

# Optimization Techniques





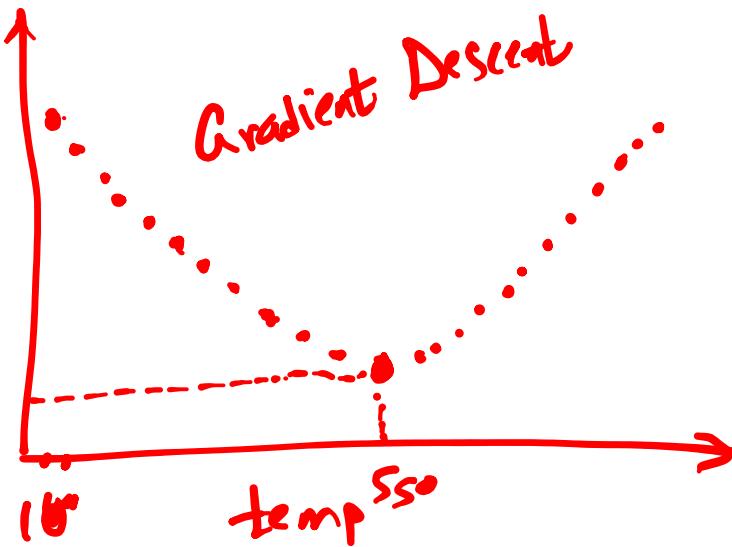
Trainer: Dr. Darshan Ingle.



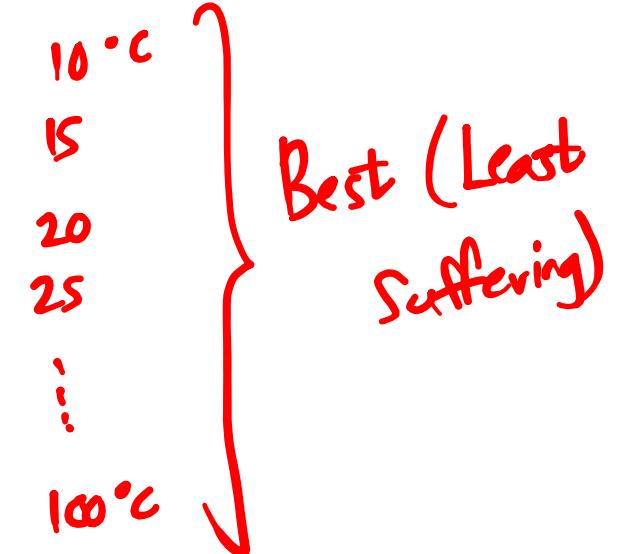
Trainer: Dr. Darshan Ingle.

## Exhaustive Search

Suffering

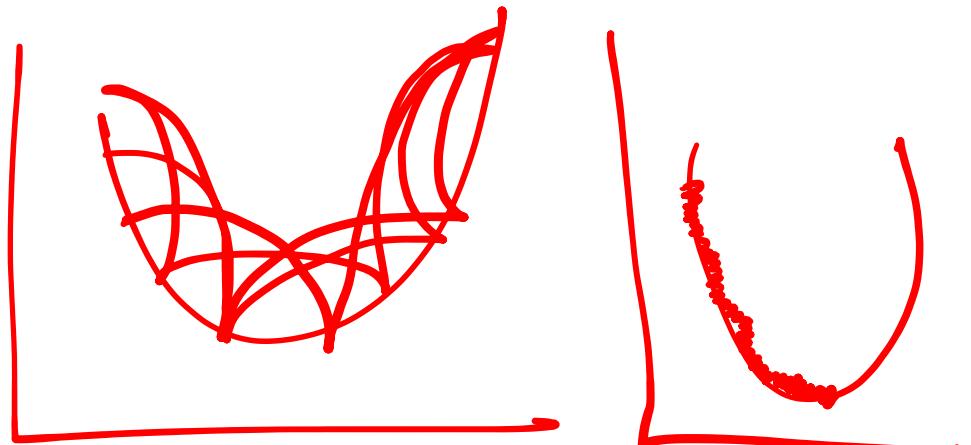
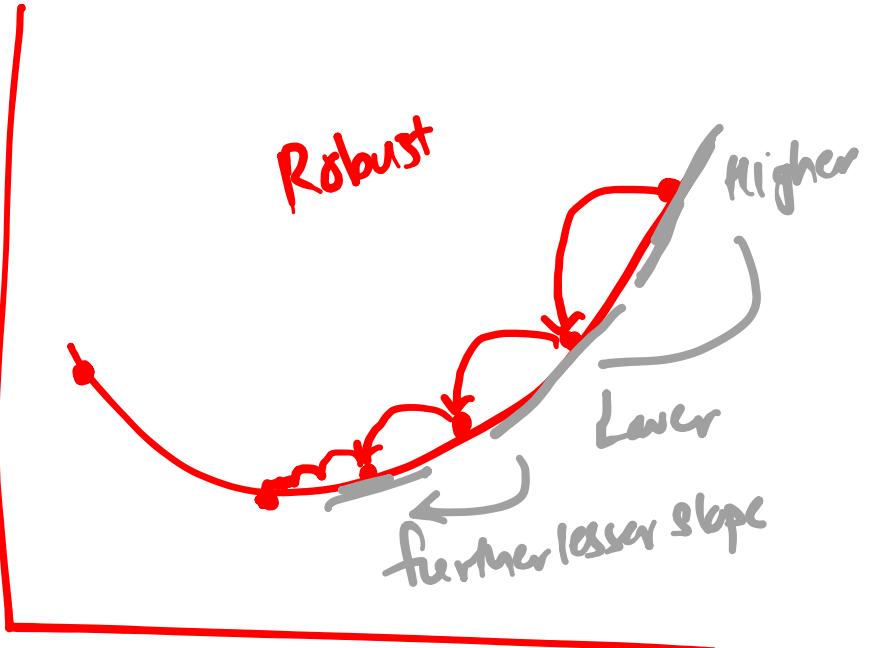


Adv: Simple Global optimum



Disadv: 1. Computationally expensive.

2. Get close to real min bcz of predefined step size, but we never achieve minimum



from sklearn.linear\_model import Logistic

obj = Logistic() #  $\downarrow$  method='Gradient'

obj.fit(x-train, y-train)

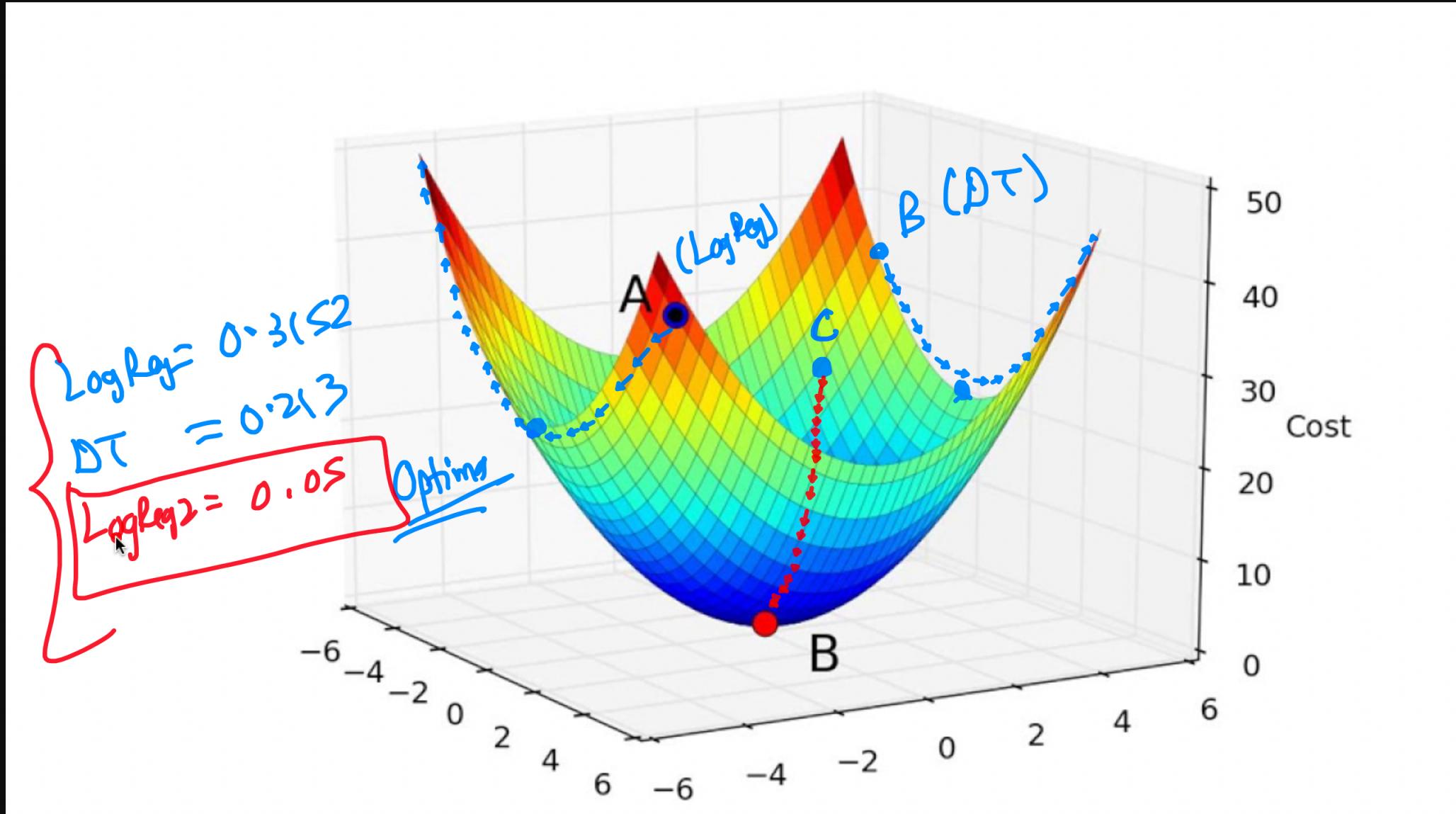


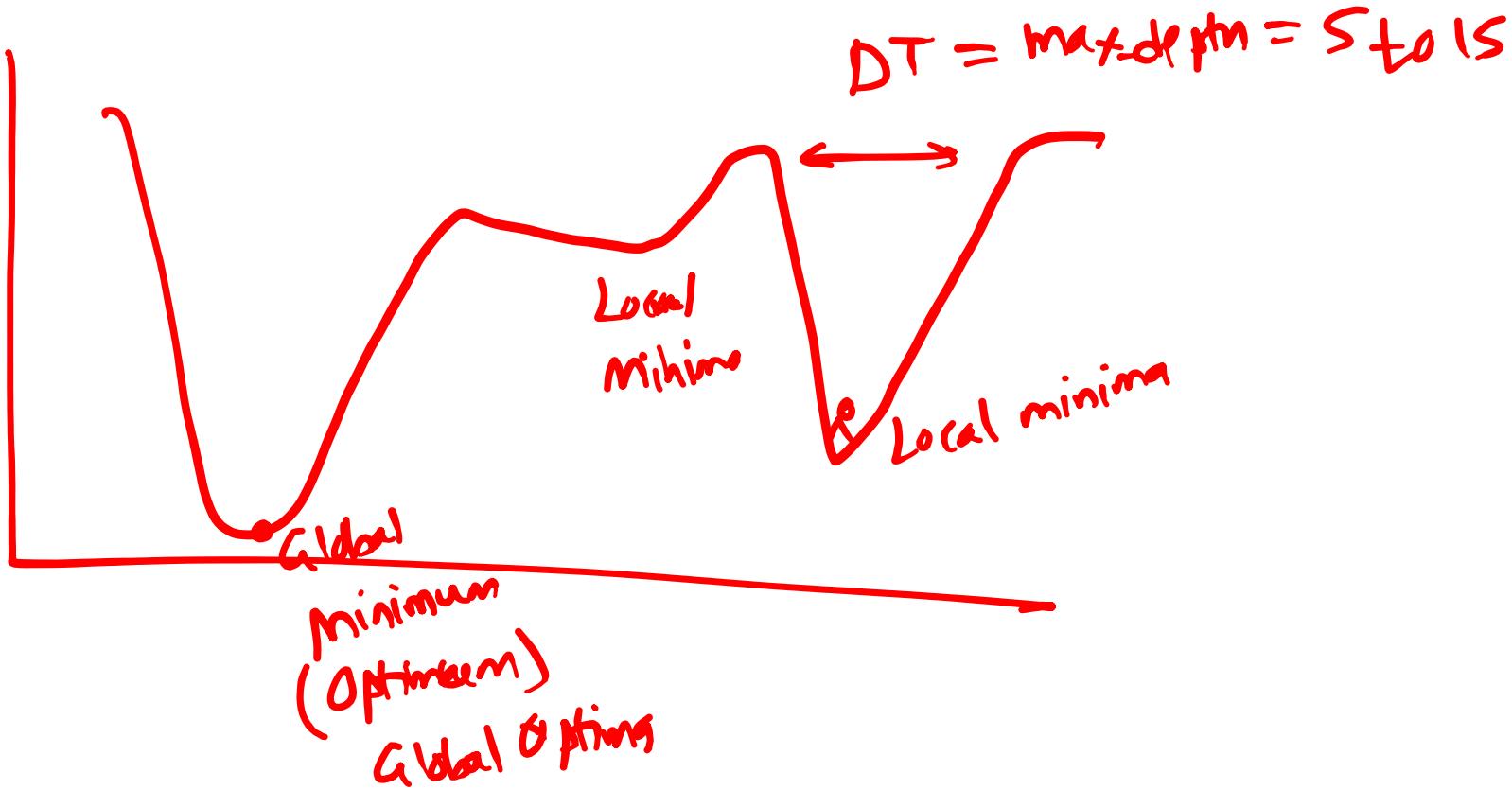
If slope is high we jump higher

If slope is less, we jump closer.

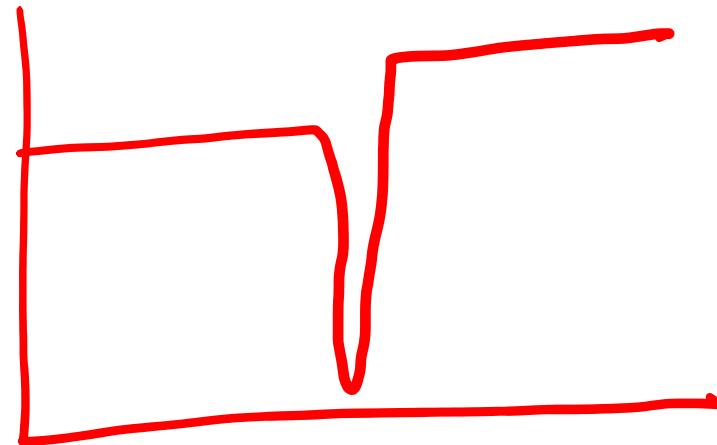
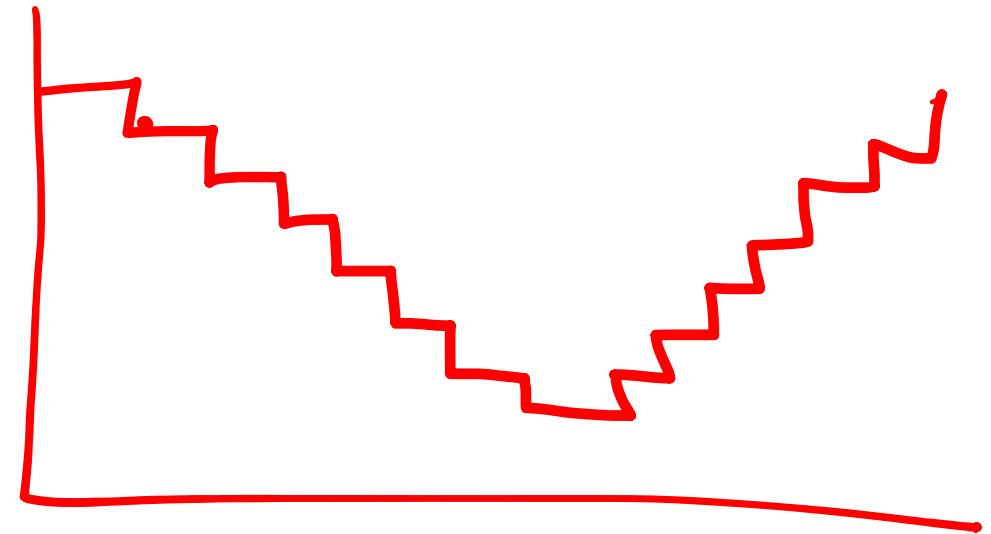
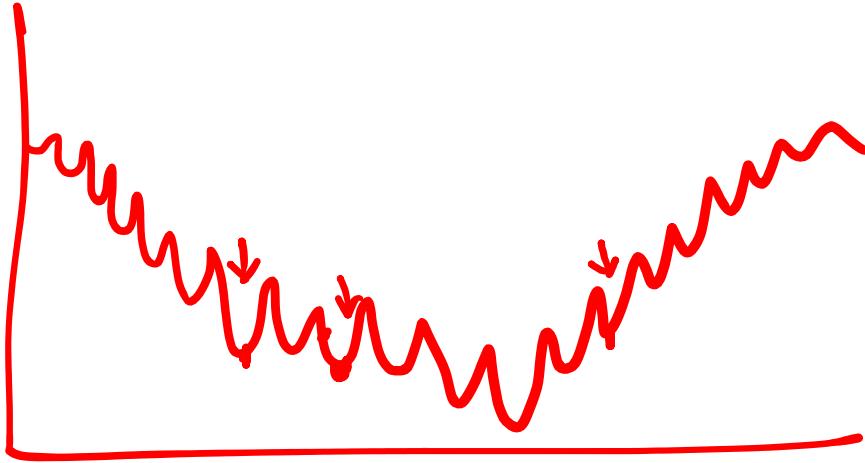
Step Size  $\propto$  Slope

$$SS = \text{Learning rate} \times \text{Slope}$$



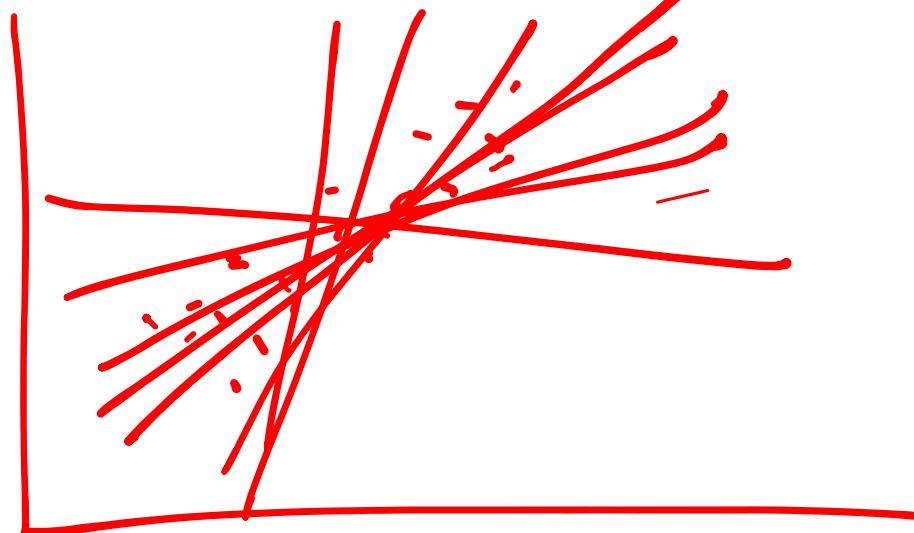


Sometimes error function is not smooth. There are lots of places a marble could get stuck.



Cost Function : We look at the loss associated with each training example & then sum these values together for an overall cost of the entire dataset.

Batch Gradient Descent:



$y$   $\left| \begin{array}{c} \text{TV} \\ 0 \end{array} \right|$   $\begin{array}{l} \text{red} \\ \text{o} = \text{error} \end{array}$   
LOBF

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

all n data points

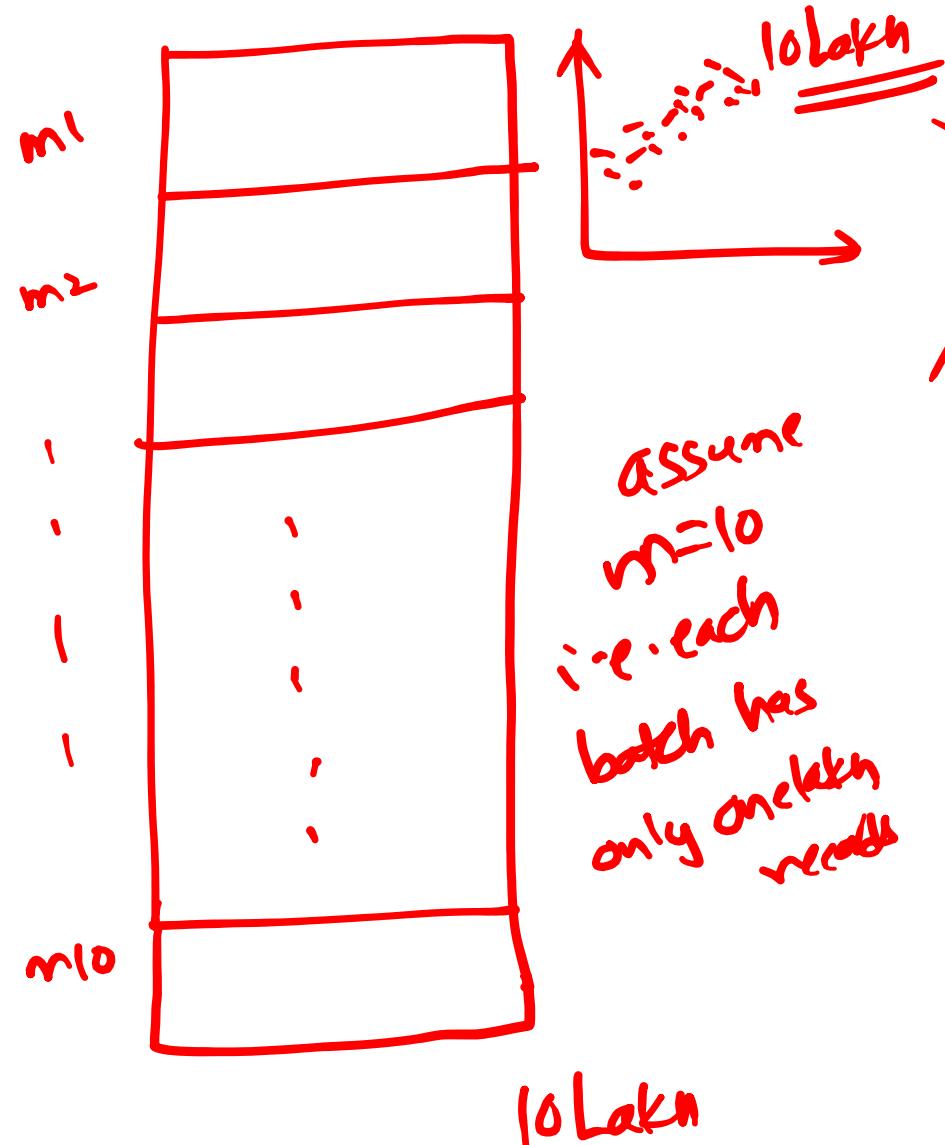
assume  $n = 10\text{baths}$

$e_1$   
 $e_2$   
 $e_3$   
 $\vdots$   
 $e_{25}$

least

Disadv: Slow & Comp. expensive.

## Mini - Batch Gradient Descent :

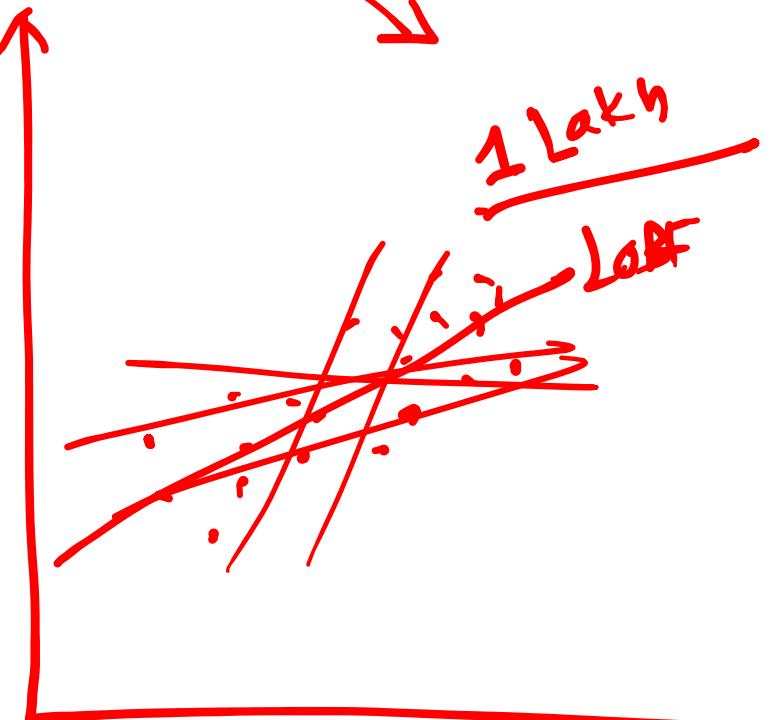


$m$  mini batches

Usually  $m = 5 \text{ to } 10$

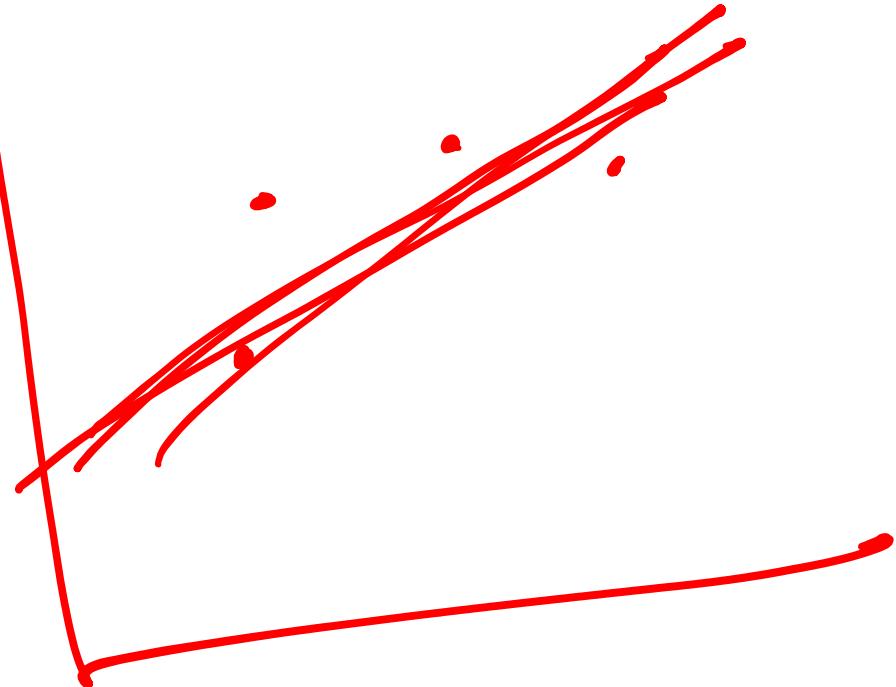
$m \geq 1$

Considering one  
batch &  
apply G.D.  
the same way.



Now the computation will  
be done on  $\frac{1}{m}$  th dataset.  
∴ its faster

Stochastic Gradient Descent : It considers any one record  
for fitting a line  $\Delta$  repeats this  
process for said no. of iterations





Trainer: Dr. Darshan Ingle.

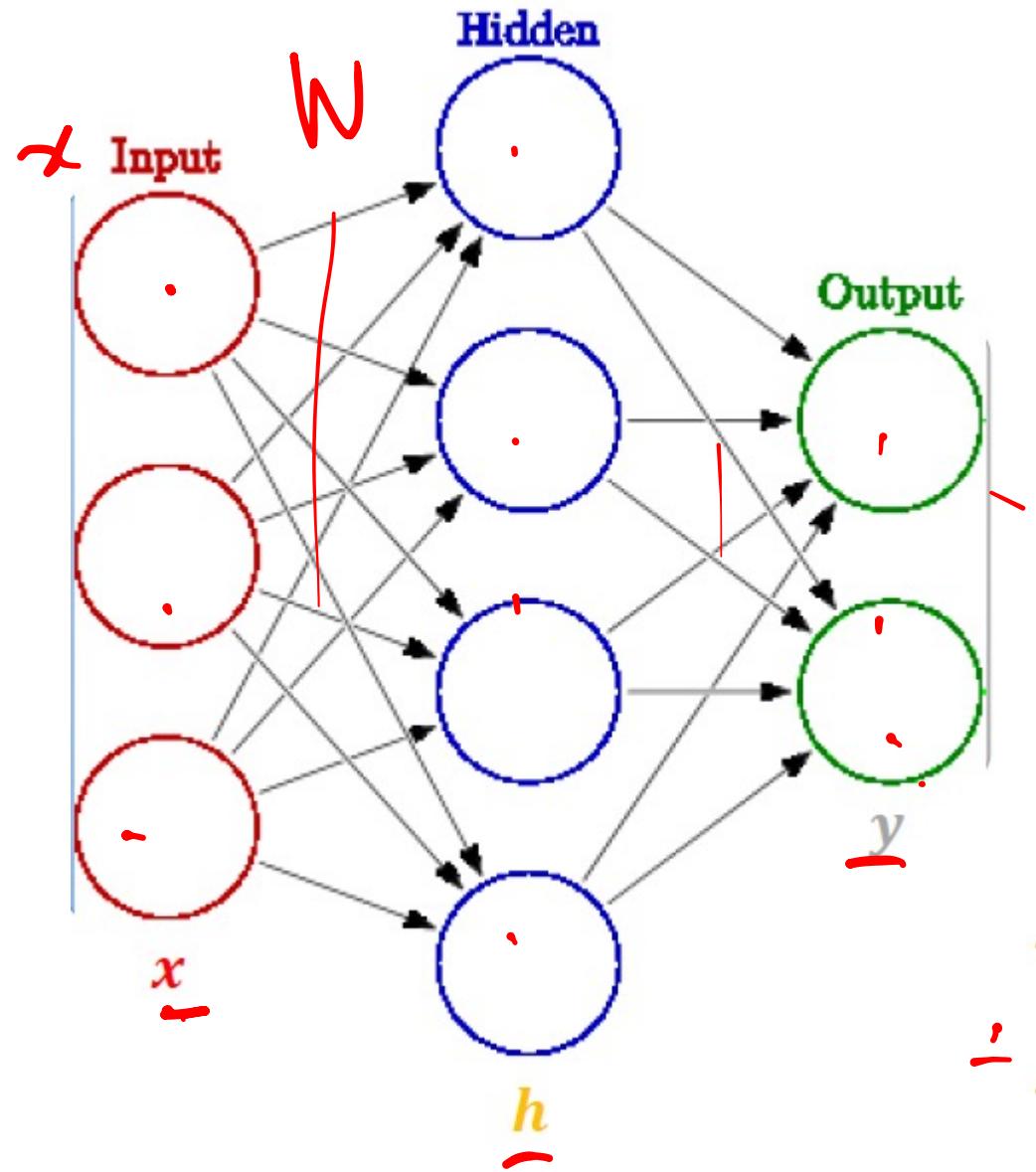


Trainer: Dr. Darshan Ingle.



Trainer: Dr. Darshan Ingle.

# Neural Network Learning Parameters



$$h = \sigma(W_1 x + b_1)$$
$$y = \sigma(W_2 h + b_2)$$

Activation function

$4+2=6$  neurons (not counting inputs)

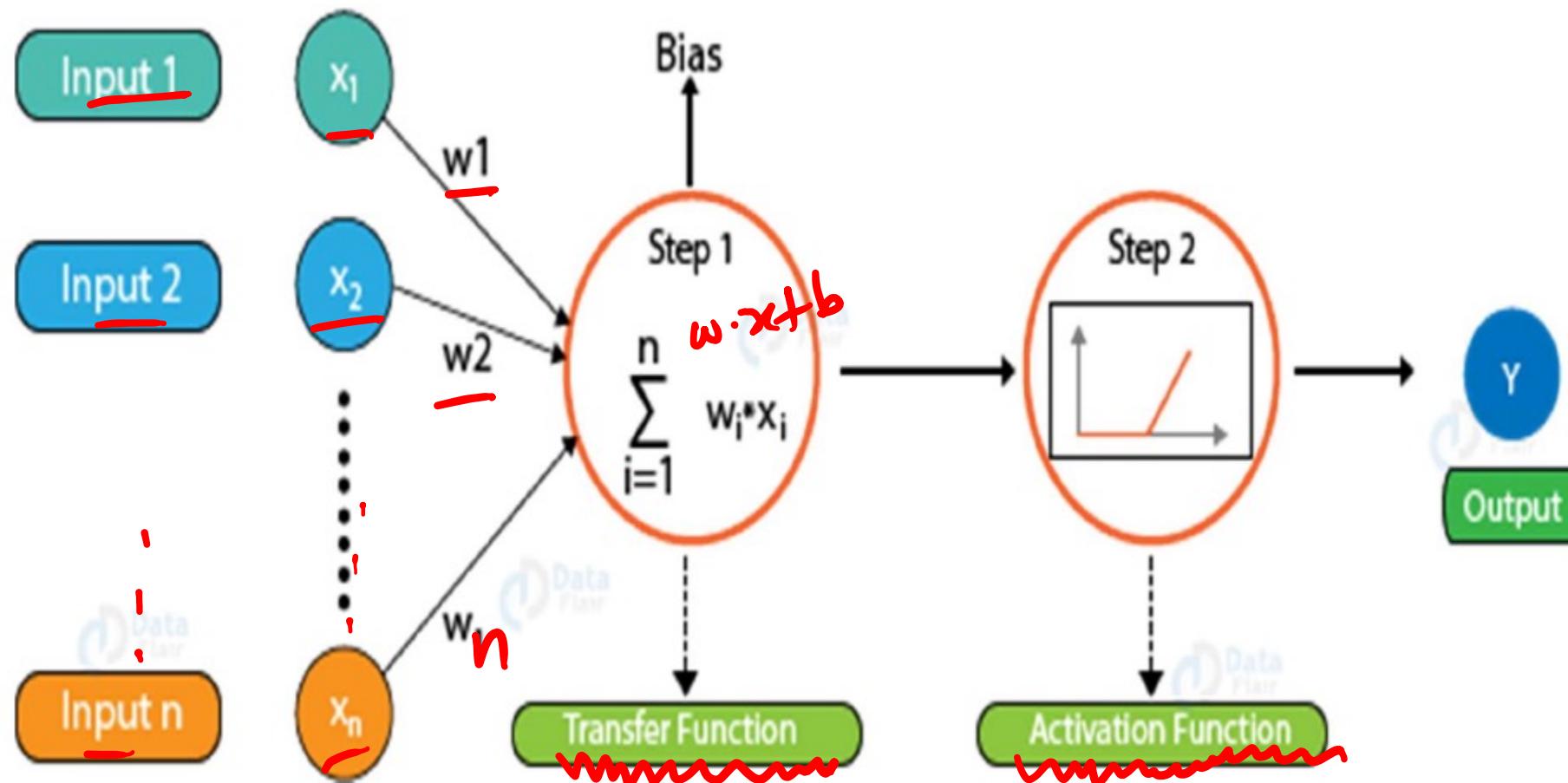
$$[3 \times 4] + [2 \times 2] = 12 + 8 = 20 \text{ weights}$$

$$4+2=6 \text{ biases}$$

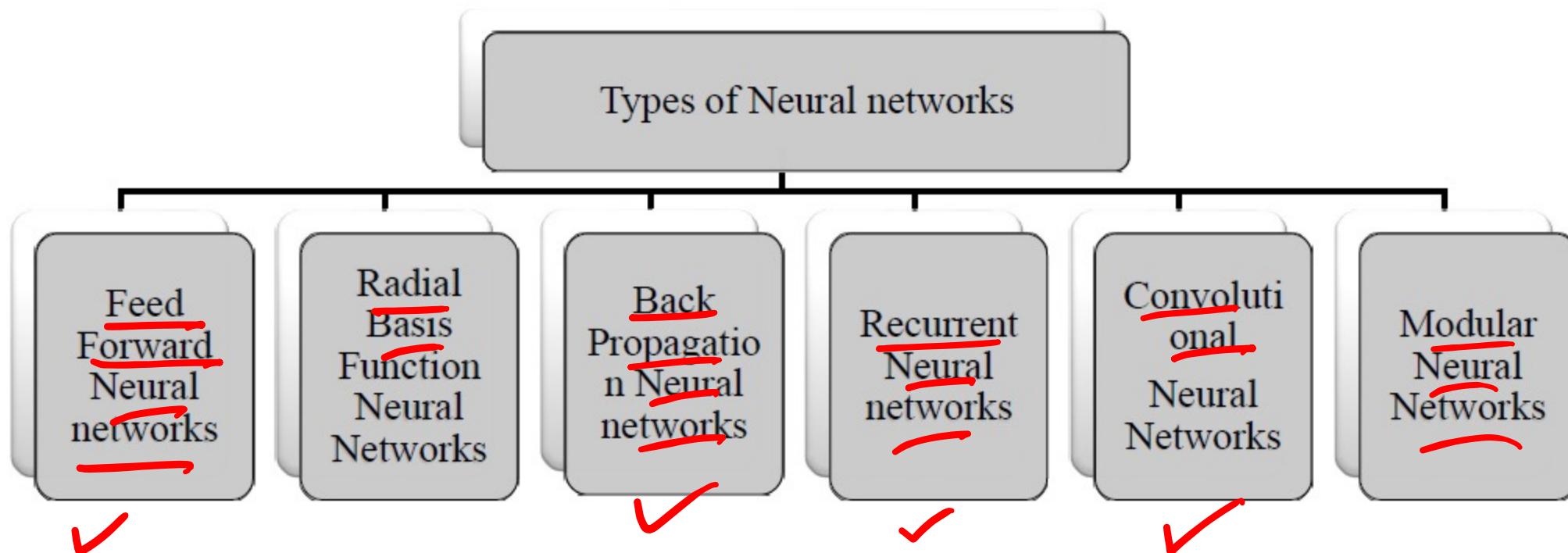
$$\therefore 20+6=26 \text{ Learnable parameters}$$

# More Terminologies of a NN

$w \cdot x + b \rightarrow$  Transfer Function

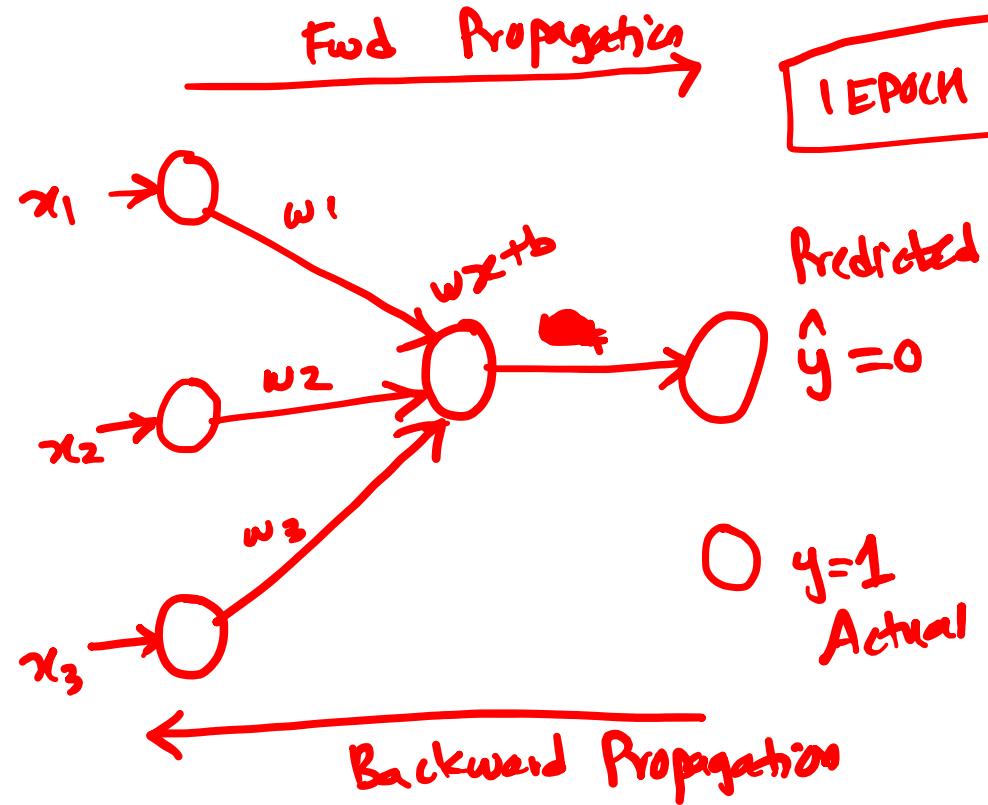


# Types of Neural Network



In real life, weights are initialized RANDOMLY.

# Forward and Backward Propagation



Note: ① Single o/p Node = Loss Function

② Multiple o/p = Cost Function

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

1 EPOCH = 1 FP + 1 BP

Min. Loss  $\Rightarrow$  Optimizers

$$\begin{aligned} \text{Loss} &= (y - \hat{y})^2 \\ &= (1 - 0)^2 \\ &= 1 \end{aligned}$$

Formula!

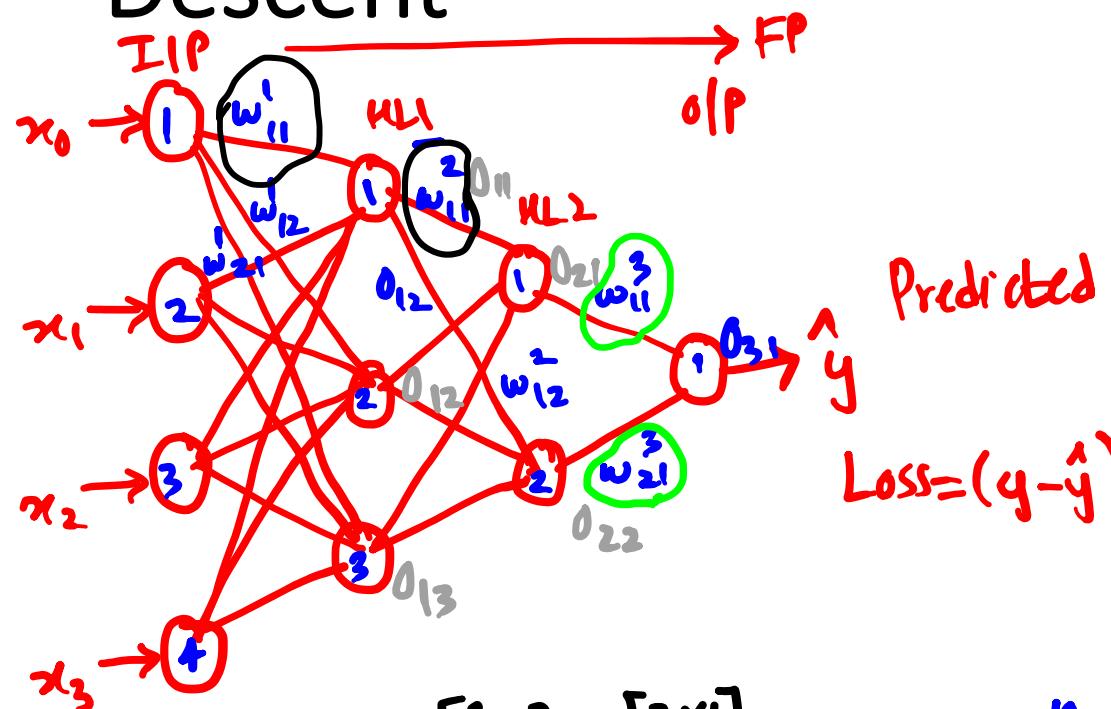
$$w_{i, \text{new}} = w_{i, \text{old}} - \eta \cdot \frac{\partial L}{\partial w_{i, \text{old}}}$$

$$\eta = \text{Learning rate. Assume } \eta = 0.01 [0-1]$$

$x_1$ Res.	$x_2$ Age	$x_3$ Score	o/p Grade
80	70	83	1

$$\begin{aligned} y &= [w_1 x_1 + w_2 x_2 + w_3 x_3] + b \\ z &= \sigma(y) \quad [\text{Sigmoid}] \end{aligned}$$

# Multi Layer Neural Network Training w.r.t Gradient Descent



Predicted

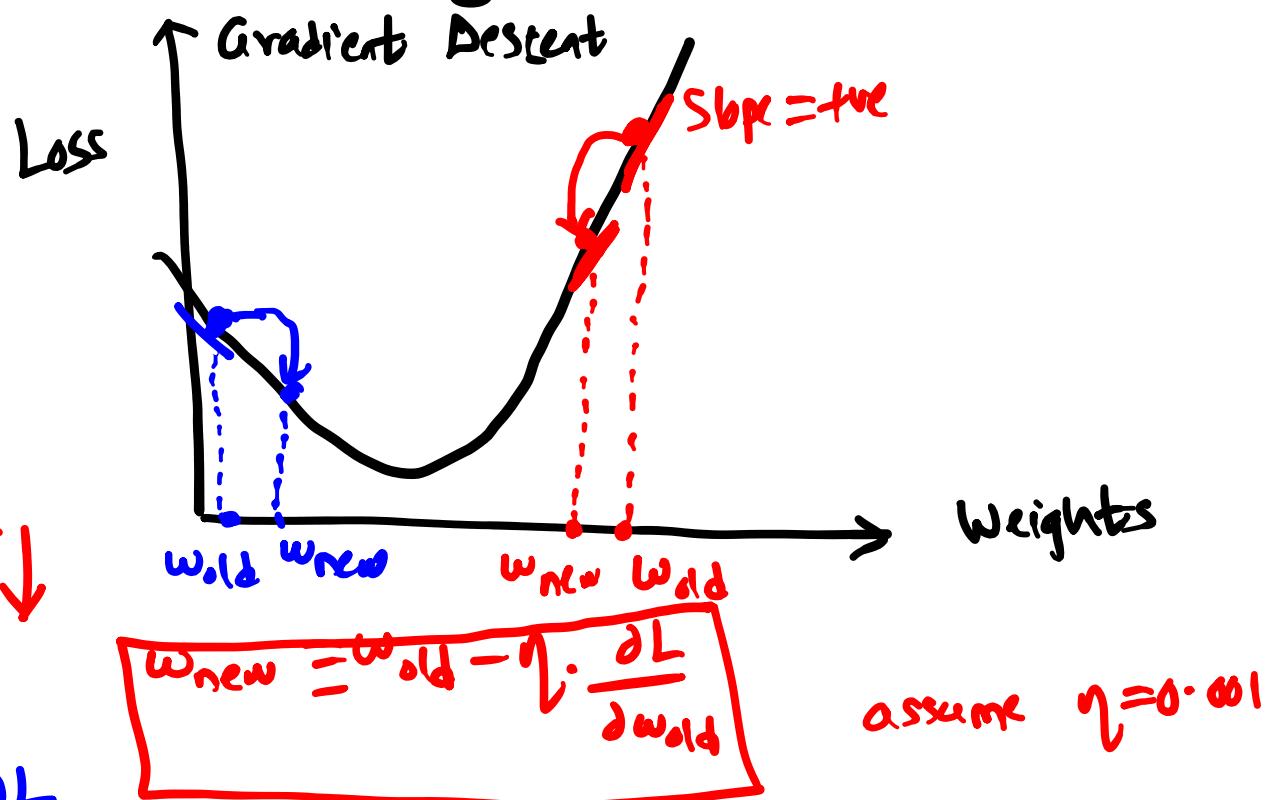
$$\text{Loss} = (y - \hat{y})^2 \downarrow$$

$$\text{Input } [4 \times 3] \quad [3 \times 2] \quad [2 \times 1] \quad w_{\text{new}} = w_{\text{old}} - \eta \cdot \frac{\partial L}{\partial w_{\text{old}}}$$

$$\begin{bmatrix} w_{11}^1 & w_{12}^1 & w_{13}^1 \\ w_{21}^1 & w_{22}^1 & w_{23}^1 \\ w_{31}^1 & w_{32}^1 & w_{33}^1 \\ w_{41}^1 & w_{42}^1 & w_{43}^1 \end{bmatrix}_{4 \times 3} \quad \begin{bmatrix} w_{11}^2 & w_{12}^2 \\ w_{21}^2 & w_{22}^2 \\ w_{31}^2 & w_{32}^2 \end{bmatrix}$$

$$w_{\text{new}} > w_{\text{old}}$$

Trainer: Dr. Darshan Ingle.



$$= w_{\text{old}} - \eta \cdot (\text{tvc})$$

$$= w_{\text{old}} - (\text{tvc})$$

$$w_{\text{new}} < w_{\text{old}}$$

$w_{11}^3_{\text{new}} = w_{11}^3_{\text{old}} - \eta \cdot \frac{\delta L}{\delta w_{11}^3_{\text{old}}}$

# CHAIN RULE IN BACKPROPAGATION

How to calculate  $w_{11}^2_{\text{new}}$ ?

 $w_{11}^2_{\text{new}} = w_{11}^2_{\text{old}} - \eta \cdot \frac{\delta L}{\delta w_{11}^2_{\text{old}}}$ , Slope
 $\frac{\delta L}{\delta w_{11}^3_{\text{old}}} = \frac{\delta L}{\delta O_{31}} \cdot \frac{\delta O_{31}}{\delta w_{11}^3_{\text{old}}}$ 

$\therefore$  we get  $w_{11}^3_{\text{new}}$  (updated weight)

||| eg, we can calculate for  $w_{21}^3_{\text{new}}$ .

$w_{11}^1_{\text{new}} = w_{11}^1_{\text{old}} - \eta \cdot \frac{\delta L}{\delta w_{11}^1_{\text{old}}}$ 
 $\frac{\delta L}{\delta w_{11}^1_{\text{old}}} = \left[ \frac{\delta L}{\delta O_{31}} \cdot \frac{\delta O_{31}}{\delta O_{21}} \cdot \frac{\delta O_{21}}{\delta w_{11}^1_{\text{old}}} \right] + \left[ \frac{\delta L}{\delta O_{31}} \cdot \frac{\delta O_{31}}{\delta O_{22}} \cdot \frac{\delta O_{22}}{\delta O_{12}} \cdot \frac{\delta O_{12}}{\delta w_{11}^1_{\text{old}}} \right]$