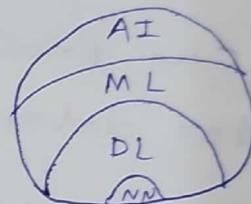


DLP

Artificial Intelligence (AI) - AI is the capability of a machine to imitate intelligent human behaviour.
AI is the tech that enables computers and machines to simulate human intelligence and problem solving capabilities.

Eg of Application

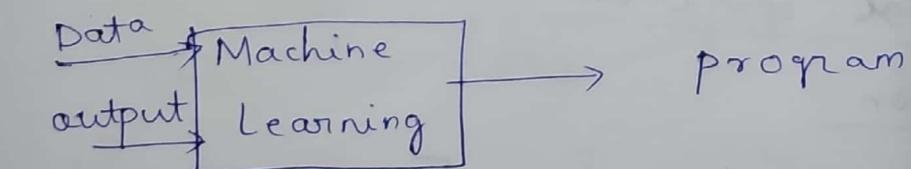
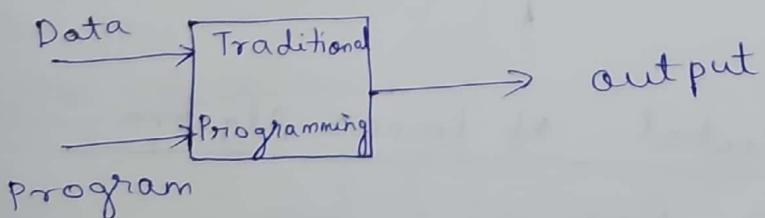
- search engines (Google)
- content Recommendations (youtube, Netflix)
- chatGPT
- self driving cars
- Facial recognition
- Computer games



Machine Learning

It is the study of computer algo that improve automatically through experience.

ML refers to the study of computer systems that learn & adapt automatically from experience without being explicitly programmed.



Deep Learning - is the subset of ML called ANN (Artificial Neural Network) which uses algo which are inspired by the structures & function of the brain. and are capable of self learning.

Difference b/w ML & DL

ML

- i) subset of AI
- ii) It can train on smaller datasets.
- iii) Require more human intervention to correct & learn.
- iv) shorter training
- v) Lower accuracy
- vi) It can train on a CPU.
- vii) Less complex

DL

- i) Subset of ML
- ii) It ^{requires} ~~can~~ train large datasets.
- iii) It learns of its own ~~for~~ from environment & part mistakes
- iv) Larger training
- v) higher accuracy
- vi) It needs specialized GPU.
- vii) more complex

ch-1 Fundamental of Linear Algebra

Matrix operations

Addition

$$\begin{bmatrix} 1 & 0 & 5 \\ 2 & 4 & 7 \end{bmatrix}_{2 \times 3} + \begin{bmatrix} 5 & 3 & 2 \\ 6 & 1 & 4 \end{bmatrix}_{2 \times 3} =$$

$$\begin{bmatrix} 6 & 3 & 7 \\ 8 & 5 & 11 \end{bmatrix}_{2 \times 3}$$

scalar Multiplication

$$2 * \begin{bmatrix} 1 & 0 & 5 \\ 2 & 4 & 7 \end{bmatrix} = \begin{bmatrix} 2 * 1 & 2 * 0 & 2 * 5 \\ 2 * 2 & 2 * 4 & 2 * 7 \end{bmatrix} = \begin{bmatrix} 2 & 0 & 10 \\ 4 & 8 & 14 \end{bmatrix}$$

$$\begin{bmatrix} 4 & 6 \\ 2 & 7 \end{bmatrix} * \begin{bmatrix} 1 & 9 \\ 5 & 2 \end{bmatrix} = \begin{bmatrix} 4 + 30 & 36 + 12 \\ 2 + 35 & 18 + 14 \end{bmatrix}$$

$$= \begin{bmatrix} 34 & 48 \\ 37 & 32 \end{bmatrix}$$

$$\begin{bmatrix} 4 & 6 \\ 2 & 7 \end{bmatrix} \begin{bmatrix} 1 & 9 \\ 5 & 2 \end{bmatrix} = \left[\begin{bmatrix} 4 & 6 \\ 2 & 7 \end{bmatrix} \begin{bmatrix} 1 \\ 5 \end{bmatrix} \begin{bmatrix} 4 & 6 \\ 2 & 7 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} \right]$$

$$= \begin{bmatrix} 4 * 1 + 6 * 5 \\ 34 & 48 \\ 37 & 32 \end{bmatrix}$$

Vector Operations

$$\begin{bmatrix} 4 & 3 & 6 \end{bmatrix}_{1 \times 3} \begin{bmatrix} 2 \\ 7 \\ 1 \end{bmatrix}_{3 \times 1} = \begin{bmatrix} 8 + 21 + 6 \end{bmatrix}_{1 \times 1}$$

$$\begin{bmatrix} 2 \\ 7 \\ 1 \end{bmatrix}_{3 \times 1} \begin{bmatrix} 4 & 3 & 6 \end{bmatrix}_{1 \times 3} = \begin{bmatrix} 35 \end{bmatrix}_{1 \times 1}$$

$$= \begin{bmatrix} 8 & 6 & 12 \\ 28 & 21 & 42 \\ 4 & 3 & 6 \end{bmatrix}_{3 \times 3}$$

Matrix vector Multiplication

$$\begin{bmatrix} 4 & 6 \\ 2 & 7 \\ 5 & 8 \end{bmatrix}_{3 \times 2} \begin{bmatrix} 1 \\ 5 \end{bmatrix}_{2 \times 1} = \begin{bmatrix} 4 + 30 \\ 2 + 35 \\ 5 + 40 \end{bmatrix} = \begin{bmatrix} 34 \\ 37 \\ 45 \end{bmatrix}_{3 \times 1}$$

$$= 1 \begin{bmatrix} 4 \\ 2 \\ 5 \end{bmatrix} + 5 \begin{bmatrix} 6 \\ 7 \\ 8 \end{bmatrix}$$

The fundamental Spaces

① column space

② Null space

Column Space

If A is a matrix, column space is represented by $C(A)$.

$C(A)$ = all possible combo of columns of A.

$$A = \begin{bmatrix} 1 & 1 \\ 2 & 1 \end{bmatrix}, c_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, c_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Feature

→ used for dimension reduction

→ used to find out relationship among features.

Null Space

The null space of matrix A is represented by $N(A)$.

$$Ax = 0$$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

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

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

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

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

$$\Rightarrow x_1 = -2x_2$$

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

Eigen Value & Eigen Vector



Vector representation of a matrix

→ It is a non-zero matrix

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

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

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

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

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

$$\Rightarrow \lambda^2 - 7\lambda + 10 = 0$$

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

6

$$\Rightarrow \lambda(\lambda-2) - 5(\lambda-2)$$

$$(\lambda-2)(\lambda-5)$$

$$\lambda_1 = 2, \lambda_2 = 5$$

considering $\lambda_1 = 5$,

$$(A - 5I)x = 0$$

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

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

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

$$\Rightarrow x_1 = x_2$$

$$\text{let } x_1 = 1, x_2 = 1$$

$$\text{eigen value } \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Types of Learning

- ① supervised Learning - All data are labelled
+ training data includes desired output. eg - classification, regression
- ② unsupervised " - no data is labelled
+ training doesn't include o/p
eg - clustering, generative models
- ③ semi-supervised " - datum singular
+ small amt of data is labelled (training data includes few desired o/p)
- ④ reinforcement "

↓

No labelled data
It tells which step is correct & which step is wrong

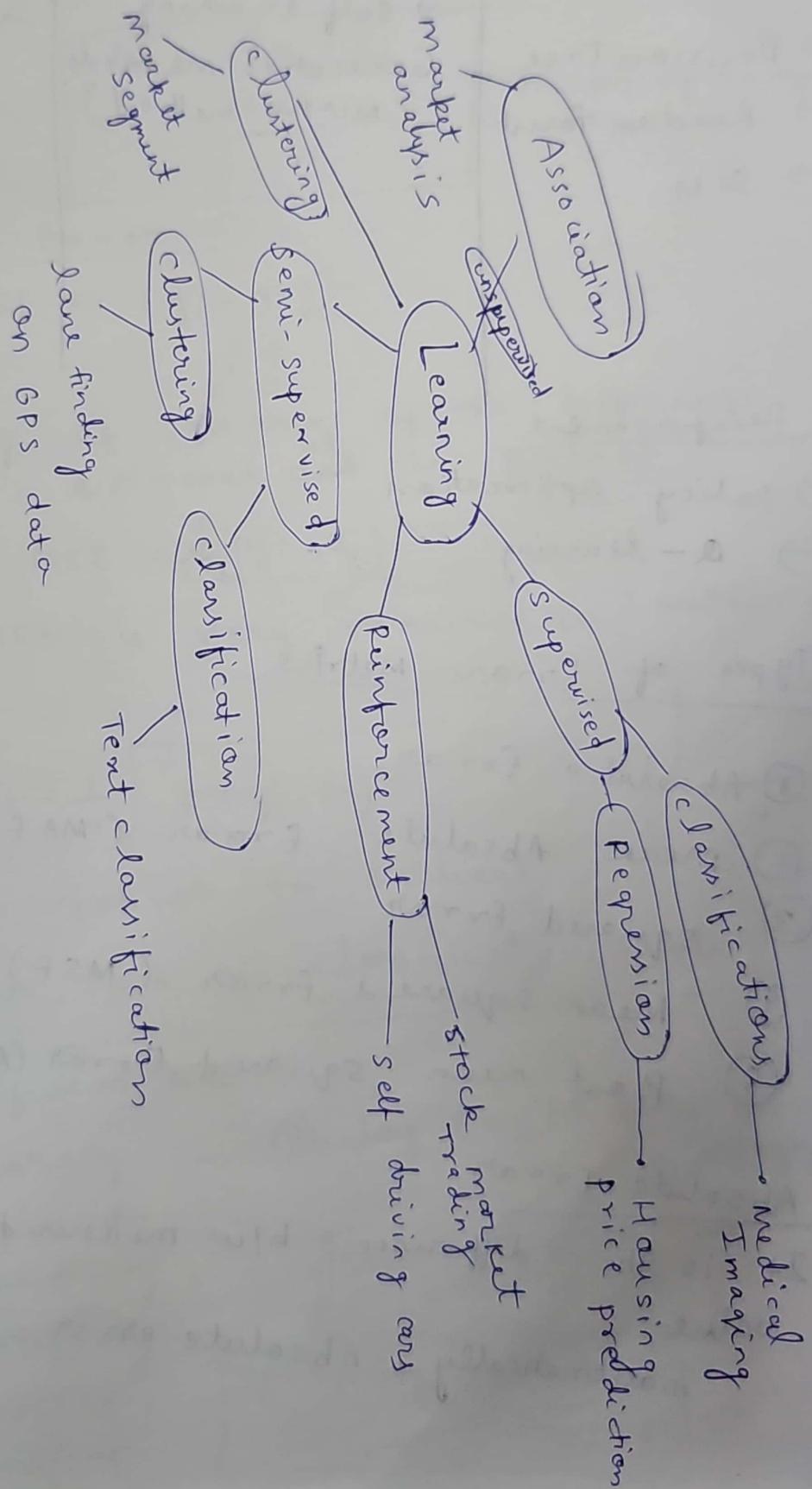
$$(5-\lambda)2 - (5-\lambda)1$$

$$(2-\lambda)(5-\lambda)$$

$$2 = \lambda, 5 = \lambda$$

$$0 = 01 + 1 + -5\lambda$$

$$0 = 01 + 1\lambda - 5\lambda - 5$$



Types of Learning platform

<u>Supervised</u>	<u>Semi-supervised</u>	<u>unsupervised</u>
<ul style="list-style-type: none">→ Logistic Regression→ Naive Bayes→ KNN→ Decision Tree→ Random forest→ SVM	<ul style="list-style-type: none">→ SSSVM (Semi-supervised SVM)→ self training Generative methods (Mixture method)	<ul style="list-style-type: none">→ k-means clustering→ t-SNE→ DBSCAN

Reinforcement

- policy optimization
- Q-learning

Types of Errors Metrics

- ① Absolute error
- ② Mean Absolute Error (MAE)
- ③ squared error
- ④ Mean squared error (MSE)
- ⑤ Root mean squared Error (RMSE)

Absolute Error

It is the difference b/w measured & actual value.

$$\text{mathematically, Absolute error} = |y_i - z_i|$$

x_i = actual value

y_i = Predicted value

Mean Absolute Error

MAE is a metric that is used to measure the arithmetic avg of deviation b/w predictions & actual value.

$$\text{MAE} = \frac{\sum_{i=1}^n |y_i - x_i|}{n}$$

n = no. of samples

Squared Error

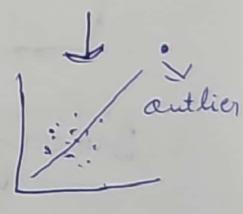
$$SE = (y_i - x_i)^2$$

MSE

The avg of squares of the difference b/w estimated ~~at~~ the actual values.

$$\text{MSE} = \frac{\sum_{i=1}^n (y_i - x_i)^2}{n}$$

The MSE is very sensitive to outliers.



RMSE

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - x_i)^2}$$

→ It is robust to outliers

Different types of classification Evaluation metrics

① confusion Matrix

② Accuracy

③ Precision

④ Sensitivity

⑤ Specificity

⑥ F1-score

↳ combo of precision & recall

⑦ Log loss / Binary cross entropy

$$\begin{aligned} \text{Sensitivity} &= \frac{TP}{TP + FN} \\ &= TPR \end{aligned}$$

$$\begin{aligned} \text{Specificity} &= \frac{TN}{TN + FP} \\ &= TNR \end{aligned}$$

Confusion Matrix

TP	FN Type-II error
FP Type-I error	TN

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

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

→ useful for skewed dataset / unbalanced

$$\text{F1-score} = \frac{2PR}{P+R}$$

Binary + entropy - It compares each of the predicted probability to ~~comp~~ actual class o/p. It is the negative avg of log of corrected predicted probabilities

$$= -\frac{1}{N} \sum_{i=1}^N y_i \log(p(y_i)) + (1-y_i) \cdot \log(1-p(y_i))$$

Entropy, cross-entropy, KL-divergence

Entropy measures uncertainty and randomness within a single probability distribution

Entropy

$$E(p) = - \sum p(x) \cdot \log p(x)$$

\sum → represents the sum over all possible events.

$p(x)$ → prob. of event x in the distr.
prob.

Cross-Entropy

It measures how well a predicted probability distribution (q) approximates a true probability distribution (p)

$$CE(p, q) = - \sum p(x) \cdot \log(q(x))$$

\sum - sum over all possible events

$p(x)$ - true probability of event x

$q(x)$ - predicted probability of event x

KL Divergence

① for discrete Distribution

② for continuous "

KL divergence measures the info lost when approximating the distribution p with distribution q .

The. KL • It quantifies the difference b/w 2 probability dist.

It results 0, iff the dist are identif.

$$KL(P||Q) = \sum p(i) \cdot \log \left(\frac{p(i)}{q(i)} \right)$$

p_i of i is the probability of event i
 q_i in the is the approximating
 denote Σ of all dist Q .
 possible outcomes \rightarrow prob. of event
 i in the predicted distr Q .

$$KL(P||Q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx$$

Properties

Non-negative

$$KL(P||Q) \geq 0$$

zero divergence

$$KL(P||Q) = 0$$

if the two dist are identical
~~and~~ i.e. $P = Q$

Ch-3

The Neural Network

The neuron - It is fundamental unit of human brain.

The avg human brain has approx 100 billions of neurons

$$= 10^{11}$$

It is a small cell that receive electrochemical signal from various sources & response by transmitting electrical impulse by to other neurons.

Each neuron is densely connected to 10^4 no. of neurons.

The weight of brain is 1.5 kg

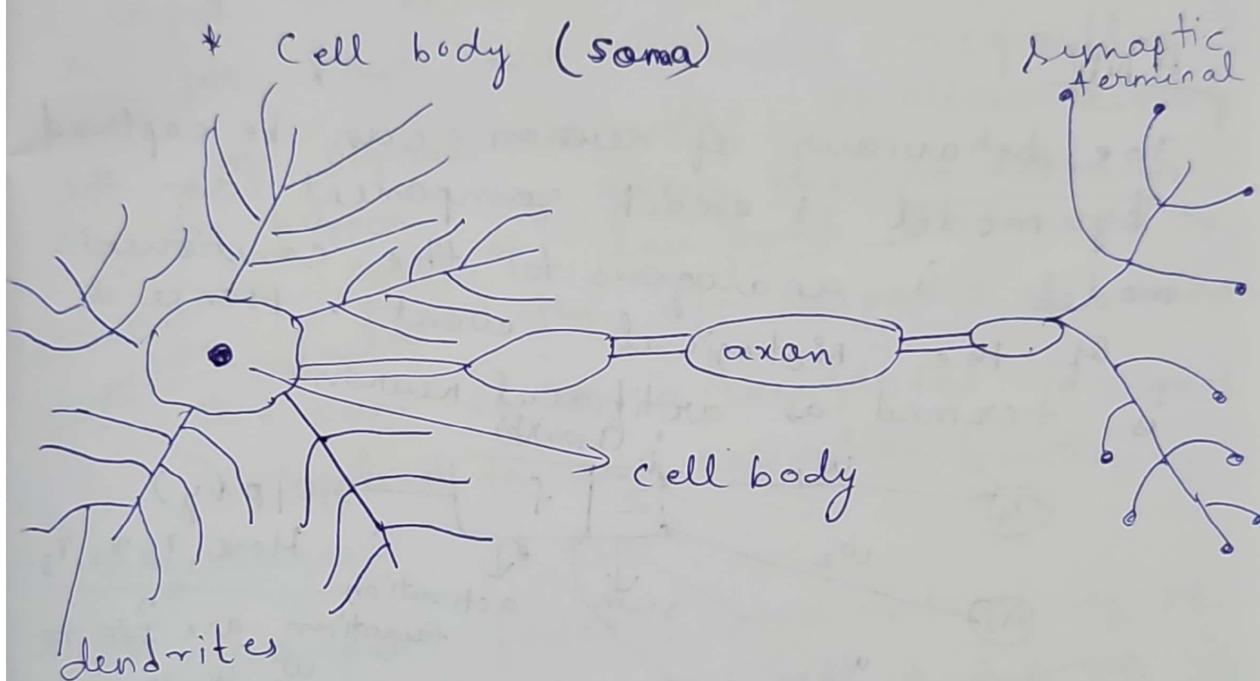
The 11^{th} neuron is 1.5×10^{-9} gm

A neuron consist of 3 main parts.

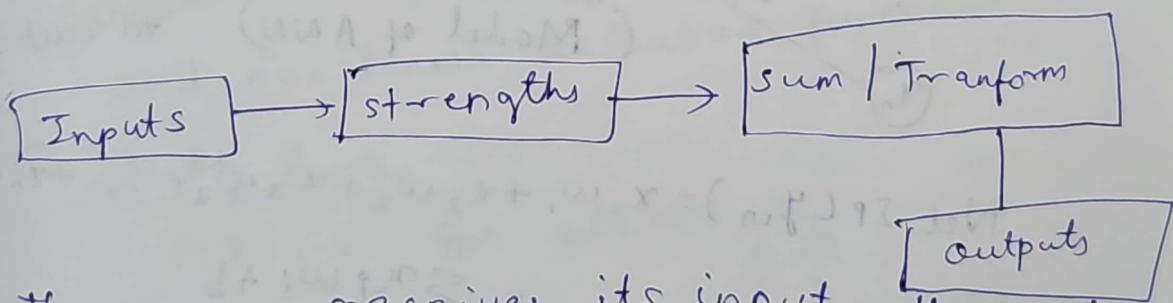
* dendrites

* Axon

* Cell body (soma)



(Structure of biological neuron)



- The neuron receives its input through antenna like structure called dendrites

- The dendrites behave as i/p chain

Each of these incoming connections is dynamically strengthen or weaken based on how often it is used.

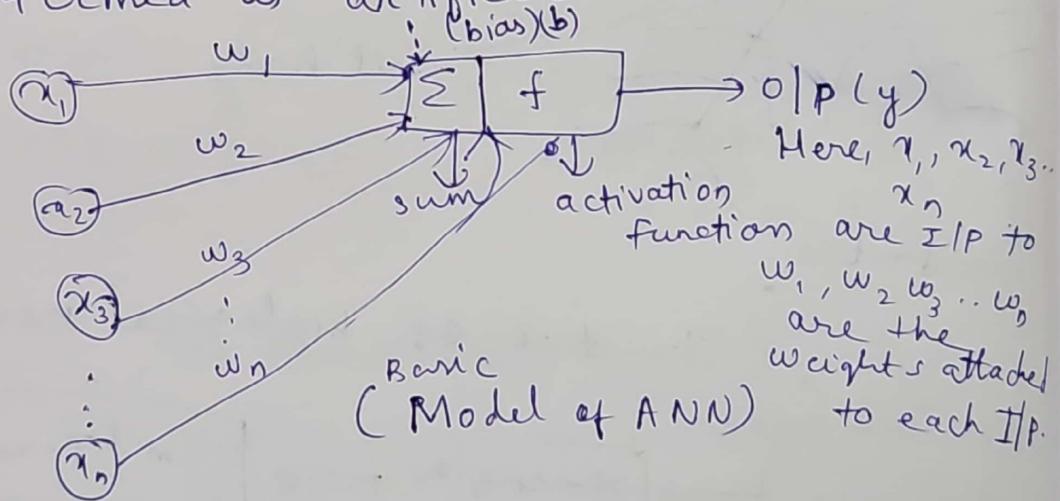
- The strength of each connection determines the contribution of i/p to n. neurons o/p.

After being waited by the strength of their respective connections, the I/P are summed together in the cell body

The sum is then transformed into a new signal that is propagated along the cell's axon & sent off to the other neurons.

ANN

The behaviour of neuron can be captured by model & each component of the model is analogous to the constituents of the biological events. Hence it is termed as artificial neuron.



$$\begin{aligned} \text{Net I/P } (y_{in}) &= x_1 w_1 + x_2 w_2 + x_3 w_3 + \dots + x_n w_n \\ &= \sum_{i=1}^n x_i w_i + b \end{aligned}$$

$$y = f(\sum w_i x_i + b)$$

Let's reformulate the I/Ps as a vector

$$x = [x_1, x_2, x_3, \dots, x_n]$$

$$\omega = [\omega_1, \omega_2, \omega_3 \dots \omega_n]$$

$$\text{Net I/p} (y_{in}) = X \cdot \omega$$

The O/p of the neuron is given by

$$y = f(X \cdot \omega + b)$$

biases are the scalar values added to the input to ensure that at least a few nodes per layer are activated regardless of signal strength. It allows learning to happen by giving the n/w action in the event of low signal. It is used to calculate the o/p response of a neuron function. The sum of a weighted i/p signal is applied to an activation function.

Advantages

a) It computes faster, speed up training, It helps n/w to use useful info & ignore irrelevant info. There are some linear and non linear activation functions.

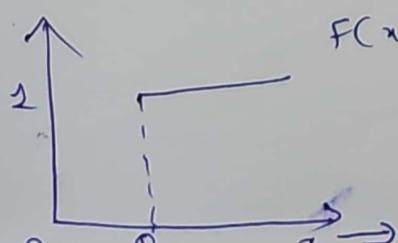
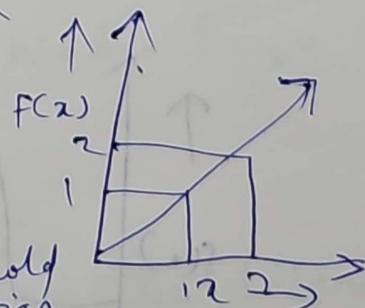
① Identity function

$$f(x) = x \text{ for all } x$$

Step function

② Binary step function

→ Also called threshold function

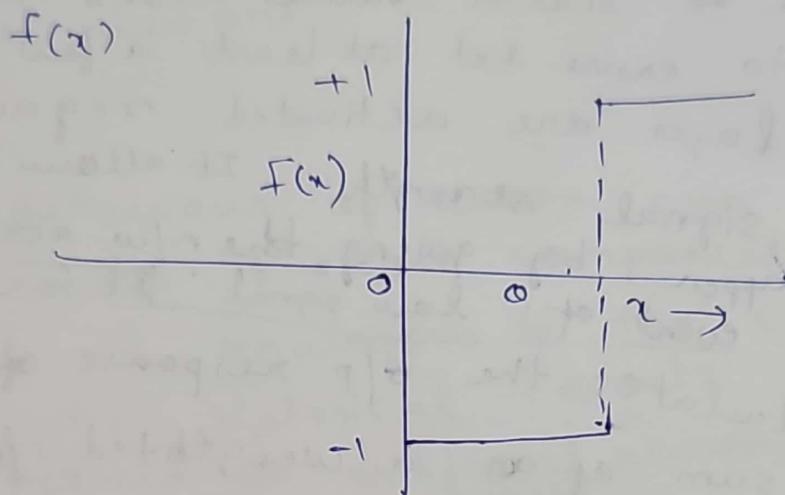


$$f(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{o/w} \end{cases}$$

Bipolar Step function

signum function,

$$f(x) = \begin{cases} +1 & \text{if } x \geq 0 \\ -1 & \text{if } x < 0 \end{cases}$$



③ Sigmoid function

there are 2

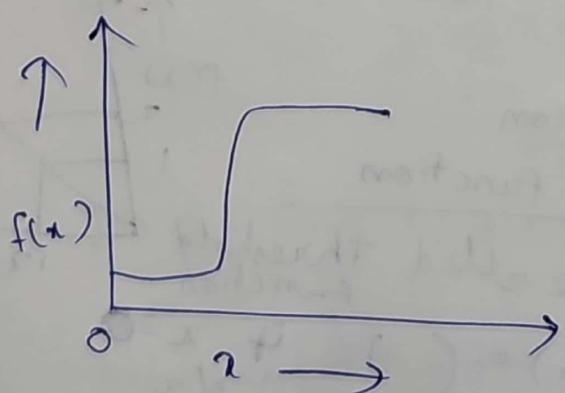
① Binary sigmoidal function - ranges b/w 0 & 1

② Bipolar

IA is ranges b/w 0.1 * 1.0

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

$$x = [-0.5, 0, 0.5, 0.7]$$



- Adv
- easy to train on small dataset
 - reduce extreme values or outliers in data without removing

Dis

- vanishing gradient problem
- Not zero centred
- It is used in O/P
(not an hidden layer)

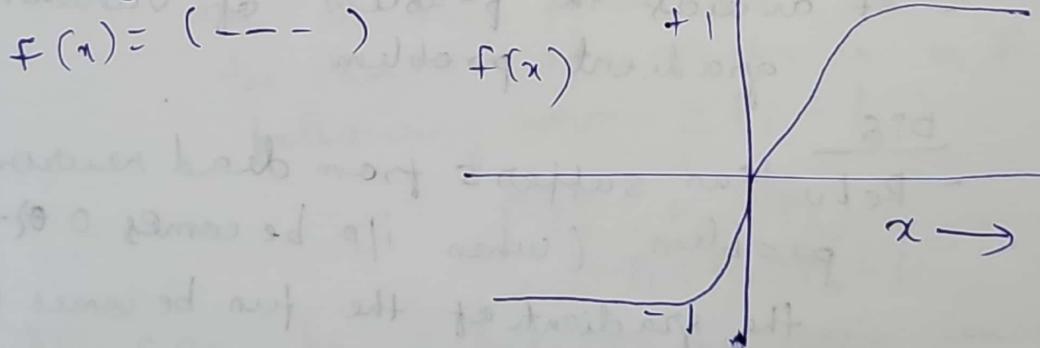
Bipolar sigmoidal fun

tangent hyperbolic function
(tanh)

→ ranges $b/\omega + 1 \approx -1$

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

$$x = [-1, -0.75, -0.5, 0, 0.25, 0.5, 0.75, 1]$$



Adv

- It is Θ -centred
- It can be used both in hidden layer & O/P layer

dis

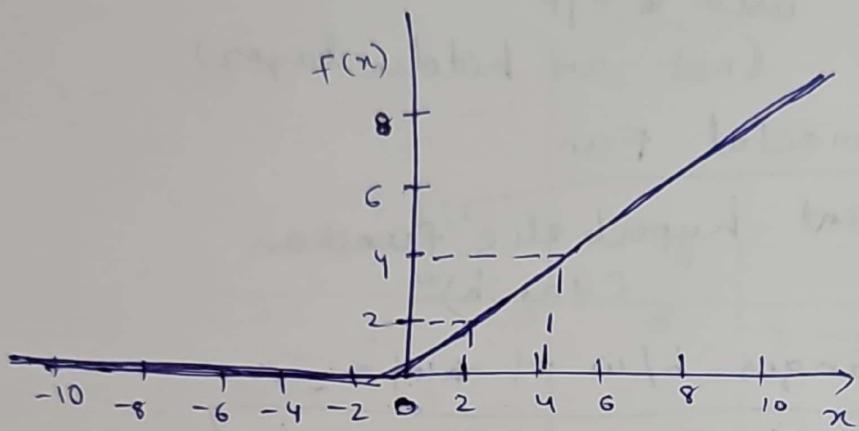
- vanishing gradient problem
- It is hard to train on smaller datasets

⑨ ReLU function

Rectified linear Unit

$$\text{ReLU}(x) = f(x) = \max(0, x)$$

ReLU fun always gives 0 for -ve numbers



- It is computationally efficient
 - (It allows the n/w to converge very quickly)
- It avoids the problem of vanishing gradient problem.

Dis

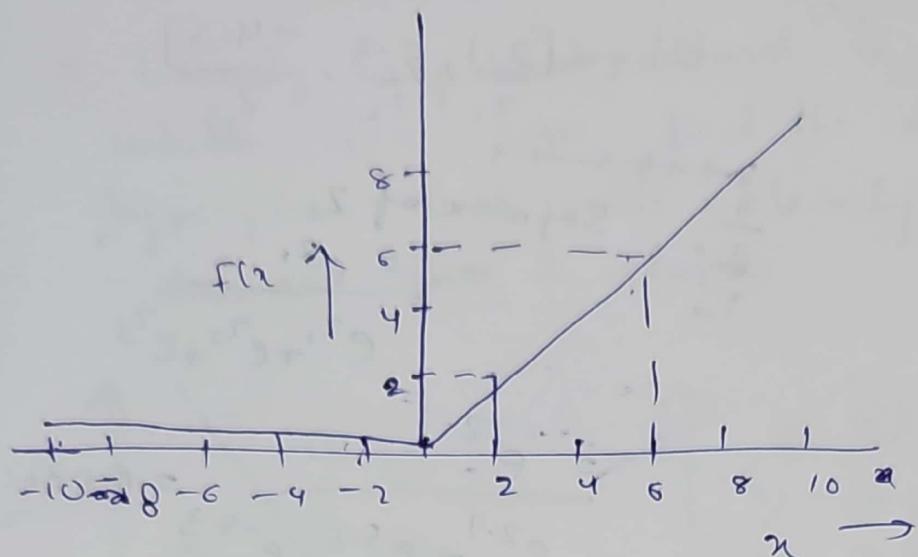
- ReLU fun suffers from dead neuron problem (when i/p becomes 0 or -ve, the gradient of the fun becomes 0). So the n/w can't perform back propagation & can't learn.

⑤ Leaky ReLU

Dead neuron problem can be solved by leaky ReLU

- when $x < 0$, the fun will have a small neg slope

$$f(x) \begin{cases} x & \text{if } x < 0 \\ 0.01x & \text{if } x \geq 0 \end{cases}$$



⑥ Softmax fun

MNIST digit dataset

sometimes we want our o/p vector \mathbf{o} to be a probability dist over a set of mutually exclusive levels.

using a probability dist gives us a better idea how confident we are in our prediction. As a result the desired o/p vector is in the following form. $\sum_i p_i = 1$

$$[p_1, p_2, p_3]$$

$$p_1 + p_2 + p_3 = 1$$

This can be achieved by ~~using~~ using a special output layer ~~as~~ called Softmax fun.

The o/p of a neuron in a softmax layer depends on the o/p of all other neurons in its layer. this is bcoz you require sum of all o/p's equal to 1

$$y_i = \frac{e^{z_i}}{\sum e^{z_i}}$$

this function is mainly used for multi-class classification data problems.

$$z = \left[\frac{2.1}{z_1}, \frac{5.5}{z_2}, \frac{-4.3}{z_3} \right]$$

$$\begin{array}{c} 2.1 + 5.5 - 4.3 \\ \hline \text{Sum of } z_i \\ \hline 7.6 \\ \hline \end{array} \quad \text{Softmax of } z_1 \\ = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$= \frac{e^{2.1}}{e^{2.1} + e^{5.5} + e^{-4.3}} = 0.0327$$

$$\cancel{\text{Sum of } z_2} = \frac{e^{5.5}}{e^{2.1} + e^{5.5} + e^{-4.3}} =$$

$$\cancel{\text{Sum of } z_3} = \frac{e^{-4.3}}{e^{2.1} + e^{5.5} + e^{-4.3}} = 0.9676$$

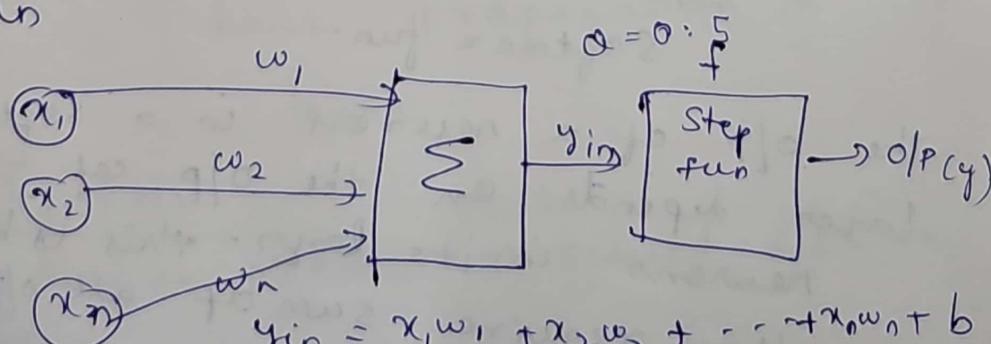
~~highest~~ 96.1%

$$= \frac{e^{-4.3}}{e^{2.1} + e^{5.5} + e^{-4.3}} = 5.3657 \times 10^{-5}$$

The perceptron Model

linear model used for binary classification

This model uses step fun as activation fun

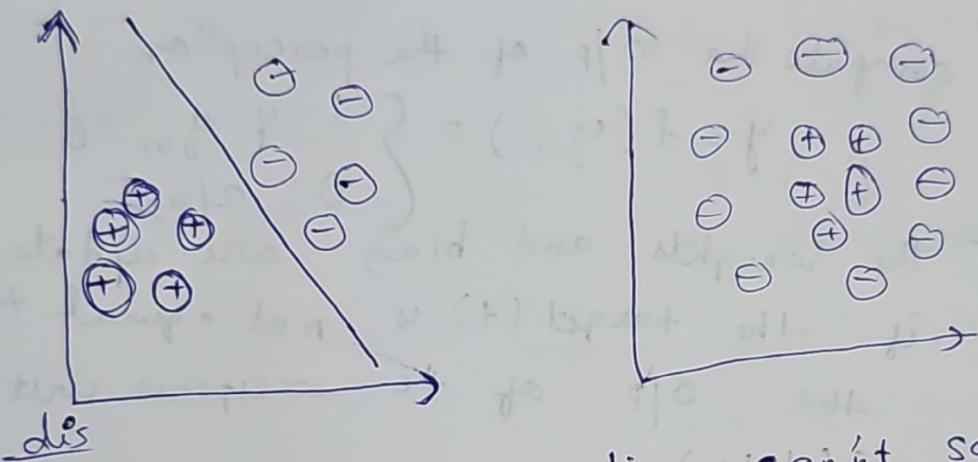


$$y_{in} = x_1 w_1 + x_2 w_2 + \dots + x_n w_n + b$$

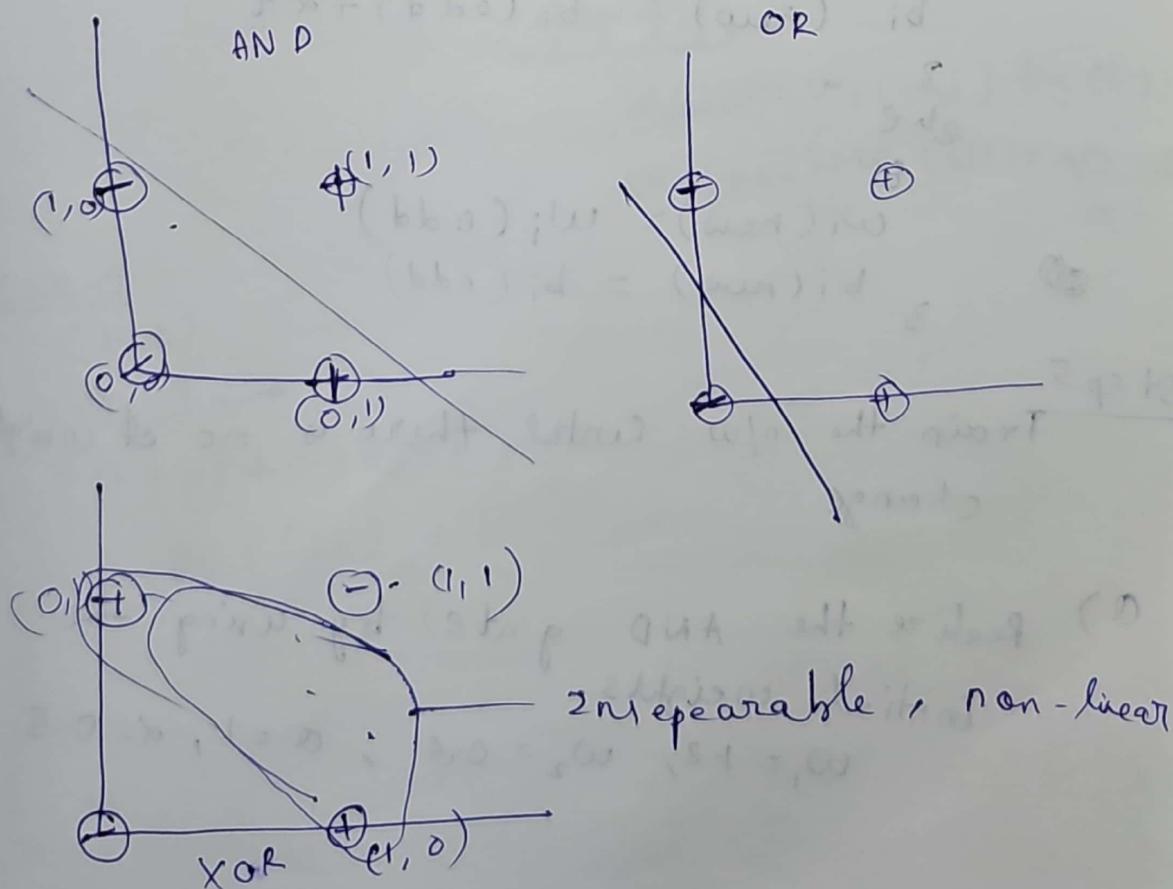
$$y = f(y_{in}) = \begin{cases} \sum_{i=1}^n x_i w_i + b & y_{in} \geq 0 \\ 0 & y_{in} < 0 \end{cases}$$

The basic perceptron model is mainly used for linearly separable dataset.

A linearly separable dataset is one for which we can find the values of hyperplane that will clearly divide the dataset into 2 classes.



A basic ~~perception~~ can't solve non-linear problem. (dataset that are not linearly separable).



non-separable, non-linear

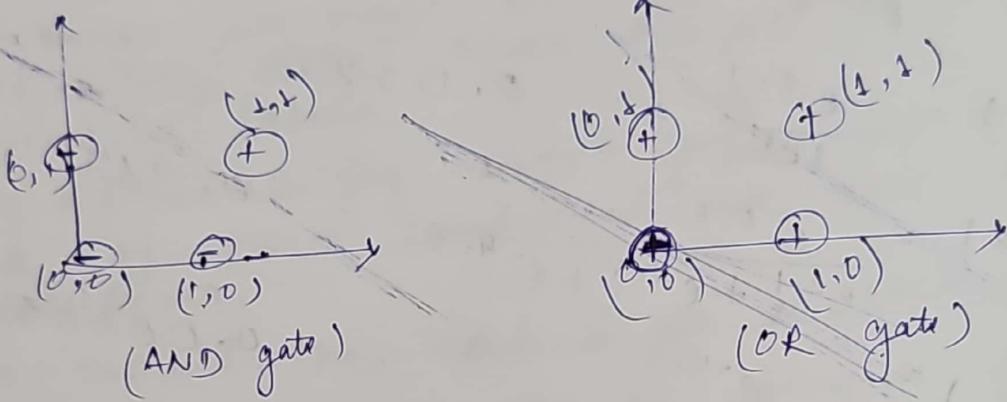
* The basic perception model is mainly used for linearly separable datasets.

23/09/25

Linearly separable dataset is one for which we can find values of hyperplane that will clearly divide the dataset into two classes.
e.g.: - OR / AND gate

Limitations

The basic perception model cannot solve non-linear problem (dataset that are not linearly separable)



But not XOR gate (i.e. we cannot draw a line that will separate the two class)

→ The Training algorithm Perception

Step 1 Initialise the weights, bias & learning rate (α)

Step 2 Compute the o/p of the response unit.
$$y_{in} = \sum w_i x_i + b$$

Step 3 Compute the o/p of the perceptron

$$y = f(y_{in}) = \begin{cases} 1 & \text{if } y_{in} > 0 \\ 0 & \text{otherwise} \end{cases}$$

Step 4 The weights and bias are updated if the target (t) is not equal to the o/p of the response unit (y).

$$\left. \begin{array}{l} w_i(\text{new}) = w_i(\text{old}) + \alpha(t-y) r_i \\ b_i(\text{new}) = b_i(\text{old}) + \alpha t \end{array} \right\} \text{use } \left. \begin{array}{l} w_i(\text{new}) = w_i(\text{old}) \\ b_i(\text{new}) = b_i(\text{old}) \end{array} \right\} \text{implies}$$

Step-5 Train the rule until there is no weight change.

Q Realize the AND gate by using the initial weights
 $w_1 = 1.2, w_2 = 0.6, \theta_{\text{threshold}} = 1, \alpha = 0.5$.

Soln

r_1	r_2	t
0	0	0
0	1	0
1	0	0
1	1	1

case 1 ($r_1 = 0, r_2 = 0$)

$$y_{\text{in}} = w_1 r_1 + w_2 r_2 \\ = 1.2 \times 0 + 0.6 \times 0$$

$$\therefore y = f(y_{\text{in}}) = 0 \\ y = t, \quad (\text{no wt. change})$$

case 2 ($r_1 = 0, r_2 = 1$)

$$y_{\text{in}} = 1.2 \times 0 + 0.6 \times 1 \\ = 0.6$$

$$y = f(y_{\text{in}}) = 0$$

$\therefore y = t$
 \Rightarrow no weight change

case 3 ($r_1 = 1, r_2 = 0$)

$$y_{\text{in}} = 1.2 \times 1 + 0.6 \times 0 \\ = 1.2$$

$$y = f(y_{\text{in}}) = 1$$

$$\therefore y \neq t$$

(as $t = 1$)

so net. change is required.

$$\therefore \text{new } w_1(\text{new}) = w_1(\text{old}) + \alpha(t-y)r_1 \\ w_2(\text{new}) = w_2(\text{old}) + \alpha(t-y)r_2$$

$$\begin{aligned} w_1(\text{new}) &= 1 \cdot 2 + 0.5(0-1) \cdot 1 \\ &= 1 \cdot 2 - 0.5 = 0.7 \\ w_2(\text{new}) &= 0.6 + 0.5(0-1) \cdot 0 \\ &= 0.6 \end{aligned}$$

Now,

$$\begin{aligned} y_{\text{in}} &= w_1 n_1 + w_2 n_2 \\ &= 0.7 \times 1 + 0.6 \times 0 \\ &= 0.7 \end{aligned}$$

$$y = f(y_{\text{in}}) = 0$$

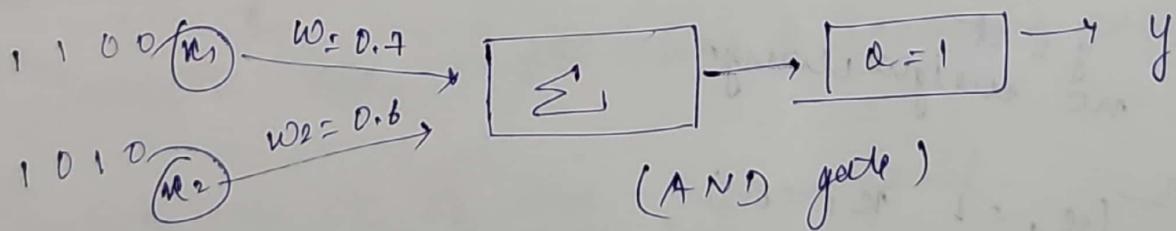
$\therefore y = t$
so no wt. change.

case-4 $n_1 = 1, n_2 = 1$

$$\begin{aligned} y_{\text{in}} &= w_1 n_1 + w_2 n_2 \\ &= 0.7 \times 1 + 0.6 \times 1 \\ &= 1.3 \end{aligned}$$

$$y = f(y_{\text{in}}) = 1$$

$y = t$
so no wt. change.



Q Realize the OR gate by using the initial weight
 $w_1 = 0.6, w_2 = 1.1, d = 1, \alpha = 0.5$

for

n_1	n_2	t
0	0	0
0	1	1
1	0	1
1	1	1

Case I $x_1 = 0, x_2 = 0$

$$\begin{aligned}y_{in} &= w_1 x_1 + w_2 x_2 \\&= 0.6 \cdot 0 + 1.1 \cdot 0 \\&= 0\end{aligned}$$

$$y = f(y_{in}) = 0$$

$$\therefore y = t$$

so no wt. change

Case II $x_1 = 1, x_2 = 0$

$$\begin{aligned}y_{in} &= w_1 x_1 + w_2 x_2 \\&= 0.6 \cdot 1 + 1.1 \cdot 0 \\&= 0.6\end{aligned}$$

$$y = f(y_{in}) = 0$$

$$\text{but } t = 1$$

$$\therefore y \neq t$$

so wt. change req.

$$w_1(\text{new}) = w_1(\text{old}) + \alpha(t - y) x_1$$

$$\Rightarrow w_1(\text{new}) = 0.6 + 0.5(1 - 0) \cdot 1 \\= 0.6 + 0.5$$

$$w_1(\text{new}) = 1.1$$

$$\begin{aligned}w_2(\text{new}) &= w_2(\text{old}) + \alpha(t - y) x_2 \\&= 1.1 + 0.5(1 - 0) \cdot 0 \\&= 1.1\end{aligned}$$

$$w_2(\text{new}) = 1.1$$

New

$$\begin{aligned}y_{in} &= 1.1 \cdot 1 + 1.1 \cdot 0 \\&= 1.1\end{aligned}$$

$$y = f(y_{in}) = 1$$

$$\therefore y = t$$

so no wt. change.

Case III $x_1 = 0, x_2 = 1$

$$\begin{aligned}y_{in} &= w_1 x_1 + w_2 x_2 \\&= 0.6 \cdot 0 + 1.1 \cdot 1 \\&= 1.1\end{aligned}$$

$$y = f(y_{in}) = 1$$

$$y = t$$

so no wt. change

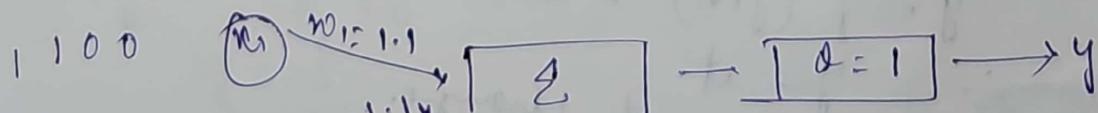
Case IV $x_1 = 1, x_2 = 1$

$$\begin{aligned}y_{in} &= w_1 x_1 + w_2 x_2 \\&= 1.1 \cdot 1 + 1.1 \cdot 1 \\&= 2.2\end{aligned}$$

$$y = f(y_{in}) = 1$$

$$y = t$$

so no wt. change.



Q Implement OR gate by using the initial weights
 $w_1 = 0$ $w_2 = 0$ $\alpha = 0.5$ $t = 1$

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

Case I

$$w_1 = 0 \quad w_2 = 0$$

$$\begin{aligned} y_{in} &= w_1 x_1 + w_2 x_2 \\ &= 0 + 0 + 0 + 0 \\ &= 0 \end{aligned}$$

$$y = f(y_{in}) = 0$$

$$\therefore y = t$$

so no wt. change

case II $w_1 = 0 \quad w_2 = 1$

$$\begin{aligned} y_{in} &= w_1 x_1 + w_2 x_2 \\ &= 0 + 0 + 0 + 1 \\ &= 1 \end{aligned}$$

$$y = f(y_{in}) = 0$$

$$\text{but } t = 1$$

$$\therefore y \neq t$$

net change req.

$$w_1(\text{new}) = w_1(\text{old}) + \alpha(t - y)x_1$$

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

$$= 0$$

$$w_2(\text{new}) = 0 + 0.5(1 - 0) * 1$$

$$= 0.5$$

Now

$$\begin{aligned} y_{in} &= 0 + 0.5 * 1 \\ &= 0.5 \end{aligned}$$

$$y = f(y_{in}) = 0.$$

$y \neq t$ net. change req.

$$w_1(\text{new}) = 0 + 0.5(1 - 0)0 = 0$$

$$w_2(\text{new}) = 0.5 + 0.5(1 - 0) * 1 = 1$$

Now

$$y_{in} = 0 + 1 * 1 = 1$$

$$y = f(y_{in}) = 1$$

$$\therefore y = t$$

so no wt. change.

Case III $x_1 = 1 \quad x_2 = 0$

$$y_{in} = w_1 x_1 + w_2 x_2 \\ = 0 * 1 + 1 * 0 \\ = 0$$

$$y = f(y_{in}) = 0$$

but $t = 1$
so wt change req.

New

$$w_1(\text{new}) = 0 + 0.5(1-0)^{*}1$$

$$w_1(\text{new}) = 0.5$$

$$w_2(\text{new}) = 1 + 0.5(1-0)^{*}0 \\ = 1$$

$$y_{in} = w_1 x_1 + w_2 x_2 \\ = 0.5 * 1 + 1 * 0 \\ = 0.5$$

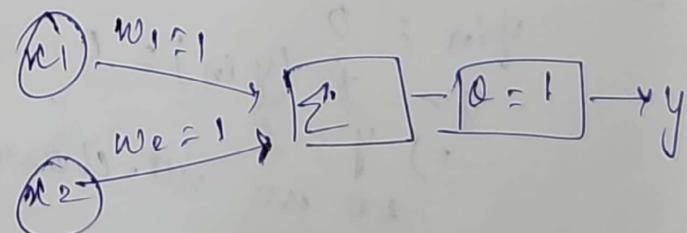
No

$$w_1(\text{new}) = 0.5 + 0.5(1-0)^{*}1 \\ = 1$$

$$w_2(\text{new}) = 1 + 0.5(1-0)^{*}0 \\ = 1$$

$$y_{in} = 1 * 1 + 1 * 0 \\ = 0 +$$

$$y = f(y_{in}) = 1 \\ \therefore \underline{y = t}$$



Case IV $x_1 = 1 \quad x_2 = 1$

$$y_{in} = 1 * 1 + 1 * 1 \\ = 2$$

$$y = f(y_{in}) = 1$$

$$\underline{y = t}$$

24/09/25

 $w_1 = 0, w_2 = 0, R = 0.5, t = 1$

Using AND gate

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

case I

$x_1 = 0, x_2 = 0$

$$\begin{aligned}y_{in} &= w_1 x_1 + w_2 x_2 \\&= 0 \cdot 0 + 0 \cdot 0 \\&= 0\end{aligned}$$

$y = f(y_{in}) = 0$

$\therefore y = t$
so no net change

case II

$x_1 = 0, x_2 = 1$

$$\begin{aligned}y_{in} &= w_1 x_1 + w_2 x_2 \\&= 0 \cdot 0 + 0 \cdot 1 \\&= 0\end{aligned}$$

$y = f(y_{in}) = 0$

$\therefore y = t$
so no net change

case III

$x_1 = 1, x_2 = 0$

$y_{in} = 0$

$y = f(y_{in}) = 0$

$\therefore y = t$
so no net change

case IV $x_1 = 1, x_2 = 1$

$y_{in} = 0$

$y = f(y_{in}) = 0$

but $y_{out} = 1$
so $y \neq t$.

so net. change req.

$$\begin{aligned}w_1(\text{new}) &= w_1(\text{old}) + \\&\quad \times (t - y) x_1 \\&\Rightarrow w_1(\text{new}) = 0 + 0.5(1 - 0) \\&= 0.5\end{aligned}$$

$$\begin{aligned}w_2(\text{new}) &= 0 + 0.5(1 - 0) \\&= 0.5\end{aligned}$$

$$\begin{aligned}y_{in} &= w_1 x_1 + w_2 x_2 \\&= 0.5 \cdot 1 + 0.5 \cdot 1 \\&= 1\end{aligned}$$

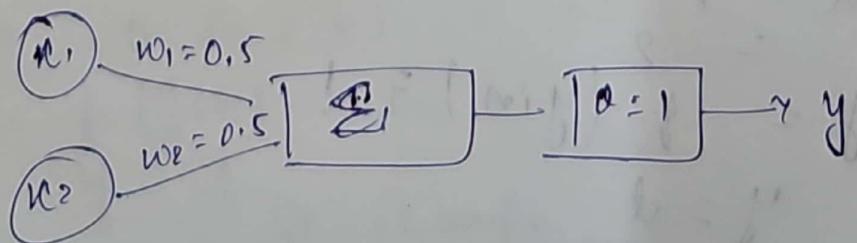
$\therefore y = f(y_{in}) = 1$

$\& t = 1$

$\therefore y = t$

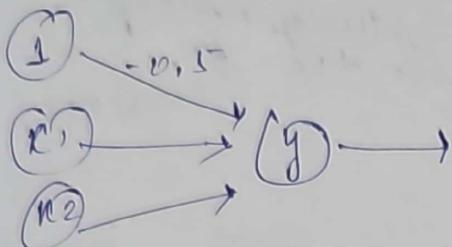
∴ no net. change req.

Q&P



Q Use a perceptron to implement OR gate with the following info. The line equation is $x_1 + x_2 - 0.5 = 0$ and the step activation fn is as follows. $f(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$

~~sofn~~ $x_1 + x_2 - 0.5 = 0$



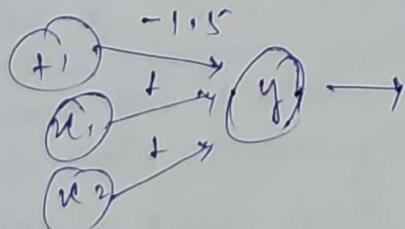
$$y = f(-0.5 + x_1 + x_2)$$

x_1	x_2	$x_1 \cup x_2$	$(-0.5 + x_1 + x_2)$	y
0	0	0	-0.5	0
0	1	1	0.5	1
1	0	1	0.5	1
1	1	1	0.5	1

Q Use perceptrons to implement AND gate. The line equation is $x_1 + x_2 - 1.5 = 0$ and the step activation fn is same as above. $f(z) = \begin{cases} 0 & z < 0 \\ 1 & z \geq 0 \end{cases}$

~~Ans~~

x_1	x_2	$\frac{k}{0}$	$\frac{z = -1.5 + x_1 + x_2}{-1.5}$	$y = f(z)$
0	0	0	-1.5	0
0	1	0	-0.5	0
1	0	1	-0.5	0
1	1	1	0.5	1

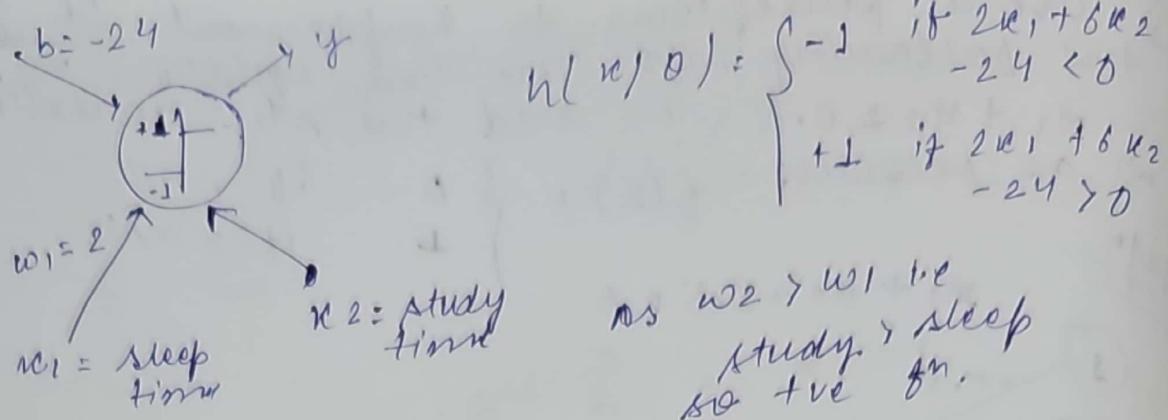


Q Build a perceptron to predict the student success based on time spent on study and sleep.

let $x_1 = \text{sleep time}$
 $x_2 = \text{study time}$

$$w_1 = 2 \quad w_2 = 6 \quad b = -24$$

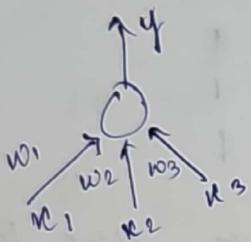
which feature contributes more towards the prediction Justify your answer.



06/10/25

Linear Neurons and Their Limitations

$$y(i) = x_1 w_1(i) + x_2 w_2(i) + x_3 w_3(i)$$



Linear Neurons are easy to compute. But they have serious limitations. Because linear neurons can be expressed as a network with hidden layers. hidden layers enable us to learn important features from the IIP layer.

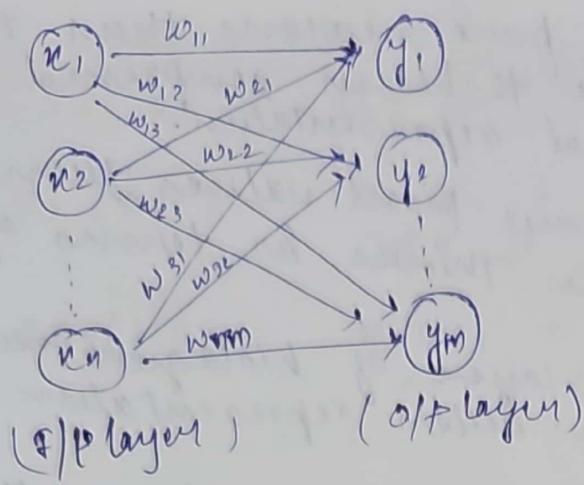
In other words, in order to learn complex relationships, we need to use neurons that employ some sort of non-linearity.

Feed Forward Neural Network

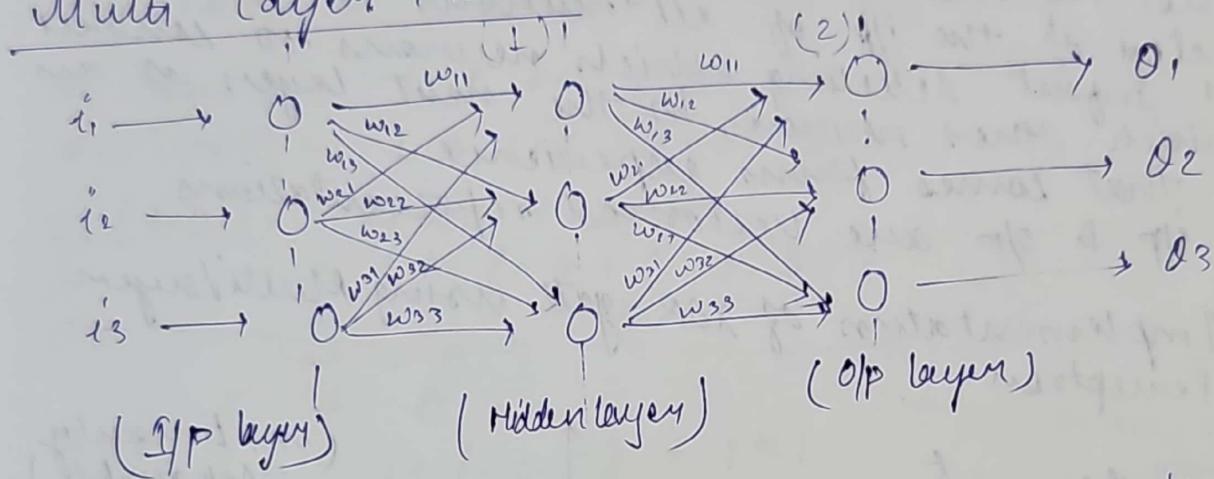
1. Single layer FNN
2. Multi layer FNN

Single Layer FNN

Despite two layers the nne is termed as single layer because all the computation being performed in the O/P layer and the IIP layer doesn't perform any computation. Hence the nne is known as single layer feed forward Neural Network.



Multi Layer FNN



This network is made up of multiple layers i.e. input layer, output layer and hidden layer.

The layers that are sandwiched b/w input layer & output layer are called hidden layer.

In this layer most of the magic happens when the neural network tries to solve problems.

e.g. In handwriting digit recognition, we have to spend a lot of time identifying useful features. But this hidden layers automate this process for us. That means the features are automatically extracted from the data through hidden layers.

(Diag :- This is a simple Feed Forward Neural Network with 3 layers and 3 neurons per layer)

Though in this example every layer has the same number of neurons this is neither necessary nor recommended.

More often hidden layers have pure neurons than the I/P layers. To force the NN to learn compressed representation of the original representation.

eg:- while obtaining raw pixel values from our surroundings, our brain thinks in terms of edges and contours.

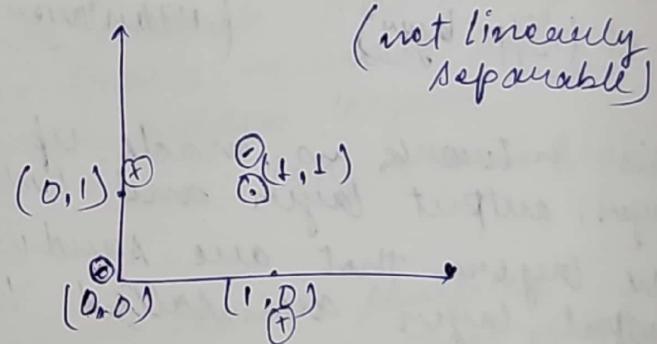
This is because the hidden layer of biological neurons force us to come up with better representation for everything we perceive.

It is not required that every neuron has its output connected to the I/P of all neurons in the next layer. In fact selecting which neurons to connect to which other neurons in the next layer is an art that comes from experience.

The I/P & O/P are vectorized representations.

→ Implementation of XOR gate using Multilayer Perceptron

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

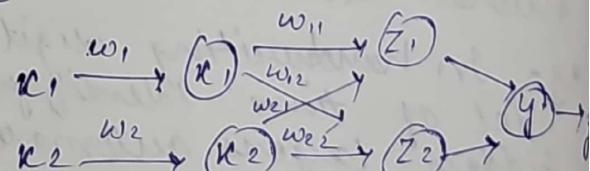


$$y = 01 + 10 \dagger = \frac{\bar{x}_1 x_2}{z_1} + \frac{x_1 \bar{x}_2}{z_2}$$

$$z_1 = \bar{x}_1 \cdot x_2 \quad (\text{1st f}^n)$$

$$z_2 = x_1 \cdot \bar{x}_2 \quad (\text{2nd f}^n)$$

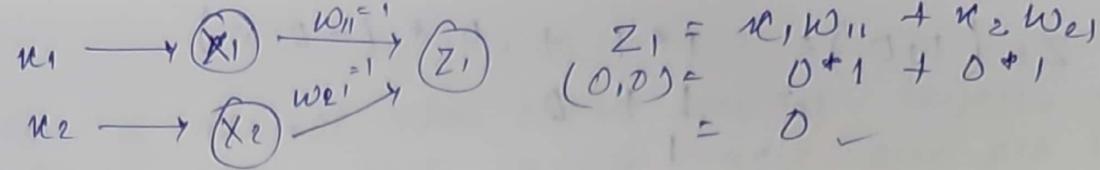
$$y = z_1 + z_2 \quad (\text{3rd f}^n)$$



Assume initial weights $w_{11} = 1$, $w_{12} = 1$, $\theta = 1$, $\eta = 0.5$

$$\text{1st f}^n \quad z_1 = \bar{x}_1 \cdot x_2$$

x_1	\bar{x}_1	x_2	$z_1(t)$
0	1	0	0
0	1	1	1
1	0	0	0
1	0	1	0



$$z_1(0,1) = \frac{\kappa_1 w_{11} + \kappa_2 w_{21}}{0+1 + 1+1} = 1 \checkmark$$

$$z_2(1,0) = \frac{\kappa_1 w_{11} + \kappa_2 w_{21}}{1+1 + 0+1} = 1 \quad (\text{mismatched, need to change})$$

$$w_{11}(\text{new}) = w_{11}(\text{old}) + \eta(t-y) \kappa_1$$

$$w_{21}(\text{new}) = w_{21}(\text{old}) + \eta(t-y) \kappa_2$$

$$w_{11}(\text{new}) = 1 + 1.5(0-1) = -0.5$$

$$w_{21}(\text{new}) = 1 + 1.5(0-1) = 1$$

$$z_1(1,0) = \kappa_1 w_{11} + \kappa_2 w_{21}$$

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

$$= -0.5 + 0$$

$$= -0.5$$

$$= 0 \checkmark$$

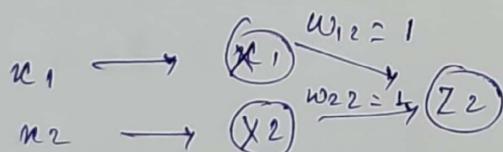
$$z_2(1,1) = \frac{1 + (-0.5) + 1+1}{-0.5 + 1} = 0.5 \approx 0$$

$$w_{11} = -0.5$$

$$w_{21} = 1$$

and B^n

\textcircled{x}_1	\textcircled{x}_2	$\bar{\textcircled{x}}_2$	$z_2(t)$
0	0	1	0
0	1	0	0
1	0	1	1
1	1	0	0



$$z_1(0,0) = \frac{\kappa_1 w_{12} + \kappa_2 w_{22}}{0+1 + 0+1} = 0 \checkmark$$

$$z_2(0,1) = \frac{\kappa_1 w_{12} + \kappa_2 w_{22}}{0+1 + 1+1} = 1 \quad (\text{mismatched})$$

$$w_{12} \text{ new} = 1 + 1.5(0-1)0$$

$$w_{22} \text{ new} = 1 + 1.5(0-1)1$$

$$z_2(0,1) = \frac{0+1 + 1+(-0.5)}{-0.5 + 1} \approx 0 \checkmark$$

$$z_1(1,0) = 1 \cdot 1 + 0 + (-0.5)$$

$$= 1 + 0$$

$$= 1$$

$$z_2(1,1) = 1 \cdot 1 + 1 + (-0.5)$$

$$= 1 - 0.5$$

$$= 0.5 \approx 0 \vee$$

$$w_{12} = 1$$

$$w_{22} = -0.5$$

3rd f^m

$$y = z_1 + z_2$$

$$= z_1 N_1 + z_2 N_2$$

$\frac{z_1}{0}$	$\frac{z_2}{0}$	$\frac{x}{0}$
0	1	1
1	0	1
1	1	1

$$y(0,0) = 0 \cdot 1 + 0 \cdot 1$$

$$= 0$$

$$y(0,1) = 0 \cdot 1 + 1 \cdot 1$$

$$= 1$$

$$y(1,0) = 1 \cdot 1 + 0 \cdot 1 \quad \therefore \text{all } y \text{ condition matched}$$

$$= 1$$

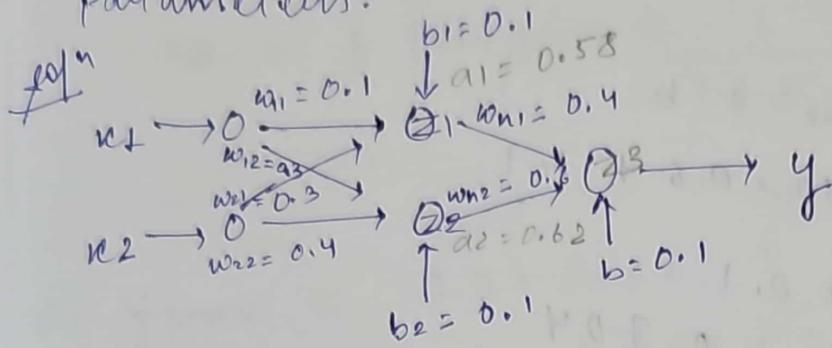
$$y(1,1) = 1 \cdot 1 + 1 \cdot 1$$

$$= 2 \approx 1 \pm 1$$

No need of weight change.

08/10/25

(i) Consider a neural net with one IP layer, one hidden-layer with two neurons and one OP layer with one neuron. Draw the architecture of neural net with given data and calculate the total no. learnable parameters.



$$\begin{aligned} \text{no. of learnable parameters} &= 9 \\ a_1 &= 0.5 \\ u_2 &= 0.6 \end{aligned}$$

(ii) calculate the OP of the net in the forward propagation.

(iii) calculate the error at OP layer for the actual OP $y = 1$.

assume neurons using sigmoid activation, f^n.

$$\begin{aligned} z_1 &= u_1 w_{11} + u_2 w_{22} + b_1 \\ &= 0.5 \times 0.1 + 0.6 \times 0.3 + 0.1 \\ &= 0.05 + 0.18 + 0.1 \\ &= 0.23 + 0.1 = \underline{\underline{0.33}} \end{aligned}$$

$$a_1 = f(z_1) = \frac{1}{1 + e^{-z_1}} \quad (\because \text{sigmoid } f^n)$$

$$= \frac{1}{1 + e^{-0.33}} = \frac{e^{0.33}}{1 + e^{0.33}} = \underline{\underline{0.58}}$$

$$\begin{aligned} z_2 &= u_1 \cdot w_{12} + u_2 \cdot w_{22} + b_2 \\ &= 0.5 \times 0.3 + 0.6 \times 0.4 + 0.1 \\ &= 0.15 + 0.24 + 0.1 \\ &= 0.39 + 0.1 = \underline{\underline{0.49}} \end{aligned}$$



$$a_2 = f(x_2) = \frac{1}{1 + e^{-x}}$$

$$= \frac{1}{1 + e^{-0.49}} = \frac{e^{0.49}}{1 + e^{0.49}} = 0.619$$

$$\approx 0.62$$

$$Z_3 = a_1 w_{01} + a_2 w_{02} + b$$

$$= 0.58 \times 0.4 + 0.62 \times 0.6 + 0.1$$

$$= 0.232 + 0.372 + 0.1$$

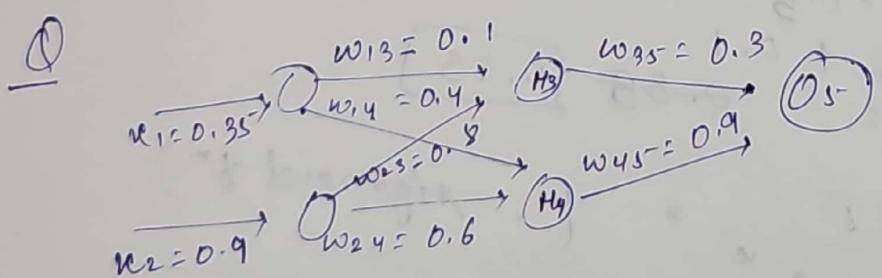
$$\approx 0.604 + 0.1 = \underline{\underline{0.704}}$$

$$y = \frac{1}{1 + e^{-0.704}} = \underline{\underline{0.668}} \approx 0.67$$

$$\therefore y_{predicted} = \underline{\underline{0.67}}$$

Also, error = $y_{actual} - y_{pred.}$

$$= 0.33 - \underline{\underline{0.67}} = \underline{\underline{-0.34}}$$



Given sigmoid fn
actual OP = 0.5

Q2

$$H_3 = x_1 w_{13} + x_2 w_{23}$$

$$= 0.35 \times 0.1 + 0.9 \times 0.8$$

$$= 0.035 + 0.720$$

$$= \underline{\underline{0.755}}$$

$$a_1 = f(H_3) = \frac{1}{1 + e^{-0.755}} = \underline{\underline{0.68}}$$

$$\begin{aligned}
 H_4 &= w_1 w_{1,4} + w_2 w_{2,4} \\
 &= 0.35 + 0.4 + 0.9 + 0.6 \\
 &= 0.140 + 0.54 \\
 &= 0.680
 \end{aligned}$$

$$a_2 = f(H_4) = \frac{1}{1 + e^{-0.680}} = 0.663$$

$$\begin{aligned}
 0.5 y_{\text{pred}} &= 0.3 + 0.68 + 0.9 + 0.66 \\
 &= 0.204 + 0.594 \\
 &= 0.798
 \end{aligned}$$

$$\text{error } = \underline{y} - \hat{y} = \frac{1}{1 + e^{-0.798}} = 0.689 \approx 0.69$$

$$\begin{aligned}
 \text{error} &= \frac{|y_{\text{actual}} - y_{\text{pred}}|}{0.5 - 0.69} \\
 &= 0.19 \approx 19\%
 \end{aligned}$$

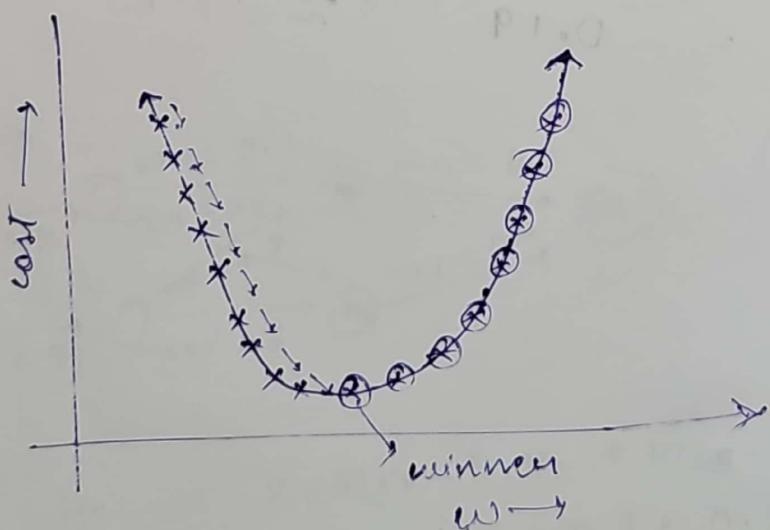
~~RECORDED~~

13/10/25 Gradient Descent & Delta Rule :-

Perceptron rule finds a successful weight vector when the training examples are linearly separable.

But it fails to converge if the examples are not linearly separable.

- The Delta Rule is designed to overcome this difficulty.
 - If the training examples are not linearly separable, the delta rule converges towards a fixed best fit approximation to the target concept.
 - The key idea behind the delta rule is to use gradient descent.
To search the hypothesis space of possible weight vector to find the ~~vector~~ that best fits the training example.
- This rule is imp. because this provides the basis for back propagation algo.



The Delta training rule is best understood by considering the task of training a single unit of a three-layer perceptron, i.e., a linear unit for which the output is given by

$$y = w_0 + w_1 u_1 + w_2 u_2 + \dots + w_n u_n$$

- or $y(\vec{w}) = \vec{w} \cdot \vec{x}$
- In order to derive weight learning rule for two linear units, let us begin by specifying a measure for the training error of a hypothesis, relative to the training examples.

Although there are many ways to define this error, one common measure is

$$E(\vec{w}) = \frac{1}{2} \sum_{i \in N} (t_i - y_i)^2$$

\downarrow
Target calculated (by output of linear unit for i th training unit)

$N \rightarrow$ set of training examples.

Ans:

Q How can we find out the direction of steepest descent along the error surface?

The direction can be found by completing the derivative of E w.r.t \vec{w} and written as

$$\nabla E(\vec{w}) = \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots, \frac{\partial E}{\partial w_n} \right]$$

\therefore the gradient specifies the direction of steepest descent of E . The training rule for gradient descent is:-

$$\vec{w} \leftarrow \vec{w} + \Delta \vec{w}$$

$$\text{where, } \Delta \vec{w} = -\eta \cdot \nabla E(\vec{w})$$

$\eta \rightarrow$ learning rate, which determines the step size in the gradient steepest descent surface.

The $-\eta$ sign is present because we want to move the weight vector in the direction that \downarrow is E .

Ans:

This training rule can also be written in component form :-

$$w_i \leftarrow w_i + \Delta w_i$$

where $\Delta w_i = -\eta \cdot \frac{\partial E}{\partial w_i}$

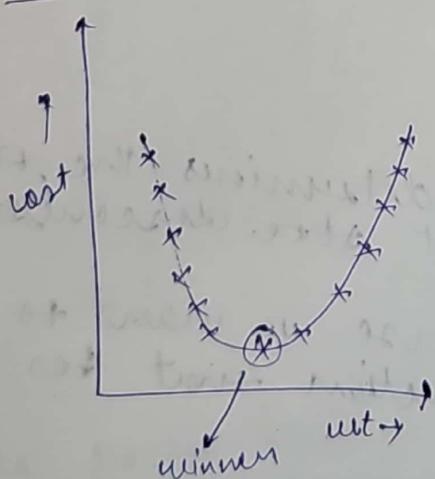
$$E(\vec{w}) = \frac{1}{2} \sum_{i \in N} (t_i - y_i)^2$$

$$\begin{aligned} \frac{\partial E}{\partial w_i} &= \frac{1}{2} \times 2(t_i - y_i) \frac{\partial}{\partial w_i} (t_i - y_i) \\ &= \sum_{i \in N} (t_i - y_i) \left[0 - \frac{\partial}{\partial w_i} (w_i \cdot n_i) \right] \\ &= \sum_{i \in N} (t_i - y_i) (-n_i) \\ \boxed{\frac{\partial E}{\partial w_i}} &= -n_i \sum_{i \in N} (t_i - y_i) \end{aligned}$$

$$\Delta w_i = -\eta \cdot -n_i \sum_{i \in N} (t_i - y_i)$$

$$= \eta n_i \sum_{i \in N} (t_i - y_i)$$

15/10/25 The delta rule and learning rate



$$E(\vec{w}) = \frac{1}{2} \sum_{i \in N} (t_i - y_i)^2$$

$$\text{where } w_i \leftarrow w_i + \Delta w_i$$

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

ideal value for learning rate = $\eta = 0.001$

If the value is too less then it will take more time or to reach converging point.

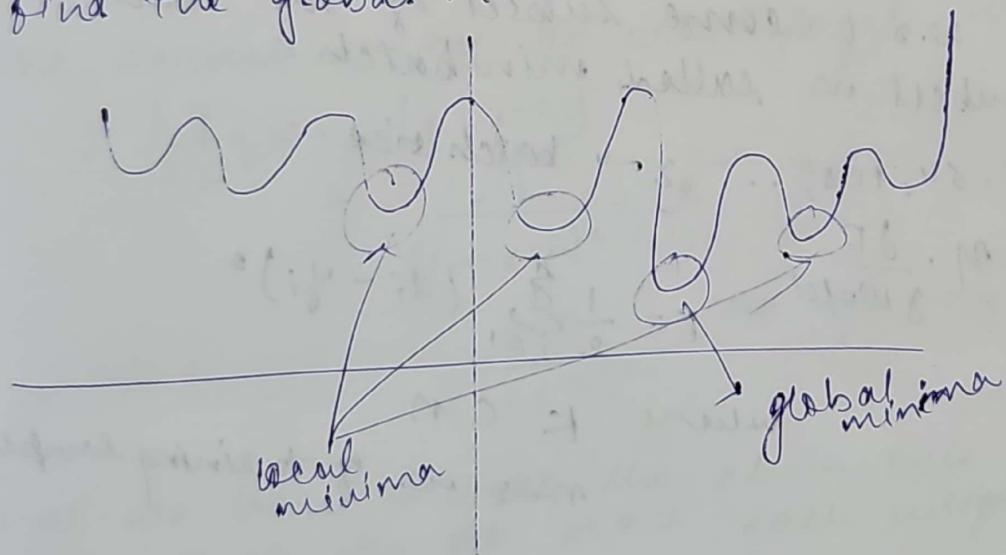
Q What is the effect of learning rate (η) or what will happen when η is too high or too low?

- If the learning rate is too high or too low then convergence is difficult ^{ie} that means we cannot may not come to the min^m point.
- If learning rate is too slow then it will take too long during the training process.

Limitations of Gradient Descent

converging to a local min^m can sometimes be quite slow [it may require min + 100 gradient descent steps]

If there are multiple local min^m in the surface, then there is no guarantee that the procedure will find the global minima.



Ans Stochastic Gradient Descent (SGD) & Minibatch Gradient Descent : -

In gradient descent, we use our entire dataset to compute the error surface and then follow the gradient to take the path of the steepest descent where

$$W_{new} = W_{old} - \eta \frac{\partial E}{\partial W_{old}}$$

$$E = \frac{1}{2} \sum_{i=1}^n (x_i - y_i)^2$$

where $n \rightarrow$ no. of training sample.

- * Hence training gradient descent algo. requires more resource and more calculation and it is very complex.
 - * gradient descent also called as batch gradient descent
- In SGD, the error surface is ~~estimated~~ estimated using single training sample w.r.t single ~~it~~ at each iteration.

$$\left\{ \begin{array}{l} w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial E}{\partial w_{\text{old}}} \\ \text{where, } E = \frac{1}{2} (t - y)^2 \end{array} \right.$$

But the major pitfall of SGD is looking at the error incurred one example at a time may not be a good enough apparent of the error surface.

Minibatch :- In minibatch SGD, we compute the error surface w.r.t some subset of the total dataset. This subset is called minibatch.

Like $\{16, 32, 64, 128, \dots\}$ \rightarrow batch size.

$$\left\{ \begin{array}{l} w_{\text{new}} = w_{\text{old}} - \eta \cdot \frac{\partial E}{\partial w_{\text{old}}} \\ E = \frac{1}{2} \sum_{i=1}^k (t_i - y_i)^2 \end{array} \right.$$

where $k \subset n$

$n \rightarrow$ no. of training sample

training, testing
(80%) (20%)

overfitting:-
If it passed training but failed in testing
e.g. - overfitting

underfitting:-
if failed in training only.

Gradient Descent with Sigmoidal Neurons

Logistic neurons compute their output values from their i/p's.

first neurons computes the weighted sum of its i/p. Then it feeds its logit into its ~~logit~~ i/p fⁿ to compute the final op $y \cdot$ ~~weight~~ ^{i/p}

$$z = \sum_k w_k r_k , \frac{\partial z}{\partial w_k} = r_k$$

$$y = \frac{1}{1+e^{-z}} , \frac{\partial z}{\partial r_k} = w_k$$

for learning we have to compute the gradient of error fⁿ w.r.t weights.
To do so we have to find the derivative of the logit w.r.t i/p and weights.

The derivative of the op w.r.t logit is given by

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

Now we have to use the chain rule to get the derivative of the op w.r.t each weight.

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

Putting all these together, we can now compute the derivatives of error fⁿ w.r.t each weight.

$$\frac{\partial E}{\partial w_k} = \sum_i \frac{\partial E}{\partial y^{(i)}} \cdot \frac{\partial y^{(i)}}{\partial w_k}$$

$$= - \sum_i \eta_k^{(i)} y^{(i)} (1 - y^{(i)}) (t^{(i)} - y^{(i)})$$

$$[\Delta w_k = \sum_i \eta_i \eta_k^{(i)} y^{(i)} (1 - y^{(i)}) (t^{(i)} - y^{(i)})]$$

$$\left\{ \therefore \Delta w = - \eta \frac{\partial E}{\partial w} \right.$$

Compute the derivative of sigmoid f^n $f(x) = \frac{1}{1+e^{-x}}$
at $x = 0$

~~soln~~

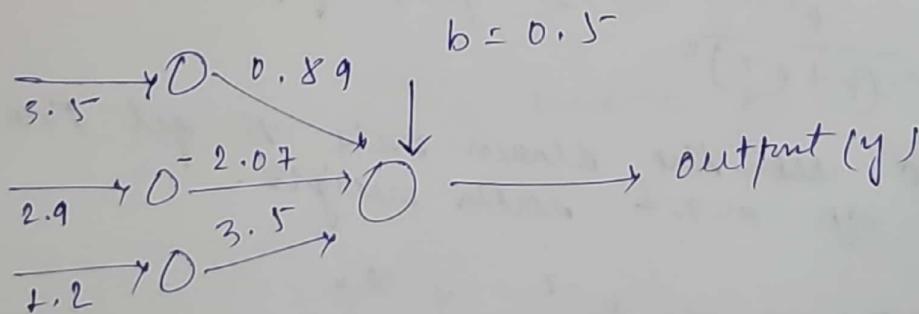
$$y = \frac{1}{1+e^{-x}}$$

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

$$= \frac{\frac{1}{e^0}}{\left(1 + \frac{1}{e^0}\right)^2} = \frac{1}{2(1+1)^2} = \frac{1}{4}$$

$$\left\{ f'(x) = f(x)(1 - f(x)) \right.$$

~~soln~~



(i) Sigmoid(z)

$$z = w_1 u_1 + w_2 u_2 + w_3 u_3 + 0.5$$

$$= 3.5 \cdot 0.89 + 2.9 \cdot -2.07 + 1.2 + 0.5$$

$$= 3.115 + (-6.003) + 4.2 + 0.5$$

$$= 3.115 - 6.003 + 4.2 + 0.5$$

$$= 7.315 - 6.003 + 0.5$$

$$= 1.312 + 0.5 = \boxed{1.812}$$

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

$$y = \frac{1}{1 + e^{-1.812}}$$

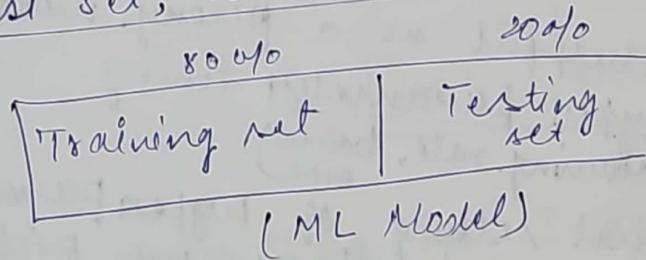
(ii) tanh

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

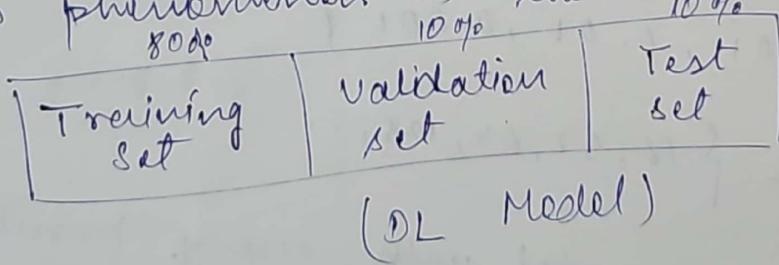
(iii) $\text{ReLU}(z) = \max(0, z)$

$$= \max(0, 1.812) = \underline{\underline{1.812}}$$

→ Test set, validation set and overfitting



Sometimes the model performs well with the training data but does not perform well with the test data. This phenomenon is called Model Overfitting.



29/10/20
In ML we split the data into training set and test set. In DL we include a validation set to prevent overfitting during the training process.

In DL we have to stop the training process as soon as model starts overfitting due to prevent poor generalization.

To do this, we divide our training process into epochs. An epoch is a single iteration over the entire training set.

Eg:- suppose we have a training set of size 'd' and we are doing minibatch AD with batch size 'b', then an epoch would be equal to size ' d/b '.

At the end of each epoch, we have to measure how well our model is generalised. To do this we use an additional validation set. The validation set will tell us how the model does on unseen data.

→ If the accuracy on the training set continuously decreases while the accuracy on the validation set remains the same or less, the training has to stop because the model is overfitting.

The validation set is helpful as a proxy measure of accuracy during hyper parameter tuning.
(learning rate, batch size)

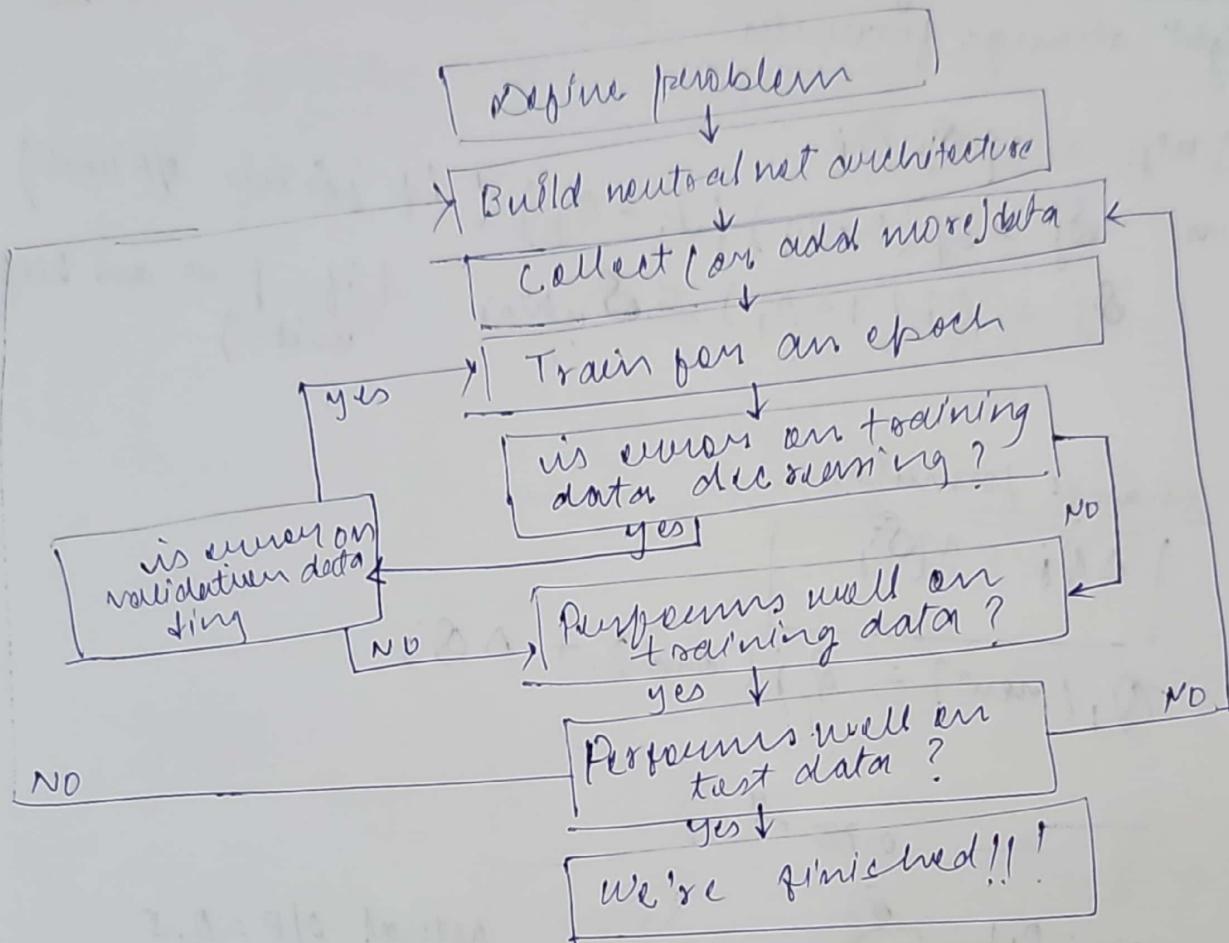
To find the optimal setting of hyper parameters grid search is applied. Where we pick a value for each hyper parameter from a finite set of possible options.

$$\eta \in \{0.01, 0.01, 0.001\}$$

$$\text{batch size}(k) \in \{16, 32, 64, 128, \dots\}$$

Then we train the model with every possible permutation of hyper parameter choices. Then we test the combination of hyper parameters with the best performance validation set and report the accuracy of the model trained with best combination on the test set.

→ Detailed workflow of Training & Evaluating a DL Model



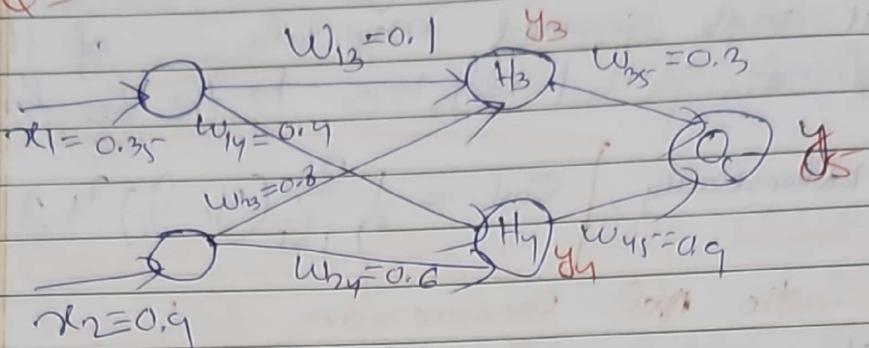
→ Preventing overfitting in Deep Neural Network:

① Regularization $\begin{cases} \text{L1} & \text{Regularization (Lasso regularization)} \\ \text{L2} & \text{Regularization (Ridge)} \end{cases}$

② Drop out

Backpropagation Algorithm :-

(1) Q -

Assume actual op $y = 0.5$ Learning rate (η) = 1Forward Pass (Complete op for y_3, y_4 and y_5)

$$f(\phi, v) = \frac{1}{1 + e^{-v}}$$

$$H_3 = x_1 * w_{13} + x_2 * w_{23}$$

$$= (0.35 * 0.1) + (0.9 * 0.8) = 0.755$$

$$y_3 = \frac{1}{1 + e^{-H_3}} = \frac{1}{1 + e^{-0.755}} = 0.68$$

$$H_4 = x_1 * w_{14} + x_2 * w_{24}$$

$$= 0.35 * 0.4 + 0.9 * 0.6 = 0.68$$

$$y_4 = \frac{1}{1 + e^{-H_4}} = \frac{1}{1 + e^{-0.68}} = 0.6637$$

$$\text{Now, } O_5 = y_3 * w_{35} + y_4 * w_{45}$$

$$= 0.68 * 0.3 + 0.6637 * 0.9$$

$$= 0.801$$

$$y_5 = \frac{1}{1 + e^{-O_5}} = \frac{1}{1 + e^{-0.801}} = 0.69$$

$$\text{Error} = Y_{\text{target}} - Y_{\text{actual}}$$

$$= 0.5 - 0.69$$

$$= -0.19$$

Weight Change Formula δ_j Error measure for unit j

$$\Delta w_{ji} = \eta \delta_j o_i$$

where -

$$\delta_j = o_j(1-o_j)(t_j - o_j) \quad (\text{if } j \text{ is an O/p unit})$$

$$\delta_j = o_j(1-o_j) \sum \delta_k w_{kj} \quad (\text{if } j \text{ is a hidden unit})$$

Backward Pass :For O/p unit :

$$\begin{aligned}\delta_5 &= y_5(1-y_5)(y_{\text{target}} - y_5) \\ &= 0.69(1-0.69)(0.5 - 0.69) = -0.0406\end{aligned}$$

For Hidden unit :

$$\begin{aligned}\delta_3 &= y_3(1-y_3) \times w_{35} * \delta_5 \\ &= 0.68(1-0.68) * 0.3 * (-0.0406) \\ &= -0.00265\end{aligned}$$

$$\begin{aligned}\delta_4 &= y_4(1-y_4)w_{45} * \delta_5 \\ &= 0.6637(1-0.6637) * 0.9 * -0.0406 \\ &= -0.0082\end{aligned}$$

Compute new weights

$$\Delta w_{ji} = \eta \delta_j o_i$$

$$\Delta w_{45} = \eta \delta_5 y_4 = 1 * -0.0082 * 0.6637 = -0.0054$$

$$w_{45}(\text{new}) = w_{45}(\text{old}) + \Delta w_{45} = -0.0269 + 0.9 = 0.8731$$

$$\Delta w_{14} = \eta \delta_4 x_1 = 1 * -0.0082 * 0.35 = -0.0287$$

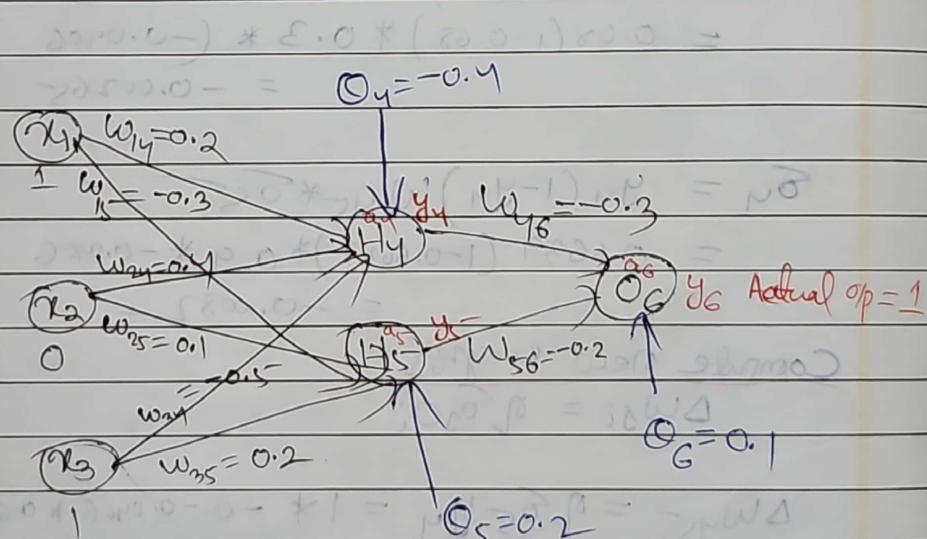
$$w_{14}(\text{new}) = w_{14}(\text{old}) + \Delta w_{14} = 0.4 - 0.0287 = 0.3971$$

Update all other weights

i	j	w_{ij}	δ_j	x_i	n	Updated w_{ij}
1	3	0.1	-0.00265	0.35	1	0.0991
2	3	0.8	-0.00265	0.9	1	0.7976
1	4	0.9	-0.0082	0.35	1	0.3971
2	4	0.6	-0.0082	0.9	1	0.5926
3	5	0.3	-0.0406	0.68	1	0.2724
4	5	0.9	-0.0406	0.68	1	0.8731

(2) Q. Assume the neurons have a sigmoid activation function, perform a forward pass and a backward pass on the network.

Assume that the actual op of y is 1 and learning rate is 0.9. Perform another forward pass.



Forward Pass

Compute Op for y_4, y_5, y_6

$$\begin{aligned}
 o_4 &= w_{14} * x_1 + w_{24} * x_2 + w_{34} * x_3 + o_1 \\
 &= (0.2 * 1) + (0.4 * 0) + (-0.5 * 1) + (-0.4) \\
 &= -0.7
 \end{aligned}$$

$$O(H_4) \Rightarrow y_4 = f(a_4) = \frac{1}{1+e^{-a_4}} = \frac{1}{1+e^{-0.332}} = 0.332$$

$$a_5 = w_{15} * x_1 + w_{25} * x_2 + w_{35} * x_3 + O_5 \\ = (-0.3 * 1) + (0.1 * 0) + (0.2 * 1) + 0.2 = 0.$$

$$O(H_5) = y_5 = f(a_5) = \frac{1}{1+e^{-a_5}} = 0.525$$

$$\text{Now, } a_6 = w_{46} * H_4 + w_{56} * H_5 + O_6 \\ = (-0.3 * 0.332) + (-0.2 * 0.525) + 0.1 \\ = -0.105$$

$$O(O_6) = y_6 = f(a_6) = \frac{1}{1+e^{-a_6}} = 0.474$$

$$\text{Error} = y_{\text{target}} - y_6 = 1 - 0.474 = 0.526$$

Backward Pass

Compute δ_4 , δ_5 and δ_6

For OP unit:

$$\delta_6 = y_6(1-y_6)(y_{\text{target}} - y_6) \\ = 0.474(1-0.474)(1-0.474) = 0.1311$$

For Hidden unit:

$$\delta_5 = y_5(1-y_5)w_{56} * \delta_6 \\ = 0.525(1-0.525)(-0.2 * 0.1311) \\ = -0.0065$$

$$\delta_4 = y_4(1-y_4)w_{46} * \delta_6 \\ = 0.332(1-0.332)(-0.3 * 0.1311) \\ = -0.0087$$

$$\Delta w_{46} = \eta \cdot \delta_6 y_4 = 0.9 * 0.1311 * 0.332 = 0.03917$$

$$w_{46}(\text{new}) = w_{46}(\text{old}) + \Delta w_{46} = -0.3 + 0.03917 \\ = -0.261$$

$$\Delta w_{14} = \eta \cdot \delta_4 x_1 = 0.9 * -0.0087 * 1 = -0.0078$$

$$w_{14}(\text{new}) = w_{14}(\text{old}) + \Delta w_{14} = -0.0098 + 0.2 = 0.192$$

Similarly, update all other weights.

i	j	w_{ij}	\bar{o}_j	x_i	η	updated w_{ij}
4	6	-0.3	0.1311	0.332	0.9	-0.261
5	6	-0.2	0.1311	0.585	0.9	-0.138
1	4	0.2	-0.0087	1	0.9	0.192
1	5	-0.3	-0.0065	1	0.9	-0.306
2	4	0.4	-0.0087	0	0.9	0.4
2	5	0.1	-0.0065	0	0.9	0.1
3	4	-0.5	-0.0087	1	0.9	-0.508
3	5	0.2	-0.0065	1	0.9	0.194

Similarly, update the biases

\bar{o}_j	Prev. \bar{o}_j	\bar{o}_j	η	Updated \bar{o}_j
\bar{o}_6	-0.14	0.1311	0.9	0.218
\bar{o}_5	0.2	-0.0065	0.9	0.194
\bar{o}_4	-0.4	-0.0087	0.9	-0.408

$$\Delta \bar{o}_j = \eta \cdot \bar{o}_j (w_{j-1})_{\text{old}} = w_j$$

$$\Delta \bar{o}_j (\text{new}) = \Delta \bar{o}_j + \bar{o}_j (\text{old})$$

Forward Pass:

Compute OIP for y_4, y_5 and y_6 .

$$\begin{aligned}
 o_4 &= (w_{14} * x_1) + (w_{24} * x_2) + (w_{34} * x_3) \\
 &= (0.192 * 1) + (0.4 * 0) + (-0.508) + o_4 \\
 &= -0.724
 \end{aligned}$$

$$y_4 = f(a_4) = \frac{1}{1+e^{-0.327}} = 0.327$$

$$\begin{aligned}a_5 &= (w_{15} * x_1) + (w_{25} * x_2) + (w_{35} * x_3) + \theta_5 \\&= (-0.306 * 1) + (0.1 * 0) + (0.191 * 1) + 0.191 \\&= 0.082\end{aligned}$$

$$y_5 = f(a_5) = \frac{1}{1+e^{-0.082}} = 0.520$$

$$\begin{aligned}a_6 &= w_{46} * h_1 + w_{56} * h_6 + \theta_6 \\&= -0.261 * 0.327 + (0.138 * 0.520) + 0.218 \\&= 0.061\end{aligned}$$

$$y_6 = f(a_6) = \frac{1}{1+e^{0.061}} = 0.515$$

$$\text{Error} = Y_{\text{target}} - y_6 = 0.485$$

Chapter - 6

Date _____

Page No. _____

Challenges with Gradient Descent

- * Vanishing gradient
- * Saddle point

⇒ Vanishing gradient

→ When training deep NN with gradient descent and back propagation, we have to calculate the partial derivatives by moving from the final output layer to the initial layer.
→ Due to this process, vanishing gradient problem arises.

⇒ Saddle point:

→ It is a point on the surface of the graph of a function where the slopes (derivatives) in orthogonal directions are all zero but which is not a local extremum of the function.

→ Saddle point injects confusion into the learning process.

→ Learning of the model becomes slow.

→ This saddle point gets into focus when the gradient descent (GD) works in multi dimensions (ie, we have more & more parameters).

Avoiding vanishing gradient

(parameters)

- Use ReLU or Leaky ReLU as activation function
- Reduce the model complexity (Reducing no. of layers)
- Better optimizers with a well-tuned learning rate

Momentum based Gradient optimization

- gt is an extension of Batch GD.
- Momentum is a strategy for accelerating the convergence of the optimization process by including a momentum element in the update rule.
- Momentum considers previous step weights gradient changes when updating current weights.
- This momentum factor assists the optimizer in continuing to go in the same direction of the search space even if the gradient changes direction or becomes zero.
- gt overcomes the oscillations of noisy gradients and coast across flat spots of the search space.
- Instead of stopping when the gradient is 0 at the 1st local minimum, momentum will continue to move forward because it considers the prior steep in the slope.

Mathematically;

$$V(t) = \beta * V(t-1) + (1-\beta) \cdot \nabla w$$

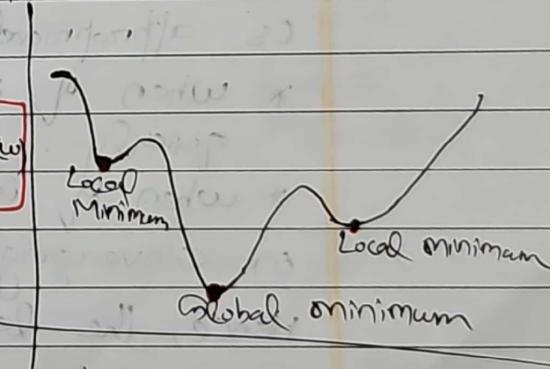
$$\nabla w = \frac{\partial E}{\partial w}$$

$$w = w - \alpha \cdot V(t)$$

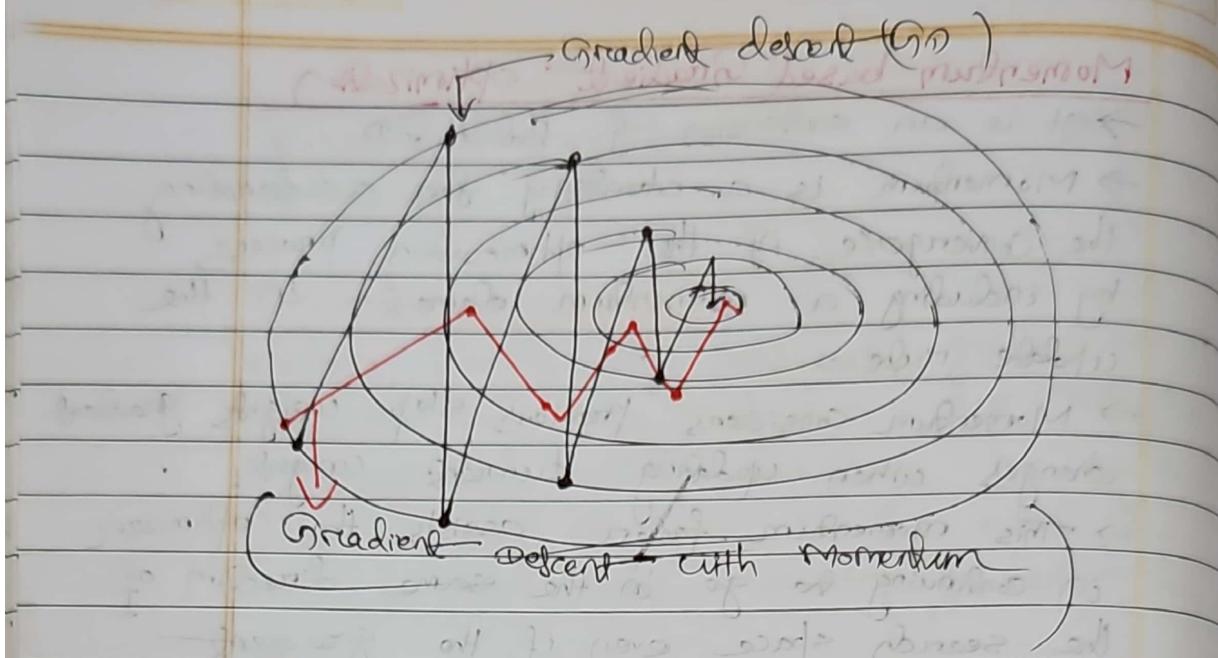
where V : velocity

β : momentum factor (ranges from 0 to 1)

α : learning rate (typically 0.9)



- the optimizer calculates the gradient of the cost function at each iteration, and updates the momentum term as the exponentially weighted average of the previous gradients.



Advantages :

- gt helps the algorithm to escape from local minima and saddle points.
- gt makes the algorithm to converge faster by avoiding oscillations.

Learning Rate Adaption

- The major challenge for training deep NN, is appropriately selecting the learning rate.
- * when η is too small, it doesn't learn quickly.
- * when η is too large, it may have difficulty in converging to a local minimum.

Hence, the learning rate adaption is very vital.

The most popular adaptive learning rate alg.s are

- ① AdaGrad (Adaptive Gradient) optimization
- ② RMS Prop (Root Mean Squared Propagation) optimization
- ③ Adam (Adaptive Moment Estimation) optimization