# Unit I :-
## Basics of Artificial Neural Network.

**Artificial Neural Network :-** Analogous to Human Brain.
imp: neuron, perceptron.

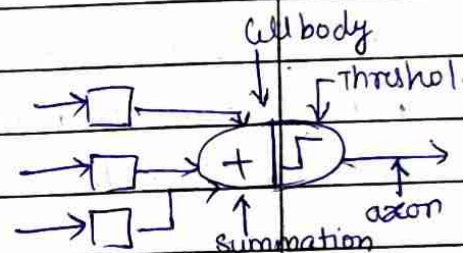Association of biological neuron with Artificial neuron,,,

→ dendrite :- receives signals from other neuron

synapse :-

## Association of biological net with artificial net.
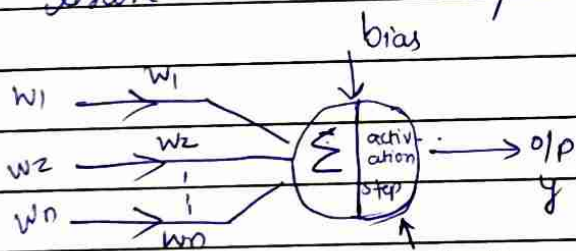
| Biological Neuron | Artificial Neuron |
|---|---|
| cell | neuron |
| Dendrites | weights or interconnectn |
| Soma | Sums the input |
| Axon | output |



i/p vector $X = [x_1, x_2 \ldots x_n]$

Weight vector $W = [W_1, W_2 \ldots w_n]$

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

This is going to decide what type of function. it is going to perform

o/p. Y is,

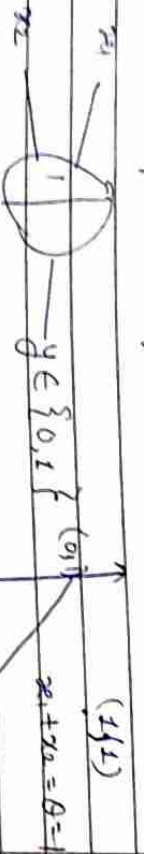$$Y = f(y_{in})$$

## MNIST Dataset

Training images ⟶ 60,000

Testing images ⟶ 10,000.     $28 \times 28 \times 1$
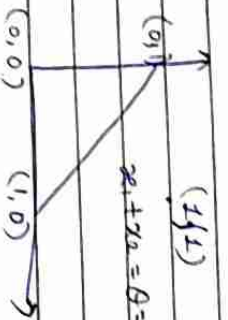
Grey scale

# Geometric Interpretation of m-p neuron



$y \in \{0, 1\}$
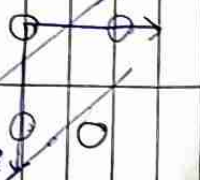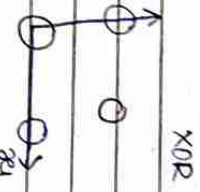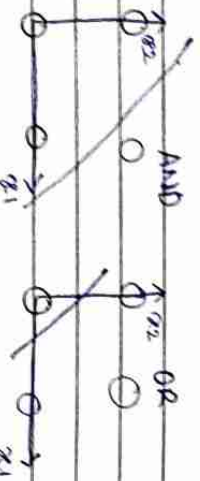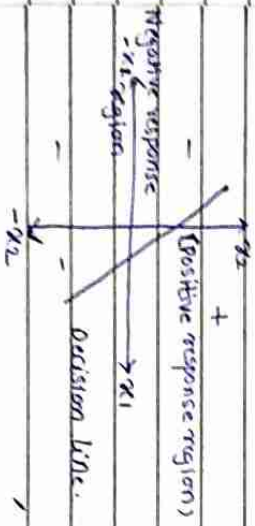
OR function

$x_1 + x_2 = \sum\limits_{i=1}^{2} x_i \geq 1$



$x_1 + x_2 = \theta = 1$

- For xor no linear $\theta$ $\phi$ $\cancel{th}$ separability

- Linear separability :- concept wherein the separation of input space into regions is based on whither the network response is positive or negative $\overline{ab} + ab$.

- A decision line is drawn to separate positive and negative responses.
- A division line is may also be called as decision making line.
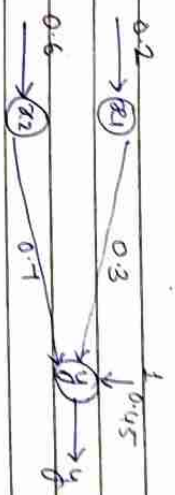- possibility was felt to classify the patterns based upon their output response.



(positive response region)

negative response region

decision line.



AND          OR          XOR

| | | No Separation is possible | Separation with 2 layers | Separation with CLU's |
|---|---|---|---|---|

output :- 0 - -
output :- 1 - -

---

$y_{in} = b + \sum\limits_{i=1}^{n} x_i w_i$   ← feature extractor

1. calculate the net input for the network shown in figure with bias included in the network?



$y_{in} = b + x_1 w_1 + x_2 w_2$
$y_{in} = 0.45 + 0.2 \times 0.3 + 0.6 \times 0.7$
$y_{in} = 0.45 + 0.06 + 0.42 = \underline{0.93}$

Activation function:- → to overcome disadvantage over its linear response → to handle non-linearity present as inputs in the network → to calculate networks output.

1. identity function.
2. Binary step function
3. Bipolar step function
4. Sigmoid function
   a. Binary Sigmoid Function
   b. Bipolar Sigmoid Function.
5. Ramp function

Activation function

Back propagation :- weight & Bias are going get updated based upon error calculation happening at the output (predicted by network).

Mean square error

Forward propagation :-

**Step activation:** $g(x) = 0$ for all $x$. for $x \to 0$ $x < 0$ $g(x) = \begin{cases} 0 & x < 0 \\ 1 & x = 1 \\ & x > 0 \end{cases}$ moguwa

**Identity function:** one to one input both a corresponding.

**Sigmoidal:** $\sigma(x) = \dfrac{1}{1 + e^{-x}}$

derivative: $g'(x) = \sigma(x)(1 - \sigma(x))$

starts approaching to 0 when we take derivative are hidden layer function so back propagation

Since the sigmoidal function is differential. So it is mostly used. And mostly used for non-linear data.

**RELU :-** Rectified Linear Unit.

**tanh Activation function:-** Differential function can be use in back propagation

$g(x) = 1 - (\tanh(x))^2$ — derivative.

$\tanh(x) = \dfrac{2}{1 + e^{-2x}} - 1$

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

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

$g(x) =$

**leaky Relu Activation function.**

$f(x) = \begin{cases} x, & \text{when } x > 0 \\ ax, & x < 0 \end{cases}$

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

**Softmax**

$f(x_i) = \dfrac{e^{x_i}}{\sum_j e^{x_j}}$

---

find out $y_{in}$ you are getting it is passed through some activation function.

RSigmoid.

**Q.** Obtain the o/p of the $H$ neuron $Y$ for the network shown in Fig 3 using activation function a) is binary sigmoid, b) Bipolar sigmoid



$y = b + x_1 w_1 + x_2 w_2$

$= -0.35 + 0.8 \times 0.1 + 0.5 \times 0.3 + 0.4 \times 0.2 = 0.5$

i. **Binary sigmoidal function,**

$y = f(y_{in}) = \dfrac{1}{1 + e^{-y_{in}}} = \dfrac{1}{1 + e^{-0.53}} = 0.62$

ii. **Bipolar sigmoidal function.**

$y = f(y_{in}) = \dfrac{2}{1 + e^{-y_{in}}} - 1 = -1 = 0.25$

**Flowchart of Training of Perception**

**Training (Single o/p classes):-**
1. Initialize the weight and bias (Initialize learning rate $\alpha$)

Random weights are assigned forward run (parameter ber)

value we choose for learning rate neural network is going to learn.

| | Single first | Hello | Pg |
|---|---|---|---|
| | | 1940 | 3.2.7 |

Binary data or be used for
Capable of handling
Capable of bipolar data

2

2. Check for stopping condition, if it is false, perform step 2-6.

3. Perform steps 3-5 for each bipolar or binary training vector pair $s:t$

4. Set activation (identity) for each input unit $i=1$ to $n$:
$$x_i = s_i$$

5. Calculate output response of each o/p unit $j=1$ to $m$: first, the net i/p is calculated.
$$y_{inj} = b_j + \sum_{i=1}^{n} x_i w_{ij}$$

6. Weight & bias adjustment for $j=1$ to $m$ & $i=1$ to $n$.

   if $t_j \neq y_j$ then
   $$w_{ij}(new) = w_{ij}(old) + \alpha \, t_j \, x_i$$
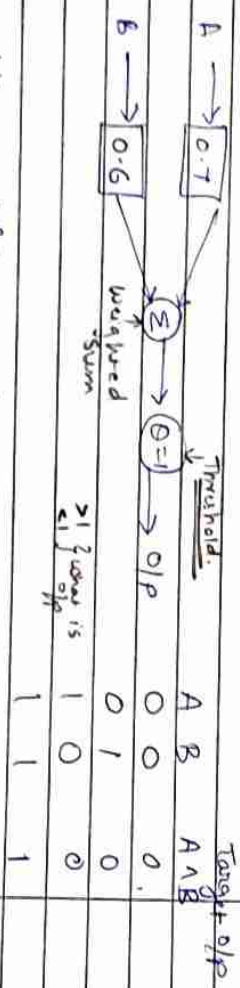   $$b_j(new) = b_j(old) + \alpha \, t_j$$
   
   else, we have
   $$w_{ij}(new) = w_{ij}(old)$$
   $$b_j(new) = b_j(old)$$

7. Test for the stopping condition i.e. if there is no change in weights then stop the training, else start again from step 2.

logic and gate perception training rule.



|   | A | B | Target o/p |
|---|---|---|---|
| | A | B | A∧B |
| | 0 | 0 | 0 |
| | 0 | 1 | 0 |
| | 1 | 0 | 0 |
| | 1 | 1 | 1 |

$>1$ ? $w_{ow}$ is 1 o/p
$<1$ ? o/p

Whenever A∧B are equal to 1 o/p must be 1, realised to have

2 weights, $w_1 = 1.2$, $w_2 = 0.6$, $\theta = 1$, $\alpha = 0.5$

∴ $A = 0$   $B = 0$, Target $= 0$

$\sum_{i=1}^{2} w_i x_i = 0 \times 1.2 + 0 \times 0.6$

— Thus is not greater than threshold ∴ calculated o/p is 0
$$y_{in} = 0 < \text{threshold} = 1. \text{ calculated o/p } y = 0$$

— Since actual calculated value is 0 & targeted o/p is 0 ∴ no need of net weight updation.

i) A = 0, B = 1,
$$y_{in} = 0 \times 1.2 + 1 \times 0.6 = 0.6$$
$$y_{in} = \sum_{i=1}^{2} w_i x_i$$

— Since $y_{in} = 0.6$ is less than the threshold = 1 ∴ the o/p is 0.
So no need of weight updation.

ii) A = 1, B = 0
$$y_{in} = 1 \times 1.2 + 0 \times 0.6 = 1.2$$

As the calculated o/p of $y$ is not match to target o/p. So we need the updation of weight.

∴ weight are updated as per perceptron training rule.
$$w_{ij}(new) = w(old) + \text{learning rate} \cdot \text{target o/p } \& \text{ actual o/p}$$
$$w_1 \text{ new} = 1.2 + 0.5(0-1) 1 = 0.7$$
$$w_2 \text{ new} = 0.6 + 0.5 (0-1) 0 = 0.6$$
various $b_{new} = b_{old} + \alpha(t-0) x_i$
(expected updation)

Taking new weights $w_1 \text{new} = 0.7$ & $w_2 = 0.6$.

i) A = 0, B = 0, Target = 0
$$y_{in} = 0 \times 0.7 + 0 \times 0.6 = 0$$
$$y_{in} < \text{threshold}.$$

ii) A = 0, B = 1, Target = 0
$$y_{in} = 0 \times 0.7 + 1 \times 0.6 = 0.6$$
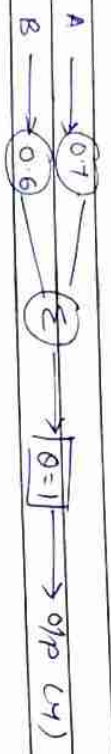∴ No updation of weight.

iii) $A=1$, $B=0$, $y_{in}=0$

$y_{in} = 1 \times 0.7 + 0 \times 0.6 = 0.7$

no need of updation of weight

iv) $A=1$, $B=1$, $y_{in}=1$

$y_{in} = 1 \times 0.7 + 1 \times 0.6 = 1.3$

No need of weight updation

we have classified all the training i/p correctly.

with $\theta$ $w_1 = 0.7$ & $w_2 = 0.6$, threshold = 1, $\theta = 1$ &

learning rate $\alpha = 0.5$

Modified perceptron And gate.

Initialize test sample.

$w_i : x_i$

A —→ 0.7

B —→ 0.6 —→ C —→ | $\theta = 1$ | —→ o/p (y)

Perceptron rule for AND gate.

$y = x_1 x_2 + x_1 \overline{x_2}$

$y = x_1 \overline{x_2} + \overline{x_1} x_2$

$y = z_1 + z_2$

$z_1 = x_1 \overline{x_2}$ — function1

$z_2 = \overline{x_1} x_2$ — function2

| $x_1$ | $x_2$ | y |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

hidden layer.

$x_1$ —→ (x1) $w_{11}$ —→ (z1) $v_1$

$x_2$ —→ (x2) $w_{21}$ $w_{12}$ $w_{22}$ —→ (z2) $v_2$ —→ (yin) —→ | ∫ | —→ y(o/p)

$w \leftarrow$ weights

$v \leftarrow$ weights

activation function.

$f(y_{in}) = 1$ if $y_{in} \geq \theta$

$= 0$ if $y_{in} < \theta$.

input = $x_i$

$w_{ij}$

---

1st function $z_1 = x_1 \overline{x_2}$

| | $x_1$ | $x_2$ | $z_1 = x_1 \overline{x_2}$ |
|---|---|---|---|
| | 0 | 0 | 0 |
| | 0 | 1 | 0 |
| | 1 | 0 | 1 |
| | 1 | 1 | 0 |

set initial weight $w_{11} = w_{21} = 1$, Threshold $\theta = 1$, learning rate = 1.5

$(x_1, x_2)$      $y_{in} = x_1 w_{11} + x_2 w_{21}$

$(0,0) = z_1$ in = $w_{ij} * x_i = (1 \times 0) + (1 \times 0) = 0$   $0 < 1$

$(0,1) = z_1$ in = $w_{ij} \times x_i = (1 \times 0) + (1 \times 1) = 1$   $1 < 1$ wrong o/p 2

$(1,0) = z_1$ in = $w_{ij} \times x_i = (1 \times 1) + (1 \times 0) = 1$   $1 < 1$

$(1,1) = z_1$ in = $w_{ij} \times x_i = (1 \times 1) + (1 \times 1) = 2$.   $2 < 1$

↳ which $z_1$ not same as target o/p ∴ weights are updated.

$w_{new} = w_{old} + \alpha (t - 0) x_i$

$w_{11} = 1 + 1.5 (0 - 1) 0 = 1$

$w_{21} = 1 + 1.5 (0 - 1) 1 = -0.5$.

   ↳ actual o/p   ↳ $x_2$

Now, assume $w_{11} = 1$ & $w_2 = -0.5$

$x_1 = 0$   $x_2 = 0$

$z_1 y_{in} = (1 \times 0) + (-0.5 \times 0) = 0$

$x_1 = 0$, $x_2 = 1$

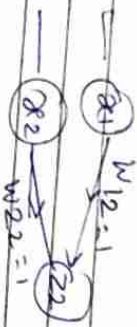$z_1 y_{in} = (0 \times 0) + (-0.5 \times 1) = -0.5$

$x_1 = 1$   $x_2 = 0$

$z_1 y_{in} = (1 \times 1) + (-0.5 \times 0) = 1$

$x_1 = 1$, $x_2 = 1$

$z_1 y_{in} = (1 \times 1) + (-0.5 \times 1) = 1 - 0.5 = 0.5$

∴ calculated o/p are values are same as target o/p $z_1$ with the updated weight $w_{11} = 1$ & $w_{21} = -0.5$

| $x_1$ | $x_2$ | $z_2$ |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 0 |

$x_1=0, x_2=0,$

$z_2 y_{in} = (1\times0)+(1\times0)=0$

$x_1=0, x_2=1$

$z_2 y_{in} = (1\times0)+(1\times1)= 1$

$x_1=1, x_2=0$

$z_2 y_{in} = (1\times1)+(0\times0)=1$

$\therefore$ we need to update the weight.

$W_{new} = W_{old} + \alpha (t-0) x_i$

$W_{12} = 1+ 1.5 (0-1)1 = -0.5$

$W_{22} = 1+ 1.5 (0-1) 0 = 1.$

| $z_1$ | $z_2$ | |
|---|---|---|
| 0 | 0 | $y_{in} = (1\times0)+(1\times0)=0$ |
| 0 | 1 | $y_{in} = (1\times0)+(1\times1)=1$ |
| 1 | 0 | $y_{in} = (1\times1)+(1\times0)=1$ |
| 0 | 0 | $y_{in} = (1\times0)+(1\times0)= 0$ |

$V_1 = V_2 = 1 \qquad \alpha = 1.5$

$y = z_1.V_1 \text{ or } z_2.V_2 \quad y_{in} = z_1V_1 + z_2V_2.$

| $x_1$ | $x_2$ | $z_1$ | $z_2$ | $y=z_1.V_1+z_2.V_2$ |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 | 1 |
| 1 | 0 | 1 | 0 | 1 |
| 1 | 1 | 0 | 0 | 0 |

$\therefore$ Perceptron learning for XOR is shown below with its weight

Back Propagation :- optimize weights / update



i) For $h_1$

Forward Propagation

$\text{out}H_1 = w_1 \times x_1 + w_2 \times I_2 + B \times 1$

$= 0.15\times 0.05 + 0.2 \times 0.1 + 0.35$

$= 0.3775$

**Binary Sigmoidal function**

o/p h1 = $\dfrac{1}{1+e^{-neth1}}$ = $\dfrac{1}{1+e^{(-0.3775)}}$ = 0.5932

ii) for neth2

neth2 = w3·I1 + w4·I2 + b1·1
= 0.25×0.05 + 0.3×0.1 + 0.35
= 1.3925

Binary Sigmoidal function

o/p h2 = $\dfrac{1}{1+e^{(-0.3925)}}$ = 0.59688431

iii) for O1

net O1 = w5×h1 + w6×h2 + b2×1
= 0.4×0.5932 + 0.45×0.59688 + 0.6×1
= 1.10596

By binary sigmoidal function

o/p O1 = $\dfrac{1}{1+e^{(-1.10596)}}$ = 0.7513 ≠ 0.01

iv) for O2

net O2 = w7×h1 + w8×h2 + b2×1
= 0.5×0.5932 + 0.55×0.59688 + 0.6×1
= 1.224884

By binary sigmoidal function

o/p O2 = $\dfrac{1}{1+e^{(-1.224889)}}$ = 0.772928 ≠ 0.99

Calculate total error.

Etotal = $\sum \dfrac{1}{2}$ (target - o/p)²

E01 = $\dfrac{1}{2}$ (0.01 - 0.7513)²

= 0.27476

---

E02 = $\dfrac{1}{2}$ (0.99 - 0.772928)² = 0.02356

E = E01 + E02
= 0.27476 + 0.02356
= 0.29832

**Back Propagation** — w5 updation

i) $\dfrac{\partial Etotal}{\partial ws} = \dfrac{\partial Etotal}{\partial out1} \times \dfrac{\partial out1}{\partial net01} \times \dfrac{\partial net01}{\partial ws}$

$\dfrac{\partial Etotal}{\partial out1}$ = -(target o1 - out o1)

= -(0.01 - 0.7513)

= 0.7513

$\dfrac{\partial out01}{\partial net01}$ = out01(1-out01)

= 0.7513(1 - 0.7513)

= 0.1868483

$\dfrac{\partial net01}{\partial ws}$ = out h1

= 0.5932

$\dfrac{\partial Etotal}{\partial ws}$ = 0.7513 × 0.1868483 × 0.5932

= 0.08316465

To decrease the error. Subtract with learning rate (x = 0.5)

ws+ = ws - x $\dfrac{\partial Etotal}{\partial ws}$

= 0.4 - 0.5 × 0.08216465

ws+ = 0.3589 ...

iv) $W6^+ = W6 - \alpha \dfrac{\partial E_{total}}{\partial W6}$

$= 0.45 - 0.5 \times 0.08216\ 4r$

$= 0.4089\ 4775$

iii) W7

$\dfrac{\partial E_{total}}{\partial out\ 2} = -(0.99 - 0.772928)^{2-1}$

$= -0.2171$

$\dfrac{\partial out\ 02}{\partial net\ 02} = out\ 02(1 - out\ 02)$

$= 0.1755142$

$\dfrac{\partial net\ 02}{\partial W7} = out\ h2$

$= 0.59688\mu$

$\dfrac{\partial E_{total}}{\partial W7} = -0.2171 \times 0.1755142 \times 0.59688\mu$

$\dfrac{\partial E_{total}}{\partial W7} = -0.02274$

$W7^+ = W7 - 0.5 \dfrac{\partial E_{total}}{\partial W7}$

$= 0.5 - 0.5\ (-0.02274)$

$= 0.51137$

iv) $W8: W8 - 0.5(-0.02274)$

$= 0.55 - 0.5(-0.02274)$

$= 0.56137$

Hidden layers:-

$\dfrac{\partial E_{total}}{\partial W1} = \dfrac{\partial E_{total}}{\partial out\ h1} \times \dfrac{\partial out\ h1}{\partial net\ h1} \times \dfrac{\partial net\ h1}{\partial W1}$

$\dfrac{\partial E_{total}}{\partial out\ h1} = \dfrac{\partial E_{01}}{\partial out\ h1} + \dfrac{\partial E_{02}}{\partial out\ h1}$

iv) $\dfrac{\partial E_{01}}{\partial out\ h1} = \dfrac{\partial E_{01}}{\partial net\ 01} \times \dfrac{\partial net\ 01}{\partial out\ h1}$

$\dfrac{\partial E_{01}}{\partial net\ 01} = \dfrac{\partial E_{01}}{\partial out\ 01} \times \dfrac{\partial out\ 01}{\partial net\ 01}$

$= 0.741385 \times 0.1126215$

$= 0.38098$

$\dfrac{\partial net\ 01}{\partial out\ h1} = W5 = 0.40$

$\dfrac{\partial E_{01}}{\partial out\ h1} = 0.38098 \times 0.40$

$= 0.55392$

iv) $\dfrac{\partial E_{02}}{\partial out\ h2} = \dfrac{\partial E_{total}}{\partial out\ 02} = -(0.99 - 0.772928)^{2-1}$

$\dfrac{\partial out\ 02}{\partial out\ 02} = 0.1755142$

$= -0.21709$

$\dfrac{\partial out\ 02}{\partial net\ 02} = 0.1755142$

$= -0.21709$

$\dfrac{\partial E_{02}}{\partial out\ h2} = -0.21709 \times 0.1755142$

$= -0.038100\mu$

$\dfrac{\partial net\ 02}{\partial out\ h2} = W1 = 0.15$

$$\frac{\partial E_{O2}}{\partial out h_2} = \frac{\partial E_{O2}}{\partial net o_2} \times \frac{\partial net o_1}{\partial out h_1}$$

$$= -0.03810 \times 0.5$$

$$= -0.0190502.$$

$$\frac{\partial E_{total}}{\partial out h_1} = \frac{\partial E_{O1}}{\partial out h_1} + \frac{\partial E_{O2}}{\partial out h_2}$$

$$= 0.0553992 + (-0.01905)$$

$$= 0.0363492.$$

$$\frac{\partial out h_1}{\partial net h_1} = out h_1 (1 - out h_1)$$

$$\frac{\partial net h_1}{\partial w_1} = 0.593269(1 - 0.593269)$$

$$= 0.241300 \neq$$

$$\frac{\partial net h_1}{\partial w_1} = i_1 = 0.05$$

$$\therefore \frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out h_1} \times \frac{\partial out h_1}{\partial net h_1} \times \frac{\partial net h_1}{\partial w_1}$$

$$= 0.0363492 \times 0.2413007 \times 0.05$$

$$= 0.000438554.$$

$$w_1^+ = w_1 - \eta \frac{\partial E_{total}}{\partial w_1}$$

$$= 0.15 - 0.5 (0.0004385)$$

$$= 0.149780$$

$$w_2^+ = w_2 - \eta \frac{\partial E_{total}}{\partial w_1}$$

$$= 0.2 - 0.5 (0.0004385)$$

$$= 0.199781$$

---

# Unit II -

$$Z = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} \\ w_{2,1} & w_{2,2} & w_{2,3} \\ w_{3,1} & w_{3,2} & w_{3,3} \\ w_{4,1} & w_{4,2} & w_{4,3} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix} =$$

Binary cross-entropy
where no of loss function currency

- Cost function - error function is for a single training example/input
  Loss function - A cost function - is average loss over the entire training
  dataset.

- Regression -
  o MSE - difference b/w actual & predicted.
  o MAE - Mean absolute error

- Classification
  o Binary cross-entropy - used in binary classification $\left(\frac{-1}{n} \sum_{i=1}^{n} y_i \log \hat{y_i} + (1-y_i)\right)$
  o Categorical cross-entropy - used in multiclass classification $- \sum_{j=1}^{m} y_j \log (\hat{y_j})$

Parameters :-
  Few layers - underfitting
  More layers - overfitting

Architecture :-

o More generally, we can calculate the activation of neuron i in layer l
  $$z_i^{[l]} = w_i^{[l]} a^{[l-1]} + b_i^{[l]}$$
  $$a_i^{[l]} = g(z_i^{[l]})$$

Similarly, we can calculate all of the activations for a given layer l by
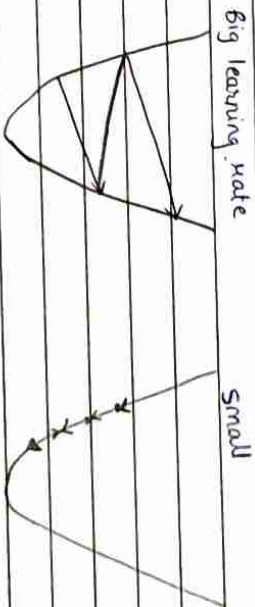using our weight matrix $w^{[l]}$

# Parameters and hyper parameters tuning

- Parameters: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$
  - No. of epochs.

- Hyperparameters are the variables which determine the network structure (eg. Number of Hidden Units) and variables which determine how the network is trained. (eg:- Learning Rate)

- Hyper parameters are set before training (helps optimizing the weights and bias)

- If hyperparameters neglected and not selected properly, they can be the cause of high error, resulting in wrong predictions.

- The number of layers, neurons, iterations, choice of activation function (ReLU, Tanh, Sigmoid, in the hidden layers) Optimizer (Adam, Adagrad, RMS prop), learning rate, choice of cost function, drop out rate, these are parameter control w's and b's

- layer size:- The second step is to tune the number of layers

- fewer layers may give an underfitting result while too many layers may make it overfitting

- Number of neurons - may be same or different in each layer. Depends on solution complexity

- An activation function is a parameter in each layer. The input values traverse from a layer to another layer keep changing according to the activation function

- Optimizer:- One of the hyperparameters in the optimizer is the learning rate

- learning rate:- Control the step size ζero so model to reach the minimum loss function (if gradient descent is working properly, it will itself lowering rate)

- A higher learning rate (alpha) makes the model learn faster, but it may miss the minimum loss function.

- A lower learning rate gives a better chance to find a minimum loss function.

- As a trade - off lower learning rate needs
  - higher epochs; or
  - more time and
  - memory capacity resources.

- **Gradient Descent :-**



  Big learning rate        Small

- The gradient vector has both a direction and a magnitude
- Gradient descent algorithm multiply the gradient by a scalar known as the learning rate (also sometimes called step size) determines the next time.
  - eg:- If the gradient magnitude is 2·5 and the learning rate is 0.0, then the gradient descent algorithm will pick the next point 0.025 (0.01 × 0.25) away from the previous point

# Gradient Descent and its variants :-

- Gradient descent algorithm updates the parameters by moving in the direction opposite to the gradient of the objective function with respect to the network parameters

  initialize w, b

  Iterate over data :

  compute ŷ

  compute $f(w, b)$ <sub>forward</sub>

  $wt+1 = wt - \eta \Delta wt$ (loss function w.r.t weight)

  $bt+1 = bt - \eta \Delta bt$

  till satisfied

- If the observation size of the training dataset is too large, it will definitely take a longer time to build the model

- To make the model learn faster, we can assign batch size so that not all of the training data are given to the model at the same time.

- Batch Size is the number of training data sub-samples for the input

  - eg:- If the training dataset has 77,500 observations and the batch size is 1000, the model will learn 77 times with 1000 training data sub-samples and another last to learning from it

    500 training data sub-samples and another last learning from this 500 training

  - The smaller batch size makes the learning process faster and bigger batch size and has a slower to learning process.

- 1 epoch :- 1 training dataset is passed forward and backward through the neural network once.

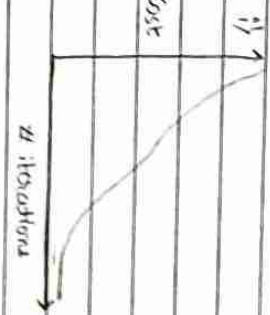- Too small no. of epochs - underfitting, because the neural network has not learned much enough

- Many epochs - Overfitting, where the model can predict the data very well, but cannot predict new unseen data well enough

- Batch gradient :- Takes the entire dataset - calculate the cost function
  - update parameter - speed , complexity & more learning cost

- Stochastic gradient :- large training examples - expensive computationally, observation and every time the weights are updated as you see as an iteration

- Mini-batch Gradient Descent :- Take a subset of data and update the parameters based on every subset

* Comparison of Gradient descent variants :-

  ii) 

  Cost

  n : iterations

  - Cost function return smoothly

  - cost keeps on decreasing over the epochs smoothly

  - Disadvantage : slow training examples

  - Computationally expensive

  - Batch Gradient

- epoch :- The no. of time a whole dataset is passed through the neural network model is called an epoch.

- Lot of variations in cost functions
- Updation is not smooth
- Updating the parameters based on a single observation, there are a lot of iterations

$$w_3^+ = w_3 - \eta \left( \frac{\partial Etotal}{\partial w_3} \right)$$

$$= 0.25 - 0.5 (0.0008746)$$

$$= 0.2495621$$

### SGD

- model starts learning noise as well

$$w_4^+ = w_4 - \eta \left( \frac{\partial Etotal}{\partial w_3} \right)$$

$$= 0.3 - 0.5 (0.0008746)$$

$$= 0.2995627$$

# iterations

- Smoother cost function as compared to SGD
- Updation of the cost function is smoother as compared to that of the cost function in SGD
- Not updating the parameters after every single observation but after every subset of the data.

### # - Batch

### Momentum Gradient Descent Algorithm :-
update rules, also include the history component. $V_t$ it stores all the previous gradient movements till this time t.

$$V_t = \gamma V_{t-1} + \eta \nabla w_t = \gamma^{t-1} \cdot \eta \nabla w_1 + \gamma^{t-2} \cdot \eta \nabla w_t \cdots \gamma^{t-n} \eta \nabla_i w_t$$

$$\frac{\partial Etotal}{\partial outh_1} = 0.036 3092$$

### Mini Batch Gradient.

- '1 epoch = one pass over the entire data.
- 1 step = one update of the parameters
- N = Number of data points
- B = Mini Batch size

### Algorithm.

| | # of steps in 1 epoch. |
|---|---|
| Vanilla (Batch) gradient. | 1 |
| Stochastic gradient. | N |
| Mini - Batch gradient | $\frac{N}{B}$ |

$$\frac{\partial outh_2}{\partial neth_2} = outh_2 \ (1- outh_2)$$

$$= 0.59688 \ (1- 0.59688)$$

$$= 0.240614$$

$$\frac{\partial neth_2}{\partial w_3} = i2 = 0.1$$

- Vanishing and Exploding gradient.
- as we are calculating gradient values from output to input
so weight values are increasing - Vanishing

$$\frac{\partial Etotal}{\partial w_3} = \frac{\partial Etotal}{\partial outh_2} \times \frac{\partial outh_2}{\partial neth_2} \times \frac{\partial neth_2}{\partial w_3}$$

$$= 0.036 349 2 \times 0.240614 u \times 0.1$$

$$= 0.0008746$$

$$\circ \xrightarrow{\hspace{2cm}} \eta \hspace{0.3cm} \text{decreasing.}$$ (exploding)

$$\boxed{\frac{\partial L}{\partial w}}$$ exploding

$$\boxed{\frac{\partial L}{\partial w}}$$

Converge.

Gradient :- derivative of loss functn w.r.t weight.
used to update the weight to minimize the loss functn during back propagation.

• Vanishing :- Derivative or slope will get smaller and smaller as we go backward with every layer during back propagative.
• update is very small or exponentially small, the training time takes too much longer, and in the worst case, this may completely Stop the neural network training. — exploding



— Occurs when the derivatives or slope will get larger and larger as we go backward.
— because of weights, not because of activation functn.
— Due to high weight values, the derivatives will also higher so that the new weight varies a lot to the older weight. and the gradient will never converge. So it may result in oscillating around the minima and the global minima.

Remedies :- Vanishing.
1. Choice of Activation function :- Vanish occurs with sigmoid and tanh activation functn because.
  because the derivatives are betn .0 to 0.25 and 0~1. Therefore, updated weight values are small, & new weight values are very similar to the old weight. This leads to vanishing.
  To avoid this problem using ReLU activation function because the gradient is 0 for negative and zero i/p, & 1 for positive i/p.

Q. Appropriate setter choice of weight.
  weight < 1     vanishing
  weight > 1     exploding. (algorithm becomes unstable for oscillate
• So choose weights randomly.
  — Intelligent choosing of Backpropagation learning algorithm

• Vanishing occurs due to sigmoid & tanh activation functn.
  exploding gradient due to large setter weight.

• For every layer, weights are commonly sampled from a normal distribution with mean = 0, standard deviation = 1
           $\mu = 0$ and $\sigma = 1$.

• "He" weight initialization.
           $\mu = 0$ and $\sigma = \sqrt{\dfrac{1}{n_i}}$, or $\sigma = \sqrt{\dfrac{2}{n_i}}$
           $n_i$ = no. of neurons in preceding layer.

— "He" initialization and "Xavier" initialization ensure that the weights or are close to 1.

— "Xavier" weight initialization.
           $\mu = 0$ and $\sigma = \sqrt{\dfrac{2}{n_i + n_{i+1}}}$

* Gradient Clipping :-
        Give the threshold value depending on that the gradient of value is more or less so for that we need to perform gradient clipping method.

  if gradient exceed, set them to max upper bound of interval
  if gradient fall below set them to min lower bound of interval

# Different architecture of Neural Network.

Bias, errors and variance.
if neural o/p model is not accurate, it can make prediction errors. and these prediction errors are usually known as Bias and variance.

- Bias is the difference b/n actual & predicted values
- High Bias :- model not captured data well during training and cannot perform well on test data.
- the instance where the model cannot find pattern in training Set and fail hence fails for both seen and unseen data, is called <u>underfitting</u>

- Neural N/w performs very well on training data, but fails as soon of it sees some new data from the problem domain and is very imp to take care of this in NN

No. of parameters in a feed-forward Neural Network

Bias                          Bias

3 i/p layer        4 hidden layer        2 output.

Assumption:-

$i$ = no. of neuron in i/p layer
$h$ = 9/- hidden layer
$o$ = -/- output layer

i) No. of connectn betn 1st & 2nd layer $(3 \times 4) = 12$ which is nothing but the product of $i$ & $h$

ii) -/- 2nd & 3rd layer $(4 \times 2) = 8$ -/- $h$ & $o$

iii) There are connections between layers via bias as well. No. of connections betn the bias of the first layer & neuron in second layer (except bias of the second layer) : $1 \times 4$ which is nothing but $h$.

iv) No. of connectn betn bias of second layer and neuron of third layer : $1 \times 2$ which is nothing but $o$

Summing up all ;

$9 \times 6 + 4 \times 2 + 1 \times 4 + 1 \times 2$

$12 + 8 + 4 + 2$

$26$ — trainable parameters.

$= i \times h + h \times 0 + h + 0$

Thus, the total number of parameters in a feed-forward neural network with one hidden layer is given by

$(i \times h + h \times 0) + h + 0$

**Scenario-2 :-** A feed-forward n/w with 3 hidden layers.
No. of units in the i/p, first hidden, third hidden, second hidden, third hidden, output layer are $h_1, h_2, h_3$.
$35, 6, 4$ and $2$

**Assumption.**

i : no. of neurons in input layer

$h_1$ : ——— $h_2$ — first hidden layer

$h_2$ : ——— second — $h_3$

$h_3$ : ——— third ; $h_1$

$0$ = no. of neurons in o/p layer

**formula for n hidden layers.**

$i \times h_1 + \sum_{k=1}^{n-1}(h_k \times h_{k+1}) + h_n \times 0 + \sum_{k=1}^{n} h_k + 0.$

o:] apply

$= (3 \times 5) + (5 \times 6) + (6 \times 4)(4 \times 2) + (1 \times 5)$

$(1 \times 6) + (1 \times 4)(1 \times 2)$

$= i \times h_1 + h_1 \times h_2 + h_2 \times h_3 + h_3 \times 0) + h_1 + h_2 + h_3 + 0.$

$= 15 + 30 + 24 + 8 + 5 + 6 + 4 + 2$

$= 94.$

③ connect of 1 & 2 = $(3 \times 5) = 15$

i] ——— $h_1 ; 2 \& 3 = (5 \times 6) = 30$

ii] ——— $; ——— 3 \& 4 = (6 \times 4) = 24$

iii] ——— $4 \& o/p = (4 \times 2) = 8$

iv] $-1 - h_1 = (1 \times 5) = 5$

v] $-1 - h_2 = (1 \times 6) = 6$

vi] $h_3 = (1 \times 4) = 4$

viii] $o/p = (1 \times 2) = 2$