Perceptron

$\uparrow \hat{y}$  y

Forward propagation.

$$\phi \sum_{i=1}^{m} w_i x_i$$

$\hat{y}$

Output ~~Actual~~ values

$\hat{y}$

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

$\longrightarrow$ cost-function

$I_1$, $I_2$, $I_m$ — inputs

$w_1$, $w_2$, $w_m$ — Weights

y  Actual value
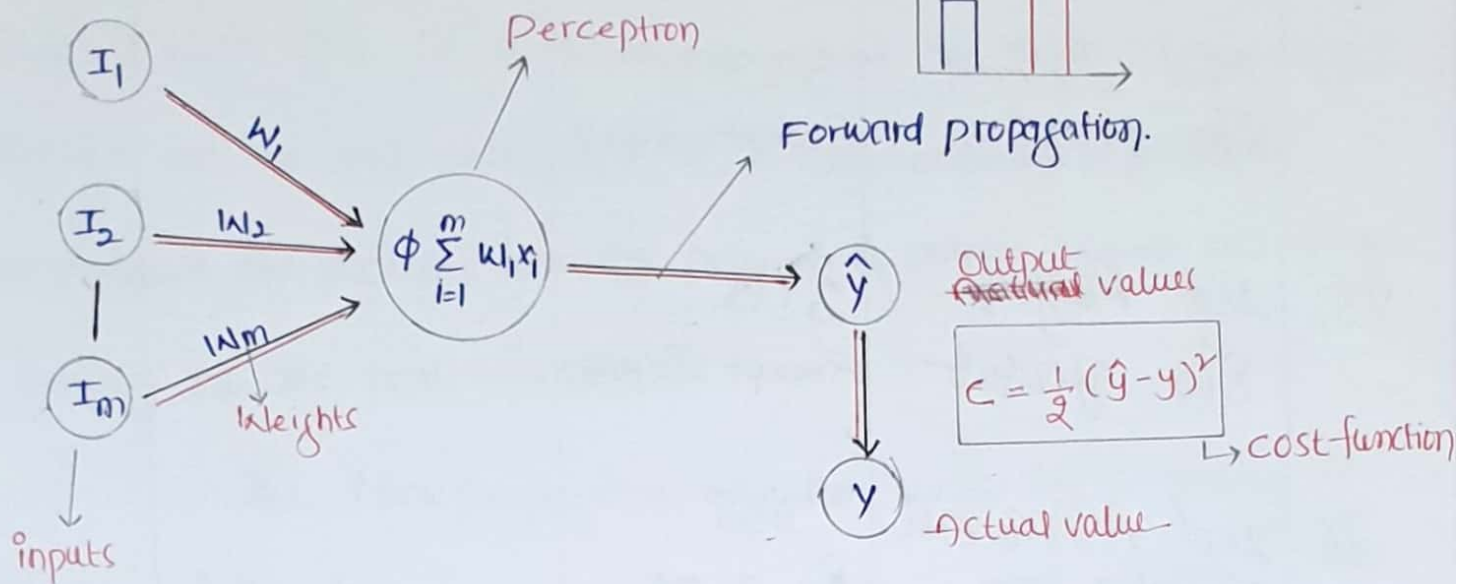
→ Feed features to the Input layers

→ Caliculate the weights of an each neuron

→ predict the output

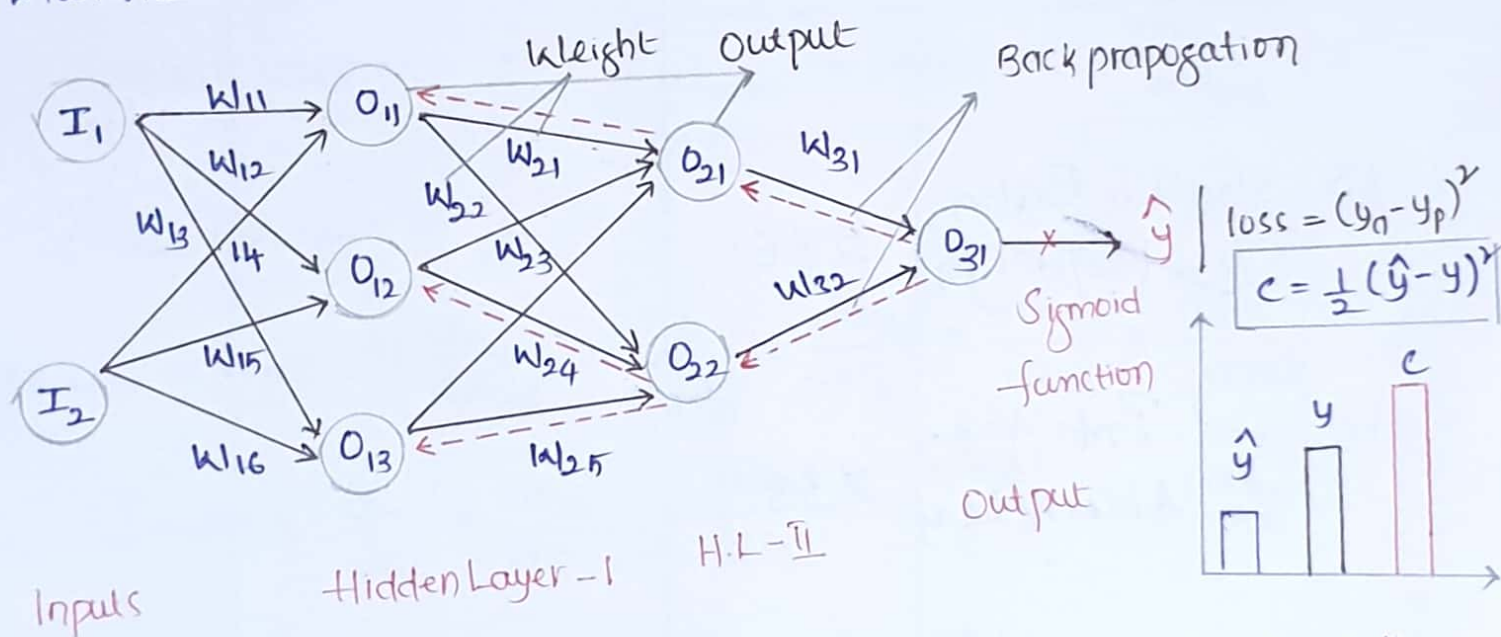→ finally compare the actual and predicted output.

Back propagation :- procedure —

i) caliculate the cost-function for predicted output in the forward Propagation.

ii) back propagate (Gradient descent) and adjust the weights of the layers so that Cost-function is minimized.

iii) predict the output for all observations. & caliculate the 'c'

iv) Adjust the weights so that Cost-function is minimized.

v) Again predict the output with new adjusted weights and repeate the procedure untill the Cost-function (c) is minimized.

# Back propagation :- [ Backword prapogation ]

It is a technique used to train a certain classes of neural networks, it is essentially a principle that allows the machine learning program to adjust it self according to looking at its past -function. to

The back prapogation algorithm works by computing the gradient of the loss-function with respect to each weight by the chain rule, computing the gradient one layer at a time, iterating backword from the last layer to avoid redundant caliculations of. Intermediate terms in the chain rule.

Weight output                                         Back prapogation

$I_1$  $W_{11}$ → $O_{11}$
$W_{12}$
$W_{13}$  $W_{21}$  $O_{21}$  $W_{31}$
$W_{22}$
14  $O_{12}$  $W_{23}$  $O_{31}$ × $\hat{y}$   | loss $= (y_a - y_p)^2$

$W_{15}$  $W_{22}$  $c = \frac{1}{2}(\hat{y} - y)^2$

$I_2$  $W_{24}$  $O_{32}$   Sigmoid   function

$W_{16}$  $O_{13}$  $W_{25}$   Output

Inputs        Hidden Layer - 1        H.L - II

Note :- Back prapogation at a time considering only one path if we observe from output we have two padbe intially it takes one path ( $O_{31} \to O_{21} \to O_{11}$ ) then other ( $O_{31} \to O_{22} \to O_{12}$ ) like that.

Activation functions $ its forms :-

1) Sigmoid $\rightarrow$ $y = \dfrac{1}{1+e^{-x}}$     2) Tanh $\rightarrow$ $y = Tanh(x)$

3) Step function $\rightarrow$ $y = \begin{cases} 0, & x < n \\ 1, & x \geqslant n \end{cases}$

4) Soft plus $\rightarrow$ $y = log(1+e^x)$

5) ReLU $\rightarrow$ $y = \begin{cases} 0, & x < 0 \\ x, & x \geqslant 0 \end{cases}$

6) Soft Sign $\rightarrow$ $y = \dfrac{x}{(1+|x|)}$

7) ELU $\rightarrow$ $y = \begin{cases} \alpha(e^x - 1), & x < 0 \\ x, & x \geqslant 0 \end{cases}$

8) log of Sigmoid $\rightarrow$ $y = log\left(\dfrac{1}{1+e^{-x}}\right)$

9) Swish $\rightarrow$ $y = \dfrac{x}{1+e^{-x}}$

10) Sinc $\rightarrow$ $y = \dfrac{Sin(x)}{x}$

11) Leaky Relu $\rightarrow$ $y = Max(0.1x, x)$

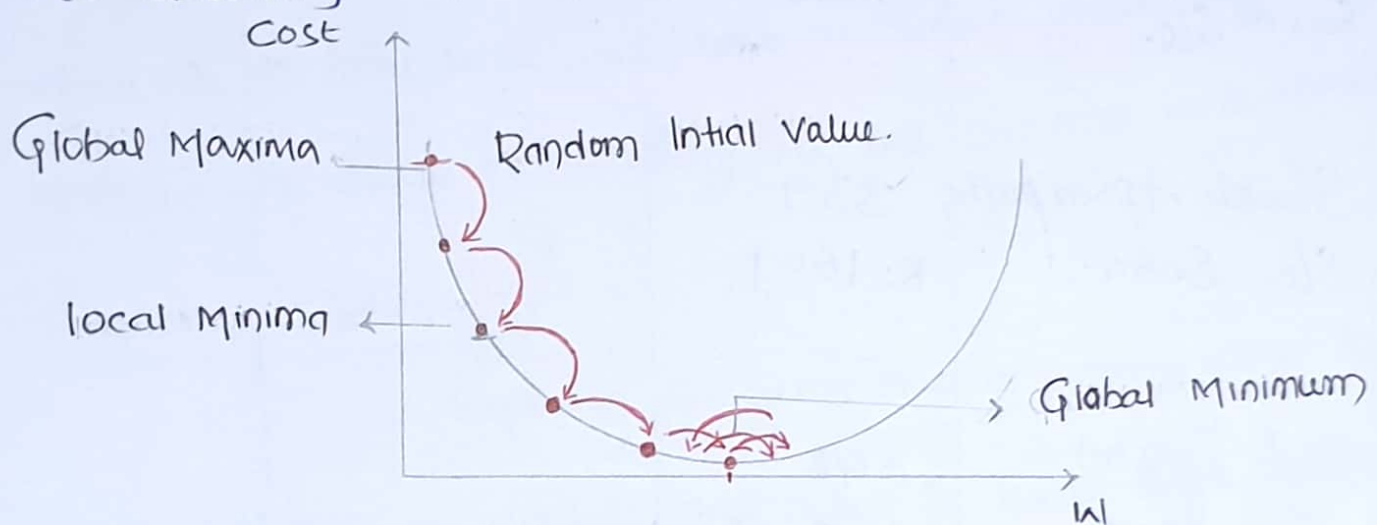12) Mish $\rightarrow$ $y = x(Tanh(Soft\ plus(x)))$

# Gradient descent :-

GD is an optimization algorithm for finding a local minimum of a differentiable function, it is used in machine learning to find a values of a functions parameters (Coefficents) that minimize a cost function as far as possible.

as well this algorithm works in neural networks, Training data helps these models learn over time and the cost function with in gradient descent specially acts as a barometer, guiding its accuracy with each itaration of parameter updates.
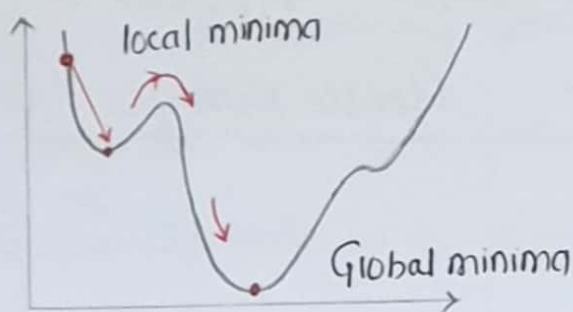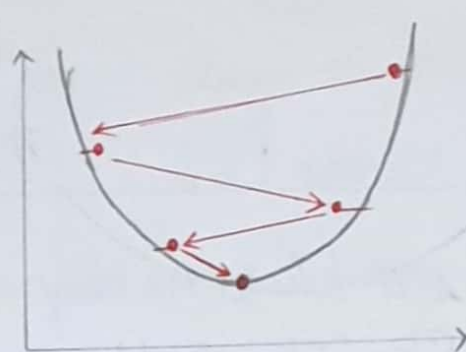


This jumps may happen if the learning rate is very high.

Vanishing Gradient descent :- it is deffect (cause) in Gradient descent, it refers to the fact that in a feed forward network (FFN). the back propagated error Signal typically decreases or increases exponantially. as a function of the distance from the final layer.

[ Vanishing G.D.]    [ Exploding Gradient ]

In G.D. vanishing G.D. is a deffect we have to overcome that

Cause of local value, when we iterate the value, then generated

new value. reaches the minimum point then after no improvement

in output, the point thinks to be the Global value reaches but

that is a local value, we have to reach global value. this

deffect is called as vanishing G.D., but our task is to overcome it.

Exploding Gradient :-

Exploding Gradients are a problem where large

error gradients accumulate, and result in heavy large updates to

neural network model weights during training. this has the effect

of our model being. unstable and unable to learn from training data.

By Normalization we explod the Gradients (stop).

→ common solution to exploding Gradients is to change the error

derivative before propagating it backward through the network

and using it to update the weights.

**Note :-** In vanishing Gradient, the weights change are negligible and hence the learning is poor.

We can identify by the change in loss for every iteration.

The Exploding Gradient makes an exponential increment in weights / higher jump in weight values, which impacts the learning and loss are infinet or losses are very very high.

Exploding Gradient Can be identified as nan in loss value.

**Important Notes on GD :-**

→ In Optimization, the main aim is to find weights that reduce loss.

→ Gradient is caliculated by Optimizing function.

→ Gradient is the change is loss with change in weights.

→ The weights are modified according to the caliculated gradients.

→ Same process is repeated untill minima is reached.

**Learning Rate and Momentum :-**

Learning Rate is a hyper parameter to what extant newly aquired weights overrides the existing weights, lies between 0 and 1.

Momentum is used to decide the weights on nodes from previous iterations, it helps in improving training speed and also avoiding local minimas.

# Types of Gradient Descent :-

They are mainly classified as 3 categories. are —

① Batch Gradient Descent

② Stochastic Gradient Descent

③ Mini - batch Gradient Descent.

## Batch Gradient Descent :-

In this we will caliculate the Cost function for all Observations and update the weights to minimize the Cost-function this is called Batch Gradient Descent.
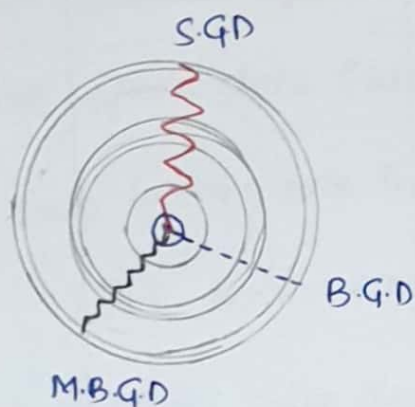
## Stochastic gradient Descent :-

Here we caliculate the Cost-function for each Obsenation and updates the weights to minimize the Cost-function this is called as Stochastic Gradient descent.

it takes the dataset in Random and negative with FWP and BP for each epochs (iterations), the main disAdvantage is it make the result more bias.

## Mini - Batch Gradient Descent :- In this we will first divide the data in to mini batches (5 or 8 Observations for batch) then caliculate the Cost -function for each mini observation and update the weights to minimize the Cost-function.

This phenominum is called as " Mini-batch G.D."

It navigates based upon steps per epoche/batch Size of the G.D.



Momentum → Improving the noise from S.G.D, Mini batch SGD /from G.D.

$$W_{new} = W_{old} - \eta \frac{dL}{dW_{t-1}}$$

Exponential weighted Average

$t_1, t_2, t_3 \cdots \cdots$   $W_{new} = W_{old} - \eta \frac{dL}{dW_{t-1}}$

$a_1, a_2, a_3$

$\qquad\qquad = (1-\beta) * a_2$

$\qquad\qquad = \beta * a_1 + (1-\beta) * a_2$

$[\beta = 0.1] = 0.1 * a_1 + (1-0.1) * a_2$

$\qquad\qquad = 0.1 * a_1 + 0.9 * a_2$

Major differences in G.D :-

Stochastic G.D.

1) Local minima problem is can be resolved
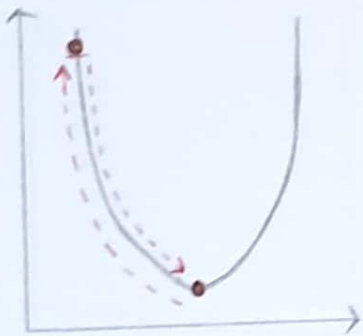
2) fast

3) Less Deterministic

Batch G.D.

1) may end up local minima.

2) Slow

3) More Deterministic

# Adaptive Gradient Descent (Adagrad) :-

It is an algorithm for gradient based Optimization

it performs Smaller updates as result, it is well Suited when dealing with Sparse data (NLP or Image recognization) each parameter has its own learning rate that improves performance of problems.

It is an extention of Gradient decent optimization algorithm that allows Step Size in each dimention used by the optimiz-ation algorithm to be automatically adapted based on the gradients Seen for the variables (partial derivatives).

In adagrade, the learning rate is not fixed and it is Subjected to change based upon the loss.

$$W_{new} = W_{old} - \eta' \frac{dL}{dw_{t-1}}$$



$$\eta' = \frac{\eta \text{ (intial learning rate)}}{(\alpha_t + e)}$$
$\quad\quad\quad\quad \hookrightarrow$ constant

$$\alpha_t = \sum_{i=1}^{t} \left(\frac{dL}{dw_t}\right)^{2}$$
$\quad\quad\quad\quad\quad\quad \hookrightarrow t = W_{11}$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad W_{12}$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad W_{13}$

$$W_{new} \approx W_{old}$$

# RMSprop Optimization :-

It is an gradient based Optimization technique used in training neural networks. this normalization balances the Step Size (Momentum) decreasing the Step for large gradients to avoid exploding and increasing the Step for small gradients to avoid vanishing.

RMSprop Stands for root mean Square prop., which can also accelarate gradient descent, it uses the Same Concept of the exponentially weighted. average of gradient descent with momentum but the differences is parameter update.

$$W_{new} = W_{old} - \eta'' \frac{dL}{dW_{old}}$$

$$\eta'' = \frac{\eta}{\sqrt{Sd_{wt} + \varepsilon}}$$

$$Sd_{wt} = \beta \cdot Sd_{wt-1} + (1-\beta)\left(\frac{dL}{dw_t}\right)^2$$

↳ previous

# Adam Optimization :— 

Adam is a replacement optimization algo for Stochastic GD. for training deep learning models. Adam Combines the best properties of the Adagrad and RMSprop. algorithms to provide an optimization algo that can handle spares gradients.

Adam is best among the adaptive optimizers in most of the cases., Good with Spare data.

Adam optimizer well Suited for large data sets and its Computationally efficient, there are few disfadvantages with it as the adam optimizer tends to coverage faster, but Other algorithms like the. Stochastic. G.D. focus on the data points and generalize in a better manner.

$$W_{new} = W_{old} - \eta''' \frac{dL}{dW_{old}}$$

where $\eta'''$ — momentum, RMSprop.

$$W_{new} = W_{old+1} - \frac{\eta * V_{dw}}{\sqrt{Sdw + \varepsilon}} \begin{array}{l} \rightarrow momentum \\ \rightarrow RMSprop. \end{array}$$

Momentum $\therefore$ $\boxed{V_{dw} = \beta V_{dw} + (1-\beta) \frac{dL}{dw}}$

Number of Weights :-



256 Weight

26 Weights

input image

16 x 16

256 x 100

100 Hidden Unit

26x100+26

Nodes

$= 25600 + 100 + 2600 + 26$

$= 28326,$