

Semester	VII	L	T	P	C	COURSE CODE
Regulation	V20	3	0	0	3	V20CSTPE11
Name of the Course	Deep Learning (Professional Elective-III)					
Branch	Common to CSE & CST					

Syllabus Details

Course Outcomes: After Successful completion of the Course, the student will be able to:

- CO1: Describe the fundamentals of deep learning. (K2)
- CO2: Illustrate the working of deep feed forward neural networks. (K2)
- CO3: Discuss regularization and optimization techniques used in deep neural networks. (K2)
- CO4: Illustrate the working of convolution neural networks. (K2)
- CO5: Explain about recurrent and recursive neural networks. (K2)

UNIT-I: Introduction: Historical Trends in Deep Learning, The Many Names and Changing Fortunes of Neural Networks, Increasing Dataset Sizes, Increasing Model Sizes, Increasing Accuracy, Complexity and Real-World Impact.

UNIT-II: Deep Feed forward Networks: Learning XOR, Gradient-Based Learning, Hidden Units, Architecture Design, Back Propagation and Other Differentiation Algorithms.

UNIT-III: Regularization for Deep Learning: Parameter Norm Penalties, Early Stopping, Dropout; **Optimization for Training Deep Models:** How Learning Differs from Pure Optimization, Challenges, Basic Algorithms, Parameter Initialization Strategies, Algorithms with Adaptive Learning Rates, Optimization Strategies and Meta-Algorithms.

Adam, RMSprop

UNIT-IV: Convolution Networks: The Convolution Operation, Motivation, Pooling, Convolution and Pooling as an Infinitely Strong Prior, Variants of Basic Convolution Functions, Structured Outputs, Data Types, Efficient Convolution Algorithms, Random or Unsupervised Features, The Neuroscientific Basis for Convolutional Networks, Convolutional Networks and the History of Deep Learning.

2x2 CNN
ReLU

UNIT-V: Sequence Modeling- Recurrent and Recursive Nets: Unfolding Computational Graphs, Recurrent Neural Networks, Bidirectional RNNs, Encoder-Decoder Sequence-to-Sequence Architectures, Deep Recurrent Networks, Recursive Neural Networks, The Challenge of Long-Term Dependencies, Echo State Networks, Leaky Units and Other Strategies for Multiple Time Scales, LSTM and Other Gated RNNs, Explicit Memory.

Textbooks:

1. Deep Learning, Ian Goodfellow, Yoshua Bengio, and Aaron Courville, MIT Press.

Reference Books:

1. Neural Networks and Deep Learning, Charu C. Aggarwal, Springer.
2. Fundamentals of Deep Learning, Nikhil Buduma, 1st Edition, O'Reilly

AI, ML, DL

★ AI :- It is the broad idea that machines can intelligently execute tasks by mimicking human behavior;

— AI recreates human intelligence @ behavior using algorithms, data, models ✓

— Ex:- self-driving cars, security systems @ medical imaging; Siri @ Alexa --

★ ML :- subset of AI, revolves around the idea that machines can learn @ adapt through experiences @ data to complete specific tasks; ✓

— ML is an AI algorithm which allow system to learn from data;

Ex:- sentiment analysis, predictive maintenance, spam filters --

structured data :- data arranged in rows @ columns

Ex:- Tables, excel sheets, databases

★ DL :- DL attempts to mimic the human brain through neural networks;

— DL is the evolution of machine learning @ NN's which uses advanced computer programming @ training to understand complex patterns hidden in large data sets; Ex:- image classification, NLP, generative models ✓

— DL is about understanding how the human brain works in different situations,

Ex:- then trying to recreate its behavior;

— used to complete complex tasks @ unstructured data → email, social media posts, audio's, images, videos..

for example :- Drawing a cat :-

AI - you tell the robot, "Draw a cat", @ it tries using rules you give it. It acts smart but doesn't really learn on its own

ML - you show the robot 100 pics of cats @ it learns from them to draw its own cat.

DL - you show it 10,000 pics, @ it uses a deep brain-like system to learn even more accurate legs --

just like a wrapper around Tensorflow ✓

① keras, Tensorflow, pytorch;
(libraries)

② mathematics + statistics; problem solving;

③ GPU - for heavy DL projects; (NVIDIA) — Titan RTx (model)

④ Why DL popular :- ① Data growth (volume of data)

② Hardware advancements — GPU, TPU → Tensorflow processing unit (google)

③ Python & open source ecosystem;

④ frameworks :- Pytorch (2016), Tensorflow 2.0 (2015), Dist Belief (2011)

snippet :-
from tensorflow import keras
model = keras.Sequential([keras.layers.Dense(10, input_shape=(784,), activation='sigmoid')])

⑤ Job ⑥ AI Boom;

⑦ High accuracy with large data → automatic feature extraction

⑧ versatility across domains — CNN, RNN, Transformer, audio

⑨ Breakthroughs in AI Capabilities — Chatbot, Deepfake, Autonomous vehicle, AlphaGo

⑩ State of art performance & scalable;

Important terms :-

① Neuron, Perceptron;

② Activation Functions :-

→ ReLU (Rectified Linear Unit)

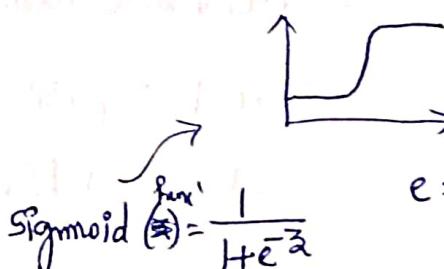
(Logit) Sigmoid, Tanh, softmax

③ Layers :- I/P, H/L, O/P; I/P = Raw data, H/L = Perform transformation, O/P = final prediction

④ Loss Functions — measures error b/t P & A O/P]

→ MSE = Regression

Cross-entropy = Classification;



$$\text{Sigmoid}(x) = \frac{1}{1+e^{-x}}$$

$e = 2.71828$
sigmoid func converts input into range 0 to 1

- ⑤ Backpropagation - It's a algorithm to minimize loss by updating weight w_j gradient
- ⑥ Gradient Descent - used to update weight to reduce error;
 (optimization) \rightarrow stochastic gradient descent, Adam, RMS prop;
-
- Model Architectures :-

① ANN, CNN, RNN ② Transformer

<u>Feature</u>	<u>ANN (Artificial NN)</u>	<u>CNN (Convolutional NN)</u>	<u>RNN (Recurrent NN)</u>
Purposes	General purpose NN	mainly for Images & Spatial data	Sequence & Time-Series data
Input type	fixed size Vector	Images	text, audio, time-series
use cases	tabular data, basic classification;	Images, video, object detection;	Language modelling, speech, time prediction (LSTM / GRU)
Realworld examples	creditcard fraud detection, medical diagnosis, stock market forecasting;	face Recognition, object detection; medical imaging, satellite image Analysis;	google Translate, audio to text, Analyzing tweets, Review, weather forecast
For	Regression / Classification	Computer vision / Images	Time series & Analysis

time
effort
money

for selection of a particular architecture

BERT

image classification \rightarrow ResNET

Text classification \rightarrow BERT

Image segmentation \rightarrow UNet

Image translation \rightarrow pix2pix

object detection \rightarrow YOLO

Speech generation \rightarrow WaveNET

Readily available Architectures

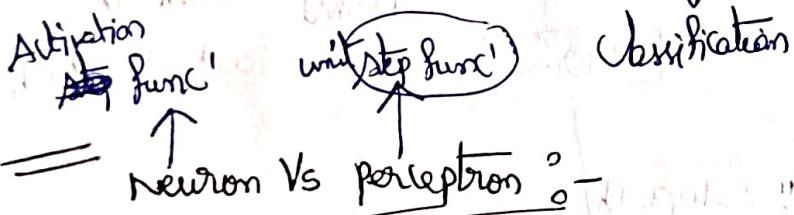
(2)

④ Training & Performance :-

- ① Epoch ;
- ② Batch size ;
- ③ overfitting & underfitting ;

④ Regularization :- prevent overfitting ; L_1/L_2 regularization, dropout, early stopping

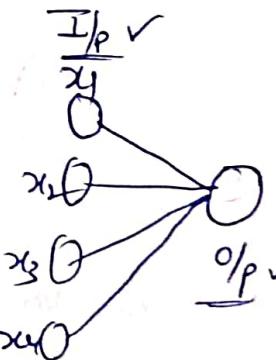
⑤ performance metrics :-



⑥ perceptron :-

Basic perceptron does not have a hidden layer ;
(single layer)

→ only ① Input layer
② O/P layer ;



— solves only linearly separable problem ;

(i) multilayer perceptron :- (MLP) :- like AND, OR ;
(MLP) :- not complex patterns ;

— includes one or more hidden layers
— uses activation functions ;

(Note) :- MLP is a subset of NN's ; i.e. ANN is the umbrella term

All MLP's are NN's, but not all NN's are MLP's

(MLP) is a type of ANN with multiple layers

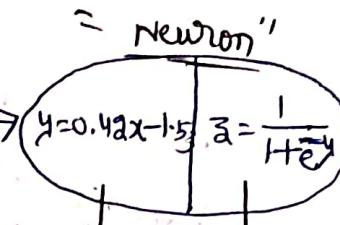
weights :- control the influence of inputs b/w neurons ;

bias :- constant added to the o/p, provides the flexibility ;
for adjustment of activation thresholds @ shifting

input layer = where data is fed into network ;
hidden layer = where all the mathematical calculations are done using weights, biases @ data
o/p layer = where meaningful results are obtained from

Neuron :-

Say ex:- Age = 43



Binary C/P

Say

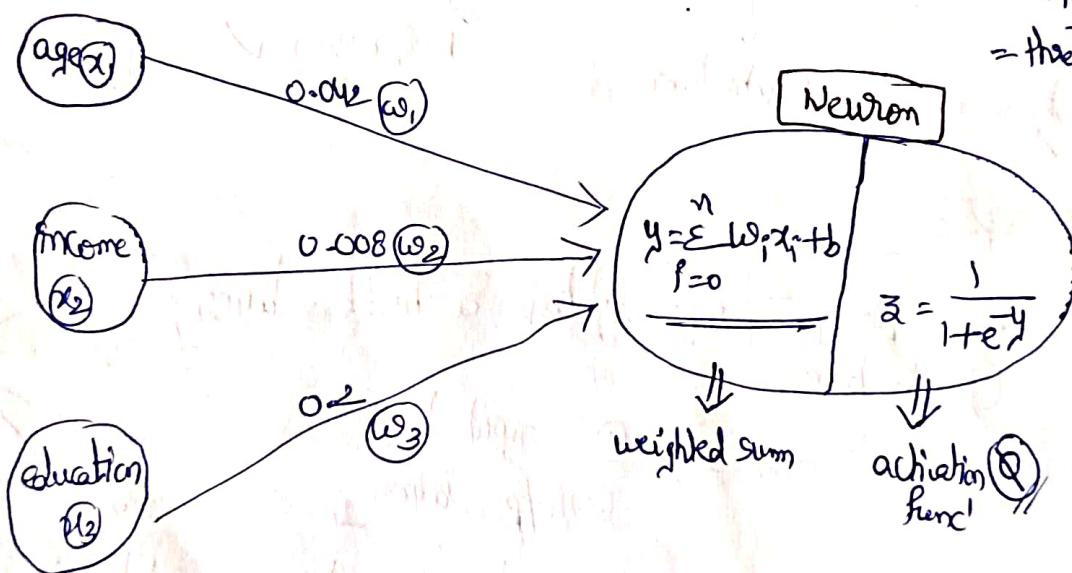
$$y = 0.042(x) - 1.53$$

$$y = 0.042(x_1) + 0.008(x_2) + 0.2(x_3) - 1.53$$

$$y = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

$$y = \sum_{i=0}^n w_i x_i + b$$

w_1, w_2, \dots = weights ✓ = importance & Ifp
 b = bias ✓ = threshold shifts the activation threshold



① Neuron means node (a) unit, is the basic computational unit of a NN;

that mimics the behavior of a biological neuron;

② Activation func! decides, if the neuron "fires" ③ passes info forward ✓

~~ANATOMY~~

(Biological) Natural NN (say)

Neuron

Neuron

dendrites

Capture info from other neuron ✓

axon

Passes the info. ✓ (like passage)

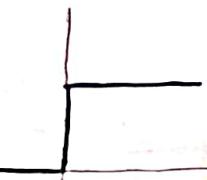
2nd neuron start point

0/p

④ Perception [Indirect] :- It is an algorithm for supervised learning of binary classifier; (T/F)

- It can be seen as a single unit of an ANN (also known as prototype for neural net);
- single layer NN with the unit step func as an activation func;

unit step function :- $g(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{otherwise;} \end{cases}$



prediction = $\hat{y} = g(f(x))$

$$= g(w^T x + b)$$

indicates "applying activation func";

w = weights
b = bias; = threshold

perception update Rule :-

for each training sample (x_i) :-

i) Initialize weights ✓

ii) $w = w + \Delta w$

$b = b + \Delta b$

$$\text{where } \Delta w = \alpha (y_i - \hat{y}_i) \cdot x_i$$

$$\Delta b = \alpha (y_i - \hat{y}_i)$$

where $\alpha = \text{learning rate}$
 $= [0, 1]$

actual	prediction
1	1
1	0
0	0
0	-1

perception simply :-

- A line in 2D which creates two regions to classify the two classes;
3D - plane makes ✓

Code :- import numpy as np

class perception:

def __init__

④ disadv :- perception not works on non linear data;

⑤ playground.tensorflow.org ✓

Plannable parameters = ? for Backpropagation

$$\begin{aligned} \text{Total weights} &= 12 + 6 + 3 = 21 \\ \text{bias} &= 3 + 2 + 1 = 6 \end{aligned} \quad (25) \quad \text{Plannable parameters}$$

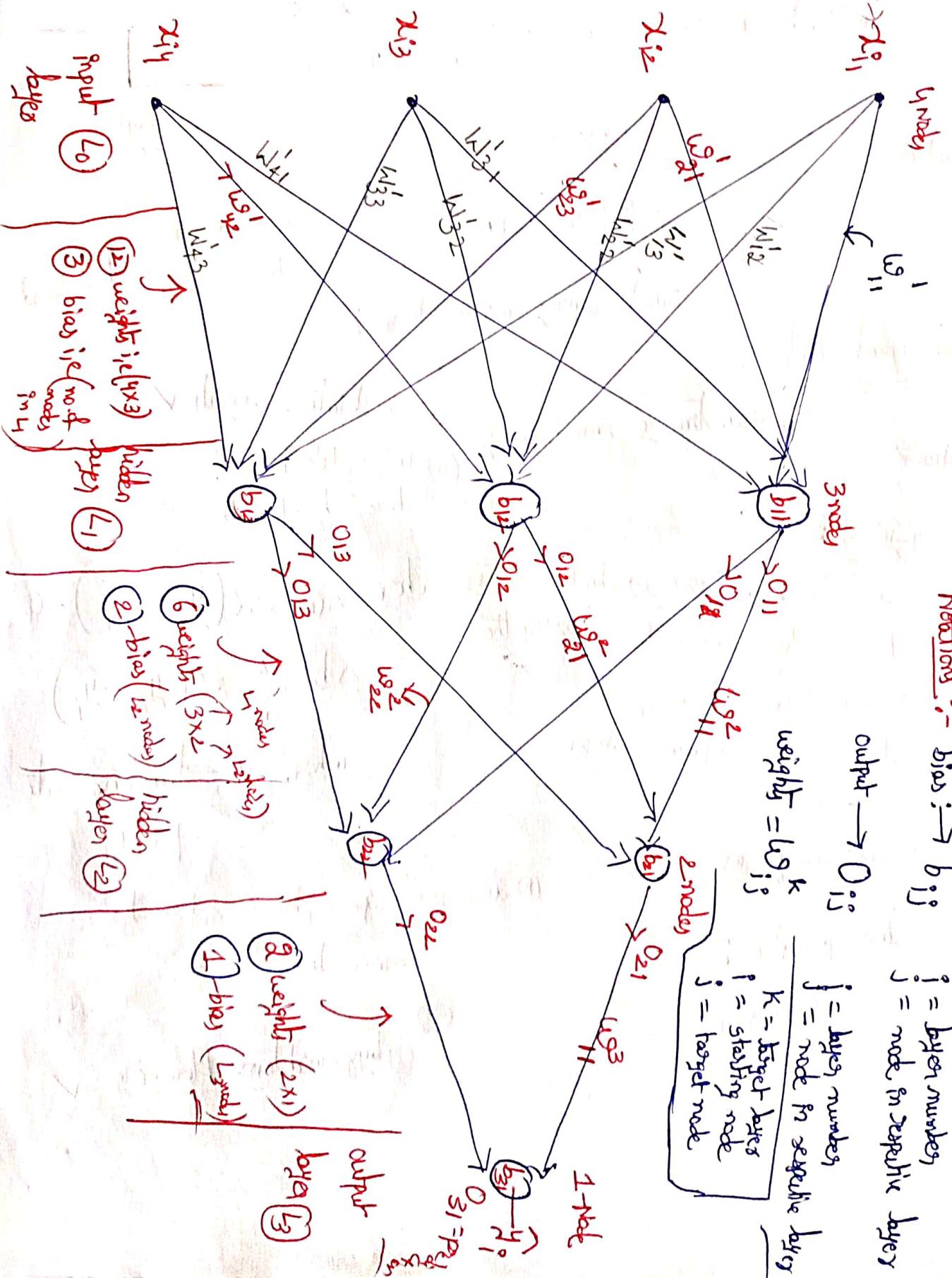
1	2	3	4	5
.
.
.

* MLP :- (1) weights (2) bias

(2) How to denote weights (3) bias ✓

MLP notations :- data = $\{n \times n\}$ $m = \text{rows}$

say $n=4 = \text{columns}$



Activation Function (AF) in DL :- $\textcircled{1}$ very important component in NN's; $\textcircled{2}$ introducing non-linearity

Definition :- In ANN, each neuron forms a weighted sum of its input $\textcircled{1}$ passes the resulting scalar value through a func' ref'd to as an activation function $\textcircled{2}$ transfer func'. If a neuron has n inputs then the output $\textcircled{3}$ activation of neuron is

$$a = g(w_1x_1 + w_2x_2 + \dots + b) \quad \text{where } \textcircled{1} \text{ ref's as AF;}$$

In short, AF if a mathematical gate b/t input & output;

Ideal AF Quality :-

$\textcircled{1}$ should be non-linear

$\textcircled{2}$ should be differentiable (to help gradient descent) $\textcircled{3}$ p'

$\textcircled{3}$ " be computationally inexpensive;

$\textcircled{4}$ zero centred; $\textcircled{5}$ normalized

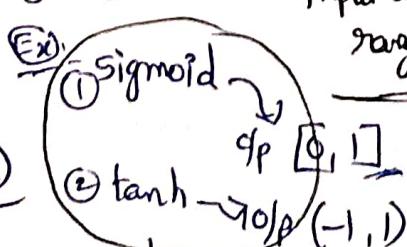
(mean = zero)

Ex: tanh

$\textcircled{6}$ Non-saturating (means should scale the input data range)

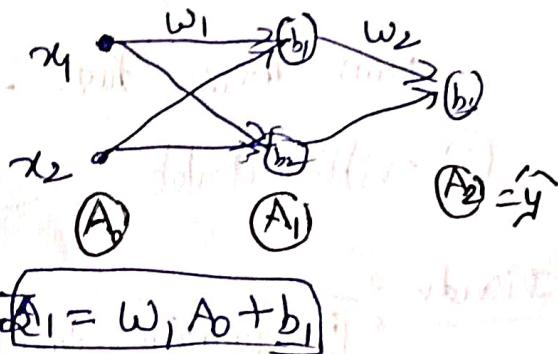
Ex:- ReLU

$$f(x) = \max(0, x)$$



Ex: $\textcircled{3}$ Saturating AF

Causes vanishing gradient problem



if AF = linear then $A_1 = z_1$

$$A_1 = g(w_1 A_0 + b_1)$$

$$= w_1 A_0 + b_1$$

$$= w_2 (w_1 A_0 + b_1) + b_2$$

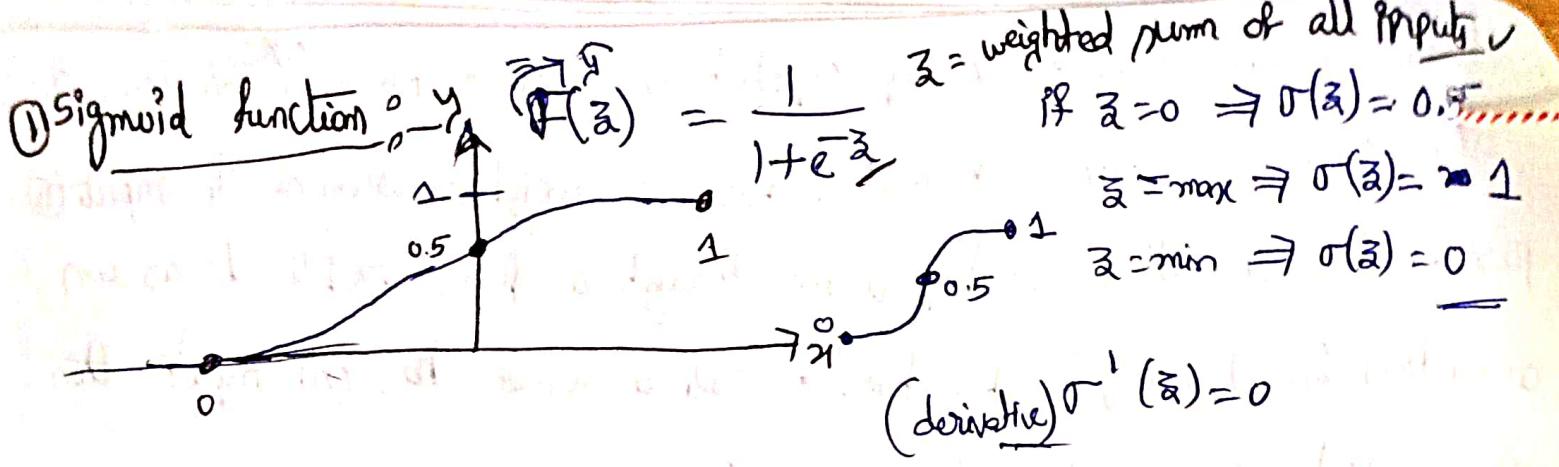
$$= \underbrace{w_2 w_1}_{w} \underbrace{A_0 + w_2 b_1 + b_2}_{b}$$

$$= w^T A_0 + b^T$$

$$A^T = w^T A_0 + b^T$$

↓
output
↓
input

its relation is linear



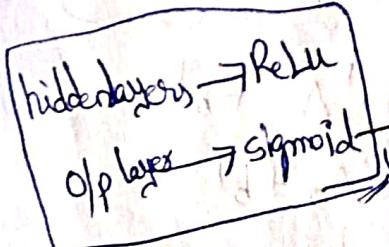
adv: ① O/p always like prob [0, 1] → feels like probability @ o/p buyers ②

user for Binary classification problem ✓ $\begin{cases} 0 = \text{No} \\ 1 = \text{Yes} \end{cases}$

② Non-linear func' → capture non-linear data;

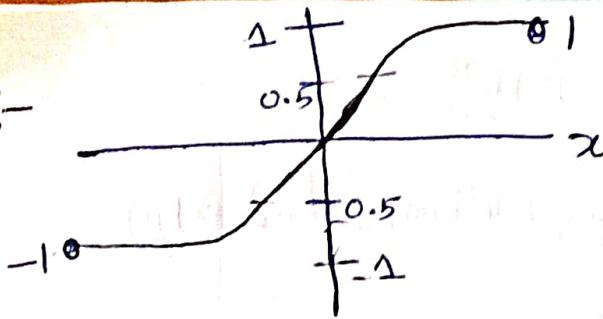
③ Differentiable → helped in Backpropagation (BP) i.e. $\left(\frac{\partial L}{\partial w} \right)$ ✓

Disadv: ① Saturating function → because of its range; it causes Vanishing gradient problem
 (i) Non-zero Centered → that's why we don't use in hidden layer;
 means $BP^1 \quad w_n = w_0 - \eta \frac{\partial L}{\partial w}$
 (ii) Causes training slow → used only in BP → means if $w_n = w_0$ no update will take place
 means normalized the values you may get either one or zero but not mean = 0
 (iii) Computationally expensive; because of component i.e. $\sigma(z) = \frac{1}{1+e^{-z}}$



binary classification

Tanh



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

(ii) range $(-1, 1)$

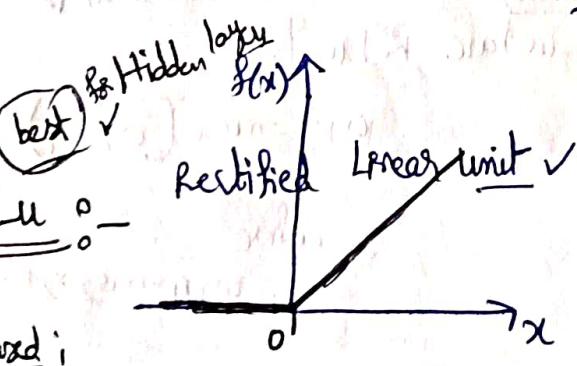
Adv :- (i) Non Linear

(ii) Differentiable

(iii) Zero-centered ; θ_w } θ_v } Training faster than Sigmoid

Disadv :- (i) saturating func! \rightarrow prone to vanishing gradient prob

(ii) computationally expensive;



(iii) mostly used;

sigmoid is used @ output layers because of its output range $(0, 1)$ that helps in binary classification.

$$f(x) = \max(0, x) \rightarrow \text{weighted sum}$$

Adv :- (i) Non Linear ; because of its $\max(0, x)$

means for θ_w values $f(x)=0$

for θ_v values $f(x)=x$

(ii) Non-saturated in θ_w region

(iv) converge, faster

(iii) computationally inexpensive

disadv :- (i) not completely differentiable;

(ii) not zero centered ; to overcome this we will use Batch normalization

⊕ Best choice for hidden layers is ReLU;

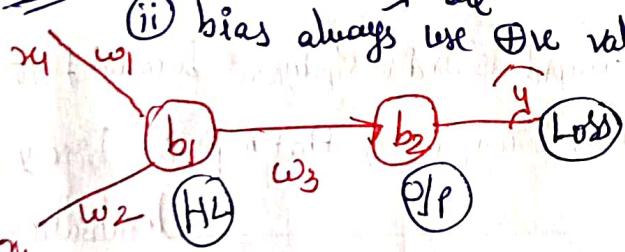
⊖ Big disadvantage of ReLU problems; (dead ReLU)

Causes :- (1) large negative bias
explanation (2) most inputs being negative;

The issue arises because the derivative of ReLU is zero for negative inputs; if a neuron's input consistently falls into the negative range due to its weight (⊖ bias), it will always o/p zero. (⊖ the zero gradient prevents the neuron from learning) (⊖ updating its weight) ✓

Solution to this :-

- ① set low learning rate (if update not happening then node said to be dead node)
- ② bias always use \oplus ve value i.e. (0.01) ✓
- ③ use its variants ✓



@HL we have ReLU like

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

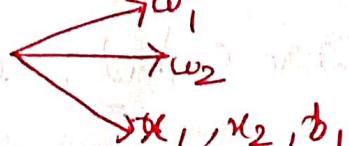
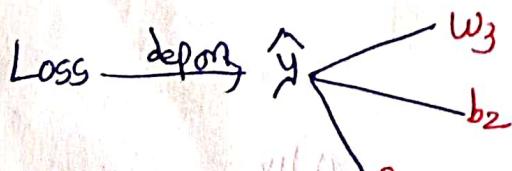
where

$$z_1 = w_1x_1 + w_2x_2 + \dots + b_1 < 0$$

then $a_1 = 0$

then $\frac{\partial a_1}{\partial z_1} = 0$ then w_1, w_2 can't be updated;

then node @HL called as dead node;



⊕ Let's take regression problem;

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1} \quad (\text{Chain Rule})$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_2}$$

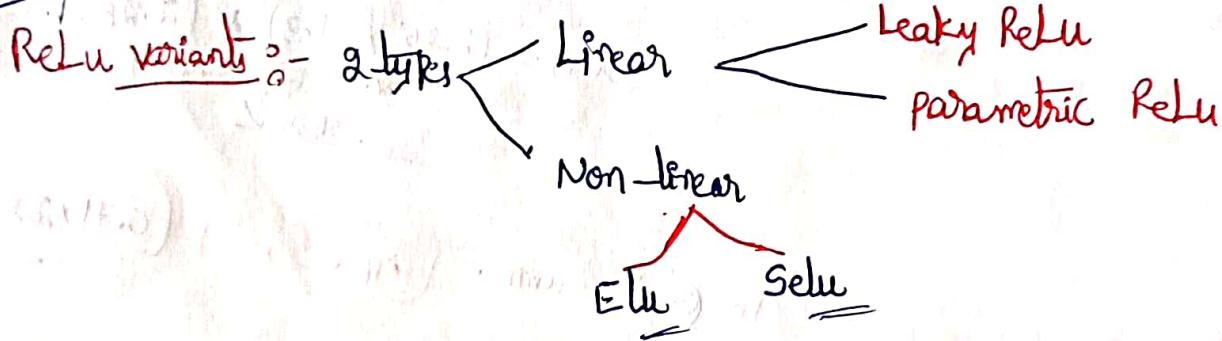
if it is zero then finally $\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2} = 0$

if $z_1 = w_1x_1 + w_2x_2 + b_1 < 0$

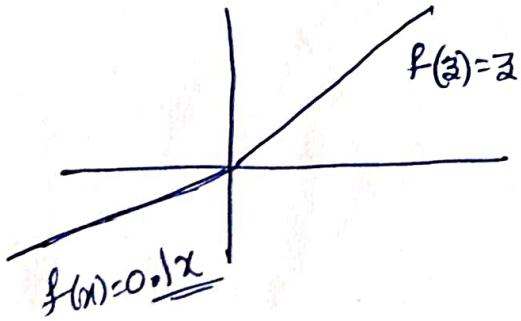
$z_1 = 0$ when ① High learning rate ②

~~This is~~ ③ High Zero bias ✓

— once a neuron dead, it never re-work



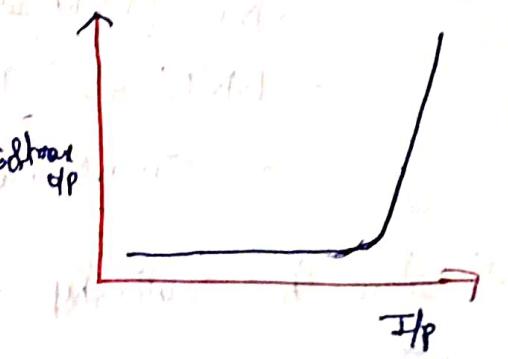
Leaky Relu :- $f(z) = \max(0.01z, z)$ if $z > 0 = z$
 $z < 0 = \frac{1}{100}z$ = fraction of z ✓



④ Neural networks learn based on connections; Forward & Backward propagation

- ⑤ ReLU → Regression → Linear unit → step func!
- Sigmoid → Classification → Non-linear → ReLU, sigmoid, Tanh
- Softmax → Multi-class → exponential Linear unit → softmax, softplus;

- ⑥ Softmax function :-
- designed to handle multi-class classification problems;
 - It transforms raw o/p scores from a NN into probabilities;
ie vector of raw scores (logits) Ex:- $p(x) = 0.5, p(y) = 0.3, p(z) = 0.2$
 - It works by squashing the output values of each class onto range 0 to 1 while ensuring that sum of all probabilities equals 1.

- ⑦ Given a vector of Real numbers $\vec{z} = (z_1, z_2, \dots, z_k)$
then softmax function computes:
- $$\sigma(\vec{z}) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}} \quad \text{for } i=1, 2, \dots, k$$
- 

Ex:- say $\vec{z} = [2, 1, 0.1]$

z_i = input score

k = Total no. of classes

$\sigma(\vec{z})$ = probability of class i

Ans:- (i) $e^{\vec{z}} = [e^2, e^1, e^{0.1}] = [7.39, 2.72, 1.11]$

(ii) $\Sigma = 7.39 + 2.72 + 1.11 = 11.22$

(iii) $\text{softmax}(\vec{z}) = \left[\frac{7.39}{11.22}, \frac{2.72}{11.22}, \frac{1.11}{11.22} \right] = [0.66, 0.24, 0.1]$

↓ ↓ ↓
66% 24% 10%

- ⑧ Typically combined with cross-entropy loss for training;
↓
loss func' used for classification purpose ✓

Softmax vs sigmoid

feature

output type

used for

formula

softmax

probability distribution (sum to 1)

multi-class classification

Involves all outputs

Sigmoid

Indep. probabilities

Binary classification

Each output computed
indep.

adv's

① Differentiable ② suitable for gradient based learning

② Ideal for tasks like Image classification

Text

③ used in output layers of classifiers like

- CNN (MNIST, CIFAR-10)

- RNN, LSTM's

- Transformers (Token classification)

④ show the relationship graph b/w sigmoid & softmax ✓

code for softmax

```
import numpy as np
```

```
def softmax(x):
```

```
e_x = np.exp(x - np.max(x))
```

```
return e_x / e_x.sum(axis=0)
```

```
logits = np.array([2.1, 0.1])
```

```
probs = softmax(logits)
```

```
print(probs) → [0.659, 0.34, 0.001]
```

Softplus :- smooth approximation of ReLU function, it allows for differentiability.

@ all points unlike ReLU which is not differentiable @ 0;

(i) $f(x) = \log(1 + e^x)$ — op always +ve

range : $(0, \infty)$ — function is smooth & differentiable everywhere

Intuition :-

For large +ve x , e^x becomes very large; so:

$$f(x) = \log(e^x) = x \checkmark$$

" " For x , e^x approaches Zero so: $f(x) = \log(1) = 0$
so simply fn behaves like:

- ReLU for large positive x
- Near Zero for large negative x
- Smooth curve fn between ✓

Comparison :-

Input

ReLU

softplus

-2

0

0.13

so, softplus always

0

0

0.69

gives a small +ve

2

2

2.13

even for negative

inputs ✓



(*) $\frac{d}{dx} [\text{softplus}] = \text{sigmoid}; \quad [\text{derivative of softplus} = \text{sigmoid}] \checkmark$

Adv :-

— smooth version of ReLU

— Differentiable everywhere; — better for optimization;

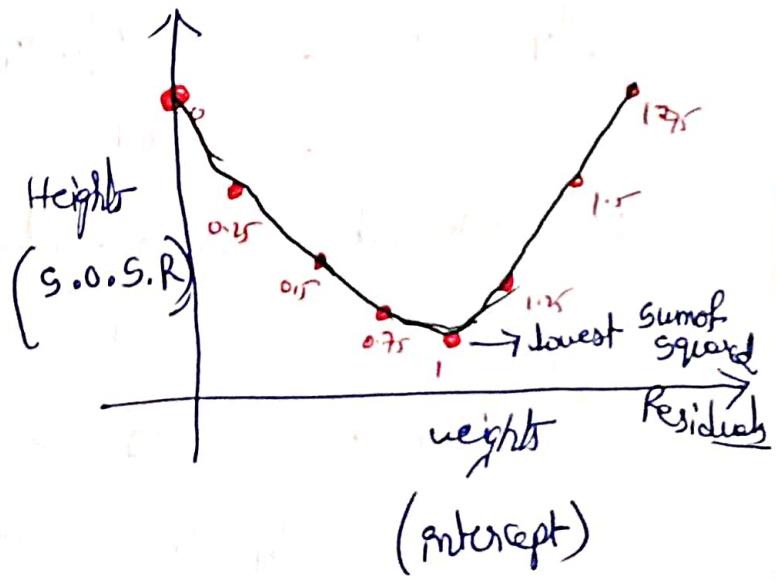
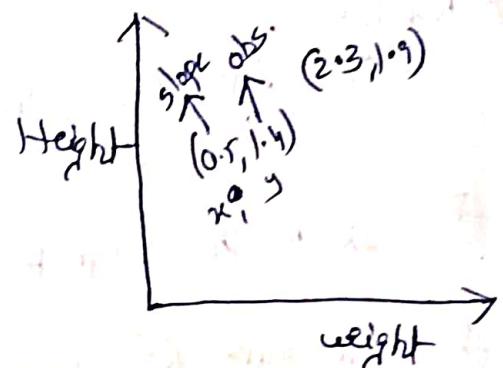
— Avoids "dying ReLU problem"

Disadv :-

— Computationally expensive than ReLU ✓

— In practice, ReLU outperforming it in terms of speed & training behavior ✓

- Gradient Descent :- Gradient descent is one of the most popular algorithms to perform optimization by far the most common way to optimise NN's;
- $J(\theta) = \text{loss func!}$ \circledast it min' the loss (or) loss func!
 - In ML & DS, DL we optimise a lot of stuff;
 - Linear Regression :- Say : Height_i = $m x_i + C$
 $= \text{slope} x \text{ weight} + \text{f(intercept)}$
 - predicted Height = intercept \oplus (slope \times weight)
 - We'll start by using GD to find the intercept;
 once we understand how GD works, we'll use it to solve for the intercept \oplus the slope;
 - Say slope = 0.64 \circledast then just estimate least squares;
 - sum of squared Residuals = type of loss func!; (say $1.2^2 + 0.4^2 + 1.3^2 = 3.1$)
 when intercept = 0
- \circledast what's the best value for intercept is b/c $(0.75, 1, 1.25) = ?$
- \circledast solution is :- GD only does a few calculations far from a optimal solution,



Sum of squared residuals = $(\text{observed} - \text{predicted})^2$

$$= (1.4 - (\text{Intercept} + 0.64 \times 0.5))^2 + (1.9 - (I + 0.64 \times 2.3))^2 \\ + (3.2 - (I + 0.64 \times 2.9))^2$$

$$\frac{d(S.O.S.R)}{d(\text{Intercept})}$$

$$= -2 [1.4(I + 0.64 \times 0.5)] - 2 [1.9(I + 0.64 \times 2.3)]$$

$$- 2 [3.2(I + 0.64 \times 2.9)]$$

The derivative, G.D will use it to find where the S.O.S.R is min or lowest;

- It where the slope of curve $\equiv 0$ then that gives best intercept; (take baby steps to get optimal)
- this makes G.D very useful when it is not possible to solve for where the derivative $\equiv 0$, & this is why G.D can be used in so many diff!

Cases:

$$- \text{put } I=0 \text{ in the derivative eq: } \frac{d(SOSR)}{d(I)} = -5.7$$

$$- \text{step size} = -5.7 \times 0.1 \rightarrow \text{learning rate} \checkmark \\ = -0.57$$

$$- \text{new intercept} = \text{old intercept} - \text{step size} \\ = 0 - (-0.57) = 0.57 \checkmark$$

2nd cycle: ~~step size~~ \rightarrow new $I = 0.57$ sub. in derivative eq: & get -2.3

$$\text{step size} = -0.23$$

$$\text{new intercept} = 0.57 - (-0.23) \\ = 0.57 + 0.23 = 0.8 \checkmark$$

Now $I(0.8) \Rightarrow -0.9$

step size $= -0.9 \times 0.1 = -0.09$

New Intercept $= \underline{0.89} \checkmark$

like next " " $= \underline{0.92} \checkmark$

" " $= \underline{0.94} \checkmark$

" " $= \underline{0.95} \checkmark$

— Now each step gets smaller & smaller the closer we get to the bottom of the curve;

(Note):- The least squares estimate for the intercept is 0.95 ✓

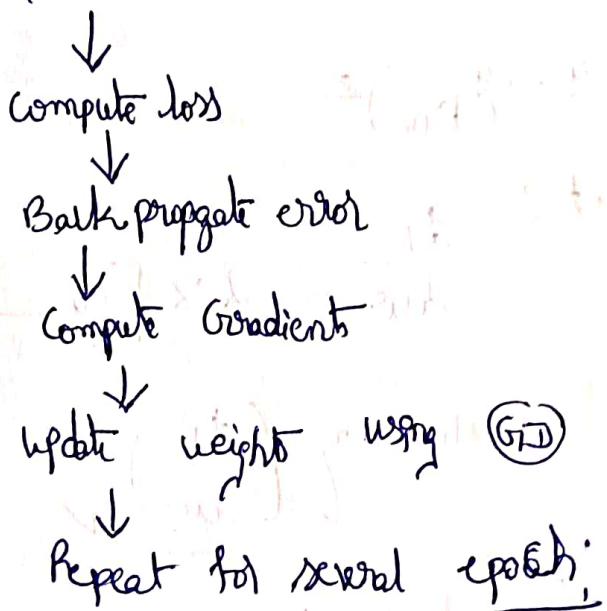
④ G.D. :- It is an optimization algorithm used to train NN by updating their weights & biases to min' loss func' [Backprop of NN]

Why do we need G.D. :-

- (i) We predict o/p using weights & biases;
- (ii) The loss func' tells us how wrong the prediction is;
- (iii) Goal :- Adjust the weights/biases so that the loss becomes as small as possible;

visual diagram :-

Input \rightarrow weights + Biases \rightarrow Activation \rightarrow o/p



How it works :-

- ① Initialize weights & biases randomly;
- ② Feed forward the I/p data through the network to get prediction;
- ③ cal' the loss (error) b/t prediction & actual o/p;
- ④ Backpropagate the error to compute gradient;
- ⑤ update the weights in the direction opposite to the gradient;

$$w \leftarrow w - \gamma \left(\frac{\partial L}{\partial w} \right) \rightarrow \text{gradient of loss}$$

LR
Learnin' rate

Challenges :-

- Local minima / Saddle point
- Learning Rate :-
 - Too high - overshoot
 - Too low - slow learning
- Vanishing / exploding gradient;

(*) prediction (\hat{y}) = wx

True value :- y

(MSE) loss function :-

- we compute the gradient :-

$$\frac{\partial L}{\partial w} = 2 \times [y - y_{\text{true}}]$$

then update :-

$$w = w - \eta \left(\frac{\partial L}{\partial w} \right)$$

(Ex:-

climbing down a hill to find the treasure;

say:-

- you are on top of a big hill, blindfolded Θ someone tells you: anything, so how do you find the treasure?

- steps :- (i) you feel which direction the hill is going down,

(ii) you take a small step downhill,

(iii) again, you feel the slope Θ take another step down;

(iv) you keep doing this until

(v) you reach the bottom where the treasure is.

hill = error (or) loss
treasure = min. error
Each step = weight update
feeling slope = (cal) gradient
step down;
walking slowly = learning rate

Weight Initialization Techniques :-

- Q) Why weight initialization is imp? :-
- (1) Initialize the parameters (w, b)
 - (2) Choose an optimization algorithm

- Out of all the points, (1) is most

Imp.

- Vanishing Gradient problem

Exploding " "

Slow Converge "

because of above issues to

- Overcome, LWT came;

③ Repeat these steps :-

- Like
GD steps
- Forward propagate an input
 - Compute the cost func' (loss func)
 - Compute the gradients of cost w.r.t. param. using BP;
 - Update each parameter using the gradient, according to optimization algorithm;

Wrong WI Techniques :- (what not to do)

- Case ① :- Zero Initialization (means all weights zero) - Regression Example :- student package how much based on GPA
-
- Linear Activation
- Test with ReLU, Tanh, Sigmoid

Say @ b_{11}, b_{12} put ReLU then

$$a_{11} = \max(0, z_{11})$$

$$\text{where } z_{11} = w_{11}^T x_1 + w_{12}^T x_2 + b_{11}$$

$$a_{12} = \max(0, z_{12})$$

where

$$z_{12} = w_{12}^T x_1 + w_{21}^T x_2 + b_{12}$$

Let us assume here, $w=0$, $b=0$ then

$$\begin{aligned} z_{11} &= 0 \\ z_{12} &= 0 \end{aligned}$$

then

$$a_{11} = 0 = a_{12}$$

equal to zero

$$\textcircled{B} \text{ BP' } w_{11}^1 = L_{11}^1 - \gamma \left(\frac{\partial L}{\partial w_{11}^1} \right)$$

↓
0 because of a_{11}

$$w_{11}^1 = w_{11}^1 \Rightarrow \text{means } \text{(no) update happening} \checkmark$$

- No training will take place

similarly with Tanh :-

$$a_{11} = \frac{e^{\bar{z}_{11}} - e^{-\bar{z}_{11}}}{e^{\bar{z}_{11}} + e^{-\bar{z}_{11}}} = \frac{1-1}{1+1} = 0 \checkmark$$

like $a_{12} = 0$ then automatically $\boxed{a_{11}, a_{12} = 0} \checkmark$

Try with

Sigmoid :- put @ HL :- $a_{11} = \sigma(\bar{z}_{11}) = 0.5$
 $a_{12} = 0.5$ } means $a_{11} = a_{12}$

No update happening ✓

$$L = w - \gamma \left(\frac{\partial L}{\partial w} \right)$$

like $\frac{\partial L}{\partial w_{11}^1} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial a_{11}} \times \frac{\partial a_{11}}{\partial \bar{z}_{11}} \times \frac{\partial \bar{z}_{11}}{\partial w_{11}^1} = x_1$ call for all 4 weights

$$\frac{\partial L}{\partial w_{11}^1} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial a_{11}} \times \frac{\partial a_{11}}{\partial \bar{z}_{11}} \times \frac{\partial \bar{z}_{11}}{\partial w_{11}^1} = x_1$$

$$\text{likewise do } \frac{\partial L}{\partial w_{12}^1} \text{ } \frac{\partial L}{\partial w_{21}^1} \text{ } \frac{\partial L}{\partial w_{22}^1} \text{ } \checkmark$$

$$\frac{\partial L}{\partial w_{12}^1} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial a_{12}} \times \frac{\partial a_{12}}{\partial \bar{z}_{12}} \times \frac{\partial \bar{z}_{12}}{\partial w_{12}^1} = x_1 \checkmark$$

$$\frac{\partial L}{\partial w_{21}^1} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial a_{21}} \times \frac{\partial a_{21}}{\partial \bar{z}_{21}} \times \frac{\partial \bar{z}_{21}}{\partial w_{21}^1} = x_2 \checkmark$$

$$\frac{\partial L}{\partial w_{22}^1} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial a_{22}} \times \frac{\partial a_{22}}{\partial \bar{z}_{22}} \times \frac{\partial \bar{z}_{22}}{\partial w_{22}^1} = x_2 \checkmark$$

say $a_{11} = a_{12}$
 $\bar{z}_{11} = \bar{z}_{12}$ } So, because of that, (all initialized value (2nd), sigmoid@HL)

even 1000 neurons in (HL) also it acts as a single node; make Linear model v

- In short, Non-linearly not existing in (HL);
- Its behavior like a perceptron;

Case ② :- Non-zero Constant Value :-

$$w = 0.5$$

$$b = 0.5$$

$$a_{11} = \max(0, \bar{z}_{11}) \neq 0$$

$$\bar{z}_{11} = w_1^T x_1 + w_2^T x_2 + b_{11} \neq 0$$

likewise $\begin{cases} a_{12} \\ \bar{z}_{12} \end{cases} \neq 0$

but

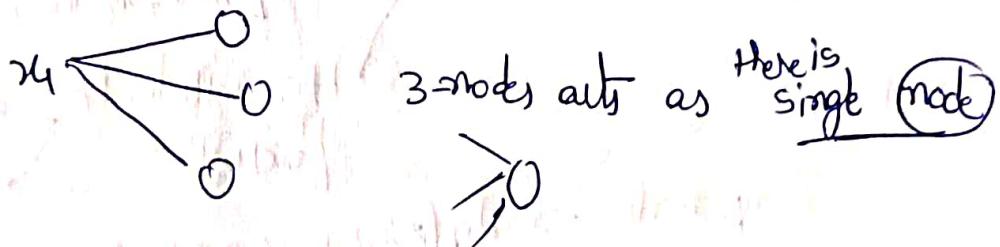
$$\begin{cases} \bar{z}_{11} = \bar{z}_{12} \\ a_{11} = a_{12} \end{cases}$$

because of

$$w = b$$

$$\frac{\partial L}{\partial w_{11}}, \frac{\partial L}{\partial w_{12}}, \frac{\partial L}{\partial w_{21}}, \frac{\partial L}{\partial w_{22}}$$

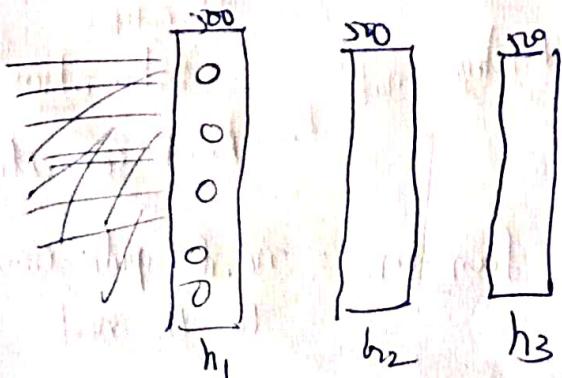
Same as above Case;



So we have only one option i.e., randomly initial the w, b;

Case ③ :- Random Initialization : ~~case 1~~ ① small Random values - Tanh/Sigmoid
 ② large R.V's;

Tanc = Rows = 1000 & Column = 500



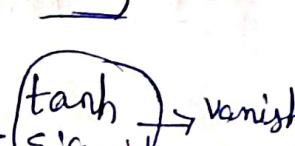
np.random.randn(500, 500) $\times 0.01$

gives small values @ kept $b=0$

activation = $\epsilon(w, x)$ where $x_i = 0$ to 1 means normalised;
 weight like 60 small $= 0.007$
 gives small number;

say Tanh is over AF' \Rightarrow means Activation function gives a very low value; i.e. o/p = 0; this causes vanishing gradient problem;

— say sigmoid is almost same, it also causes, ...;

— In summary as you used small (b, w) like 0.0007  \rightarrow vanishing problem

website? www.deeplearning.ai/ai-notes/initialization

What Can be done

- Heuristics (Jugad) → practical solution

Intuition :- take small weights \rightarrow

卷之三

O O O O

250
notes

For large weights:-

np. random, $\sigma_{\text{rand}}(250, 250) \times 0.01$

np.random.rand(250,250) \times 0.01 ✓ updated uniform

n = A hidden layer node

$$SD = \frac{1}{\sqrt{250}}$$

Yavien Smith

(normal)

✓ Jan

fan_in = no. of inputs coming
to the node

④ If random variable $(\bar{x}, \hat{\sigma})$, $\sqrt{\frac{1}{n-2}}$ ✓

~~GO~~
~~GO~~

(d)

Sum. resp. mode like $\rightarrow \sqrt{\frac{2}{f_{\text{min}} g_n + f_{\text{max}} h_n}}$

⑤ He Normal :- $\sqrt{\frac{2}{f_{\text{min}} g_n}} \rightarrow$ If you are walking with Pell ✓ (test)

uniform version means :- based on uniform distribution;

① Xavier Uniform :- $[-\text{limit}, \text{limit}]$ where limit = $\sqrt{\frac{6}{f_{\text{min}} g_n + f_{\text{max}} h_n}}$

② He Uniform :- $[-\text{limit}, \text{limit}]$ where limit = $\sqrt{\frac{6}{f_{\text{min}} g_n}}$ ✓