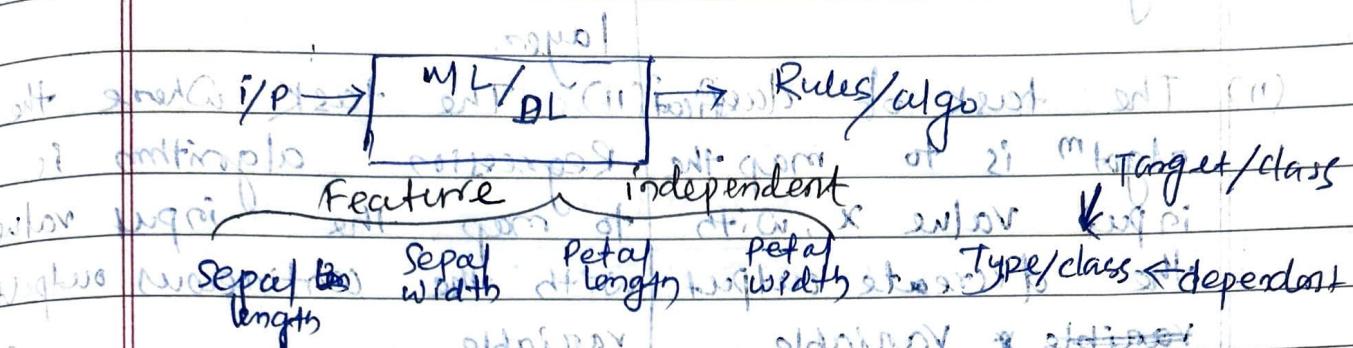
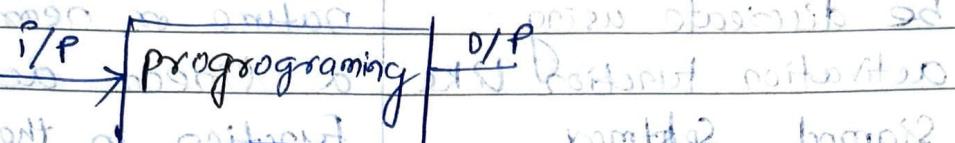
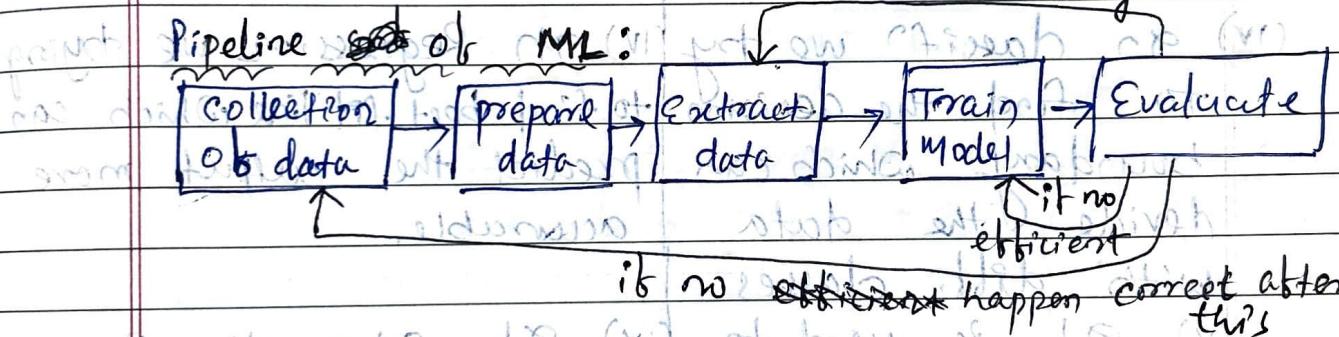


Text Book - Nikhil Buduma, "Fundamentals of Deep learning : Designing next-generation Machine Learning algorithms", O'Reilly, 2017



Sample / Observation



Datatypes → Structured data vs unstructured data

Types of ML → unsupervised learning, supervised learning, reinforcement learning

(a) unsupervised Learning → clustering

→ Dimensionality Reduction

(b) Reinforcement learning →

→ Agent need a environment to make a decision.

(c) supervised learning →

→ classification

→ Regression

Dift. b/w classification and Regression.

Classification

(i) In classification, the output variable must be discrete using activation function like Sigmoid, Softmax.

(ii) The task of classification is to map the input value x , with the discrete output variable.

(iii) Classification algorithms use discrete data.

(iv) In classification we try to find the decision boundary which can divide the data with diff. classes.

(v) It is used to solve the classification problem such as identification of spam emails, speech recognition.

(vi) Classification algorithm can be divided into binary and multiclass classifier.

Regression

(i) In Regression the output variable must be continuous nature or near value using a linear activation function in the output layer.

(ii) The task where the algorithm is to map the regression algorithm is to map the input value with the continuous output variable.

(iii) Regression algorithms use continuous data.

(iv) In Regression we trying to find best fit line which can predict the outputs more accurate.

(v) It is used to solve the regression problem such as weather prediction, house price pred' etc.

(vi) Regression algorithm further divided into linear and non-linear Regression.

Metrics for classification:

True Positive : Yes
True Negative : NO
False Positive : NO
False Negative : YES

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

Accuracy = $\frac{\text{TP} + \text{TN}}{\text{Total no. of input sample}}$

Sensitivity = $\frac{\text{TP}}{\text{TP} + \text{FN}}$

Specificity = $\frac{\text{TN}}{\text{TN} + \text{FP}}$

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

Recall / Sensitivity = $\frac{\text{TP}}{\text{TP} + \text{FN}}$

F1-Score = $\frac{2 \cdot [\text{Precision} \times \text{Recall}]}{\text{Precision} + \text{Recall}}$

(Harmonic Mean between precision and recall)

Confusion matrix :

Actual Class	+ve	TP	FN
-ve		FP	TN
	+ve		-ve

predicted class

Specificity / TNR = $\frac{\text{TN}}{\text{TN} + \text{FP}}$

FPR = $\frac{\text{FP}}{\text{FP} + \text{TN}} = 1 - \text{TNR}$

ROC curve (Recall vs operating characteristics) ?
Plot using TPR vs FPR

	sepal length	sepal width	prediction	target
1	0.192	0.9085	Setosa	Setosa
2	1.132	2.1726	virginica	virginicolor
3	-0.958	2.1738	Setosa	Setosa
4	-2.952	2.138	virginica	virginicolor
5	-0.5055	-2.1499	virginicolor	virginica
6	0.8187	0.6222	virginica	virginica
7	-0.938	-2.1703	virginicolor	virginicolor
8	0.8777	0.0536	virginica	virginica
9	-4.388	-0.239	virginicolor	Setosa
10	-1.419	-0.6867	Setosa	Setosa

Overall classification : Accuracy = $\frac{6}{10} = 0.6$

For a class Setosa : TP = 3, TN = 6 \Rightarrow $\text{Accuracy} = \frac{3+6}{3+6+0+1} = 0.9$

precision = $\frac{TP}{TP+FP} = \frac{3}{3+0} = 1$

Recall = $\frac{TP}{TP+FN} = \frac{3}{3+2} = 0.7$

F1-score = $2 \times \left[\frac{1 \times 0.7}{1 + 0.7} \right] = 2 \times \frac{3}{4} = \frac{3}{2}$

$= 2 \times \frac{\frac{3}{4}}{\frac{4+3}{4}} = 2 \times \frac{3}{7} \times \frac{4}{7} = \frac{6}{7}$

For a class virginicolor = $\frac{3}{3+2} = 0.6$

TP = 1, TN = 5, FP = 2, FN = 2 \Rightarrow $\text{Accuracy} = \frac{5+1}{5+1+2+2} = 0.5$

For a class virginica = $\frac{2}{2+1} = 0.6667$

TP = 2, TN = 5, FP = 1, FN = 1 \Rightarrow $\text{Accuracy} = \frac{2+5}{2+5+1+1} = 0.75$

versicolor : Accuracy = $\frac{3+1}{3+1+2+2} = 0.5$

f1 score = $2 \times \frac{1}{\frac{1}{3} + \frac{1}{2}} = \frac{4}{5}$

Precision = $\frac{1}{3} = \frac{1}{3+0+1+2} = \frac{1}{6}$

Recall = $\frac{1}{3} = \frac{1}{3+2} = \frac{1}{5}$

virginica : Accuracy = $\frac{7+2}{7+2+2+2} = 0.75$

f1-score = $2 \times \frac{2}{\frac{2}{3} + \frac{2}{3}} = 2$

Precision = $\frac{2}{3} = \frac{2}{2+0+1+1} = \frac{2}{6}$

Recall = $\frac{2}{3} = \frac{2}{2+2} = \frac{2}{4} = 0.5$

Regression : (for metrices)

(i) Mean absolute error (MAE) = $\frac{1}{n} \sum_{i=1}^n |y_{\text{actual}} - y_{\text{predicted}}|$

(ii) Mean squared error (MSE) = $\frac{1}{n} \sum_{i=1}^n (y_{\text{actual}} - y_{\text{predicted}})^2$

(iii) Root mean square error (RMSE) = $\sqrt{\frac{1}{n} \sum_{i=1}^n (y_{\text{actual}} - y_{\text{predicted}})^2}$

- You need to predict how close you come to the prediction, for this use MAE
- for DL and ML, for large data set use RMSE.
- for standard error/less error, use RMSE.

$$A = 10 \quad S = 97$$

Q. Actual Value = (3, 5, 7)

Predicted Value = (2, 5, 8)

$$MAE = \frac{|3-2| + |5-5| + |7-8|}{3} = \frac{1+0+1}{3} = \frac{2}{3}$$

$$MSE = \frac{(3-2)^2 + (5-5)^2 + (7-8)^2}{3} = \frac{2+0+1}{3} = \frac{3}{3} = 1$$

$$RMSE = \sqrt{\frac{2+0+1}{3}} = \sqrt{\frac{3}{3}} = \sqrt{1} = 1$$

Actual Value = (5, 6, 9)

Predicted Value = (3, 5, 7)

$$MAE = \frac{|5-3| + |6-5| + |9-7|}{3} = \frac{2+1+2}{3} = \frac{5}{3}$$

$$MSE = \frac{(5-3)^2 + (6-5)^2 + (9-7)^2}{3} = \frac{4+1+4}{3} = \frac{9}{3} = 3$$

$$RMSE = \sqrt{MSE} = \sqrt{3} = \sqrt{3}$$

logistic loss function (cross entropy): * entropy ↑ for always 2 value, $y \in [0, 1], p = P(y=1)$

$$\text{logloss}_{N=2} = -(y \log(p) + (1-y) \log(1-p))$$

$$\text{logloss} = -\frac{1}{N} \sum_{i=1}^N (y_i \log(p_i) + (1-y_i) \log(1-p_i))$$

It's use of Multiclass, in order to capture the model confidence.

Advantage of ML:

→ Automation of repetitive task

→ Improved Decision making

→ Pattern Recognition

→ Continuous Improvement

→ Scalability

→ Cost efficiency

→ Innovation enablement

Disadvantage of ML:

→ Data dependency

→ High computational cost

→ Complexity and interpretability

→ Risk of bias

→ Job Replacement

Q. Find the null set of matrix A = $\begin{bmatrix} 2 & 3 \\ 6 & 9 \end{bmatrix}$

Q. find the basis of the column space of the matrix

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

Q. find the eigen values and corresponding eigen vector of the matrix, B = $\begin{bmatrix} 1 & 2 \\ 5 & 4 \end{bmatrix}$

Q. Explain the significance of the eigenvalues and eigen-vectors in the data sequence application?

Q. Find the eigenvalues and corresponding eigen vectors of the matrix $A = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 2 & 3 \\ 0 & 0 & 3 \end{bmatrix}$ and

(a) $C_1 = \begin{bmatrix} 1 & 1 & -2 \\ -1 & 2 & 1 \\ 0 & 1 & -1 \end{bmatrix}$

and $C_2 = \begin{bmatrix} 1 & 1 & -2 \\ -1 & 2 & 1 \\ 0 & 1 & -1 \end{bmatrix}$

(c) $\begin{bmatrix} 1 & 2 & 2 \\ 0 & 2 & 1 \\ -1 & 2 & 2 \end{bmatrix}$ to establish

Answer:

Given matrix = $\begin{bmatrix} 2 & 3 \\ 6 & 9 \end{bmatrix}$ with data 2
format in 2x2 matrix

take, $Ax = 0$ for solution

$$\begin{bmatrix} 2 & 3 \\ 6 & 9 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 2x_1 + 3x_2 \\ 6x_1 + 9x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\begin{aligned} 2x_1 + 3x_2 &= 0 \times 3 \\ 6x_1 + 9x_2 &= 0 \times 3 \end{aligned}$$

$$2x_1 + 3x_2 = 0$$

$$2x_1 = -3x_2 \Rightarrow \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -3 \\ 2 \end{bmatrix}$$

$$\Rightarrow x_1 = -\frac{3}{2}x_2$$

(3)

Given matrix = $\begin{bmatrix} 1 & 2 \\ 5 & 4 \end{bmatrix}$ for rot

we know $\Rightarrow 0 \in \sigma(A) \Rightarrow \lambda = 0$

$$\begin{bmatrix} 1 & 2 \\ 5 & 4 \end{bmatrix} - \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 2 \\ 4 & 3 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 \\ 5 & 4 \end{bmatrix} - \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} -1 & 2 \\ 3 & 3 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 \\ 5 & 4 \end{bmatrix} - \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 2 \\ 4 & 3 \end{bmatrix}$$

$$\text{L.C.} \Rightarrow (1-\lambda)(4-\lambda) - 12 = 0$$

$$\Rightarrow 4\lambda^2 - 5\lambda + 10 = 0$$

$$\text{To factorize } \Rightarrow \lambda^2 - 5\lambda + 6 = 0$$

$$\Rightarrow \lambda^2 - (6-1)\lambda + 6 = 0$$

$$\text{To solve eqn } \Rightarrow \lambda^2 - 6\lambda + 6 = 0$$

$$\Rightarrow (\lambda - 6)(\lambda + 1) = 0$$

$$\Rightarrow \lambda = 6, \lambda = -1$$

For $\lambda = 6$, put in eq? (P)

$$\begin{bmatrix} 1 & 2 \\ 5 & 4 \end{bmatrix} - \begin{bmatrix} 6 & 0 \\ 0 & 6 \end{bmatrix} = \begin{bmatrix} -5 & 2 \\ 5 & -2 \end{bmatrix} = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

$$(A - 6I)x = 0$$

$$\Rightarrow \begin{bmatrix} -5 & 2 \\ 5 & -2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} -5 & 2 \\ 5 & -2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0$$

$$\Rightarrow -5x_1 + 2x_2 = 0 \Rightarrow -5x_1 = -2x_2$$

$$\Rightarrow 5x_1 - 2x_2 = 0 \Rightarrow \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2 \\ 5 \end{bmatrix}$$

For $\lambda \geq 1$ (eigenvalue condition)

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

$$\Rightarrow \begin{bmatrix} 1 & 2 \\ 5 & 4 \end{bmatrix} - \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \lambda = 0 \Leftrightarrow \text{Want } \lambda$$

$$\Rightarrow \begin{bmatrix} 2 & 2 \\ 5 & 3 \end{bmatrix} - \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \lambda = 0$$

$$\Rightarrow 2x_1 + 2x_2 = 0$$

$$\Rightarrow \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

(4) → Eigenvalues represent the magnitude of variance in the data direction.

→ Eigenvector represent the direction of maximum variance.

→ In PCA eigenvectors are principal components, eigenvalues (importances) of those components.

(5) Given matrix, $A = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 2 & 3 \\ 0 & 0 & 3 \end{bmatrix}$

We know, $|A - \lambda I| = 0$.

$$\Rightarrow \begin{vmatrix} 1 & 2 & 3 \\ 0 & 2 & 3 \\ 0 & 0 & 3 \end{vmatrix} - \begin{vmatrix} \lambda & 0 & 0 \\ 0 & \lambda & 0 \\ 0 & 0 & \lambda \end{vmatrix} = 0$$

$$\Rightarrow \begin{vmatrix} 1-\lambda & 2 & 3 \\ 0 & 2-\lambda & 3 \\ 0 & 0 & 3-\lambda \end{vmatrix} = 0$$

$$\Rightarrow (1-\lambda)(2-\lambda)(3-\lambda) = 0 \Rightarrow \lambda_1 = 1, \lambda_2 = 2, \lambda_3 = 3$$

$$(1-\lambda)(3-\lambda)(2-\lambda) = 0 \Rightarrow \lambda_1 = 1, \lambda_2 = 2, \lambda_3 = 3$$

$$(\lambda-1)(\lambda-2)(\lambda-3) = 0 \Rightarrow \lambda_1 = 1, \lambda_2 = 2, \lambda_3 = 3$$

$$(\lambda-1)(\lambda-2)(\lambda-3) = 0 \Rightarrow \lambda_1 = 1, \lambda_2 = 2, \lambda_3 = 3$$

$$\Rightarrow \lambda^2 - (2+1)\lambda + 2 = 0$$

$$\Rightarrow \lambda^2 - 2\lambda - 1 = 0$$

$$(1-\lambda)(2-\lambda) = 0 \Rightarrow \lambda_1 = 1, \lambda_2 = 2$$

$$(\lambda-1)(\lambda-2) = 0 \Rightarrow \lambda_1 = 1, \lambda_2 = 2$$

$$\text{For } \lambda = 1, \lambda = 1$$

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

$$\Rightarrow 0x_1 + 0x_2 + 2x_3 = 0 \Rightarrow 2x_3 = 0 \Rightarrow x_3 = 0$$

$$(1-\lambda)(2-\lambda)(3-\lambda) = 0 \Rightarrow \lambda_1 = 1, \lambda_2 = 2, \lambda_3 = 3$$

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ -2 \end{bmatrix}$$

Q: A fair coin has probability $P(H) = 0.5$ and $P(T) = 0.5$, find its entropy $H(x)$.

Ans: For a fair coin $H(x) = -(0.5 \log_2(0.5) + 0.5 \log_2(0.5))$

(entropy)

For single value

For a biased coin,

$$P(H) = 0.8, P(T) = 0.2$$

$$H(x) = -(0.8 \log_2 0.8 + 0.2 \log_2 0.2)$$

Q: For two distribution how to find cross entropy?

$$\text{P} = (0.7, 0.3), \text{Q} = (0.6, 0.4)$$

$$H(P, Q) = -[(0.7 \log_2(0.6)) + (0.3 \log_2(0.4))]$$

(cross entropy)

~~Why 2 values? 2 values~~

$$KL\text{-Divergence } D_{KL}(P||Q) = H(P, Q) - H(P)$$

$$H(P) = -(0.7 \log_2(0.7) + 0.3 \log_2(0.3))$$

$$H(P, Q) - H(P) = -0.1$$

Binary cross Entropy:

Single example, $y=1, \hat{y}=0.8$

$$H(P) = -(y \log_2 \hat{y} + (1-y) \log_2(1-\hat{y}))$$

Batch of 3: $y = (1, 0, 1), \hat{y} = (0.9, 0.25, 0.6)$

$$H_1(P) = -\log_2 0.9 + 0$$

$$H_2(P) = 0 - \log_2(0.75)$$

$$H_3(P) = +\log_2(0.6)$$

$$H(P) = H_1(P) + H_2(P) + H_3(P)$$

$$\frac{1}{3} [-(1 \cdot \log_2 0.9 + 0 + 0 - \log_2 0.75 + 0 + 1 \cdot \log_2 0.6)] = -0.14$$

Cost function is

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

$$((0.9 - 1)^2 + (0.25 - 0)^2 + (0.6 - 1)^2) = 0.44$$

Suppose we have a classification problem with 3 classes x, y, z , given that for the probability distribution P ,

$$P(x) = 0.7, P(y) = 0.2, P(z) = 0.1$$

Find the cross entropy of given probability distribution.

Suppose you have two probability distribution function P and Q , also $P(x) = 0.8, P(y) = 0.2, P(z) = 0.1$, $Q(x) = 0.6, Q(y) = 0.3, Q(z) = 0.1$. Find the cross entropy $H(P||Q)$. Answer: 0.33

$$H(P||Q) = -[(0.8 \log_2(0.6) + 0.2 \log_2(0.3) + 0.1 \log_2(0.1))]$$

$$= -1.156$$

$$H(P, Q) = -[(0.7 \log_2(0.6) + 0.2 \log_2(0.3) + 0.1 \log_2(0.1))]$$

$$= -1.189$$

Definition of Deep learning:

Deep learning is the type of machine learning that uses Artificial Neural Network that multiple layers forward to learn and make decision.

It replicates just like a human brain, as all the neural network are connected in the brain which exactly fits the concept of deep learning.

In Deep NNs that holds various layers of interconnected nodes elaborated with this hierarchical pattern and data feature.

→ An automated learning system, that enhances the data without the intervention of manual feature.

→ Success in various fields such as image recognition, natural language processing, speech recognition and recommendation system.

Dift. b/w ML and DL in
Machine learning

Deep learning

→ ML is a study that uses statistical methods enabling the machine to improve (without imitation) the functionalities just like a human brain.

→ ML is the subset of AI

→ ML allows the system to learn from the data and provides the outputs accordingly.

→ The aim is to increase accuracy to carrying much about the success. Ratio of trained with longer amount of data, do nothing dominant with few.

→ DL is the deep neural networks to analyze the data.

→ DL attempts the highest rank in terms of accuracy when it is trained with long amount of data.

→ If you have the clear & clear about the logic involved in problem and can visualize the complex functionalities, then it defines the most aspect.

→ Types of ML → SL, RL, VSL, and Reinforcement learning.

→ Some is less effective than deep learning. ML as it can easily fail it cannot work for larger sets of data.

→ ML is a subset of higher amount of data.

→ ML is a subset of AI, that focuses on developing algorithm that can learn from the data and improve the performance over time without being explicitly programmed.

→ Ex: Virtual personal Assistant, Siri, Alexa, Google etc., email spam.

→ DL you have clear idea about the logic involved in it, but I don't have idea about the feature so you break the complex functionalities into low dimensional features by adding more layers than it defines the DL aspects.

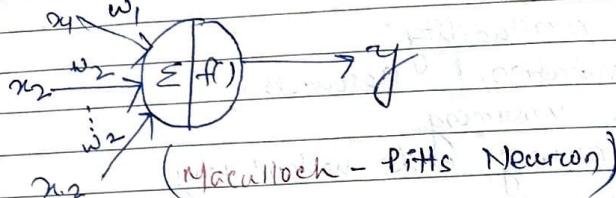
→ Types of DL → unsupervised learning, Recurrent NN and RNN.

→ DL is powerful than deep learning. ML as it can easily work for smaller sets of data.

→ DL is the subset of ML that focuses on developing DNN, that can automatically learn extracted features from the data.

Ex: Sentiment based news aggregation, image analysis, caption generation etc.

ANN:



Maculloch - Pitts Neuron works as follows:-
 The inputs are multiplied by the corresponding weights, the weighted sum of which is then summed up.

→ If the sum value is greater than or equal to threshold value then the neuron fires and the output is 'y', otherwise the neuron does not fire and output is zero.

Characteristics of ANN:-

- A neuron is the mathematical function model on the working of biological neuron.
- It is an elementary unit in an artificial neural network.
- One or more inputs are separately weighted.
- Inputs are summed up and passed through a non-linear function to produce the outputs.
- Each neuron holds one internal state, called activation signal.
- Each connection link is carrying information about the input signal.
- Every neuron is connected to another

neurons to the connection link.

→ Information is transmitted via the connection link.

Differences between biological and Artificial Neural network

- | | |
|--|--|
| <ul style="list-style-type: none"> → Biological NN → Structure is → Dendrites → cell body → Axon → Axon terminal → Memory happens → Learning can handle ambiguous noise input. → The speed of processing information is slow. → Storage allocation to new processes can be done easily by adjusting the interconnection strength. → Massive parallel operation can be conducted simultaneously. → Information can be retrieved from the sub nodes even if it is corrupted. → The information is being transmitted and processed into the network is not monitored. | <ul style="list-style-type: none"> → ANN → Structure is → Input → hidden layer → output → Speed of processing is faster compared to biological NN. → Storage allocation into a new process is not possible as each location is dedicated to a specified process. → Only sequential operations can be conducted. → Once corrupted the information can not be restored. → A control unit monitors all information and activities in the network. |
|--|--|

→ Our brain contains about 86 billion and more than 100 synapses. → The no. of neurons in ~~the~~ cortical neurons is much less than biological neural network.

→ The brain works asynchronously.
→ Brain have complicated topology. BNN can learn completely new task.

→ ANN works synchronously.
→ ANN are often in tree structure.
→ DNN are specialized and they can perform one task.

I/P value, $y = [w_1x_1 + w_2x_2 + \dots + w_nx_n]$

Advantages of ANN

- Pattern Recognition
- Parallel processing
- Adaptability
- Generalization

- Disadvantages of ANN
- Large data and computation
- Overfitting
- Hyperparameter tuning
- Slow convergence
- Time consuming
- Limited data efficiency
- Categorical data handling

Perception is a configuration proposed by Frank Rosenblatt

- 1957
- Q) What is the perception ?
- Q) What are the components of a perception?
- Q) What is meant by the data is linearly separable
- Q) Can perceptron handle non linearly separable data?
- Q) What are the limitations of perceptron?

① A perceptron is the simplest type of artificial neural network used for binary classification, it takes multiple inputs applies weights sums them, passes the result through an activation function and outputs.

② Components of perceptron -
 (i) input layer
 (ii) weight
 (iii) Activation function
 (iv) output

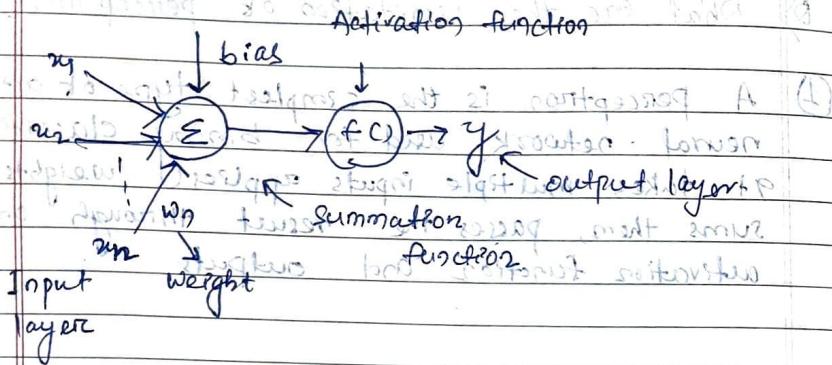
③ Data is linearly separable if a straight line, plane, hyperplane can separate the data points of different classes without error.

④ No, a simple perceptron cannot handle non-linearly separable data. for that

multilayer perceptrons are used.

- (5) Limitations of perception :-
- works only for linearly separable data.
- outputs are binary.
- learning may be slow for large datasets.
- cannot solve complex problem like XOR.

Perception



- Output is produced only in binary mode.
- perception supervised learning, not simple binary classifier.

Training algorithm :-

- If weight value = +ve, & target value = +ve, then it is +ve - reinforcement of targeted value.
- + If weight value = -ve, & target value = +ve, then it is -ve - punishment of targeted value.
- + If weight value = +ve, learned more & much.

Step-1 :-

$$w = x_1 w_1 + x_2 w_2 + \dots + x_n w_n + b$$

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

$$= x^T w + b$$

$$x = [x_1 \ x_2 \ \dots \ x_n]$$

$$w = [w_1 \ w_2 \ \dots \ w_n]$$

$$b = \text{bias}$$

$$= [x_1 \ x_2 \ \dots \ x_n] \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

$$= x^T \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

$$= x^T w$$

$$= x^T w + b$$

$$= y = f(z)$$

$$\Rightarrow \text{step function, if } y = 1, w > 0$$

$$0, w < 0$$

Activation function :-

- (1) Linear activation function

- (2) non-Linear activation function

Binary step Activation function :-

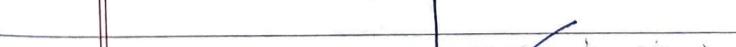
$$f(z) = 1, \text{ if } z > 0$$

$$0, \text{ if } z \leq 0$$



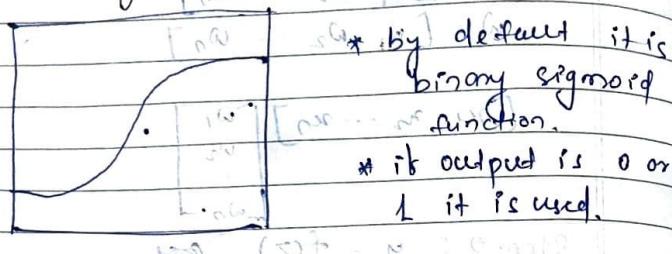
Linear activation function :-

$$f(z) = \sum_{i=1}^n -\infty \text{ to } +\infty$$



pros :-

- It gives + other than 0 or 1 to neuron.
 - It can't connect to multi neuron.
 - changes are const. for error correct.
- (3) No-1 linear activation function :-
- (a) Sigmoid or Logistic activation function :-



→ Bipolar sigmoid AF :-

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

* standard derivative form

→ Binary sigmoid AF :-

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

* -2 to 2 range

(4) Tan h AF :-

$$= \frac{e^z - e^{-z}}{1+e^{-2z}} - 1$$

$$= 2 \left[\frac{1}{1+e^{-2z}} \right] - 1 = 2 \times \text{sigmoid}(2z) - 1$$

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

$$\rightarrow -1 to 1 = \text{range}$$

(5)

ReLU function :- (Rectify Linear Unit)

$$y = f(z) = \max(0, z)$$

→ It ranges from 0 to ∞ [0, ∞)

non-linearizing nature

less computational expensive

(6)

Heaviside function :-

$$\text{heaviside function} = \begin{cases} 0, & \text{if } z < 0 \\ 1, & \text{if } z \geq 0 \end{cases}$$

(7)

Softmax function :-
 → It is used in multi-class classification function.

$$\sigma(z_i) = \frac{e^{z_i}}{\sum e^{z_j}}$$

z_i = output vector
 z_j = output vector

→ If output is $\sum e^{z_j}$ then softmax is used
 If output is e^{z_j} then softmax is not used
 If output is $e^{z_j}/\sum e^{z_j}$ then softmax is used

* If output is binary, then sigmoid AF is used

If output is multi-class, then softmax AF is used

→ Multiclass \rightarrow Branch \rightarrow CSC, ME, Engineering

→ Multilevel \rightarrow CSC-A, CSC-B, CSC-C

(1) $0.2 \times 0.3 \rightarrow 0.3$
 $0.6 \rightarrow 0.4 \rightarrow 0.3$ \rightarrow y value

→ Obtained net output for the network using AF, Binary SF, Bipolar SF, Tan-h, ReLU

→ Activation function \rightarrow $y = f(z)$

where $z = x_1 w_1 + x_2 w_2 + b_1$ and $f(z) =$

$$= 0.2 \times 0.3 + 0.6 \times 0.4 + 0.1 = 0.38$$

$$= 0.06 + 0.24 = 0.38$$

Binary SF : $y = f(z) = \frac{1}{1 + e^{-z}}$

$$\begin{aligned} & (w_0) \text{ mult } (e^z) = 8 \\ & (w_1) \text{ mult } (-1) \text{ leads to } -0.381 \\ & \text{another mult } (+0.68) \\ & \text{sum} = \frac{8.21}{1 + 0.68} = 0.595 \end{aligned}$$

Bipolar SF : f_y

$f_y = \frac{1 - e^{-|z|}}{1 + e^{-|z|}}$ constant sigmoid

Linear neuron :-

- 4 linear neurons is a neuron network that uses only linear transformation in its layers.
- It is a linear function: $y = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n$
- It can't learn complex relationships between the input and output vector.
- If a neural network uses a linear activation function all network after the network converges into one.
- It can't be used the backpropagation technique, bcoz the derivative of the linear activation function is a constant.
- It can only classify the objects into linearly separable objects.
- It can't use the hidden layer.

* What is perceptron function?

→ It is a function that map its input 'x' which is multiplied with the learned weight co-efficient and the output value $f(z)$ is generated.

$$f(z) = \begin{cases} 1, & \text{if } w_0 + w_1 x_1 + \dots + w_n x_n > 0 \\ 0, & \text{otherwise} \end{cases}$$

bias - that element that adjust the boundary away from the origin without any dependencies on the input value.

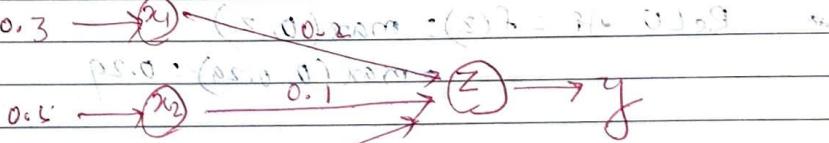
x = vector of input value

bias = That element that adjust the boundary away from the origin without any dependencies on the input value.

Characteristics of perceptron model :-

- It is a ML algorithm for the supervised learning of binary classifier.
- The weight coefficient is automatically learned.
- Initially weight w and input features are multiplied in bipolar mode.
- The decision is made whether the neuron should be dismissed or not.
- That activation function in any perceptron algorithm examples employees are rule to determine whether the weight function value is higher than zero.
- If the sum of all the input values is higher than the threshold value then it should have an output signal else no output is displayed.
- The linear decision boundary is plotted and it separates the distinction between the 2 linearly separable classes +1 and -1.

(1)



$$0.6 \rightarrow x_3 \xrightarrow{w_3} 0.3 \xrightarrow{f(z)} 1.0$$

Obtained the outputs of the network using binary sigmoid AF & Bipolar sigmoid AF.

Ans:

$$z = \sum x_i w_i + b$$

$$= 0.3 \times 0.2 + 0.5 \times 0.1 + 0.6 \times 0.3$$

$$= (0.06 + 0.05 + 0.18) = 0.29$$

$$= 0.29$$

Binary Sigmoid AF = $\frac{1}{1+e^{-z}}$

$$\rightarrow \text{Binary Sigmoid AF}(y) = f(z) = \frac{1}{1+e^{-z}}$$

$$= \frac{1}{1+e^{-0.29}} = 0.572$$

* Bipolar Sigmoid AF = $f(z) = \frac{1-e^{-|z|}}{1+e^{-|z|}}$

$$= \frac{1-0.748}{1+0.748} = -0.272$$

* $\tan h$ AF = $f(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

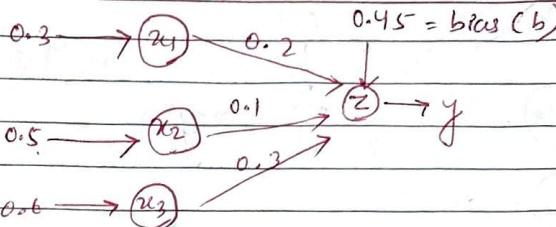
$$= \frac{e^{0.29} - e^{-0.29}}{e^{0.29} + e^{-0.29}} = 0.412$$

$$= \frac{1.336 - 0.748}{1.336 + 0.748} = \frac{0.588}{2.084} = 0.287$$

* ReLU AF = $f(z) = \max(0, z)$

$$= \max(0, 0.29) = 0.29$$

(3Q)



Find same as previous question?

$$z = \sum x_i w_i + b$$

$$= 0.29 + 0.45 = 0.74$$

$$\rightarrow \text{Binary Sigmoid AF}(y) = f(z) = \frac{1}{1+e^{-z}}$$

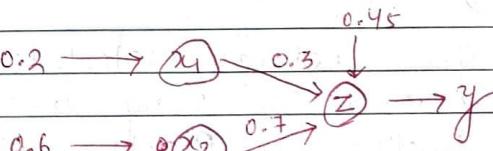
$$= \frac{1}{1+e^{-0.74}} = 0.67$$

\rightarrow Bipolar Sigmoid AF = 0.35

$$\rightarrow \tan h = 0.62$$

$$\rightarrow \text{ReLU} = \max(0, z) = \max(0, 0.74) = 0.74$$

(4Q)



$$z = \sum x_i w_i + b$$

$$= (0.2 \times 0.3 + 0.6 \times 0.7) + 0.45$$

$$= (0.06 + 0.42) + 0.45$$

$$= 0.48 + 0.45 = 0.93$$

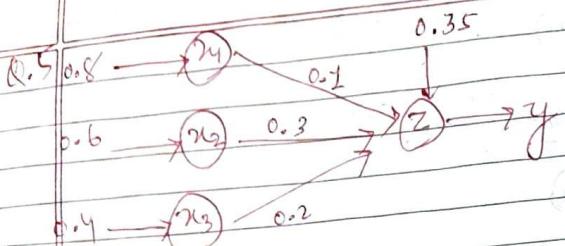
\rightarrow Binary Sigmoid AF = 0.71

\rightarrow Bipolar Sigmoid AF = 0.71

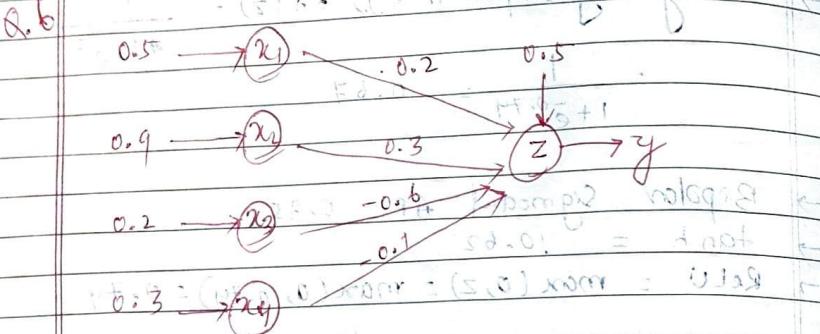
$$\rightarrow \tan h =$$

$$\rightarrow \text{ReLU} = f(z) = \max(0, z) = \max(0, 0.93)$$

$$= 0.93$$



→ obtain the output of the neuron 'y' for the network using binary sigmoidal function, Bipolar sigmoidal function, and tanh.



→ For the given network calculate output of the neuron 'y' using binary sigmoidal function, Bipolar sigmoidal function and tanh.

Q.7

The perceptron is having weights $w = (0.4, -0.3, 0.1)^T$, and a bias is 0. If the input to the perceptron is $x = (0.2, 0.6, 0.25)$ then find the output of the perceptron using binary SF, Bipolar SF and Tanh function.

NOTE: IF output is binary classification, then by default we need to find binary Sigmoidal SF function.

Ans:

$$\Rightarrow [0.4 \ 0.3 \ 0.1] [0.2 \ 0.6 \ 0.25] + 0$$

Let's calculate it step by step. Suppose the weight considered as figure has shown corresponding values together 3 inputs have the following values $w_1 = 0.4, w_2 = 0.3, w_3 = 0.25, b = -1$

and $x_1 = 0.2, x_2 = 0.6, x_3 = 0.25$. Now calculate what will be the output value if you put the configuration for each of the following four input patterns?

(i) If $v = 0$, then $f(v) = 1$, if $v > 0$, then $f(v) = 1$, if $v < 0$, otherwise $f(v) = 0$. Now calculate what will be the output value if you put the configuration for each of the following four input patterns?

(ii) If $v = 1$, then $f(v) = 1$, if $v > 0$, then $f(v) = 1$, if $v < 0$, otherwise $f(v) = 0$. Now calculate what will be the output value if you put the configuration for each of the following four input patterns?

	p_1	p_2	p_3	p_4
x_1	1	0	1	0
x_2	0	1	0	1
x_3	0	1	1	1

Ans: $p_1 = v = \sum x_i w_i + b$, $b = 0$

$$\Rightarrow p_1 = 0.4 \cdot 1 + 0.3 \cdot 0 + 0.1 \cdot 0.25 + (-1) = 0.4 + 0 + 0.025 - 1 = -0.575$$

$$\Rightarrow p_1 = 2, f(v) = 1$$

$$\Rightarrow p_2: x_1 = 0, x_2 = 1, x_3 = 1$$

$$\Rightarrow p_2 = 0.4 \cdot 0 + 0.3 \cdot 1 + 0.1 \cdot 1 + (-1) = -0.4 + 0.3 + 0.1 - 1 = -1$$

$$\Rightarrow p_2 = -1, f(v) = 0$$

$$\Rightarrow p_3 = 3, f(3) = 1$$

$$\Rightarrow p_4 = 1, x_1 = 1, x_2 = 1, x_3 = 1$$

$$\Rightarrow p_4 = -1, f(-1) = 0$$

Feed Forward Neural Network :-

- feed forward neural network is one of the foundations of a neural network architecture particularly in applications where prediction machine learning algorithm face limitation.

- They facilitate the task such as simple classification like face recognition, computer vision and speech recognition through their unidirectional flow of data.

- It consists of input and output layers with optional hidden layers in between.
- FFN networks are efficient for handling noisy data, are relatively straight forward to implement making them versatile tool in various AI applications.

Advantage of FFN:

- Less complex
- Easy to design and maintain
- It is first and speedy, due to one way propagation.
- It is highly responsive to noisy data.
- Disadvantage of FFN:
- It can not be used in deep learning due to the absence of dense layers and back propagation.

(7)

$$w = [0.4, -0.3, 0.1]$$

$$x_2 = [0.2, 0.6, 0.5]$$

$$b_2 = 0$$

$$\text{Q. net} = 0.4 \times 0.2 + -0.3 \times 0.6 + 0.1 \times 0.5 + 0$$

$$\therefore \text{net} = 0.08 - 0.18 + 0.05 = -0.05$$

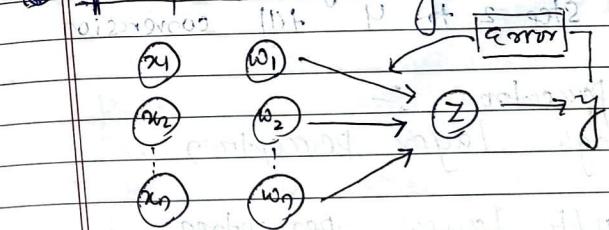
$$f(x) = 0, x = -0.05$$

$$\text{Binary SF} = \frac{1}{1+e^{-x}} = \frac{1}{1+e^{-(-0.05)}} = 0.51$$

$$\text{Bipolar SF} = \frac{1-e^{-x}}{1+e^{-x}} = -0.025$$

$$\tan h(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = 0.0499$$

perception and learning Rule:



Error \rightarrow perception :-

In a perception learning rule the predicted output is compared with actual output, if does not match the error is propagated backward to allow weight adjustment to happen.

$$w_{\text{new}} = w_{\text{old}} + \eta (y_{\text{actual}} - y_{\text{predicted}}) x_i$$

η = learning rate (choose like the minimum iteration is required)

x_i = input value

y_{actual} = output actual

$y_{\text{predicted}}$ = output predicted

- formula :-
- $$b_{\text{new}} = b_{\text{old}} + \eta (Y_{\text{actual}} - Y_{\text{predicted}}) \quad [\because x_0 = 1 \text{ (fixed)}]$$
- perceptron algorithm : $b = x_0 w_0 = w_0$
- steps to find error reduce :-
- ① Initialize the weight and bias value randomly.
 - ② calculate the weighted sum z by the given input value.
 - ③ Applied activation function of the weighted sum and find the predicted output.
 - ④ If $Y_{\text{predicted}}$ is not equal to Y_{actual} then no weight updation is required. else update the weight using formula 1.
 - ⑤ Repeat step 2 to 4 till convergence

- Types of perceptron :-
- Single layer perceptron
 - Multi layer perceptron.

Single layer perceptron :-

SLP can learn only linearly separable objects.

Step-1: Initialize weight and bias.

For simplicity set weight = 0, bias = 0

② Set the learning rate in the range of 0 to 1. ($\eta = [0-1]$)

③ While stopping condition is false do Step 2 to 6.

such that

④ For each training pair, do step 3 to 5.

- ⑤ Set the activation function of the input unit x_i .
- ⑥ calculate the weighted sum value.
- ⑦ compute the response of the output unit based on the activation function.
- ⑧ update the weight and bias if an error occurs in one of these patterns; update the weights and bias using w_{new} and b_{new} ; else the w_0 and b remains same.
- ⑨ Test the stopping condition.

Limitation of SLP :-

→ A single layer SLP can only classify data that is linearly separable which means that only separate a data in a straight line or a hyper plane.

→ It is not suitable for complex ML problem.

→ Ex: It can not go for XOR problem.

→ It is a batch learning process so it can't learn from a new data points incrementally.

→ It is susceptible to local minima which means it can not find the global optimal solution to a problem.

→ It is prone to overfitting and noise bcoz it tries to fit a straight line to a data.

Implementation of OR using perceptron :-

x_1	x_2	y_{actual}
0	0	0
0	1	1
1	0	1
1	1	1

Suppose, $w_1 = 1, w_2 = 1, b = 0.5$

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

case-1: $x_1 = 0, x_2 = 0$ (Actual)

$$z = x_1 w_1 + x_2 w_2 + b = 0 \times 1 + 0 \times 1 + 0.5 = 0.5 > 0$$

$y_{pred} = f(z) = 1 = y_{actual}$ (No weight update)

case-2: $x_1 = 1, x_2 = 0$ (Actual)

$$z = x_1 w_1 + x_2 w_2 + b = 1 \times 1 + 0 \times 1 + 0.5 = 1.5 > 0$$

$$y_{pred} = f(z) = 1 = y_{actual}$$

case-3: $x_1 = 0, x_2 = 1$ (Actual)

$$z = x_1 w_1 + x_2 w_2 + b = 0 \times 1 + 1 \times 1 + 0.5 = 1.5 > 0, f(z) = 1$$

$y_{pred} = f(z) = 1 = y_{actual}$ (no weight update)

case-4: $x_1 = 1, x_2 = 1$ (Actual)

$$z = 1 + 1 - 0.5 = 2 - 0.5 = 1.5 > 0$$

$y_{pred} = f(z) = 1 = y_{actual}$ (no weight update)

Q. Suppose $w_1 = 0.6, w_2 = 0.6, b = 0.25$

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

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

Ans case-1: $x_1 = 0, x_2 = 0$

$$z = x_1 w_1 + x_2 w_2 + b = 0 \times 0.6 + 0 \times 0.6 + 0.25 = 0.25 < 1$$

$y_{pred} = f(z) = 0 = y_{actual}$ (no weight updated)

case-2:

$$x_1 = 0, x_2 = 1$$

$$z = 0.6 + 0 = 0.6 \quad (\text{otherwise})$$

$y_{pred} = f(z) = 0 \neq y_{actual}$

$$\begin{aligned} w_1' &= w_1 + \eta (y_{actual} - y_{pred}) x_1 \\ &= 0.6 + 0.5 (1 - 0) \times 0 = 0.6 \end{aligned}$$

$$\begin{aligned} w_2' &= w_2 + \eta (y_{actual} - y_{pred}) x_2 \\ &= 0.6 + 0.5 (1 - 0) \times 1 \\ &= 0.6 + 0.5 = 1.1 \end{aligned}$$

$$b' = b + \eta (y_{actual} - y_{pred})$$

$$= 0.25 + 0.5 (1 - 0) = 0.5$$

* case-1: $x_1 = 0, x_2 = 0$

$$z = 0 + 0 + 0.5 = 0.5$$

$y_{pred} = f(z) = 0 = y_{actual}$ (no weight update)

case-2: $x_1 = 0, x_2 = 1$

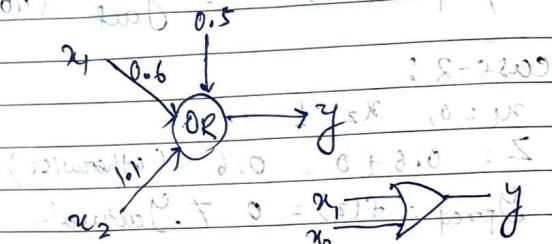
$$z = 0 + 1 \times 0.6 + 0.5 = 1.1 > 1 = f(z) = 1$$

$y_{pred} = f(z) = 1 = y_{actual}$ (no weight update)

case 3: $x_1=1, x_2=0$
 $z = 0.6x_1 + 0 + 0.5 = 1.1 \geq 1, f(z) = 1$
 $y_{pred} = f(z) = 1 = y_{actual}$ (no weight update)

case 4: $x_1=1, x_2=1$
 $z = 0.6 + 1.1 + 0.5 = 2.2 > 1, f(z) \approx 1$

$y_{pred} = f(z) = 1 = y_{actual}$ (no weight update)



How to represent input in a vector format :-

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

$$w^T = \begin{bmatrix} w_0 & w_1 & w_2 \end{bmatrix} = \begin{bmatrix} 0.6 & 0.5 & 1 \end{bmatrix}$$

but $x.w$ not possible so convert to w^T .

$$w^T = \begin{bmatrix} -0.5 \\ 1 \\ 2.0 \end{bmatrix} \quad \text{so } x.w^T = 0 + 0.5 + 2.0 = 2.5$$

$$x^T = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix} \quad w^T = \begin{bmatrix} -0.5 \\ 1 \\ 2.0 \end{bmatrix} \quad \text{so } x^T.w^T = 0 + 0.5 + 2.0 = 2.5$$

$$= \begin{bmatrix} -0.5 \\ 1 \\ 2.0 \\ 2.5 \end{bmatrix} \quad y_{pred}$$

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

$$w^T = \begin{bmatrix} 0.6 \\ 0.5 \\ 1 \end{bmatrix}$$

$$x^T.w^T = \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0.6 \\ 0.5 \\ 1 \end{bmatrix} = \begin{bmatrix} 1.6 \\ 1.5 \\ 1 \end{bmatrix}$$

$y_{pred} = \begin{bmatrix} 1.6 \\ 1.5 \\ 1 \end{bmatrix}$

$$w_1^T = 0.6 + 0.5(1-0) \times 0 = 0.6 + 0 = 0.6$$

$$w_2^T = 0.6 + 0.5(1-0) \times 1 = 1.1 \neq 1$$

$$b^T = 0 + 0.5(1-0) = 0.5 \neq 1$$

new data :- Step-2 (iteration - 2)

$$\begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix} \quad \begin{bmatrix} 0.5 \\ 0.4 \\ 0.5 \\ 1.0 \end{bmatrix}$$

$$y_{pred} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \end{bmatrix} = y_{actual}$$

(Q3)

$w_1 = 0.7$, $w_2 = 0.7$, $\eta = 0.4$, $f(z) = \text{heavy side AF}$
 Implement 'OR' gate using
 perceptron?

Date _____
 Page _____

16 not given

Ans:

$$Y = \begin{bmatrix} b & x_1 & x_2 \end{bmatrix} \cdot \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

$$\therefore w = \begin{bmatrix} 0 \\ 0.7 \\ 0.7 \end{bmatrix}, \text{Yactual} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

$$xw_2 = \begin{bmatrix} 0 \\ 0.7 \\ 0.7 \\ 1.4 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \neq \text{Yactual}$$

$$w_1^T = 0.7 + 0.4(1-0) \times 0 \\ = 0.7$$

$$w_2^T = 0.7 + 0.4(1-0) \times (1-0+0.0) = 1.0 \\ = 0.7 + 0.4 = 1.1$$

$$b^T = 0 + (1-0)0 = 0.0$$

$$X = \begin{bmatrix} -0.4 & 0 & 0 \\ 0.4 & 0 & 1 \\ 0.4 & 1 & 0 \\ 0.4 & 1 & 1 \end{bmatrix}, w = \begin{bmatrix} 1 \\ 0.7 \\ 1.1 \end{bmatrix}$$

$$xw_2 = \begin{bmatrix} 0.4 \\ 1.5 \\ 1.1 \\ 2.2 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \text{Yactual}$$

*

NOTE: If z . Activation function $f(z)$ not given,
 then heavy side AF is used.

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

Q1 What is deep learning? How does
 deep learning differ from ML.

Q2 What ANN? How does biological neurons
 similar to ANN? $f(z) = \text{heavy side AF}$.

Q3 Application of DL.

Q4 What are challenges of DL.

③ Ans: * Image Recognition: Used in face unlock
 and photo tagging. $f(z) = \text{heavy side AF}$.

* Self-driving cars: Helps detect objects, lanes,
 road signs and traffic signs.

* Healthcare: Used to detect diseases from
 X-rays or scans.

* Chatbots: Helps computer talk like human.

④

Find the challenges of DL are :-

→ Needs lots of data

→ Takes long time to train.

→ overfitting

→ High cost.

Implement of AND using Perceptron :-

x_1	x_2	$x_1 \text{ AND } x_2$	Yactual
0	0	0	0
0	1	0	0
1	0	0	0
1	1	1	1

Suppose $w_1 = 1, w_2 = 1, b = -1.5$

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

case-1:

$$x_1 = 0, x_2 = 0$$

$$z = x_1 w_1 + x_2 w_2 + b = -1.5 < 0$$

$$Y_{\text{pred}} = f(z) = 0 = Y_{\text{actual}} \quad (\text{no weight update})$$

case-2:

$$x_1 = 0, x_2 = 1$$

$$z = 0 \times 1 + 1 \times 1 + (-1.5) = -0.5 < 0$$

$$Y_{\text{pred}} = f(z) = 0 = Y_{\text{actual}} \quad (\text{no weight update})$$

case-3:

$$x_1 = 1, x_2 = 0$$

$$z = -0.5 < 0$$

$$Y_{\text{pred}} = f(z) = 0 = Y_{\text{actual}} \quad (\text{no weight update})$$

case-4:

$$x_1 = 1, x_2 = 1$$

$$z = 1 + 1 - 1.5 = 2 - 1.5 = 0.5 > 0 \Rightarrow 1$$

$$Y_{\text{pred}} = f(z) = 1 = Y_{\text{actual}} \quad (\text{no weight update})$$

Vector format :-

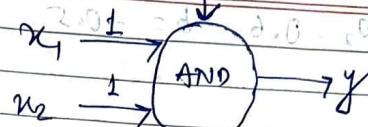
$$\begin{matrix} x_0 & x_1 & x_2 \\ \hline x_1 & \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \end{bmatrix} & \end{matrix}$$

$$w = \begin{bmatrix} 1 \\ 0 \\ -1.5 \end{bmatrix}$$

4×3

3×1

~1.5



$$\text{Ques: } w_1 = 1, w_2 = 0.6 \Rightarrow b = 0, \eta = 0.5$$

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

$$z = 1.0 - 2.0 + 0.6 = -0.4$$

$$\text{Ans: } w_1 = 1, w_2 = 0 \Rightarrow b = 0, \eta = 0.5$$

weight

$$Y_{\text{pred}} = f(z) = 0 = Y_{\text{actual}} \quad (\text{no update})$$

$$\text{Ques: } w_1 = 0, w_2 = 1$$

$$z = 0 + 0.6 + 0 = 0.6$$

$$Y_{\text{pred}} = f(z) = 0 = Y_{\text{actual}} \quad (\text{no weight update})$$

case-3:

$$x_1 = 1, x_2 = 0$$

$$z = 1 \times 1 + 0 \times 0 = 1$$

$$Y_{\text{pred}} = f(z) = 1 \neq Y_{\text{actual}}$$

$$w_{\text{new}} = w_1 + \eta (Y_{\text{actual}} - Y_{\text{predicted}}) \times x_1$$

$$= 1.0 + 0.5 (0 - 1) \times 1$$

$$= 1.0 - 0.5 = 0.5$$

$$w_2 = 0.6 + 0.5 (0 - 1) \times 0$$

$$= 0.6$$

$$b = 0 + 0.5 (0 - 1) = -0.5$$

Date _____
Page _____

 $w_1 = 0.7, w_2 = 0.6, b = +0.5$

case-1:

$x_1 = 0, x_2 = 0$

$z = -0.5 < 0$

$y_{pred} = f(z) = 0 = Y_{actual}$ (no weight update)

case-2:

$x_1 = 0, x_2 = 1$

$z = 0 + 0.6 - 0.5 = 0.1 < 0$

$y_{pred} = f(z) = 0 = Y_{actual}$ (no weight update)

case-3:

$x_1 = 1, x_2 = 0 = 0 = (0) + 0$

$z = 0.7 + 0 + -0.5 = 0.2 < 0$

$y_{pred} = f(z) = 0 = Y_{actual}$ (no weight update)

case-4:

$x_1 = 0 + 0.6 + 0 = 0$

$x_2 = 0 = (0) + 0$

$z = 0.7 + 0.6 - 0.5 = 0.8 > 1$

$y_{pred} = f(z) = 0 \neq Y_{actual}$

$w_1 \leftarrow$

$w_1 + \eta (Y_{act} - Y_{pred}) x_1$

$= 0.7 + 0.5 (1-0) \times 1 = 1.2$

$w_2 \leftarrow$

$w_2 + \eta (Y_{act} - Y_{pred}) x_2$

$b \leftarrow -0.5 + 0.5 (1-0) = 0$

case-1:

$x_1 = 1 - 0 = 2.0 + 0.1 =$

$x_2 = 0, x_2 = 0 = 2.0 - 0.1 =$

$z = 0$

$y_{pred} = f(z) = 0 \neq Y_{actual}$ (no weight update)

$2.0 - 0 = (1-0) 2.0 + 0 = 2.0$

case-2:

$x_1 = 0, x_2 = 1$

$z = 1.0$

$y_{pred} = f(z) = 1 \neq Y_{actual}$

$w_1 \leftarrow 1.2 + 0.5 (0-1) \times 0 = 1.2$

$w_2 \leftarrow 1.1 + 0.5 (0-1) \times 1 = 0.6$

$b \leftarrow 0.5 + 0.5 (0-1) = -0.5$

case-1:

$x_1 = 0, x_2 = 0$

$z = -0.5$

$y_{pred} = f(z) = 0 \neq Y_{actual}$ (no weight update)

case-2:

$x_1 = 0, x_2 = 1$

$z = 0.6 - 0.5 = 0.1$

$y_{pred} = f(z) = 0 = Y_{actual}$ (no weight update)

case-3:

$x_1 = 1, x_2 = 0 = 0 + 0$

$z = 0.7 - 0.5 = 0.2 > 0$

$y_{pred} = f(z) = 0 = Y_{actual}$ (no weight update)

case-4:

$x_1 = 2.0 + 0.1 = 2.1 + 0.1 =$

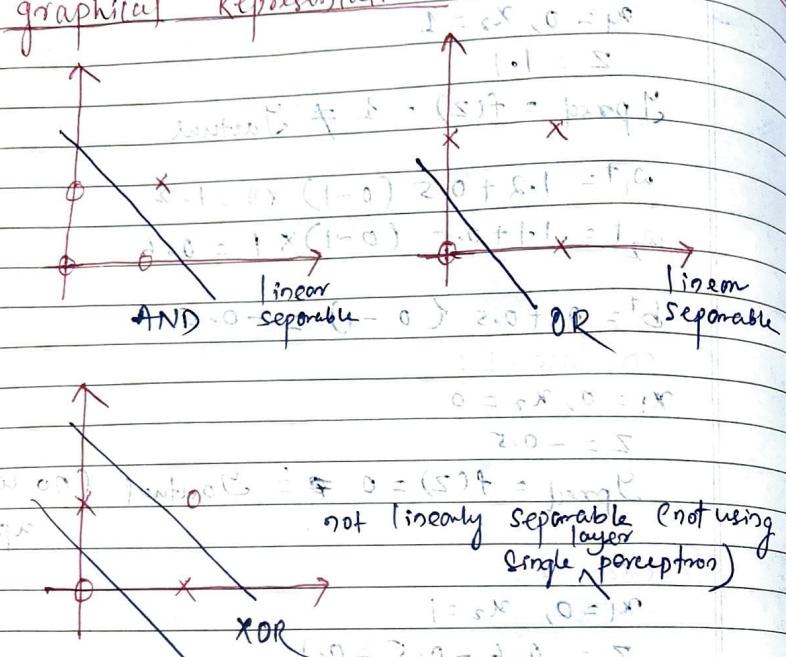
$x_2 = 0, x_2 = 0 = 2.0 - 0.1 =$

$z = 0$

$y_{pred} = f(z) = 1 \neq Y_{actual}$ (no weight update)



graphical Representation :-



Linear Separability :-

→ Linear separability of the point is the ability to classify the data points, in the hyper plane by avoiding the overlapping of the classes in the plane.

Only Single layer perceptron use for linearly separable data.

→ For 'XOR' It is a multiple layer used, so we use multi-layer perceptron for non-linearly separable data.

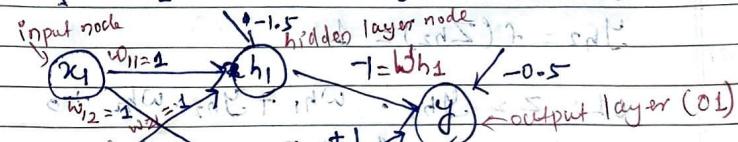
q-learnable parameters, no. of hidden layer
increases = no. of learnable parameter increases.

XOR : (Architecture of Multilayer perceptron)

$$Y_{\text{actual}} = x_1 w_1 + x_2 w_2 + b = 1.5$$

$$\begin{matrix} x_1 & x_2 \\ 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{matrix} \quad H_1 = (H_2) \quad Y_{\text{act}} = H_1 \cdot H_2$$

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



$$\begin{aligned} & \text{input node } x_1 \rightarrow h_1 \rightarrow o \\ & \text{input node } x_2 \rightarrow h_1 \rightarrow o \\ & \text{input node } x_1 \rightarrow h_2 \rightarrow o \\ & \text{input node } x_2 \rightarrow h_2 \rightarrow o \end{aligned}$$

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

$$f(z) = 0 \text{ otherwise}$$

$$w_{ij} = i = \text{input node}$$

$$j = \text{hidden node}$$

$$\text{case-1: } z = 2.0 + 2.0 + 1.0 = 5.0$$

$$x_1 = 0, x_2 = 0$$

$$h_1:$$

$$z = x_1 \cdot w_{11} + x_2 \cdot w_{21} + b$$

$$= -0.5 \cdot 2.0 + 2.0 + 1.0 = 1.0$$

$$y_{h1} = f(z_{h1}) = 0 \quad (z_{h1} < 0)$$

$$h_2: z = x_1 \cdot w_{12} + x_2 \cdot w_{22} + b$$

$$= -0.5 \cdot 2.0 + 2.0 + 1.0 = 1.0$$

$$y_{h2} = f(z_{h2}) = 0$$

$$o_1: z = y_{h1} \cdot w_{h1} + y_{h2} \cdot w_{h2} + b$$

$$= -0.5 < 0 \quad (\text{final})$$

$$y_{\text{pred}} = f(z) = 0 = Y_{\text{actual}} \quad (\text{no weight update})$$

Case-2:

$$x_1 = 0, x_2 = 1$$

$$h_1 = Z_{h_1} = x_1 \cdot w_{11} + x_2 \cdot w_{21} + b \\ = 1 \times 1 - 1.5 = 1 - 1.5 = -0.5 < 0$$

b) $y_{h_1} = f(Z_{h_1}) = 0$

$$h_2 = Z_{h_2} = x_1 \cdot w_{12} + x_2 \cdot w_{22} + b \\ = 0 + 1 - 0.5 = 0.5 > 0$$

$$y_{h_2} = f(Z_{h_2}) = 1$$

$$O_1: z = y_{h_1} \cdot w_{h_1} + y_{h_2} \cdot w_{h_2} + b \\ = 0 + 1 - 0.5 = 0.5 > 0$$

$$y_{pred} = f(z) = 1 = y_{actual} \quad (\text{no weight update})$$

case-3:

$$x_1 = 1, x_2 = 0$$

$$h_1 = Z_{h_1} = 1 - 1.5 = -0.5 < 0$$

$$y_{h_1} = f(Z_{h_1}) = 0$$

$$h_2 = Z_{h_2} = 0.1 - 0.5 = 0.5 > 0$$

$$y_{h_2} = f(Z_{h_2}) = 1 = y_{actual}$$

~~Ypred = f(z) = 1 = Yactual~~

$$O_1: 0 + 1 - 0.5 = 0.5 > 0$$

$$y_{pred} = f(z) = 1 = y_{actual}$$

(Case-2) $O_2 = 1.0$

(Step 3) $O_3 = 1.0$

Case-4:

$$x_1 = 1, x_2 = 1$$

$$h_1 = Z_{h_1} = 1 + 1 - 1.5 = 0.5 > 0$$

b) $y_{h_1} = f(Z_{h_1}) = 1$

$$h_2 = Z_{h_2} = 1 + 1 - 0.5 = 2 - 0.5 = 1.5 > 0$$

b) $y_{h_2} = f(Z_{h_2}) = 1$

$$O_1: 1 + 1 - 0.5 = 1.5 > 0$$

$y_{pred} = f(z) = 1 = y_{actual}$ ("no weight update")

What is multilayer perceptron?

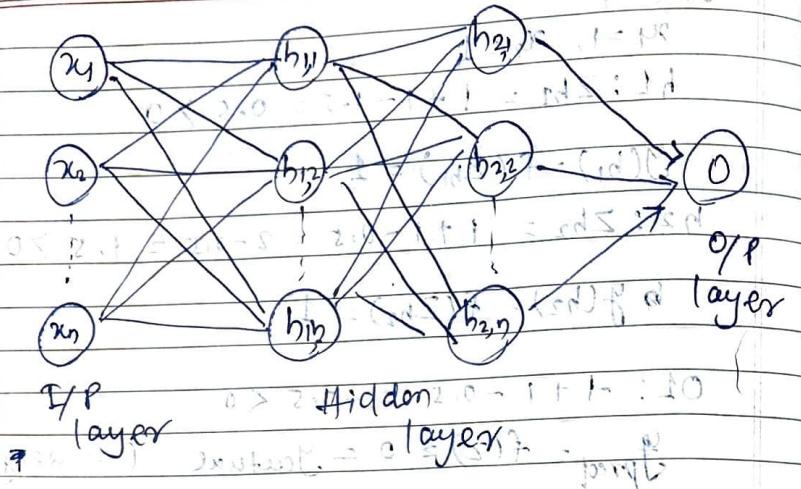
→ It is a type of artificial neural network

→ In ~~single~~ which all the nodes were interconnected with nodes of different layers.

→ Likewise a single-layer perceptron model or multilayer perceptron model also has the same model structure, but greater no. of hidden layers.

→ A multilayer perceptron neural network belongs to other feed forward neural networks because it works in the ~~backward~~ forward direction.

→ It has been widely used in various fields including the image recognition, speech recognition and natural language processing.



→ The multilayer perceptrons also known as backpropagation algorithm which executes using two main stages as follows -

- (i) Forward Stage :- Activation function starts from the first layer and terminates in the output layer.
- (ii) Backward Stage :- In this stage weight and bias values are modified as per the model requirements.
- In this stage the error between the actual output and the predicted output originates from the output layer and depends on the input.

- ① Multilayer perceptron algorithm :-
- ② Step-1: Initialize the weights and bias to small random values near to '0'.
- ③ Set the learning rate η in the range

- ④ of 0.01 to 0.1 of no. of iterations
 - ⑤ check for the stop condition if stop condition is false do step 3 to 7
 - ⑥ From 20 to 20 stop after 0.05
 - ⑦ for each training pair do step 4 to 7.
 - ⑧ set activation function of the output units in the 2nd last layer
 - ⑨ calculate the weighted sum.
- $$\text{weighted sum } (z) \text{ range} \leq x_1 w_1 + \dots + x_n w_n$$
- Apply the activation function of the weighted sum.

NOTE: * For multilayer network based on the no. of layers step 6 and 7 are repeated.

⑩ If the target is predicted output is not equal to the actual output then update the weight and bias based on the perceptron learning rules.

$$\eta = \text{Eta}$$

$$w_{i(\text{new})} = w_{i(\text{old})} + \eta (Y_{\text{actual}} - Y_{\text{predicted}}) \times x_i$$

$$b_{(\text{new})} = b_{(\text{old})} + \eta (Y_{\text{actual}} - Y_{\text{predicted}})$$

- ⑪ Test for stop condition.

- ⑫ Advantages of Multilayer perceptron :-
- This model can be used to solve complex non-linear problems.
- It works well with small and large input data.

→ It helps us to obtain quick prediction after the training.

→ It helps to obtain same accuracy ratio with large as well as small data.

* Disadvantages of Linear Function

(1) Computation are difficult and time consuming.

(2) It is difficult to predict how much the dependent variables affect each independent variable.

(3) The model's functioning depends on the quality of the training.

Gradient descent :-

at A

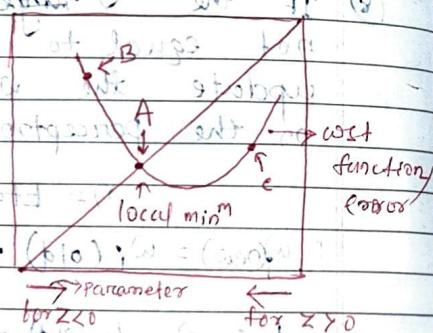
A. $z = 0$, $f'(z) > 0$

B. $f'(z) = 0$, $f''(z) < 0$

C. $z < 0, f'(z) < 0$

at e,

c. $z > 0, f'(z) > 0$



→ Gradient descent is a fundamental optimization algorithm in deep learning which is used to minimize the model's cost function by iteratively adjusting its parameters. It then moves in the direction of the negative steepest descent guided by the negative of the gradients.

* Hidden layer we use ReLU function.

→ It is an algm used to find the best soln. to a problem by making small adjustment in the right direction.
Learning rate:-

→ It is defined as a size of the steps towards the desired location or to reach the minimum cost value.

→ It is important hyper parameter in gradient decent that controls how big or small the step should be when going towards ip gradient for updating the model's parameter.

→ To achieve the min cost we have diff. techniques to prevent overfitting problem.

(a) weight regularization

(b) gradient clipping

→ we need to choose the right learning rate.

* Local min :-

→ If we move towards the -ve gradient or away from the gradient of a function at a current point then it will give the local min of the function.

* Local max :-

→ When we move towards +ve gradient or towards the gradient of the function at a current point then we will get the local max of the function.

$$\frac{d}{dx} = 2x \frac{d}{dx}$$

Minimizing Squared Error Loss on 2nd term

$$J(\omega, b) = \frac{1}{n} \sum_{i=1}^n (y_{\text{act}} + y_{\text{pred}})^2$$

$$y_{\text{pred}} = \sum w_i x_i + b$$

Gradient W.r.t. w_i :- unit vector of

$$\frac{\partial J(\omega, b)}{\partial \omega} = \frac{\partial}{\partial \omega} \left[\frac{1}{n} \sum_{i=1}^n (y_{\text{act}} - y_{\text{pred}})^2 \right]$$

$$= \frac{1}{n} \times 2 \sum_{i=1}^n (y_{\text{act}} - y_{\text{pred}}) \frac{\partial}{\partial \omega} (y_{\text{act}} - y_{\text{pred}})$$

$$= \frac{2}{n} \sum_{i=1}^n (y_{\text{act}} - y_{\text{pred}}) \left[\frac{\partial}{\partial \omega} y_{\text{act}} - \frac{\partial}{\partial \omega} (\sum w_i x_i + b) \right]$$

$$= \frac{2}{n} \sum_{i=1}^n (y_{\text{act}} - y_{\text{pred}}) (0 - n - 0) = \frac{-2n}{n} \sum_{i=1}^n (y_{\text{act}} - y_{\text{pred}})$$

$$\text{so gradient wrt } w_i \text{ is } \frac{-2n}{n} (y_{\text{pred}} - y_{\text{act}})$$

* Gradient W.r.t. b :-

$$\Delta b = \frac{\partial J(\omega, b)}{\partial b} = \frac{\partial}{\partial b} \left[\frac{1}{n} \sum_{i=1}^n (y_{\text{act}} - y_{\text{pred}})^2 \right]$$

$$= \frac{1}{n} \times 2 \sum_{i=1}^n (y_{\text{act}} - y_{\text{pred}}) \frac{\partial}{\partial b} (y_{\text{act}} - y_{\text{pred}})$$

$$= \frac{2}{n} \sum_{i=1}^n (y_{\text{act}} - y_{\text{pred}}) \left[\frac{\partial}{\partial b} y_{\text{act}} - \frac{\partial}{\partial b} (\sum w_i x_i + b) \right]$$

$$= \frac{2}{n} \sum_{i=1}^n (y_{\text{act}} - y_{\text{pred}}) (0 - 0 - 1) = \frac{-2}{n} \sum_{i=1}^n (y_{\text{act}} - y_{\text{pred}})$$

$$= \frac{2}{n} \sum_{i=1}^n (y_{\text{pred}} - y_{\text{act}})$$

Date _____
Page _____

+ve gradient :- when gradient is +ve

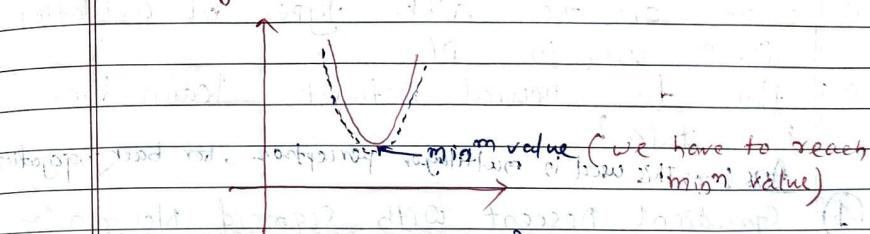
derivative of cost function

initial value $f(z) > 0$, so it should decrease to left

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial}{\partial w} J(w, b)$$

$$b_{\text{new}} = b_{\text{old}} - \eta \frac{\partial}{\partial b} J(w, b)$$

-ve gradient :-



$$w_{\text{new}} = w_{\text{old}} + \eta \frac{\partial}{\partial w} J(w, b)$$

$$b_{\text{new}} = b_{\text{old}} + \eta \frac{\partial}{\partial b} J(w, b)$$

Step of gradient descent :-

1 Initialize the parameter of the model randomly.

2 compute the gradient of the cost function wrt each parameter. it involves partial differentiation of cost function with respect to the parameter.

3 update the parameter of the model by taking steps in the opposite direction of the model here we choose the

a hyperparameter i.e. learning rate which helps in deciding the step size of the gradient.

→ Repeat Step 2 and 3 iteratively to get the best parameter for the define model.

Q.1 What are the diff. layer in ANN?

What ~~node~~ is the notation representing a node of a particular layer.

Q.2 What is cost function in DL.

Q.3 What are the activation functions in DL and where it is used.

Q.4 What are the diff. types of activation func used in DL.

Q.5 How do neural network learn from the data?

~~This used is multilayer perceptron, for backpropagation~~

Gradient Descent with sigmoid Neuron in

(1) $\sigma(z)$ (sigmoid function) :- (binary sigmoid function used)

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

(2) Sigmoid neuron's output (\hat{y}) :-

$y = \sigma(z) = \sigma(wx+b)$ (for single value)

(3) Cost function (J) :-

$J = \frac{1}{2} \sum (y - \hat{y})^2$ (stages)

Step-1 :- Derivative of cost function w.r.t w

$$\frac{\partial c}{\partial w_j} = \frac{\partial c}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \left(\frac{\partial z}{\partial w_j} \right)$$

i) $\frac{\partial z}{\partial w_j} = \frac{\partial}{\partial w_j} (wx_j + b) = x_j \cdot \frac{\partial w}{\partial w_j} = x_j$

ii) $\frac{\partial c}{\partial \hat{y}} = \frac{\partial}{\partial \hat{y}} (\frac{1}{2}(y-\hat{y})^2)$

$$= \frac{1}{2} x^2 (y-\hat{y}) \frac{\partial}{\partial \hat{y}} (\hat{y}-y)$$

$$= (y-\hat{y}) = \hat{y}-y$$

iii) $\frac{\partial \hat{y}}{\partial z} = \frac{\partial}{\partial z} \sigma(z) = \sigma(z)(1-\sigma(z))$

we can write $\frac{\partial \hat{y}}{\partial z} = \hat{y}(1-\hat{y})$

$\frac{\partial c}{\partial w_j} = (\hat{y}-y) \hat{y}(1-\hat{y}) x_j$

Step-2 : find the derivative of cost func w.r.t bias.

$$\frac{\partial c}{\partial b} = \frac{\partial c}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial b}$$

① derivative of 'cost func'

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

② $\frac{\partial \hat{y}}{\partial z} = \hat{y}(1-\hat{y})$

$$\textcircled{8} \quad \frac{\partial z}{\partial b}, \frac{\partial (w_i x_i)}{\partial b} = 0 + 1 = 1$$

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

update the weight and bias :-

$$w_j(\text{new}) = w_j(\text{old}) - \eta \frac{\partial c}{\partial w_j}$$

$$= w_j(\text{old}) - \eta [(y - \hat{y}) \cdot \hat{y}(1-\hat{y}) \cdot x_j]$$

$$\rightarrow b(\text{new}) = b(\text{old}) - \eta \frac{\partial c}{\partial b}$$

$$= b(\text{old}) - \eta [(y - \hat{y}) \cdot \hat{y}(1-\hat{y})]$$

 challenges of Gradient descent algm:-

 Local minima and Saddle points:-

 Local min: represent the point, where the loss function reaches the locally min value, compared to nearby points.

 The gradient of the loss function is zero at local minima.

 saddle point, represents the points where the gradient is zero but not all direction are flat potentially causing optimization difficulties.

 The gradient of the loss function is also zero at saddle points but not all directions are flat.

 optimization algm may encounter difficulties as the gradient is '0', making it challenging to escape and continue descending towards the global min.

 Vanishing gradient and exploding gradient:-

 vanishing gradient :-

 occurs if vanishing gradient is the phenomenon that occurs during the training the deep neural network where the gradient that are used to update the network become extremely small or vanish as they backpropagate from the output layer to the early layer.

 exploding gradient :-

 it happens when the gradient is too large causing unstable model, in this case the model weights will grow too large and they will eventually be represented as NaN.

 so one solution to this issue is to leverage or dimensionality reduction technique, which can help to minimize the complexity of the model.

 Advantage of gradient descent :-

 Ease of Implementation

 convergence guarantee

 Scalability

 Versatility

 Flexibility

Disadvantages of gradient descent

- (1) Sensitivity to learning rate (more, more complex)
- (2) Sensitivity to initialization (weight and bias)
- (3) Convergence speed (depends on learning rate)
- (4) Termination of algorithm (takes long time)
- (5) Correctness (best fit)

Assignment :-

Q.1) The input to a single input neuron is 2.0, its weight is 2.3 and its bias is -3.2. Present steps to find output.

- (i) What is the net input to the transfer function?
- (ii) What is the neuron output for the following transfer function using Sigmoid activation, tanh activation function and ReLU activation function.

Q.2) Given a 2 input neuron with the following parameters bias = 1.2, $w_2 [3 \ 2]$ and the input $(x) = [-5 \ 6]^T$, calculate the neuron output for the following tanh transfer function using bipolar sigmoidal func, binary sigmoidal func and tanh activation function.

Q.3) A single layer neurons network is to have 6 input and 2 outputs, the output are to be limited to and continuous over the range of 0 to 1. What can you tell about the network architecture?

- (i) Specifically how many neurons are required
 - (ii) What are the dimensions of the weight matrix.
 - (iii) What kind of transfer function could be used.
 - (iv) Is Bias required?
- (4) The final layer of the network produces logits for 3 classes, given logits are $Z = [2.0 \ 1.0 \ 0.1]$, True class is class 1, use softmax to get probability P_i and cross entropy loss function as $L = -\log P_{true}$, compute the probability and loss. (upto 4 decimal places), use numerically stable softmax (use padding starting from 0).

class work :-

- (1) What is overfitting and How to avoid it.
- (2) Define learning rate in DL.
- (3) Why is bias value added.
- (4) How weights are initialized in neural network.
- (5) What is cross entropy loss function.
- (6) What is Single layer perceptron? What's multilayer perceptron and how it is different from single layer perceptron.
- (7) How are the number of hidden layers and neurons in each hidden layer selected.

Answer to class work Question :-

- (1) → **Overfitting** :- When a deep learning model learns the training data too well, including noise and minor details, it performs poorly on new (test) data.
- ways to avoid overfitting :
- use more training data - improves generalization
 - Early stopping - Stop training when validation loss increases.
 - use simpler models - reduce model complexity.

- (2) → **Learning rate** : It is a hyperparameter that controls how much the model's weights are updated during training after each iteration.
- When high learning rate, faster learning but may skip optimal points, and more accurate when low learning rate.
- (3) → Bias value is added to give the model flexibility to fit data better by allowing the activation function to shift left or right.
- Without bias, all activations would pass the origin (0,0).

- (4) → **Weights initialization** means setting the starting value of weights before training begins.

- common methods -
- Zero initialization** : All weights = 0 (not used - because it causes symmetry problem)
 - Random initialization** - small random values (breaks symmetry)

Date _____
Page _____

~~0°C = 273 K~~ **meaningless value**
no temp.

Date _____
Page _____

The response variable in a dataset is the output variable.

Ex: Insulin level → Diabetic

Age	BMI	Weight	Insulin level	Diabetic
38	42.5	58	96	No
42	38.6	62	150	Yes
55	30.5	59	94	No
69	48.6	65	194	Yes
72	57.5	70	202	Yes
38	38.1			
60	38.1			

Q. The following are the attributes of a student. State if the data are nominal, ordinal, interval or ratio scale.

- Gender → Male → nominal
- class → Third year → ordinal
- Heart rate → 75 beats per minute → ratio
- Blood pressure → 180/120 mm of Hg → ratio
- Blood type → A positive → nominal
- Blood sugar → 40 → ratio
- Known allergies → Nil → nominal
- Body temperature → 36°C → ratio

→ cross entropy loss measures the difference between predicted probabilities and actual labels in classification problem.

→ formula for binary,

$$L = -(y \log p + (1-y) \log (1-p))$$

→ y = true label, p = predicted probability
→ it is used in logistic regression (softmax)

6 Single layer perception

- (i) It has only one layer of weights.
- (ii) It has only input and output.
- (iii) It can solve linearly separable problem like AND, OR.

Multilayer Perception

- (i) It has one or more hidden layers between input and output.
- (ii) It has hidden layers between input and output.
- (iii) It can solve non-linearly problems like XOR.

Difference:

→ SLP has one layer and linear decision boundary but MLP has multiple layers and non-linear decision boundaries.

7 There is no fixed rule; it depends on the problem's but in general below:

- (a) For image or speech more than 2 layers for complex data.
- (b) For simple task we need 1 or 2 layers.
- (c) Neurons per layer usually betn input size and output size.
- (d) Too many neurons causing overfitting.

8 What are Feed forward neural network.

9 When should you choose Deep learning over Machine learning solution.

10 How do you choose the right deep learning approach for your problem and data.

11 What are some challenges some common

Challenges found in DL model and how would you overcome them.

12 Input (x) = $2, 0, 0, 0, 0, 0, 0, 0$
Weight (w) = $0, 7, 0, 0, 0, 0, 0, 0$
Bias (b) = $-0, 1, 0, 0, 0, 0, 0, 0$

Find the output of the neuron (value should be upto 4 decimal places)

8 → feed-forward neural network: a type of neural network where data moves only one direction from input to hidden layer then output with no loops.

→ It is used for tasks like classification and regression.

→ ex: multilayer perceptron (MLP)

9 → To choose deep learning, when we have

- (a) large amount of data
- (b) problem is complex
- (c) high computing powers

10 To choose the right deep learning approach -

- (a) understand the data type
- (b) consider dataset size
- (c) Define the task - Regression or generation
- (d) check computational resources

11 common challenges in deep learning are -

- (a) overfitting - use regularization, dropout, more data
- (b) Long training time - use GPU, batch normalization.

- (c) Underfitting - increase model complexity.
(d) Data scarcity and poor generalization.

(12) Given, $x = 2.0, w = 0.7, b = -0.1$
and activation function (z) = $wx + b$
 $= 0.7 \times 2.0 + (-0.1)$
 $= 1.4 - 0.1 = 1.3$

- (13) The perceptron has 2 input and one output where $w_1 = 0.2, w_2 = 0.3, b = -0.3$, $b_{\text{pred}} = 0.4$, input $x = [1, 1]$, $y_{\text{actual}} = 1$, $\eta = 0.1$. Find the updated weights after 1 iteration if the perceptron mis-classifies the input.

Ans: Given, $w_1 = 0.2, w_2 = 0.3, b = 0.4$,

$x = [1, 1]$, $y_{\text{actual}} = 1, \eta = 0.1$ (i)

Input is from $x = [1, 1]$ (i)

$z = w_1 x_1 + w_2 x_2 + b$ is $z = 0.2 \times 1 + 0.3 \times 1 + 0.4$ (i)

$= 0.2 + 0.3 + 0.4 = 0.9 < 1$ (i)

\therefore so $y_{\text{pred}} = f(z) = 0 \neq Y_{\text{actual}}$ (i)

\therefore input is not classified (i)

so, \rightarrow update weights (i)

$w_1^{+2} = w_1 + \eta (Y_{\text{act}} - Y_{\text{pred}}) \times x_1$ (i)

$= 0.2 + 0.1 (1 - 0) \times 1 = 0.3$ (i)

$= 0.2 + 0.1 = 0.3$ (i)

$w_2^{+2} = w_2 + \eta (Y_{\text{act}} - Y_{\text{pred}}) \times x_2$ (i)

$= 0.3 + 0.1 (1 - 0) \times 1 = 0.4$ (i)

$= 0.3 + 0.1 = 0.4$ (i)

$b^{+2} = b + \eta (Y_{\text{act}} - Y_{\text{pred}}) \times 1$ (i)

$= 0.4 + 0.1 (1 - 0) \times 1 = 0.5$ (i)

$= 0.4 + 0.1 = 0.5$ (i)

so, so updated weight values are $w_1^+ = 0.3, w_2^+ = -0.2$ (i)

Types of GD: (i) stochastic gradient descent

(ii) Batch training (iii) mini-batch training

(iv) Stochastic gradient descent with momentum

(v) Batch Gradient descent without momentum

(vi) Gradient descent with cross entropy loss

(vii) Gradient descent with cross entropy loss for entire dataset

(viii) Gradient descent with momentum

(ix) Gradient descent with adaptive learning rate

(x) Gradient descent with adaptive learning rate with momentum

(xi) Gradient descent with adaptive learning rate with momentum and adaptive learning rate

(xii) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum

(xiii) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate

(xiv) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate with momentum

(xv) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate

(xvi) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate

(xvii) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate

(xviii) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate

(xix) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate

(xx) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate

(xxi) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate

(xxii) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate

(xxiii) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate

(xxiv) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate

(xxv) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate

(xxvi) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate

(xxvii) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate

(xxviii) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate

(xxix) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate

(xxx) Gradient descent with adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate with momentum and adaptive learning rate

- (4) → Repeats Step 1 to 3, unless the loss function is minimised.
- Advantages of backtracking batch GD :-
- It produces less noise in comparison to other gradient descent.
 - It follows the true gradient of the loss function which can lead to a most stable convergence.
 - It is computationally efficient as all resources are used for all training samples.
 - It consider the entire dataset for each update providing a global perspective of the loss landscape.

Disadvantage :-

- For large dataset computing the gradient for the entire dataset can be computationally expensive and slow.
- This requires storing and processing the entire dataset in memory which can be impractical for very large dataset.

Stochastic gradient descent :- (commonly we use)

$$\theta_{new} = \theta_{old} - \eta \nabla_{\theta} (\theta, x^{(i)}, y^{(i)})$$

↓
It is a i^{th} training example

Gradient of loss function
w.r.t θ , for i^{th}
training dataset

- Stochastic GD is an alternative approach that updates the weight for each training

example rather than the entire dataset. This means the gradient is calculated and the weight and bias are updated for each training example.

- It is easier to store in allocated memory as it requires only one training example at a time.
- Due to frequent updates
- It shows some computational efficiency losses as compared to batch gradient descent as it shows frequent updates that requires more detail and speed.
- It is treated as a noisy gradient due to frequent updates and it converges faster when the dataset is large.

- Steps :- (update rule)
- ① Selects a single training example from the dataset.
 - ② Compute the gradient of the loss function w.r.t each weight and bias using only that training example.
 - ③ Update the weight by taking a step in the direction of -ve gradient.
 - ④ Repeat Step 1 to 3 for each training example.

Advantages :

- It is easier to allocate in the desire memory.
- It is relatively faster to compute than batch gradient descent.
- It is more efficient than large dataset as it requires less memory.

- It may get new minima.
- It is sometimes helpful in finding global minima and also in escaping from local minima.
- Disadvantages:**
 - It can suffer from high variance and noisy updates as each update is based on single example, which may not be representative of the entire data set.
 - It takes more iterations as compared to other gradient descent.
 - The frequent updates can cause the algorithm to overshoot the minimum, especially with the high learning rate.
 - It requires careful tuning of the learning rate to ensure stable and efficient convergence.
 - E.g. CNN (Where image processing) and training computational power is limited.
- iii) Mini Batch GD:**
 - $Q_{\text{new}} = Q_{\text{old}} - \eta \nabla c(\theta, \{x_i, y_i\}_{i=1}^m)$
 - Involves computing gradient of loss function w.r.t θ , for small subset of m , of the dataset.
 - It is a compromise b/w the batch gradient decent and stochastic GD, instead of using the entire dataset or a single training example, it computes the gradient and updates the weight using

- Small subset or mini batches of the entire dataset. This is because splitting the training data set into smaller data set makes a better balance to maintain the computation balancing of batch GD and Speed of stochastic GD.
- Steps of Mini GD :-**
 - ① → Select a mini batch from the training data set.
 - ② → Compute the gradient of the loss function w.r.t each weight and bias using only for the mini batches.
 - ③ → Update the weight by taking a step in the direction of -ve of gradient.
 - ④ → Repeat steps 1 to 3 from different mini batches until the entire data set has been traversed.
- Advantages :-**
 - It strikes a balance b/w stability of batch gradient and the efficiency of stochastic gradient descent.
 - It helps to reduce the variance of stochastic gradient decent by still being computationally efficient, especially when leveraging vectorized operations on modern hardware.
 - Mini batches allow for better utilization of hardware resources such as GPU, which are optimised for parallel computation.
- Disadvantages :**
 - It may get stuck at local minima.

- It requires careful tuning of the mini batch size and learning rate to ensure optimal performance.
- While most stable than stochastic GD, it can still be less stable than batch gradient descent especially if mini batch size is too small.
- Selecting an inappropriate mini batch size can lead to sub optimal performance and conversion issue.
- Ex: architecture size is huge and model is more complex and huge in amount of dataset.

- Q1:** Why to use Stochastic GD :-
- Stochastic GD brings several benefits to business and plays a crucial role in machine learning and AI.

Computational Efficiency

- » Computational efficiency
- » Faster convergence
- » Memory efficiency
- » Escapes local minima

Ques 2: Minimize the function $f(w) = (w-3)^2$

Given: $w_0 = 0$, $\eta = 0.1$, no. of iteration = 3

By using the GD, minimizing to find the value of w_3 to minimize $f(w)$.

Given: $w_0 = 0$, $\eta = 0.1$, no. of iteration = 3

Minimize the function, $f(x, y) = x^2 + y^2$

Given, Initial point: $(x_0, y_0) = (1, 1)$

Given η is 0.1, No. of iteration = 3

24/01/2013
Assignment
Topic (gradient descent)

Ans: we know, $\frac{\partial f}{\partial w}(w)$

$$w_{t+1} = w_t - \eta \cdot \frac{\partial f}{\partial w}(w)$$

$$t=0, w_1 = w_0 - \eta \cdot \frac{\partial f}{\partial w}(w_0)$$

$$= 0 - 0.1(2(w_0 - 3))$$

$$= -0.1(2w_0 - 6)$$

$$= -0.2w_0 + 0.6 = 0.6$$

$$t=1, w_2 = w_1 - \eta \cdot \frac{\partial f}{\partial w}(w_1)$$

$$= 0.6 - 0.1(-4.8) \Rightarrow 0.6 - 0.1(2w_1 - 6)$$

$$= 0.6 + 0.48 = 0.6 - 0.2w_1 + 0.6$$

$$= 1.08 \quad \cancel{= -0.4w_1 + 1.2 + 0.6}$$

$$t=2, w_3 = w_2 - \eta \cdot \frac{\partial f}{\partial w}(w_2)$$

$$= 1.08 - 0.4w_2 + 1.2 + 0.6$$

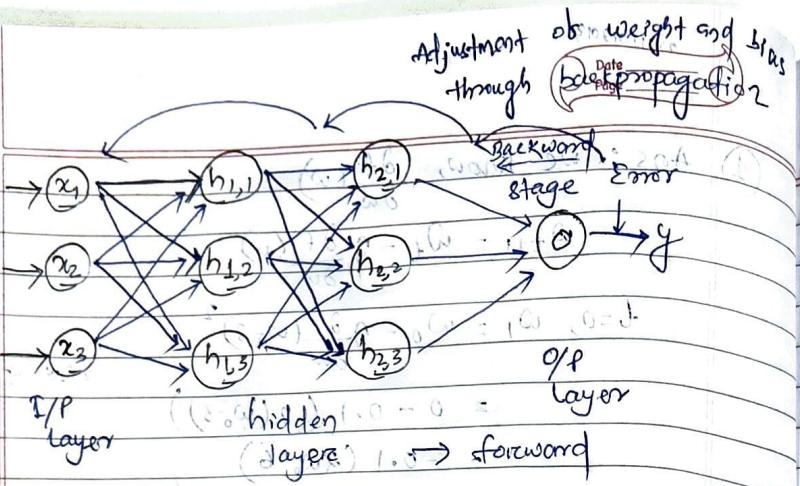
$$= 2 - 0.6w_2 + \cancel{0.4w_2} + 1.8 \quad \cancel{= 1.2 + 0.6}$$

$$\approx 1.08 + 0.1(3 \times 1.08 - 6)$$

$$= 0.696$$

$$t=0; w_1 = w_0 - \eta \cdot \frac{\partial f}{\partial w}(w_0)$$

Backpropagation algorithm :-



→ Backpropagation is a method used to train artificial neural network by minimizing the error. It is extensively used to train feed forward neural network such as CNN, RNN.

→ It enables the use of gradient methods to train multilayer network and update the weights to minimize the error.

Why Backpropagation is Important?

→ It plays a crucial role in how neural improves over time, the reason for this

(a) Efficient weight update

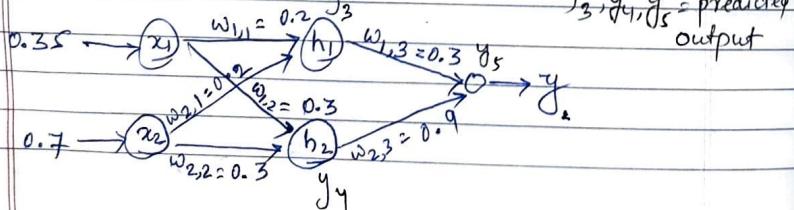
(b) Scalability

(c) Automated learning

→ The backpropagation algorithm involves

(a) forward pass

(b) backward pass



Assume sigmoid activation function for forward pass and backward pass.
 $y_{act} = 0.5, \eta = 1$

Ans: Forward calculation :: backpropagation

$$\text{hidden layer } z_j = \sum_i w_{ij} \times x_i = 0.4$$

$$\text{at } h_1 : z_1 = w_{11} x_1 + w_{12} x_2$$

$$= 0.2 \times 0.35 + 0.3 \times 0.7$$

$$= 0.21$$

$$y_1 = f(z_1) = \frac{1}{1 + e^{-z_1}} = \frac{1}{1 + e^{-0.21}} = 0.55$$

$$\text{at } h_2 : z_2 = w_{1,2} x_1 + w_{2,2} x_2$$

$$= 0.3 \times 0.35 + 0.3 \times 0.7 + 0.0$$

$$= 0.315$$

$$y_2 = f(z_2)$$

$$= \frac{1}{1 + e^{-z_2}} = \frac{1}{1 + e^{-0.315}} = 0.58$$

At O (Output layer):

$$z_3 = w_{1,3} \times y_1 + w_{2,3} \times y_2$$

$$= 0.3 \times 0.55 + 0.9 \times 0.58$$

$$= 0.687$$

Sigmoid function of z_3 :

$$y_3 = \frac{1}{1 + e^{-0.687}} = 0.67$$

but actual output is $y_{act} = 0.5$, so

δ = delta of gradients)

02/19/2012

Date _____
Page _____

Error calculation:

$$\text{error} = \frac{1}{2} \times \sum_{j=1}^m (y_{\text{pred.}} - y_{\text{act.}})^2 = 0.67 - 0.5 = 0.17$$

Backpropagation:

* δ_j = error term j for each unit

η = learning rate

To compute, e.g. gradient.

$$[\Delta w_{ij} = \eta \times \delta_j \times o_j]$$

for o_3 :

$$\begin{aligned}\delta_5 &= y_5 (1-y_5) \cdot (y_{\text{act.}} - y_{\text{pred.}}) \\ &= 0.67 (1-0.67) (0.5 - 0.67) \\ &= -0.0376\end{aligned}$$

for h_1 :

$$\begin{aligned}\delta_3 &= y_3 (1-y_3) (\Delta w_{1,3} \times \delta_5) \\ &= 0.55 (1-0.55) (0.3 \times -0.0376) \\ &= 0.55 (1-0.55) (-0.3 \times -0.0376) \\ &= -0.00279\end{aligned}$$

For h_2 :

$$\begin{aligned}\delta_4 &= y_4 (1-y_4) (\Delta w_{2,3} \times \delta_5) \\ &= 0.58 (1-0.58) (0.9 \times -0.0376) \\ &= -0.00824\end{aligned}$$

$\delta_5 = \text{delta}$

$$w_{\text{new}} = w_{\text{old}} + \Delta w_{ij}$$

Date _____
Page _____

Weight update

$$\text{for output layer } (0) * \Delta w_{2,3} = \eta \times \delta_5 \times y_4$$

$$= 0.1 \times (-0.0376) \times 0.58 = -0.0218$$

new weight,

$$w_{2,3}(\text{new}) = 0.9 - 0.0218 = 0.878$$

$$* \Delta w_{1,3} = \eta \times \delta_3 \times y_3$$

$$= 0.2 \times (-0.0376) \times 0.55 = 0.023$$

$$w_{1,3}(\text{new}) = 0.3 - 0.02 = 0.28$$

for hidden layer:

$$* \Delta w_{1,1} = \eta \times \delta_3 \times x_1$$

$$= 0.2 \times (-0.00279) \times 0.35 = 0.000976$$

$$w_{1,1}(\text{new}) = 0.2 - 0.000976 = 0.199$$

$$* \Delta w_{1,2} = \eta \times \delta_3 \times x_2$$

$$= 0.2 \times (-0.00279) \times 0.35 = -0.00288$$

$$w_{1,2}(\text{new}) = 0.3 - 0.00288 = 0.297$$

$$* \Delta w_{2,1} = \eta \times \delta_4 \times x_2$$

$$= 0.1 \times (-0.00824) \times 0.35 = -0.00288$$

Forward calculation:

At b_1 :

$$z_1 = w_{1,1} \times x_1 + w_{2,1} \times x_2 + b_1$$

$$= 0.2082 \times 2 + 0.55$$

$$y_1 = \sigma(z_1) = \frac{1}{1 + e^{-0.2082}} = 0.55$$

At b_2 :

$$z_2 = w_{1,2} \times x_1 + w_{2,2} \times x_2 + b_2$$

$$= 0.30975$$

$$y_2 = \sigma(z_2) = \frac{1}{1 + e^{-0.30975}} = 0.58$$

At b_3 :

$$z_3 = w_{1,3} \times y_1 + w_{2,3} \times y_2$$

$$= 0.571$$

$$y_3 = \sigma(z_3) = \frac{1}{1 + e^{-0.571}} = 0.64$$

Iteration - 2

do same process; again calculate back pass then

Algorithm:

Training algo:-

Step-1

→ Initialize the weight to small random values ; where input training vector is $x = [x_1, x_2, \dots, x_n]$.

Step-4

→ While the steps stopping condition is to be false, do step 3 & go to step 4.

Step-3

→ For each training pair do step 4 to

Step-4

→ Each input unit receives the signal x_i to all the units.

Step-5

→ For each hidden unit calculate the weight sum $z_j = \sum w_{ij} \times x_i + \text{bias}$, then

apply the activation function over the weighted sum and sends this signal to all the units in the layer about.

Step-6

calculation of propagation error for each output unit y_k , where k is 1 to n receives the target pattern to an input pattern and error is calculated.

$$y_k = f(z_j)$$

5.

- (1) Initialize the weights to small random values
- (2) → The input is modelled using true weights for the inputs that arrive through the pre-connected path.
- (3) calculate the output to each neuron from the input layer to the hidden layer then to the output layer etc.
- (4) Calculate the error in the output

$$\text{Backpropagation error} = Y_{\text{act}} - Y_{\text{pred}}$$

- (5) From the output layer, go back to the hidden layer and adjust the weights in order to reduce the error.
- (6) Repeat the process until the desired output is achieved.

Advantages:

- (i) Ease to Implementation.
- (ii) It is simple and flexible.
- (iii) Efficiency
- (iv) Generalization
- (v) Scalability

challenges
Disadvantages:-

- Vanishing gradient problem
- Exploding gradient.
- Overfitting