# Today's Agenda

## Perceptron Model

$$x \longrightarrow \boxed{wx + b} \longrightarrow \hat{y}$$

Neuron

$\hat{y}$ = prediction

$y$ = actual

$$x = \begin{Bmatrix} x_0 \\ x_1 \\ x_2 \end{Bmatrix}$$

$$x_1 \xrightarrow{w} \bigcirc \nearrow wx + b$$

$$\uparrow b$$

## Loss Calculation

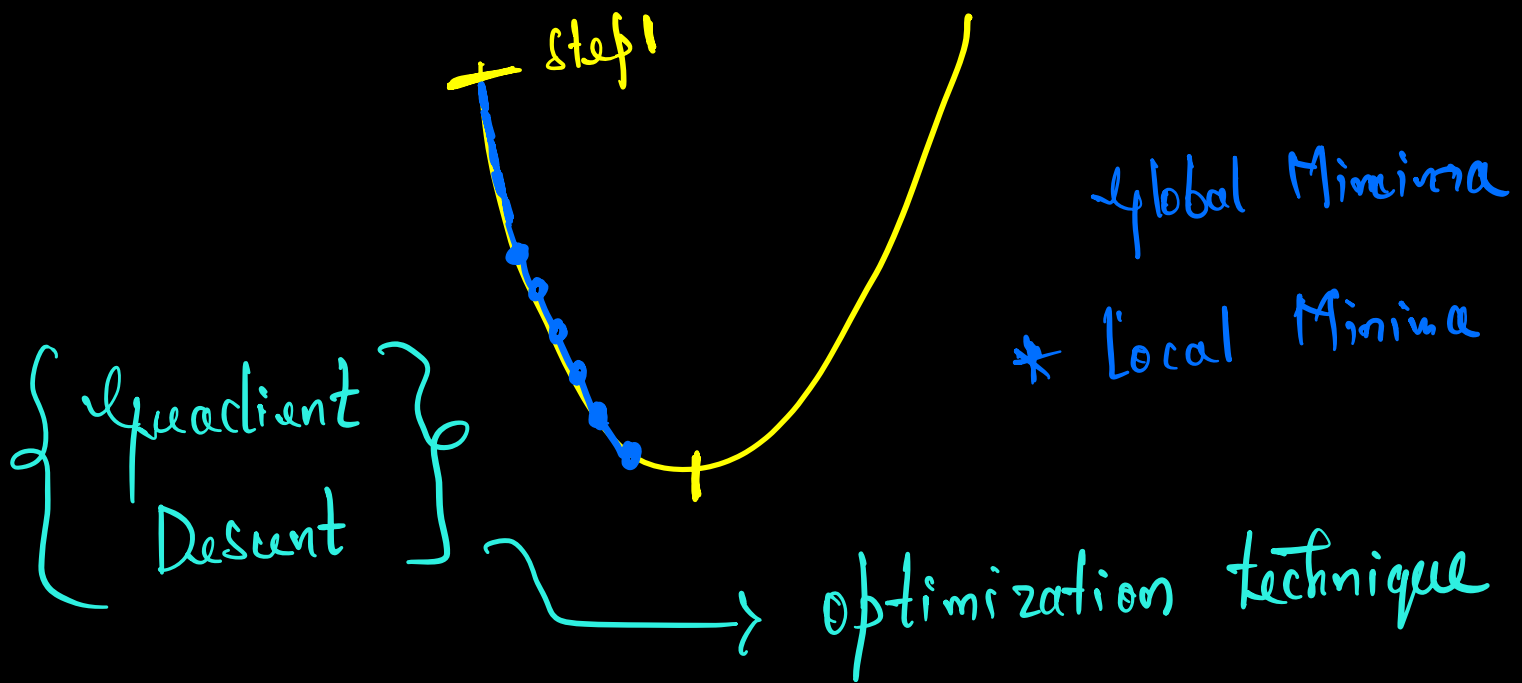$$\text{Squared Error} = \left( \hat{y} - y \right)^2 = \text{Loss}$$

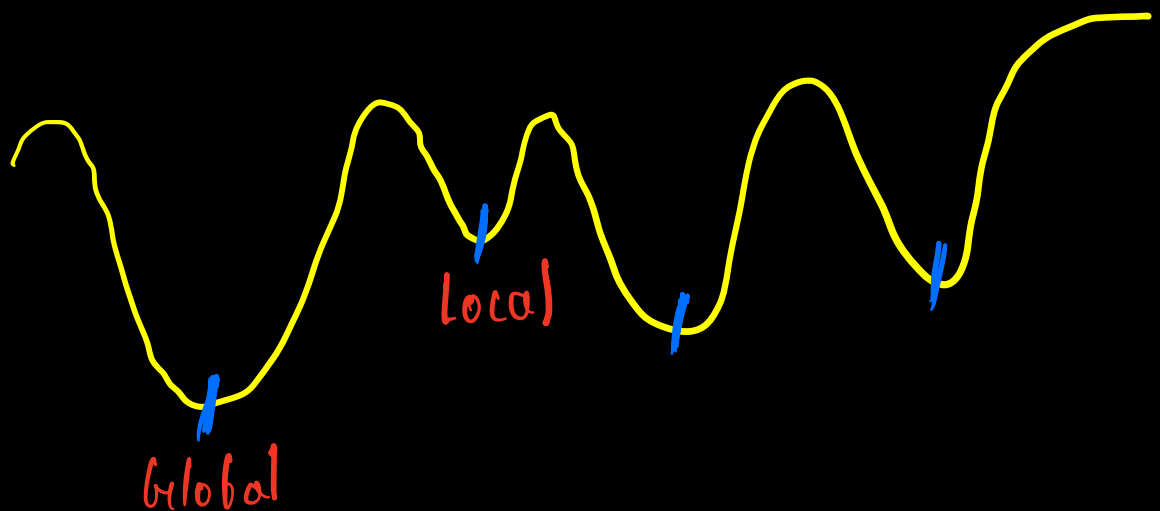Error $\longrightarrow$ Reducing the error close to O

$x \rightarrow$ Fixed

$\{w, b\} \longrightarrow$ change / learnable params

Tune $\{w, b\}$ so that the loss comes

close to Zero.

Update $w$ & $b$
_____

step1

Global Minima

* local Minima

$\{$ Gradient Descent $\}$ $\longrightarrow$ optimization technique

$w$ & $b \approx loss (0)$ close

local

Global

## W & b Update Rule

$L = E$
Loss = Error

$$W_x = W_{x'} - \eta \left\{ \frac{\partial L}{\partial W_x} \right\}$$

L·R → Learning Rate
(0.01, 0.001, 0.000$

$W_x$ → New weight

$W_{x'}$ → old weight

$$\{ \cdot 0001 \}$$

$$\frac{\partial L}{\partial W_{x'}} \quad \text{ore} \quad \frac{\partial E}{\partial W_{x'}}$$
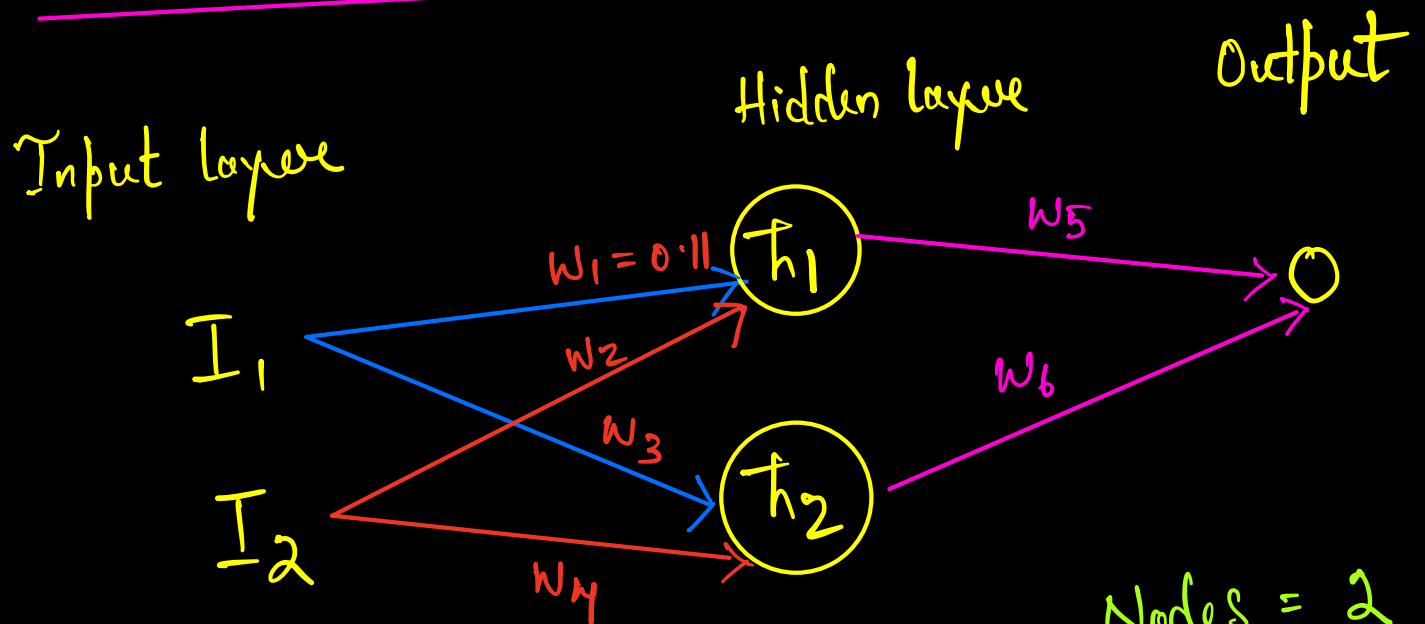
→ Derivative of Error wort weight

## Bias

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

L·R → hyperparameter

## Weight Update Rule

# Forward Propagation

Input layere          Hidden layere          Output



$W_1 = 0.11$

$\overline{h_1}$

$W_5$

$O$

$W_2$

$W_3$

$\overline{h_2}$

$W_6$

$I_1$

$I_2$

$W_4$

Nodes = 2

$W_4 = 0.08$

$W_1 = 0.11$

$W_5 = 0.14$

$W_2 = 0.21$

$W_6 = 0.15$

$W_3 = 0.12$

No of nodes
⇓
hyperparameter

$$X = [i_1, i_2] = [2, 3]$$

$$y = 1$$

No of layers
⇓
hyperparameter

## Forward Propagation

(i)    Input ⟶ Hidden layere

$$X \qquad W$$

$$[2, 3] \cdot \begin{bmatrix} \cdot 11 & \cdot 12 \\ \cdot 21 & \cdot 08 \end{bmatrix}$$

$h_1 = \quad 2 \times \cdot 11 + 3 \times \cdot 21 = 0.85$

$h_2 = \quad 2 \times \cdot 12 + 3 \times \cdot 08 = 0.48$

$$= \quad [0.\overset{h_1}{85} \quad \overset{h_2}{0.48}]$$

(ii) <u>Hidden Layer</u> $\longrightarrow$ Output layer

$$= \quad [0.85 \quad 0.48] \begin{bmatrix} 0.14 \\ 0.15 \end{bmatrix}$$

$$= \quad 0.85 \times 0.14 + 0.48 \times 0.15$$

$$= \quad [\cdot 191]$$

<u>Calculate Errors</u>

Input diagram: nodes 2 and 3 connect via weights $W_1$, $W_2$, $W_3$, $W_4$ to hidden nodes $h_1$ and $h_2$. $h_1 \to 0.85$, $h_2 \to 0.48$. Weights $W_5$, $W_6$ connect to Out node. $(0.19) = \hat{y}$

## Error Calculation

$$\text{Loss Function} = \frac{1}{2}(\hat{y} - y)^2$$

$$= \frac{1}{2}(0.19 - 1)^2$$

$$\text{Error} = 0.327$$

## Reducing the Error

$$\text{Prediction} = \text{output}$$

$$= (h_1)W_5 + (h_2)W_6$$

Phase, Output layer $\longrightarrow$ Hidden layer

Reverse

$$h_1 = i_1 w_1 + i_2 w_2$$
$$h_2 = i_1 w_3 + i_2 w_4$$

$$\text{Prediction} = \left( i_1 w_1 + i_2 w_2 \right) w_5 + \left( i_1 w_3 + i_2 w_4 \right) w_6$$

To change prediction value, weights.

How to change $\longrightarrow$ Backpropagation

1) Reduce the error

2) Change the weight

Weight Update $\longrightarrow$ Gradient Descent

$\left\{ \begin{array}{l} \text{Iterative Optimization algorithm for} \\ \text{finding the minimum of a function.} \end{array} \right\}$

# Backprop with Gradient Descent

Weight Update :-

$$W_{new} = W_{old} - \eta \left( \frac{\partial E}{\partial W_{old}} \right)$$

As per our Network

$$* \; W_6 = W_6 - \eta \left( \frac{\partial E}{\partial W_6} \right)$$

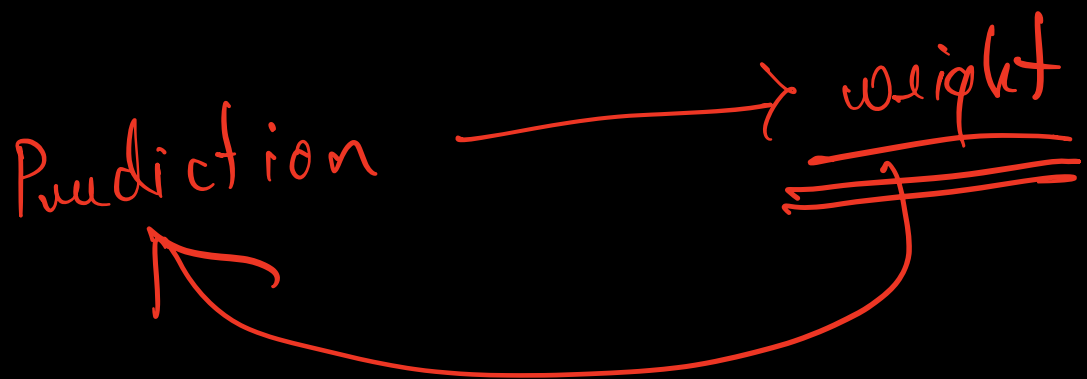$$* \; W_5 = W_5 - \eta \left( \frac{\partial E}{\partial W_5} \right)$$

$$W_6 =$$

$$\eta = .001$$

* Chain Rule *

$$\frac{\partial Error}{\partial w} = \frac{\partial Error}{\partial prediction} * \frac{\partial prediction}{\partial W_6}$$

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

$$\underbrace{\qquad\qquad}_{\{\text{prediction}\}}$$

Prediction $\longrightarrow$ weight

## loss vs Cost Function

1) MSE (Regression)

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

$\frac{1}{n}$ → no of samples.

$y$ → actual
$\hat{y}$ → predicted

2) log loss (Binary Cross Entropy)

$$\text{Logistic Regression}$$

$$loss = -\frac{1}{n} \sum_{i=1}^{n} \left[ Y \log(\hat{y}) + (1 - Y_i) \log(1-\hat{y}) \right]$$

## Loss Function

1) Sample

## Cost Function

1) Dataset

Average loss on the entire dataset

## Gradient Descent

Optimization Algorithm

Update
$$\theta = \theta - \alpha \frac{\partial J(\theta)}{\partial \theta}$$

$\theta(w, b)$

$\alpha - L.R$

$J(\theta)$ : cost function

$w_0 = 0.6$

$\alpha = 0.001$

$w_{new} = w_{old} - 0.001 \frac{\partial E}{\partial w}$

$$= 0.6 - (0.001 \times 5)$$

$$= 0.6 - ( \qquad )$$

$$w_{new} = ?$$

Gradient Descent Types
_____

* 1) Batch Gradient Descent : Full Dataset
for every update

* 2) SGD ( Stochastic Gradient Descent )

* 3) Mini Batch Gradient Descent

SGD → One data point at a time

(SLOW) Very slow

* Mini - Batch :- A small batch of dataset
_____

Batch size :- $(2^n)$ 8, 16, 32, 64, 128, 256,
512, 1024

↓

hardware

Framework :- Keras / Pytorch
_____

→ Mini Batch Gradient Descent.

(SGD)

# Data Terminologies

1) Batch :- Subset of your data

2) Iteration :- 1 update = 1 batch
   weight

Batch :- 1000 training sample

Batch Size :- 100

Total No of batches :- 10

Iterations :- $\dfrac{10 \text{ iterations}}{}$

3) Epoch :- full cycle of the entire
   dataset

All students taught in a class → 1 epoch

1 epoch = Network has seen every

Total Samples :- 1000

Batch size :- 100

Total Batch :- 10

Total Iteration :- 10

1000

Epochs = 1 , 10 Iteration / 10 Batches

1000

Epoch = 2 ) 20 Iteration / 10 Batch

Epoch = 10 , 100 iterations

4) Step

1 step = weight update

DL Frameworks :- Tensorflow & Pytouch

1 weight update :- At every batch, the weight update takes place.