

Practicals

Lab 1

derivative

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

$$\textcircled{4} \quad \text{For } \tanh = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

more like between 0 & 1

approx zero

$$\tanh(x) = 2 \times \text{sigmoid}(2x) - 1$$

$$\text{now with } \frac{d}{dx} \left(\frac{f(x)}{g(x)} \right) = \frac{f'(x)g(x) - f(x)g'(x)}{[g(x)]^2}$$

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

$$(e^x + e^{-x})^2$$

$$\text{ReLU}(x) = x^+ = \max(0, x) = x + |x| = \begin{cases} x & \text{if } x > 0, \\ 0 & \text{if } x \leq 0 \end{cases}$$

$$\textcircled{4} \quad e^x + e^{-x} \frac{d(e^x - e^{-x})}{dx} - e^x - e^{-x} \frac{d(e^x + e^{-x})}{dx}$$

now we can calculate derivative

$$(e^x + e^{-x})^2$$

$$\Rightarrow (e^x + e^{-x})^2 - (e^{-x} - e^x)^2 \quad \text{Leaky ReLU}$$

$$\text{this is good at } (e^x + e^{-x})^2$$

smoothness

$$\text{smooth function} \Rightarrow 1 - \frac{(e^x - e^{-x})^2}{(e^x + e^{-x})^2}$$

but not

$$\text{smooth function} \Rightarrow 1 - \frac{[e^x + e^{-x}]^2}{[e^x - e^{-x}]^2}$$

$$\text{Leaky ReLU}(x) = \frac{(1+\alpha)x + (1-\alpha)x}{2}$$

gradient example

mathematical example
numerical example

$$\text{Softmax} \quad \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \Rightarrow \begin{bmatrix} s_1 \\ s_2 \\ s_3 \end{bmatrix}$$

$$\alpha = \begin{bmatrix} 2.0 \\ 1.0 \\ 0.1 \end{bmatrix} \Rightarrow s = \begin{bmatrix} s_1 \\ s_2 \\ s_3 \end{bmatrix} \quad s_i = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

- 1) Compute Softmax output
- 2) Compute the Jacobian

$$s_1 = \frac{e^{x_1}}{\sum_{j=1}^3 (e^{x_1} + e^{x_2} + e^{x_3})}$$

$$s_1 = \frac{e^2}{(e^2 + e^1 + e^0)} = \frac{7.3}{10.0}$$

$$s_2 = \frac{e^{x_2}}{\sum_{j=1}^3}$$

$$s_2 = \frac{e^1}{10.0} = \frac{2.7}{10.0} = 0.24$$

$$= \frac{7.3}{10.0}$$

$$s_3 = \frac{e^{x_3}}{\sum_{j=1}^3} = \frac{e^0}{10.0} = \frac{1.0}{10.0} = 0.099$$

$$= \frac{7.3}{10.0} + \frac{2.7}{10.0} + \frac{1.0}{10.0} = 1.0$$

$$\frac{ds_i}{dx_i} = s_i(1-s_i)$$

$$(s_1, s_2, s_3) = (0.66, 0.24, 0.1)$$

$$S = \begin{bmatrix} s_1 \\ s_2 \\ s_3 \end{bmatrix} = \begin{bmatrix} 0.66 \\ 0.24 \\ 0.1 \end{bmatrix}$$

$$J_{ij} = \frac{\partial s_i}{\partial x_j}$$

$$J = \begin{bmatrix} \frac{\partial s_1}{\partial x_1} & \frac{\partial s_1}{\partial x_2} \\ \frac{\partial s_2}{\partial x_1} & \frac{\partial s_2}{\partial x_2} \\ \frac{\partial s_3}{\partial x_1} & \frac{\partial s_3}{\partial x_2} \end{bmatrix}$$

Case 1: $i=j$ (Diagonal element)

$$\frac{\partial s_i}{\partial x_i} = s_i(1-s_i) = 0.66$$

$$\text{Case 2: } \frac{\partial s_i}{\partial x_j} = -s_i s_j$$

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

$$= 0.66(1-0.66)$$

$$= 0.22$$

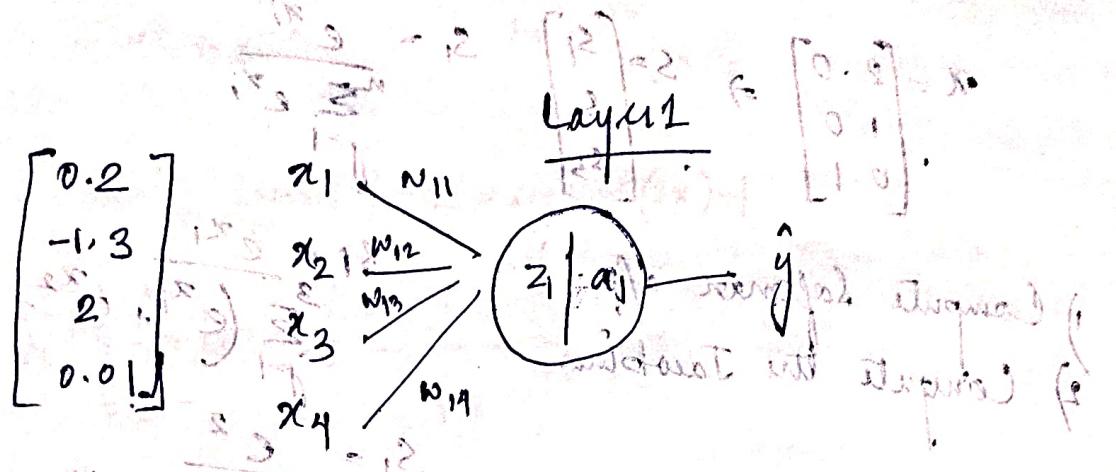
$$J = \begin{bmatrix} 0.66 & -0.66 & -0.66 \\ -0.66 & 0.66 & -0.66 \\ -0.66 & -0.66 & 0.66 \end{bmatrix}$$

$$\frac{\partial s_2}{\partial x_1} = -s_1 s_2$$

$$= -0.66 \times 0.24 = -0.15$$

2.C. Saturates in the case of Sigmoid

Lab 2



$$w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + w_{14}x_4 - n_{11} = 0.001$$

$$\text{P.S.O.} = \frac{1}{100} = 0.01$$

$$z_1 = (0.001 \times 0.2) + (0.01 \times -1.3) + (-0.005 \times 2) + (-1.2 \times 0.01) = -0.005$$

$$= 0.0002 + (-0.013) + (-0.01)$$

$$2 + (-0.012) \quad (2-1)2 = \frac{2}{2}$$

~~-0.021~~ -0.0348

$$a_1 = \text{softmax}(z)$$

$$z = -0.024$$

$$\begin{aligned}
 & \text{①) } e^{0.2} \\
 & \frac{e^{-1.3}}{e^{0.2}} + e^{0.2} + e^2 + e^{0.01} \\
 & = \frac{(e^{2-1})^2}{e^{1.22}} = \frac{2^2}{e^{1.22}}
 \end{aligned}$$

$$\left\{ \begin{array}{l} 850.0 - \frac{e}{e^{-0.0348}} \\ 20.0 - \frac{21.0}{21.0} \\ 10.0 \end{array} \right.$$

$$(12-1)12 = 12b$$

$$(11.0-1)11.0 = 11b$$

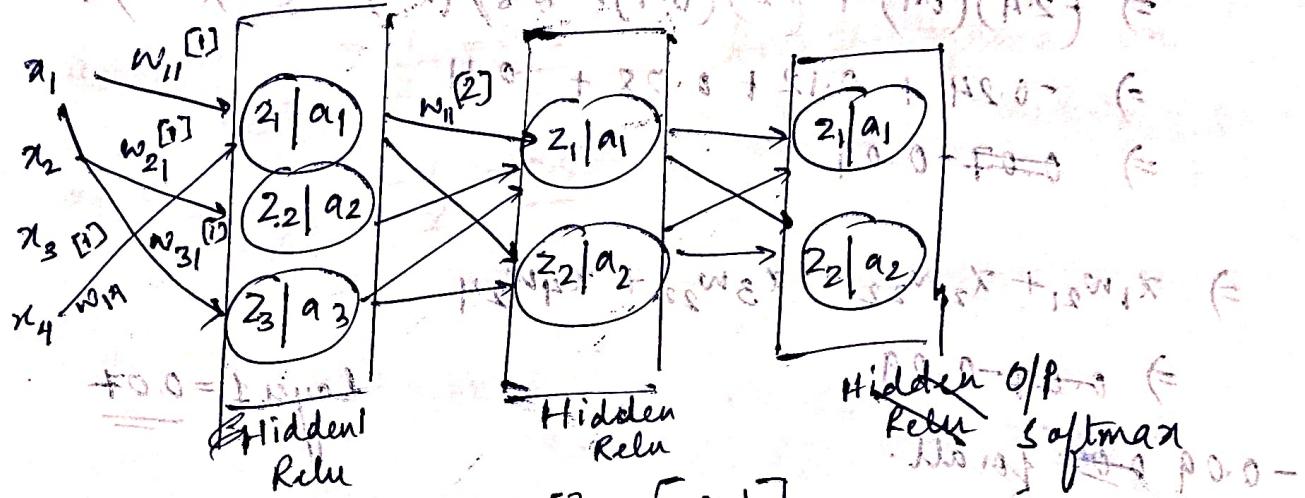
$$\frac{d^2 P}{dt^2} = \frac{d^2 p}{dt^2}$$

$$\begin{bmatrix} 0.2 \\ -1.3 \\ 1.2 \\ 0.01 \end{bmatrix} + \begin{bmatrix} 0.001 \\ 0.01 \\ -0.005 \\ -1.2 \end{bmatrix} = \begin{bmatrix} 0.2 \\ -1.3 \\ 1.2 \\ 0.01 \end{bmatrix} + e^{W_1 x} + e^{W_2 x} + e^{W_3 x}$$

$$= 0.002 + 0 - 0.013 + 0.01 - 0.012$$

$$= \cancel{-0.024} - 0.0348.$$

$$f((1.0)(1.1) + (1.0)(2.3) + (1.0)(8.1) + (1.0)(1.0))$$



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

$$w_{11}^{[1]} = \begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \end{bmatrix}$$

$$x_1 = -2.4 + 2.4 = 0$$

$$x_2 = \frac{1.2 + 1.2}{2} = \frac{2.4}{2} = 1.2$$

$$x_3 = \frac{-0.8 + 0.8}{2} = 0$$

$$x_4 = \frac{1.1 + 1.1}{2} = \frac{2.2}{2} = 1.1$$

$$z_1 = w_{11} x_1 + x_1 w_{21} + x_1 w_{31} + x_4 w_{41}$$

$$= (0.1)(-2.4) + (-2.4)(0.1) + (-2.4)(0.1) + (1.1)(0.1)$$

$$x_1 w_{11} + x_1 w_{21} + x_1 w_{31} + x_2 \overbrace{w_{12}}^{100.0} + x_2 \overbrace{w_{22}}^{100.0} + x_3 \overbrace{w_{32}}^{100.0} \\ + x_3 \overbrace{w_{13}}^{100.0} + x_3 \overbrace{w_{23}}^{100.0} + x_3 \overbrace{w_{33}}^{100.0} + x_4 \overbrace{w_{14}}^{100.0} + x_4 \overbrace{w_{24}}^{100.0} \\ + x_4 \overbrace{w_{34}}^{100.0}$$

$$\Rightarrow x_1 w_{11} + x_2 w_{12} + x_3 w_{13} + x_4 w_{14}.$$

$$\Rightarrow (-2.4)(0.1) + (-1.2)(0.1) \neq (-0.8)(0.1) + (1.1)(0.1)$$

$$\Rightarrow -0.24 + 0.12 + 0.08 + 0.11$$

$$\Rightarrow \text{O} = 0.06$$

$$\Rightarrow x_1 w_{21} + x_2 w_{22} + x_3 w_{23} + x_4 w_{24}$$

$$\Rightarrow \cancel{0.07} - 0.09$$

-0.09 ~~0.07~~ for all

$$\text{Layer 1} = \underline{\underline{0.07}}$$

$$a_1 = \text{ReLU}(0.09)$$

$$= 0 \quad \left(\frac{x+1}{x} \right)$$

$$\begin{aligned} & \left. \begin{array}{l} 1.0 \\ -0.09 \\ \hline 0.91 \end{array} \right\} \quad \left. \begin{array}{l} 0.09 \\ +0.09 \\ \hline 0.18 \end{array} \right\} \\ & = \cancel{e^{-0.09}} \quad \cancel{e^{0.09}} \\ & \frac{e^{0.09} + e^{-0.09}}{2} \\ & = \frac{1.09 + 0.91}{2} \\ & = \underline{\underline{1.00}} \end{aligned}$$

Layer 1

$$z_1 = 0 - 0.09$$

$$a_1 = 0$$

$$Z_2 = -0.0$$

α, γ

$$Z_3 = -0.9$$

$$\alpha_3 = 0$$

$$(100)(1.8) + (100)(1.8)(1.8) + (100)(1.8)(1.8)(1.8) = 18 \text{ m}^3$$

Layer 2

0 for all z_1 and z_2 values

New α values are ~~0~~ $a_1, a_2 \& a_3 = 0$ is violated

Hidden value = 0.

Number of iterations = 1000

Layer 3

$$a_1 = \frac{e^0}{e^0 + e^0} = \frac{1}{2} = 0.5$$

doing the softmax on

a_1, a_2

$$a_2 = \frac{e^0}{e^0 + e^0} = \frac{1}{2} = 0.5$$

Iteration 1000 is completed

Structure of network after iteration 1000 is laid

The minimum loss value is 11

rate of steps

(Coordinate, loss(1)), stepsize (stepsize), w

no step size settings for perpendicular and distance



Rate of loss

initial weight & initial loss value are 0 steps

number of step

with overshoot for

loss, x, p, w

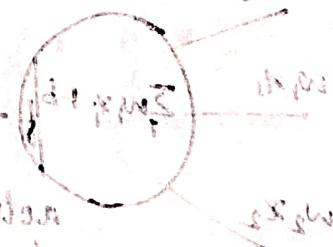
initial loss & final loss

loss value

constant

stepsize

f(x) = (x, p, w)



negative

antiderivative

loss

loss value is added to initial loss

and loss is the antiderivative benefit gained

steps

$$\int f(x) dx = F(x)$$

$$(f-1) \text{ fat } (f-1)$$

Representation power of NN

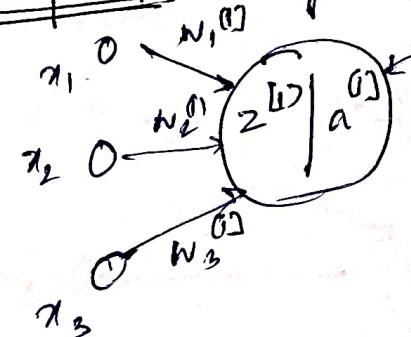
- NN with at least one hidden layer are universal approximators
- any continuous fn, $f(x)$, and some $\epsilon > 0$, there exists a neural network, $g(x)$, with at least one hidden layer (with some non-linearity g), such that

$$\forall x, |f(x) - g(x)| \leq \epsilon$$

both functions mimic hence difference is zero.

Neural Network are non-convex

Perception \rightarrow single layer NN.



$$z^{(1)} = w_1 x_1 + w_2 x_2 + w_3 x_3 + b_1$$

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

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

AND

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



$$\text{e.g. } w_1 x_1 + w_2 x_2 + (-1)$$

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

$$= 0 \quad g(-1) = \frac{e^{-\frac{-1}{2}}}{e^{\frac{-1}{2}}} = e^{-2}$$

$$w_{11} = 1$$

$$w_{12} = 1$$

for $x_1=0, x_2=1$
 ~~$\frac{w_1}{w_2}$~~ $0 + 1 - 1 = 0 \quad g(z)$
 however $x_1=0$ if bias is $b=1$ so 1
 for $x_1=1, x_2=0$
 $= 1 + 0 - 1 = 0 \quad g(z)$
 and $0 < 1$ if bias is 0 , $1 + 0 + 0$
 we want to say, (x) , $g(z) = 1$
 for $x_1=1, x_2=1$
 $\Rightarrow 0 + 1 - 1 \Rightarrow 2 - 1$
 $\Rightarrow 1 \quad g(z) = 1 \neq 0$
 so $g(z) = 1 \neq 0$
 one is non-affine and other is affine

AND OR

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

XOR

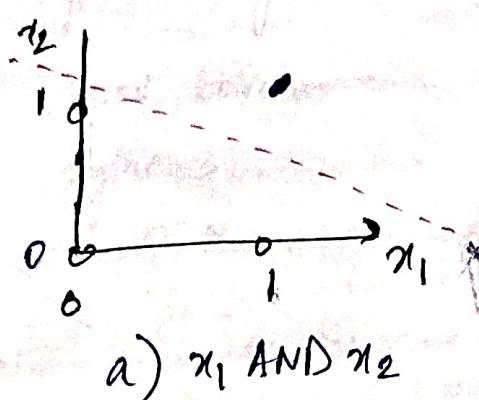
x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0

- 1) $0 + 0 + 0 = 0 \quad g(z) = 0$
- 2) $0 + 1 + 0 = 1 \quad g(z) = 1$
- 3) $1 + 0 + 0 = 1 \quad g(z) = 1$
- 4) $1 + 1 + 0 = 2 \quad g(z) = 1$

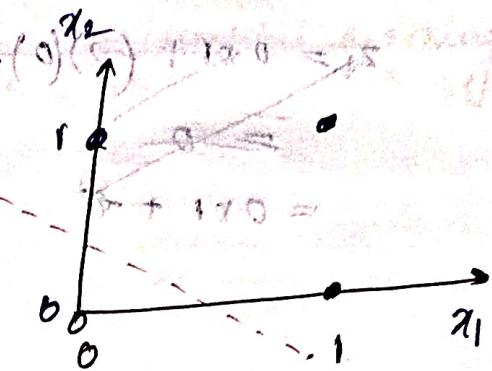
$1 + 1 + 0 = 2$

$1 + 1 + 0 = 2$

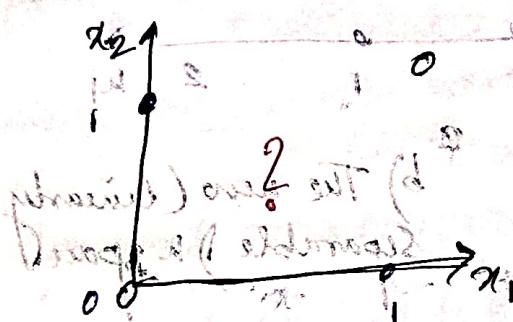
Decision Boundaries



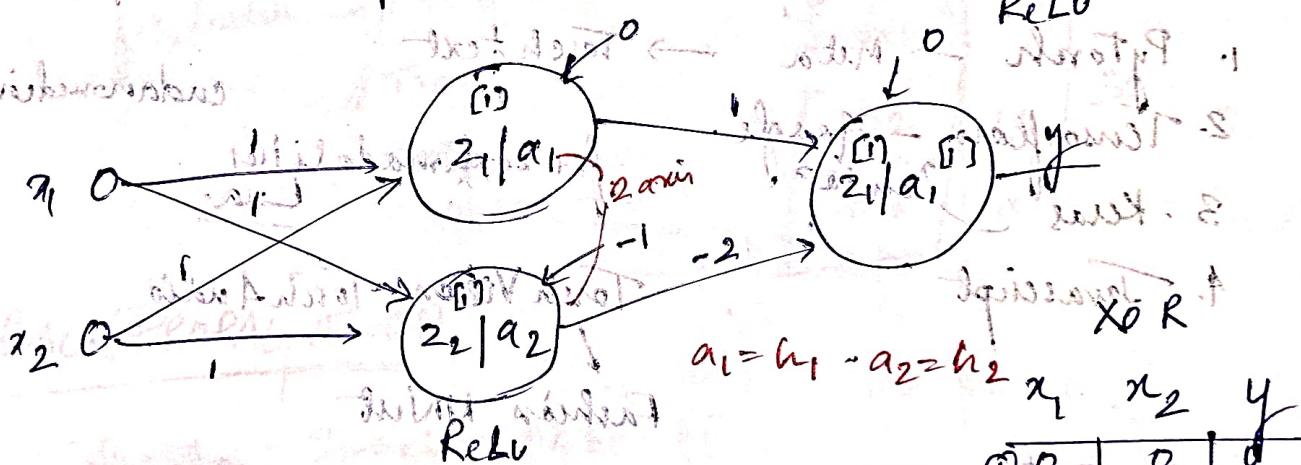
a) $x_1 \text{ AND } x_2$



b) $x_1 \text{ OR } x_2$



c) $x_1 \text{ XOR } x_2$



$$z_1 = w_{11}x_1 + w_{12}x_2 + b_1 = 1 \cdot 0 + (-1) \cdot 0 + 0 = 0$$

$$= 0 \quad \text{and } g(0) = 0$$

$$z_2 = 0 + 0 + 0 = 0 \quad g(0) = 0$$

$$\text{total } z = 0 \quad g(z) = 0$$

$$z_1 = w_{11}x_1 + w_{12}x_2 + b_1 = 1 \cdot 1 + (-1) \cdot 0 + 0 = 1$$

$$= 1 \quad g(1) = 1$$

$$z_2 = 0 + 0 + 0 = 0 \quad g(0) = 0$$

$$\text{total } z = 1 \quad g(z) = 1$$

$$z_1 = 0 + 0 - 1 = -1 \quad g(-1) = 0$$

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

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

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

$$\text{total } z = 4 \quad g(4) = 1 \quad a_1 = 1$$

$$\text{total } z = 1 - 1 = 0 \quad g(0) = 0 \quad a_2 = 0$$

Computing the equations through several layers has computational cost $\Theta(n)$. Random initialization is done at 0 will for nets \rightarrow lead to zero in the equations.

Computational Graph

net initialization

i) $w_{ij} \sim \text{uniform } [-0.1, 0.1]$

$$\text{let } f = (x + y)^2$$

multivariable

func.

x has the impact on q .

but q has negative impact on f .

$$x = -2 \\ y = -4$$

$$z = -4 \\ q = 2$$

$$f = q^2$$

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

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

$$\frac{\partial f}{\partial q} = 2$$

forward pass output

$$\left[\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z} \right] = [2, 2, 2]$$

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

$$= 2 \cdot 1 \\ = 2$$

$$= 2 \Rightarrow -4$$

global gradient.

local gradient

output q respect to x — local gradient

then global gradient \rightarrow value that came to me

\rightarrow uniform distribution

$$ii) w_{ij} \sim N(0, \frac{2}{n(p) + n(q)})$$

faster convergence

- start from back
- put the gradients
- compute the local & global gradients
- multiply them.

① $\frac{\partial f}{\partial z} \rightarrow$ local grad \times global grad

~~20/07/25~~ $f(x, y, z) = x + yz$

$$x = -2, y = 5, z = -4$$

1) Identify additional functions
1.1 Compute local derivatives.

2) Draw computational graph

3) Perform forward pass

4) Perform backward pass

starting from the end of the circuit.

$$f = x + yz$$

let

$$y = \frac{5}{-4}$$

$$x = -2$$

$$z = -4$$

$$f = -22$$

$$\text{let } q = yz, f = q + x \rightarrow \frac{\partial f}{\partial q} = 1, \frac{\partial f}{\partial x} = 1$$

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

$$\Rightarrow -4 \rightarrow 5$$

~~$\frac{\partial q}{\partial x}$~~

② $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

$$\Rightarrow \cancel{\frac{\partial f}{\partial x}} = \cancel{\frac{\partial f}{\partial q}} \cdot \cancel{\frac{\partial q}{\partial x}} \Rightarrow 0$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \cdot \frac{\partial q}{\partial y} \Rightarrow (1) \cdot (-4) \equiv -4$$

$$\frac{\partial f}{\partial z} = \frac{\partial f}{\partial q} \cdot \frac{\partial q}{\partial z} \Rightarrow (1) \cdot 5 \equiv 5$$

$$f(w, x) = \frac{1}{c_0 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

$$w_0 = 2$$

$$w_1 = -1$$

$$w_2 = 3$$

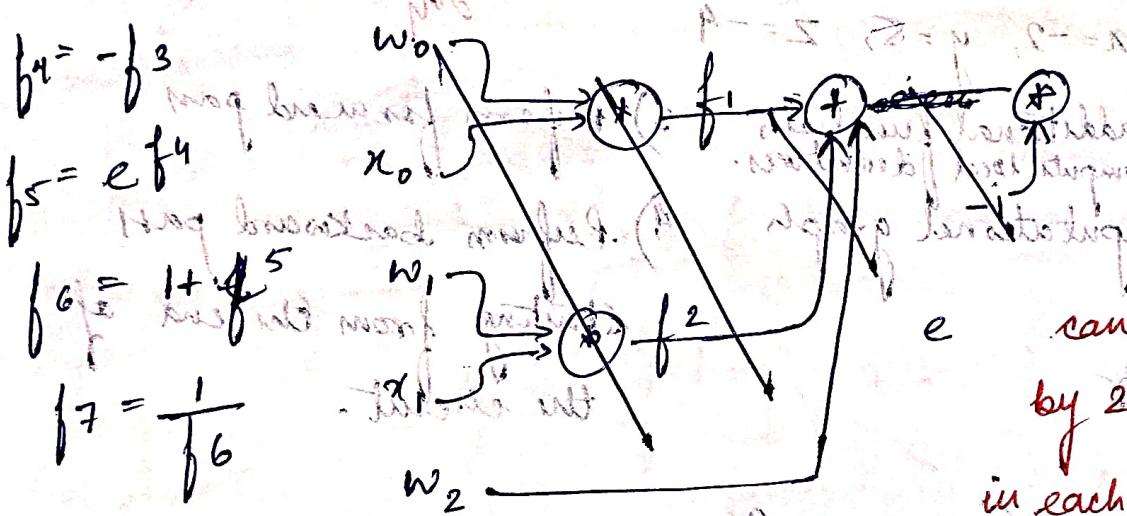
$$x_1 = -2$$

$$w_0 = -3$$

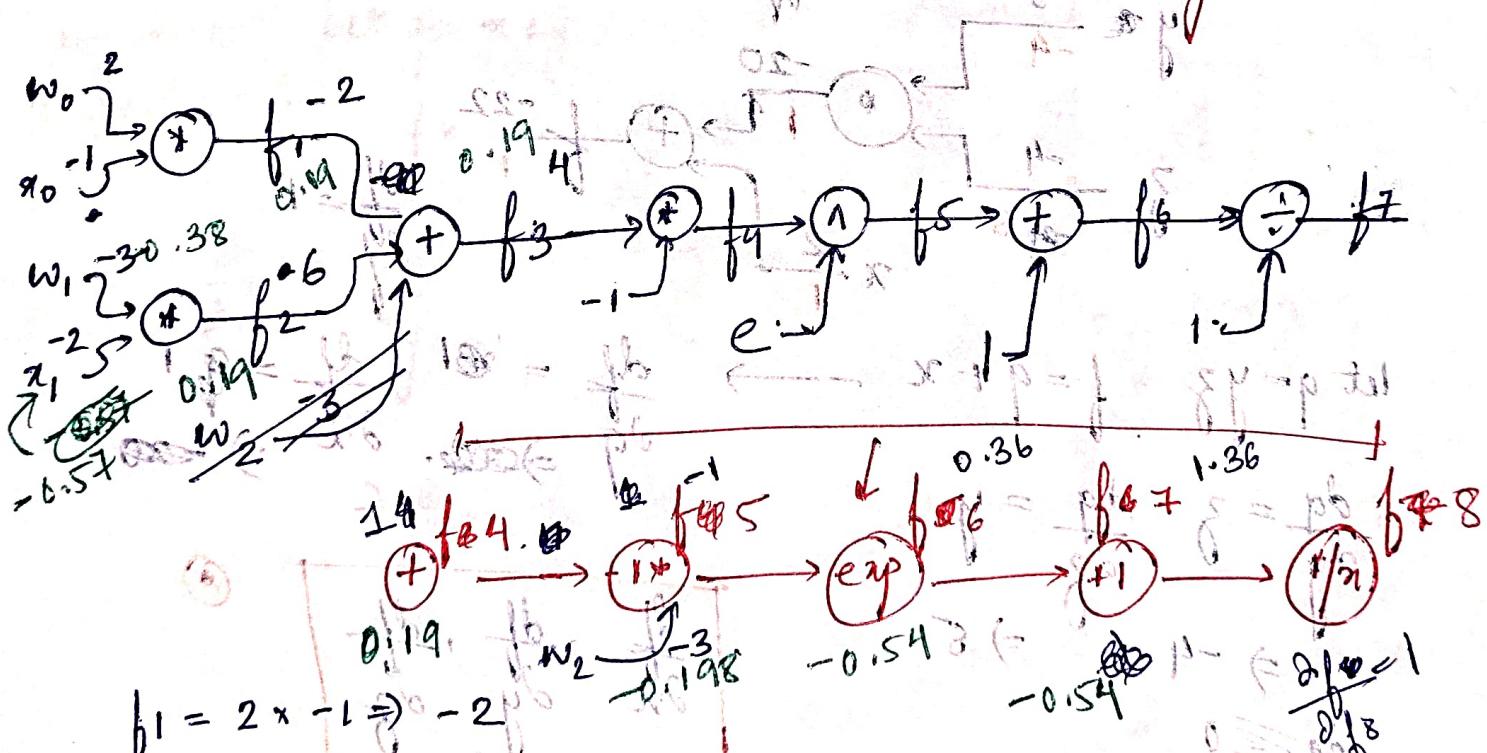
$$\text{let } q = w_0 x_0 + w_1 x_1 + w_2$$

$$f = \frac{1}{1+e^{-q}}$$

$$\text{let } f_1 = w_0 x_0, f_2 = w_1 x_1, f_3 = \frac{1}{1+e^{-q}} f_1 + f_2 + w_2$$



e can be done
by 2 inputs
in each gate



$$(2) \quad f_8 = p_5 f_7 = \frac{1}{2.20} \approx 0.45$$

$$f_3 = 4 \quad f_4 = w_2 + f_3 \quad f_5 = -1$$

$$f_6 = e^{f_5} = 0.36 \quad f_7 = 1 + 0.36 = 1.36 \quad f_8 = \frac{1}{1.36(f_7)} = 0.735$$

$$f_1 = w_0 x_0 \rightarrow f_2 = w_1 x_1 \quad f_3 = f_1 + f_2 \quad f_4 = f_3 + w_2$$

$$f_5 = -f_4, f_6 = e^{f_5} \quad f_7 = 1 + f_6, f_8 = \frac{1}{f_7}$$

$$f_1 \Rightarrow \frac{\partial f_1}{\partial x_0} \Rightarrow w_0 \quad \frac{\partial f_1}{\partial w_0} \Rightarrow -1 \quad f_2 \Rightarrow \frac{\partial f_2}{\partial x_1} = x_1 \rightarrow -2$$

$$\frac{\partial f_2}{\partial w_1} = w_1 \rightarrow -3$$

$$f_3 = \frac{\partial f_3}{\partial f_1} = 1 \quad f_4 = \frac{\partial f_4}{\partial f_3} = 1 \quad f_5 \Rightarrow \frac{\partial f_5}{\partial f_4} = -1$$

$$f_6 = \frac{\partial f_6}{\partial e^{f_5}} \Rightarrow e^{f_5} \quad f_7 = \frac{\partial f_7}{\partial f_6} = 1$$

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

~~$\frac{\partial f_8}{\partial f_7} = 1.36$~~ local \times global

~~$\frac{\partial f_8}{\partial f_8} = \frac{df}{df} = 1$~~ Global

$$\frac{\partial f_8}{\partial f_7} = \text{local} \times \text{global} \quad \frac{\partial f_8}{\partial f_7} = \frac{-1}{f_7^2} \times 1 \Rightarrow \left(\frac{-1}{1.36}\right)^2 \times 1 \Rightarrow -0.54$$

$$\frac{\partial f_8}{\partial f_6} = \frac{\partial f_7}{\partial f_6} \times \frac{\partial f_8}{\partial f_7}$$

$$= 1 \times -0.54$$

$$= -0.54$$

$$= -0.54 \times 0.36$$

$$= -0.198$$

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

$$\frac{\partial f_8}{\partial f_4} \Rightarrow \frac{\partial f_5}{\partial f_4} \times \frac{\partial f_8}{\partial f_5}$$

$$\Rightarrow -1 \times -0.198$$

$$= 0.19$$

$$= 1 \times 0.19$$

$$= 0.19$$

$$\cancel{\frac{\partial f_8}{\partial f_1}} = \cancel{\frac{\partial f_2}{\partial f_1}} \times \cancel{\frac{\partial f_8}{\partial f_2}}$$

$$\cancel{\frac{\partial f_8}{\partial f_2}} = \cancel{\frac{\partial f_3}{\partial f_2}} \times \cancel{\frac{\partial f_8}{\partial f_3}}$$

$$\cancel{\frac{\partial f_8}{\partial f_2}} = \cancel{\frac{\partial f_3}{\partial f_2}}$$

$$= 0.19$$

$$= 0.19 \times \cancel{\frac{\partial f_3}{\partial f_3}}$$

$$\cancel{\frac{\partial f_8}{\partial f_4}} = \cancel{\frac{\partial f_5}{\partial f_4}} \times \cancel{\frac{\partial f_8}{\partial f_5}}$$

$$= 0.19.$$

$$\frac{\partial f_8}{\partial f_3} = \frac{\partial f_4}{\partial f_3} \times \frac{\partial f_8}{\partial f_4}$$

$$= 0.19$$

$$\frac{\partial f_8}{\partial w_1} = \frac{\partial f_2}{\partial w_1} \frac{\partial f_8}{\partial f_2}$$

$$= -2 \times 0.19$$

$$= -0.38$$

~~1 x 0.19~~

$$\frac{\partial f_8}{\partial w_2} = \frac{\partial f_2}{\partial w_2} \frac{\partial f_8}{\partial f_2}$$

$$= 1 \times 0.19$$

$$= 0.19$$

$$\frac{\partial f_8}{\partial x_1} = \frac{\partial f_2}{\partial x_1} \cdot \frac{\partial f_8}{\partial f_2}$$

$$= w_1 \cdot 0.19$$

$$= -3 \cdot 0.19$$

$$\frac{\partial f_8}{\partial w_0} = \frac{\partial f_3}{\partial f_1} \cdot \frac{\partial f_8}{\partial f_3}$$

$$= 1 \cdot 0.19$$

$$\frac{\partial f_8}{\partial w_0} = \frac{\partial f_1}{\partial w_0} \cdot \frac{\partial f_8}{\partial f_1}$$

$$= -1 \cdot 0.19$$

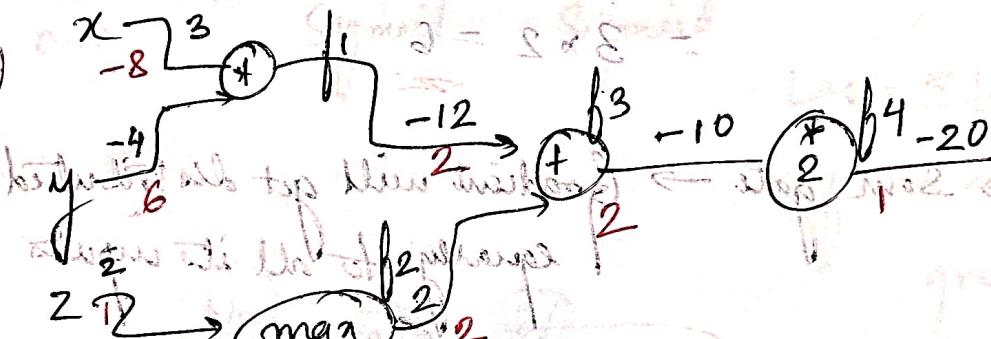
$$= -0.19$$

$$\frac{\partial f_8}{\partial x_0} = \frac{\partial f_1}{\partial x_0} \cdot \frac{\partial f_8}{\partial f_1}$$

$$= 2 \cdot 0.19$$

$$= 0.38$$

②



equations

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

$$\frac{\partial f_1}{\partial w} = -4$$

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

$$\frac{\partial f_2}{\partial z} = 1 \rightarrow \text{if } z > w$$

$$\frac{\partial f_2}{\partial w} = 0 \quad \text{if } z \leq w$$

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

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

undefined if $z = w$

$$\frac{\partial f_2}{\partial w} = 1 \quad \text{if } w > z$$

$$0 \quad \text{if } w \leq z$$

$$\text{undefined if } w = z$$

$$\frac{\partial f_4}{\partial f_3} = \frac{\partial f_4}{\partial x^1} \cdot \frac{\partial x^1}{\partial f_3}$$

$$\frac{\partial f_4}{\partial f_1} = \frac{\partial f_3}{\partial f_1} \cdot \frac{\partial f_4}{\partial f_3}$$

$$= 1 \times 2$$

$$\frac{\partial f_4}{\partial f_2} = \frac{\partial f_3}{\partial f_2} \cdot \frac{\partial f_4}{\partial f_3}$$

$$= 1 \times 2$$

$$P_{f_4} \times 2 = 2$$

$$\frac{\partial f_4}{\partial x} = \frac{\partial f_1}{\partial x} \cdot \frac{\partial f_4}{\partial f_1}$$

$$= -4 \cdot 2$$

$$= P_{f_4} \times 3$$

$$\frac{\partial f_4}{\partial y} = \frac{\partial f_1}{\partial y} \cdot \frac{\partial f_4}{\partial f_1}$$

$$= 3 \times 2 = 6$$

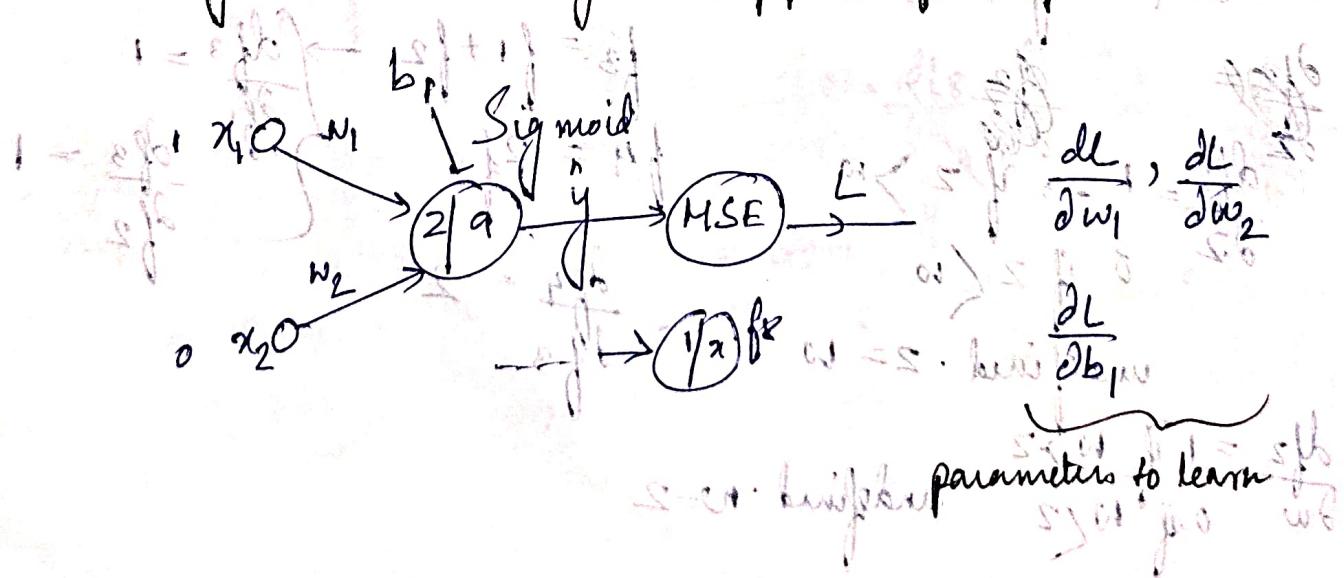
Observation:

The gradient \rightarrow Sum gate \rightarrow Gradient will get distributed equally to all its inputs

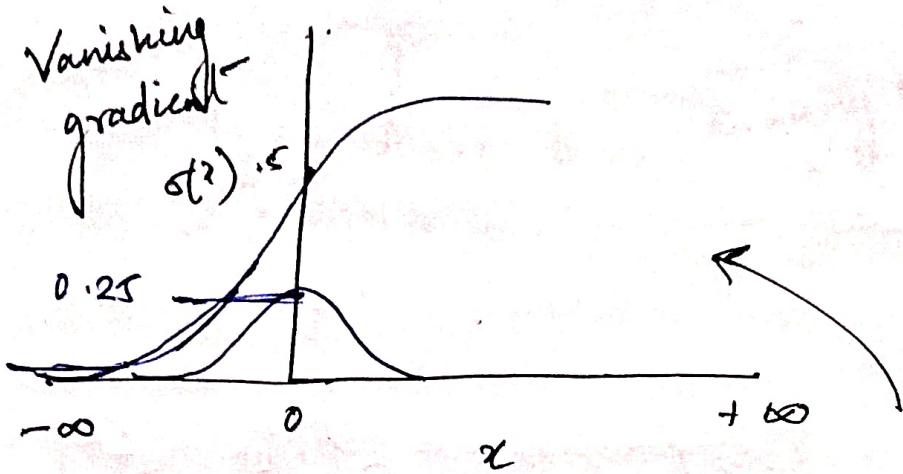
Max gate \rightarrow

Impact of w on f_2 is zero. \rightarrow gradient will flow where activation values are higher (Routes the gradient to the higher values)

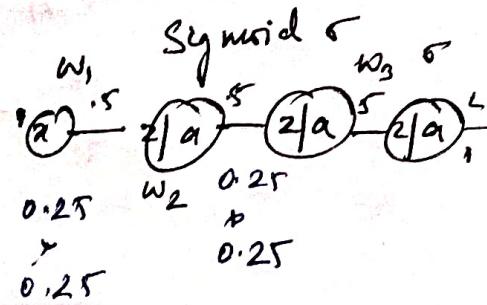
* gate \rightarrow values gets swapped. (of the input values)



Vanishing and exploding Gradient Problems



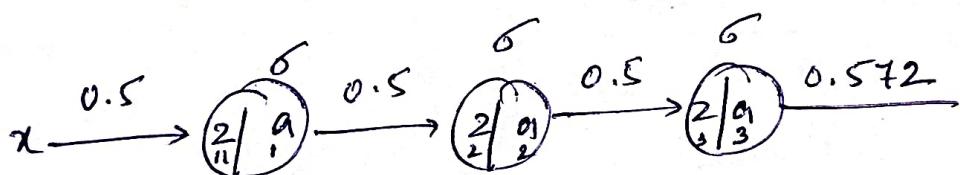
No update as gradients will be non-existent



gradient is vanished as more decimal number will get multiplied - and no learning will be there.

ReLU can be solved by \rightarrow Clipping the gradients \rightarrow threshold.

ReLU \rightarrow Exploding gradient problem.
can happen potentially to ReLU



$$z_1 = .1 \times 0.5 + 0.5$$

$$a_1 = \sigma(z_1) = 0.25 \approx 0.622$$

$$\frac{1}{1+e^{-z}} = \frac{1}{1+e^{-0.5}} = \frac{1}{1+0.606} = \frac{1}{1.606}$$

$$L = \frac{1}{2}(a_3 - y)^2, \quad \frac{dL}{da_3} = a_3 - y$$

$$\frac{dL}{da_2} = 0.572$$

$$\frac{dL}{da_3} = \frac{1}{2} \times 2(a_3 - y) = a_3 - y$$

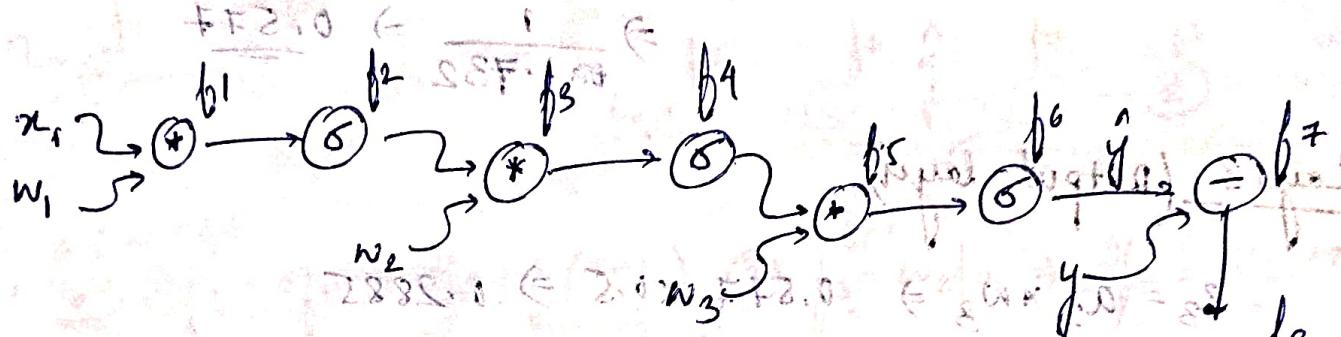
$$a_1 = 0.622$$

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

$$a_3 - y = 0.572 - 1 \\ = -0.428$$

$$\frac{\partial L}{\partial a_3} = -0.572$$

$$\frac{\partial L}{\partial a_3} = -0.428$$



$$f_1 = \sigma(x_1 w_1) \\ = 1 \times 0.5 \\ = 0.5$$

$$f_2 = \sigma(0.5) \\ = \frac{1}{1+e^{-2}} \\ = 0.6$$

$$f_3 = 0.622 \times 0.5 (w_2) \\ = 0.311$$

$$f_4 = \sigma(0.311) \\ = \frac{1}{1+e^{-0.311}} \\ = 0.577$$

$$f_5 = \sigma(0.577 \times 0.5) \\ = 0.577$$

Step 1 Forward Pass

$$f_6 = \sigma(f_5) \\ = \frac{1}{1+e^{-0.577}}$$

$$= 0.574$$

$$z_1 = x_1 w_1$$

$$\Rightarrow 1 \times 0.5 \Rightarrow 0.5$$

$$8.57 \times 0.5 \times 0.5$$

Layer 1

$$a_1 = \sigma(z_1) \\ = \frac{1}{1+e^{-2}} \\ = \frac{1}{1+0.606} \\ = \frac{1}{1.606}$$

For Vanishing/Use Symmetric

Sigmoid

Layer 2

$$z_2 = a_1 \times w_2 \Rightarrow 0.622 \times 0.5 \Rightarrow 0.311$$

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

Layer 3 (Output layer)

$$z_3 = a_2 \times w_3 \Rightarrow 0.577 \times 0.5 \Rightarrow 0.2885$$

$$a_3 = \sigma(z_3) \Rightarrow \frac{1}{1+e^{-z_3}} \Rightarrow \frac{1}{1+e^{-0.2885}} \Rightarrow 0.572$$

Loss: $L = \frac{1}{2} (a_3 - y)^2$ need to check
 $= \frac{1}{2} (0.572 - 1)^2 \Rightarrow 0.092$

Gradient Calculation $\rightarrow \frac{dl}{da_3} = a_3 - y \Rightarrow 0.572 - 1 \Rightarrow -0.428$

Output layer:

$$\frac{dl}{dz_3} = \frac{dl}{da_3} \cdot \sigma'(z_3) = \frac{dl}{da_3} \cdot a_3(1-a_3)$$

~~$$l \cdot z_3(1-z_3) \Rightarrow 0.2885(1-0.2885) \Rightarrow -0.428 \cdot 0.2885$$~~

~~$$\frac{dl}{dz_3} = (-0.428)(0.2885) \Rightarrow -0.09$$~~

~~$$a_3(1-a_3) = 0.572(1-0.572) \Rightarrow 0.572 \times 0.428$$~~

~~$$\Rightarrow 0.244$$~~

~~$$\frac{dl}{dz_3} = 0.244 \times -0.428 \Rightarrow -0.104$$~~

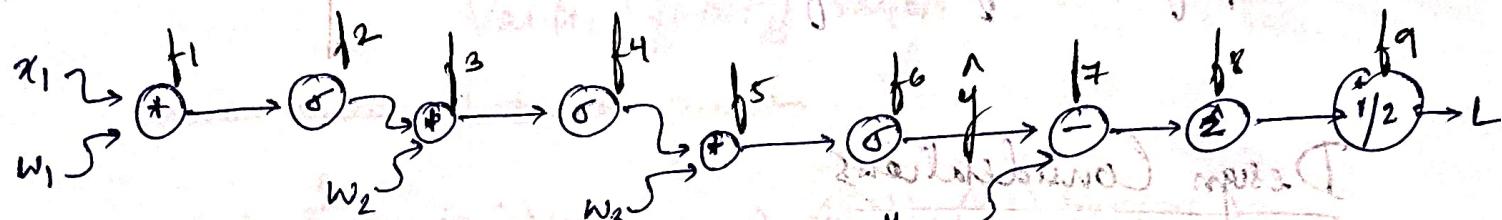
~~2nd layers~~

$$\frac{dl}{da_2} = \frac{\partial l}{\partial a_2}$$

Use computational graph.

$$\frac{\partial f_2}{\partial q} = 1 \quad \frac{\partial f_2}{\partial q} = \frac{\partial l}{\partial f_2}$$

$$\frac{\partial f_2}{\partial l} = \frac{\partial l}{\partial l} = 1 \times \frac{\partial f_2}{\partial l}$$



- ① $f_1 = \sigma(x_1)$ ② $f_2 = \sigma(f_1) \Rightarrow \frac{\partial f_2}{\partial f_1} = \sigma'(f_1) = 0.5$ ③ $f_3 = w_2 f_2 \Rightarrow \frac{\partial f_3}{\partial f_2} = w_2 = 0.5$
- $\frac{\partial f_1}{\partial x} = w_1$ ④ $\frac{\partial f_2}{\partial f_1} = \sigma'(f_1) = 0.5$ ⑤ $\frac{\partial f_3}{\partial f_2} = f_2 = 0.25$
- $\frac{\partial f_1}{\partial w_1} = x$ ⑥ $f_4 = \sigma(f_3) \Rightarrow \frac{\partial f_4}{\partial f_3} = \sigma'(f_3) = 0.25$ ⑦ $f_5 = \sigma(f_4) \Rightarrow \frac{\partial f_5}{\partial f_4} = \sigma'(f_4) = 0.25$
- ⑧ $f_6 = w_5 f_5 \Rightarrow \frac{\partial f_6}{\partial f_5} = w_5 = 0.5$ ⑨ $f_7 = \sigma(f_6) \Rightarrow \frac{\partial f_7}{\partial f_6} = \sigma'(f_6) = 0.25$ ⑩ $f_8 = (f_7)^2 \Rightarrow \frac{\partial f_8}{\partial f_7} = 2f_7 = 0.5$
- ⑪ $f_9 = f_8 - y \Rightarrow \frac{\partial f_9}{\partial f_8} = 1$

Backpropagation

~~$$\frac{dl}{df_1} = \frac{\partial f_1}{\partial q} = 1 \Rightarrow \frac{dl}{df_1} = 0.25$$~~

~~$$\frac{dl}{df_9} = \frac{\partial f_9}{\partial q} = \frac{\partial f_9}{\partial f_8} \times \frac{\partial f_8}{\partial q} = 0.5 \times 0.25 = 0.125$$~~

~~$$\frac{dl}{df_7} = \frac{dl}{df_8} \times \frac{\partial f_8}{\partial f_7} = 0.25 \times 2(f_7) = 0.25 \times 0.428 \times 2 = 0.428$$~~

~~$$= 0.125 \times 0.428 = 0.053125$$~~

~~$$\frac{dl}{df_6} = \frac{dl}{df_7} \times \frac{\partial f_7}{\partial f_6} = 0.428 \times 0.25 = 0.107$$~~

~~$$= 0.107 \times 0.25 = 0.027$$~~

~~$$\frac{dl}{df_5} = \frac{dl}{df_6} \times \frac{\partial f_6}{\partial f_5} = 0.027 \times 0.25 = 0.00675$$~~

~~$$= 0.00675 \times 0.25 = 0.0016875$$~~

~~$$= 0.0016875 \times 0.25 = 0.000421875$$~~

Batch norm → Batch Normalization

Normalization per batch. \rightarrow (i) nn. Linear (ii) nn. Batchnorm (iii) nn. ReLU (iv) nn. Sigmoid

Normalizing the inputs first \rightarrow (i) nn. Batchnorm (ii) nn. ReLU (iii) nn. Sigmoid (iv) nn. Linear

Trained during any loop.

Introduce before Activation functions

$$\text{Input} \rightarrow \text{BatchNorm} \rightarrow \text{ReLU}$$

Mini Batch GD.

- - - (measured first) - - - we backpropagate through the network with

Backpropagation

$$\frac{df_9}{df_9} = 1 \quad \frac{df_9}{df_8} = 1 \times 0.5 = 0.5$$

$$\frac{df_9}{df_7} = \frac{df_8}{df_7} \times \frac{df_9}{df_8}$$

= $2(f_7) \times 0.5$

$$f_1 = 0.8 \quad f_2 = \sigma f_1 \Rightarrow \sigma f_1(1 - \sigma f_1) = 0.622(1 - 0.622) = 0.428$$

$$f_3 = w_2 f_2 \Rightarrow \frac{df_3}{df_2} = 0.5 \Rightarrow 0.235$$

$$\frac{df_3}{dw_2} = 0.235 \quad \frac{df_9}{df_6} = \frac{df_7}{df_6} \times \frac{df_9}{df_7}$$

$$f_5 = w_3 = 0.5 \quad f_4 = 0.244 \quad f_6 = \sigma_5(1 - \sigma_5) = 0.574(1 - 0.574) = 0.244$$

$$f_7 = 0.428 \quad f_8 = 2(f_7) = 0.856$$

$$f_9 = 0.5$$

$$\frac{df_9}{df_5} = \frac{df_6}{df_5} \times \frac{df_9}{df_6}$$

$$= 0.245 \times 0.1831$$

$$= 0.0448$$

$$\approx 0.045$$

$$\frac{df_4}{df_4} = \frac{df_5}{df_4} + \frac{df_9}{df_5} \Rightarrow 0.5 \times 0.045 \Rightarrow 0.0225$$

~~$$\frac{df_9}{df_3} = \frac{df_4}{df_3} \times \frac{df_9}{df_4} \Rightarrow 0.244 \times 0.0225 \Rightarrow 0.00544 \approx 0.0055$$~~

~~$$\frac{df_9}{df_2} = \frac{df_3}{df_2} \times \frac{df_9}{df_3} \Rightarrow 0.5 \times 0.0055 \Rightarrow 0.002745$$~~

~~$$\frac{df_9}{df_1} = \frac{df_2}{df_1} \times \frac{df_9}{df_2} \Rightarrow 0.235 \times 0.002745 \Rightarrow 0.0006450$$~~

~~$$\frac{df_9}{df_6} = \frac{df_7}{df_6} \times \frac{df_9}{df_7} \Rightarrow 0.1 \times 0.428 \Rightarrow 0.428$$~~

~~$$\frac{df_9}{df_5} = \frac{df_6}{df_5} \times \frac{df_9}{df_6} \Rightarrow 0.245 \times 0.428 = 0.1048 \approx 0.105$$~~

~~$$\frac{df_9}{df_4} = \frac{df_5}{df_4} \times \frac{df_9}{df_5} \Rightarrow 0.5 \times 0.105 = 0.0525 \Rightarrow 0.0525$$~~

~~$$\frac{df_9}{df_3} = \frac{df_4}{df_3} \times \frac{df_9}{df_4} \Rightarrow 0.0525 \times 0.244 \Rightarrow 0.01281 \approx 0.013$$~~

~~$$\frac{df_9}{df_2} = \frac{df_3}{df_2} \times \frac{df_9}{df_3} \Rightarrow 0.5 \times 0.013 \Rightarrow 0.0064$$~~

~~$$\frac{df_9}{df_1} = \frac{df_2}{df_1} \times \frac{df_9}{df_2} \Rightarrow 0.235 \times 0.0064 \Rightarrow 0.00150$$~~

Vanishing Gradient Function

Exploding Gradient

Taking the weights as 5

$$f_1 = x w_1 \Rightarrow 1 \times 5$$

$$f_2 = \text{ReLU}(f_1) \Rightarrow \frac{5+|5|}{2} = 5$$

$$f_3 = w_2 f_2 \Rightarrow 5 \times 5 = 25$$

$$f_4 = \text{ReLU}(f_3) = \frac{25+25}{2} = 25$$

$$f_5 = w_3 f_4 \Rightarrow 5 \times 25 = 125$$

$$f_6 = \text{ReLU}(f_5) = 125$$

$$f_7 = f_6 - y \Rightarrow 125 - 1 = 124$$

$$f_8 = (f_7)^2 \Rightarrow (124)^2 = 15376$$

$$f_9 = 15376 \times 0.5 \Rightarrow 7688$$

$$df_9 = 0.5$$

$$df_8$$

$$\frac{df_2}{df_1} = 1$$

$$\frac{df_3}{df_2} = w_2 = 5$$

$$\frac{df_4}{df_3} = 1$$

$$\frac{df_5}{df_4} = w_3 = 5$$

$$\frac{df_6}{df_5} = 1 \quad \frac{df_7}{df_6} = 1$$

$$\frac{df_8}{df_7} = 2f_7 = 2 \times 124 = 248$$

Calculations:-

$$\frac{df_9}{df_8} = 1 \quad \frac{df_9}{df_8} = \frac{df_9}{df_8} \times 0.5 \Rightarrow 0.5$$

$$\frac{df_9}{df_7} = \frac{df_8}{df_7} \times \frac{df_9}{df_8} \Rightarrow 248 \times 0.5 \Rightarrow 124$$

$$\frac{df_9}{df_6} = \frac{df_7}{df_6} \times \frac{df_9}{df_7} \Rightarrow 1 \times 124 \Rightarrow 124 \quad \frac{df_9}{df_5} = \frac{df_6}{df_5} \times \frac{df_9}{df_6}$$

$$\frac{df_9}{df_5} = 1 \times 124$$

$$\frac{df_9}{df_4} = \frac{df_5}{df_4} \times \frac{df_9}{df_5} \Rightarrow 5 \times 124 \quad \frac{df_9}{df_3} = \frac{df_4}{df_3} \times \frac{df_9}{df_4} \Rightarrow 25 \times 124$$

$$\frac{df_9}{df_3} = 124$$

$$\frac{df_1}{df_3} = \frac{df_4}{df_3} \times \frac{df_1}{df_4} \Rightarrow 1 \times 6.20 \Rightarrow 6.20$$

$$\frac{df_1}{df_2} = \frac{df_3}{df_2} \times \frac{df_1}{df_3} \Rightarrow 5 \times 6.20 \Rightarrow 31.00$$

$$\frac{df_1}{dx} = \frac{df_2}{dx} \times \frac{df_1}{df_2} \Rightarrow 1 \times 31.00 \Rightarrow 31.00$$

$$\frac{df_1}{dx} \Rightarrow \cancel{31.00} \times \frac{df_1}{df_1} \times \frac{df_1}{df_1} \Rightarrow 5 \times 31.00 \Rightarrow \underline{\underline{155.00}}$$

Final Exploding Gradient

Batch Norm Continued

Scale and shift operations on top of variance and other operations
mean of gaussian is something else

Shifting of space might happen wrongly. The datapoint might
a shift somewhere else.

- Normalize each layer

- Apply to μ^2

- scale and shift \rightarrow Learnable parameters.

(γ) (β)

(becomes prominent with each layer)

$\gamma^{[l]}$ $\beta^{[l]}$

Example:- $z = [4, 8, 6, 10]$.

$$-\mu = 4+8+6+10 \Rightarrow 7$$

$$\sigma^2 = \frac{1}{4} \sum_{i=1}^4 (z_i - \mu)^2$$

$$= 5$$

Convolutional Neural Network

1 sample \rightarrow Image \rightarrow $64 \times 64 \times 1$ input layers huge

made of R G B channels passed as a vector

(channels)

join and give colour image.

Histogram of Oriented Gradient (HOG)

- " " Flow (HOF)

- Optical flow features

$[43 | 240 | 1 | 28]$ X

feature
detector

adjacent pixels will have

homogeneous property.

Convolution Operation.

Input 6x6 image					
3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	3	7	7	8
4	2	1	6	2	8
2	4	8	2	3	9

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

-5	-1	0	1	8
-10	-2	2	3	5
0	-2	-4	-7	1
-3	-2	-3	-15	2

Convolution operation
of
kernel
 3×3

filter

kernel

$$\text{Convolution operation} \Rightarrow (3 \times 1) + (0 \times 0) + (1 \times -1) + (1 \times 1) + (5 \times 0) + (8 \times -1) \\ + (2 \times 1) + (7 \times 0) + (2 \times -1) \\ \Rightarrow 3 - 1 + 1 + 0 - 8 + 2 + 0 - 2 \\ \Rightarrow -5$$

$$\Rightarrow (0 \times 1) + (-1 \times 0) + (2 \times -1) + (5 \times 1) + (8 \times 0) + (9 \times -1) + 1 \\ + (7 \times 1) + (2 \times 0) + (5 \times -1) \\ \Rightarrow 0 - 0 - 2 + 5 + 0 - 9 + 7 + 0 - 5 \\ \Rightarrow -4$$

$$\Rightarrow (1 \times 1) + (2 \times 0) + (7 \times -1) + (8 \times 1) + (9 \times 0) + (3 \times 0) + (2 \times 1) \\ + (5 \times 0) + (1 \times -1) \\ \Rightarrow 1 + 0 - 7 + 8 + 0 - 3 + 2 + 0 - 1 \\ \Rightarrow 0$$

Similarly $\Rightarrow 2 + 0 - 4 + 9 + 0 - 1 + 5 + 0 - 3$

$$\Rightarrow 1 + 0 - 8 + 2 + 0 - 2 + 0 + 0 - 3 \Rightarrow -10$$

$$\Rightarrow 5 + 0 - 9 + 7 + 0 - 5 + 1 + 0 - 1 \Rightarrow -2$$

$$\Rightarrow 8 + 0 - 3 + 2 + 0 - 1 + 3 + 0 - 7 \Rightarrow 2$$

$$\Rightarrow 9 + 0 - 1 + 5 + 0 - 3 + 1 + 0 - 8 \Rightarrow 3$$

$$\Rightarrow 2 + 0 - 2 + 0 + 0 - 3 + 4 + 0 - 0 \Rightarrow 0$$

$$\Rightarrow 7 + 0 - 5 + 1 + 0 - 1 + 2 + 0 - 6 \Rightarrow -2$$

$$\Rightarrow 2 + 0 - 1 + 3 + 0 - 0 + 1 + 0 - 2 \Rightarrow -4$$

$$\Rightarrow 5 + 0 - 3 + 1 + 0 - 8 + 6 + 0 - 8 \Rightarrow -7$$

$$\Rightarrow 0 + 0 - 3 + 4 + 0 - 1 + 2 + 0 - 0 \Rightarrow -3$$

$$\Rightarrow 1 + 0 - 1 + 2 + 0 - 6 + 4 + 0 - 2 \Rightarrow -2$$

$$\Rightarrow 3 + 0 - 7 + 1 + 0 - 2 + 5 + 0 - 0 \Rightarrow -3$$

$$\Rightarrow 1 + 0 - 8 + 6 + 0 - 8 + 2 + 0 - 9 \Rightarrow -16$$

The filter goes horizontally
and then slides

down by one box and then

again goes horizontally

down by one box and then

again goes horizontally

down by one box and then

again goes horizontally

Vertical filter \rightarrow

1	0	-1
0	0	1
1	0	-1

Vertical edges

Horizontal filter \rightarrow

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

Horizontal edges

edges are taken as pixels and passed it on through SVM.

Least int Padding

$$\begin{aligned}
 & \text{Input: } 6 \times 6 \\
 & \text{Kernel: } 3 \times 3 \\
 & \text{Output: } 3 \times 3 = 3
 \end{aligned}$$

$$n = \text{input} + 2 \times \text{padding} - \text{kernel} + 1 = 6 + 2 \times 1 - 3 + 1 = 6$$

$$p = \frac{n - \text{kernel} + 1}{2} = \frac{6 - 3 + 1}{2} = 2$$

After applying too many layers it might shrink too much.

and might vanish

less no. of hits \rightarrow edges

more no. of hits \rightarrow middle

To overcome — shrinking problem, equal attention to all pixels.

$$p=1 \quad n-f+1 \times n-f+1$$

$$= (6-3+1) \times (6-3+1)$$

when $n=6, p=1$

$$f=3 \quad o/p = 6 \times 6$$

for Padding $n+2p-f+1 \times n+2p-f+1$

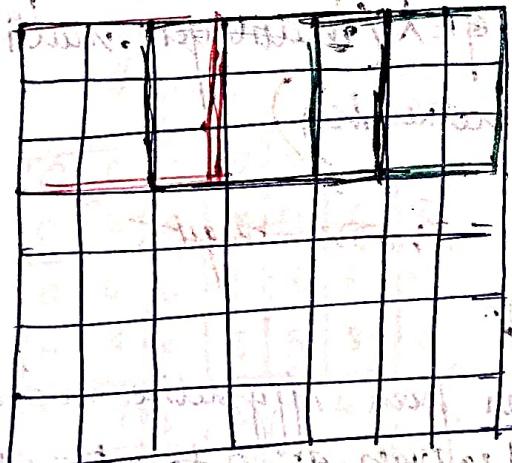
valid convolution \Rightarrow i/p size \geq o/p size

"same" convolution \Rightarrow i/p size $=$ o/p size

see what should be the value of p to get "same" conv?

$p = \frac{f-1}{2}$ (new term written w/ [conv structure])

Strided convolution: stride s is not applied at every step



3x3 kernel \leftarrow stride 2

Stride $s=2$ (hyperparameter)

faster processing \rightarrow leads to shrunken image (even rows x cols)

Padding will solve the issue.

1) GED

2) GTEx

3. 1000G

4. L1000

$\frac{1}{10} \text{GB}$

for training - 943
output - 9520

GEO data \rightarrow 80% training \rightarrow 88807

val \rightarrow 1101

test \rightarrow 1101

ARIA

Last layer is simple linear layer

hence linear activation $z = a$ or $a = z$

19x23

LM-1D }
TG-1D } Split

Input (11 genes)

Target-genes

pull them out

462 samples

1 sample \rightarrow 22K genes \rightarrow index into 1000 file.

for every sample

102 genes input

100 total

target

all the 1000 output

100 target

1000 total

2 3 2 3

2 3 2 3

2 3 2 3

P = 11

S = 21

1	-1
0	-1

f = 2x2

-2	-7	-6
-9	-14	-7
11	4	-3

0	0	0	0	0	0
0	2	3	4	5	6
0	6	6	9	8	7
0	3	4	8	3	8
0	7	8	3	6	6
0	11	2	1	8	3
0	0	0	0	0	0

n = 5x5

= Convolution operation = $100 \cdot 0 + 0 + 0 - 2 = -2$

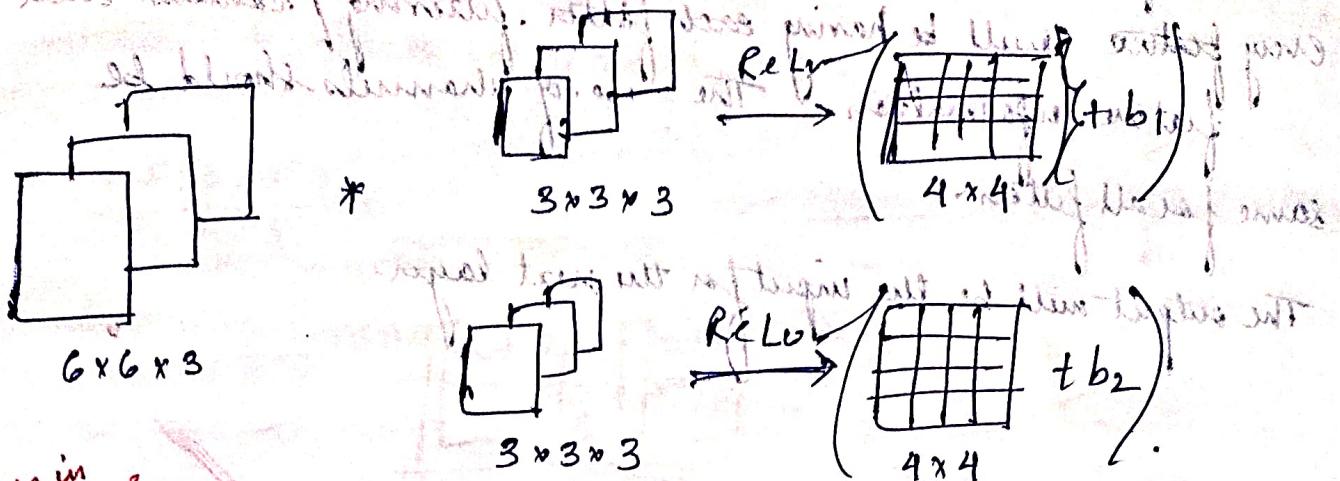
$$-6 - 3 = -9$$

$$\Rightarrow 6 - 9 + 0 - 8$$

$$\Rightarrow 5$$

$$= 0 + 0 + 0 - 7 = -7$$

$$= 0 + 0 + 0 - 6 = -6$$



~~filters in CNN~~

$$z^{[l+1]} = w^{[l+1]} a^{[l]} + b^{[l+1]}$$

$$a^{[l+1]} = \text{ReLU}(z^{[l+1]})$$

we apply the filter to a portion of the input and bias.

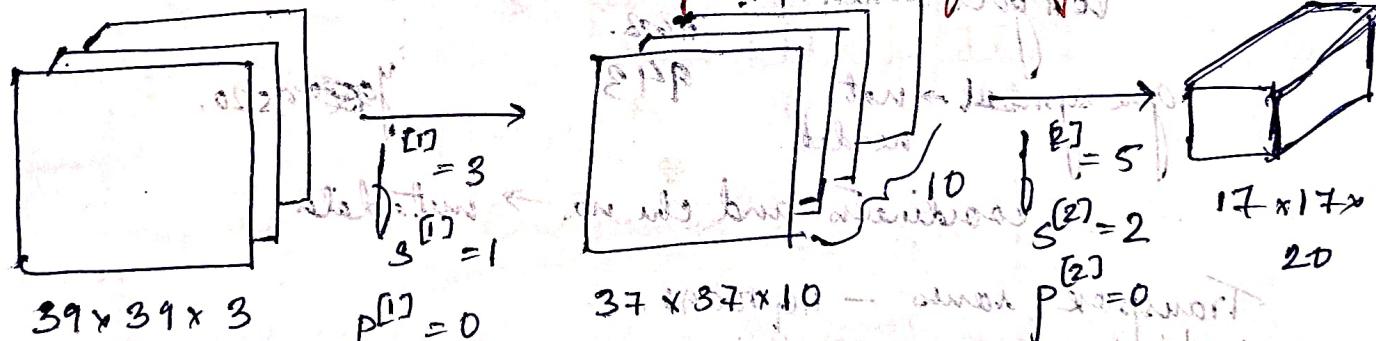
~~Sparsity of connections~~

Weight Sharing

Weights (filter) - same weight (filter) used across the entire image.

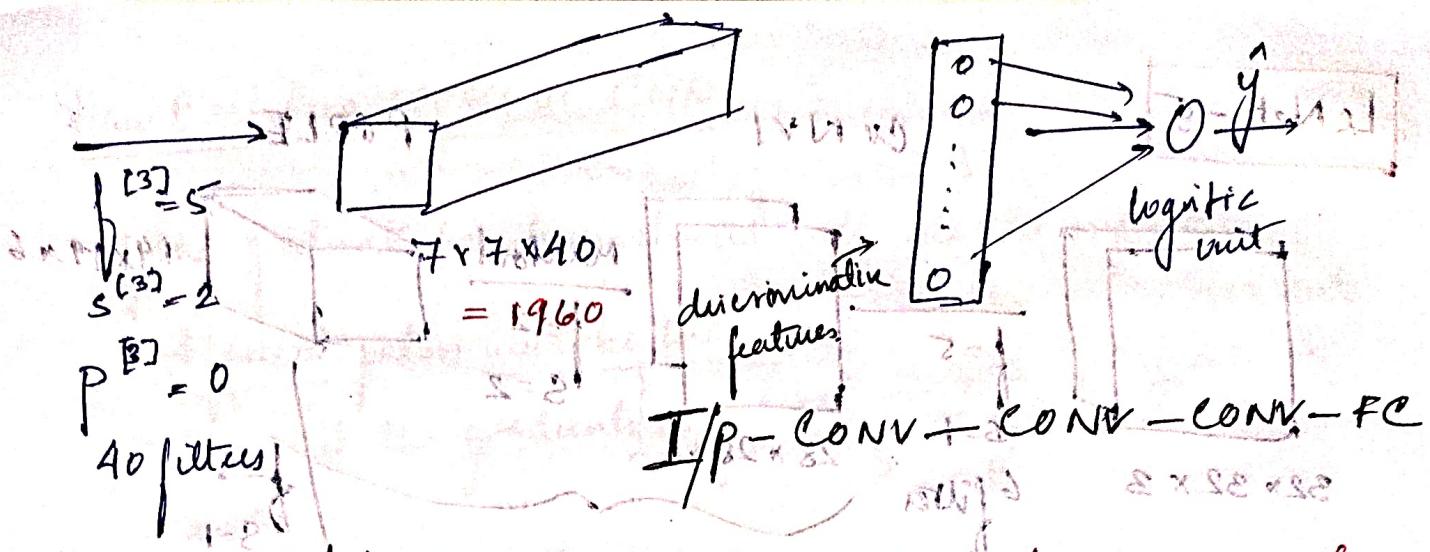
~~Sparsity of connections~~

These 2 properties make CNN efficient and will make it less prone to overfitting.



$$\frac{n+2p-f+1}{s} = \frac{(39+2(0)-3)}{3} + 1 \\ \Rightarrow 36 + 1 \\ \Rightarrow 37$$

$$\frac{n+2p-f+1}{s} = \frac{37-5}{2} + 1 \\ = \frac{32}{2} + 1 = 16 + 1 \\ \Rightarrow 17$$



$$\frac{n+2p-f}{s} + 1$$

$$\Rightarrow \frac{17-5}{2} + 1$$

$$\Rightarrow \frac{12+1}{2} \Rightarrow 6+1$$

The whole conv is only for one sample.

This is just one row.

(every channel picks up feature)

40 different types of features.

It is combined and flattened and passed through logits.

10 filters } in one layer of conv net

$3 \times 3 \times 3$ } based on how many channels are there

the size of kernel

How many parameters does this layer have?

$$3 \times 3 \times 3 = 27 \times 10 \text{ filters} \Rightarrow 270 + 10 \text{ biases}$$

$$\Rightarrow 280$$

~~total x 31 A~~

1. Conv. layer
 2. Pooling layers
 3. Fully connected (FC) layer.
- Max Pooling (Most commonly used) computationally expensive.

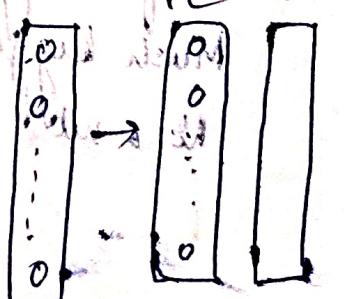
1	3	2	1	3
2	9	1	1	5
1	3	2	3	2
8	3	5	1	0
5	6	1	2	9

Pooling

9	9	5
9	9	5
18	6	9

$3 \times 3 \times n_c$

Flattened dimension

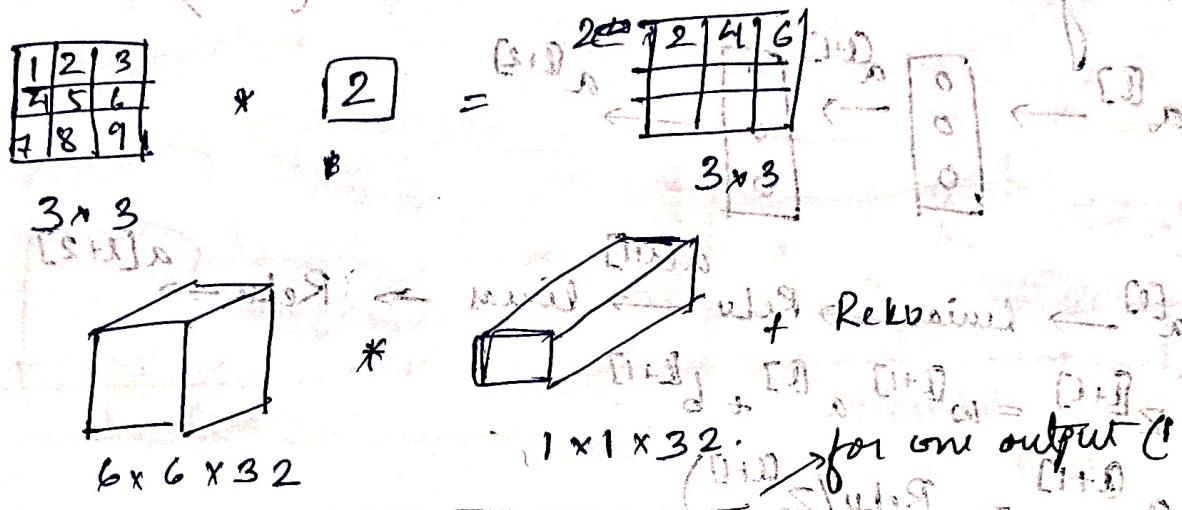


$$a^{[l+2]} = \text{ReLU} (z^{[l+2]} + a^{[l]})$$

$64 \times 1 \quad 32 \times / + W \cdot 64 \times 1$

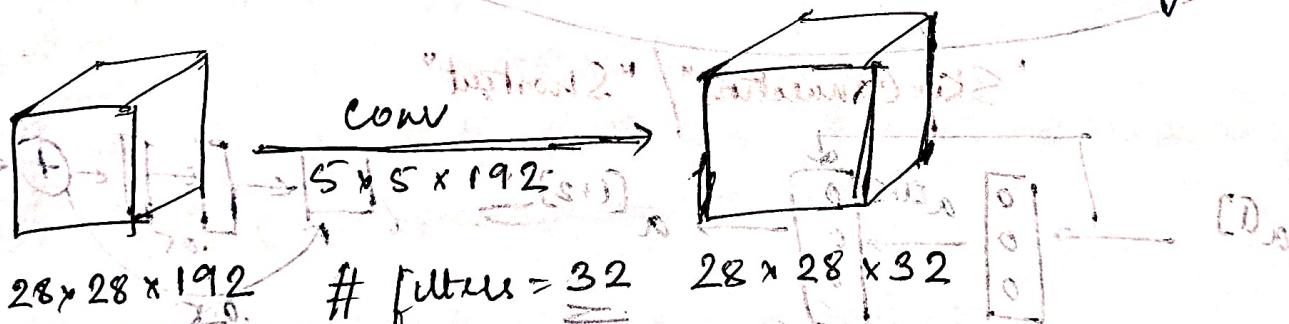
1x1 Convolution

(slow) operations with most of them are null products

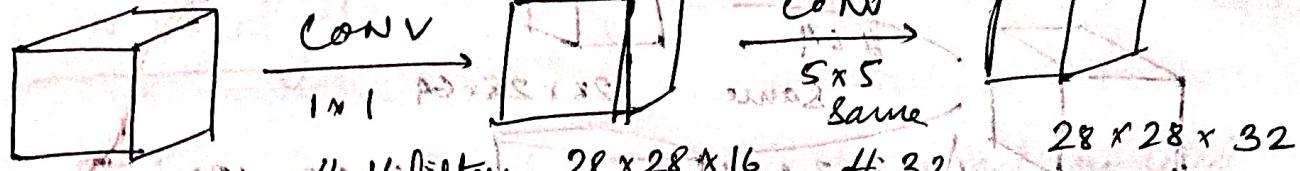


Also called Network in a network \rightarrow if this 1x1x32

is overlapped on $6 \times 6 \times 32$ then it just becomes a
MLP \rightarrow contains weights



Same \rightarrow $28 \times 28 \times 192$ each filter $\rightarrow 5 \times 5 \times 192$
 Total operations = $28 \times 28 \times 192 \times 5 \times 5 \times 192 = 4800 \times 28 \times 28$
 (mult.) $\rightarrow 153600 \times 28 \times 28$
 $\rightarrow 120,422,400$



$28 \times 28 \times 192$

16 filters

$28 \times 28 \times 16$

5×5

Same

$28 \times 28 \times 32$

32 filters

$$28 \times 28 \times 192 = 150,528$$

$$150,528 \times 16 \text{ filters} = 2,408,448$$

~~$$2,408,448 \times 5 \times 5 \times 32 = 19,267,520,000$$~~

$$\begin{array}{r}
 \cancel{\text{from } 192 \text{ to } 16} \\
 + \cancel{240,844,8} \\
 + \cancel{627,200} \\
 \hline
 \text{Bottleneck layer} \\
 \text{from } 192 \text{ to } 32. \\
 \hline
 240,844,8 \\
 + 10,035,200 \\
 \hline
 12,443,648
 \end{array}$$

1) Modity height widths

→ Padding

→ Striding, Pool

2) Modity channels

→ filters

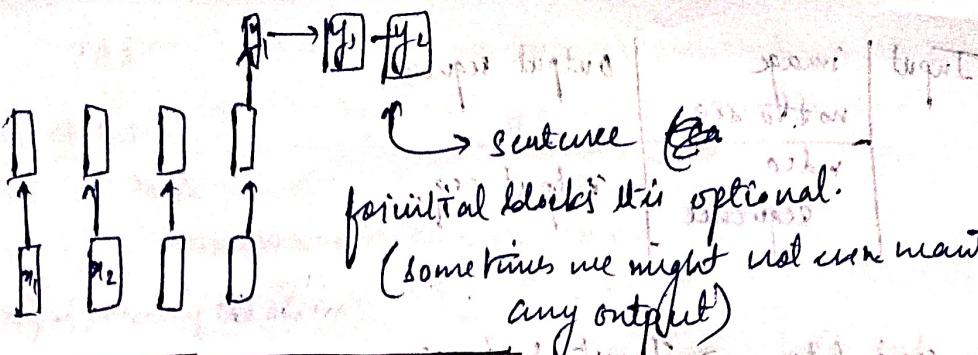
→ 1x1 CONV

INCEPTION NETWORK

Just put all filters and let the machine learn from it.

Inception Module

[Inception module]



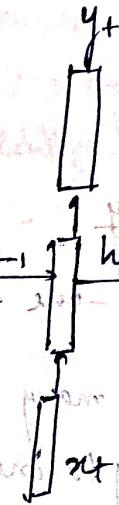
h_t is a function = $f(h_{t-1}, x_t)$

$$= \tanh(W_{1h} h_{t-1} + W_{2h} x_t)$$

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \tanh\left(\begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \alpha\right)$$

W_{1h} h_{t-1}
 4×4 4×1

W_{2h} x_t
 4×3 3×1



$$y_t = W_y h_t$$

(dimensions of W_y can be decided
(column is fixed))

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

2×3 3×1

weight sharing for all
(Whole shared across whole thing)

Should avoid a lot of work if we do it sequentially.

Vanilla
RNN

Character-level

2×1

$T=2$ (two time steps)

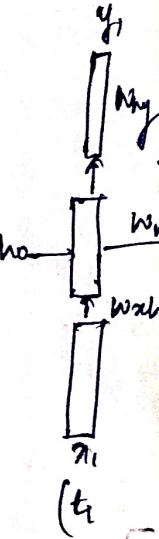
$$x^{(1)}: x_1 \in \mathbb{R}^2 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, x_2 \in \mathbb{R}^2 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

Input dimensions = 2 $\Rightarrow x_t \in \mathbb{R}^2$ has two neurons

Input dimensions = 2 $\Rightarrow x_t \in \mathbb{R}^2$ has two neurons

(all neurons are fully connected) probably not what we want

$$h_t = \tanh$$



$$\frac{e^{-1/2}}{e^{2/2} e^{-2}}$$

$$\begin{aligned} 0.9 &= 1.105 \\ 0.9 &= 1.105 \\ 0.3 &= 3.3 \end{aligned}$$

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

$$h_0 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, W_{xh} = \begin{bmatrix} 0.5 & -0.3 \\ 0.8 & 0.2 \\ 0.1 & 0.9 \end{bmatrix}, W_{hh} = \begin{bmatrix} 0.1 & 0.4 & 0.0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.2 \end{bmatrix}$$

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

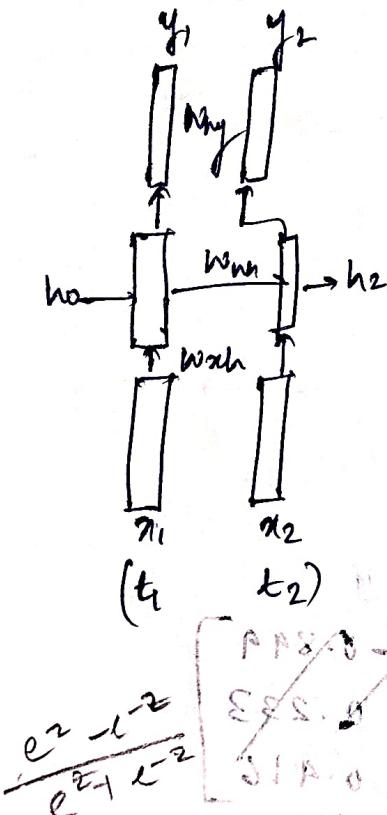
$$y_t = W_{yh} h_t \Rightarrow 2 \times 2 = [2 \times 3] \times (3 \times 1) \rightarrow h_t$$

$$h_t = 3 \times 1$$

$$y_t = 2 \times 1$$

$$\begin{bmatrix} 1.0 & -1.0 & 0.5 \\ 0.5 & 0.5 & -0.5 \end{bmatrix} \begin{bmatrix} 0.1 & 0.4 & 0.0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.2 \end{bmatrix} = \begin{bmatrix} 0.28 & -0.28 & 0.1 \\ 0.15 & 0.15 & -0.1 \end{bmatrix}$$

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



$$h_1 = \tanh(W_{hh} h_{t-1} + W_{xh} x_t)$$

$$h_2 = \tanh(W_{hh} h_1 + W_{xh} x_2)$$

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

$$= \tanh \left(\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0.5 & -0.3 & (2) \\ 0.8 & 0.2 & (2) \\ 0.1 & 0.9 & (2) \end{bmatrix} \right)$$

$$= \tanh \begin{bmatrix} -0.1 \\ 1.2 \\ 0.9 \end{bmatrix} = \begin{bmatrix} 0.099 \\ 0.833 \\ 0.716 \end{bmatrix}$$

Practical \rightarrow time steps of input.

$$\begin{bmatrix} 0.0 & 0.0 & 1.0 \\ 2.0 & 0.0 & 0.0 \\ 3.0 & 1.0 & 0.0 \end{bmatrix} \begin{bmatrix} y_1 = W_{11}x_1 + b_1 \\ y_2 = W_{21}x_1 + b_2 \\ y_3 = W_{31}x_1 + b_3 \end{bmatrix} = \begin{bmatrix} 1.0 & -1.0 & 0.5 \\ 0.5 & 0.5 & -0.5 \end{bmatrix} \begin{bmatrix} x \\ -0.099 \\ 0.833 \\ 0.716 \end{bmatrix}$$

$$= 1.0x + 0.099 + -1.0 \times 0.833 + 0.5 \times 0.716 \\ 0.5x - 0.099 + 0.5 \times 0.833 + -0.5 \times 0.716$$

$$= \begin{bmatrix} -0.099 + -0.833 + 0.358 \\ -0.0495 + 0.4165 + -0.358 \end{bmatrix}$$

$$= \begin{bmatrix} -0.579 \\ 0.009 \end{bmatrix}$$

$$h_2 = \tanh(W_{22}x_1 + W_{23}x_2)$$

~~$$= \tanh \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} + \begin{pmatrix} -0.5 & -0.3 \\ -0.8 & +0.2 \\ -0.1 & +0.4 \end{pmatrix}$$~~

~~$$= \tanh \begin{pmatrix} -0.099 \\ 0.833 \\ 0.716 \end{pmatrix}$$~~

~~$$\Rightarrow \tanh \begin{pmatrix} -0.87 \\ -0.8 \\ -0.6 \\ -0.6 \\ -0.3 \end{pmatrix} = \begin{pmatrix} -0.66 \\ -0.5 \\ 0.29 \\ -0.899 \\ 0.233 \\ 0.416 \end{pmatrix}$$~~

$$\begin{bmatrix} 0.0 \\ 2.0 \\ 3.0 \end{bmatrix} \Rightarrow \begin{bmatrix} -0.3 \\ 0.2 \\ 0.0 \end{bmatrix} \tanh \begin{pmatrix} -0.7 \end{pmatrix} = \begin{bmatrix} 0.0 \\ 0.1 \\ 0.0 \end{bmatrix}$$

hence output with unit weight

$$\tanh \left(\begin{bmatrix} 0.1 & 0.1 & 0.0 \\ -0.2 & 0.3 & 0.2 \\ 0.05 & -0.1 & 0.2 \end{bmatrix} \begin{bmatrix} -0.099 \\ 0.833 \\ 0.716 \end{bmatrix} \right) + \begin{bmatrix} -0.8 \\ -0.6 \\ 0.3 \end{bmatrix}$$

Given \Rightarrow $\begin{bmatrix} -0.099 + 0.332 + 0.000 \\ 0.0198 + 0.2499 + 0.1432 \\ -0.0099 + -0.0833 + 0.1432 \end{bmatrix} + \begin{bmatrix} -0.8 \\ -0.6 \\ 0.3 \end{bmatrix}$

$$\Rightarrow \begin{bmatrix} 0.3221 - 0.8 \\ 0.4129 - 0.6 \\ 0.055 + 0.3 \end{bmatrix} \Rightarrow \begin{bmatrix} -0.4779 \\ -0.1871 \\ 0.355 \end{bmatrix}$$

\Rightarrow

$$\begin{bmatrix} 0.5 & 0.3 & 0.1 \\ 0.2 & 0.2 & 0.1 \\ 0.1 & 0.1 & 0.1 \end{bmatrix}$$

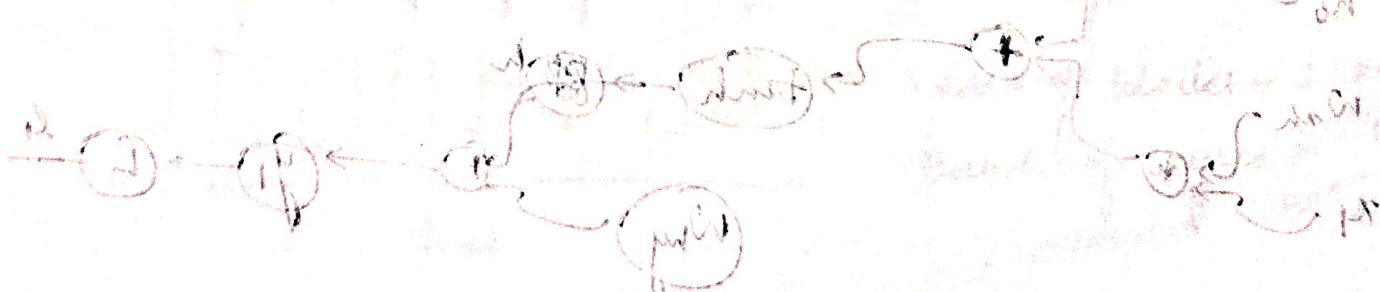
minimum attained

What we wanted

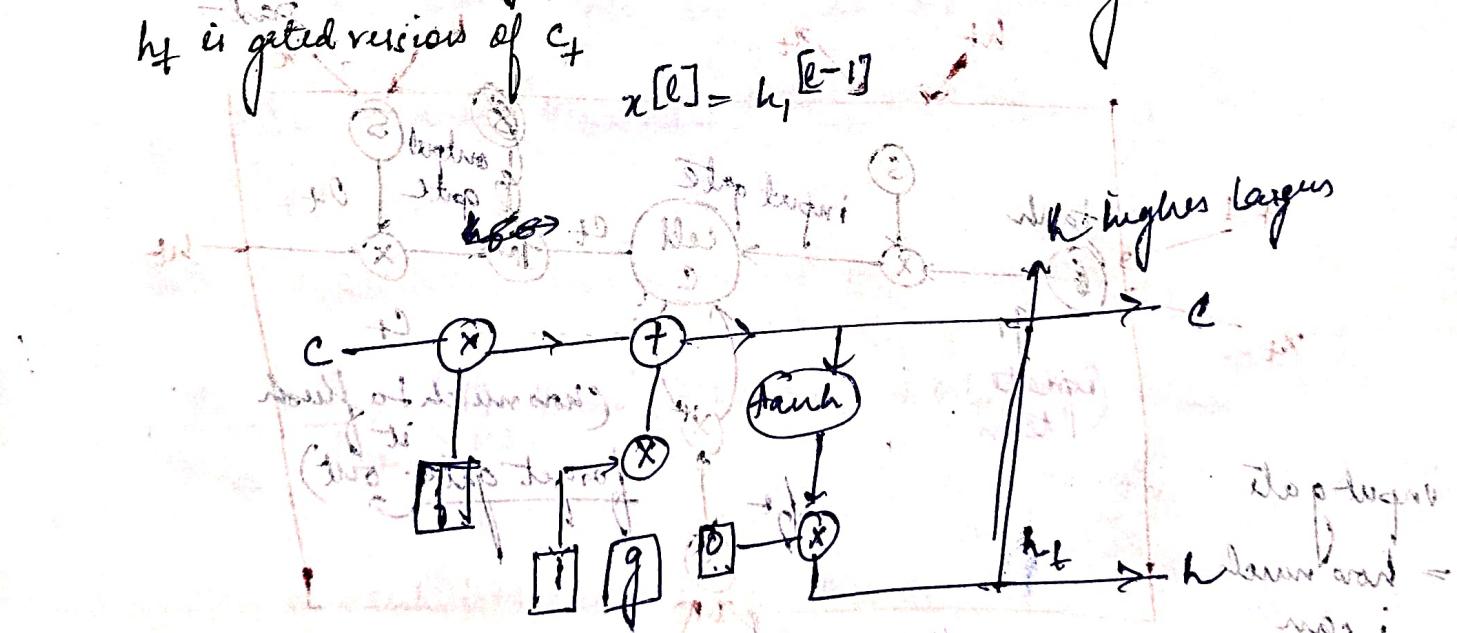
$$\begin{bmatrix} 0.5 \\ 0.3 \\ 0.1 \end{bmatrix} \leftarrow \text{as ad below}$$

$$\begin{bmatrix} 1.0 \\ 0.0 \\ 0.0 \end{bmatrix} \leftarrow \text{as ad}$$

$$+ \frac{\partial f}{\partial w_1} = \text{min}$$



element
 next $c_t = f_t \odot c_{t-1} + i_t \odot g_t$ Memory controlled by
 $h_t = o_t \odot \tanh(c_t)$ current/news.
 gate $f_t = \sigma(w_{af} h_{t-1} + w_{af} x_t)$
 gate $i_t = \sigma(w_{ai} h_{t-1} + w_{ai} x_t)$
 gate $o_t = \sigma(w_{ao} h_{t-1} + w_{ao} x_t)$
 note $g_t = \tanh(w_{hg} h_{t-1} + w_{gx} x_t)$ it is a boolean thing.
 gate $\times g_t = \tanh(w_{hg} h_{t-1} + w_{gx} x_t)$
 if f_t is a vector with ranges from 0 to 1.
 g_t is not a gate but is an update term.
 because of gates some or few will be allowed through the gate.
 c_t = previous memory.



want - stop first
 Stratified modelling (female only / male only)
 for demographic fl.
 with a new feature
 prediction

Example

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

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

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

$$w_{hi}^T = \begin{bmatrix} -0.4 & 0.2 \\ 0.3 & 0.3 \end{bmatrix}$$

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

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

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

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

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

$$h_t = ?$$

$$ct = ?$$

At t

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

$$= \sigma \left(\begin{bmatrix} 0.15 & 0.05 \\ 0.1 & -0.2 \end{bmatrix} \begin{bmatrix} 0.0 \\ 0.1 \end{bmatrix} \right) + \begin{bmatrix} 0.5 & 0.3 \\ 0.4 & 0.1 \end{bmatrix} \begin{bmatrix} 0.5 \\ 0.1 \end{bmatrix}$$

wrong calculation

$$= (0 + 0.005) + (0 + -0.02) + (0.25 + -0.03) (0.2 + 0.01)$$

$$\Rightarrow [0.01 - 0.02] + [0.16 - 0.14]$$

$$\Rightarrow [0.17 - 0.16] \Rightarrow [0.005 - 0.02] + [0.22 - 0.19]$$

$$\Rightarrow \sigma [0.225 \quad 0.17]$$

$$\text{OR } \frac{1}{1 + e^{-2.2}} \left[\frac{0.17}{0.17 + 0.78} \right] \Rightarrow \frac{1}{1 + e^{-2.2}} \left[\frac{0.17}{0.95} \right] \Rightarrow 0.56$$

$$[0.17 + 1] \Rightarrow \frac{1}{1 + 0.84} \Rightarrow \frac{1}{21.89} \Rightarrow 0.54$$

$$i_t = \boxed{0.78 \quad 0.54} \cdot \boxed{0.56 \quad 0.59}$$

$$o_t = \sigma \left(w_{hy} h_{t-1} + w_{xy} x_t \right)$$

= did by mistake in if:

$$\begin{bmatrix} 1.0 & -0.1 \\ -0.1 & 1.0 \end{bmatrix} = 0.5 \times \begin{bmatrix} 0.005 & -0.02 \\ -0.02 & 0.005 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 0.15 & -0.025 \\ -0.025 & 0.15 \end{bmatrix} = \begin{bmatrix} 0.125 & -0.12 \\ -0.12 & 0.125 \end{bmatrix}$$

$$\neq \begin{bmatrix} 0.005 & -0.002 \\ -0.002 & 0.005 \end{bmatrix} + \begin{bmatrix} 0.125 & -0.12 \\ -0.12 & 0.125 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 0.13 & -0.122 \\ -0.122 & 0.13 \end{bmatrix}$$

$$\sigma \begin{bmatrix} 0.13 & -0.122 \end{bmatrix}$$

$$\Rightarrow \frac{1}{1 + e^{-0.13}} = \frac{1.0}{1 + 0.87} = \frac{1.0}{1.87} = \frac{0.53}{1.87} \Rightarrow 0.53$$

$$\Rightarrow \frac{1}{1 + e^{+0.122}} = \frac{1}{1 + 1.29} = \frac{1}{2.29} = 0.43$$

$$\theta_t \Rightarrow \begin{bmatrix} 0.53 & 0.43 \\ 0.43 & 0.53 \end{bmatrix}$$

$$g_t = \tanh(w_{hy} h_{t-1} + w_{xy} x_t)$$

$$= \begin{pmatrix} 0.2 & 0.1 \\ -0.1 & 0.05 \end{pmatrix} \begin{bmatrix} 0.0 \\ -0.1 \end{bmatrix} + \begin{pmatrix} -0.5 & 0.4 \\ 0.2 & -0.3 \end{pmatrix} \begin{bmatrix} 0.5 \\ -0.1 \end{bmatrix}$$

$$= [0 + -0.01 \quad 0 + -0.005] + [-0.25 - 0.189 + 0.1 + 0.03]$$

$$\Rightarrow [-0.01 \quad -0.005] + [-0.29 \quad 0.13]$$

$$\Rightarrow [-0.35 \quad 0.125]$$

$$\tanh[-0.35 \quad 0.125]$$

$$\Rightarrow \frac{e^{-0.35} - e^{+0.35}}{e^{-0.35} + e^{+0.35}} \Rightarrow \frac{0.74 - 1.41}{0.74 + 1.41} \Rightarrow \frac{-0.67}{2.15}$$

$$\Rightarrow \frac{e^{0.125} - e^{-0.125}}{e^{0.125} + e^{-0.125}} = \frac{1.133 - 0.88}{1.133 + 0.88} \Rightarrow -0.311$$

$$\Rightarrow \frac{0.253}{2.013} \Rightarrow 0.125 \approx 0.13$$

Similarity score \rightarrow Score.

Content is a

unidirectional vector.

Alignment score \rightarrow

$$h_1 = [1, 0, 1] \text{ and } h_3 = [1, 1, 0] \quad s_{t-1} = [1, 0, 1]$$

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

Score of h = dot product.

$$c_t = ?$$

$$c_t = \alpha \sum_j^T \alpha_{t,j} h_j$$

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

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

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

$$= \frac{\exp(1, 0, 1)}{e^2 + e^1 + e^1}$$

$$= \frac{e^2}{e^2 + e^1 + e^1} \Rightarrow \frac{7.4}{7.4 + 2.7 + 2.7} \Rightarrow \frac{7.4}{12.8} \Rightarrow 0.57$$

$$\alpha_{t,2} = \frac{e^1}{12.8} \Rightarrow \frac{2.7}{12.8} \Rightarrow 0.21 \quad \alpha_{t,3} = \frac{e^1}{12.8} \Rightarrow 0.21$$

$$c_t = 0.57 + 0.21 + 0.21$$

$$\Rightarrow 0.99$$

Continuation of prev class

Thinking

Machines

① $q_1 = \text{Thinking}$

$K_1 = \text{Thinking}$

$v_1 = \text{Thinking}$

① $q_1 = \text{Machines}$

$K_1 = \text{Thinking}$

$q_2 = \text{Machines}$

② $q_1 = \text{Thinking}$

$K_2 = \text{Machines}$

$K_2 = \text{Machines}$

$1/p = \text{Playing Outside} \Rightarrow z_1 \text{ and } z_2$

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

$$K_1 = [0.31, 0.84, 0.963, 0.57]^T$$

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

outside

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

$$K_2 = [0.45, 0.99, 0.73, 0.58]^T$$

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

$\rightarrow A(\dim k)$

$$q_1 \cdot K_1 = [0.06572, 0.0336, 0.60664, 0.2052]$$

$$\Rightarrow 0.91116 / \sqrt{4} \Rightarrow 0.45558 \approx 0.45 \Rightarrow \text{Softmax}$$

$$q_1 \cdot v_1 = [0.0959, 0.0376, 0.4599, 0.2088]$$

$$\Rightarrow 0.8017 / [0.40085] \approx 0.48 \Rightarrow \frac{e^{0.45}}{e^{0.45} + e^{0.4}}$$

$$\frac{1.4}{2.9} \Rightarrow 0.48 \leq$$

$$\frac{1.5}{1.5+1.4} \Rightarrow 0.5 \leq$$

softmax x value $\Rightarrow 0.5 [0.36 \ 0.83 \ 0.12 \ 0.38]$

$$\Rightarrow [0.18 \ \cancel{0.415} \ 0.05 \ 0.19]$$

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

$$\Rightarrow [0.148 \ 0.1728 \ 0.091 \ 0.84]$$

add

$$\Rightarrow [0.328 \ \cancel{0.2628} \ 0.141 \ 0.53]$$

for 22

$$q_2 \times k_2 = \cancel{0.02} + 0.0212$$

$$\text{mindest} = 0.2$$

$$= [0.031 + 0.1176 + 0.828 + 0.4389]$$

$$\text{mindest} = 0.2$$

$$= 1.41155 / \sqrt{4}$$

$$\Rightarrow 0.70775$$

$$\approx 0.7 \quad \text{softmax} = \frac{2.01}{2.01 + 1.8} \Rightarrow \frac{2.01}{3.81} \Rightarrow 0.52$$

$$q_2 \times k_2 \Rightarrow [0.45 + 0.1316 + 0.6278 + 0.4466]$$

$$\Rightarrow \cancel{1.65} 1.251 \Rightarrow 0.62 \Rightarrow \text{softmax} \Rightarrow \frac{1.8}{3.81}$$

$$\Rightarrow 0.472$$

softmax

softmax x value $\Rightarrow 0.52 [0.36 \ 0.83 \ 0.12 \ 0.38]$

$$0.4 [0.31 \ 0.36 \ 0.19 \ 0.72]$$

$$\Rightarrow [0.1872 \ 0.4316 \ 0.052 \ 0.1976]$$

$$\checkmark \Rightarrow [0.3112 \ 0.5756 \ 0.0128 \ 0.4856]$$

$$P_E = \frac{\sin(\text{pos})}{\text{pos}_i} \quad P_{\text{par2}ii} = \cos\left(\frac{\text{pos}_i}{100 \cdot (2i/R^2)}\right)$$

for
this
do this
with 100
 $i = 0$

$$P_{0,0} = \sin\left(\frac{0}{100 \cdot 0}\right) \quad \sin \frac{0}{100 \cdot 0} \Rightarrow \sin 0 = 0$$

$$\Rightarrow \sin 0 = 0$$

$$\Rightarrow 0 = 0$$

$$P_{0,1} = \sin\left(\frac{0}{100 \cdot 1/4}\right) \cos\left(\frac{0}{100 \cdot 1/4}\right)$$

$$\Rightarrow \cos \frac{0}{1/4} = 1$$

parallel lines
at 90 degrees

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

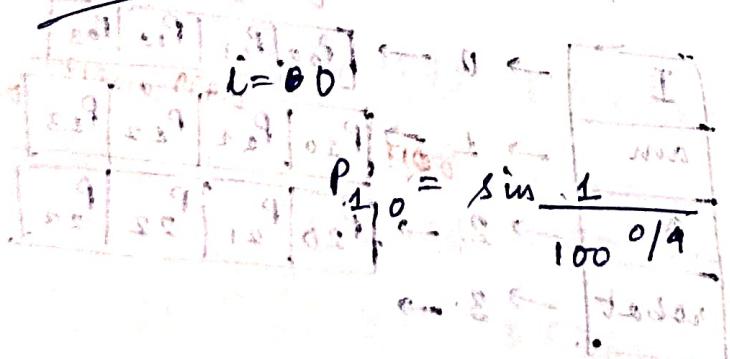
$$\Rightarrow \sin \frac{0}{10} = 0$$

$$\Rightarrow 0 = 0$$

$$P_{0,3} = \cos\left(\frac{0}{100 \cdot 3/4}\right)$$

$$\Rightarrow \cos \frac{0}{3/4} = 1$$

Position 1



$$P_1^0 = \sin \frac{1}{100 \cdot 0/4} \Rightarrow \sin \frac{1}{1} = 0.017$$

sliding window $P_{1,1} = \cos \frac{1}{100 \cdot 1} \Rightarrow \cos \frac{1}{1 \cdot 2} \Rightarrow 0.999$

$$i=1 \quad P_{1,1} = \sin \frac{1}{100 \cdot 1} \Rightarrow \sin \frac{1}{10} \Rightarrow 0.0017$$

$$P_{1,2} = \cos \frac{1}{100 \cdot 2} \Rightarrow 0.999$$

$$P_{1,3} = \sin \frac{1}{100 \cdot 3} \Rightarrow 0.0017$$

Position 2

$$i=0 \quad P_{2,0} = \sin \frac{2}{100 \cdot 0} = \sin \frac{2}{1} = 0.039$$

$$P_{2,1} = \cos \frac{2}{100 \cdot 1} = \cos \frac{2}{1 \cdot 2} \Rightarrow 0.999$$

$$P_{2,2} = \sin \frac{2}{100 \cdot 2} \Rightarrow \sin \frac{2}{10} = 0.0008$$

$$P_{2,3} = \cos \frac{2}{100 \cdot 3} \Rightarrow \cos \frac{2}{1.73} = 0.0006$$

Dot product of initial positional vector will be high

$$P_{500} \quad \text{dot} = \boxed{}$$

$$P_{501} \quad \text{prod}$$

Continuation

Multihed Attention

$$\boxed{y_1 \ y_2 \ y_3} \quad \boxed{x_1 \ x_2 \ x_3}$$

in multihead attention

the self attention will

look at everything but to predict y_3 suppose

we need to mask

the further one.