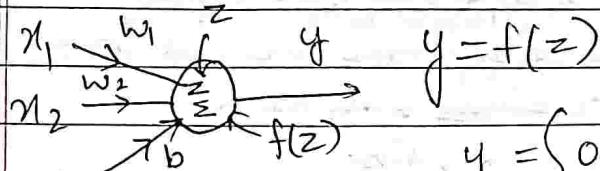


13/08

Deep Learning

Page No.	
Date	



$$y = \begin{cases} 0, & f(z) \leq 0 \\ 1, & f(z) > 0 \end{cases}$$

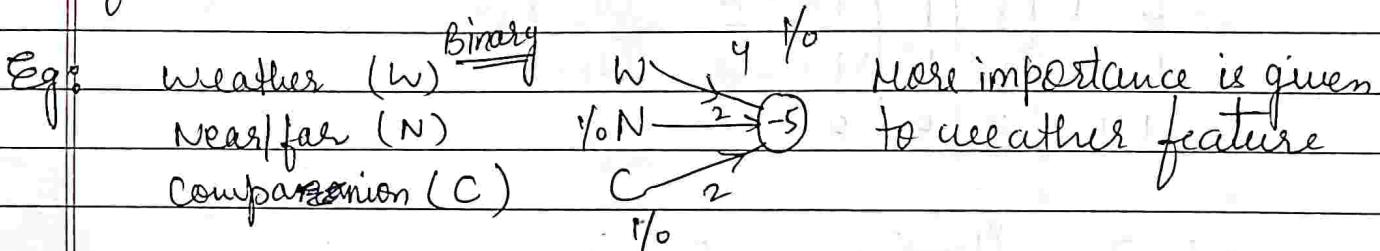
$$z = w^T u + b$$

w = weight

b = bias

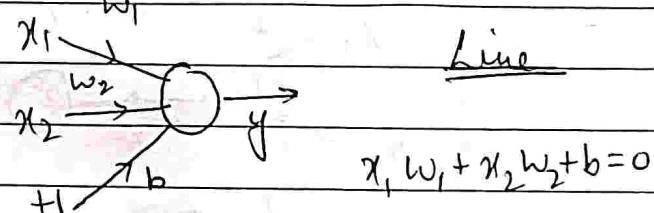
Linear classifier

More value of weight, more important is the feature



→ OR gate

x_1	x_2	O/p
0	0	0
0	1	1
1	0	1
1	1	1



If $x_1 w_1 + x_2 w_2 + b \geq 0 \rightarrow y = 1$
 $\leq 0 \rightarrow y = 0$

$x_1, x_2 = y$

for $w_1 = 1$ $0 \ 0 \ 0 \ 0$

$w_2 = 1$ $0 \ 1 \ 1 \ 1$

$b = 0$ $1 \ 1 \ 2 \ 1$

→ AND gate

x_1	x_2	O/p
0	0	0
0	1	0
1	0	0

1 1 1 Class 2

→ XOR

x_1	x_2	O/p
0	0	0
0	1	1
1	0	1
1	1	0

C₁ C₂

C₂ C₁

C₁

NAND \rightarrow Universal gate

$x_1 \ x_2$ O/P

0 0 1

0 1 1

1 0 1

0 1 1

$$x_1 w_1 + x_2 w_2 + b = y z$$

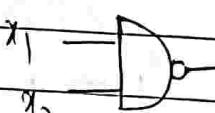
$$z = 0w_1 + 0w_2 + b \Rightarrow z = b$$

$$z = w_1 + w_2 + b$$

$$z = w_1 + b$$

$$z = w_1 + w_2 + b$$

$$w_1 = -1, w_2 = -1, b = 2$$



$x_1 \ x_2$ O/P z

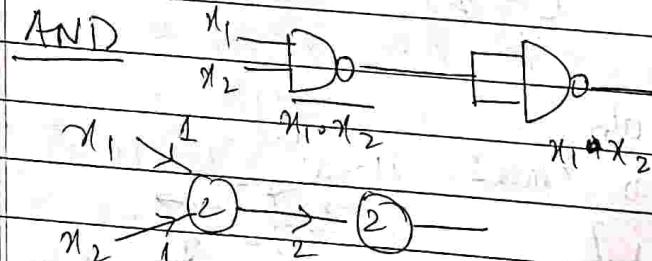
0 0 1 2 1

0 1 1 1 1

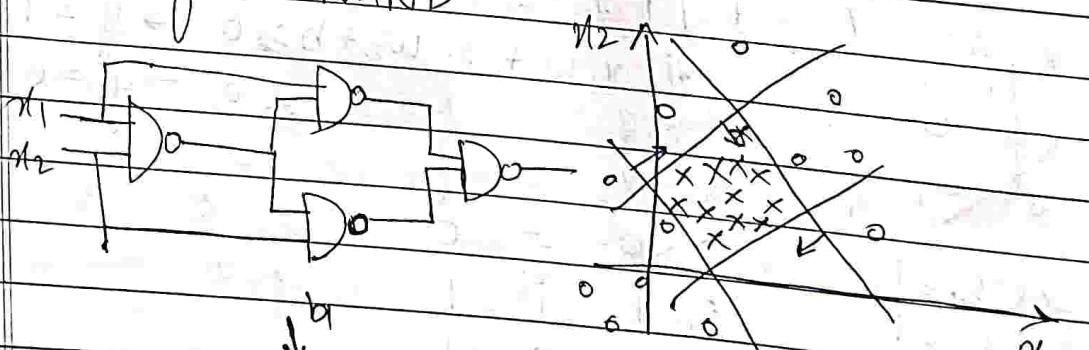
1 0 1 1 1

1 1 0 0 0

AND

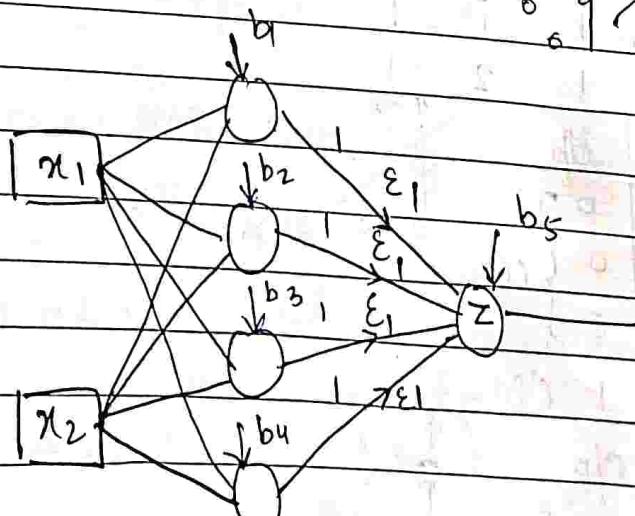


XOR from NAND



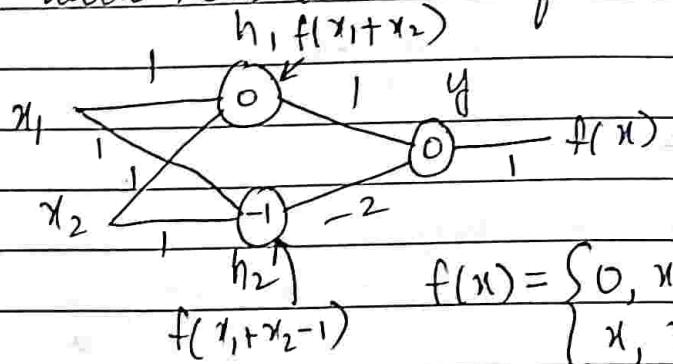
0 - class 1

* - Class 2



Piece wise
linear

XOR with ReLU activation function

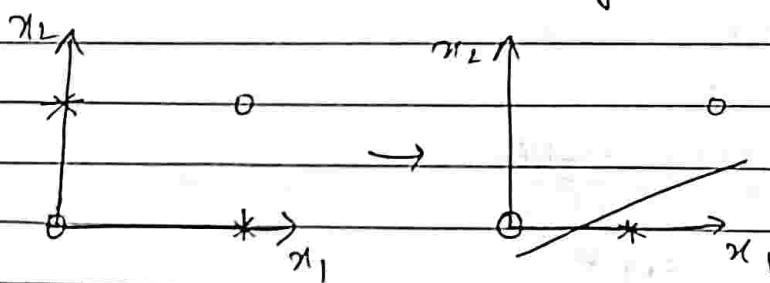


$$x_1, x_2 (0, 0) \quad h_1, h_2 (0, 0) \quad y = f(h_1 - 2h_2) = 0$$

$$x_1, x_2 (0, 1) \quad h_1, h_2 (1, 0) \quad y = f(1 - 0) = 1$$

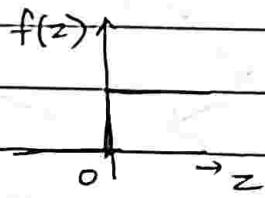
$$x_1, x_2 (1, 0) \quad h_1, h_2 (1, 0) \quad y = f(1 - 0) = 1$$

$$x_1, x_2 (1, 1) \quad h_1, h_2 (2, 1) \quad y = f(2 - 2) = 0$$



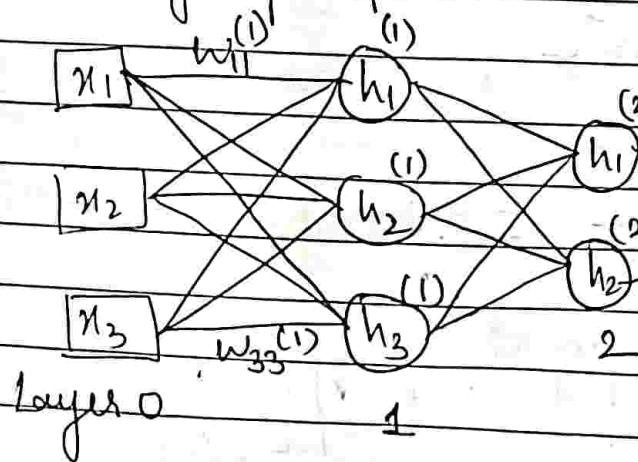
* Activation function

(1) Step function



(2) Sigmoid function

Multi layer perceptron



Wab
(c)

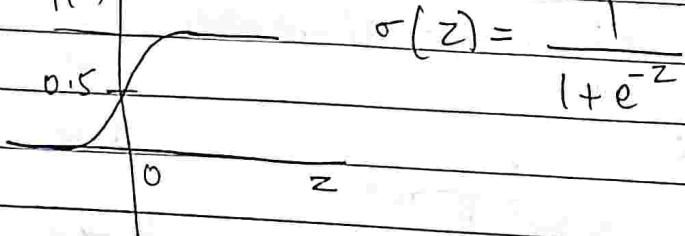
a → coming from
b → going to
c → layer no.

3

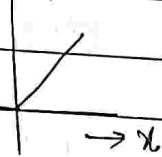
Layer 0

1

2

$f(n)$ 

③ ReLU function

 $f(n)$ ReLU

$$f(n) = \begin{cases} 0, & n \leq 0 \\ n, & n > 0 \end{cases}$$

27/8

MLP

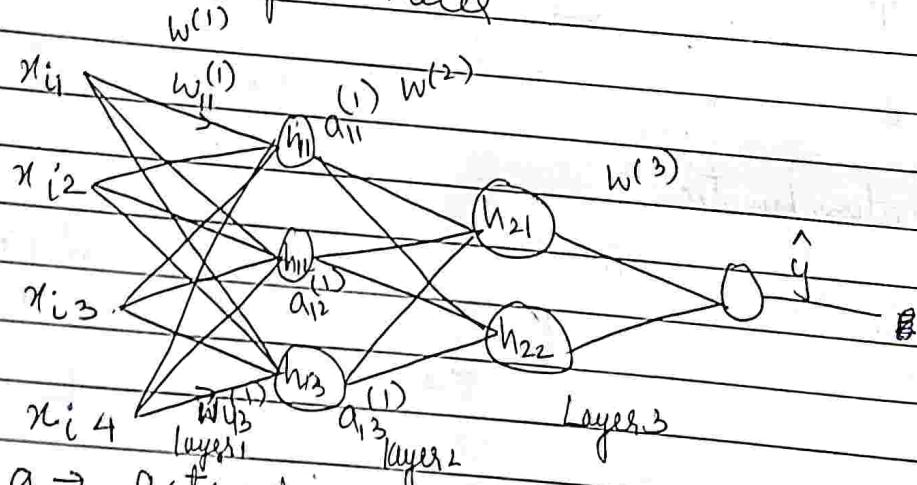
→ Feed forward propagation

→ loss

→ Back propagation

→ Activation

→ Trainable parameters

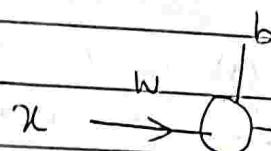


a → activation

$z \rightarrow w^T x + b$

x → input

w → weight b → bias



$z = w^T x + b$

$a = \sigma(z)$

$z^{(1)} = w^{(1)T} x + b^{(1)}$

$a^{(1)} = \sigma(z^{(1)})$

$$z^{(1)} = \begin{bmatrix} w_{11}^{(1)} & w_{12}^{(1)} & w_{13}^{(1)} \\ w_{21}^{(1)} & w_{22}^{(1)} & w_{23}^{(1)} \\ w_{31}^{(1)} & w_{32}^{(1)} & w_{33}^{(1)} \\ w_{41}^{(1)} & w_{42}^{(1)} & w_{43}^{(1)} \end{bmatrix}^T \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \\ x_{i4} \end{bmatrix} + \begin{bmatrix} b_{11}^{(1)} \\ b_{12}^{(1)} \\ b_{13}^{(1)} \end{bmatrix}$$

$(4 \times 3)^T \quad (4 \times 1) + (3 \times 1) =$

21/8

DL

$$a^{(2)} = \sigma(z^{(2)}) \quad \text{layer 2}$$

$$z^{(2)} = w^{(2)T} a^{(1)} + b^{(2)}$$

$$z^{(2)} = \begin{bmatrix} w_{11}^{(2)} & w_{12}^{(2)} \\ w_{21}^{(2)} & w_{22}^{(2)} \\ w_{31}^{(2)} & w_{32}^{(2)} \end{bmatrix} \begin{bmatrix} a_{11}^{(1)} \\ a_{12}^{(1)} \\ a_{13}^{(1)} \end{bmatrix} + \begin{bmatrix} b_{11}^{(2)} \\ b_{12}^{(2)} \end{bmatrix}$$

$$(3 \times 2)^T (3 \times 1) + (2 \times 1) \rightarrow 2 \times 1 \rightarrow z^{(2)}$$

$$z^{(3)} = \begin{bmatrix} w_{11}^{(3)} \\ w_{12}^{(3)} \end{bmatrix}^T \begin{bmatrix} a_{11}^{(2)} \\ a_{12}^{(2)} \end{bmatrix} + b_{11}^{(3)}$$

$$\sigma \left(w^{(2)T} \left[\sigma \left(w^{(1)T} a^{(0)} + b^{(1)} \right) \right] + b \right)$$

- Loss functions

	x_1	x_2	class/value	$w^{(1)}$	$w^{(2)}$	$y = 6.7$
1	6.9	7.2	7	6.7	$x_{11} \rightarrow$	
2			8	7.2	$x_{12} \rightarrow$	
3			9	9.2	$x_{12} \rightarrow$	
4					$\hat{y} = 6.7$	

$$L = y_i - \hat{y}_i = \text{loss for single input}$$

cost for entire dataset

$$\text{loss} = f(w, b)$$



Regression

- Mean Squared error (MSE)
- Mean absolute error (MAE)
- Huber loss

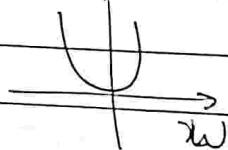
Classification

- Binary cross entropy
- Categorical cross entropy

MSE (Mean Squared Error)

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{mean}(e^2)$$

$(y_i - \hat{y}_i)^2$ \Rightarrow quadratic \rightarrow differentiable \rightarrow 1 local minima



More / elastic impact of outlier

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad \begin{array}{l} \text{loss} \\ \text{Non} \\ \text{differentiable} \end{array}$$

\rightarrow low impact of outliers

29/8. Huber Loss

Parameter that decides whether it is an outlier or not

$$L = \begin{cases} \frac{1}{2} (y_i - \hat{y}_i)^2 & ; \text{ if } |y_i - \hat{y}_i| \leq s \\ s |y_i - \hat{y}_i| - s^2 & ; \text{ otherwise } \end{cases} \quad \begin{array}{l} \text{if } y_i \text{ is not} \\ \text{outlier} \\ \text{when } y_i \text{ is outlier} \end{array}$$

s = trainable parameter

Decides whether datapoint y_i is not or not an outlier

MSE $\text{If } L = (y_i - \hat{y}_i)^2 \rightarrow \text{cost} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$ Entire set

Single data point \uparrow guide the n/w used for test & validation set

Training data 1 \rightarrow FP \rightarrow loss
 \downarrow Update w, b
 data 2 \rightarrow FP \rightarrow loss
 \downarrow Update w, b

MAE $L = |y_i - \hat{y}_i|$
 $\text{cost} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$

29/8

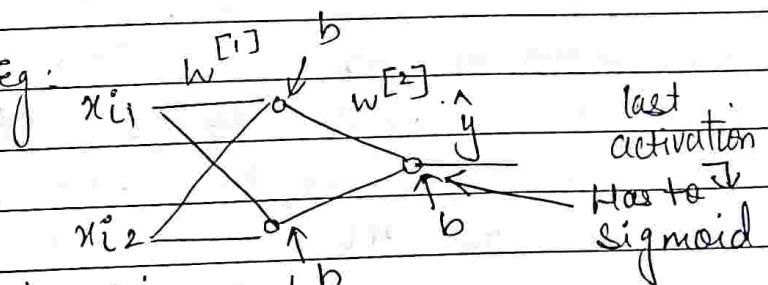
DL

• Classification

i. Binary Cross Entropy

x_1	x_2	y	\hat{y}
70	90	1	0.6
50	40	0	0.2
60	75	1	0.4

Eg:



For binary class classification only

$$L = -y_i \log(\hat{y}_i) - (1-y_i) \log(1-\hat{y}_i)$$

$$\text{Cost} = -\frac{1}{n} \sum_{i=1}^n y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)$$

For 1st set of input, $y_i = 1$ but $\hat{y}_i = 0.6$

$$\Rightarrow L = -\log(0.6)$$

2nd case $\rightarrow y_i = 0, \hat{y}_i = 0.2$

$$\Rightarrow L = -\log(0.8)$$

3rd case $\rightarrow y_i = 1, \hat{y}_i = 0.4 \Rightarrow L = -\log(0.6)$

ii. Categorical Cross Entropy

For multi-class classification

$$L = \sum_{j=1}^k -y_j \log(\hat{y}_j) \text{ where } k = \text{no. of classes}$$

$$\text{Eg: If } k=3 \Rightarrow L = -y_1 \log(\hat{y}_1) - y_2 \log(\hat{y}_2) - y_3 \log(\hat{y}_3)$$

x_1	x_2	y	OHE	Ground truth	Last activation has. to be	Softmax
70	90	Yes	[1 0 0]			[1 0 0]
60	40	No	[0 1 0]			[0 1 0]
40	40	Maybe	[0 0 1]			[0 0 1]
60	75	Yes	[1 0 0]			[1 0 0]

$$\sum_{i=1}^k f(z_i) = 1$$

$$f(z_i) = e^{z_i} \text{ activation}$$

$$y_1=1, y_2=0, y_3=0$$

$$\sum_{j=1}^k e^{z_j} \text{ for last layer}$$

$$x_{i1} \quad x_{i2} \quad \text{---} \quad y_1=0, y_2=1, y_3=0$$

$$x_{i1} \quad x_{i2} \quad \text{---} \quad y_1=0, y_2=0, y_3=1$$

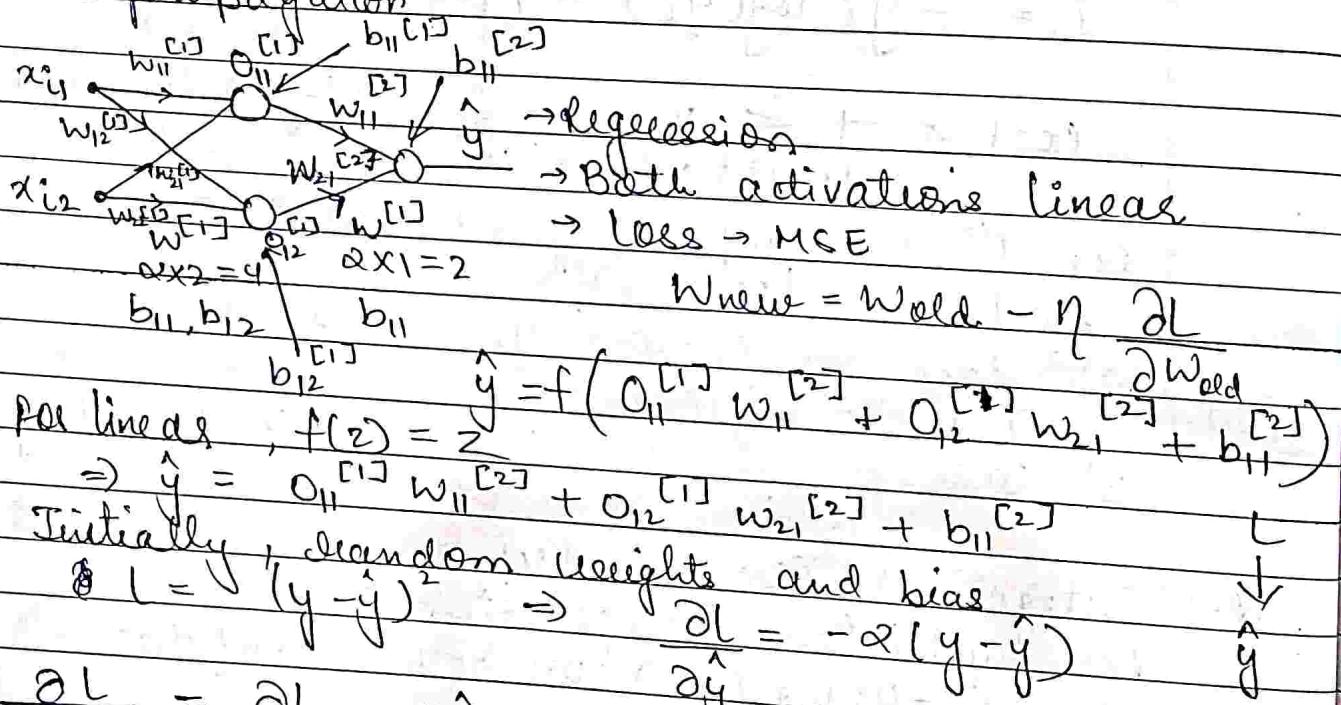
DL

3. Sparse Categorical Entropy

x_1	x_2	y	OHE	Pred	L	O/P
70	90	Yes	1 0 0	[0.6, 0.2, 0.2]	$\log 0.6$	1
60	40	No	0 1 0	[0.3, 0.5, 0.2]	$\log 0.5$	2
40	40	Maybe	0 0 1			3
60	75	Yes	1 0 0			1

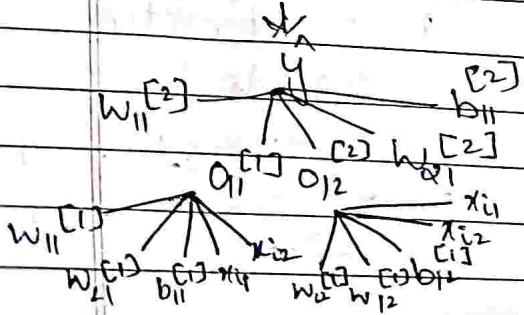
2/9/24

Backpropagation



$$\frac{\partial L}{\partial w_{11}^{[2]}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_{11}^{[2]}}$$

$$\frac{\partial L}{\partial w_{21}^{[2]}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_{21}^{[2]}}$$



$$\frac{\partial L}{\partial b_{11}^{[2]}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b_{11}^{[2]}}$$

$$\frac{\partial \hat{y}}{\partial w_{11}^{[2]}} = o_{11}^{[2]}$$

$$\Rightarrow \frac{\partial L}{\partial w_{11}^{[2]}} = -\alpha(y - \hat{y}) o_{11}^{[2]}$$

$$\frac{\partial L}{\partial w_{21}^{[2]}} = -\alpha(y - \hat{y}) o_{12}^{[2]}$$

To update $w_{21}^{[2]}$

$$\frac{\partial L}{\partial b_{11}^{[2]}} = -\alpha(y - \hat{y}) \cdot 1$$

for $b_{11}^{[2]}$

2/9.

DL

Backpropagation

$$w_{ii}^{[2]}_{\text{new}} = w_{ii}^{[2]}_{\text{old}} - \eta \cdot \frac{\partial L}{\partial w_{ii}^{[2]}} \quad \leftarrow \text{To update } w_{ii}^{[2]}$$

$$O_{ii}^{[1]} = f(x_{ii} w_{ii}^{[1]} + x_{i2} w_{i2}^{[1]} + b_{ii}^{[1]})$$

$$O_{ii}^{[1]} = x_{ii} w_{ii}^{[1]} + x_{i2} w_{i2}^{[1]} + b_{ii}^{[1]} \quad \text{Here for linear}$$

$$\frac{\partial L}{\partial w_{ii}^{[1]}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial O_{ii}^{[1]}} \cdot \frac{\partial O_{ii}^{[1]}}{\partial w_{ii}^{[1]}} =$$

$$\frac{\partial L}{\partial w_{i2}^{[1]}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial O_{ii}^{[1]}} \cdot \frac{\partial O_{ii}^{[1]}}{\partial w_{i2}^{[1]}}$$

$$\frac{\partial L}{\partial w_{ii}^{[1]}} = -\alpha(y - \hat{y}) \cdot w_{ii}^{[2]} \cdot x_{ii} \quad \frac{\partial L}{\partial b_{ii}^{[1]}} = -2(y - \hat{y}) w_{ii}^{[2]}$$

$$\frac{\partial L}{\partial w_{i2}^{[1]}} = -\alpha(y - \hat{y}) \cdot w_{ii}^{[2]} \cdot x_{i2}$$

$$O_{i2}^{[1]} = x_{ii} w_{i2}^{[1]} + x_{i2} w_{i2}^{[1]} + b_{i2}^{[1]}$$

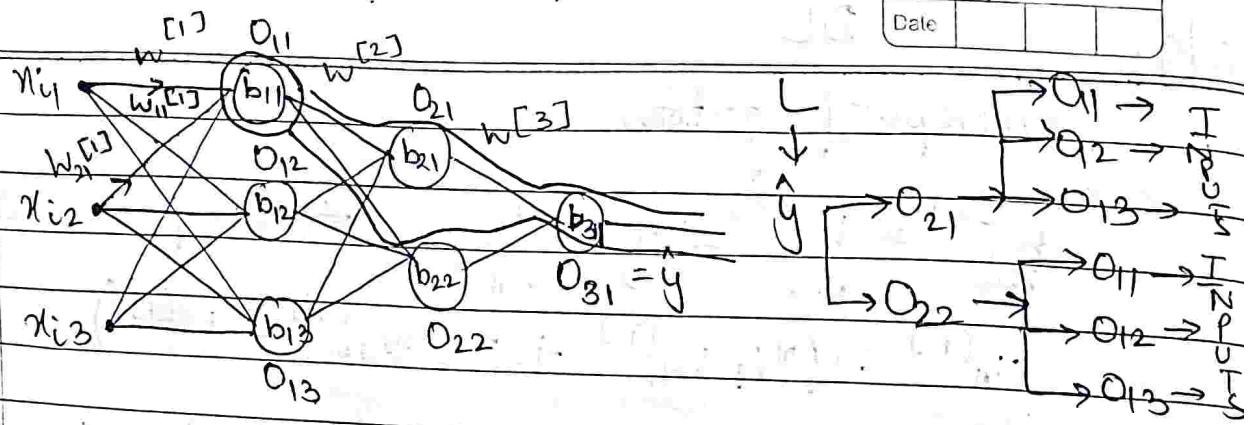
$$\frac{\partial L}{\partial w_{i2}^{[1]}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_{i2}^{[1]}} \frac{\partial O_{i2}^{[1]}}{\partial w_{i2}^{[1]}} = -\alpha(y - \hat{y}) w_{i2}^{[2]} x_{ii}$$

$$\frac{\partial L}{\partial w_{i2}^{[1]}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_{i2}^{[1]}} \frac{\partial O_{i2}^{[1]}}{\partial w_{i2}^{[1]}} = -\alpha(y - \hat{y}) w_{i2}^{[2]} x_{i2}$$

$$\frac{\partial L}{\partial b_{i2}^{[1]}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_{i2}^{[1]}} \frac{\partial O_{i2}^{[1]}}{\partial b_{i2}^{[1]}} = -\alpha(y - \hat{y}) w_{i2}^{[2]}$$

2/9 → Lab → TSNE

3/9 Memorisation in MLP (Multilayer Perception)



$$x \xrightarrow{f(x)} y \rightarrow h(f(x), g(x))$$

$$\frac{\partial L}{\partial w_{11}^{[1]}} = \frac{\partial L}{\partial y} \left[\frac{\partial y}{\partial O_1} \frac{\partial O_1}{\partial O_{11}} \frac{\partial O_{11}}{\partial w_{11}^{[1]}} + \frac{\partial y}{\partial O_2} \frac{\partial O_2}{\partial O_{11}} \frac{\partial O_{11}}{\partial w_{11}^{[1]}} \right]$$

$$\frac{\partial L}{\partial w_{11}^{[1]}}$$

Activation function

→ hidden layer] present at
→ output layer]

Gaussian Bernoulli Multinoulli / multinomial
distribution distribution distribution

$$f_x(x)$$

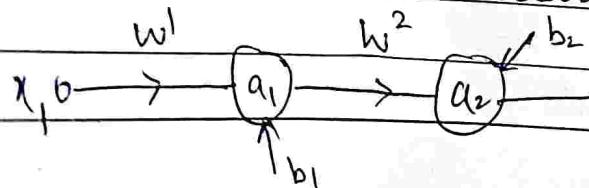
$$p(x=0)$$

$$P(x=0)$$

Softmax

i. Linear activation $\rightarrow f(z) = z$

Linear combination of linear activations
will result into linear plane.



$$z_1 = x, w^1 + b^1$$

$$a' = f(z_1) = x, w^1 + b^1$$

3/9

DL

$$z_2 = a^{[2]} w^{[2]} + b^{[2]} \rightarrow a^{[2]} = f(z_2) = z_2$$

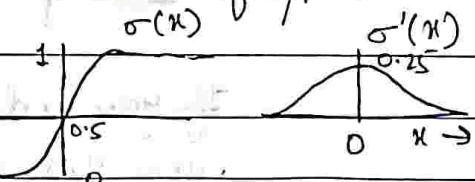
$$a^{[2]} = [x_1 w^{[1]} + b^{[1]}] \cdot w^{[2]} + b^{[2]}$$

$$a^{[2]} = x_1 w^{[1]} w^{[2]} + b^{[1]} w^{[2]} + b^{[2]} = n, w + b$$

9/9 Characteristics of activation function

1. Non linear : Universal approximation theorem
2. Differentiable : weight update / back propagation : exception = ReLU
3. Computationally inexpensive
4. Non saturating
5. Non centered zero : o/p of the activation should be centred to zero . Mean of o/p is 0.

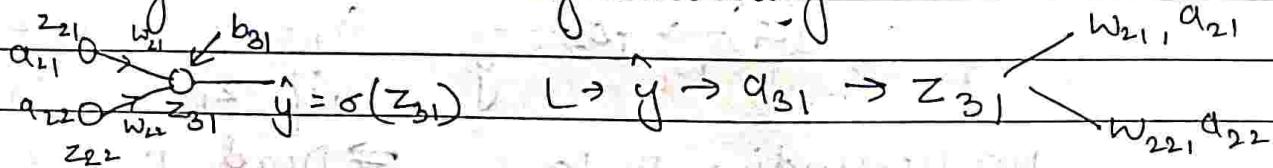
1. Sigmoid $\Rightarrow \sigma(x) = \frac{1}{1+e^{-x}}$



Advantages \rightarrow 1. Non linear & Differentiable 3. o/p is b/w [0, 1] consider o/p as probabilities

4. used for output layer

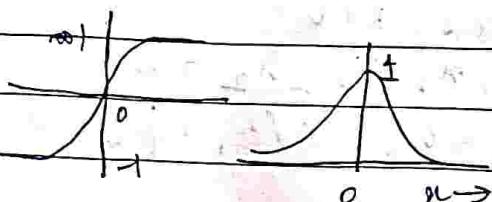
Disadvantages \rightarrow 1. Saturating



$$\frac{\partial L}{\partial w_{21}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_{31}} \frac{\partial z_{31}}{\partial w_{21} (\text{+ve})} \quad z_{31} = a_{21} w_{21} + a_{22} w_{22} + b_{31}$$

$$\frac{\partial L}{\partial w_{22}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_{31}} \frac{\partial z_{31}}{\partial w_{22} (\text{+ve})} \quad w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}}$$

2. Tanh =



DL lab

→ Grid search & Random search

No. of hidden layers

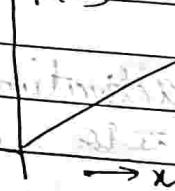
Learning rate

Activation func

Optimizer

10/9 • ReLU [Rectilinear Unit]

$$f(x)$$



Advantages

Non linear

$f(x) = \max(0, x) \rightarrow$ computationally inexpensive

$$f'(x)$$

$$1$$

If you add 2 linear ReLU
you get non linear

Disadvantage: Output is not centered to zero as there are no negative value



$$z_1 = x_1 w_1 + x_2 w_2$$

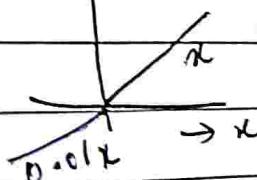
$$a_1 = \max(0, z_1)$$

Not responding to changes in input. \Rightarrow Dead neuron

z_1 can be negative if w is high and if bias has high -ve value.

Dis → dead neurons and no response to -ve values

• Leaky ReLU



$$f(x) = \begin{cases} 1, & x > 0 \\ 0.01x, & x \leq 0 \end{cases}$$

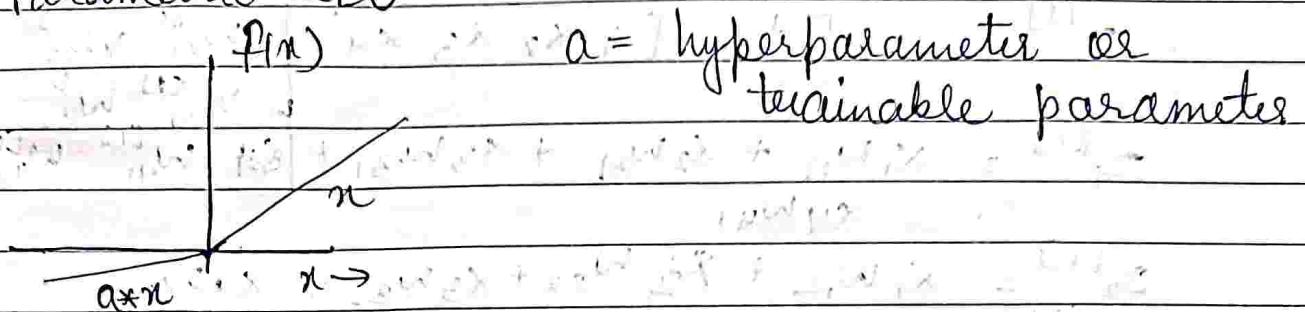
$$z_1 = w_1 x_1 + w_2 x_2 + b,$$

$$\text{if } z_1 < 0, a(z_1) \neq 0$$

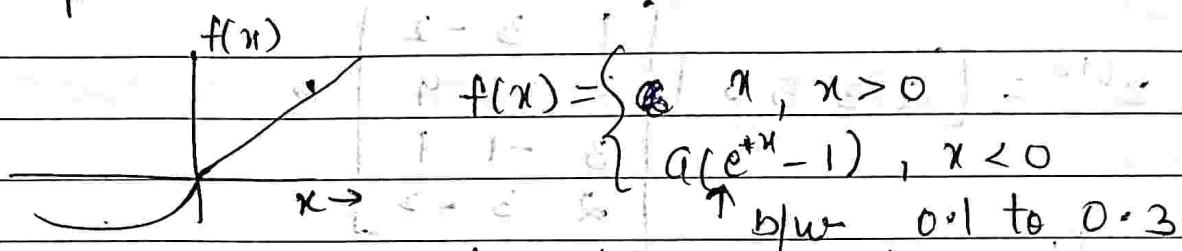
10/9

DL

- Parametric ReLU



- Exponential Linear Unit



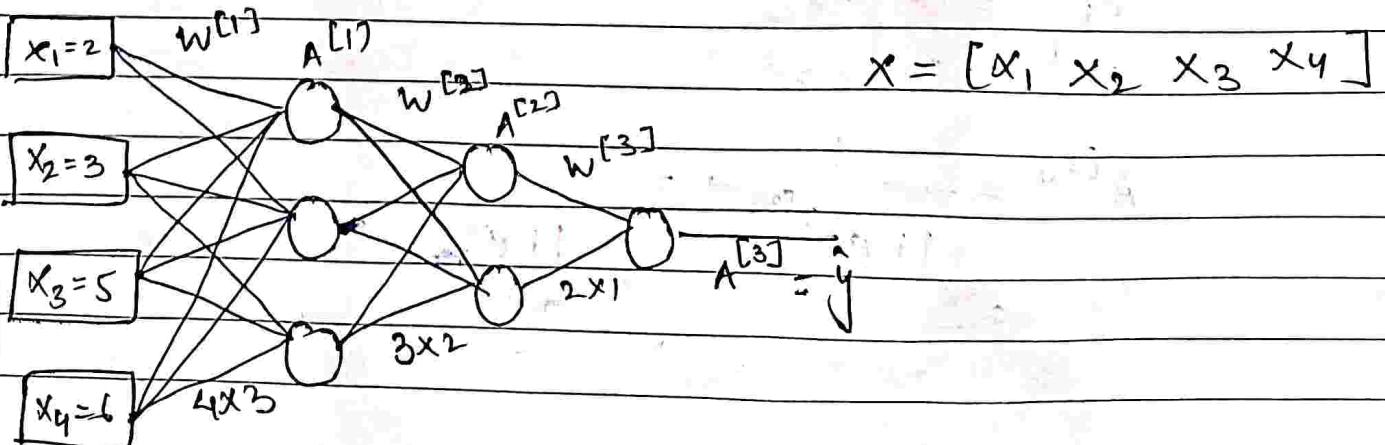
- Scalable Exponential Linear Unit

$$f(n) = \begin{cases} n, & n > 0 \\ a(e^{n-x} - 1), & n \leq 0 \end{cases}$$

12/9 Consider a 2×2 image applied at the input of the following MLP. Perform forward propagation and weight update. consider sigmoid activation on all neurons and all layers.

$$x = \begin{bmatrix} 2 & 3 \\ 5 & 6 \end{bmatrix}, w^{[1]} = \begin{bmatrix} 1 & 3 & -2 \\ 2 & -3 & 4 \\ 3 & -1 & 1 \end{bmatrix}, w^{[2]} = \begin{bmatrix} 1 & -2 \\ -3 & 4 \\ -1 & 1 \end{bmatrix}$$

$$w^{[3]} = \begin{bmatrix} 4 \\ 1 \end{bmatrix}, y = 1, \text{ consider all bias terms are } 0, \text{ and } \eta = 0.01$$



$$A^{[1]} = f^{[1]}(z^{[1]})$$

$$z^{[1]} = xw^{[1]} = [x_1 \ x_2 \ x_3 \ x_4]$$

$$z_1^{[1]} = x_1 w_{11} + x_2 w_{21} + x_3 w_{31} + x_4 w_{41}$$

$$\begin{matrix} w_{11} & w_{12} & w_{13} \\ w_{21}^{[1]} & w_{22}^{[1]} & w_{23}^{[1]} \\ w_{31}^{[1]} & w_{32}^{[1]} & w_{33}^{[1]} \\ w_{41}^{[1]} & w_{42}^{[1]} & w_{43}^{[1]} \end{matrix}$$

$$z_2^{[1]} = x_1 w_{12} + x_2 w_{22} + x_3 w_{32} + x_4 w_{42}$$

$$z_3^{[1]} = x_1 w_{13} + x_2 w_{23} + x_3 w_{33} + x_4 w_{43}$$

$$z^{[1]} = \begin{bmatrix} 2 & 3 & 5 & 6 \end{bmatrix} \begin{bmatrix} 1 & 3 & -2 \\ 2 & -3 & 4 \\ 3 & -1 & 1 \\ 2 & 3 & -2 \end{bmatrix}$$

$$z_1^{[1]} = 2+6+15+12 = 35$$

$$z_2^{[1]} = 6+(-9)-5+18 = 10 \Rightarrow z^{[1]} = \begin{bmatrix} 35 & 10 & 1 \end{bmatrix}$$

$$z_3^{[1]} = -4+12+5-12 = 1$$

$$A^{[1]} = \frac{1}{1+e^{-z^{[1]}}} = \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix}$$

$$A^{[1]} = \begin{bmatrix} e^{35} & e^{10} & e^1 \\ e^{35}+1 & e^{10}+1 & 1+e^1 \end{bmatrix} = \begin{bmatrix} 1 & 0.9 & 0.73 \end{bmatrix}$$

$$A^{[2]} = f^{[2]}(z^{[2]}) \text{ & } z^{[2]} = A^{[1]} \cdot w^{[2]}$$

$$z^{[2]} = \begin{bmatrix} 1 & 0.9 & 0.73 \end{bmatrix} \begin{bmatrix} 1 & -2 \\ -3 & 4 \\ -1 & 1 \end{bmatrix} \Rightarrow \begin{bmatrix} 1-0.7-0.73 & -2+3.6 \\ -3+3.6 & 4+0.73 \\ -1+1 & 1+0.73 \end{bmatrix} = \begin{bmatrix} 0.43 & 2.33 \end{bmatrix}$$

$$z_1^{[2]} = -0.43$$

$$z_2^{[2]} = 0.33$$

$$A^{[2]} = \frac{1}{1+e^{-z^{[2]}}} = \frac{1}{1+e^{-(0.43)}} = \frac{1}{1+e^{-2.43}} = \frac{1}{1+11.35888} = \frac{1}{11.27794} = 0.08091$$

$$= [0.08 \ 0.91]$$

$$A^{[3]} = f^{[3]}(z^{[3]}) \quad \text{and} \quad z^{[3]} = A^{[2]} w^{[3]}$$

$$z^{[3]} = [0.08 \ 0.91] \begin{bmatrix} 4 \\ 1 \end{bmatrix} = 0.32 + 0.91 \\ = 1.23$$

$$A^{[3]} = \frac{1}{1+e^{-z^{[3]}}} = \frac{1}{1+e^{-1.23}} = 0.77$$

$$\hat{y} = 0.77 \quad \text{and } y = 1 \quad \Rightarrow \hat{y} - y = 0.77 - 1 \\ = -0.23$$

$$L = -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$$

$$w_n^{[3]} = w_0^{[3]} - \eta \frac{\partial L}{\partial w_0^{[3]}}$$

$$\frac{\partial L}{\partial w_0^{[3]}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z^{[3]}} \frac{\partial z^{[3]}}{\partial w_0^{[3]}}$$

$$\frac{\partial L}{\partial \hat{y}} = \frac{-y}{\hat{y}} - \frac{(1-\hat{y})(-1)}{(1-\hat{y})} = \frac{-y}{\hat{y}} + \frac{1-\hat{y}}{1-\hat{y}}$$

$$\frac{\partial \hat{y}}{\partial z^{[3]}} = \frac{\sigma'(z^{[3]})}{\sigma(z^{[3]})} \Rightarrow \hat{y} = \sigma(z^{[3]}) = \frac{1}{1+e^{-z^{[3]}}}$$

$$\frac{\partial \hat{y}}{\partial z^{[3]}} = \sigma(z^{[3]}) (1 - \sigma(z^{[3]}))$$

$$\frac{\partial z^{[3]}}{\partial w_0^{[3]}} = \hat{y}(1-\hat{y})$$

$$\frac{\partial z^{[3]}}{\partial w_0^{[3]}}$$

$$\frac{\partial L}{\partial w_0^{[3]}} = - \left(\frac{y}{\hat{y}} - \frac{1-\hat{y}}{1-\hat{y}} \right) = \hat{y}(1-\hat{y}) A^{[2]}$$

$$= - [y(1-\hat{y}) - \hat{y}(1-y)] A^{[2]}$$

$$= [\hat{y}(1-y) - \hat{y}(1-\hat{y})] A^{[2]}$$

$$= [\hat{y} - y\hat{y} - \hat{y} + y\hat{y}] A^{[2]} = (\hat{y} - y) A^{[2]}$$

$$w_n^{[3]} = \begin{bmatrix} 4 \\ 1 \end{bmatrix} - 0.01(\hat{y} - y) A^{[2]}$$

But Cap

DL by Ankit Das Book for numericals

classmate

Date _____

Page _____

$$w_n^{[2]} = w_0^{[2]} - \eta \frac{\partial L}{\partial w_0^{[2]}}$$

$$\frac{\partial L}{\partial w_0^{[2]}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial A^{[2]}} \frac{\partial A^{[2]}}{\partial w_0^{[2]}}$$

17 | 9 | 24

Weight initialization

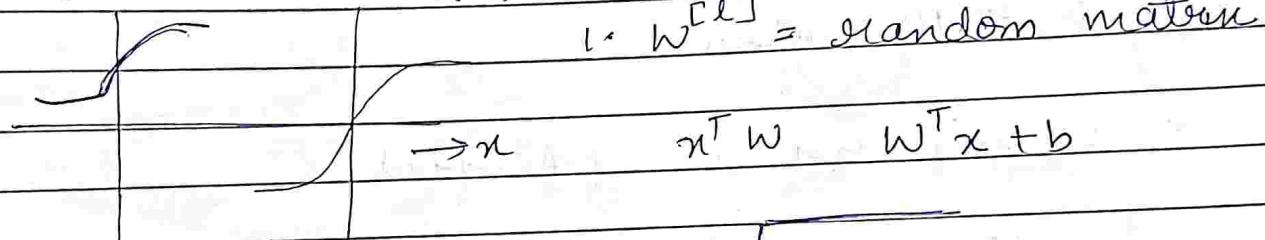
$w^{[l]}$ initial value of weight matrix

0 is not the correct way to initialize

He initialization & Xavier initialization
of w)

tanh(x) Normal

1. $w^{[l]} = \text{random matrix}$



He 1) $w^{[l]} = \text{Normal random matrix} *$

size $[l-1]^2$

Xavier 2) $w^{[l]} = \text{Uniform random matrix} *$

size $[l-1]$

→ Gradient descent (GD)

$$w_n = w_0 - \eta \frac{\partial L}{\partial w_0}$$

→ Stochastic GD → epochs = $\frac{n}{B}$ for no. of epochs:

no of weight update
= epochs × data point

Shuffle

for data point

predict using forward prop.

loss

$$\frac{\partial L}{\partial w}$$

update weights

end

Take lot of time
Not guaranteed
to reach somewhere
or converge to
near solution

- Batch GD : no. of times weight updates = epochs
Faster
- Mini batch : Data is randomly divided into patches
No. of batches × No. of epochs
- Adaptive learning rate

$$w = w - \eta \frac{\partial L}{\partial w} \text{ (without momentum)}$$

$$m w_i = \beta_1^i m w_{i-1} + (1 - \beta_1) \frac{\partial L}{\partial w} \Rightarrow w_i = w_{i-1} - \eta_m w_i$$

↑
(With momentum)

* Adaptive Gradient (Adagrad)

$$\alpha_i = \frac{\alpha_{i-1}}{\sqrt{\delta + \sum_{k=1}^{i-1} \partial w_k^2}} \quad \delta \approx 10^{-5} \text{ to } 10^{-1}$$

• Optimizer

$$\alpha_i = \frac{\alpha}{\sqrt{S + \sum_{k=1}^{i-1} \partial W_k^2}}$$

↑
fixed 10^5 to 10^7

$$w_i = w_{i-1} - \alpha_i \partial w$$

- Root mean square propagation

More weightage is given to current gradient.

$$S_{Wi} = \beta_2 S_{Wi-1} + (1-\beta_2) \Delta w_{i-1}$$

\leftarrow initial LR

$$\alpha_i = \frac{\alpha}{\sqrt{S_{ii}}}$$

- Gradient descent with momentum

$$m w_i = \beta_i m w_{i-1} + (1 - \beta_i) \frac{\partial L}{\partial w} \Rightarrow m w_i = B$$

$$w_i = w_{i-1} - \alpha_i \partial w$$

$$w_i = w_{i-1} - \alpha m w_v$$

- ### • Adam optimizer (AO)

$$m_{WI} = \beta_1 m_{W1} + (1 - \beta_1) m_{W2} \quad] \text{GD}_M^{\text{Wett}}$$

$$w_i = w_{i-1} - \alpha m w_i$$

$$A \left\{ \begin{array}{l} w_i = w_{i-1} + \alpha_i \Delta w_i \quad \& \quad \alpha_i = \sqrt{8 + s_i} \\ s_i = \beta_2 s_{i-1} + (1-\beta_2) \Delta w_i \\ w_i = w_{i-1} - \alpha_i \Delta w \\ m w_i^{\text{corrected}} = \frac{m w_i}{(1-\beta_2)} \end{array} \right.$$

$$x_i^{\text{corrected}} = \frac{(1-\beta_i) x_i}{\sqrt{s + s_{w,i}^{\text{corrected}}}}$$

$$S_{Wi}^{\text{corrected}} = \frac{S_{Wi-1}}{1-\beta_i}$$

$$w_i = w_{i-1} - \alpha_i u_i$$

Q For a NN, initial LR is 0.02 and squared deviate of weight in 1st iteration is 6. If SD of w ↑ by 5% in each iteration, find the LR of adagrad after 10 iterations. Assume $\delta = 10^{-5}$

$$x_i = \frac{\alpha}{\sqrt{\delta + \partial w_i^2}}$$

$$x_{i_2} = \frac{0.02}{\sqrt{10^5 + 6}} = \frac{0.02}{\sqrt{\frac{1}{10^5} + 6}} = \frac{0.02}{\sqrt{1 + 6 \times 10^5}} = 174.5973$$

$$= 0.008$$

$$\alpha_i =$$

Iteration	∂w^2	$\sum \partial w^2$	α_i
1	6	6	0.008
2	$6 \cdot 3$	$6 + 6 \cdot 3 = 12 \cdot 3$	0.006
3	$12 \cdot 3 \times 0.05 + 6 = 6.62$	18.92	

$1 \rightarrow 6 \rightarrow 6 \rightarrow 0.08$
 $2 \rightarrow 6 + \frac{6 \times 0.05}{6 \cdot 3} \rightarrow 6 + 6 \cdot 3 = 12 \cdot 3 \rightarrow$
 $3 \rightarrow \underbrace{12 \cdot 3 + 12 \cdot 3 \times 0.05}_{a} \rightarrow 12 \cdot 3 + a = \rightarrow 25.219$
 12.915

Iteration	∂w^2	$\sum \partial w^2$	α_i
1	6	6	0.008
2	$6 + 6 \times 0.05 = 6 \cdot 3$	$6 + 6 \cdot 3 = 12 \cdot 3$	0.006
3	$12 \cdot 3 \times 0.05 + 6 = 6.62$	$12 \cdot 3 + 6.62 = 18.92$	
4	$18.92 \times 0.05 + 6 = 6.94$	$18.92 + 6.94 = 25.86$	
5	$25.86 \times 0.05 + 6 = 7.293$	$7.293 + 25.86 = 33.153$	
6	$33.153 \times 0.05 + 6 = 7.65$	$7.65 + 33.153 = 40.8$	
7	$40.8 \times 0.05 + 6 = 8.04$	$8.04 + 40.8 = 48.84$	
8	$48.84 \times 0.05 + 6 = 8.44$	$8.44 + 48.84 = 57.282$	
9	$57.282 \times 0.05 + 6 = 8.864$	$8.864 + 57.282 = 66.14$	
10	$66.14 \times 0.05 + 6 = 9.307$	$9.307 + 66.14 = 75.447$	

• Regularization

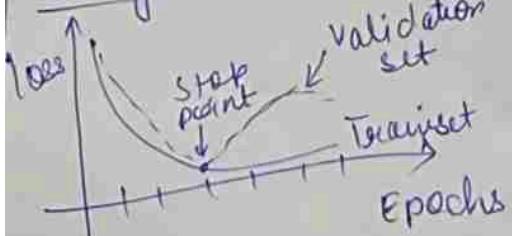
Type $\rightarrow L_1 \& L_2$

$$L_1 \rightarrow L'(w, b) = L(w, b) + \lambda w, \leftarrow \text{sum of weights}$$

$$L_2 \rightarrow L''(w, b) = L(w, b) + \lambda w^2, \leftarrow \text{squared sum of weights}$$

Dropout

Early stopping



Data augmentation

- related to images
- resize
- rotate
- shift

Normalization

$$x_n = x_i - \mu x_i$$

→ Normalize data by making it zero mean & 1 variance

28/9 Unit 2 - CNN

0 - black, 255 black

Convolution : For input image $n \times n$ & filter $m \times m$, dimension of the output matrix is $(n-m+1) \times (n-m+1)$

Eq: $\begin{matrix} 6 \times 6 \\ \text{Image} \end{matrix} \quad \begin{matrix} 3 \times 3 \\ \text{Filter} \end{matrix} \Rightarrow \text{Output} \rightarrow \begin{matrix} (6-3+1) \times (6-3+1) \\ 4 \times 4 \end{matrix}$

Padding & stride 1/10

size of padding depends on size of filter

$$n + 2p - f + 1 \rightarrow \text{for padding} \\ = n \Rightarrow p = \frac{f-1}{2}$$

With Stride $\rightarrow \frac{n-f+1}{s}$

without padding, $s \Rightarrow \frac{n-f+1}{s} \times \frac{n-f+1}{s}$

padding \rightarrow valid \rightarrow no padding
Same \rightarrow with padding o/p shape \Rightarrow i/p shape

with padding, $s \Rightarrow \frac{n+2p-f+1}{s} \times \frac{n+2p-f+1}{s}$

Q. I/p image $\Rightarrow 28 \times 28$

layer 1 } 32 filter 3×3
2 } no stride
3 } valid padding

Layer 4 \rightarrow flatten

Layers 5 \rightarrow 128 neuron, Fully connected
6 \rightarrow 10 neurons fully connected

Find o/p shape

$$\rightarrow \text{Layer 1} = 28 \times 28 \times 32 \quad (28-3+1 = 26)$$

$$\text{Layer 2} = 26 \times 26 \times 32 \quad (26-3+1 = 24)$$

$$\text{Layer 3} = 24 \times 24 \times 32 \quad (24-3+1 = 22)$$

$$\text{Layer 4} = 22 \times 22 \times 32$$

$$\text{Layer 5} = 128$$

$$\text{Layer 6} = 10$$

a. Layer 1 } 32 filter 3×2
2 } Stride = 2
3 } same padding

$$1 \rightarrow 14 \times 14 \times 32$$

$$2 \rightarrow 7 \times 7 \times 32$$

$$3 \rightarrow 4 \times 4 \times 32$$

$$4 \rightarrow 4 \times 4 \times 32$$

$$5 \rightarrow 128$$

$$6 \rightarrow 10$$

$$\text{Layer 1} = \frac{14 \times 14 \times 32}{28 \times 28} = \frac{(28+2-3+1)}{2} = 7$$

$$\text{Layer 2} = \left(\frac{14+2-3+1}{2} \right) = \left(\frac{13+1}{2} \right) = 6 \cdot 5 + 1 = 7$$

$$\text{Layer 3} = \left(\frac{7+2-3+1}{2} \right) = \left(\frac{6+1}{2} \right) = 3 + 1 = 4 \times 4 \times 32$$

Pooling

- 1. Image storage
- 2. Translation invariance

$$\text{Layer 2} = 24 \times 24 \times 32 \quad (24 - 3 + 1 = 24)$$

$$\text{Layer 3} = 22 \times 22 \times 32 \quad (24 - 3 + 1 = 22)$$

$$\text{Layer 4} = 20 \times 20 \times 32$$

$$\text{Layer 5} = 128$$

$$\text{Layer 6} = 10$$

1. Layer 1 } \rightarrow 32 filter 3×2
 2 } $\text{Stride} = 2$
 3 } same padding

$$1 \rightarrow 14 \times 14 \times 32$$

$$2 \rightarrow 7 \times 7 \times 32$$

$$3 \rightarrow 4 \times 4 \times 32$$

$$4 \rightarrow 4 \times 4 \times 32$$

$$5 \rightarrow 128$$

$$6 \rightarrow 10$$

• Pooling

1. Image storage

2. Translation invariance

12/11 Popular CNN Architectures

1. LeNet-5

3. VGG-16

2. AlexNet

4. ResNet

$$\text{Layer 1} = 14 \times 14 \times 32$$

$$\text{Layer 2} = \left(\frac{14+2-3}{2} + 1 \right) = 13 \cdot 5 + 1 = 14$$

$$= \left(\frac{13}{2} + 1 \right) \cdot 5 + 1 = 7 \cdot 5 + 1 = 36$$

$$\rightarrow 7 \times 7 \times 32$$

$$\text{Layer 3} = \left(\frac{7+2-3}{2} + 1 \right) = \left(\frac{6}{2} + 1 \right) = 3 + 1 = 4 \times 4 \times 32$$

LiNet-5 - for recognizing handwritten and machine-printed characters



Output shapes for =
 $(32-5+1) \times (32-5+1)$

~~Architecture 2. AlexNet~~

3. VGG-16

Traditional vs Transfer learning

94/11 Normalization

$$x \rightarrow \text{input}, \mu_x = \text{mean}, \sigma_x = \text{variance/SD}$$

$$\hat{x}_N = \frac{x - \mu_x}{\sigma_x} \quad x_{\min} \text{ & } x_{\max} \Rightarrow \hat{x}_N = \frac{x_i - \min}{\max - \min}$$

mini-batch: $\square \square \square \square$ i/p & hidden layers

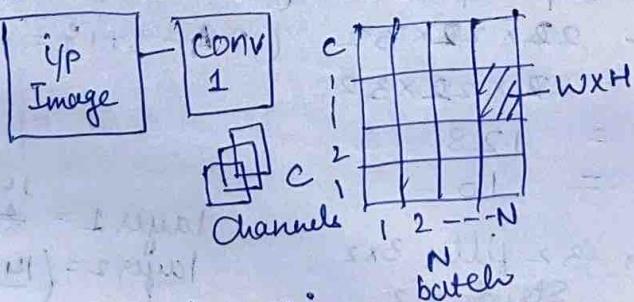
Batch normalization \rightarrow i/p

Classifier

$x \rightarrow \text{mini batch data}, \mu_B = \text{mean } x \quad \sigma_B = \text{std of } x$

$$\hat{x}_B = \frac{x - \mu_B}{\sqrt{\epsilon + \sigma_B^2}}$$

i/p Image $\xrightarrow{\text{conv} 1}$ channels



$$F_1 \quad F_2 \quad F_3 \quad F_4$$

$$\begin{matrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \\ 7 & 8 \\ 9 & 10 \\ 11 & 12 \\ 13 & 14 \\ 15 & 16 \end{matrix}$$

$$\xrightarrow{\text{I1}} \begin{matrix} 1 & 5 \\ 9 & 14 \end{matrix}$$

$$\xrightarrow{\text{I2}} \begin{matrix} 2 & 6 \\ 10 & 14 \end{matrix}$$

$$\hat{x} = \frac{x - \mu_c}{\sqrt{\sigma_c^2 + \epsilon}}$$

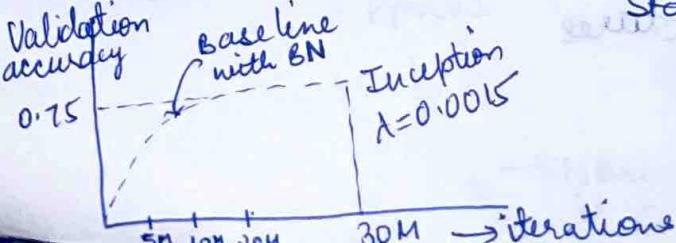
Reparameterization: learn mean & variance of a class through data channel

$$\rightarrow \mu_c = \frac{1}{N \times W \times H} \sum_{i=1}^N \sum_{j=1}^W \sum_{k=1}^H x_{ijk} \quad \text{index}$$

$$\sigma_c^2 = \frac{1}{N \times W \times H} \sum_{i=1}^N \sum_{j=1}^W \sum_{k=1}^H (x_{ijk} - \mu_c)^2 \quad \text{batch (row, col)}$$

$$\rightarrow \gamma + \beta$$

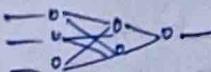
$$y_i = \gamma \hat{x}_i + \beta \quad \gamma \text{ & } \beta \text{ are trainable}$$



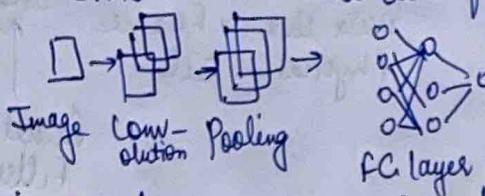
18/11 RNN

DL

ANN

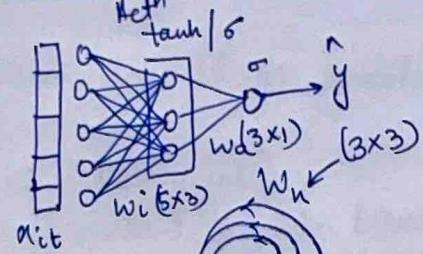


DL \rightarrow MLP \rightarrow size of input is fixed
 \rightarrow CNN Not suitable for image

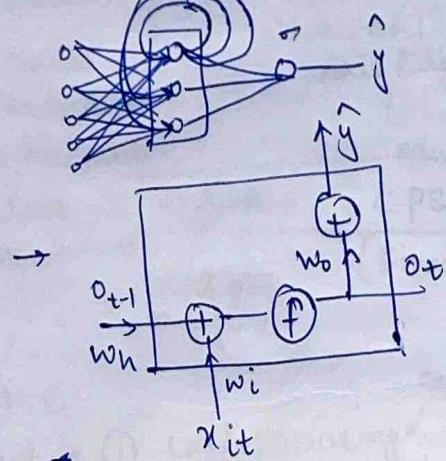
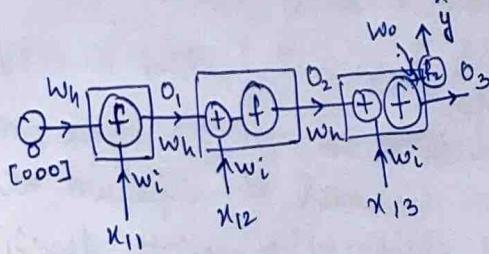


RNN is used for sequential data.
 Which specifies the timestamps as well.

$$\begin{array}{cccccc} x_1 & x_{12} & x_3 \\ S_1 & v_1 & v_2 & v_3 & y & 1 \\ x_2 & v_1 & v_2 & v_3 & v_4 & v_5 & 0 \\ S_3 & x_{21} & x_{22} & x_{23} & x_{24} & x_{25} \end{array}$$



$$f_2([f_1(x_{it} \cdot w_i)] w_o) \rightarrow \hat{y}$$



$$[f_1(x_{it} \cdot w_i)] = o_t$$

$$\hat{y} = [f[(o_{t-1} \cdot w_h) + w_i x_{it}] w_o]$$

$$\begin{aligned} w_i &= w_i - \eta \frac{\partial L}{\partial w_i} \\ w_o &= w_o - \eta \frac{\partial L}{\partial w_o} \\ w_h &= w_h - \eta \frac{\partial L}{\partial w_h} \end{aligned} \quad \begin{matrix} \text{Back-} \\ \text{Prop-} \\ \text{Through} \\ \text{Time} \end{matrix}$$

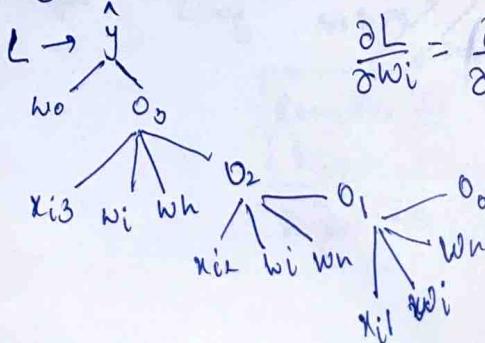
$$o_0 \rightarrow o_1 \rightarrow o_2 \rightarrow o_3$$

$$\frac{\partial L}{\partial w_o} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_o}$$

$$\frac{\partial L}{\partial \hat{y}} = - \left[\frac{y}{\hat{y}} - \frac{1-y}{1-\hat{y}} \right]$$

$$L = -y \log \hat{y} - (1-y) \log(1-\hat{y})$$

$$\frac{\partial \sigma(o_3 w_o)}{\partial w_o} = \sigma(o_3 w_o) (1 - \sigma(o_3 w_o)) \cdot w_o \left[\frac{\partial \sigma(z)}{\partial z} = \sigma(z) (1 - \sigma(z)) \right]$$

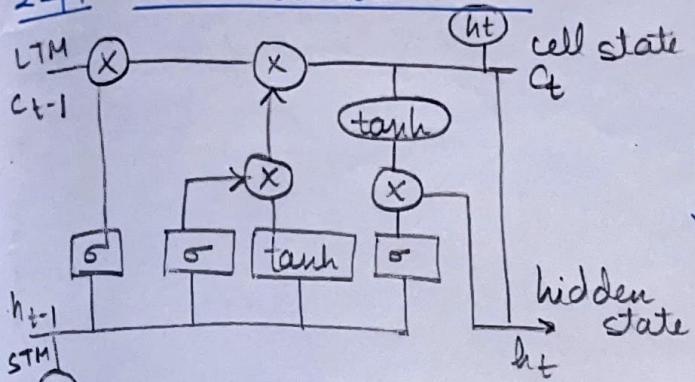


$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial w_i} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_2} \frac{\partial o_2}{\partial w_i} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_1} \frac{\partial o_1}{\partial w_i}$$

$$\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial w_i} = \sigma(o_3 w_o) (1 - \sigma(o_3 w_o)) \cdot w_o \cdot \sigma'(o_3 w_o) \cdot w_o$$

$$\frac{\partial L}{\partial w_h} = \frac{\partial L}{\partial g} \frac{\partial g}{\partial o_3} \frac{\partial o_3}{\partial w_h} + \frac{\partial L}{\partial g} \frac{\partial g}{\partial o_2} \frac{\partial o_2}{\partial w_h} + \frac{\partial L}{\partial g} \frac{\partial g}{\partial o_1} \frac{\partial o_1}{\partial w_h}$$

22/11 LSTM Architecture



long term memory is cell state

Short term memory is in hidden state

Content - long term context
short term

Inputs to LSTM \rightarrow Previous cell state $- c_{t-1}$

Previous hidden state $- h_{t-1}$

i/p at current time index $- x_t$

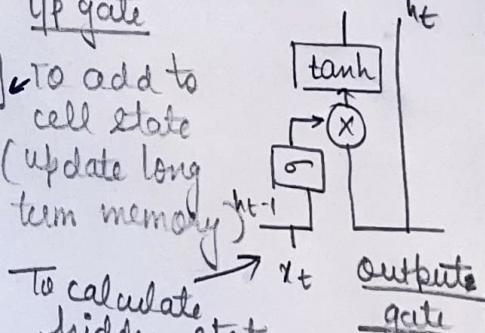
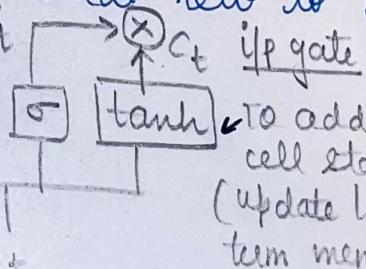
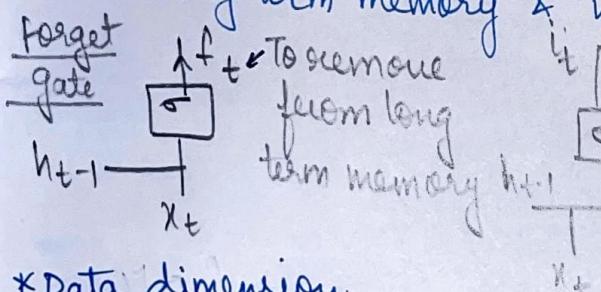
Outputs generated by LSTM \rightarrow current hidden state $- h_t$

current cell state $- c_t$

Process of LSTM \rightarrow 1. Update cell state i.e. long term content
2. from c_{t-1} to c_t

Calculate current hidden state i.e. short term content $\rightarrow h_{t-1}$ to h_t

* c_t \rightarrow Based on current i/p decide what to remove from long term memory & what new to be added.



Data dimension

$x_t \rightarrow$ vector

$h_t \& c_t \rightarrow$ also vector, same size

Data \rightarrow Sent 1 $t=0$ $t=1$ $t=2$ $t=3$
 Sent 2
 Sent 3

x_t one word vector

Dimension of $x_t \& c_t, h_t$ need not be same

* output of forget gate is f_t

* Other o/p vectors in LSTM are i_t, o_t, c_t
 \uparrow are of same dimension as that of $c_t \& h_t$

\otimes element wise multiplication

\oplus element wise addition

\tanh element wise tanh

$$\text{Eg: } [1, 2, 3] \otimes [1, 4, 7] = [1 \times 1, 2 \times 4, 3 \times 7]$$

$$\tanh(1, 4, 3) = [\tanh(1), \tanh(4), \tanh(3)]$$