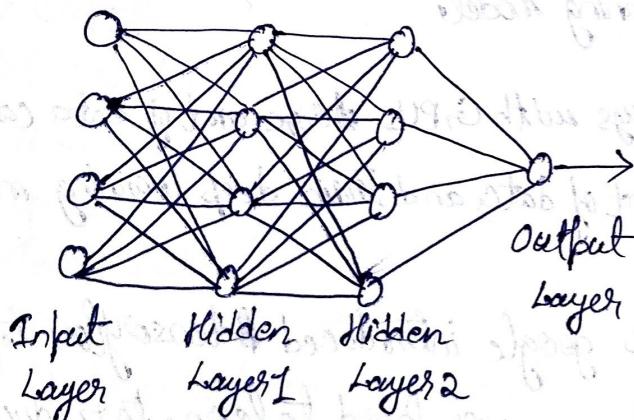
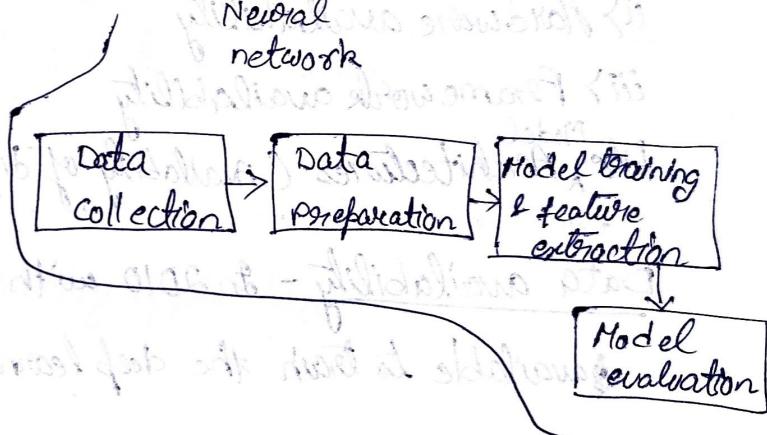
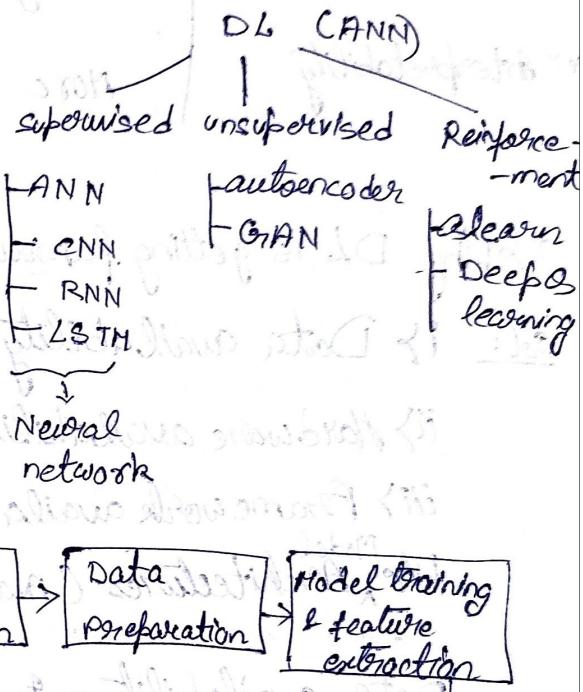
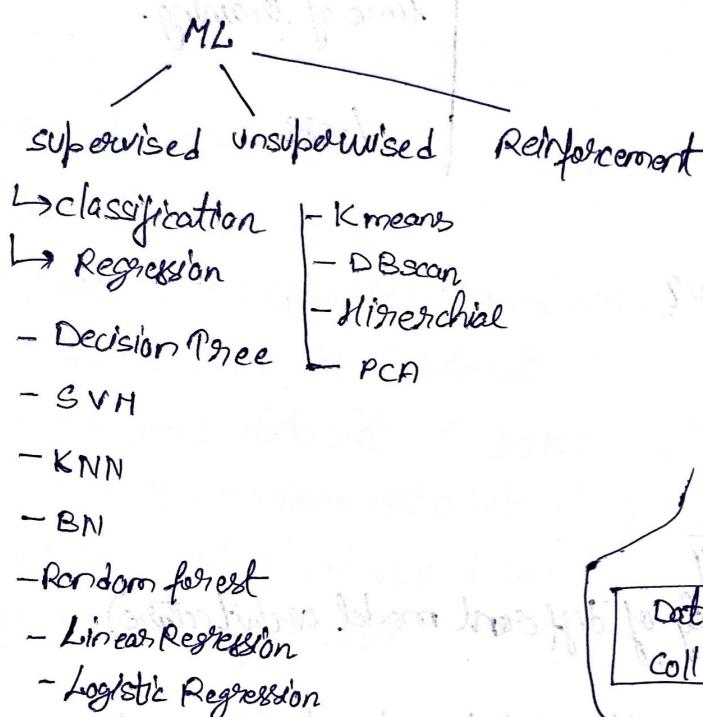


Deep Learning Using Python

6/9

Deep learning - is a subfield of ML that uses neural networks with multiple layers to learn complex patterns and make predictions from data.

- It enables computers to automatically learn representations and features from raw data like images, videos, audio or text often with minimal human ~~interaction~~ interaction.



- Each hidden layer extracts some feature of the input data.

ML vs DL

ANN → works with numerical data

sequential data

CNN → " " image data

" " Video

RNN → " " sequential data

LSTM → " "

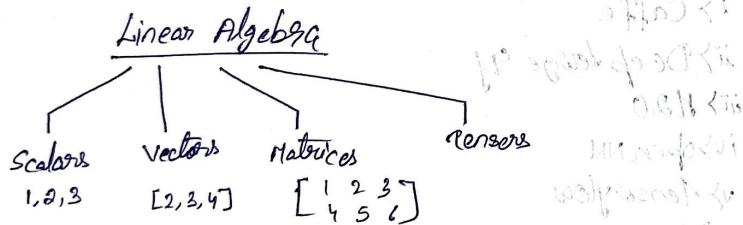
Application :-

- i) Autonomous car
 - ii) X-ray classification in health care (CNN)
 - iii) Advanced facial recognition (CNN)
 - iv) Visual question and answer (CNN ~~and RNN~~)
 - v) Advanced voice assistance tools and devices (RNN & LSTM)
 - vi) Automatic machine translation (RNN)
 - vii) Image recognition and image capturing (CNN) & (RNN)
 - viii) Anatomy detection & classification (ANN)
 - ix) Predicting earthquakes (Regression can be used & ANN) \rightarrow (DL or ML algo)
 - x) Language translation (RNN)
 - xi) Place restoration (CNN)
 - xii) Photo descriptions (CNN)

819

Linear Algebra

- Linear Algebra is a branch of mathematics that deals with the study of linear systems which are sets of equations involving linear functions of variables.



Q3) Why linear algebra for DL?

Sol: i) Good generalization ability.

ii) Good data representation.

Movie Recommendation System

```

graph LR
    movie[movie] --> matrix[matrix]
    matrix --> m1["m1"]
    matrix --> m2["m2"]
    matrix --> m3["m3"]
    matrix --> m4["m4"]
    matrix --> m5["m5"]
  
```

movie data

<u>movie name</u>	<u>Summary</u>
m1	S1
m2	S2
m3	S3
m4	S4

51 → The weather is good today

S2 → My name is Ram

$S_2 \rightarrow$ New or \neq yes

S₄ → I T E R is a good college

Unique words → The weather is ^{today} good _{my name}

Ram how are you I T E R a college

$$s_1 \rightarrow 1111100009990000$$

$s_2 \rightarrow 00100111000000$
 $\quad\quad\quad (is)$

- In hyperplane we can represent a point of n coordinates.

Dot product

$$a_1x_1 + a_2x_2 + \cdots + a_n x_n$$

$$a_1b_1 + a_2b_2 + \dots + a_nb_n$$

$$A = \begin{bmatrix} a_1 \\ a_n \end{bmatrix}, B = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$$

$$A \cdot B = A^T B$$

Fundamental spaces

column space - of a matrix is the set of all linear combinations of its columns.

$$A = \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix} \quad q_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad q_2 = \begin{bmatrix} 2 \\ 4 \end{bmatrix}$$

$$\overline{1} \alpha_2 = \alpha_1$$

Rank of a matrix is the dimension of its column space, the max no. of linearly independent columns or rows.

Rank of $A = 1$

$$B = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \text{ rank } = 3$$

$$C = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \\ 2 & 6 & 9 \end{bmatrix} \text{ rank } = 2$$

Use of column space in neural network - In a fully connected linear neural network equation $[y = w\bar{x} + b]$, where $w \rightarrow$ weight matrix, $\bar{x} \rightarrow$ input matrix

If w has low rank, its column space is small meaning the output lives in a low dimensional space limiting the networks capacity to represent complex patterns. Model expressiveness \propto full rank, the more the rank the more the expressive the model is.

Null space - The null space of a matrix A is the set of all vectors \bar{x} such that $A\bar{x} = 0$

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

$$A\bar{x} = \begin{bmatrix} 1 & 2 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = 0 \Rightarrow \begin{aligned} 1x_1 + 2x_2 - x_3 &= 0 \Rightarrow x_1 = -2x_2 + x_3 \\ 2x_1 + 4x_2 - 2x_3 &= 0 \quad x_2 = s, x_3 = t \end{aligned}$$

$$\bar{x} = \begin{bmatrix} -2s+t \\ s \\ t \end{bmatrix} = s \begin{bmatrix} -2 \\ 1 \\ 0 \end{bmatrix} + t \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

$$\text{Null space of } A = \left(\begin{bmatrix} -2 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \right)$$

$$\text{Find out null space of } A = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 4 & 3 & 2 & 1 \end{bmatrix}$$

$$\text{Sol: } A\bar{x} = 0$$

$$\Rightarrow \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 4 & 3 & 2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = 0 \Rightarrow \bar{x}$$

(To find this type of ans we need row echelon form).

Row Echelon form -

Condition - i) Leading element of each row must be unity or 1.

ii) If any 0 row is present then it must be at last.

iii) No. of zeros in each row must be more than previous row.

Condition for reduced row echelon form -

Below and above of leading element, they must be 0.

$$\text{Eg: } \begin{bmatrix} 1 & 0 & 1 & 2 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 2 & 3 & 4 \end{bmatrix}$$

$$A = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 4 & 3 & 2 & 1 \end{bmatrix} \xrightarrow{\begin{array}{l} R_2 = R_2 - R_1 \\ R_3 = R_3 - 4R_1 \end{array}} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 2 & 3 \\ 0 & -1 & -2 & -3 \end{bmatrix} \xrightarrow{\begin{array}{l} R_1 = R_1 - R_2 \\ R_3 = R_3 + 2R_2 \end{array}} \begin{bmatrix} 1 & 0 & -1 & -2 \\ 0 & 1 & 2 & 3 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Now, eqn is

$$x_1 - x_3 - 2x_4 = 0 \Rightarrow x_1 = x_3 + 2x_4$$

$$x_2 + 2x_3 + 3x_4 = 0 \Rightarrow x_2 = -2x_3 - 3x_4$$

$$\text{Let, } x_3 = s, x_4 = t$$

$$\bar{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} s+2t \\ -2s-3t \\ s \\ t \end{bmatrix} = s \begin{bmatrix} 1 \\ -2 \\ 1 \\ 0 \end{bmatrix} + t \begin{bmatrix} 2 \\ -3 \\ 0 \\ 1 \end{bmatrix}$$

Q) Find the eigenvalues and eigenvectors of $A = \begin{bmatrix} 5 & 4 \\ 1 & 2 \end{bmatrix}$

$$\text{Sol: } A = \begin{bmatrix} 5 & 4 \\ 1 & 2 \end{bmatrix}$$

$$|A - \lambda I| = 0 \Rightarrow \begin{vmatrix} 5-\lambda & 4 \\ 1 & 2-\lambda \end{vmatrix} = 0 \Rightarrow (2-\lambda)(5-\lambda) - 4 = 0 \Rightarrow 10 - 2\lambda - 5\lambda + \lambda^2 = 4 = 0$$

$$\text{For, } \lambda = 1 \quad \begin{bmatrix} 4 & 4 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 0 \Rightarrow 4x + 4y = 0 \Rightarrow x = -y \quad \therefore \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

$$\text{For, } \lambda = 6 \quad \begin{bmatrix} -1 & 4 \\ 1 & -4 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \Rightarrow -x + 4y = 0 \Rightarrow x = 4y \quad \therefore \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} -4 \\ 1 \end{bmatrix}$$

19

Entropy - The ~~order~~ of ent of is order, uncertainty, randomness associated with the event.

$$H(P) = - \sum_i P_i \log(P_i)$$

Cross entropy - Cross entropy is the ~~good~~ goodness of fit, the closeness of the predictive distribution to the actual one.

Eg: 100 times a coin is tossed
75 heads

case I: 75 heads
25 tails
~~more heads, so~~ ~~low cross entropy~~

case II: 50 heads
50 tails

$$\text{Now, } H(p, q) = - \sum_i p_i \log(q_i)$$

For case I:

$$H(p, q) = -0.75 \log(0.75) - 0.25 \log(0.25) = 0.819$$

For case II:

$$H(p, q) = -0.75 \log(0.75) - 0.25 \log(0.30) = 0.829$$

Note: The cross entropy is smaller for case I than case II as it is closer to the actual distribution.

* In ML problem when we are trying to fit a model with a training data and data is already labelled as +ve and -ve examples, we would like to know how our model is generalizing itself by calculating cross entropy loss, i.e. for a +ve labelled data how good the model is predicting it as positive and for -ve labelled data how good the model is predicting it as negative.

Note - We need to minimize the cross entropy loss.

KL-Divergence - is a difference between cross entropy and entropy.

$$\begin{aligned} \text{KL-Divergence } (P||Q) &= H(p, q) - H(P) \\ &= - \sum_i p_i \log(q_i) + \sum_i p_i \log(P_i) \end{aligned}$$

$$= \sum_i p_i \log(P_i/q_i)$$

$$\text{For case I: } 0.75 \log(0.75/0.75) + 0.25 \log(0.25/0.25) = 0.0032$$

$$\text{For case II: } 0.75 \log(0.75/0.70) + 0.25 \log(0.25/0.30) = 0.0032 - 0.0088$$

• Distance is less in case I than in case II as the case I is closer to the actual distribution.

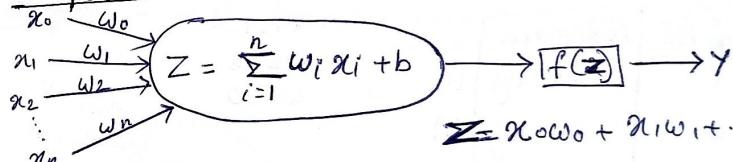
* If a neural network is regression then there will be only one neuron in output layer.

* If a neural network is classification. ~~then there will be min and it's~~ ^{neuron} binary classification. ~~or min of one~~ ^{neuron} and max of 2 ^{neuron} are possible in output layer.

* If a neural network is multi-classification then no. of neural network depends on no. of different inputs.

Architecture of artificial neural network

Perception -



$$z = x_0 w_0 + x_1 w_1 + \dots + x_n w_n$$

inputs = x_0, \dots, x_n

weights = w_0, \dots, w_n

z = summation function

$f \rightarrow$ activation function

Difference b/w biological neuron and perceptron -

- Nervous system is a collection of billions of neurons and neuron is a biological cell consisting of a nucleus where all complex processing tasks takes place.

- The cells consisting of ~~tentacles~~ ^{dendrites} to receive input and the input processed in the nucleus to generate output which is passed through ~~axom~~ axon.

Perception

- In perception there are multiple inputs, for each input there is an assigned weight. The weighted input is compared to dendrites. The sum and activation as a whole can be compared with nucleus and output is compared with axon.

Biological Neuron

i) It has very complex architecture.

ii) The process is?

iii) What processing task is done in nucleus is not known exactly.

Q) Design a perception to predict the performance of a student based on no. of study hours and no. of sleep hours.

Sol: 1) Take dataset

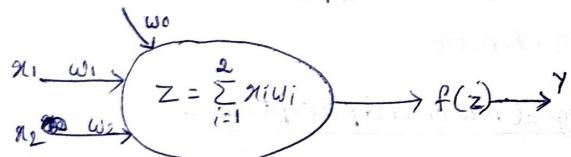
	Study hrs	Sleep hrs	Performance
S1	12	6	Above avg
S2	18	4	Above avg
S3	6	10	Below avg
S4	10	10	Below avg

$$x_0 = 1$$

$$w_1 = 1$$

$$w_2 = 2$$

$$w_0 = 0$$



$$\begin{aligned} Z &= x_0 w_0 + x_1 w_1 + x_2 w_2 \\ &= x_0 w_0 + 12 w_1 + 6 w_2 = 0 + 12 + 12 = 24 \end{aligned}$$

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

$$f(24) = 1 \text{ (as } z \geq 0\text{)}$$

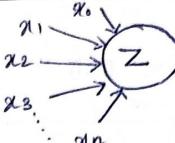
Perception

i) It's simple architecture

ii) Here, we are performing the summation and followed by an activation function to get the output.

- Perception is a mathematical model or a mathematical function, which undergoes 2 no. of processes, 1st is training and 2nd is prediction.
- In the process of training we are finding out the optimal parameter i.e. weight or (weight and Bias).

Geometric Interpretation



$$Z = x_0 w_0 + x_1 w_1 + x_2 w_2 + \dots + x_n w_n$$

$$Z = x_0 w_0 + \sum_{i=1}^n x_i w_i$$

$$Z = \sum_{i=1}^n x_i w_i + b \quad \text{line eqn}$$

$$Z = XW + b$$



$$Ax + By + b = 0$$

$$w_1 x_1 + w_2 x_2 + b = 0$$

$$Ax + By + b = 0$$

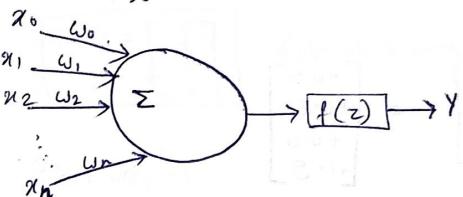
- Perception is a binary classifier because it divides the data into 2 classes by creating regions.

- A perception acts as a line in case of 2D, in 3D its a plane and in 4D onwards it will act as a hyperplane.

159

Perception - is the basic building block of a artificial neural network.

(ANN).



- Perception can be used in logic gates. (AND, OR, XOR)

$$x_1 + x_2 = 1 \cdot 5$$

$$x_0 w_0 + x_1 w_1 + x_2 w_2 = 0$$

$$w_1 = 1, w_2 = 1 \quad b = -1 \cdot 5 = x_0 w_0$$

$$X = \begin{bmatrix} x_0 & x_1 & x_2 \\ 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix}$$

$$w = [-1 \cdot 5 \quad 1 \quad 1]$$

AND Gate

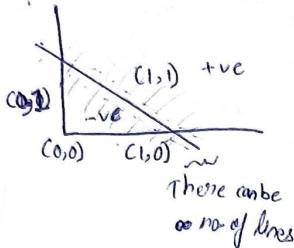
x1	x2	y
0	0	0
0	1	0
1	0	0
1	1	1



$$w x^T = [-0.5 \ 1] \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} -0.5 \\ -0.5 \\ -0.5 \end{bmatrix}$$

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

$$f(z) = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \quad (\text{o/p of AND Gate})$$



OR Gate

$$x_1 + x_2 - 0.5 \text{ (chosen)} \Rightarrow \text{there can be } \infty \text{ no. of lines}$$

$$\begin{aligned} w_1 = 1, w_2 = 1, b = -0.5 \\ x_0 = 1 \text{ (always)} \\ w_0 = -0.5 \end{aligned}$$

TRUTH Table (OR)

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

$$w = [-0.5 \ 1 \ 1]$$

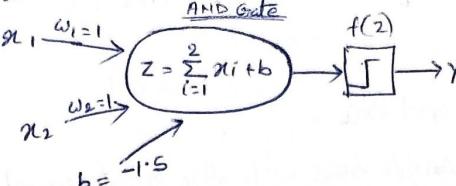
$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$w x^T = [-0.5 \ 1 \ 1] \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} -0.5 \\ +0.5 \\ +0.5 \\ +1.5 \end{bmatrix}$$

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

$$f(z) = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} \quad (\text{o/p of OR Gate})$$

Architecture of AND Gate



$$(i) \ x_1 + x_2 - 1 \text{ for AND Gate}$$

$$(ii) 2x_1 + 2x_2 - 1 \text{ for OR Gate}$$

$$\text{Sol: i) } w_1 = 1, w_2 = 1, w_0 = -1$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$w = [-1 \ 1 \ 1]$$

$$w x^T = [-1 \ 1 \ 1] \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix}$$

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

$$f(z) = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$$

$$\text{(AND Gate) } \neq \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

$$(ii) 2x_1 + 2x_2 - 1$$

$$\text{Sol: ii) } w_1 = 1, w_2 = 1, w_0 = -1$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$w = [-1 \ 2 \ 2]$$

$$w x^T = [-1 \ 2 \ 2] \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \\ 3 \end{bmatrix}$$

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

$$f(z) = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} \quad (\text{o/p of OR Gate})$$

linear activation function in backpropagation:

- If we apply the linear activation function, all layers will collapse to a single layer in multi-layer perceptron.

- Essentially, the linear activation func turns the neural network into just one layer and due to its limited power this doesn't allow the model to create complete mapping b/w input and output.

Non-linear activation function -

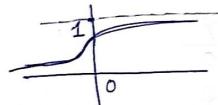
- i) Sigmoid activation function
- ii) tanh
- iii) ReLU
- iv) Softmax
- v) Leaky ReLU
- vi) Parameter ReLU
- vii) Switch
- viii) GELU
- ix) SELU

Sigmoid activation function -

$$f(z) = \frac{1}{1+e^{-z}}, z \rightarrow \text{obj of affine transformation}$$

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

$$f'(z) = y(1-y) = \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}}\right) = \frac{e^{-z}}{(1+e^{-z})^2}$$



Sigmoid function value is from (0 to 1) for $(-\infty, \infty)$

(Q) A perceptron with i/p vector $(x_1, x_2) = (1, 1)$ has weights $(w_1, w_2) = (0.5, 0.5)$ and bias = 0.6. Check whether it is correctly classify the AND function for input (1, 1).

Sol: Truth Table

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

Q) ~~Q2~~ ~~Q3~~

$$w_1 = 0.5, w_2 = 0.5$$

$$x_1 = 1, x_2 = 1$$

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

$$z = w_1 x_1 + w_2 x_2 + b = 0.5 \times 1 + 0.5 \times 1 - 0.6 = 0.5 + 0.5 - 0.6 = 0.4$$

$$f(0.4) = 1 \quad (\text{Yes}) \quad \text{Ans}$$

(Q) A perceptron with learning rate $\eta = 0.1$ and current weights are

$w_1 = 0.2, w_2 = -0.3, b = 0.1$ for input $(1, 2)$ and target output = 1. If perception output is $y = 0$, update the weights and bias.

Sol: $w_1 = 0.2, w_2 = -0.3, b = 0.1, \eta = 0.1$

$$x_1 = 1, x_2 = 2$$

$$\begin{aligned} z &= w_1 x_1 + w_2 x_2 + b \\ &= 0.2 \times 1 + (-0.3) \times 2 + 0.1 = 0.2 - 0.6 + 0.1 = -0.3 \end{aligned}$$

$$f(-0.3) = 0 \quad (\text{update required})$$

$$w_{1\text{ new}} = 0.2 + 0.1(1 - 0) = 0.2 + 0.1 = 0.3$$

$$w_{2\text{ new}} = -0.3 + 0.1(1 - 0) = -0.3 + 0.1 = -0.2$$

$$b = 0.1 + 0.1(1 - 0) = 0.1 + 0.1 = 0.2$$

$$z_{\text{new}} = 0.3 \times 1 + (-0.2) \times 2 + 0.2 = 0.3 - 0.4 + 0.2 = 0.1$$

$$f(0.1) = 1 \quad (\text{Ans})$$

(Q) Give perceptron weights $w_1 = 1, w_2 = -1, b = -0.2$. Write the eqn of decision boundary.

b) Classify the points $(0.5, 0.5)$ and $(0.8, 0.2)$.

Sol: $w_1 = 1, w_2 = -1, b = -0.2$

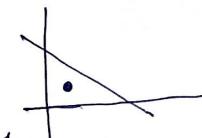
$$\begin{aligned} a) z &= w_1 x_1 + w_2 x_2 + b = 1 \cdot x_1 + (-1) x_2 + (-0.2) \\ &= x_1 - x_2 - 0.2 \quad \text{Ans} \end{aligned}$$

$$b) x_1 = 0.5, x_2 = 0.5$$

$$z = 0.5 - 0.5 - 0.2 = -0.2$$

$$x_1 = 0.8, x_2 = 0.2$$

$$z = 0.8 - 0.2 - 0.2 = 0.4, f(0.4) = 1$$



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



$f(0.2) = 1$ (it will lie in the +ve region).

Q3) Train a perceptron for NAND gate for the input $(x_1, x_2) = (0, 1)$.

Target (y) = NOT ($x_1 \cdot x_2$). Initial weights are $w_1 = 0$, $w_2 = 0$, $b = 0$ and $\eta = 1$. Perform weight update step by step until convergence.

Sol: $\theta = 0$ (if not given)

$$x_0 = 0, x_1 = 0$$

$$z = 0 \times 0 + 0 \times 0 + 0 = 0$$

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

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

$$w_{1,\text{new}} = w_{1,\text{old}} + \eta(t^i - y)x_1$$

$$w_{1,\text{new}} = 0 + 1(1-0)0 = 0$$

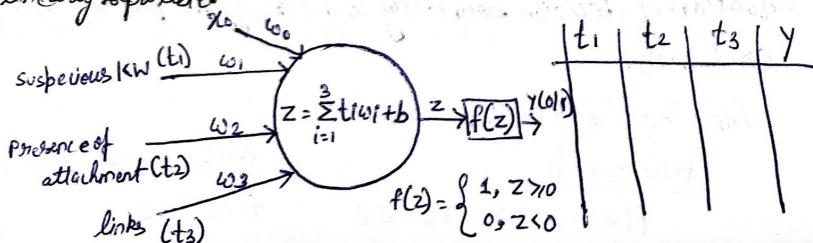
$$w_{2,\text{new}} = 0 + 1(1-0)0 = 0$$

$$b_{\text{new}} = b_{\text{old}} + \eta(t^i - y) = 0 + 1(1-0) = 1$$

$$x_0 = 0, x_1 = 0$$

Q4) You are using a perceptron to classify whether an email is spam or not spam. Input features are no. of suspicious keywords, presence of attachments, no. of links. Target is whether spam (1) or not spam (0). What are the advantages of using a perceptron for this task? What are the limitations and suggest an improvement if the data is not linearly separable?

Sol:



Q5) Define a perceptron and explain its role as a binary classifier.

Q6) Explain the purpose of weights or bias in perceptron.

Q7) Why is a step activation function commonly used in a perceptron?

Q8) What is linear separability and why it is important for a perceptron?

Q9) Explain why a single layer perceptron can't solve the XOR problem.

Q10) How does the choice of activation func affect the learning process of perceptron?

Activation function -

An activation function decides whether a neuron should be activated or not.

3 types -

i) Binary step function

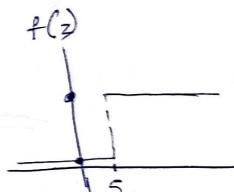
ii) Linear activation function

iii) Non linear activation function.

Binary step function - Depends on a threshold value that decides whether a neuron should be activated or not.

$$f(z) = \begin{cases} 1, & z \geq \theta = 5 \\ 0, & z < 0 \end{cases} \rightarrow \text{threshold}$$

z	$f(z)$	$\theta = 5$	y
0	$f(0)$	$\theta = 5$	0
1	$f(1)$	$\theta = 5$	0
4	$f(4)$	$\theta = 5$	0
14	$f(14)$	$\theta = 5$	1



\Rightarrow Binary activation function is not differentiable.
 $\frac{d(f)}{dz} = 0$ (always constant).

Linear activation function -

- Also known as no activation or identity function. It's where the activation is proportional to input.

$$\boxed{f(x) = x}$$

- If the value of activation function is 0, the o/p value is also 0.

Limitation - The derivative of the func is constant and has no relation to input 'x'. Therefore, we can't use the

Q1) Let, $w_1 = 1.2, w_2 = 0.6, b = 0, m = 0.5, \theta = 1$. Find the optimal line equation for AND Gate.

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

Sol: $Z_1 = w_1 x_1 + w_2 x_2 + b$
 $w_0 = 0, x_0 = 1$

$$Z_1 = +\cancel{2x_0} + 0x_1 \cdot 2 + 0x_0 \cdot 6 + 0 = 0 \text{ (no update)}$$

$$Z_2 = 0x_1 \cdot 2 + 1x_0 \cdot 6 + 0 = 0.6 = f(0.6) = 0 \text{ (no update)}$$

$$Z_3 = 1x_1 \cdot 2 + 0x_0 \cdot 6 + 0 = 1.2 = f(1.2) = 1 \text{ (update required)}$$

$$w_{1\text{new}} = 1.2 + 0.5(0-1) = 1.2 - 0.5 = 0.7$$

$$w_{2\text{new}} = 0.6 + 0.5(0-1) = 0.6$$

$$b_{\text{new}} = 0 + 0.5(0-1) = -0.5$$

$x_0, x_1 = \underline{x_0 = 0, x_1 = 0}$

$$Z_1 = 0x_0 \cdot 7 + 0x_1 \cdot 6 - 0.5 = -0.5 = f(-0.5) = 0$$

$\underline{x_0 = 0, x_1 = 1}$

$$Z_2 = 0x_0 \cdot 7 + 1x_1 \cdot 6 - 0.5 = 0.6 - 0.5 = 0.1 = f(0.1) = 0$$

$\underline{x_0 = 1, x_1 = 0}$

$$Z_3 = 1x_0 \cdot 7 + 0x_1 \cdot 6 - 0.5 = 0.2 = f(0.2) = 0$$

$\underline{x_0 = 1, x_1 = 1}$

$$Z_4 = 1x_0 \cdot 7 + 1x_1 \cdot 6 - 0.5 = 0.7 + 0.6 - 0.5 = 0.8$$



update

Ans

$$w_{1\text{new}} = 0.7 + 0.5(1-0) = 0.7 + 0.5 = \underline{\underline{1.2}}$$

$$w_{2\text{new}} = 0.6 + 0.5(1-0) = 0.6 + 0.5 = \underline{\underline{1.1}}$$

$$b_{\text{new}} = -0.5 + 0.5(1-0) = 0$$

$\underline{x_0 = 0, x_1 = 0}$

$$Z_1 = 0x_1 \cdot 2 + 0x_1 \cdot 1 + 0 = 0 = f(0) = 0$$

$\underline{x_0 = 0, x_1 = 1}$

$$Z_2 = 0x_1 \cdot 2 + 1x_1 \cdot 1 + 0 = 1.1 = f(1.1) = 1 \text{ (update required)}$$

$$w_{1\text{new}} = 1.2 + 0.5(0-1) = 1.2 - 0.5 = 0.7$$

$$w_{2\text{new}} = 1.1 + 0.5(0-1) = 1.1 - 0.5 = 0.6$$

$$b_{\text{new}} = 0 + 0.5(0-1) = -0.5$$

$\underline{x_0 = 0, x_1 = 0}$

$$Z_1 = 0x_1 \cdot 2 + 0x_1 \cdot 6 - 0.5 = f(0.5) = 0$$

$\underline{x_0 = 0, x_1 = 1}$

$$Z_2 = 0x_1 \cdot 2 + 1x_1 \cdot 6 - 0.5 = f(0.1) = 0$$

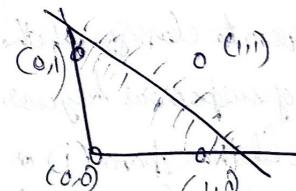
$\underline{x_0 = 1, x_1 = 0}$

$$Z = 0x_1 \cdot 2 + 0x_1 \cdot 6 - 0.5 = 1.2 - 0.5 = 0.7 = f(0.7) = 0$$

$\underline{x_0 = 1, x_1 = 1}$

$$Z = 1x_1 \cdot 2 + 1x_1 \cdot 6 - 0.5 = 1.2 + 0.6 - 0.5 = 1.2 + 0.1 = 1.3 = f(1.3) = 1$$

$\therefore \boxed{1.2 x_1 + 0.6 x_2 - 0.5 = 0}$



Q2) Given perceptron weights $w_1 = 1, w_2 = 1, b = -0.2$. Write the equation of decision boundary and classify the points $(0.5, 0.5)$ and $(0.8, 0.2)$.

Sol: Line eqn = $x_1 + x_2 - 0.2$

for $(0.5, 0.5)$

$$f(x) = x_1 + x_2 - 0.2 = 0.8$$

for $(0.8, 0.2)$

$$f(x) = 0.8$$

Q) Design a perceptron which will act as an OR gate with line eq. $x_1 + x_2 - 0.5$.

Perceptron Learning Rule

Step 1) Randomly initialize the weights and bias.

Step 2) Randomly select a training sample bias with the replacement.

Step 3) Perform the affine transformation.

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

Step 4) Pass the value of z to the step activation function to get y .

Step 5) If y doesn't match with your target then update the weights and bias with the formula

$$\text{if } y \neq t \quad w_{\text{new}} = w_{\text{old}} + m \left(\frac{t - y}{x_i} \right) x_i$$

$$b_{\text{new}} = b_{\text{old}} + m \left(\frac{t - y}{x_i} \right)$$

else

$$w_{\text{new}} = w_{\text{old}}, \quad b_{\text{new}} = b_{\text{old}}$$

Step 6) Repeat this steps 2 to 5 till convergence.

Q) Find the line equation of a perceptron which will act as an OR Gate by application of perceptron learning rule.

The initial value of $w_1 = 0.6, w_2 = 0.6, b = 0, m = 0.5, \theta = 1$

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

Sol: Truth table

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

1) $x_1 = 0, x_2 = 0$.

$$z_1 = 0 \times 0.6 + 0 \times 0.6 = 0$$

$$f(z_1) = f(0) = Y = 0$$

Here, $Y = t$ (no updation)

$$x_1 = 0, x_2 = 1$$

$$z_2 = 0 \times 0.6 + 1 \times 0.6 + 0 = 0.6$$

$$f(0.6) = Y = 0$$

Here, $(Y \neq t)$

So, update weights and bias

$$w_{1\text{new}} = w_{1\text{old}} + 0.5 \left(\frac{1-0}{0.6} \right) 0 = 0.6 + 0 = 0.6$$

$$w_{2\text{new}} = w_{2\text{old}} + 0.5 \left(\frac{1-0}{0.6} \right) 1 = 0.6 + (0.5)1 = 1.1$$

$$b_{\text{new}} = b_{\text{old}} + 0.5 \left(\frac{1-0}{0.6} \right) = 0 + (0.5) = +0.5$$

$$x_1 = 0, x_2 = 0$$

$$z_1 = 0.6 \times 0 + 0 \times 1 + 0 = 0.6 - 0.6 = 0$$

$$f(z_1) = 0 \quad (Y \neq t) \rightarrow \text{no updation}$$

$$x_1 = 0, x_2 = 1$$

$$z_2 = 0 \times 0.6 + 0 \times 1 + 0.5 = 0.5$$

$$f(z_2) = f(-0.4) = 0$$

$$x_1 = 0, x_2 = 1$$

$$z_3 = 0.6 \times 0 + 1 \times 1 + 0.5 = 1.1$$

$$f(1.1) = 1 \quad (Y = t) \rightarrow \text{update}$$

$$x_1 = 1, x_2 = 0$$

~~$$z_4 = 0 \times 0.6 + 1 \times 0 + 0.5$$~~

$$z_4 = 1 \times 0.6 + 1 \times 1 + 0.5$$

$$= 2.2$$

$$f(2.2) = 1 \quad (Y = t)$$



$$[0.6x_1 + 1.1x_2 + 0.5] \text{ Ans}$$

$$x_1 = 0, x_2 = 0$$

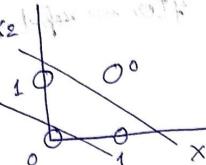
$$z_1 = 0.6 \times 0 + 0 \times 1 + 0.5 = 0.5$$

$$z_2 = 0.6 \times 0 + 0 \times 1 + 0.5 = 0.5$$

$$f(z_2) = f(0.5) = 0 \quad (Y \neq t) \rightarrow \text{no updation}$$

XOR

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



$$XOR = (x_1 + x_2) \cdot (\bar{x}_1 + \bar{x}_2)$$

OR AND NAND

x_1	x_2	y	$(x_1 + x_2)$	$(\bar{x}_1 + \bar{x}_2)$	$h_1 \cdot h_2$
0	0	0	0	1	0
0	1	1	1	1	1
1	0	1	1	1	1
1	1	0	1	0	0

$$OR = x_1 + x_2 - 0.5$$

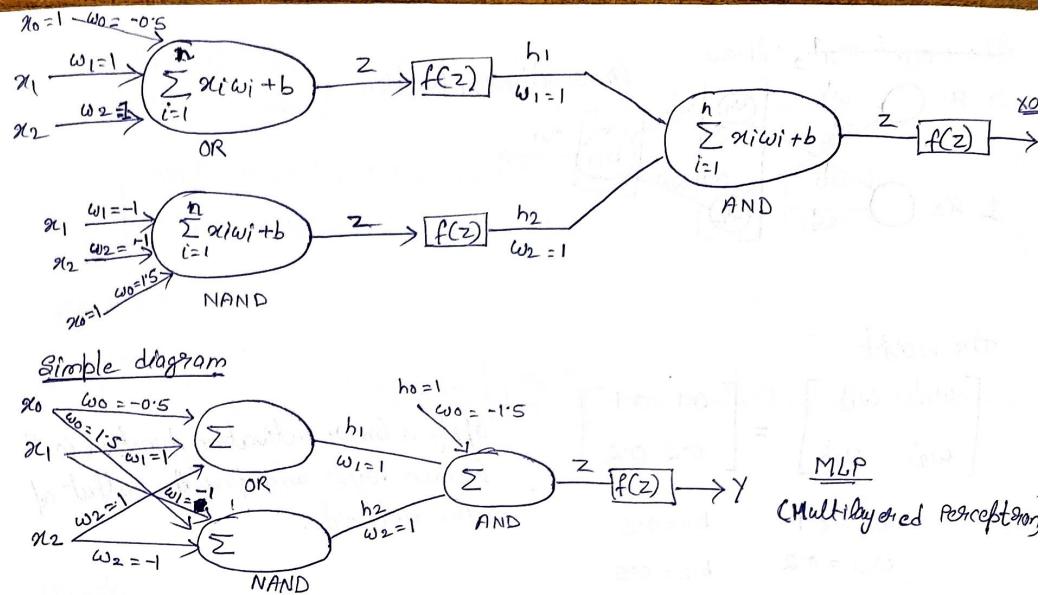
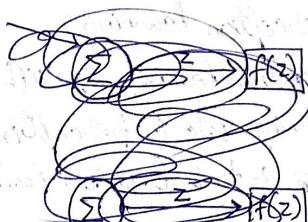
$$AND = x_1 + x_2 - 1.5$$

$$NAND = -x_1 - x_2 + 1.5 \quad (\text{NOT AND})$$

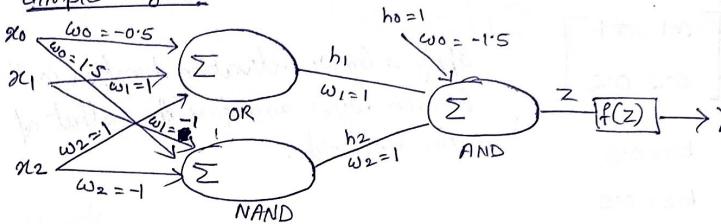
$$\begin{bmatrix} x_0 & x_1 & x_2 \end{bmatrix} \begin{bmatrix} -0.5 \\ 1.5 \\ 1 \\ -1 \end{bmatrix} = \begin{bmatrix} -0.5 & 1.5 \\ 0.5 & 0.5 \\ 0.5 & 0.5 \\ 1.5 & -0.5 \end{bmatrix} f(z) = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

x w z

$$\text{Now, } \begin{bmatrix} x_0 & b_1 & b_2 \end{bmatrix} \begin{bmatrix} -1.5 \\ 1 \\ 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} -0.5 \\ 0.5 \\ 0.5 \\ -0.5 \end{bmatrix} f(z) = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} \neq XOR$$



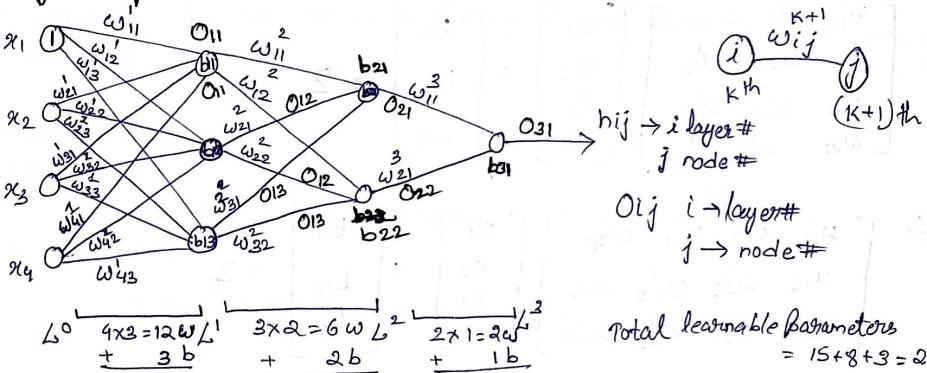
Simple diagram



MLP

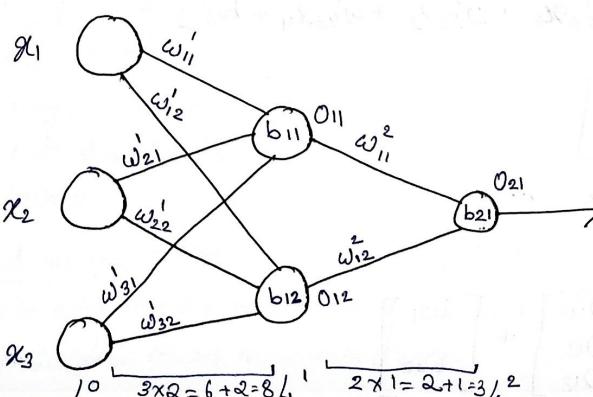
(Multi-layered Perceptron)

Multi-layer Perceptron -



$$\text{Total learnable parameters} = 15 + 8 + 3 = 26$$

Q) Annotate the MLP with weight, bias, output.



$$\text{Total learnable parameters} = 8 + 3 = 11$$

$$C_{\text{new}} = C_{\text{old}} - \text{slope}$$

when,

$$C_{\text{old}} = -10$$

$$\text{slope} = -40$$

$$C_{\text{new}} = -10 - (-40) = -10 + 40 = 30$$

$$\boxed{C_{\text{new}} = C_{\text{old}} - \eta \cdot \text{slope}}$$

$$C_{\text{old}} = -10$$

$$\text{slope} = -40$$

$$C_{\text{new}} = -10 - 0.01(-40) = -10 + 0.40 = -9.6$$

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

$$z = x_1 w_1 + x_2 w_2 + b_2$$

$$L = \frac{1}{2} (y_i - \hat{y}_i)^2 \rightarrow \hat{y}_i \text{ has introduced to make calculation easy (given in book)}$$

$$L = (y_i - x_1 w_1 - x_2 w_2 - b)^2$$

$$L(w_1, w_2, b) = (y_i - x_1 w_1 - x_2 w_2 - b)^2$$

$$w_{1\text{new}} = w_{1\text{old}} - \eta \frac{\partial L}{\partial w_{1\text{old}}}$$

$$w_{2\text{new}} = w_{2\text{old}} - \eta \frac{\partial L}{\partial w_{2\text{old}}}$$

$$b_{\text{new}} = b_{\text{old}} - \eta \frac{\partial L}{\partial b_{\text{old}}}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}_i} \times \frac{\partial \hat{y}_i}{\partial z} \times \frac{\partial z}{\partial w_1} = -\frac{1}{2} (y_i - \hat{y}_i) \times 1 \times x_1$$

$$w_{1\text{new}} = w_{1\text{old}} - \eta (-(y_i - \hat{y}_i) x_1)$$

$$= w_{1\text{old}} + \eta (y_i - \hat{y}_i) x_1$$

$$w_{2\text{new}} = [w_{2\text{old}} + \eta (y_i - \hat{y}_i) x_2]$$

$$b_{\text{new}} = [b_{\text{old}} + \eta (y_i - \hat{y}_i)]$$

Weight update formula for linear and binary neuron.

9/10

Q3) A sigmoid perceptron is initialized with weights $(w_1, w_2) = (0.7, 0.8)$ and bias $b = 0$. For given input $(x_1, x_2) = (1, 0)$, target $t = 0$ and learning rate $\eta = 0.2$. What will be the new weights and bias after one step update?

For sigmoid function

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_1}$$

$$\frac{\partial}{\partial z} \left(\frac{1}{1+e^{-z}} \right) = \frac{e^{-z}}{(1+e^{-z})^2} = \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}} \right) = \hat{y}(1-\hat{y})$$

$$\therefore \frac{\partial L}{\partial w_1} = -(y - \hat{y}) \hat{y}(1-\hat{y}) x_1$$

$$\frac{\partial L}{\partial w_2} = -(y - \hat{y}) \hat{y}(1-\hat{y}) x_2$$

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

$$w_{1\text{new}} = w_{1\text{old}} + \eta (y - \hat{y}) \hat{y}(1-\hat{y}) x_1$$

$$w_{2\text{new}} = w_{2\text{old}} + \eta (y - \hat{y}) \hat{y}(1-\hat{y}) x_2$$

$$b_{\text{new}} = b_{\text{old}} + \eta (y - \hat{y}) \hat{y}(1-\hat{y})$$

11/10

Q4) What's a perceptron? What are the main components of a perceptron?

Q5) Explain the working principle of a perceptron. What is the role of weights and bias in a perceptron? Why is an activation function used in a perceptron?

Q6) What type of problem can a single layer perceptron solve?

Q7) What's the difference between perception and neuron in a biological brain?

Q8) How does the perceptron learning rule work?

Q9) Explain the perceptron update rule.

Q10) What's meant by convergence in perceptron learning?

Q11) Under what condition does the perceptron learning algorithm converges?

Q12) What's the effect of learning rate in perceptron training? What happens if the learning rate is too high or too low?

Q13) What's the stopping cond'n for perceptron learning?

- (Q) What for a 3-class perception the raw output are $z = [1.2, 0.8, -0.5]$. Compute the softmax probabilities for each class.
- Sol: Softmax can be used only for last layer of classification.

$$\begin{aligned} z_1 &= 1.2 \rightarrow \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}} \\ z_2 &= 0.8 \rightarrow \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}} \\ z_3 &= -0.5 \rightarrow \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}} \end{aligned}$$

$$z_i = e^{z_i} + e^{z_2} + e^{z_3} = 6.132$$

(Probabilities)

$$\begin{aligned} z_1 &= 0.54, z_2 = 0.36, z_3 = 0.09 \\ 0.999 &\approx 1 \\ \text{upon adding} \end{aligned}$$

As $z_1 = 0.54$ (highest). Sample belongs to class 1.

- (Q) $z = [1000, 1001, 1002]$. Calculate the softmax probabilities for each class.

Sol: Softmax can be used only for last layer of classification.

$$\begin{aligned} z_i &= \frac{e^{z_i}}{\sum e^{z_i}} \\ z_1 &= \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}} \\ z_2 &= \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}} \\ z_3 &= \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}} \end{aligned}$$

$$z_i = e^{1000} + e^{1001} + e^{1002} = 2.1882 \times 10^{435}$$

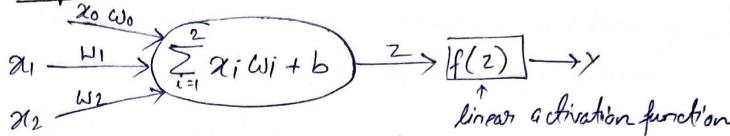
$$z_1 = \frac{e^{1000}}{2.188 \times 10^{435}} = 7.48 \times 10^{-432}$$

$$z_2 = \frac{e^{1001}}{2.188 \times 10^{435}} = 2.90 \times 10^{-435} \quad 2.955 \times 10^{-435}$$

$$z_3 = \frac{e^{1002}}{2.188 \times 10^{435}} = 1.503 \times 10^{-436}$$

Highest value is z_1 . So sample belongs to class 1.

6/10/25
Perception as a Regression -



$$\begin{aligned} x_1 w_1 + x_2 w_2 + b &= z \\ f(z) &= z \end{aligned}$$

$$\boxed{f(x) = x} \rightarrow \text{linear activation function}$$

Perception training (Linear)

$$w_{1\text{new}} = w_{1\text{old}} + \eta(t-y)x_1$$

$$w_{2\text{new}} = w_{2\text{old}} + \eta(t-y)x_2$$

$$b_{\text{new}} = b_{\text{old}} + \eta(t-y)$$

$$\eta = 0 \text{ to } 1 \text{ (always)}$$

$$\text{ideal } \eta = 0.01$$

\hookrightarrow loss (error)

$$L = 80$$

$$L = 60$$

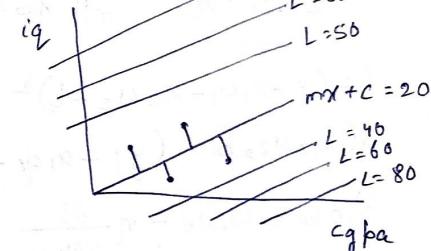
$$L = 50$$

$$mx + c = 20$$

$$L = 40$$

$$L = 60$$

$$L = 80$$



$$L = (y_i - \hat{y}_i)^2$$

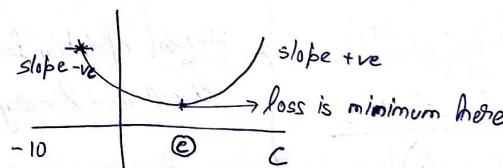
\downarrow Target \hat{y}_i Predicted

$$y_i = mx_i + c$$

$$L(m, c) = (y_i - mx_i - c)^2$$

$$L(c) = (y_i - 78.5x_i - c)^2, m = 78.5 \text{ (constant)}$$

$$L = c^2 \text{ (as } L \text{ depends on } c \text{ only).}$$



$$\text{slope} = \frac{dL}{dt}$$

$$f(x) = x^2$$

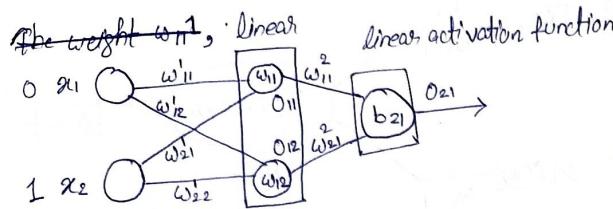
$$\frac{d}{dx} f(x) \Big|_{x=2} = \frac{d}{dx} x^2$$

$$= 2x$$

$$= 4$$

Slope is -ve (value is increasing)

Slope is +ve (value is decreasing)



The weight

$$\begin{bmatrix} w_{11}^1 & w_{12} \\ w_{21}^1 & w_{22}^1 \end{bmatrix} = \begin{bmatrix} 0.1 & 0.1 \\ 0.2 & 0.2 \end{bmatrix}$$

$$w_{11}^1 = 0.1 \quad b_{11} = 0.5$$

$$w_{12}^1 = 0.2 \quad b_{12} = 0.5$$

apply a linear activation function in the hidden layer and find the output of the network.

Layer #1

$$W = \begin{bmatrix} w_{11}^1 & w_{12}^1 & w_{13}^1 \\ w_{21}^1 & w_{22}^1 & w_{23}^1 \\ w_{31}^1 & w_{32}^1 & w_{33}^1 \\ w_{41}^1 & w_{42}^1 & w_{43}^1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix} = Z$$

$$\Rightarrow W^T = \begin{bmatrix} w_{11}^1 & w_{21}^1 & w_{31}^1 & w_{41}^1 \\ w_{12}^1 & w_{22}^1 & w_{32}^1 & w_{42}^1 \\ w_{13}^1 & w_{23}^1 & w_{33}^1 & w_{43}^1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix} = Z$$

$$\Rightarrow \begin{bmatrix} w_{11}^1 x_1 + w_{21}^1 x_2 + w_{31}^1 x_3 + w_{41}^1 x_4 + b_{11} \\ w_{12}^1 x_1 + w_{22}^1 x_2 + w_{32}^1 x_3 + w_{42}^1 x_4 + b_{12} \\ w_{13}^1 x_1 + w_{23}^1 x_2 + w_{33}^1 x_3 + w_{43}^1 x_4 + b_{13} \end{bmatrix} = Z$$

$$\Rightarrow f(Z) = \begin{bmatrix} O_{11} \\ O_{12} \\ O_{13} \end{bmatrix}$$

Layer #2

$$W = \begin{bmatrix} w_{11}^2 & w_{12}^2 \\ w_{21}^2 & w_{22}^2 \\ w_{31}^2 & w_{32}^2 \end{bmatrix} \begin{bmatrix} O_{11} \\ O_{12} \\ O_{13} \end{bmatrix} + \begin{bmatrix} b_{21} \\ b_{22} \end{bmatrix}$$

$$\Rightarrow W^T = \begin{bmatrix} w_{11}^2 & w_{21}^2 & w_{31}^2 \\ w_{12}^2 & w_{22}^2 & w_{32}^2 \end{bmatrix} \begin{bmatrix} O_{11} \\ O_{12} \\ O_{13} \end{bmatrix} + \begin{bmatrix} b_{21} \\ b_{22} \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} w_{11}^2 O_{11} + w_{21}^2 O_{12} + w_{31}^2 O_{13} + b_{21} \\ w_{12}^2 O_{11} + w_{22}^2 O_{12} + w_{32}^2 O_{13} + b_{22} \end{bmatrix} = f(Z) = \begin{bmatrix} O_{21} \\ O_{22} \end{bmatrix}$$

Layer #2

$$\begin{bmatrix} w_{11}^3 & w_{21}^3 \\ w_{21}^3 & w_{31}^3 \end{bmatrix} \begin{bmatrix} O_{21} \\ O_{22} \end{bmatrix} = \begin{bmatrix} w_{11}^3 & w_{21}^3 \\ w_{21}^3 & w_{31}^3 \end{bmatrix} \begin{bmatrix} O_{21} \\ O_{22} \end{bmatrix} + b_{31} = w_{11}^3 O_{21} + w_{21}^3 O_{22} + b_{31}$$

$$f(Z) = O_{31}$$

Answers

- 1) A perceptron is the simplest type of Artificial Neural Network a computational model inspired by the way biological neurons work. It's mainly used for binary classification.

Main Components -

i) Input Layer

Take multiple input values (features)

ii) Weights (w_1, w_2, \dots)

Each function is assigned a weight showing its importance in decision.

iii) Summation function

iv) Activation function

v) Output

2) Working of perceptron

i) Input Reception

ii) Weighted sum calculation

iii) Activation function

iv) Learning

Role of weights

- Weights determine the importance of each input feature on the output.
- Larger weights means more importance.

- v) Weight updates - Use an optimization algorithm (like gradient descent) to adjust weights and biases to minimize the loss.
- v) Repeat the above steps ~~repe~~ iteratively over many epochs until the network converges to a good solution.

Forward propagation in MLP - Forward propagation is the process of passing input through the network to get output.

Step 1: For each neuron compute a weighted sum of inputs and bias.

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

Step 2: Apply an activation function to get the neuron output, ~~except~~ output of 1 layer becomes input to the next layer.

Continue layer by layer until final output layer produces predictions.

Loss function - is a method of evaluating how well your algorithm is modelling your data. If the loss function value is high your algorithm is performing poorly and if the loss function value is small your algorithm is performing great.

- Loss function is a mathematical func of the parameters of Machine learning parameter. If we adjust the parameters the loss function value changes.

Categories of loss function -

~~Classification~~

Regression	Classification	Autoencoder	GAN
MSE	Primary cross entropy	KL Divergence	Discriminator loss
MAE	Categorical cross entropy		
huber loss	hinge loss		min max GAN loss.

23/10

Regression

① Mean Squared Error (MSE) / L2 loss / squared loss -

Difference b/w loss function and cost function

Loss function

is calculated for a single data sample

$$\text{loss} = (y_i - \hat{y}_i)^2$$

Cost function

is calculated for a batch (entire data)

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

cgpa	iq	package (lpa)	Predicted Value (\hat{y})	error	error ²
6.3	100	6.3	6.1	0.2	0.04
7.1	91	4.1	4	0.1	0.01
8.5	83	3.5	3.7	-0.2	0.04
9.2	102	7.2	7	0.2	0.04
		90	8	3.2	10.24

Advantages and Limitations of MSE -

Advantage

1. It is easy to interpret
2. MSE is differentiable
3. It has only one local minima

$$\frac{\partial L}{\partial y_i} = -2(y_i - \hat{y}_i)$$

Disadvantage

1. error unit is different
2. It is not robust to outliers

MLP setup for regression problem where MSE is used or loss function is used.

- output layer should contain one neuron and for output neuron the linear activation function is required.

- We can use MSE as a loss function if the dataset is free of outliers.

② Mean Absolute Error (MAE)

Loss function

$$L = |y_i - \hat{y}_i|$$

Cost function

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

For layer 2

$$\begin{bmatrix} w_{11}^2 & w_{12}^2 \\ w_{21}^2 & w_{22}^2 \\ w_{31}^2 & w_{32}^2 \end{bmatrix} = \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.1 \\ 0.2 & 0.2 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 0.1 & 0.1 & 0.2 \\ 0.1 & 0.1 & 0.2 \end{bmatrix} \begin{bmatrix} 0.5 \\ 0.6 \\ 0.7 \end{bmatrix} =$$

$$\Rightarrow \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.1 \\ 0.2 & 0.2 \end{bmatrix} \begin{bmatrix} 0.22 \\ 0.23 \\ 0.24 \end{bmatrix} + \begin{bmatrix} 0.5 \\ 0.6 \end{bmatrix} = \begin{bmatrix} 0.593 \\ 0.693 \end{bmatrix}$$

For layer 3

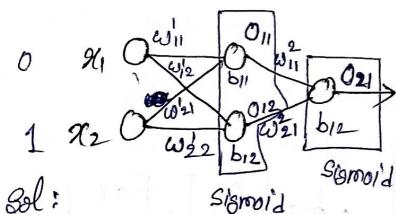
$$\begin{bmatrix} w_{11}^3 & w_{12}^3 \end{bmatrix} = \begin{bmatrix} 0.5 & 0.6 \end{bmatrix}, b_{31} = 0.3$$

~~0.5~~

~~0.6~~

$$\begin{bmatrix} 0.5 & 0.6 \end{bmatrix} \begin{bmatrix} 0.593 \\ 0.693 \end{bmatrix} + [0.3] = \boxed{1.0123}$$

Q) Annotate the given network and output generated



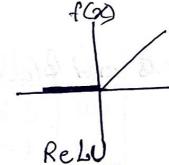
$$\begin{bmatrix} w_{11}^1 & w_{12}^1 \\ w_{21}^1 & w_{22}^1 \end{bmatrix} = \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.1 \end{bmatrix}, \begin{bmatrix} b_{11} \\ b_{12} \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.6 \end{bmatrix}$$

$$\sigma \left(\begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} + \begin{bmatrix} 0.5 \\ 0.6 \end{bmatrix} \right) = \sigma \left[\begin{bmatrix} 0.6 \\ 0.7 \end{bmatrix} \right] = \begin{bmatrix} \frac{1}{1+e^{-0.6}} \\ \frac{1}{1+e^{-0.7}} \end{bmatrix} = \begin{bmatrix} 0.6456 \\ 0.6681 \end{bmatrix}$$

$$\begin{aligned} \text{Layer 2} \quad \sigma \left(\begin{bmatrix} 0.6456 & 0.6681 \\ 0.6681 & 0.693 \end{bmatrix} \begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix} + 0.6 \right) &= \begin{bmatrix} 0.6456 & 0.6681 \\ 0.6681 & 0.693 \end{bmatrix} \begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix} + 0.6 \\ &= 0.3963 + 0.6 = \boxed{0.99636} \end{aligned}$$

16/10 ReLU activation -

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



- It is differentiable
 $f'(x) = \begin{cases} 0, & x \le 0 \\ 1, & x > 0 \end{cases}$

x	f(x)
100	100
-80	0
5	5
-4	0

- ReLU stands for Rectified Linear Unit. If the input 'x' is (+ve) then $f(x) = x$. If, input is (+ve) then $f(x) = x$.

- Although it gives an impression of linear function, ReLU has a derivative function and allows for back propagation making it computationally efficient.

- ReLU doesn't make all the neurons active at the same time. The neuron will only be deactivated if the output of the linear transformation is less than 0.

- The limitation is, ReLU faces a Dying ReLU problem.

Training process of MLP -

General Training process of MLP -

Training process of MLP involves the following steps -

i) Forward propagation - Here, input data passes through the network layer by layer producing an output.

ii) Loss calculation - Compare the output to the true labels using a loss function (MSE, cross entropy).

iii) Back propagation - calculate gradients of the loss with respect to each weight by propagating error backward through the network.

Role of perception

- Acts as a decision making unit.
- Classifies data into one of two categories based on weighted combination of inputs.

3) The activation function is used in a perceptron to introduce non-linearity and decide whether the ~~neuron~~ neuron should fire (activate) or not.

~~It also uses~~

4) A single layer perceptron can only solve linearly separable problems.
Eg: AND Gate, OR Gate can be solved but XOR can't.

5) The perceptron Learning rule is a method to adjust the weights and bias of a perceptron so it can correctly classify training examples.

1) Initialise weights and bias

2) Compute output

3) Calculate error ($t - y$)

4) Update weights and bias

5) Repeat until max no. of iterations is reached.

6) It's a simple error correction mechanism that allows the perceptron to learn linearly separable patterns by gradually adjusting weights and bias towards correct classification.

7) In perceptron learning, convergence refers to the point at which the ~~perceptron~~ perceptron stops updating its weights and bias ~~so it has correctly classified in all training examples~~.

10) Learning rate is a small +ve constant that controls how much weight and bias are adjusted during each update.

Too high

- The perceptron may overshoot the correct solution.
- Can lead to ~~no~~ oscillation and failure to converge.

Too low

- Learning becomes slow.
- May take long time to converge.

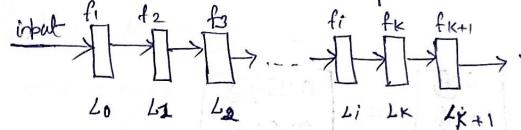
11) Stopping condition for perceptron learning are -

i) All training examples are correctly classified.

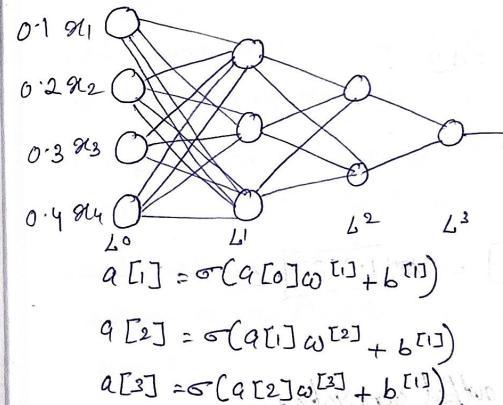
ii) Max no. of iterations (epochs) Reached.

13/10

Neural network is a cascade of non-linear functions



$f_k(f_i(\dots f_3(f_2(f_1(x)))) \dots)$ → cascade of non-linear functions



$$a[1] = \sigma(a[0]w^{[1]} + b^{[1]})$$

$$a[2] = \sigma(a[1]w^{[2]} + b^{[2]})$$

$$a[3] = \sigma(a[2]w^{[3]} + b^{[3]})$$

For 1st layer

12 weights, 3 bias

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{43} \end{bmatrix} = \begin{bmatrix} 0.1 & 0.1 & 0.1 \\ 0.2 & 0.2 & 0.2 \\ 0.1 & 0.1 & 0.1 \\ 0.2 & 0.2 & 0.2 \end{bmatrix}, b = \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix} = \begin{bmatrix} 0.6 \\ 0.7 \\ 0.5 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{41} \\ w_{21} & w_{22} & w_{23} & w_{42} \\ w_{31} & w_{32} & w_{33} & w_{43} \end{bmatrix} \begin{bmatrix} 0.1 \\ 0.2 \\ 0.3 \\ 0.4 \end{bmatrix} + \begin{bmatrix} 0.6 \\ 0.7 \\ 0.5 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 0.1 & 0.2 & 0.1 & 0.2 \\ 0.1 & 0.2 & 0.1 & 0.2 \\ 0.1 & 0.2 & 0.1 & 0.2 \end{bmatrix} \begin{bmatrix} 0.1 \\ 0.2 \\ 0.3 \\ 0.4 \end{bmatrix} + \begin{bmatrix} 0.6 \\ 0.7 \\ 0.5 \end{bmatrix} = \boxed{\begin{bmatrix} 0.22 \\ 0.23 \\ 0.24 \end{bmatrix}}$$

			Bredicted value	error	error
cgpa	iq	package	lpa		
				0.2	0.2
				0.1	0.1
				-0.2	0.2
				0.2	0.2
				3.2	3.2

Advantages

1. It's intuitive to understand.
2. The error unit is same.
3. It's robust to outliers.
- In MAE each error contributes linearly to the total loss.

Disadvantages

1. It's not differentiable.

② Huber loss

$$L = \begin{cases} \frac{1}{2}(\gamma - \hat{\gamma})^2 & \text{for } |\gamma - \hat{\gamma}| \leq \delta \\ \delta(1|\gamma - \hat{\gamma}| - \frac{1}{2}\delta) & \text{for } |\gamma - \hat{\gamma}| > \delta \end{cases}$$

→ acts as MSE
→ acts as MAE

Q3 Given $\gamma = 5$, $\hat{\gamma} = 3$, $\delta = 1$. Find the huber loss

Sol: $e = |\gamma - \hat{\gamma}| = |5 - 3| = 2$

$e > \delta$, so linear branch will be chosen

$$L = 1(2 - \frac{1}{2}(1)) = 2 - 0.5 = \boxed{1.5}$$

Q4 Given $\gamma = 1$, $\hat{\gamma} = 0.5$, $\delta = 1$. Find the huber loss

Sol: $e = |\gamma - \hat{\gamma}| = |1 - 0.5| = 0.5$

$e < \delta$, so quadratic branch will be chosen

$$L = \frac{1}{2}(0.5)^2 = \frac{1}{2} \times 0.25 = \boxed{0.125}$$

Advantage

- 0 It is also differentiable (when 25% data is outliers)

Classification -

1) Binary Cross entropy / log loss -

- It is used for binary classification problem.

$$\text{Loss } L = -\gamma_i \log(\hat{\gamma}_i) - (1-\gamma_i) \log(1-\hat{\gamma}_i)$$

$$\text{cost} = -\frac{1}{n} \sum_{i=1}^n [\gamma_i \log(\hat{\gamma}_i) + (1-\gamma_i) \log(1-\hat{\gamma}_i)]$$

cgpa	iq	performance (γ_i)	$\hat{\gamma}_i$
80	8	1	0.73
60	6	0	0.27
50	5	0	0.27

25/10

Q1) What's role of hidden layers in MLP?

Q2) How does an MLP learn from data explain the general training process?

Q3) What are common activation functions used in MLPs and what are their pros and cons?

Q4) How does the no. of hidden layers and neurons affect model capacity and performance?

Q5) What is the loss function in the MLP training?

Q6) How can MLP be used for classification and regression?

Q7) MLP best suited for what kind of problems?

2) Categorical cross entropy - (Used in softmax regression)

cgpa	iq	placed	one hot encoding
8	80	Yes	1 0 0
6	60	No	0 1 0
7	70	may be	0 0 1

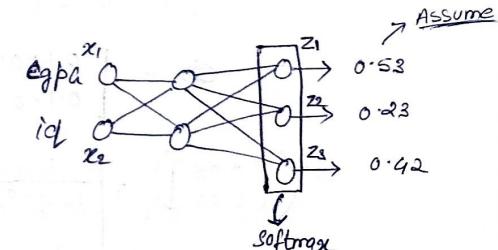
- It is used in a loss function in multiclass classification problem.

$$L = -\sum_{j=1}^K \gamma_j \log(\hat{\gamma}_j)$$

* K = no. of classes

$\gamma_j \rightarrow$ true class

$\hat{\gamma}_j \rightarrow$ predicted class



$$\hat{\gamma} = [0.53, 0.23, 0.42]$$

$$\gamma = [1 \ 0 \ 0]$$

$$L = -\gamma_1 \log(\hat{\gamma}_1) - \gamma_2 \log(\hat{\gamma}_2) - \gamma_3 \log(\hat{\gamma}_3)$$

* Now for the multiclass classifier loss function applied categorical cross entropy

Sol: a) Layer #1

$$\rightarrow \left(\begin{bmatrix} 0.2 & -0.4 \\ 0.7 & 0.1 \end{bmatrix} \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix} \right) = \sigma \begin{bmatrix} 0.06 \\ 0.66 \end{bmatrix} = \begin{bmatrix} \frac{1}{1+e^{-0.06}} \\ \frac{1}{1+e^{-0.66}} \end{bmatrix}$$

~~$\begin{bmatrix} 0.515 \\ 0.66 \end{bmatrix}$~~ = $\begin{bmatrix} 0.515 \\ 0.66 \end{bmatrix} = \begin{bmatrix} 0.11 \\ 0.12 \end{bmatrix}$

Layer #2



$$\begin{bmatrix} 0.3 & -0.5 \end{bmatrix} \begin{bmatrix} 0.515 \\ 0.66 \end{bmatrix} + \begin{bmatrix} 0.05 \end{bmatrix} = \begin{bmatrix} -0.125 \end{bmatrix} = \hat{y}$$

b) $w_{11}^2, w_{21}^2, b_{21}$, loss function is squared loss $\Rightarrow L = (\hat{y} - y)^2$

$$w_{11}^2_{\text{new}} = w_{11}^2 - \eta \frac{\partial L}{\partial w_{11}^2} = 2(\hat{y} - y)0_{11} = 2(-0.125 - 1)0_{11}$$

$$w_{21}^2_{\text{new}} = w_{21}^2 - \eta \frac{\partial L}{\partial w_{21}^2} = 2(\hat{y} - y)0_{12} = 2(-0.125 - 1)0_{12}$$

$$b_{21}_{\text{new}} = b_{21}_{\text{old}} - \eta \frac{\partial L}{\partial b_{21}} = 2(\hat{y} - y) = 2(-0.125 - 1)$$

$$w_{11}^2_{\text{new}} = 0.3 - 0.1 \times 2(-0.125 - 1)0.515 = 0.415$$

$$w_{21}^2_{\text{new}} = -0.5 - 0.1 \times 2(-0.125 - 1)0.66 = -0.3515$$

$$b_{21}_{\text{new}} = -0.05 - 0.1 \times 2(-0.125 - 1) = 0.275$$

c) find the ~~that~~ updated weights and biases for the input to hidden layer.

110

Training the MLP -

- Forward propagation ✓
- Back propagation ✓

Neural network performance enhancement -

- i) Tuning the hyper parameters
 - a) hidden layer
 - b) batch
 - c) No. of neurons
 - d) learning rate
 - e) activation function
 - f) optimizers
 - g) epochs

ii) Solving neural network problems

- a) vanishing gradient
- b) ~~verifying~~ problem
- c) overfitting problem.

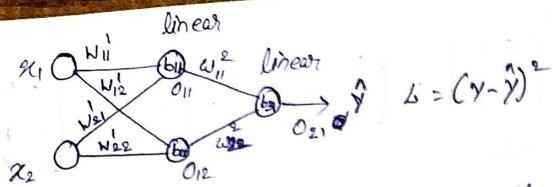
Variants of Gradient descent -

- i) Batch GD
- ii) stochastic GD
- iii) Mini Batch GD

Batch GD

epoch (1) Training Dataset = 100 samples

Batch → NN → $\hat{y}_{\text{batch}} = y - b$, update weights and biases



9 learnable parameters are there in the above diagram

minimize the loss

$$L = (y - \hat{y})^2$$

$$L(\hat{y}) = (\hat{y} - y)^2$$

$$z = \hat{y} = O_{11}w_{11}^2 + O_{12}w_{21}^2 + b_2$$

$$\hat{y} = \left[\frac{x_1 w_{11}^1 + x_2 w_{21}^1 + b_{11}}{O_{11}} \right] w_{11}^2 + \left[\frac{x_1 w_{12}^1 + x_2 w_{22}^1 + b_{12}}{O_{12}} \right] w_{21}^2 + b_{21}$$

1st

$$w_{11\text{new}} = w_{11\text{old}} - \eta \frac{\partial L}{\partial w_{11}^2} = 2(\hat{y} - y)O_{11}$$

$$w_{21\text{new}} = w_{21\text{old}} - \eta \frac{\partial L}{\partial w_{21}^2} = 2(\hat{y} - y)O_{12}$$

$$b_{21\text{new}} = b_{21\text{old}} - \eta \frac{\partial L}{\partial b_{21}} = 2(\hat{y} - y)$$

$$w_{11\text{new}} = w_{11\text{old}} - \eta \frac{\partial L}{\partial w_{11}^1}$$

$$w_{12\text{new}} = w_{12\text{old}} - \eta \frac{\partial L}{\partial w_{12}^1}$$

$$w_{21\text{new}} = w_{21\text{old}} - \eta \frac{\partial L}{\partial w_{21}^1}$$

$$w_{22\text{new}} = w_{22\text{old}} - \eta \frac{\partial L}{\partial w_{22}^1}$$

$$b_{11\text{new}} = b_{11\text{old}} - \eta \frac{\partial L}{\partial b_{11}}$$

$$b_{12\text{new}} = b_{12\text{old}} - \eta \frac{\partial L}{\partial b_{12}}$$

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

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

$$\frac{\partial L}{\partial b_{21}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial b_{21}}$$

$$\frac{\partial}{\partial \hat{y}} (\hat{y} - y)^2 = -2(\hat{y} - y) = 2(y - \hat{y})$$

$$\frac{\partial \hat{y}}{\partial z} = \frac{\partial f(z)}{\partial z} = \frac{\partial z}{\partial z} = 1$$

$$\frac{\partial z}{\partial w_{11}^2} = O_{11}$$

2nd

$$\frac{\partial L}{\partial w_{11}^1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_f} \cdot \frac{\partial z_f}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial z_{11}} \cdot \frac{\partial z_{11}}{\partial w_{11}^1} = 2(\hat{y} - y) w_{11}^2 x_1$$

$$\frac{\partial L}{\partial w_{21}^1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_f} \cdot \frac{\partial z_f}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial z_{11}} \cdot \frac{\partial z_{11}}{\partial w_{21}^1} = 2(\hat{y} - y) w_{21}^2 x_2$$

$$\frac{\partial L}{\partial b_{11}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_f} \cdot \frac{\partial z_f}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial z_{11}} \cdot \frac{\partial z_{11}}{\partial b_{11}} = 2(\hat{y} - y) b_{11} = 2(\hat{y} - y) w_{11}^2$$

27/10

a) Design an MLP with the following setup

Input (2) \rightarrow hidden (2) \rightarrow output (1)
(sigmoid) (linear)

a) Find the \hat{y} of the model for the following iteration

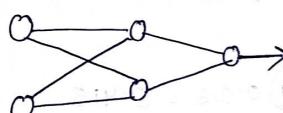
$$x = \begin{bmatrix} 0.6 \\ 0.9 \end{bmatrix}, y_{\text{true}} = 1, \eta = 0.1$$

$$w^{[1]} = \begin{bmatrix} 0.2 & -0.4 \\ 0.7 & 0.1 \end{bmatrix}, b^{[1]} = \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix}$$

$$w^{[2]} = \begin{bmatrix} 0.3 & -0.5 \end{bmatrix}, b^{[2]} = 0.05$$

b) find the updated weights, and bias for hidden to output layer by choosing a suitable loss function for the given setup.

Sol:



* If a MLP solves a regression problem the output layer should contain only 1 neuron and on that neuron linear activation function need to be applied.

* If an MLP acts as binary classifier the output layer should contain one neuron and sigmoid activation function need to be applied on that neuron.

* If an MLP acts as a multiclass classifier the output layer should contain K no. of neurons where, K = no. of classes and on the output layer softmax activation function need to be applied.

* Now, for regression function loss function applied MSE, MAE, huber loss.

* " " Binary function binary cross entropy.

$$\Rightarrow L = -1 \log(0.53) - 0 \log(0.23) - 0 \log(0.42) \\ = -\log(0.53) = 0.275$$

b) Design a MLP with the following information.

input(2) \rightarrow hidden(2) \rightarrow output(3)

(sigmoid) (softmax)

c) find the predicted values using forward propagation

$$w^{[1]} = \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.1 \end{bmatrix}, b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$w^{[2]} = \begin{bmatrix} 0.2 & 0.2 \\ 0.1 & 0.1 \\ 0.3 & 0.3 \end{bmatrix}, b^{[2]} = \begin{bmatrix} 0.2 \\ 0.2 \\ 0.4 \end{bmatrix}$$

$$x = \begin{bmatrix} 0.2 \\ 0.3 \end{bmatrix}$$

b) Find the categorical cross entropy loss for this sample having true value $y = [1, 0, 0]$

Sol:



a) Layer #1

$$\left(\begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{21} \end{bmatrix} \right) = \begin{bmatrix} 0.11 \\ 0.12 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.1 \end{bmatrix} \begin{bmatrix} 0.2 \\ 0.3 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.1 \end{bmatrix} \begin{bmatrix} 0.2 \\ 0.3 \end{bmatrix}$$

$$\Theta = \begin{bmatrix} 0.1 \times 0.2 + 0.1 \times 0.3 \\ 0.1 \times 0.2 + 0.1 \times 0.3 \end{bmatrix} = \begin{bmatrix} 0.05 \\ 0.05 \end{bmatrix}$$

$$\sigma \left(\begin{bmatrix} 0.05 \\ 0.05 \end{bmatrix} \right) = \frac{1}{1+e^{-0.05}} = \begin{bmatrix} 0.5125 \\ 0.5125 \end{bmatrix}$$

Layer #2

$$\begin{bmatrix} 0.2 & 0.2 \\ 0.1 & 0.1 \\ 0.3 & 0.3 \end{bmatrix} \begin{bmatrix} 0.5125 \\ 0.5125 \end{bmatrix} + \begin{bmatrix} 0.2 \\ 0.3 \\ 0.4 \end{bmatrix} = \begin{bmatrix} 0.2048 \\ 0.1024 \\ 0.3072 \end{bmatrix} + \begin{bmatrix} 0.2 \\ 0.3 \\ 0.4 \end{bmatrix} = \begin{bmatrix} 0.4048 \\ 0.4024 \\ 0.7072 \end{bmatrix}$$

softmax

$$z_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}, z_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}, z_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$= 0.298 \quad = 0.297 \quad = 0.403$$

class of the sample is 3.

$$b) \hat{y} = [0.298 \quad 0.297 \quad 0.403]$$

$$L = -y_1 \log(\hat{y}_1) - y_2 \log(\hat{y}_2) - y_3 \log(\hat{y}_3)$$

$$= -0.298 \log(0.298)$$

$$= -1 \log(0.298) = 0.525 \text{ Ans}$$

Backpropagation learning - is an algorithm for supervised learning of artificial neural networks using gradient descent.

- Given an artificial neural network and an error function the methods calculate the gradient of the error function with respect to neural network, weights, and biases.

- Backpropagation is an algorithm to train a neural network and training a neural network means for a given data it finds the optimum values of weights and biases.

for i in range(epochs):

 for j in range(x.shape[0]):

 select row1 (random)

 predict using forward propagation

 calculate loss (using a suitable loss function)

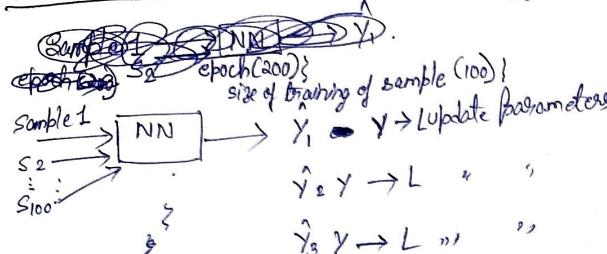
 update weights and bias using GD (Gradient descent)

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

Calculate the average loss.

so, no. of updates = no. of epochs (Batch or D)

Stochastic Gradient Descent



∴ In 1 epoch \rightarrow 100 updates happens in (SGD)
 So, in 200 epochs, $200 \times 100 = 20,000$ updates

Batch GD	Stochastic GD
i) Computational speed is high.	i) Computational convergence speed is high.
ii) Convergence speed is low.	ii) Computational speed is high.
iii) Stable	iii) Noise and unstable.

Mini Batch Gradient Descent -

100 training samples
into no. of mini batches
batch size = 20
no. of batches = 5

$\therefore 1 \text{ epoch} = 5 \text{ updates}$

$$200 \text{ epochs} = 200 \times 5 = 1000 \text{ updates}$$

overfitting data - for training data the loss is decreasing.

for test data the loss is increasing.

* Training accuracy is high, but test accuracy is low.

Solutions to combat overfitting

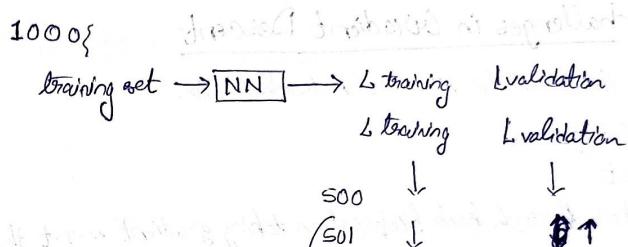
- i) Early stopping
 - ii) Dropout
 - iii) Regularization
 - iv) Data augmentation

Early stopping -

Before testing the model splitting ratio is -

Training set = 50% validation set = 20%

Test set = 30%.



→ Comes out of loop when loss of training decreasing & loss of validation increasing.
(i.e) Early stopping.

* No. of hidden layer increases with increasing complexity of the model.

No. of neurons

$$n_h = \frac{n_{\text{inlet}} + n_{\text{out}}}{2} \quad (\text{linear})$$

↓
no. of hidden layer

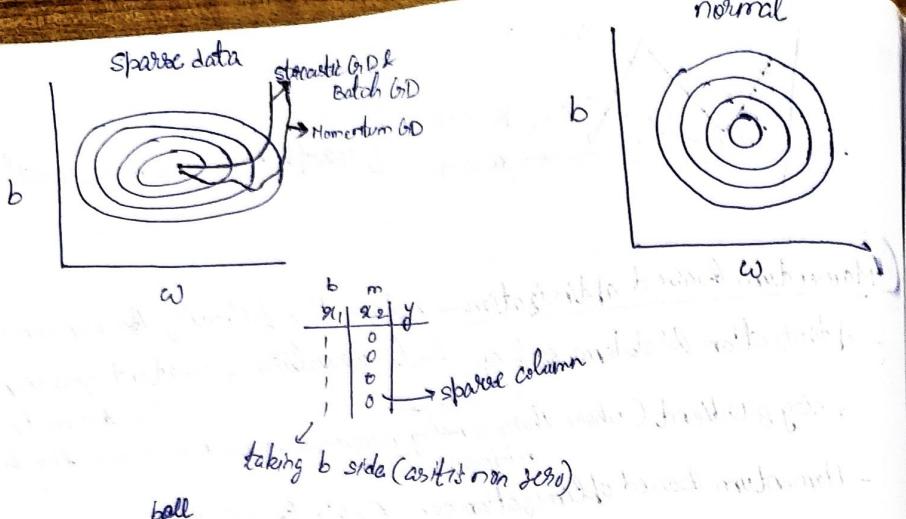
* hidden layer activation function \rightarrow sigmoid / tanh / ReLU

Vanishing gradient problem

If we have a lot of hidden layer (e.g.: 30, 40) it is known as deep neural network.

•

- $0.1 \times 0.1 \times 0.1 \times 0.1 = 0.0001$ ~~so~~ (vanishing)
 - ReLU activation funⁿ is used to avoid vanishing gradient problem.
 - If η is high it will keep oscillating in a the range



- In elongated ~~ball~~ ~~loss~~ problem the slope changes in one axis but is fixed in one axis, but is fixed in another axis. If the data is sparse the graph plotted for ~~loss~~ v.r.t. parameters will be elongated otherwise it is circular.
- In elongated graph the slope changes in 1 axis and ~~is fixed in other axis~~ other axis. In this case our normal optimization algo doesn't solve well, to solve this problem the adagrad algo is used.

$$\begin{aligned}
 L &= (y - \hat{y}) \\
 \text{updated weight} &w_t = w_{t-1} - \eta \frac{\partial L}{\partial w_{t-1}} \\
 \text{updated bias} &b_t = b_{t-1} - \eta \frac{\partial L}{\partial b_{t-1}}
 \end{aligned}$$

- To overcome this disadvantage adagrad is being introduced.
- In adagrad we set different learning rates for different parameters.

~~$w_t = w_{t-1} - \eta \frac{\partial L}{\partial w_{t-1}}$~~

~~$w_t = w_{t-1} - \eta \frac{\partial L}{\partial w_{t-1}} + \frac{(\nabla w_t)^2}{V_t + \epsilon}$~~

$w_{t+1} = w_t - \eta \frac{\nabla w_t}{\sqrt{V_t + \epsilon}}$ → it will become small if it's very large
 $b_{t+1} = b_t - \eta \frac{\nabla b_t}{\sqrt{V_t + \epsilon}}$ → as it's not sparse its value will increase so, V_t will also increase.
 $V_t = V_{t-1} + (\nabla w_t)^2$
 ↓ current velocity old velocity $(\nabla w_t)^2$
 $V_t = V_{t-1} + (\nabla b_t)^2$

- We generally don't use adagrad in NN, we can use adagrad ~~in~~ linear regression but not for complex neural network.
- Adagrad able to reach near to the solution but doesn't converge to the solution the reason is that we are decreasing the learning rate.

for, $w \rightarrow$ it is increasing the learning rate
 $b \rightarrow$ it is decreasing the learning rate

- Learning rate is divided by $\sqrt{V_t + \epsilon}$ and where epochs increases the learning rate becomes smaller due to which the updates become small and when no update happens it stops at a point and can't reach the minima or we can't reach to the global minima.

- By calculating ~~last~~ exponentially weighted moving average we are giving more importance to the current data than the previous data.
- So, we have to focus on 2 points in EWMA, the first point is the
 - i) In due time the weightage of the data will reduce,
 - ii) It's a technique to find patterns or sudden trends in time series data.

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

where, θ_t = current value (loss of gradient)

V_t = EWMA at time t .

V_{t-1} = EWMA at time $t-1$.

β = constant ($0 < \beta < 1$ (smoothing parameter)).

	task
D ₁	25
D ₂	13
D ₃	17
D ₄	23
D ₅	5

$$\left. \begin{array}{l} V_0 = 0 \\ V_1 = 0.99 \times 25 + 0.1 \times 13 = 25 \\ V_2 = 0.99 \times 25 + 0.1 \times 17 \\ V_3 = 0.99 \times 23 + 0.1 \times 17 \end{array} \right\}$$

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

$$V_0 = 0$$

$$V_1 = \beta V_0 + (1-\beta) \theta_1 = (1-\beta) \theta_1$$

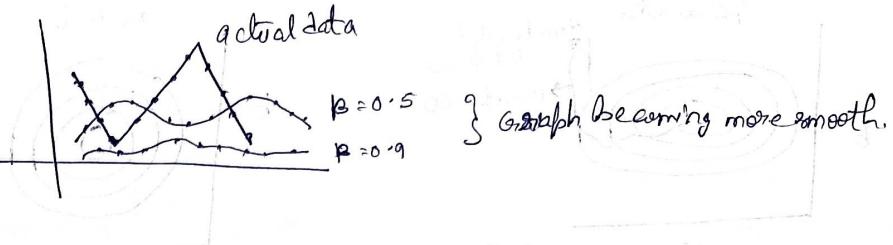
$$V_2 = \beta V_1 + (1-\beta) \theta_2 = \beta (1-\beta) \theta_1 + (1-\beta) \theta_2$$

$$\begin{aligned} V_3 &= \beta V_2 + (1-\beta) \theta_3 = \beta [\beta (1-\beta) \theta_1 + (1-\beta) \theta_2] + (1-\beta) \theta_3 \\ &= \beta^2 (1-\beta) \theta_1 + \beta (1-\beta) \theta_2 + (1-\beta) \theta_3 \end{aligned}$$

$$V_4 = \beta V_3 + (1-\beta) \theta_4$$

$$= \beta [\beta^2 (1-\beta) \theta_1 + \beta (1-\beta) \theta_2 + (1-\beta) \theta_3] + (1-\beta) \theta_4$$

$$= \beta^3 (1-\beta) \theta_1 + \beta^2 (1-\beta) \theta_2 + \beta (1-\beta) \theta_3 + (1-\beta) \theta_4$$



Graph becoming more smooth.

Momentum based optimization - solves the following non convex

- optimization problems such as high curvature, constant gradients & local noisy gradient (where there is noisy minima & we can't overcome the local minima).

- Momentum based optimization can navigate all these 3 problems

$$W_{t+1} = W_t + V_t \rightarrow V_t = \beta V_{t-1} - \eta \frac{\partial L}{\partial W_t}$$

↓ momentum

Adagrad optimization

Adaptive gradient optimization -

- In this optimization technique we don't fix the learning rate rather we change it or we adopt the learning rate acc. to the situation.

- The cases where adagrad works well in comparison to other optimization techniques

Case I: The input feature scales are different.

X ₁	X ₂	Y = 10
cgi pa	package	
10	1000	K-L

Case II: If the data is sparse then we can use adagrad optimization b/c

the momentum and SGD can't sparse the data

cgpa	iq	IT	package
1	1	1	0

The question is if there is sparsity in data what will happen.

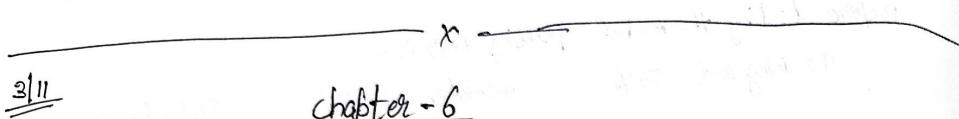
$$w_{ii\text{ new}} = w_{ii\text{ old}} - \eta \frac{\partial L}{\partial w_{ii}}$$

→ minimize the loss

(derivation of sigmoid activation function, tanh activation function).

$$\frac{\partial L}{\partial w_{ii}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z_i} \times \frac{\partial z_i}{\partial w_{ii}} \times \frac{\partial w_{ii}}{\partial w_{ii}}$$

- To solve this we need to use another activation function ReLU.



chapter - 6

challenges in Gradient Descent

- Step 1) Randomly initialize the parameter of neural network (w, b).
- Step 2) Perform forward propagation.
- Step 3) Calculate the cost.
- Step 4) Update the parameters through back propagation taking gradient w.r.t the parameters.

Challenges in GD -

The problems in the current optimizers are

i) Deciding the correct value for the learning rate.

(η value will be same across all epochs).

$$w_t = w_{t-1} - \eta \frac{\partial L}{\partial w_t}$$

ii) So, to overcome this particular problem, we can use learning rate scheduling.

- Can't be applied to all neural networks.

iii) We have to find the optimal weight and bias (for several parameters). If there are 2 parameters w_1 & w_2 for both we can't keep separate learning rates in the existing gradient descent.

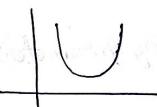
iv) local minima problem - In non convex surface we have multiple minima but our solution is the global minima, while at the time of convergence mostly we stuck in local minima.

v) Saddle point - It is a point on the surface of a function where slopes are $\neq 0$ in all orthogonal directions but it is neither a local minima nor a local maxima.

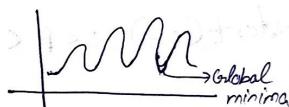
- To remove this problem we can use some other optimizers -

- i) Momentum based optimizer
- ii) Adagrad
- iii) RMS prop
- iv) Adam

Convex vs Non Convex



convex



non convex.



contour plot

is a 2D curve where the loss function is a function of single parameter.



surface at point A if we set to know somehow that my global minima is at lower direction.

Then at point B, when I got to know confirmly that my global minima is there then I should take less step means have to increase the learning rate but this is not possible in our all three Gradient Descents.

(Batch, stochastic, miniBatch). Based GD.

Therefore, η should be adaptive.

Exponentially weighted moving average -

EWMA is a technique based on which we are capable of finding trends in a time series based data. It is basically used in time series forecasting, financial forecasting, signal processing and in deep learning to build cool optimizers.