \cos\left(\frac{103(0.01)}{100}\right) = 0.999$$

$$P_{1,2(0)+1} = \begin{array}{|c|c|c|c|} \hline & 2 & 1 & 2 \\ \hline 2 & 1 & 9 & 1 \\ \hline 0 & 9 & 1 & 0 & 9 \\ \hline \end{array} = \sin\left(\frac{1}{10}\right)$$

$$\begin{array}{|c|c|c|c|} \hline & 2 & 1 & 2 \\ \hline 2 & 1 & 9 & 1 \\ \hline 0 & 9 & 1 & 0 & 9 \\ \hline \end{array} = \sin(0.1) = 0.001$$

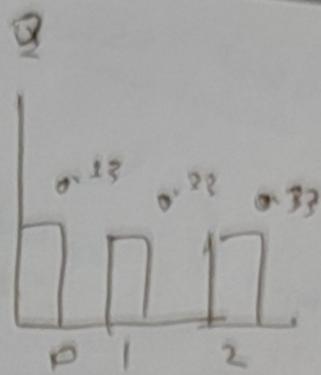
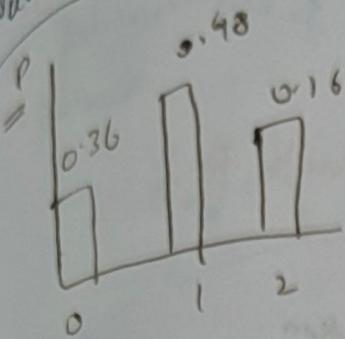
$$\begin{array}{|c|c|c|c|} \hline & 2 & 1 & 3 \\ \hline 2 & 1 & 9 & 1 \\ \hline 0 & 9 & 1 & 0 & 9 \\ \hline \end{array} = \cos(0.1) =$$

$$P_{1,2(0)+1} = \begin{array}{|c|c|c|c|} \hline & 2 & 0 & 9 \\ \hline 2 & 1 & 9 & 1 \\ \hline 0 & 9 & 1 & 0 & 9 \\ \hline \end{array} = 0.0001745$$

R

$$\text{Exponent } (0)^{\text{index}} = (0.0001745)^{\text{index}} =$$

Numerical



Binomial dist

distribution

$$P(x)$$

$$(x=0, 1, 2)$$

$$Q(x)$$

Q. compute  $D_{KL}(P||Q)$

$$D_{KL}(P||Q) = \sum_{x \in \{0, 1, 2\}} P(x) \ln \frac{P(x)}{Q(x)}$$

$$\begin{aligned} & \approx \sum_{x=0}^2 P(x) \ln \left( \frac{P(x)}{Q(x)} \right) \\ & \approx 0.36 \ln \left( \frac{0.36}{0.13} \right) + 0.48 \ln \left( \frac{0.48}{0.22} \right) + 0.16 \ln \left( \frac{0.16}{0.33} \right) \\ & \approx 0.0852986 \end{aligned}$$

$$D_{KL}(Q||P) = \sum_{x \in \{0, 1, 2\}} Q(x) \ln \frac{Q(x)}{P(x)}$$

$$= 0.13 \ln \left( \frac{0.13}{0.36} \right) + 0.22 \ln \left( \frac{0.22}{0.48} \right) + 0.33 \ln \left( \frac{0.33}{0.16} \right)$$

$$= 0.13(-0.087) + 0.22(-0.374) + 0.33(0.724)$$

$$= -0.0261 - 0.1122 - 0.2172 = 0.09245$$

Note:  $D_{KL}(P||Q) \neq D_{KL}(Q||P)$