

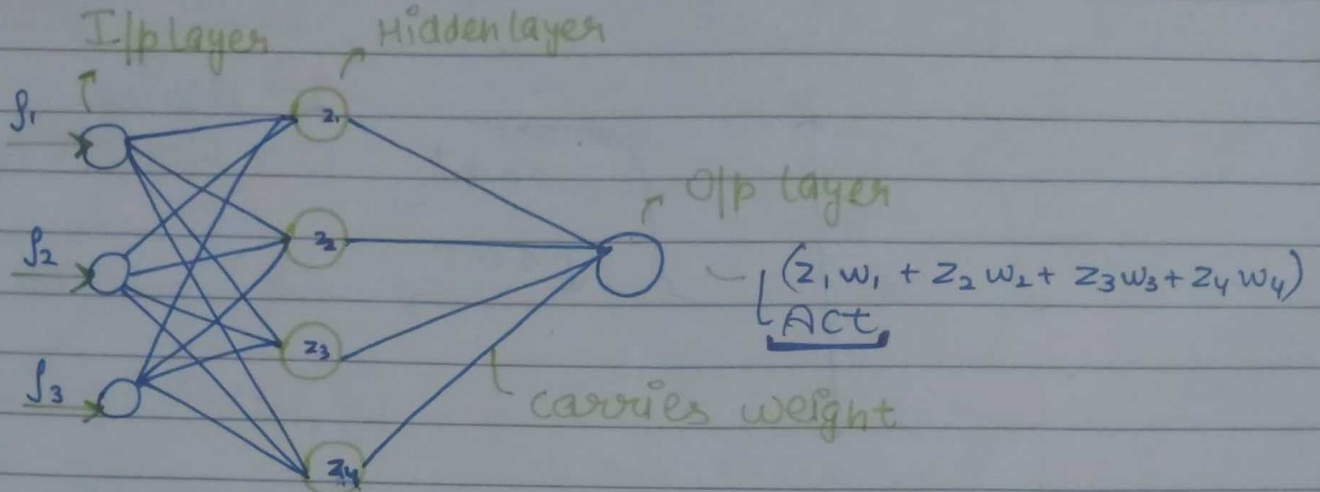
# NN

5

29/09/22

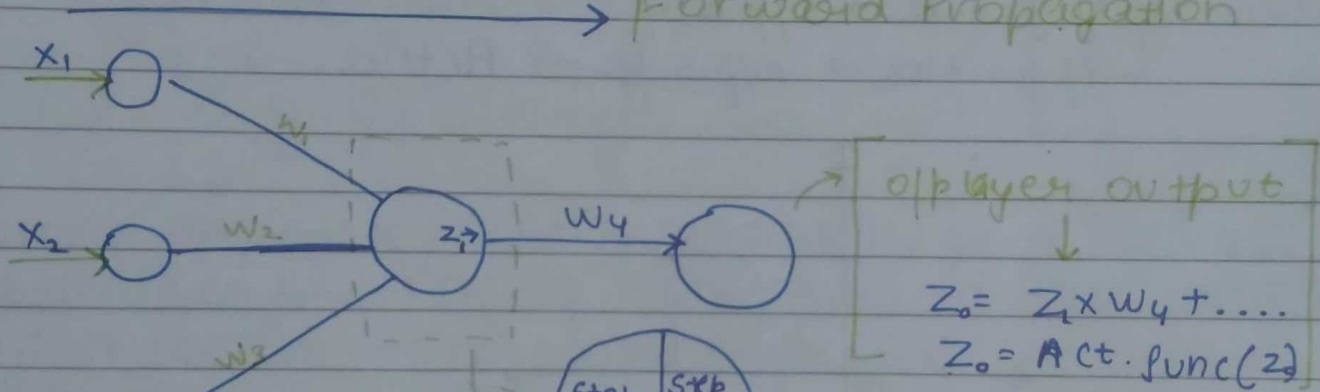
Nodes

• eyes, skin, ear → I/p layer | 'Neurons' → Hidden layer



↓ Simplified version to understand Better

→ 'Forward Propagation'



similar steps

Step 1 →  $y = w_1 x_1 + w_2 x_2 + w_3 x_3 + \text{Bias}$

Step 2 →  $Z_i = \text{Act. func}(y)$

ex: sigmoid

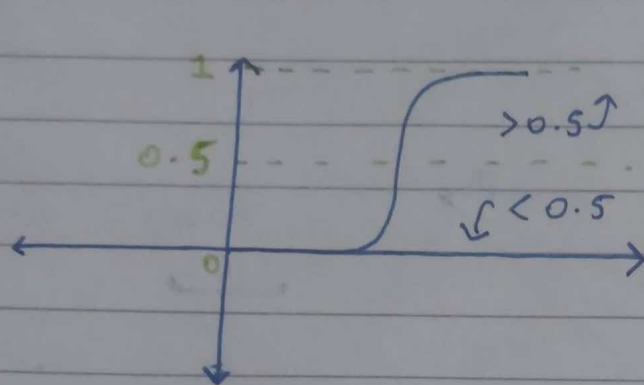
$\frac{1}{1+e^{-y}}$

The value → active or not

Act func → To activate which Neuron will forward the Response (like → hot piece on left hand will activate Neurons of left hand and not Right hand)

- 1) Sigmoid AF  $\rightarrow$  B/w 0 to 1 (To activate and act as output)

$\rightarrow \frac{1}{1+e^{-y}}$  [Set as Threshold]

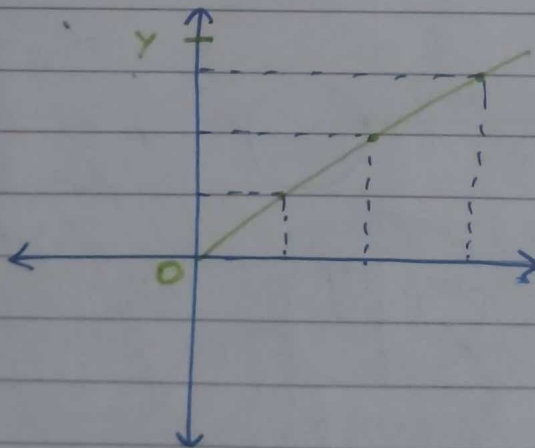


1  $\rightarrow$  Active

0  $\rightarrow$  Active-Non.

- 2) ReLU AF  $\rightarrow \max(y, 0)$  [No -ve values]

$\begin{cases} \text{if } y = -ve \rightarrow o/p = 0 \rightarrow \text{Non-Active} \\ \text{if } y = +ve \rightarrow o/p = y \rightarrow \text{Active.} \end{cases}$



X ——— X ——— X ——— X ——— X ———

★ NN TRAINING  $\rightarrow$  Back Prop + Loss func

Play  
2h

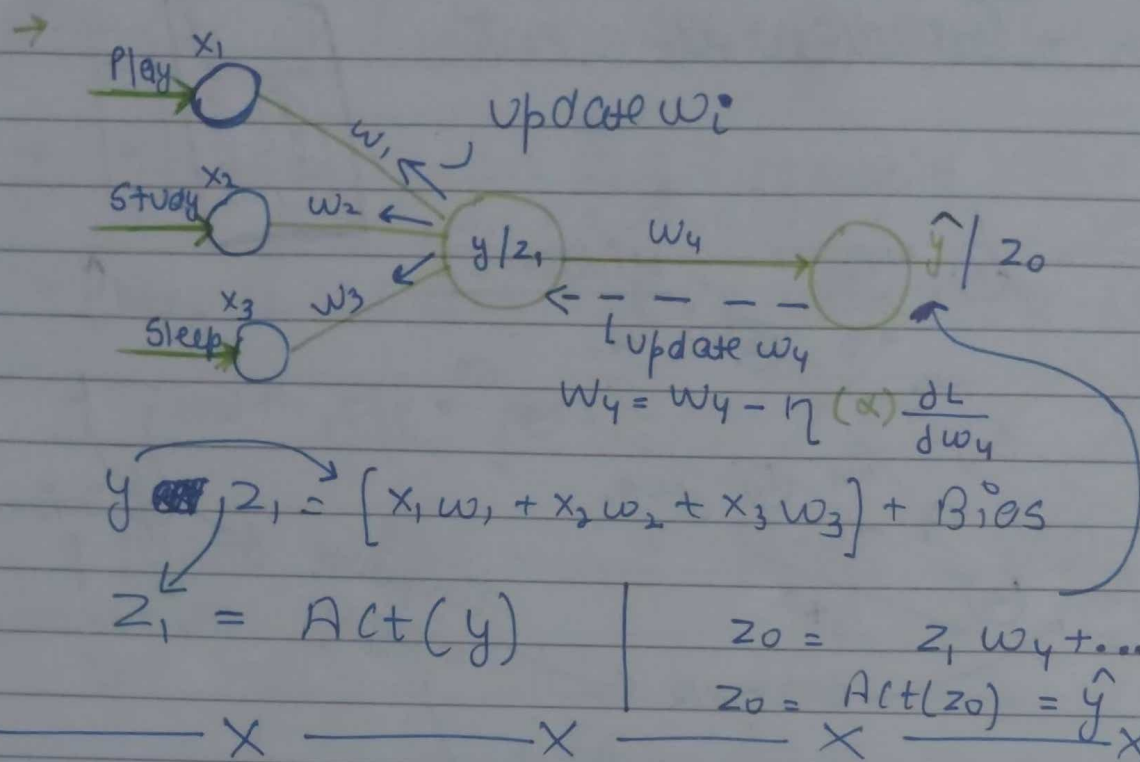
Study  
4h

Sleep  
8h

O/p

1 - pass  $\rightarrow y$

- Forward Prop  $\rightarrow$  Act func
- Backward Prop  $\rightarrow$  optimizer + loss func



★ To check the accuracy of Result  $\rightarrow$  compare  $y$  and  $\hat{y}$

By Loss func  $\rightarrow$  optimize  $\rightarrow$  minimize loss

Loss func  $\rightarrow \sum (y - \hat{y})^2$  (Normal loss func) (1)

• if wrong propagate / Classify / error huge  $\rightarrow$  Define optimizer func  $\rightarrow$  Define loss func  $\rightarrow$  Back track  $\rightarrow$  Adjust weight  $\rightarrow$  Forward Propagate.  $\uparrow$  cycle repeat

$\uparrow \uparrow \uparrow$  until epochs.

• updated  $w_4 = w_4 - \alpha \cdot \frac{dL}{dw_4}$

$\downarrow$  small ' $\alpha$ '  $\rightarrow$  To reach global minima

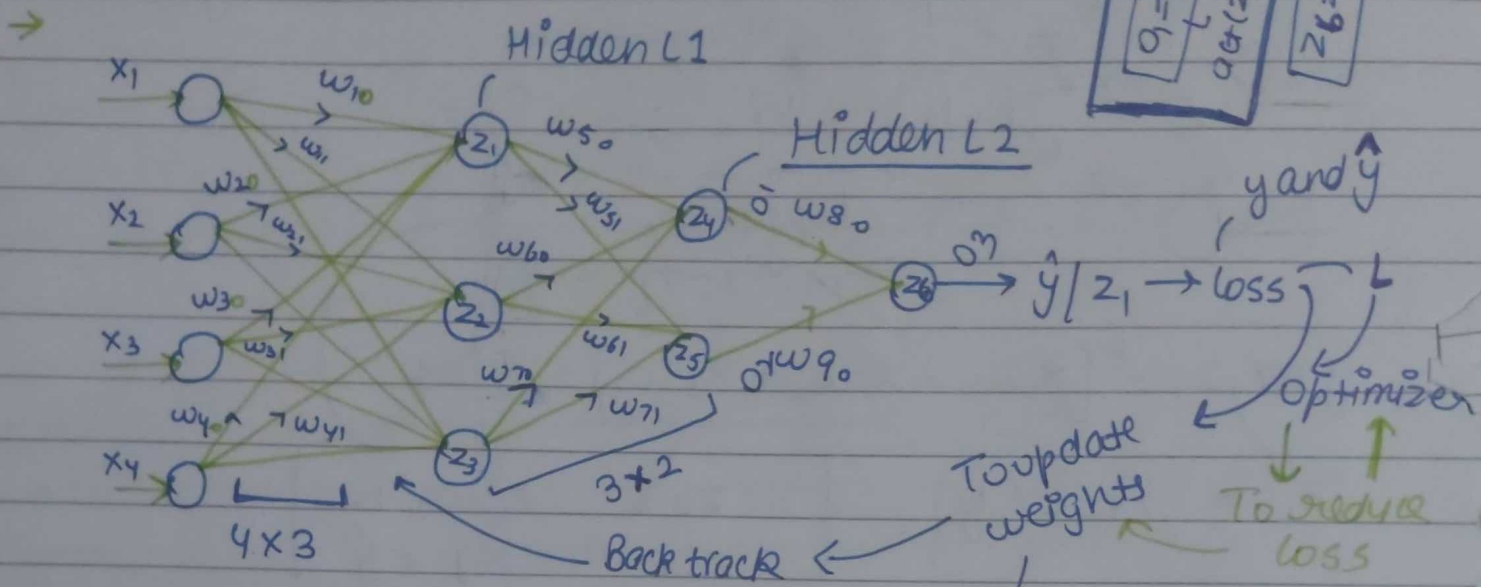
• updated  $w_3 = w_3 - \alpha \cdot \frac{dL}{dw_3}$



$w_{50} = w_{51}$   
singular  
But  
Diff

Output  $\rightarrow Z_i \rightarrow \text{Act}(y) \cdot \text{Weight Assign}$

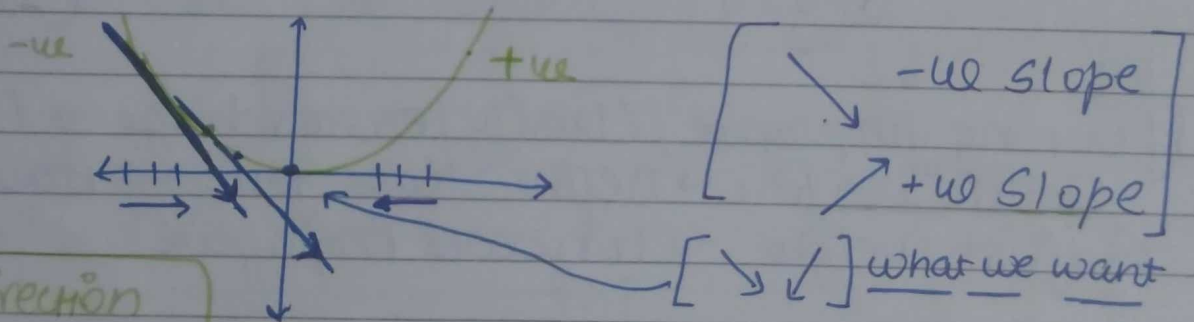
★ Optimizer  $\rightarrow$  Gradient descent



$\rightarrow \text{Updated } w_i^o = w_i^o - \alpha \cdot \frac{\partial L}{\partial w_i^o} \quad (1)$

slope

GD  $\rightarrow$  update weights  $\rightarrow$  To reach global optimal sol<sup>n</sup>



Slope  $\rightarrow$  Direction  
L.R  $\rightarrow$  Movement

L.R determine the movement and new updated weight.

if Slope = -ve  $\rightarrow$  updated  $w_i^o > w_i^o$  (from ①)  $\rightarrow$  Move Down  $(\rightarrow)$

if Slope = +ve  $\rightarrow$  updated  $w_i^o < w_i^o$  (from ①)  $\rightarrow$  Move Down  $(\leftarrow)$

- After some epoch  $\rightarrow$  If loss  $\downarrow \rightarrow \checkmark$
- After some epoch  $\rightarrow$  If loss  $\uparrow$  or does not  $\downarrow \rightarrow \times \rightarrow$  prob with  $\alpha$ .

## Chain Rule in Back Propagation

- $\rightarrow z_1 \cdot w_{50} \rightarrow \text{Output } O_1 \rightarrow z_1 \cdot w_{51}$  (  $w_{50}$  similar same to  $w_{51}$  )  
 $+ z_2 \cdot w_{60} + z_3 \cdot w_{60}$  Ka whole act from
- $\rightarrow z_4 \cdot w_{80} \rightarrow \text{Output } O_3$  (  $O_3$  ~~act~~ weight  $\rightarrow$  act )  $\in$  of  $z_1, z_2$  and  $z_3$  and  $z_4$  and  $z_5$
- $\downarrow$  Loss func  $\rightarrow L$

updated weight  $w_{80} \rightarrow \text{old } w_{80} - \alpha \frac{\partial L}{\partial w_{80}}$  [  $w_{80}$  updates or affects the Result  $z_6$  ]

$$\frac{\partial L}{\partial w_{80}} = \frac{\partial L}{\partial z_6} \cdot \frac{\partial z_6}{\partial w_{80}}$$

$$\frac{\partial L}{\partial w_{90}} = \frac{\partial L}{\partial z_6} \cdot \frac{\partial z_6}{\partial w_{90}}$$

To update  $w$  of  $z_1$   
 $\hookrightarrow \frac{\partial L}{\partial w_{50}} + \frac{\partial L}{\partial w_{51}}$

$$\frac{\partial L}{\partial w_{50}} = \frac{\partial L}{\partial z_6} \cdot \frac{\partial z_6}{\partial z_4} \cdot \frac{\partial z_4}{\partial w_{50}} \rightarrow \text{Chain Rule (Piche sai)}$$

$$\frac{\partial L}{\partial w_{51}} = \frac{\partial L}{\partial z_6} \cdot \frac{\partial z_6}{\partial z_5} \cdot \frac{\partial z_5}{\partial w_{51}} \rightarrow \text{Chain Rule (Piche sai)}$$

$\frac{\partial L}{\partial w_{50}} + \frac{\partial L}{\partial w_{51}} = \text{up/slope of } z_1$

$$\frac{\partial L}{\partial w_{61}} = \frac{\partial L}{\partial z_6} \cdot \frac{\partial z_6}{\partial z_4} \cdot \frac{\partial z_4}{\partial w_{61}}$$

$$\frac{\partial L}{\partial w_{71}} = \frac{\partial L}{\partial z_6} \cdot \frac{\partial z_6}{\partial z_5} \cdot \frac{\partial z_5}{\partial w_{71}}$$

~~$$\frac{\partial L}{\partial w_{10}} = \frac{\partial L}{\partial z_6} \cdot \frac{\partial z_6}{\partial z_4} \cdot \frac{\partial z_4}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_{10}}$$~~

~~$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z_6} \cdot \frac{\partial z_6}{\partial z_5} \cdot \frac{\partial z_5}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_{11}}$$~~

$U_1 + U_2 + U_3 \dots U_n$  ( $n = \text{no. of paths}$ )  
 from  $z_1$  there is 2 paths

~~$$\frac{\partial L}{\partial z_6} \cdot \frac{\partial z_6}{\partial z_4}$$~~

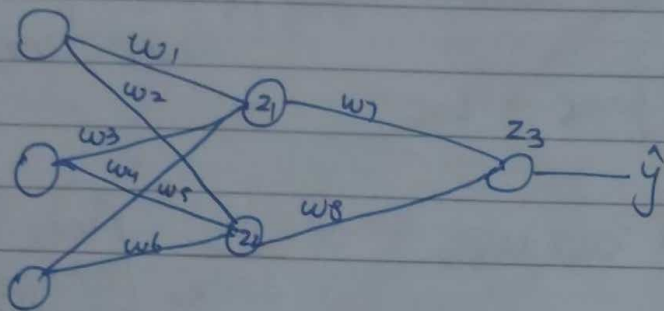


chain Rule  $\rightarrow$  affect + K<sup>th</sup> path

$w_{old} = w_{new} =$  No Change  
in movement

## # Vanishing Gradient Problem

- occurs (may) using Sigmoid AF. ( $z \rightarrow$  due to Sigmoid A



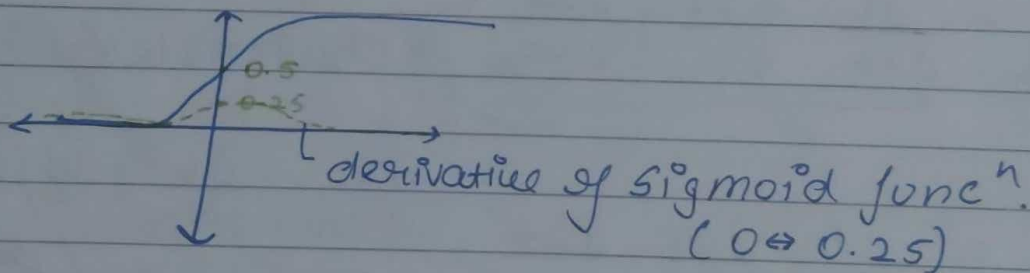
$$w_{1, new} = w_1 - \alpha \frac{dL}{dw_1}$$

③  $\rightarrow$  That's why don't use Sigmoid AF in all Hidden layers.

Result

$$\frac{dL}{dw_1} = \frac{dL}{dz_3} \cdot \frac{dz_3}{dz_1} \cdot \frac{dz_1}{dw_1}$$

★ Derivative of Sigmoid  $\rightarrow$  Blue 0 and 0.25



- from ①  $\rightarrow$  The product we get will get smaller and the further we Backtrack the smaller it will get.

+ Product with  $\alpha$  may make it more smaller  
So in ②

③  $\rightarrow$

$$w_{new} = w_{old}$$

No change

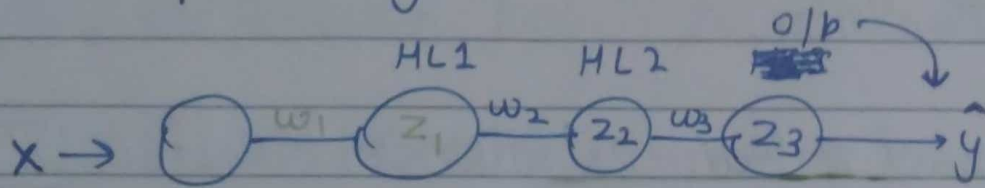
The Result will not be THAT  
Diff from the old value

(No significance diff in new and old value).

happens because of weights.

$$\begin{aligned} 0.21 &= z_2 \\ 0.11 &= z_1 \end{aligned}$$

## # Exploding Gradient Problem



$$\rightarrow w_1 = w_1 - \alpha \frac{\partial L}{\partial w_1} \left[ \frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z_3} \cdot \frac{\partial z_3}{\partial z_2} \cdot \frac{\partial z_2}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1} \right]$$

125.100  
= 12500

•  $z_2 = w_2 \cdot z_1 + b \rightarrow z_2 = \text{Act}(z_2)$   
 $z_2 = \phi(z_2)$   
Sigmoid (0-1)

$0 \leftrightarrow 0.25$

$$\frac{\partial z_2}{\partial z_1} = \frac{\partial \phi(z_2)}{\partial z_2} \cdot \frac{\partial z_2}{\partial z_1} \rightarrow \frac{\partial (w_2 \cdot z_1 + b)}{\partial z_1} \text{ Const}$$

$0 \leq \partial \phi(z_2) \leq 0.25 \cdot w_2$

$0.25 \cdot w_2 \rightarrow 0.25 \times 500 = 125$

larger value of derivative. (Result) Because of weights

old - large value = negative value (-ve)  
 new value

• Old and new value will vary a lot  $\rightarrow$  GD will never converge



1) Fix  $\rightarrow$  Overfitting

Multi-layer

- $\rightarrow$  No underfitting
- $\rightarrow$  But Overfitting

# DROP-OUT Layers in Multi-Neural Networks.

# And Regularization.

$\rightarrow$  Similar to 'subset of features' given to Multiple DTs in a Random Forest.

1) Select Drop-out Ratio  $\rightarrow (p) \rightarrow 0 \leq p \leq 1$

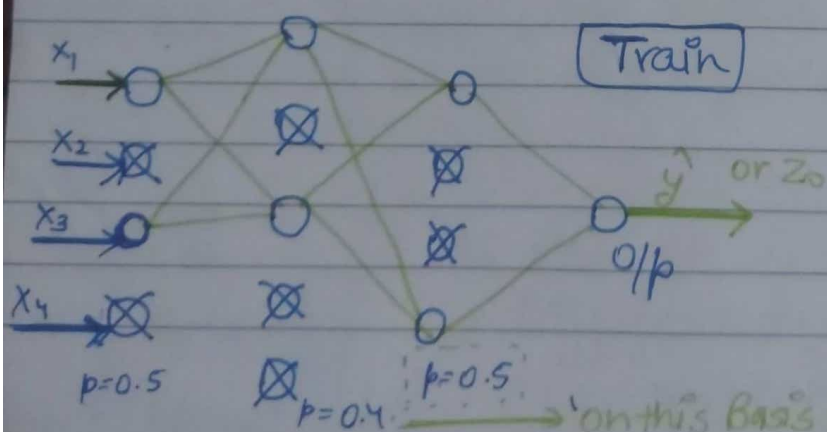
2) Select subset of features from I/p layer.

$\rightarrow$  Subset of features/activation func neurons from Hidden Layer

$\rightarrow$  selecting activated neurons from Hidden Layer.

$\rightarrow$  Depending on Drop-out Ratio some nodes and neurons are activated and deactivated. In Backtracking, weights of activated neurons are updated, in next iteration randomly some nodes and neurons are selected/activated and deactivated and cycle repeats itself.

For test data - There is no random selection, everything is connected but each and every weights are multiplied by drop-out Ratio  $(p)$  which was selected at the time of training.



Test

$\rightarrow$  everything is connected

But  $\rightarrow$   $w = w \times p$

every weights

$\rightarrow$  on this Basis Nodes are activated and deactivated.



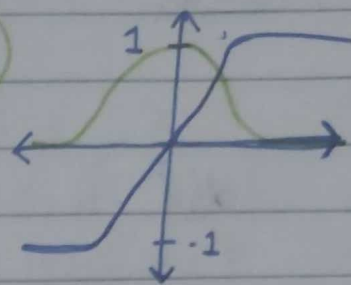
↓ solves Vanishing/Exploding G.D problem.

## # ReLU and Leaky ReLU

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

→ Threshold Activation function →  $f_{th}$  (-1 to 1)

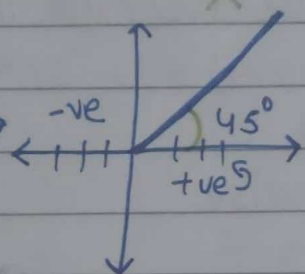
$$\frac{\partial J_{th}}{\partial z} = (0 \text{ to } 1) \quad L < 1$$



↳ Give prob to Vanishing and Exploding G.D problem.

• → ReLU →  $\max(y, 0)$  or  $\max(z, 0)$

$$\begin{cases} z & z \geq 0 \\ 0 & z < 0 \end{cases}$$



$$\tan 45^\circ = 1$$

→  $J_{rel}$

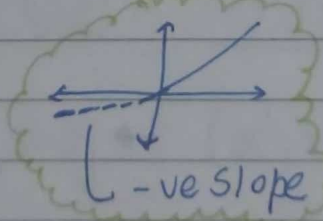
→ Derivative of ReLU Act (in Backtracking)

$$\frac{\partial J_{rel}}{\partial z} \rightarrow \begin{cases} 1, & z > 0 \\ 0, & z < 0 \end{cases} \quad \left| \begin{array}{l} \text{only 2 values} \\ \text{i.e., 0 or 1} \end{array} \right|$$

!! Prob - ? → one of the derivative become zero, then again vanishing G.D problem will arise

↓ Fix by 'Leaky ReLU'

→  $J_{lrl}$



• → Leaky ReLU → when  $z$  is negative, do not keep it zero, add some values. (small value addition)

$$J_{lrl} = \begin{cases} z, & z > 0 \\ 0.01(z), & z < 0 \end{cases}$$

↳ negative slope

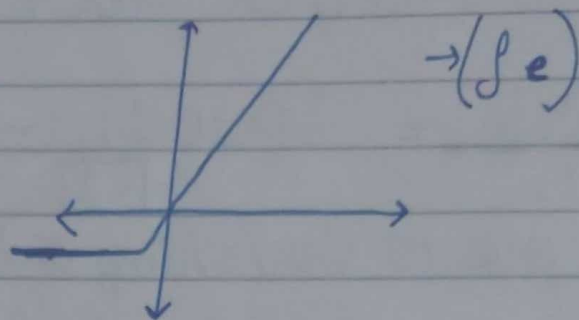
$$\frac{\partial J_{lrl}}{\partial z} \rightarrow \begin{cases} 1, & z > 0 \\ 0.01, & z < 0 \end{cases}$$

↳ -ve slope But not 0.

# # More Act functions ( $x=z$ )

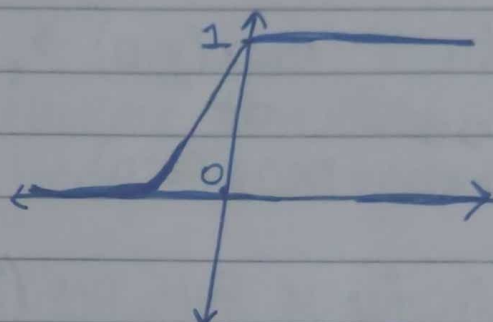
→ 1) ELU → (exponential Linear Unit) func<sup>n</sup>

$$= \begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{otherwise} \end{cases}$$



\* Don't have to find derivative of zero.

$$\frac{d \text{le}}{dR} \rightarrow \begin{cases} 1 & x > 0 \\ \alpha e^x & \text{otherwise} \end{cases}$$



• Because of  $e^x \rightarrow$  Time  $\uparrow\uparrow$

2) PReLU (Parametric Relu) → (P<sub>PR</sub>)

$$\rightarrow \begin{cases} x & x > 0 \\ \alpha x & \text{otherwise} \end{cases}$$

$$\begin{cases} \text{if } \alpha = 0.01 = \text{LReLU} \\ \alpha = 0 = \text{ReLU} \end{cases}$$

$\frac{\partial \text{loss}}{\partial x} \rightarrow \begin{cases} 1 & \text{if } x > 0 \\ \alpha & \text{if } x \leq 0 \end{cases}$   
No zero, solves problem

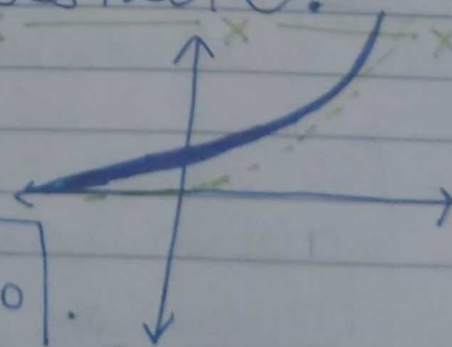
3) Swish Function (Use in LSTM) (NN > 40 layers)

$$y = x \cdot \text{sigmoid}(x) \rightarrow y = x \cdot \sigma(x)$$

→ solves dead ReLU as it passes near 0.

4) Softplus

$$y = \ln(1 + e^x) \quad \left( \begin{array}{l} \text{No need} \\ \text{of der of 0} \end{array} \right)$$





works on probab and real values

5) softmax  $\rightarrow S(x_j) = \frac{e^{x_j}}{\sum_{k=1}^K e^{x_k}}, j=1, 2, 3, \dots, K!$

$\rightarrow$  for an arbitrary real vector of length  $K$ , softmax can compress it into real vector of length  $K$  with a value in the Range  $(0, 1)$  and the sum of the elements in the vector is 1

ex 2

$\rightarrow$  Before the input to the o/p layer or i/p of o/p layer  
 $\{x_1, x_2, x_3, x_4\}$

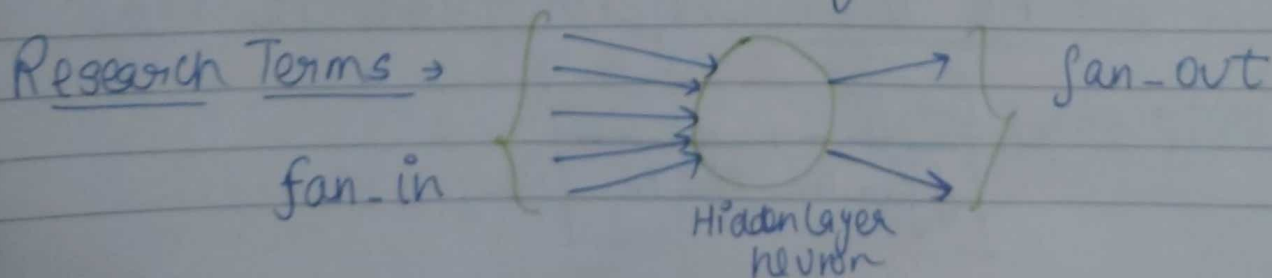
$$\frac{e^{x_1}}{e^{x_1} + e^{x_2} + e^{x_3} + e^{x_4}} + \dots + \frac{e^{x_4}}{e^{x_1} + e^{x_2} + e^{x_3} + e^{x_4}} \rightarrow \{p_1, p_2, p_3, p_4\}$$

highest P  
Ans.  
L o/p

\* Sigmoid and Softmax act probab and kept at last layer.  
( $>0.5, <0.5$ ) (Range of Probab)

## \* WEIGHT INITIALIZATION

- 1) weight should be small.
- 2) weight should not be same.
- 3) weight should have good Variance.



### METHOD-1: Uniform Distribution

$$w_{ij} \sim \text{Uniform} \left[ \frac{-1}{\sqrt{fan\_in}}, \frac{1}{\sqrt{fan\_in}} \right]$$

→ good with SIGMOID ACTIVATION FUNCTION

### METHOD-2: Xavier / Glorot Normal DISTRIBUTION

$$w_{ij} \sim N(0, \sigma) \quad (\text{mean}=0 \mid \text{sd}=\sigma)$$

$$\sigma = \sqrt{\frac{2}{fan\_in + fan\_out}}$$

### METHOD-2.1 → Xavier / Glorot Uniform Distribution

$$w_{ij} \sim U \left[ \frac{-\sqrt{6}}{\sqrt{fan\_in + fan\_out}}, \frac{\sqrt{6}}{\sqrt{fan\_in + fan\_out}} \right]$$

→ good with SIGMOID ACTIVATION FUNCTION

### METHOD-3: He init Uniform Distribution

$$w_{ij} \sim U \left[ \frac{-\sqrt{6}}{\sqrt{fan\_in}}, \frac{\sqrt{6}}{\sqrt{fan\_in}} \right]$$



$$GD \rightarrow \text{All points} \rightarrow \text{Loss} = \frac{1}{2} \sum_{i=1}^n (y - \hat{y})^2$$

$$SGD \rightarrow \text{one by one} \rightarrow \text{loss} = \frac{(y - \hat{y})^2}{2}$$

METHOD 3.1  $\rightarrow$  Heinit Normal Distribution

$$w_{ij}^0 \sim N(0, \sigma)$$

$$\sigma = \sqrt{\frac{2}{\text{fan-in}}}$$

$\rightarrow$  good for ReLU, LReLU, elu, P-ReLU and others.

\* STOCHASTIC GD [SGD]

$\rightarrow$  Find error  $\rightarrow$  find partial derivative  $\rightarrow$  finding slope  $\rightarrow$  finding weights

$\hookrightarrow$  one by one  $\rightarrow$  SGD  $\cdot$  Iterations / Epochs will  $\downarrow$   $\hookrightarrow$  F+Backward

$$w_{\text{new}} = w_{\text{old}} - \alpha \frac{dL}{dw_{\text{old}}} \rightarrow \frac{dL}{dw_{\text{old}}}$$

Req more comp. power  
To load  $\uparrow$

$$GD \rightarrow \text{Loss}(L) \rightarrow \sum_{i=1}^n (y - \hat{y})^2$$

$$SGD \rightarrow \text{loss}(L) \rightarrow (y - \hat{y})^2 \rightarrow \text{like in L.R}$$

$$\text{m.B.SGD} = \text{loss}(L) \rightarrow \sum_{i=1}^K (y - \hat{y})^2$$

$\hookrightarrow$  n data points (all)  $\rightarrow$  GD   
 very slow  $\uparrow$

$\hookrightarrow$  1 data point  $\rightarrow$  SGD

$\hookrightarrow$  K data points  $\rightarrow$  mini Batch SGD   
 ( $K < n$ )

## 3D - 2D representation of GD and MB SGD



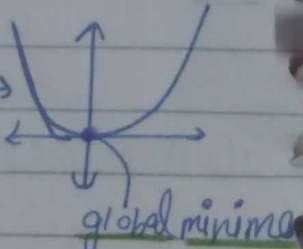
The zig-zag movement of MB SGD sol<sup>n</sup> is known as NOISY DATA

↳ We use MB SGD with Momentum to eliminate this or remove this (Noise).

$$\left[ \frac{\partial L}{\partial w_{old}} \right]_{\text{MS SGD}} \approx \left[ \frac{\partial L}{\partial w_{old}} \right]_{\text{GD}}$$

↳ sample
↳ population

## \* Global and local minima

- 1) Not every loss func<sup>n</sup> gives GD graph → 
- 2) some loss func<sup>n</sup> may give graph