

Null space

- Kernel of a matrix
- when a matrix A is multiplied with a vector x it gives

$$A\vec{x} = 0$$

- Null space is the subspace of a matrix.

Q. $A = \begin{bmatrix} 1 & 1 \\ 2 & 2 \end{bmatrix}$

Solve for null space using $Ax=0$.

$$x_1 + x_2 = 0 \quad \text{--- (1)}$$

$$2x_1 + 2x_2 = 0 \quad \text{--- (2)}$$

from this $x_1 = -x_2$

Let $x_1 = c$, $x_2 = -c$

$$\vec{x} = \begin{bmatrix} c \\ -c \end{bmatrix}$$

- Q. Find out null space of A .

$$A = \begin{bmatrix} -1 & 3 & 2 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix}$$

Ans.

$$\begin{bmatrix} -1 & 3 & 2 & | & x_1 \\ 1 & 1 & 0 & | & x_2 \\ 1 & 1 & 0 & | & x_3 \end{bmatrix} = 0$$

$$R_3 \rightarrow R_3 - R_2, \quad R_2 = R_2 + R_1$$

$$\begin{bmatrix} -1 & 3 & 2 & | & 0 \\ 0 & 4 & 2 & | & 0 \\ 0 & 0 & 0 & | & 0 \end{bmatrix}$$

$$R_1 = -R$$

$$R_2 = R_2/4$$

$$\begin{bmatrix} 1 & -3 & -2 \\ 0 & 1 & 0.5 \\ 0 & 0 & 0 \end{bmatrix}$$

$$R_1 = R_1 + 3R_2$$

$$\begin{bmatrix} 1 & 0 & -0.5 \\ 0 & 1 & 0.5 \\ 0 & 0 & 0 \end{bmatrix}$$

$$x_1 - 0.5x_3 = 0$$

$$x_2 + 0.5x_3 = 0$$

Eigen value & Eigen vector
 ↓
 feature selection

PCA (Principal component analysis)

$$\text{Eg: } A = \begin{bmatrix} 10 & 5 & -10 \\ 0 & 22 & 16 \\ 0 & -9 & -2 \end{bmatrix}$$

$$x_A = \begin{bmatrix} -5 \\ -4 \\ 3 \end{bmatrix}, \quad x_B = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad x_C = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$$\text{Soln: } x_A = \begin{bmatrix} -50 \\ -40 \\ 30 \end{bmatrix}, \quad x_B = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}; \quad x_C = \begin{bmatrix} -5 \\ 34 \\ 11 \end{bmatrix}$$

$$AX_A = 10x_A \quad ; \quad AX_B = 0x_B$$

$$Ax = \lambda x$$

→ eigen value

$$A \neq x \quad (\vec{x} \neq 0)$$

$\lambda = 2$

Ch-2

* Probab
-finc→ App
hean
→ Bene
mou* Rand
real
rand

→ Role

→ Set

* For the first two cases we can say write $AX = \lambda x$, where λ is a scalar, when this eqⁿ holds for x & λ , the scalar is called as the eigen value of A .

* For case 1 $\lambda = 10$, and

$$\text{case 2 } \lambda = 0$$

* when $A\vec{x} = \lambda \vec{x}$, $\vec{x} \neq 0$, we can say \vec{x} is an eigen vector of matrix A .

Q) Find out the eigen vector for A :

$$A = \begin{bmatrix} 4 & 1 \\ 2 & 3 \end{bmatrix}$$

Sol.:

$$Ax = \lambda x$$

$$(A - \lambda I)x = 0$$

$$\Rightarrow \begin{vmatrix} 4-\lambda & 1 \\ 2 & 3-\lambda \end{vmatrix} - x \begin{vmatrix} 1 & 0 \\ 0 & 1 \end{vmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0$$

$$\Rightarrow -(4-\lambda)(3-\lambda) - 2 = 0$$

$$\Rightarrow 12 - 4\lambda - 3\lambda + \lambda^2 = 0$$

$$\Rightarrow \lambda^2 - 7\lambda + 12 = 0$$

$$\Rightarrow (\lambda - 3)(\lambda - 4) = 0$$

$$\Rightarrow \lambda = 3, 4$$

$$\left(\begin{bmatrix} 4 & 1 \\ 2 & 3 \end{bmatrix} - \begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix} \right) \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0$$

$$\begin{bmatrix} -1 & 1 \\ 2 & -2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0$$

$$+x_1 = +x_2$$

$$\text{let } x_1 = 1 \quad \vec{x} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$x_2 = 1$$

$$\begin{aligned} \lambda = 2 \\ \left\{ \begin{bmatrix} 4 & 1 \\ 2 & 3 \end{bmatrix} - \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \right\} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = 0 \\ = \begin{bmatrix} 2 & 1 \\ 2 & 1 \end{bmatrix} \end{aligned}$$

$$2x_1 + x_2 = 0$$

$$\Rightarrow x_1 = -x_2/2.$$

Let $x_2 = 1$

$$x = \begin{bmatrix} -\frac{1}{2} \\ 1 \end{bmatrix}$$

Ch-2 Fundamentals of probability

- * Probability: field of mathematics that quantifies uncertainty w.r.t event.
- App? : image processing, NLP, medical data handling.
- Benefits: building robust model, handling noisy ~~data~~ I/P, decision making.
- * Random variable: mathematical rep" of a variable, whose value is generated by randomness.
- Role of random no. in NN:
 - weight initialization
 - stochastic gradient descent
 - dropout.
- It improves convergence of model.

expectation

$E(x)$ - let x be input distr"

$E(x) = \text{input distribution multiplied by its prob.}$

$$= \{x_1 \cdot p_1\} + \{x_2 \cdot p_2\} + \dots + \{x_n \cdot p_n\}.$$

(Q)

Ans.

- (Q). A pair of 6 side dice be rolled. Let x be the random variable representing the no. off on dice.

Ans.

$$x = 1, 2, 3, 4, 5, 6$$

$$P(x) = \frac{1}{6}$$

$$E(x) = \frac{1}{6} \cdot 1 + \frac{1}{6} \cdot 2 + \dots + \frac{1}{6} \cdot 6 = \frac{21}{6} = \frac{7}{2}$$

→

→

- (Q) A game involves rolling a fair 6 sided die. If the no. is even you win \$10 else loose \$5. Let x be random variable representing amt. won or loss. Find out $E(x)$.

Ans.

$$x = \{10, 5\}$$

$$E(x) = \frac{1}{2} \times 10 + \frac{1}{2} \times 5 = 7.5$$

Ans

$$E(\text{win}) = 5 \quad ; \quad E(\text{loss}) = -5$$

$$E(x) = 5 - 2.5 = 2.5$$

Variable & standard deviation

$$\text{var}(x^2) = E(x - \bar{x})^2 = (x_1 - \bar{x})^2 + \dots$$

where, $x = \text{input distr"}$

, $\bar{x} = \text{mean}$

standard (σ) = $\sqrt{\text{var}}$

Q) Find out popⁿ variance of data 5, 7, 9, 10, 14, 15. Just find var & std.

Ans.

$$\bar{x} = \frac{1}{6}(5+7+9+10+14+15)$$

$$\sigma^2 = 12.660 ; \sigma = 3.55$$

* Entropy ($H(x)$)

$$= -P(x) \log_2 P(x)$$

$P(x)$ = prob of event x .

↓ uncertainty of event x .

- It measures avg. level of uncertainty present in random variable placed outcome.
- The more unpredictable an event is the higher will be entropy.

coin toss = $-[P(H) \log_2 P(H) + P(T) \log_2 P(T)]$

Q). $P(H) = 0.6 , P(T) = 0.4$

Find $H(x)$.

Ans

$$H(x) = 0.9709$$

Cross Entropy Loss

It measures the diff b/w true label & predicted class label probability during data classification.

$$L(y, y') = -\sum y * \log(y')$$

y = actual class

y' = predicted class.

Ques.

Consider a 2 class classification problem. The true label & predicted label of classifier are given as:

True label (y)	Predicted label (y')
Sample 1 [1, 0]	[0.7, 0.3]
Sample 2 [0, 1]	[0.4, 0.6]

Find out the cross entropy loss.

Ans. Cross entropy = $-\sum y * \log(y')$

Sample 1 = $-(1 * \log(0.7) + 0 * \log(0.3)) = 0.5145$

Sample 2 = $0 * \log(0.4) + 1 * \log(0.6) = 0.7369$

$L = L_1 + L_2 = 0.5145 + 0.7369$

Cross entropy loss = 1.2514

Q) There are 3 containers containing shapes such as triangle & circle. The prob. of picking a shape from these 3 containers are given as:

Container No.	Picking A	Picking O
C1	$\frac{26}{30}$	$\frac{4}{30}$
C2	$\frac{14}{30}$	$\frac{16}{30}$
C3	$\frac{1}{30}$	$\frac{29}{30}$

Apply cross entropy loss to find out from which container the picking of shapes are container & why?

Ans.

Container C1 = 2.84.

$$C2 = -0.089$$

$$C3 = -0.161$$

C3 will be preferred.

Q).

consider a 3 class classification problem, considering cat, dog & bird. There are 2 cases for which the labels are: True label [1, 0, 0]

case 1: Predicted label [0.7, 0.2, 0.1]

case 2: Predicted label [0.9, 0.05, 0.05].

Evaluate which will be predicted.

Ans.

$$L(Y, Y') = 0.5145$$

$$L_e(Y, Y') = 0.1520$$

L_e loss $L_2 < L_1$. Thus, M_2 is better.

KL divergence

it is the difference b/w two prob. distribution. Represent as:

$$D(P||Q) = \sum P(x) \log_2 \left(\frac{P(x)}{Q(x)} \right)$$

where P & Q are prob. distributions, of random variable.

Q)

Let P & Q are given as

$$P(x) = [0.4, 0.3, 0.3]$$

$$Q(x) = [0.2, 0.5, 0.3]$$

Find out KL divergence D .

Ans.

~~Ans~~ $D = 0.179$

3. The Neural Networks (NN)

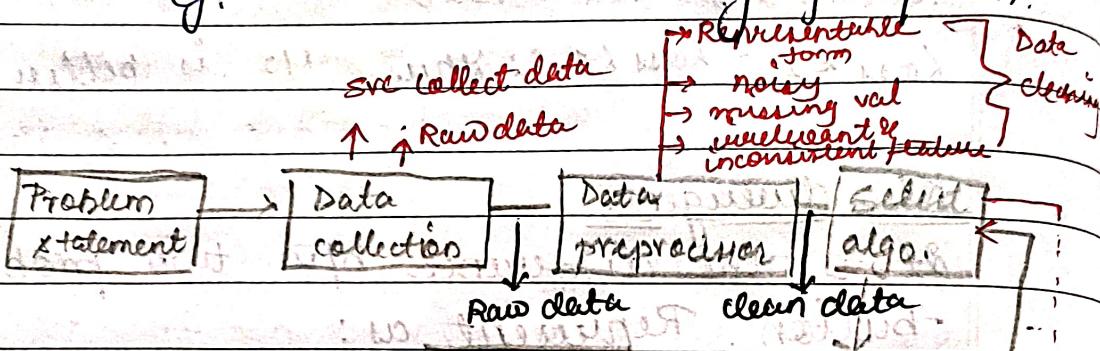
Defⁿ of Intelligence

- learn
- understand
- adapt

limitation of traditional system

- set of defined "instincts"
- efficient in performing v. op.ⁿ.
- arithmetic

eg: understand handwriting of a person.



Q: Identify hidden pattern

→ 1. classificⁿ prop.

if accuracy is less.

→ Support vector machines

→ Decision tree

→ Random forest

→ Linear Regression

→ ANN

→ Naive Bayes classifier

eg:

* A machine learning technique is an algo. that is able to learn from data.

* ML enables us to tackle tasks that are too diff. to solve with the fixed program written by humans.

- * ML tasks are usually described in terms of how a model should process an example.

Types of Machine Learning

- ① Supervised learning - labelled data.
- ② Unsupervised learning - unlabelled data.
- ③ Semi-supervised learning - few of both.
- ④ Reinforcement learning - completely diff. from one ① & ②.

① Supervised Learning

- It is the study of algo. that use labelled data in which every sample is associate with class label.
- Used to discover hidden patterns.

Eg: Classification, Regression

② Unsupervised Learning:

- Involves difficult tasks as pre defined labels are not availables

Eg: Clustering, Anomaly Detection, Dimensionality Reduction

③ Semi-supervised learning

- few samples of labelled & few of unlabelled.
- It fills the gap b/w ① & ②. It leverages the strength of both approaches by training using labelled data with unlabelled.

④ Reinforcement Learning

- It is completely different from first two.
- It adopts an agent as it interacts with surroundings & learns.
- The agent learning is through experiment. Thus maybe rewarded when performs desired task.

eg:

Recommend model; gaming AI.

Confusion matrix

Actual outcome

Predicted outcome

1	1 (TP)
---	--------

0 (TN)

0 (FP)

0	1 (FN)
---	--------

⇒ True positive (TP) - +ve sample correctly identified

⇒ True negative (TN) - -ve sample correctly identified

⇒ False positive (FP) - -ve sample incorrectly identified

⇒ False negative (FN) - +ve sample incorrectly identified

		Predicted class	
		0	1
Actual class	0	TP	FN
	1	FP	TN

① High value of TP & TN = Good performance

② High value of FP = Over +ve prediction (unnecessary action)

③ High value of FN = under predict of positive samples.

(a) A confusion matrix w.r.t. COVID dataset is as follows : $TP = 30$, $TN = 950$, $PP = 15$ and $PN = 5$. How to interpret this outcome.

Ans.

$$\text{Total no. of samples} = 1000$$

$$\begin{aligned}\text{Total no. of positive prediction} &= TP + PP \\ &= 30 + 15 = 45\end{aligned}$$

$$\begin{aligned}\text{Total no. of negative prediction} &= TN + PN \\ &= 950 + 5 = 955\end{aligned}$$

(b) Find out precision, recall, accuracy, sensitivity, F1 score.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} = \frac{980}{1000} = 0.98$$

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{30}{45} = 0.667$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{30}{35} = 0.857$$

$$\text{F1 score} = 0.75 = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\text{Error} = 1 - \text{accuracy} = 0.02$$

Deep Learning

- * It is a sub-field of machine learning.
- * It is a type of ANN in which automatic feature selection takes place.
- * Explicit feature selection is not reqd.

Advantages

- * DL can learn complex patterns, especially high dimensional prob.
- * Has the ability to learn from raw data such as image or audio.

Requirements

- * Huge amt. of data for model training.
- * Requires GPU for modelling.
- * Requires more training & processing time.

types of DL model (wrt the type of data it investigates)

① Convolutional Neural Network (CNN)

→ image & video data

② Recurrent Neural Network (RNN)

→ used for sequential data learning
e.g.: time series data (temporal data)

Data

temporal long term
(time comp) (No time comp).

③ Long short term memory (LSTM)

→ Advance version of RNN → temporal data.
→ used for prediction → NLP

④ Generative Adversarial Network (GAN)

→ generate synthetic data (new data)

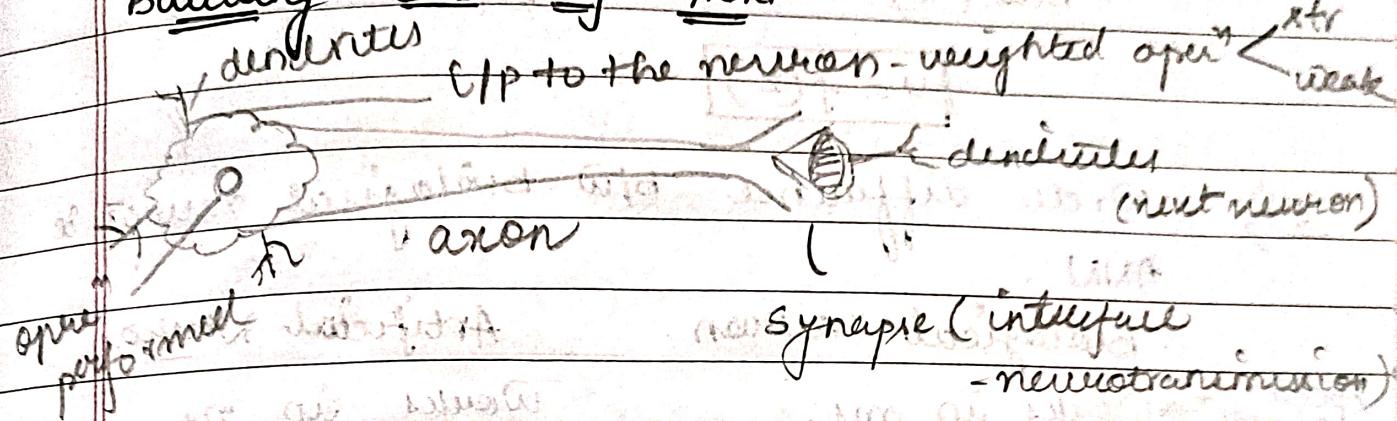
⑤ Auto encoder

→ noise reduction → produce new dimension
→ compress data

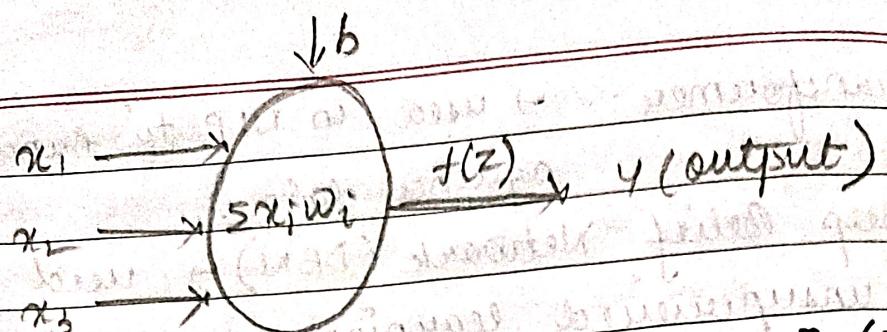
⑥ Transformer: → used in RNN to handle sequence data.

⑦ Deep Belief Network (DBN) → used for unsupervised learning.
Eg: clustering, anomaly detection.

Building Blocks of ANN



- * ~~Revered~~ Neuron receives i/p through mini dendrites.
- * Input is either strengthened or weakened depending on how & which dendrite will participate in info. sharing.
- * After being weighted by the strength of respective connecⁿ, the i/p are summed in the cell body.
- * The summed value is transformed into a new signal & propagates to axon to other neurons.
- * Here, synapses are present, which act as an interface b/w two neurons & helps in neurotransmission.



→ input (x) → activation funcⁿ (z)

→ weight (w)

→ bias (b)

$$z = \sum_{i=1}^n x_i w_i + b$$

$$y = f(z)$$

Q.) Write difference b/w biological neuron & ANN.

Biological Neuron.

Artificial Neuron

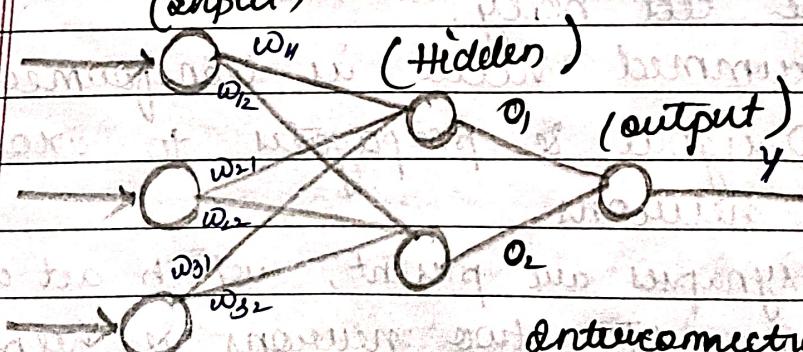
speed of oper ⁿ	Works in ms.	Works in ns.
weighted i/p proc.	Parallel	Parallel
storage cap.	unlimited	limited.
size	Huge	comparatively less.
complexity	More	comparatively less.
Fault tolerance	hardly any loss	data loss, n/w failure

Layered architecture of ANN

(input)

(Hidden)

(output)



Interconnection

(takes i/p from dataset) ① Input layer

(weights are defined) ② Hidden layer (optional)

③ Output layer

① Input Layer : Top layer in ANN. Here, the input data in the form of image pixels, numbers, alphabets is received.

② Hidden Layer : Is located in centre of ANN. Runs to variety of mathematical operations over input data to identify patterns.

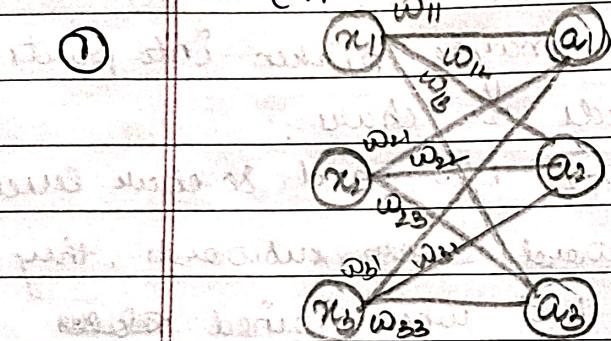
③ Output Layer : The computation in middle layer enable us to generate desired output in the output layer.

- * ANN employs a variety of mathematical processing stages. It has a no. of units arranged in no. of layers. A single unit is known as neuron.
- * The input layers input units take a variety of input from the outside environment.
- * The data then travels to the hidden units, which transforms it into a format that output units may use it.
- * "Interconnection" maybe defined as the way of processing elements in ANN, are connected to each other.
- * These arrangements have always two layers common in every n/w, that are : input layer (buffers the i/p signal) & output layer (generates o/p of n/w).

DL UP

- * The third layer is ~~out~~ hidden layer, in which neurons are neither kept in input layer nor in output layer.
- * These layers are hidden from people who are interfacing from the system. Thus, it is completely a blackbox.
- * No. of hidden layers & no. of neurons are used in ANN to introduce non-linearity.

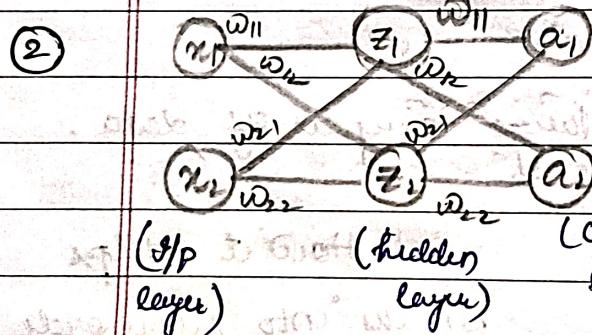
(S/P Layer) (O/P Layer)



single layer feed forward
neural n/w.

(1 layer FFNN)

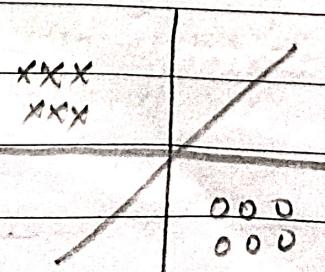
- linear neuron



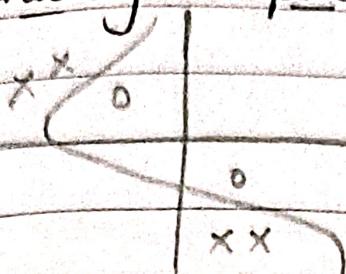
Multilayer feed forward
neural n/w

m-layer FFNN

- * Types of problem
- ① linearly separable problem

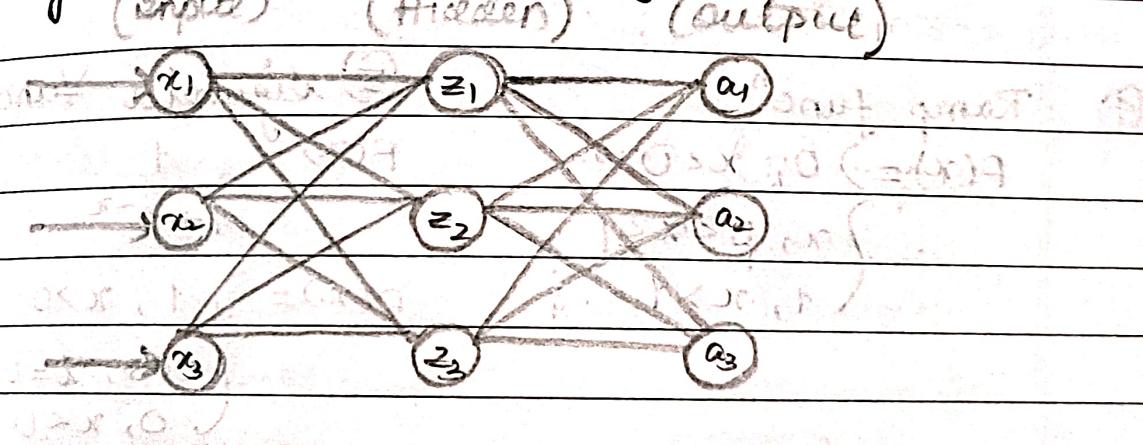


Non linearly separable problem

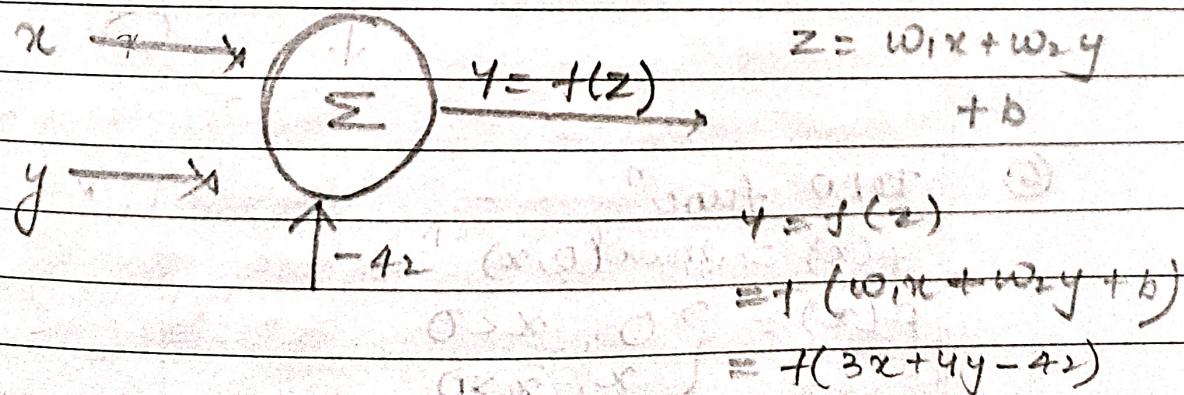


- * linear neurons never use hidden layer in NN.
- * in realtime "itua" problems are non-linear.
- * hidden layers are helpful to extract imp. features from data.
- * it is used to learn complex reln'ships & direct the output towards the desired o/p.

Q) Design a feed forward NN with 3-layers. Each layer contains 3 no. of neurons.



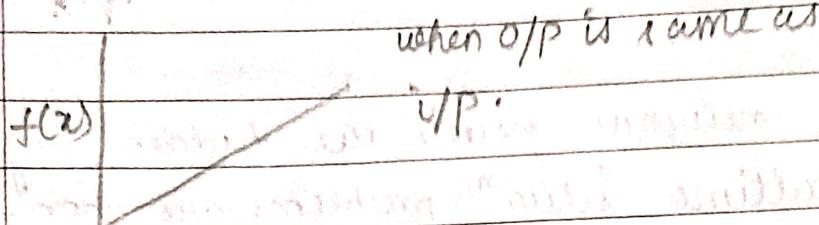
Q) Let $z = 3x + 4y - 42$. Design a NN for this.



Activation function

① Identity func"

$$f(x) = x$$

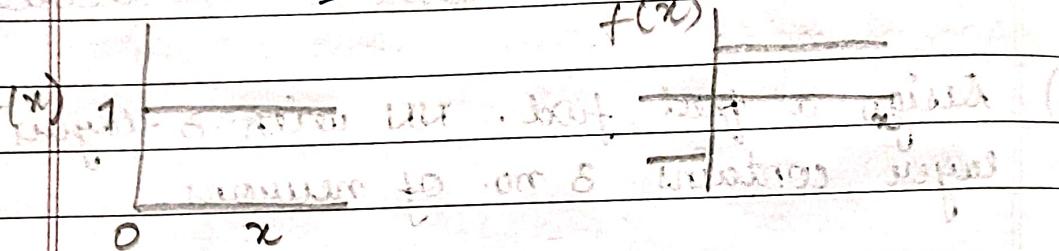


② Binary func"

$$f(x) = \begin{cases} 0, & x < 0 \\ 1, & x \geq 0 \end{cases}$$

③ Bipolar step func"

$$f(x) = \begin{cases} 1, & x \geq \theta \\ -1, & x < \theta \end{cases}$$



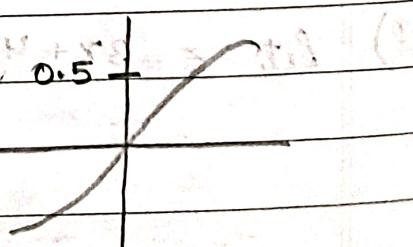
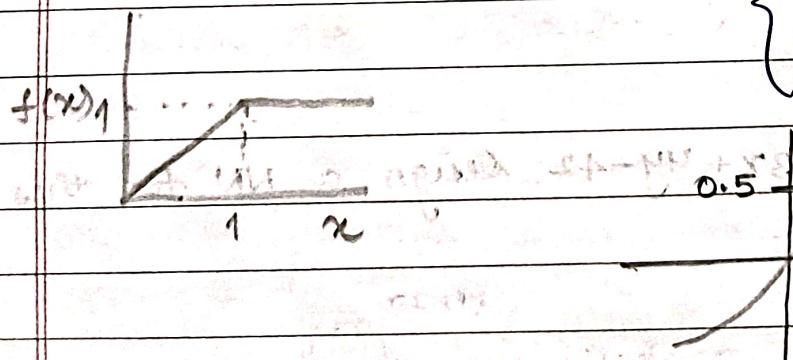
④ Ramp func"

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

⑤ Sigmoid Func

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

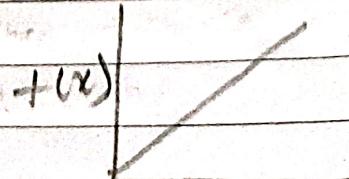
$$f(x) = \begin{cases} 1, & x > 0 \\ 0.5, & x = 0 \\ 0, & x < 0 \end{cases}$$



⑥ ReLU func"

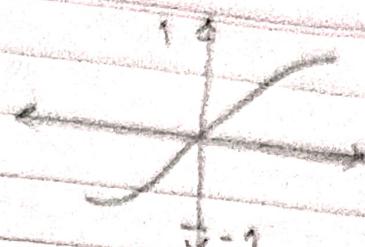
$$f(x) = \max(0, x)$$

$$f(x) = \begin{cases} 0, & x < 0 \\ x, & x \geq 0 \end{cases}$$



③ Tanh function

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



- ④ Consider a label dataset of size 5×4 . The data descrip " is as follows:

Data =	1	0	0	1
	1	0	1	1
	0	1	1	0
	0	1	0	0
	1	1	1	1

The last row represents class label. given weights of NN are (w_0, w_1, w_2, w_3) . given activation func " $P(z) = \begin{cases} 0, & z \leq 0.5 \\ 1, & z > 0.5 \end{cases}$

Find out classific " performance of given n/w.

Ans. $z_1 = 0.5 \times 1 = 0.5 = 0 = 0$

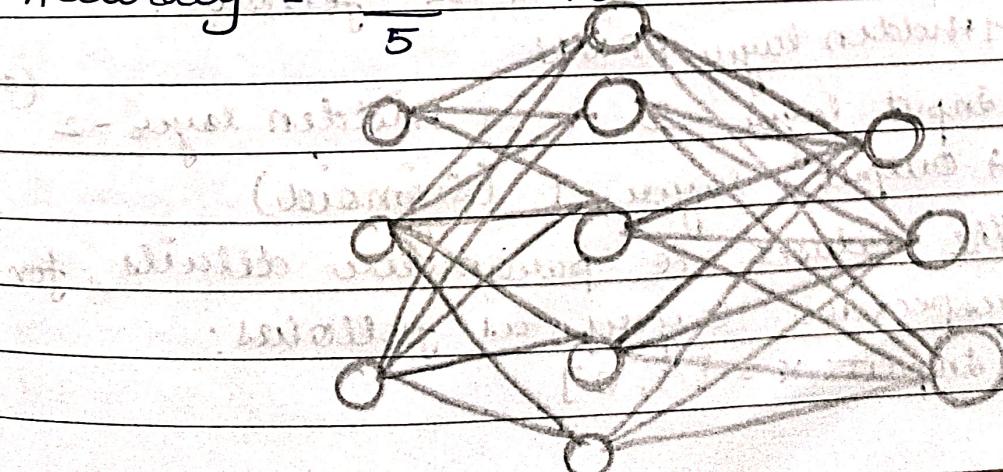
$z_2 = 0.5 + 0.6 = 1.1 = 1$

$z_3 = 0.6 + 0.4 = 1 = 1$

$z_4 = 0.6 = 1$

$z_5 = 1.5 = 1$

Accuracy = $\frac{2}{5} = 40\%$



0.5, -0.6, 0.4



$x_1 \rightarrow \text{circle} \rightarrow \text{circle} \rightarrow y = f(z)$

$$w_1 = 0.5$$

$$w_2 = -0.6 \quad f(z) = \begin{cases} 0, & z \leq 0 \\ 1, & z > 0 \end{cases}$$

$$w_3 = 0.4$$

(Q) Solve the given "classification" problem & find
actual precision & recall. $y = f(z) = \begin{cases} 1, & z > 0 \\ 0, & z \leq 0 \end{cases}$

Input (x) Actual class $\frac{PC}{O}$

0.1	$2x-1 (-0.8)$	0	0
-----	---------------	---	---

0.4	-0.2	0	0
-----	------	---	---

0.5	1	1	1
-----	---	---	---

0.6	0.2	1	1
-----	-----	---	---

0.9	0.8	1	1
-----	-----	---	---

Any Precision = 1 ; Recall = 1.

(Q) Evaluate the performance of given n/w.

The n/w structure is as follows:

1-hidden layer FFNN

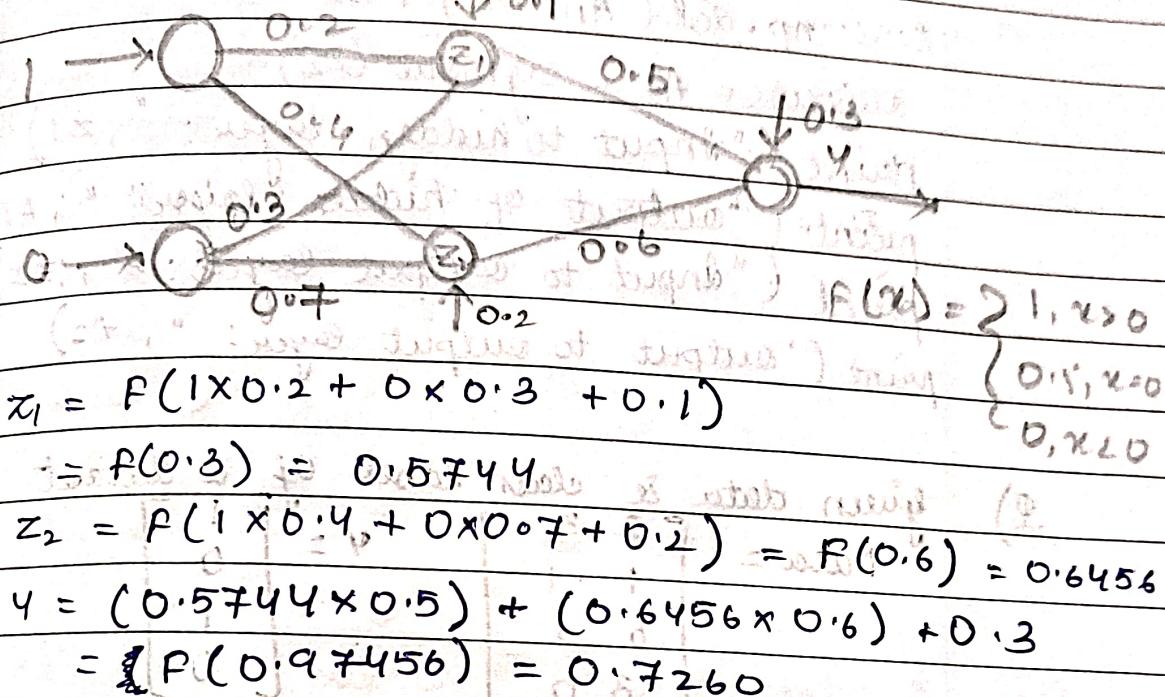
→ Input layer - 2 → Hidden layer - 2 (sigmoid)

→ Output layer - 1 (sigmoid)

The dataset & parameter details for
respective layers as follows:

→ Input $x = [1, 0]$

- weights [input → hidden layer] : $w^1 = \begin{bmatrix} 0.2 & 0.4 \\ 0.3 & 0.7 \end{bmatrix}$
- Bias for hidden layer, $b^1 = [0.1, 0.2]$.
- weights [hidden → output] : $w^2 = [0.5, 0.6]$
- Bias for output layer, $b^2 = 0.3$.



Python code

```

import numpy as np
# define sigmoid func
def sigmoid(x):
    return 1 / (1 + np.exp(-x))
# initialize data
x = np.array([1, 0])
# define weights & bias
w1 = np.array([[0.2, 0.4], [0.3, 0.7]])
# weights [input → hidden]
b1 = np.array([0.1, 0.2])
w2 = np.array([0.5, 0.6])
b2 = 0.3

```

```

# Initialization over
# Step 1 :  $z_1 = xw_1 + b_1$ 
 $z_1 = \text{np.dot}(x, w_1) + b_1$ 
# Step 2 :  $A_1 = \text{sigmoid}(z_1)$ 
 $A_1 = \text{sigmoid}(z_1)$ 
 $z_2 = \text{np.dot}(A_1, w_2) + b_2$ 
def  $A_2 = \text{sigmoid}(z_2)$ 
print("Input to hidden layer: ", z1)
print("Output of hidden layer: ", A1)
print("Input to output layer: ", z2)
print("Output to output layer: ", A2)

```

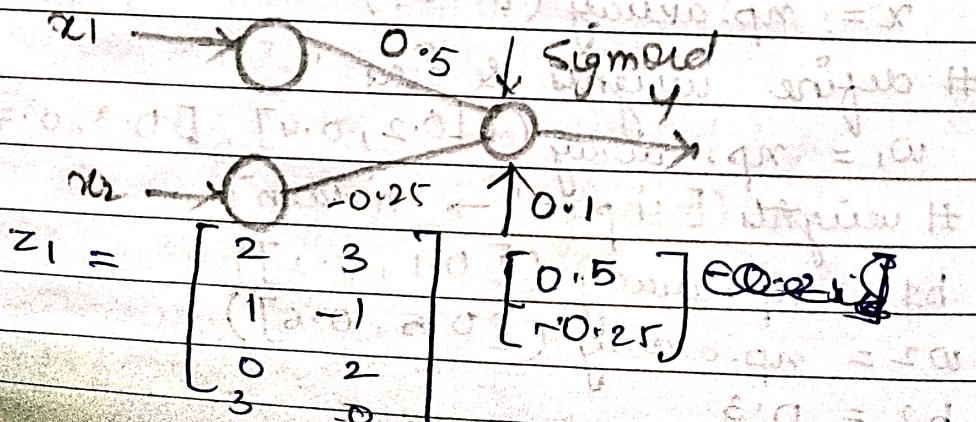
(Q) Given data & class label of a dataset

$$\text{Data} = \begin{bmatrix} 2 & 3 \\ 1 & -1 \\ 0 & 2 \\ 3 & 0 \end{bmatrix} \quad \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

a) Design the PFNN with the following n/w details. $w = \begin{bmatrix} 0.5 \\ -0.25 \end{bmatrix}$, $b = 0.1$, AF in O/p layer = Sigmoid

b) Evaluate the performance of binary classifier in form of eval criteria such as accuracy, error, precision & recall.

Ans.



$$= \begin{bmatrix} 0.25 \\ 0.75 \\ -0.5 \\ 1.5 \end{bmatrix} + 0.1 = \begin{bmatrix} 0.35 \\ 0.85 \\ -0.4 \\ 1.6 \end{bmatrix}$$

$A(z) = \begin{bmatrix} 0.5866 \\ 0.7005 \\ 0.4613 \\ 0.8320 \end{bmatrix}$

$\text{TP} = 1, \text{FP} = 2$
 $\text{TN} = 0, \text{FN} = 1$
 $\text{PP} = 1, \text{PN} = 1$

$$\text{Accuracy} = \frac{\text{TP}}{\text{TP} + \text{FP} + \text{TN} + \text{FN}} = \frac{1}{4}$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{1}{3}$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{1}{2}$$

$$\text{Error} = 1 - 0.25 = 0.75$$

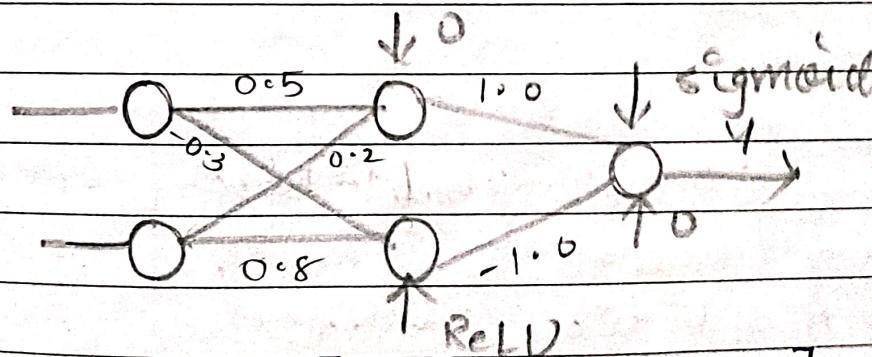
(Q) Consider 2 layer FFNN. Input, $x = [1, 2]$

Input layer details as follows:

$$w_F \begin{bmatrix} 0.5 & -0.3 \\ 0.2 & 0.8 \end{bmatrix}, b_1 = 0, \text{AF} = \text{ReLU}$$

$$w_2 = \begin{bmatrix} 1.0 \\ -1.0 \end{bmatrix}, b_2 = 0, \text{AF} = \text{Sigmoid}$$

Ans.



$$z_1 = \begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} 0.5 & -0.3 \\ 0.2 & 0.8 \end{bmatrix} = \begin{bmatrix} 0.9 & 1.3 \end{bmatrix}$$

$$z_2 = \begin{bmatrix} 0.9 & 1.3 \end{bmatrix} \begin{bmatrix} 1.0 \\ -1.0 \end{bmatrix} = \begin{bmatrix} -0.4 \\ 0.404 \end{bmatrix}$$

Q Binary classification problem

* Q over single layer PFNN

a) Design a PFNN with the following n/w details

b) evaluate the performance of binary classification problem in terms of evaluation criteria like accuracy, error, precision, F1 score.

$$x = \begin{bmatrix} 2 & 1 \\ 1 & -1 \\ 0 & 2 \\ 3 & 0 \end{bmatrix}, AC = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

$$w_1 = \begin{bmatrix} 0.2 & -0.5 \\ 0.3 & 0.8 \end{bmatrix}, b_1 = [0.1, -0.2]$$

$$w_2 = \begin{bmatrix} 0.4 \\ -0.6 \end{bmatrix}, b_2 = 0.05$$

$$AF_1 = \text{ReLU}, AF_2 = \text{sigmoid}$$

Ans.

$$\text{Input} = \begin{bmatrix} 2 & 1 \\ 1 & -1 \\ 0 & 2 \\ 3 & 0 \end{bmatrix}$$

Q.) $x = \begin{bmatrix} 1 & 2 \\ 2 & 0 \\ 0 & 1 \end{bmatrix}$ $y = \begin{bmatrix} 1 \\ 3 \\ 2 \end{bmatrix}$

hidden layer wt. & bias.

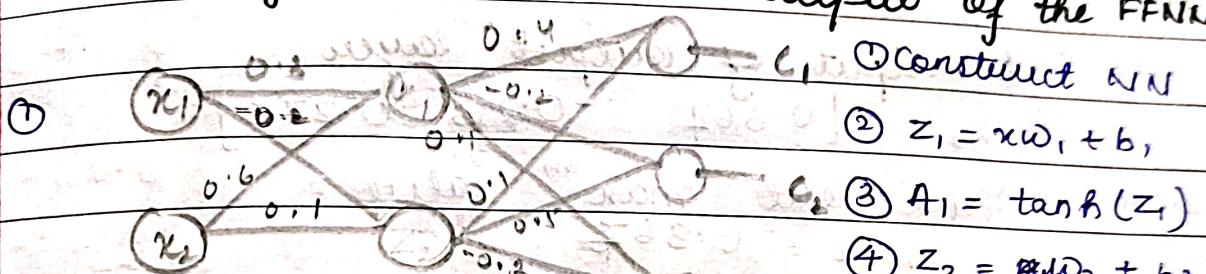
$$w_1 = \begin{bmatrix} 0.3 & -0.2 \\ 0.6 & 0.1 \end{bmatrix} \quad b_1 = \begin{bmatrix} 0 \\ 0.5 \end{bmatrix} \quad AF_1 = \tanh$$

output layer wt. & bias.

$$w_2 = \begin{bmatrix} 0.4 & -0.2 & 0.1 \\ 0.1 & 0.5 & -0.3 \end{bmatrix} \quad b_2 = \begin{bmatrix} 0 \\ -0.1 \\ 0.2 \end{bmatrix}$$

$AF_2 = \text{Softmax}$. Evaluate output of the FFNN.

An.



$$\textcircled{2} \quad z_1 = \begin{bmatrix} 1 & 2 \\ 2 & 0 \\ 0 & 1 \end{bmatrix} \quad \begin{bmatrix} 0.3 & -0.2 \\ 0.6 & 0.1 \end{bmatrix} + \begin{bmatrix} 0 \\ 0.5 \end{bmatrix}$$

$$= \begin{bmatrix} 1.5 & 0 \\ 0.6 & -0.4 \\ 0.6 & 0.1 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 0.5 \end{bmatrix} \quad \begin{bmatrix} 0. \\ 0.05 \end{bmatrix}$$

$$= \begin{bmatrix} 1.5 & 0.05 \\ 0.6 & -0.35 \\ 0.6 & 0.15 \end{bmatrix} \quad f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\textcircled{3} \quad A_1 = \begin{bmatrix} \tanh(1.5) & -\tanh(0.05) \\ \tanh(0.6) & \tanh(-0.35) \\ \tanh(0.6) & \tanh(0.15) \end{bmatrix} \quad \begin{bmatrix} 0.905 \\ 0.537 \\ 0.537 \end{bmatrix} \quad \begin{bmatrix} 0.049 \\ -0.336 \\ 0.148 \end{bmatrix}$$

output of hidden layer.

$$④ z_2 = \begin{bmatrix} 0.905 & -0.049 \\ 0.537 & -0.336 \\ 0.537 & 0.148 \end{bmatrix} + \begin{bmatrix} 0.4 & -0.2 & 0.1 \\ 0.1 & 0.5 & 0.3 \\ 0 & -0.1 & 0.2 \end{bmatrix}$$

$$\begin{aligned} z_2' &= \begin{bmatrix} 0.8669 & -0.176 & 0.075 \\ 0.1812 & -0.275 & 0.154 \\ 0.229 & -0.033 & 0.009 \end{bmatrix} + \begin{bmatrix} 0 \\ -0.1 \\ 0.2 \end{bmatrix} \\ &= \begin{bmatrix} 0.366 & -0.256 & 0.275 \\ 0.181 & -0.375 & 0.354 \\ 0.229 & -0.133 & 0.209 \end{bmatrix} \end{aligned}$$

output of output layer

$$z_2 = [0.366 \quad -0.256 \quad 0.275]$$

Step 1: Find out \max^m value
 $z_{\max} = 0.366$.

Step 2: Shift z_2 , $[z_2'] = z_2 - z_{\max}$
 $z_2' = [0, -0.623, -0.09]$

Step 3: Exponentiate z_2'

$$z_2' = e^{z_2'} = [1, 0.536, 0.912]$$

Step 4: Normalize z_2' ,

$$\text{Sum} = 2.44$$

$$P_1 = \left[\frac{1}{2.44}, \frac{0.536}{2.44}, \frac{0.912}{2.44} \right] = [0.408, 0.218, 0.372]$$

For Row 2:

$$z_2 = [0.181 \quad -0.375 \quad 0.354]$$

$$z_{\max} = 0.354$$

$$z_2' = [-0.173 \quad -0.729 \quad 0]$$

$$e^{z_2'} = [0.8413, 0.482, 1], \text{ Sum} = 2.323$$

$$P_2 = [0.362 \quad -0.207 \quad 0.4304]$$

For Row 3:

$$\begin{aligned} z_2' &= [0, -0.362, -0.02] \\ e^{z_2'} &= [1, 0.696, 0.980] \\ p_3 &= [0.373, 0.260, 0.366], \text{ sum} = 2.676 \end{aligned}$$

Class 1	Class 2	Class 3
0.408	0.219	0.372
0.362	0.207	0.430
0.373	0.260	0.366

* Evaluate which col. has "highest contrib" towards class label.

$$y = \begin{bmatrix} 1 \\ 3 \\ 2 \end{bmatrix}, \quad y' = \begin{bmatrix} 1 \\ 3 \\ 1 \end{bmatrix}$$

Solving Regression Problem

Q.) Find out the performance of given n/w in terms of MSE.

① $x = 2.0, y = 3.5, w = 1.2, b = 0.1$

$$MSE = \frac{1}{2} \sum_{i=1}^n (y_i - y'_i)^2$$

$$x = xw + b = 2 \times 1.2 + 0.1 = 2.5$$

$$MSE = \frac{1}{2} (3.5 - 2.5)^2 = 0.5$$

$$x = \underbrace{\sum_{i=1}^n}_{w=1.2} \times y = 3.5$$

② $x = [1, 2], y = 4, w_1 = [0.5, -0.3]$

$$b_1 = [0.1], AF = \text{ReLU}$$

$$w_2 = \begin{bmatrix} 1 \\ -0.5 \end{bmatrix}, b_2 = 0.2$$

Ques.

$$Z_1 = \begin{bmatrix} 1 & 2 \end{bmatrix}^T \begin{bmatrix} 0.5 & -0.3 \\ 0.2 & 0.4 \end{bmatrix} + \begin{bmatrix} 0.1 \end{bmatrix}$$

$$= \begin{bmatrix} 0.9 & 0.5 \end{bmatrix} + \begin{bmatrix} 0.01 \end{bmatrix} = \begin{bmatrix} 0.9 & 0.6 \end{bmatrix}$$

$$Z_2 = \begin{bmatrix} 0.9 & 0.6 \end{bmatrix} \begin{bmatrix} 1 \\ -0.5 \end{bmatrix} + 0.2$$

$$= \underline{0.8}$$

$$MSE = \underline{0.12}$$

Multi-output based Regression

$$x = [1, 2] \quad y = [3, 1]$$

$$w_1 = \begin{bmatrix} 0.3 & 0.7 \\ -0.2 & 0.4 \end{bmatrix} \quad b_1 = \begin{bmatrix} 0 \\ -0.1 \end{bmatrix}$$

$$w_2 = \begin{bmatrix} 0.5 & -0.1 \\ 0.2 & 0.8 \end{bmatrix} \quad b_2 = \begin{bmatrix} 0.1 & -0.1 \\ -0.2 \end{bmatrix}$$

$$\text{Ans. } Z_1 = x w_1 + b_1$$

$$= \begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} 0.3 & 0.7 \\ -0.2 & 0.4 \end{bmatrix} + \begin{bmatrix} 0 \\ -0.1 \end{bmatrix}$$

$$= \begin{bmatrix} -1.01 & 1.04 \\ 0.099 & 0.885 \end{bmatrix}$$

$$AF_1 = \begin{bmatrix} -0.099 & 0.885 \end{bmatrix}$$

$$Z_2 = A_1 w_2 + b_2$$

$$= \begin{bmatrix} -0.099 & 0.885 \end{bmatrix} \begin{bmatrix} 0.5 & -0.1 \\ 0.2 & 0.8 \end{bmatrix} + \begin{bmatrix} 0.1 & -0.1 \end{bmatrix}$$

$$= \begin{bmatrix} 0.1275 & 0.7179 \\ 0.2275 & 0.5179 \end{bmatrix} + \begin{bmatrix} 0.1 & -0.1 \end{bmatrix}$$

$$MSE = \underline{\frac{1}{2}(3.959)}$$

Thur \Rightarrow Fri \Rightarrow Outcome

Example of regress

How to analyse time series data using FFNN?

Price

Gold Prediction

Eg:

Mon \rightarrow y_{t-2} Tue \rightarrow y_{t-1} Wed \rightarrow y_t Thu \rightarrow y_{t+1} Fri \rightarrow Sat \rightarrow Wed \rightarrow

1

0.5

0.15

0.148

0.1511

$$x = [y_{t-1}, y_{t-2}] = [0.5, 1].$$

$$w_1 = [0.2, -0.1]$$

$$b_1 = 0, w_2 = 0.3, b_2 = 0$$

AF = Sigmoid

$$z_1 = x w_1^\top + b_1$$

$$= [0.5 \ 1] \begin{bmatrix} 0.2 \\ -0.1 \end{bmatrix} + 0 = 0$$

$$A_1 = \sigma(0) = 0.5$$

$$z_2 = A_1 w_2 + b = 0.5 \times 0.3 = 0.15.$$

$$\begin{aligned} \text{Thu } \Rightarrow z_1 &= x w_1 + b & x = [y_t, y_{t-1}] \\ &= [0.15 \ 0.5] \begin{bmatrix} 0.2 \\ -0.1 \end{bmatrix} & = 0.025 \ 0.03 - 0.01 \\ & & = 0.02 \end{aligned}$$

$$A_1 = \sigma(0.02) = 0.495.$$

$$z_2 = 0.495 \times 0.3 = 0.1485$$

$$\begin{aligned} \text{Fri } \Rightarrow z_1 &= [0.148 \ 0.15] \begin{bmatrix} 0.2 \\ -0.1 \end{bmatrix} & = 0.0296 - 0.015 \\ & & = 0.0147 \end{aligned}$$

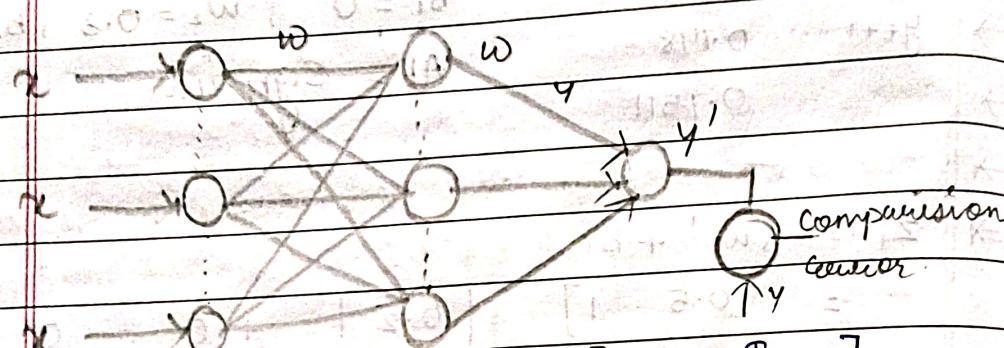
$$A_1 = \sigma(0.0147) = 0.4996 \ 0.503$$

$$z_2 = 0.503 \times 0.3 = 0.1511$$

Outcome: For the next 3 days the output (gold price) varies mostly fixed & a minor fluctuation is observed wrt 0.15.

4. Training FFNN

- Any type of learning in FPNN
- i) Forward propagation (no training is possible).
 - ii) Backward propagation



Forward pass

- linear transformation $[z = xw^T + b]$.

- non linear transformation [Activⁿ Funcⁿ].

⇒ Classification problem - Cross Entropy loss.

⇒ In regression problem - MSE/RMSE.

$loss = \text{compare}(\text{true o/p}, \text{predicted o/p})$

→ loss occurs due to improper selecⁿ of weight & bias. To reduce it take derivatives wrt weights ($\frac{\partial L}{\partial w}$)

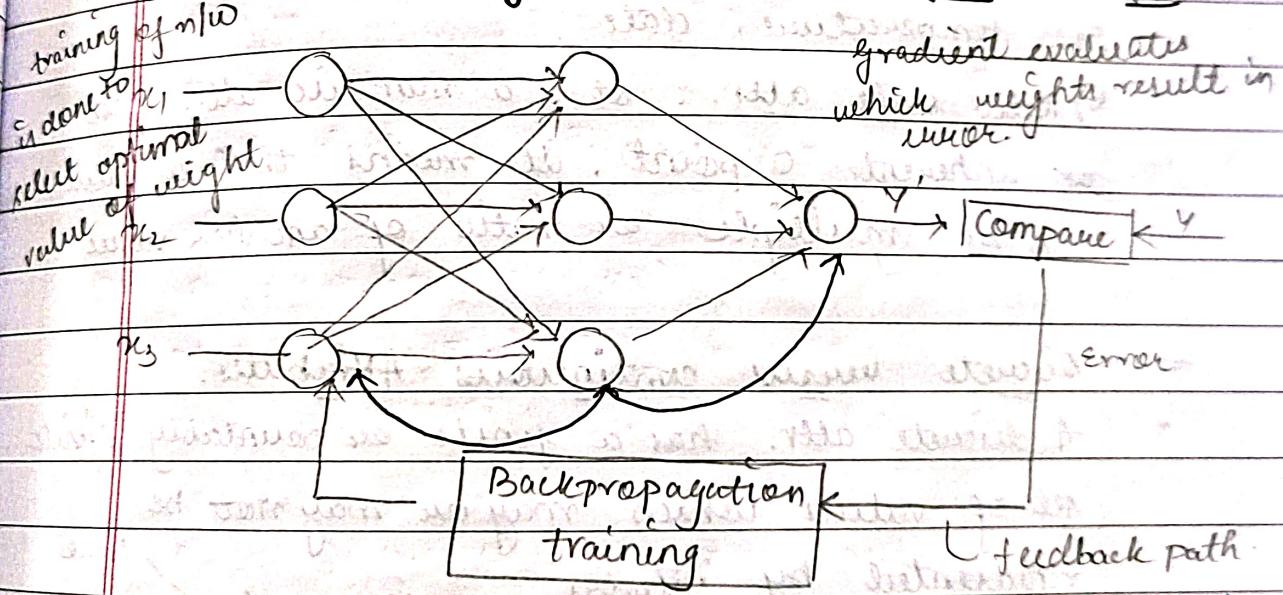
- * Neural n/w is a n/w of interconnected neurons having weight, bias & activⁿ funcⁿ.
- * learning starts from the linear transformⁿ of i/p to non linear transformⁿ using AF.
- * There are 2 types of learning

- * Forward Propagⁿ: i/p are passed to hidden layer with some randomly initialized weights alongwith bias at linear transforⁿ
i.e. $z = xw + b$.

possible).

- * To solve complex prob. non linear transferⁿ is used through AF. The output of linear transformⁿ is supplied to AF i.e.
 $A = f(Z)$ layer
- * everywhere applies these linear followed by non-linear transformⁿ. Hence, o/p of applying AF in o/p layer is our final outcome.
- * In forward propⁿ, as i/p travels from i/p through hidden to o/p layer, no training is possible.

• Backward Propagation Neural Network (BPNN)



- * BP means propagⁿ of error in the backward direction in NN.
- * Its objective is to reduce the diff b/w actual & predicted o/p by adjusting weight & bias.
- * It operates iteratively to adjust weight & bias.
- * In each EPOCH the models adapt these parameters by reducing loss following error gradient.

* BPNN has two phases:

- ① Forward Pass - signal travels from input through hidden to o/p layer following linear & non linear transformation.
- ② Backward pass - error b/w predicted & actual output is evaluated as it propagates back through the w/b to adjust weight & bias.
→ It continues layer by layer ensuring that the w/b learns & improves performance

when to stop learning?

- ① Define threshold
- ② Set no. of epochs

- ③ Construct a perceptron which acts as an AND gate for the given straight line $x_1 + x_2 - 1.5 = 0$.

Ans

$x_1 \quad x_2 \quad y$

0 0 0

0 1 0

1 0 0

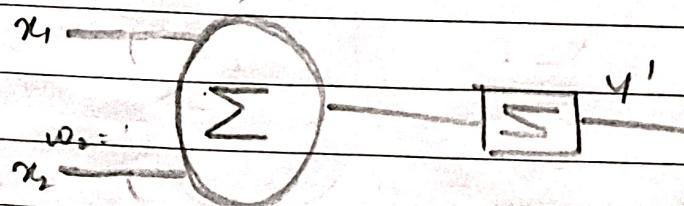
1 1 1

$$x = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \end{bmatrix} \quad w = [1, 1]$$

$$b = -1.5$$

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

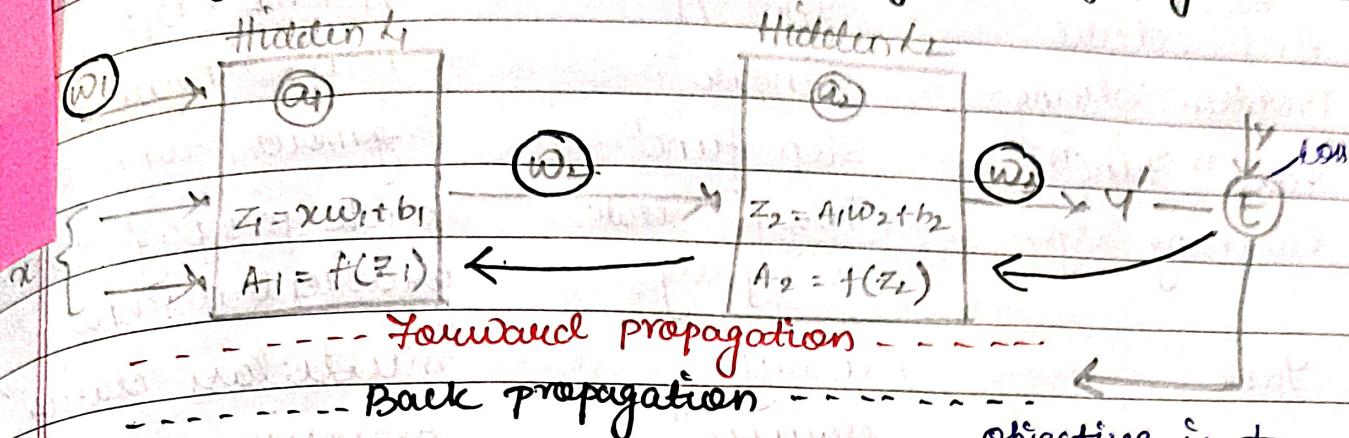
$$w_1 = 1 \quad \downarrow b = -1.5$$



$$\text{Af} = \text{Step func}^n = f(z) = \begin{cases} 1, & z > 0 \\ 0, & z \leq 1 \end{cases}$$

$$A_1 = f(z) = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

Working principle of backpropagation of algo.



$\frac{\partial E}{\partial w}$ → gradient funcⁿ.

objective is to
reduce error & select
optimal weights.
 E = actual o/p - predicted o/p.

$$\frac{\partial A_1}{\partial w_{11}} \times \frac{\partial a_1}{\partial A_1} \times \frac{\partial z_1}{\partial a_1} \times \frac{\partial y'}{\partial z_1} \times \frac{\partial E}{\partial y'} = \frac{\partial E}{\partial w_{11}}$$

- * A set of i/p neurons with input x connected to neurons of next layer having certain weights are multiplied & passed ^{to the activⁿ} function giving a certain o/p.
- * It is calculated keeping in mind that the actual o/p is calculated using derivative of cost funcⁿ.
- * Basically back propagⁿ is to update the weight to for reducing loss in the next iterⁿ.
- * We need to have derivative of cost funcⁿ i.e. partial derivative of cost funcⁿ wrt weight follows a chain rule i.e.
- * Aggregating all derivatives the change of 'w' is evaluated.

Criterias

Architecture

Problem Solving

Active" Func"

Learning algo.

Task

Types of app"

N/W complexity

SLP

I/P + O/P

Linear

Step func"

Percep" law
- ming algo.

Binary

classific"

Simple

simple

MLP

I/P + H/L + O/P

Non - linear

sigmoid, tanh, Relu

Back propag" +
gradient descentmulticlass classific",
regression

Complex

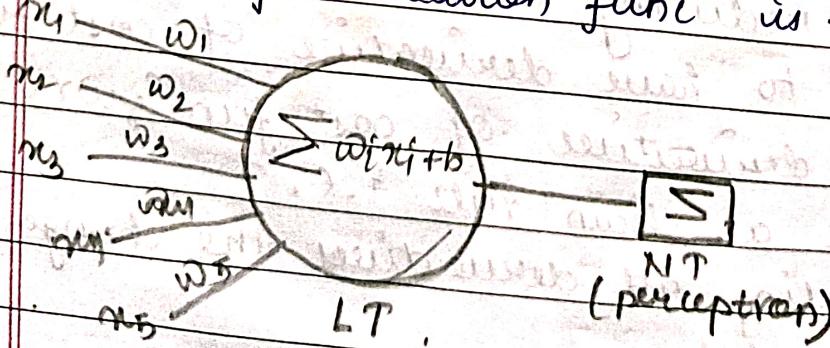
Complex

Perception

→ When neurons are connected through a network it is called perceptron.

→ They have diff n/w parameters - weight, bias,
active" func".

→ Output of activation func" is the predicted o/p.



Types of perceptron

- ① Single layer perceptron (No hidden layer)
- ② Multilayer perceptron (1 or more hidden layers)

Perceptron Training Rule

- It is an algo. used to train perceptron.
- It adjust the weights assigned to the i/p features based on errors made in predict.
- This allow the perceptron to improve its ability to distinguish b/w two classes.

- ① Initialization: Start initializing (w_0) & (b) to a very small random value.
- ② Set learning rate: Choose a learning rate that controls the magnitude of weight updates. It is set to very small value like 0.01, 0.001...

Training process :

- i) present each training example at a time.
- ii) For each i/p vector x calculate the weighted sum $z = x\bar{w} + b$ — ①
- iii) Apply AF, "typically a step func" to determine the predicted o/p.

$$Y_{pred} = \begin{cases} 0, & \text{if } z < 0 \\ 1, & \text{if } z \geq 0 \end{cases} — ②$$

- iv) Calculate error by comparing predicted o/p with actual o/p.

$$\text{Error} = y - Y_{pred} — ③$$

v) If there is error, then update weight

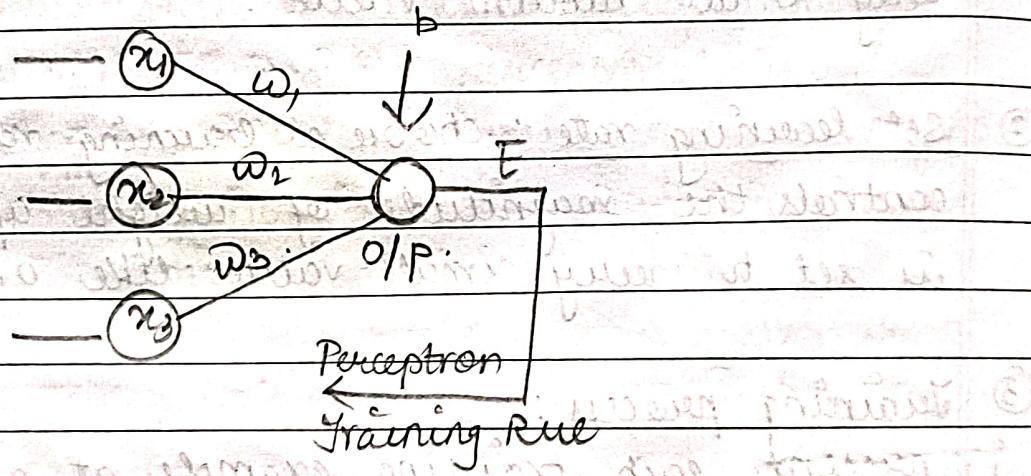
i.e. $y_{pred} \neq y$

$$w \leftarrow w + \eta * \text{Error} * x \quad \text{--- (4)}$$

$$b \leftarrow b + \eta * \text{Error} \quad \text{--- (5)}$$

This adjustment moves the decision boundary in a "direction", that reduces the classification error for the current sample.

④ Repeat the training process EPOCHS until (i) the perceptron correctly classifies all training example or (ii) until a max. EPOCHS is reached.



Q) $x = \begin{bmatrix} 2 & 1 \\ 1 & -1 \\ -1 & -2 \end{bmatrix}, t = \begin{bmatrix} +1 \\ -1 \\ -1 \end{bmatrix}, w = [0, 0], b = 0, \eta = 1$

$$y = \text{sign}(z) = \begin{cases} 1, & z > 0 \\ 0, & z = 0 \\ -1, & z < 0 \end{cases}$$

Ans. $z = x^T w + b = \begin{bmatrix} 2 & 1 \\ 1 & -1 \\ -1 & -2 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} + [0] = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$

consider 1st sample

$$x = [2, 1], w = [0, 0], b = [0], \eta = 1, t = +1$$

$$z = x w^T + b = 2 \times 0 + 0 \times 1 + 0 = 0$$

$$y = \text{sign}(z) = 0$$

$$\text{Error} = 1 - 0 = 1$$

$$w_1' = 0 + 1 \cdot 1 \times 2 = 2; w_{1,1} = 0 + 1 \cdot 1 \times 1 = 1$$

$$w_1' = [2 \ 1], b' = b + \eta * E = 1$$

consider 2nd sample.

$$x = [1 \ -1], w_1' = [2 \ 1], b_1 = 1, t = -1$$

$$z = x w^T + b = [1 \ -1] \begin{bmatrix} 2 \\ 1 \end{bmatrix} + [1]$$

$$= [1] + 1 = 2$$

$$\text{Error} = 1 - 2 = -1$$

$$e = t - y = -1 - 1 = -2$$

$$w_2' = 2 + 1 \times -2 \times 1 = [0 \ 3]$$

$$= 1 + 1 \times -2 \times -1$$

$$b = 1 + 1 \times (-2) = -1$$

consider 3rd sample

$$x = [-1, -2], w = [0, 3], b = 1, t = -1, \eta = 1$$

$$z = -7; y = -1$$

$$E = -1 + 1 = 0 \ (\text{no update})$$

$$w = [0, 3]; b = -1$$

Now verify for all samples

Perception learning for AND/OR gate

(Q) Apply the perceptron training rule to find optimal weights. Given:

$$w_1 = 1.2, w_2 = 0.6, \eta = 0.5, \text{threshold} = 1.$$

AND

x_1	x_2	t	$z = w \cdot x$	
0	0	0	0	b - no update
0	1	0	-1	$E = t - y \rightarrow 0 - \text{update wt.}$
1	0	0	1.2	$w \leftarrow w + \eta E x$
1	1	1	2.2	

Ans.

For sample 1:

$$x = [0, 0]; w = [1.2, 0.6]; \eta = 0.5$$

$$t = 1; z = 0.$$

Step 1:

$$z = w \cdot x = 0.$$

$$y = z < 1 = 0; E = 0 - 0 = 0$$

No weight update.

For sample 2:

$$x = [0, 1]; w = [1.2, 0.6]$$

$$z = [0, 1] \times \begin{bmatrix} 1.2 \\ 0.6 \end{bmatrix} = 0.6$$

$$y = z < 1 = 0; E = 0 - 0 = 0.$$

No weight update.

For sample 3:

$$x = [1, 0] \times \begin{bmatrix} 1.2 \\ 0.6 \end{bmatrix} = 1.2$$

$$y = 1; E = 0 - 1 = -1$$

$$w_1' = 1.2 + (0.5 \times (-1)) \times 1 = 0.7$$

$$= 0.6$$

step 2:

for sample 1:

$$\underline{x} = [0, 0] ; \underline{w} = [0.7, 0.6]$$

$$z = \underline{w} \cdot \underline{x} = [0.7 \ 0.6] \begin{bmatrix} 0 \\ 0 \end{bmatrix} = 0$$

$$E = 0$$

No update.

for sample 2:

$$\underline{x} = [0, 1]$$

$$z = [0 \ 1] \begin{bmatrix} 0.7 \\ 0.6 \end{bmatrix} = 0.6$$

$$\Rightarrow y = 0 ; E = 0$$

No update.

for sample 3:

$$\underline{x} = [1, 0] ; \underline{w} = [0.7 \ 0.6]$$

$$z = (\underline{x} \cdot \underline{w}) = [1 \ 0] \begin{bmatrix} 0.7 \\ 0.6 \end{bmatrix} = 0.7$$

$$y = z < 1 = 0 ; E = 0$$

No update.

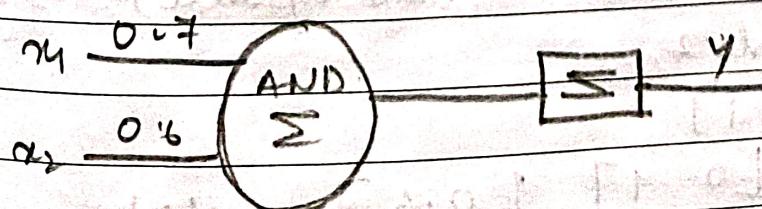
for sample 4:

$$\underline{x} = [1, 1]$$

$$z = [1 \ 1] \begin{bmatrix} 0.7 \\ 0.6 \end{bmatrix} = 0.7 + 0.6 = 1.3$$

$$y = z > 1 = 1 ; E = 0$$

No update.



(Q) Apply perceptron training rule over OR gate to find optimal weights.

$$w_1 = 0.6, w_2 = 0.6, \text{ thresh} = 1; \eta = 0.5$$

OR

x_1	x_2	t
0	0	0
0	1	1
1	0	1
1	1	1

For sample 1:

$$x = [0, 0]; w = [0.6, 0.6]$$

$$z = xw = [0 \ 0] \begin{bmatrix} 0.6 \\ 0.6 \end{bmatrix} = 0$$

$$y = z < 0 = 0. \text{ (No update)}$$

For sample 2:

$$x = [0, 1]$$

$$z = [0 \ 1] \begin{bmatrix} 0.6 \\ 0.6 \end{bmatrix} = 0.6$$

$$y = z < 0 = 0.; E = 1 - 0 = 1.$$

$$w = 0.6 + 0 = 0.6$$

$$0.6 + 0.5 \times 1 \times 1 = 1.1$$

$$\text{New weight } = [0.6, 1.1]$$

Step 2: $x = [0, 0]$

$$z = 0; E = 0 \text{ (No update)}$$

For sample 2

$$x = [0, 1]$$

$$z = [0 \ 1] \begin{bmatrix} 0.6 \\ 1.1 \end{bmatrix} = 1.1$$

$$y = 1; E = 0 \text{ (No update)}$$

for sample 3:

$$x = \underline{\underline{[1 \ 0]}}$$

$$z = [1 \ 0] \begin{bmatrix} 0.6 \\ 1.1 \end{bmatrix} = 0.6$$

$$y = z < 1 = 0; E = 1 - 0 = 1.$$

$$w = 0.6 + 1 \times 0.5 \times 1 = 1.1$$

$$= 1.1 + 0 = 1.1$$

new weight $\begin{bmatrix} 1.1 & 1.1 \end{bmatrix}$

step 3: $x = \underline{\underline{[0 \ 0]}}$

$z = 0, y = 0, E = 0$ (No update)

for sample 2

$$x = \underline{\underline{[0 \ 1]}}$$

$$z = 1.1; y = 1; E = 0$$

for sample 3

$$x = \underline{\underline{[1 \ 0]}}$$

$$z = 1.1; y = 1; E = 0$$

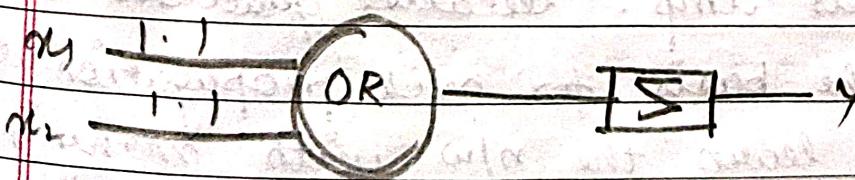
for sample 4

$$x = \underline{\underline{[1 \ 1]}}$$

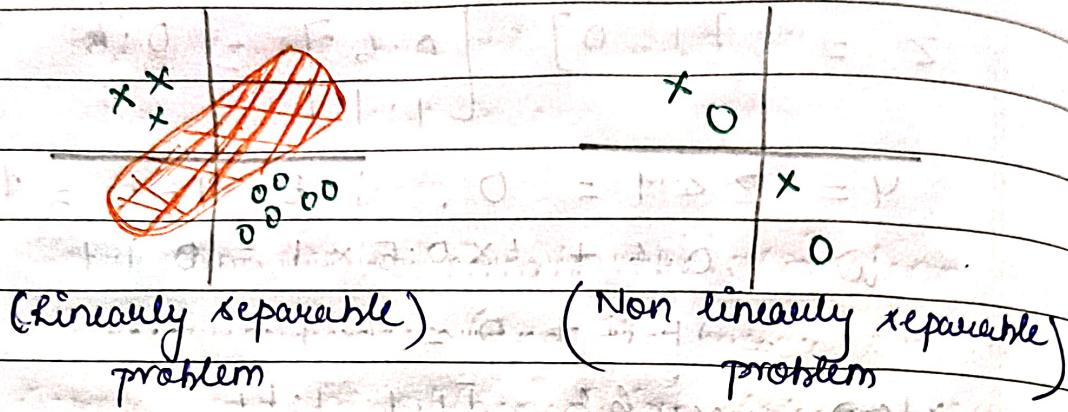
$$z = 1.1 + 1.1 = 2.2$$

$$y = z > 1 = 1$$

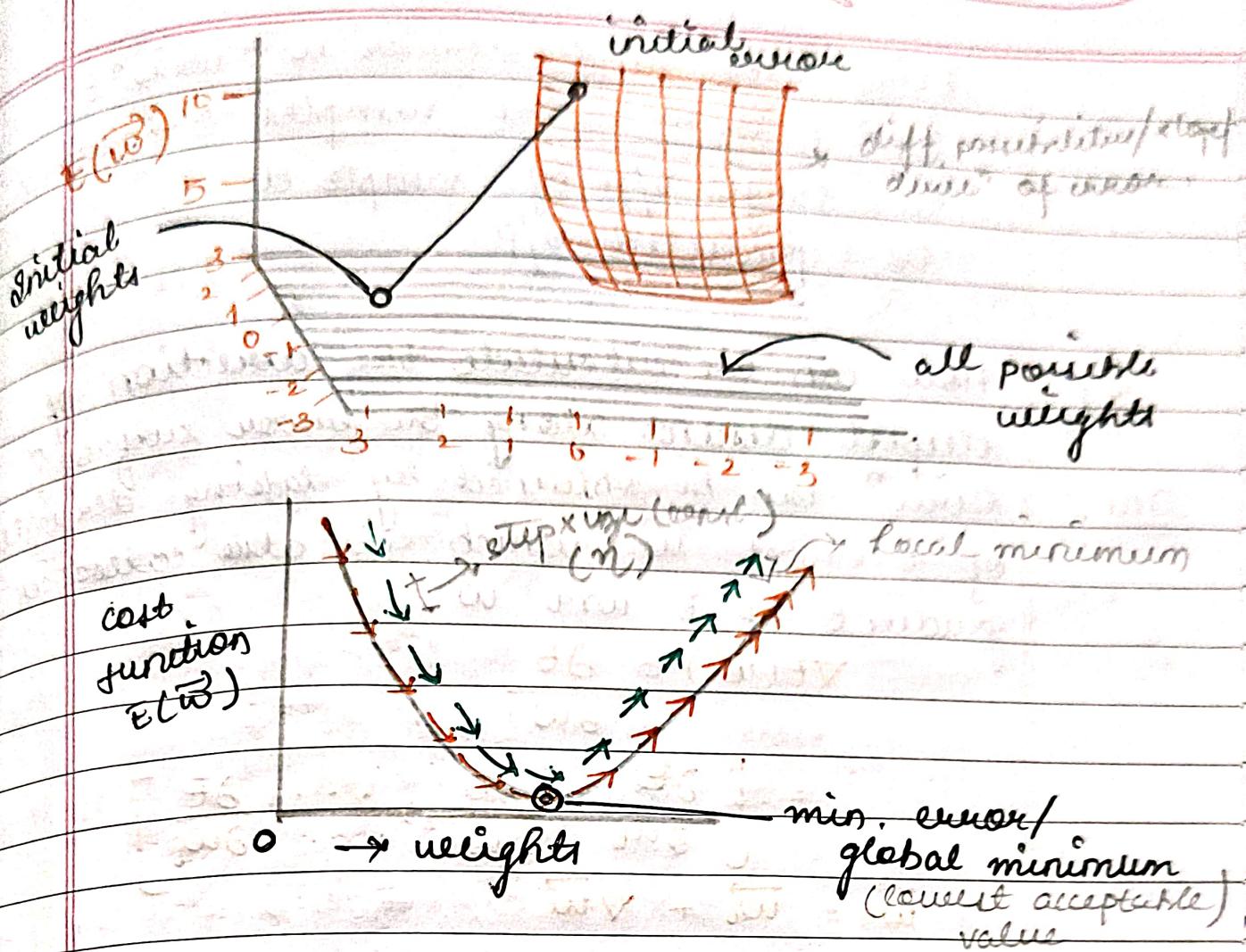
$E = 0$ (No update).



gradient descent & delta rule



- * Perceptron training rule finds a successful weight vector when training examples are linearly separable, but it fails to converge if examples are non linearly separable.
- * A 2nd rule called as delta rule is designed to overcome this issue.
- * If the training examples are not linearly separable, delta rule converges towards a least fit approx to a target o/p.
- * The key idea behind delta rule is to use gradient descent to search the hypothetical space of all possible weight vector, to find the weights that best fit the training example.
- * This rule is imp. because gradient descent provides the basis for back propagation algo, which can learn the n/w with many inter-connected units.



How delta rule works?

- * The delta training rule is best understood for considering the task of training an unthresholded perceptron which is a linear unit.

Here, O/p 'O' is defined as :

$$O = \omega_0 x_0 + \omega_1 x_1 + \omega_2 x_2 + \dots$$

$$\Rightarrow O(\vec{w}) = \vec{x} \cdot \vec{w}$$

- * In order to derive a wt. learning rule for linear learning unit, let us begin by specifying a measure for training error.

- * There are many ways to identify this error, one such way is considering MSE.

$$E(\vec{w}) = \frac{1}{2d} \sum_{d=1}^D (t_d - O_d)^2$$

$E(\vec{w})$ = cost func' = error wrt weight.

$d \in D = d \Rightarrow$ no. of samples $D \rightarrow$ dataset

t_d = target o/p of sample d

O_d = predicted o/p.

(Q) How can we calculate the direction of steepest descent along the error surface?

Ans. "Diric" can be obtained by taking derivative of 'E' wrt \vec{w} , which is also called a gradient of E wrt \vec{w} .

$$\nabla E(\vec{w}) = \frac{\partial E}{\partial \vec{w}}$$

$$= \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_i} \right]$$

$$\vec{w}_n = \vec{w}_0 + \nabla \vec{w}$$

$$\nabla \vec{w} = -\eta \frac{\partial E}{\partial \vec{w}} \quad \eta \rightarrow \text{learning rate}$$

a +ve constraint which determines step size in gradient descent search.

* The -ve sign is present because we want to move the \vec{w} in diric' that decreases E.

$$\frac{\partial E}{\partial w_i} = \frac{1}{2} \frac{\partial}{\partial w_i} \sum_{d \in D} (t_d - O_d)^2$$

$$= \frac{1}{2} \sum_{d \in D} (t_d - O_d)^2$$

$$= \frac{1}{2} \times \sum (t_d - O_d) \frac{\partial}{\partial w_i} (t_d - O_d)$$

$$= (t_d - O_d) \frac{\partial}{\partial w_i} (t_d) - \frac{\partial}{\partial w_i} (O_d)$$

Step-1:

Step-2:

Step-3:

Step-4:

(Q).

Ans.

$$= (t_d - o_d) - \frac{\partial}{\partial w_i} (o_d) \rightarrow \frac{\partial}{\partial w_i} (w_n + w_0 x_0 + w_i x_i)$$

$$= -(t_d - o_d) x_i$$

$$\therefore \nabla \vec{w} = \eta [-(t_d - o_d) x_i]$$

$$\nabla \vec{w} = \eta (t_d - o_d) x_i$$

$$w = w + \nabla \vec{w}$$

$$= w + \eta (t_d - o_d) x_i$$

Delta rule over linear perceptron

Step-1: $y = w_n$ (Apply forward pass).

Step-2: Find out error

$$E = t - y, t = \text{actual class.}$$

Step-3: If $E \neq 0$, then apply delta rule

$$\nabla \vec{w} = \eta (t - y) x_i$$

Step-4: Update weights.

$$w_n = w_0 + \nabla \vec{w}$$

Continue till all samples have $E = 0$.

Q). $x = 0.5, t = 0.8, w = 0.4, \eta = 0.1$. Apply delta rule & find out E for 2 iterations. update weights. & show the trend.

Ans. $y = w_n = 0.5 \times 0.4 = 0.2$

$$E = t - y = 0.8 - 0.2 = 0.6$$

$$\nabla \vec{w} = \eta (t - y) x_i = 0.1 \times (0.8 - 0.2) \times 0.5 = 0.1 \times 0.6 \times 0.5$$

$$= 0.03$$

$$w_n = 0.4 + 0.03 = 0.43$$

Step-2: $y = 0.43 \times 0.5 = 0.215$

$$E = 0.8 - 0.215 = 0.585$$

$$\nabla \vec{w} = 0.1 \times 0.585 \times 0.5 = 0.029$$

$$w = 0.43 + 0.029 = 0.459$$

<u>Round</u>	<u>Error</u>	<u>Weight</u>
1.	0.6	0.43
2.	0.585	0.459

Q). $x_1 = 0.6$; $x_2 = 0.8$; $w_1 = 0.2$; $w_2 = 0.4$
 $t = 0.75$; $b = 0.1$; $\eta = 0.1$. apply delta rule
over the given linear perceptron, find
weight & bias update for 2 iterations

Ans.

$$y = xw + b$$

$$= [0.6 \quad 0.8] \begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix} + 0.1$$

$$= 0.24 + 0.1 = 0.34$$

$$E = t - y = 0.41$$

$$\nabla w = 0.2 \times 0.6 \quad 0.41 \times 0.6 = 0.024$$

$$0.1 \times 0.41 \times 0.3 = 0.012$$

$$w_n = [0.224 \quad 0.412] \quad \nabla b = \eta t = 0.041; b = 0.41$$

iter²: $y = [0.6 \quad 0.8] \begin{bmatrix} 0.224 \\ 0.412 \end{bmatrix} + 0.1$

$$= 0.258 + 0.141 = 0.399$$

$$E = 0.75 - 0.399 = 0.351$$

$$\nabla w = 0.1 \times 0.351 \times 0.6 = 0.021$$

$$0.1 \times 0.351 \times 0.3 = 0.010$$

$$w_n = [0.245 \quad 0.422]$$

$$\nabla b = \eta E = 0.1 \times 0.1 = 0.0351$$

$$b_n = 0.41 + 0.0351 = 0.4451$$

If non-linear definitely Hidden layer

classmate

Date _____
Page _____

Derivation of learning rule for a perceptron
with sigmoid activation function

Step 1:

Step 2:

Initialize weights & bias.

apply forward pass & calculate 'z'.

$$z = \sum_{i=1}^n x_i w_i \text{ OR } \sum_{i=1}^n x_i w_i + b$$

Step 3:

Apply activation func sigmoid

$$y' = \frac{1}{1 + e^{-z}}$$

Step 4: If $y' \neq y$, apply delta rule & update wt.

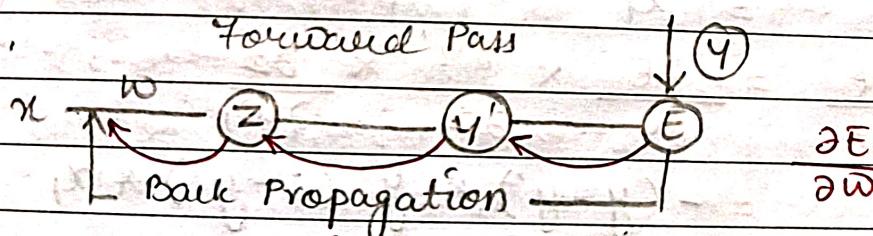
$$w_i = w_i - \eta \frac{\partial E}{\partial w_i}$$

y' → predicted o/p ; y → actual o/p

η → learning rate ; $\frac{\partial E}{\partial w_i}$ → gradient of error wrt weight.

Step 5: Repeat step 2 to step 4 until loss func converges or stopping criteria reached.

$$E(w) = \frac{1}{2} (y - y')^2$$



$$\frac{\partial E}{\partial w} = \frac{\partial E}{\partial y'} \times \frac{\partial y'}{\partial z} \times \frac{\partial z}{\partial w} = \frac{\partial E}{\partial y'}$$

$$\frac{\partial E}{\partial w} = \frac{\partial E}{\partial y'} \times \frac{\partial y'}{\partial z} \times \frac{\partial z}{\partial w}$$

$$\frac{\partial E}{\partial y'} = \frac{\partial}{\partial y'} \left[\frac{1}{2} (y - y')^2 \right]$$

$$= \frac{1}{2} \times 2 (y - y') \frac{\partial}{\partial y'} (y - y')$$

$$\begin{aligned}
 &= (y - y') x - 1 \\
 &= -(y - y') \\
 \frac{\partial y'}{\partial z} &= \frac{\partial}{\partial z} (1 + e^{-z})^{-1} \\
 &= -(1 + e^{-z})^{-2} x - \frac{e^{-z}}{(1 + e^{-z})^2} \\
 &= \frac{e^{-z}}{(1 + e^{-z})^2} \\
 &= \frac{e^{-z}}{(1 + e^{-z})^2}
 \end{aligned}$$

$$\begin{aligned}
 y(1-y) &= & y' &= (1 + e^{-z})^{-1} \\
 \frac{1}{1+e^{-z}} \times e^{-z} &= & 1 - y' &= 1 - \frac{1}{1+e^{-z}} \\
 = \frac{e^{-z}}{(1+e^{-z})^2} &= & = \frac{e^{-z}}{1+e^{-z}} \\
 \therefore \frac{\partial y'}{\partial z} &= y(1-y')
 \end{aligned}$$

$$\frac{\partial z}{\partial w_i} = \frac{\partial}{\partial w_i} [w_0 x_0 + w_1 x_1 + w_2 x_2 + \dots] = x_i$$

$$\begin{aligned}
 \frac{\partial E}{\partial y} &= \frac{\partial E}{\partial y} \times \frac{\partial y'}{\partial z} \times \frac{\partial z}{\partial w_i} \\
 &= -(y - y') * y' (1 - y') * x_i \\
 w_i &= w_i - \eta \frac{\partial E}{\partial w_i} \\
 &= w_i - \eta [- (y - y') y' (1 - y') x_i] \\
 &= w_i + \eta (y - y') y' (1 - y') x_i
 \end{aligned}$$

$$\left. \begin{array}{l} \text{Linear} \\ w = w + \eta (t_{di} - o_{di}) x_i \end{array} \right\} \quad \left. \begin{array}{l} \text{Non-linear} \\ w = w + \eta (y - y') y' (1 - y') x_i \end{array} \right\}$$

X

4P
Hidden prob

linearly separable non linearly separable.

Q) $x_1 = 0.6, w_1 = 0.8, b = 0.2, t = 1, \eta = 0.1$
 AP = sigmoid. Apply delta rule & update weight for 2 iterations.

$$z = xw + b$$

$$\text{Q.M.} \quad = (0.6 \times 0.8) + 0.2 = 0.68$$

$$y' = \text{sigmoid}(z) = \frac{1}{1+e^{-z}} = 0.663$$

$$E = t - y' = 1 - 0.663 = 0.337$$

$$\nabla w = \eta (y - y') y' (1 - y') x \\ = 0.068$$

$$\nabla b = \eta (y - y') y' (1 - y') = 0.007$$

$$w_n = w_0 + \nabla w = 0.808 + 0.068 = 0.868$$

$$b_n = b_0 + \nabla b = 0.207 + 0.007 = 0.207$$

$$\text{Step 2: } z = 0.6 \times 0.868 + 0.207 = 0.6902$$

$$y' = \frac{1}{1+e^{-z}} = \frac{1}{1+e^{-0.6902}} = 0.6659$$

$$E = 1 - 0.6659 = 0.3341$$

$$\nabla w = 0.0044$$

$$\nabla b = 0.0074$$

$$w_n = 0.808 + 0.0044 = 0.8089$$

$$b_n = 0.207 + 0.0074 = 0.2149$$

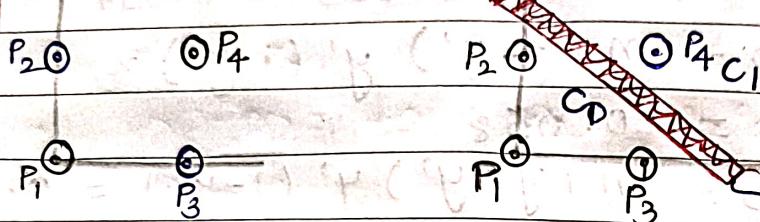
Perception training for XOR gate

AND

XOR

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

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



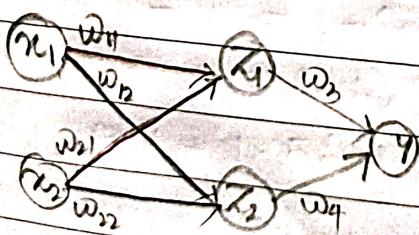
- * If we want to draw a straight line wrt XOR gate in 2D plane, to differentiate the data objects, we cannot draw.
- * This is a different type of logical gate which is nonlinearly separable.
- * So, we cannot apply simple perceptron rule to find or implement XOR gate.
- * Hence, we will use a different eqⁿ which represent a non-linear solⁿ to XOR gate.

$$y = x_1 \bar{x}_2 + \bar{x}_1 x_2$$

$$z_1 = x_1 \bar{x}_2 - \text{1st function}$$

$$z_2 = x_2 \bar{x}_1 - \text{2nd function}$$

$$y = z_1 \text{ OR } z_2$$



$$z_{11} = x_1 w$$

$$y' = f(z_{11})$$

$$\epsilon = y - y'$$

If $\epsilon \neq 0$, update w

$$w = w + \Delta \nabla w = w + \eta \epsilon x$$

$$\omega_1 = [\omega_{11}, \omega_{21}] = [1, 1]$$

$$\omega_2 = [\omega_{12}, \omega_{22}] = [1, 1] \quad \eta = 1.5$$

1st function

$$\begin{array}{cccc} x_1 & x_2 & z_1 & f(y') \\ \hline 0 & 0 & 0 & 1, y' \geq 0 \\ 0 & 1 & 0 & 0, y' < 0 \\ 1 & 0 & 1 & \\ 1 & 1 & 0 & \end{array}$$

$z_1 = x_1 \bar{x}_2$

Step 1: $S^1: x_1 = [0, 0], y = 0$

$$z_{11} = [0 \ 0] \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 0$$

$$f(z_{11}) = 0, \text{ no update}$$

$$S^2: x_2 = [0, 1], \omega = [1, 1], \eta = 1.5, y = 0$$

$$z_{11} = [0 \ 1] \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1.$$

$$E = 0 - 1 = -1; \omega = 1 + 1.5 \times (-1) \times 0 = 1$$

Step 2: $S^1: x = [0, 0], y = 0$

$$z = [0 \ 0] \begin{bmatrix} 1 \\ -0.5 \end{bmatrix} = 0$$

$$E = 0, \text{ no update.}$$

$$S^2: x = [0, 1]$$

$$z = [0 \ 1] \begin{bmatrix} 1 \\ -0.5 \end{bmatrix} = -0.5$$

$$f(z) = 0; E = 0, \text{ no update}$$

$$S^3: x = [1, 0]$$

$$z_{11} = [1 \ 0] \begin{bmatrix} 1 \\ -0.5 \end{bmatrix} = 1.$$

$$E = 0, \text{ no update.}$$

$$S^4: \quad x = [1 \ 1] \\ z = [1 \ 1] \begin{bmatrix} 1 \\ -0.5 \end{bmatrix} = 0.5 = 0$$

$$E = 0 - 0 = 0, \\ w = 1 + 1.5(-1)1 \neq -0.5 \\ -0.5 + 1.5(-1)1$$

function \bar{w}

$$z = \bar{w}_1 x_1$$

$$x_1 \quad x_2 \quad z_2$$

$$\begin{bmatrix} 0 & 0 & 0 \end{bmatrix} = \bar{w}_2 = [\bar{w}_{12}, \bar{w}_{22}]$$

$$\begin{bmatrix} 0 & 1 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 & 1 & 1 \end{bmatrix}$$

$$\underline{\text{step 1:}} \quad S^1: x = [0 \ 0]$$

$$z = [0 \ 0] \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 0$$

$E = 0$, no update

$$S^2: x = [0, 1]$$

$$z = [0 \ 1] \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1$$

$E = 0$, no update

$$S^3: x = [1, 0]$$

$$z = [1 \ 0] \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1$$

$$E = 0 - 1 = -1$$

$$w = 1 + 1.5(-1)1 = -0.5$$

$$\bar{w} = 1 + 1.5(-1)0 = 1$$

$s^1: s' = [0, 0]$

$z = [0 \ 0]$

$$\begin{bmatrix} -0.5 \\ 1 \end{bmatrix} = 0$$

$E = 0 - 0 = 0$, no update.

$s^2: x = [0 \ 1]$

$z = [0 \ 1]$

$$\begin{bmatrix} -0.5 \\ 1 \end{bmatrix} = 1$$

$E = 1 - 1 = 0$, no update

$s^3: x = [1 \ 0]$

$z = [1 \ 0]$

$$\begin{bmatrix} -0.5 \\ 1 \end{bmatrix} = -0.5$$

$f(z) = 0$.

$E = 0 - 0 = 0$, no update.

$s^4: x = [1 \ 1]$

$z = [1 \ 1]$

$$\begin{bmatrix} -0.5 \\ 1 \end{bmatrix} = 0.5$$

$f(z) = 0$.

$E = 0 - 0 = 0$, no update.

3rd Function

$$\begin{array}{c} x_1 \\ \hline 0 & 0 & 0 \end{array}$$

$$\begin{array}{c} 0 & 1 & 1 \end{array}$$

$$\begin{array}{c} 1 & 0 & 1 \end{array}$$

$$\begin{array}{c} 0 & 0 & 0 \end{array}$$

$s^1: s' = [0 \ 0] \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 0$

$E = 0$, no update.

$$S^2: w = [0 \ 1]$$

$$z = [0 \ 1] \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1.$$

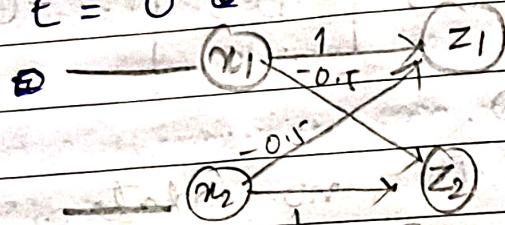
$$E = 1 - 1 = 0, \text{ no update.}$$

$$S^3: x = [1 \ 0] \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1.$$

$$E = 1 - 1 = 0, \text{ no update.}$$

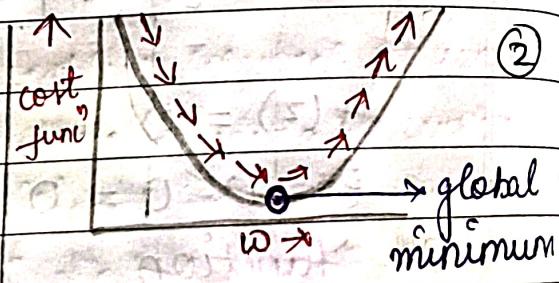
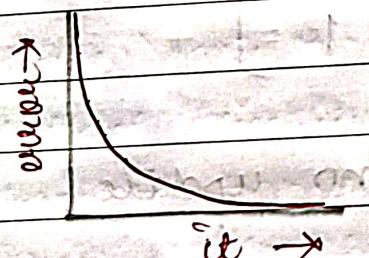
$$S^4: [0 \ 0] \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 0.$$

$$E = 0 - 0 = +0, \text{ no update}$$

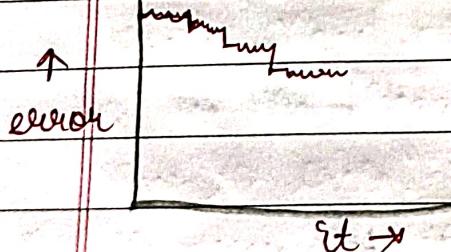


Problem with gradient descent

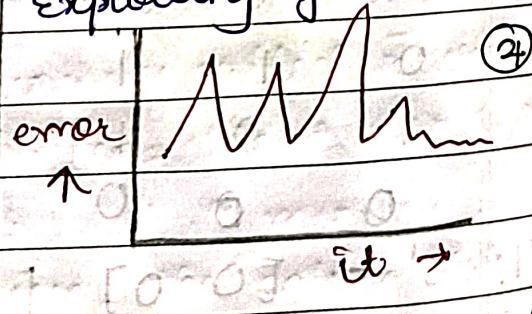
① Ideal curve



③ Vanishing gradient



Exploding gradient



* Graph b/w cost func & weight demonstrate (fig-2) how gradient descent achieves its aim of reaching global cost minima point.

- * Each step towards local minima pt is determined by gradient or slope ($\frac{\partial E}{\partial w}$) & step size i.e. learning rate (η).
- * Choosing inappropriate AF, learning leads to gradient problem.

① vanishing gradient: It is a scenario in the learning process of NN, where model doesn't learn at all.

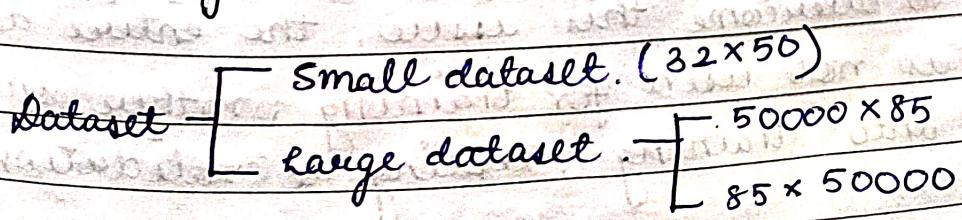
→ This is due to, when gradient becomes too small, almost vanishes leads weight getting stuck, never reaching the optimal value. This n/w is not able to learn & converge.

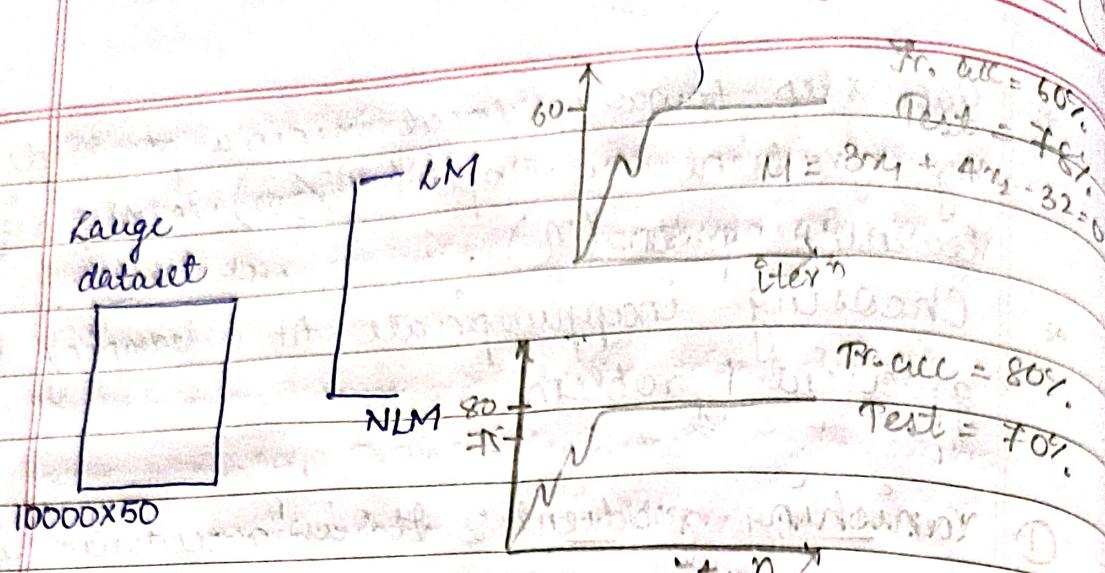
② Exploding gradient: Is exactly opposite to vanishing gradient, where model keeps on learning, weight keeps updating but model never converges.

→ It computes gradient wrt weight which becomes extremely large in the earlier layer of NN in such a way that it explodes or oscillates.

Overfitting

construct
its
point





- * Given with a bunch of data points on a flat plane, our goal is to find a curve, that best describes the data.
- * Let us consider two models to train these data points:

① Linear model ② Non linear model

- * By building a more complex model it's quite easy to perfectly fit our training dataset because we are giving enough freedom to our model to control itself to fit the "observed" data.
- * But when we evaluate such complex model over new dataset, model fails. This phenomena is called overfitting & is the biggest challenge that ML & DL model fails.

* To overcome this issue, the entire dataset is not used for training, rather split into training, testing & validation set.

Data, split
25 samples

Training Set
(20 samples)

Test set
(5 samples)

Training set

Validation set

16

4



Case 1 (✓)

Case 2 (Overfitting)

tr v

tr v

B1 98 96

B1 98 96

B2- 96 95

B3 - 98 95

23 23

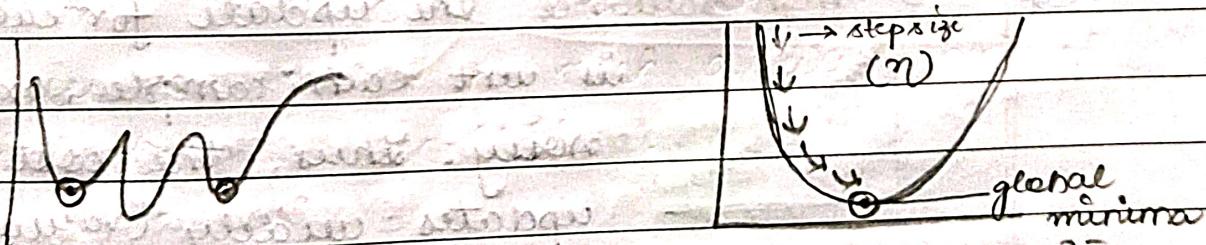
Page 86 of 88

卷之三

B1 98 95

Challenges with gradient descent

* Local minima & Saddle Point:



- The point where gradient is 0 & cost is lower than its immediate neighbour, higher than global minima.
 - Point where gradient is 0, but funcⁿ increases in some direcⁿ & decreases in some other. Thus, algo slow down significantly around these points.

- ② Choosing learning rate.
- * Too large - overstepped: The algo may overshoot the minima, fails to converge or may diverge.
 - * Too less - slow convergent: The algo. will converge very slowly. Potentially will take long time to reach min. point.

③ Computational time

- * gradient descent (GD)
 - i) Batch gradient descent - entire data 
 - ii) Stochastic GD - single sample 
 - iii) Minibatch GD - a small batch sample. 

- * Batch GD is computationally very expensive.
- * storing the entire data in the memory & computing avg. gradient is infeasible.

Variance stochastic GD & Minibatch GD.

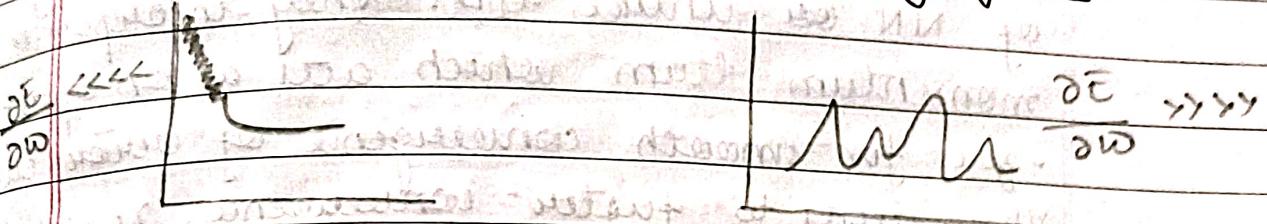
- (i) SGD - High Variance The update for each step is w.r.t each sample, maybe noisy. Thus, the cost func' updates widely (variance) & thus may not reach global minima.

- (ii) Selection of batch size in minibatch.
- * One extra parameter i.e. batch size is added for efficient learning.

Scaling issue - feature value - [0-1]

- * If the feature value are scaled within a large interval, GD may follow a zig-zag pattern for convergence.

Vanishing gradient & exploding gradient



- * When gradient becomes extremely small i.e. close to 0 in earlier layers, making the weights in those layer learn very slowly or stop learning. In this scenario, model may not reach ~~go~~ global min-point.

- * Gradient becomes extremely large, leading to unstable n/w behaviour. Thus, leads to large weight updates & prevents convergence.

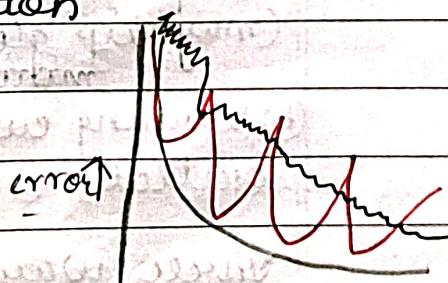
Solution to gradient descent

- * Momentum based optimization

$$v_{t+1} = \beta v_t + (1-\beta) \frac{\partial E}{\partial w_t}$$

$$w_{t+1} = w_t - \eta v_{t+1}$$

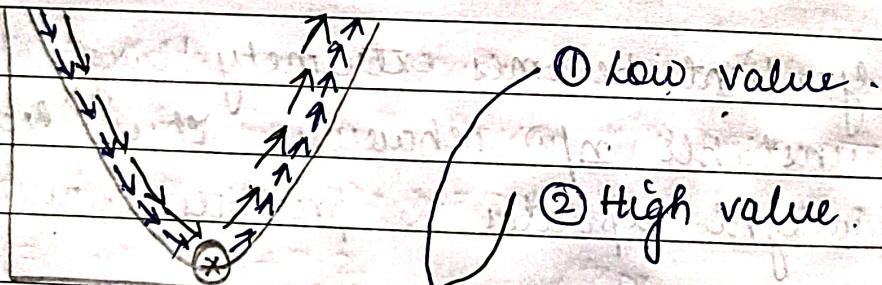
$$w_{t+1} = w_t - \eta \frac{\partial E}{\partial w}$$



- * Here, v_t is the velocity i.e. running avg. of GD.

- * $\beta = \text{momentum factor}$ ranges from 0 to 1.
- * $\eta = \text{learning rate}$
- $\frac{\partial E}{\partial W} = \text{gradient of loss func"}$
- * Momentum base GD optimizer are advanced techniques used to enhance the training of NN ~~as~~ unlike GD. They incorporate momentum term which acts as an optimizer for smooth convergence of error curve.
- * This leads to faster convergence, reduced oscillation & improved performance particularly while training a large dataset.

Impact of learning rate (η) over gradient



- | | |
|---|--|
| <ul style="list-style-type: none"> → ① Overshooting
(may skip global minima) ⑥ learning curve may oscillate | <ul style="list-style-type: none"> → ② convergence is slow ⑤ more epoch ④ computatⁿ cost is high |
|---|--|

Three ways to update weight

- ① Adagrad (Adaptive gradient algo.)
- ② RMSprop (Root mean square propagation)
- ③ ADAM (Adaptive moment optimizer).

① Adagrad

- * It adapts learning rate per parameter based on squared gradient.

$$\text{let}_t = \frac{n}{(g_t + \epsilon)^{1/2}} \quad n \rightarrow \text{const}$$

$g_t \rightarrow \text{sum of squared grad.}$
 $\epsilon \rightarrow \text{small const. to avoid divided by 0 error. } \omega = \omega - n \frac{\partial E}{\partial \omega}$

$$\theta_{t+1} = \theta_t - \text{let}_t \nabla \theta \quad \begin{matrix} \text{current} \\ \text{parameter} \end{matrix} \quad \begin{matrix} \nabla \theta \\ \rightarrow \text{gradient of cost func}^n \end{matrix}$$

→ Sparse data (α^2)

→ Features have diff. importance

→ change of η is too small. May decay [$0-10^{-5}$]

② RMSProp

- learning rate is updated considering running sum of average squared gradient

$$= \left(\left(\frac{\partial E}{\partial \omega_1} \right)^2 + \left(\frac{\partial E}{\partial \omega_2} \right)^2 \dots \right)^{1/2}$$

→ SGD & Adagrad:

④ gradient at step t

$$g_t = \nabla \theta \quad (\text{grad. of cost func}^n \text{ wrt } \theta)$$

⑤ update moving avg of gradient.

$$E[g^2]_t = \gamma E[g^2]_{t-1} + (-g_t)$$

$\gamma = \text{decay rate} = \text{considering } [0-1], \text{ mostly is } 0.9$

⑥ update parameter

$$\theta_{t+1} = \theta_t - \frac{\eta}{(E[g^2]_t + \epsilon)^{1/2}} \quad \begin{matrix} \eta = 0.0001 \\ \epsilon = 0.9 \end{matrix}$$

③ Adam

- * computationally efficient (in terms of error) &
- * high performance speed up convergence.

* RMSProp + momentum concept (technique).

@ first moment estimate

$$m_t = \beta_1 * m_{t-1} + (1 - \beta_1) g_t \quad [1^{\text{st}} \text{ moment vector}]$$

(b) 2nd moment estimate

$$v_t = \beta_2 * v_{t-1} + (1 - \beta_2) g_t^2 \quad [2^{\text{nd}} \text{ moment vector}]$$

$$\hat{m}_t = m_t / (1 - \beta_1^t)$$

$$\hat{v}_t = v_t / (1 - \beta_2^t)$$

@ weight update

$$w_t = w_{t-1} - \alpha * \hat{m}_t \quad m_t \text{ & } v_t = 1^{\text{st}} \text{ & } 2^{\text{nd}}$$

$$(\sqrt{v_t} + \epsilon) \quad \text{moment vector.}$$

$$\alpha = 0.0001; \beta_1 = 0.9; \beta_2 = 0.999; \epsilon = 10^{-8}$$

Q) Apply delta rule & gradient descent to update weight for 2nd iter. Given n/w parameters are $x=2$; $y=0.8$; $w=0.5$; $b=0.2$; $\eta=0.1$

Ans.

$$z = 2 \times 0.5 + 0.2 \quad \text{steps}$$

$$= 1.2 \quad ① z = xw + b$$

$$\text{sigmoid} = 0.768 \quad ② y' = \text{sigmoid}(z)$$

$$E = 0.0314$$

$$③ E = y - y'$$

$$\nabla w = 0.00114$$

④ if $E \neq 0$, update w_t &

$$\nabla b = 0.000057$$

bias AND restart

$$w = 0.50114$$

$$\nabla w = \eta (y - y') y' (1 - y') x$$

$$b = 0.200057$$

$$\nabla b = \eta (y - y') y' (1 - y')$$

$$⑤ w = w + \nabla w$$

$$b = b + \nabla b$$

Iterⁿ2: $x = 1.202$

$$y' = 0.7688$$

$$E = 0.0312$$

$$\nabla w = 0.00109$$

$$\nabla b = 0.0000554$$

$$w = 0.50223 ; b = 0.2001124$$

$$w=0.5, b=0.2$$

$$x=2 \quad y'$$

$$\left(\sum_{z=1}^z \right) y'$$

$f = \text{Sigmoid}$

Q) $[x_1, x_2] = [2, 3] ; [w_1, w_2] = [0.4, 0.6]$

$$y = 0.8 ; b = 0.2 ; \eta = 0.1$$

Iterⁿ1: $x = 2.8 ; y' = 0.942$

$$E = 0.8 - 0.942 = -0.142$$

$$\nabla w = [-0.00155, -0.00232].$$

$$\nabla b = -0.000775$$

$$w' = [0.39845, 0.59768]$$

$$b' = 0.19925.$$

Iterⁿ2: $x = 2.78919$

$$y' = 0.94208 ; E = -0.14208$$

$$\nabla w = [-0.00155, -0.00232]$$

$$\nabla b = -0.0007752$$

$$w'' = [0.3969, 0.59536]$$

$$b'' = 0.19847.$$