



5 DAYS !  
**DEEP LEARNING**

**LIVE COMMUNITY  
SESSION**

# Day 1 : Deep Learning

7pm - 9pm IST

Motive : { Clear their Basics, Maths, Interview Preparation }

## Agenda

- ① Deep Learning → Perceptron { AI VS ML VS DL VS DS }
- ② Forward Propagation
- ③ Backward Propagation
- ④ Loss function
- ⑤ Activation functions
- ⑥ Optimizers.

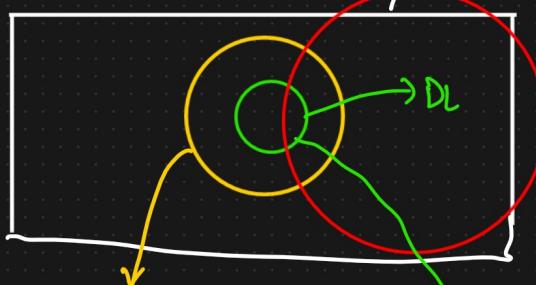
## Prerequisite

- ① Python
- ② ML
- ③ Stats

→ Deep learning

AI VS ML VS DL VS DS

{ ML is a subset  
of AI }



Goal

{ application which can do its own task  
Without any human Intervention }

① Netflix App

② Self Driving Cars

③ Amazon Application

④ Sofia

RTX  
3090

ML → Statstool to analyze  
the data, visualize the data,  
predictions, forecasting, clustering

Researchers (1958)

Multi Layered  
Neural Net  
{ Mimic the human brain }

{ Perceptron }

- ⑥ Why Deep Learning Is Becoming Popular?

① 2005 → ORKUT, FACEBOOK, Instagram, WhatsApp, LinkedIn, Twitter

DATA → Exponentially ↑↑↑

2008 → of BIG DATA} → Efficiently

Years {2013-2014}

Netflix

2013 → Company had huge amount of Data

↓ {AI → popular}

Scans → Products

Panasonic : AC's, TV's, Refrigerator {Data}



Model → Reduce the Electricity Bills ↓

↓  
Subscription Basis

{RTX Titan}



{Current Revenue, Better Decisions} {RTX 3090} ↑

② Hardware Advancement (Nvidia) → GPU's → TRAINING THE MODEL.

GPU's → Cost ↓↓↓

④ Perception {Single Layered Neural N/W}.

I/P layers

Study  
 $x_1$

play  
 $x_2$

Sleep  
 $x_3$

HL 1

O/p Layer

Output  
Neuron

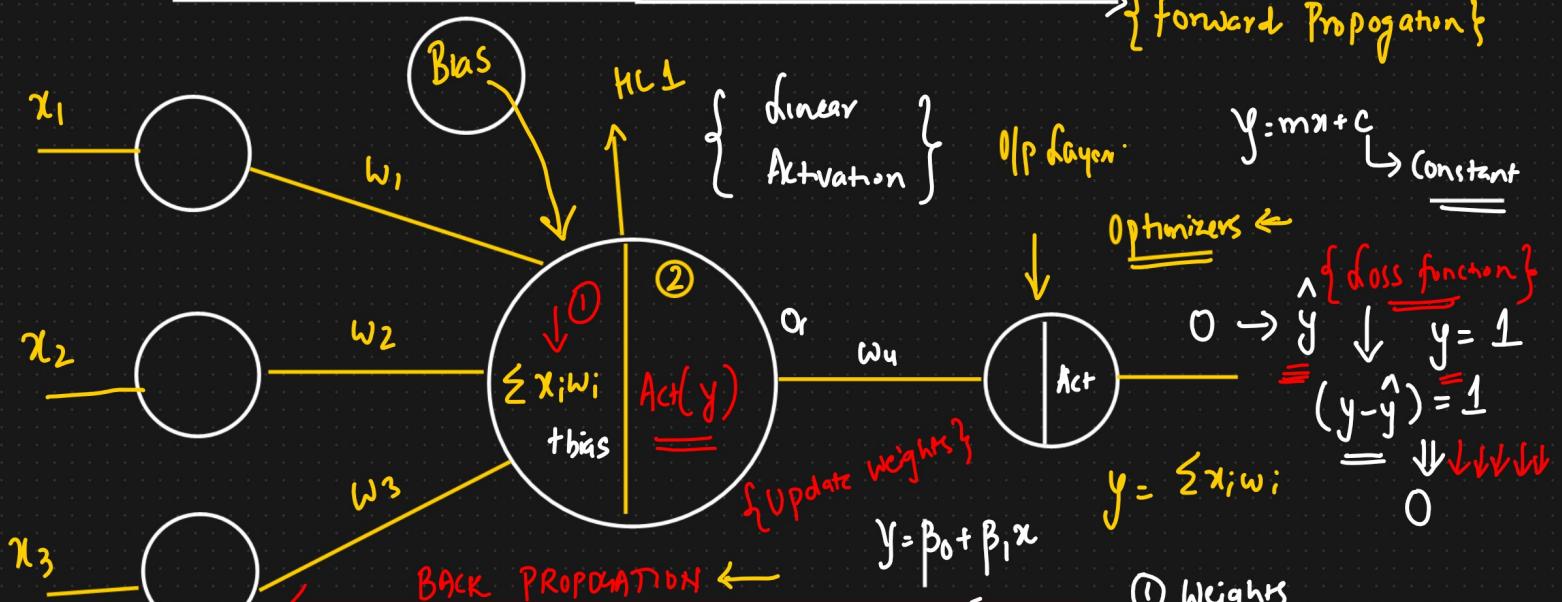
Binary Classification

Dataset



	Study	Play	Sleep	Pass/Fail
Study	7	3	7	1 ←
Play	2	5	8	0
Sleep	4	3	7	1

Forward Propagation



$$\text{Binary Classification} = \hat{y}$$

$$\text{Sigmoid} = \frac{1}{1 + e^{-\hat{y}}} = \frac{1}{1 + e^{-(\sum x_i w_i + b)}}$$

0 to 1

$$\begin{cases} \hat{y} \geq 0.5 \rightarrow 1 \\ \hat{y} < 0.5 \rightarrow 0 \end{cases}$$

$$\begin{cases} 0 \text{ or } 1 \end{cases}$$

① ANN

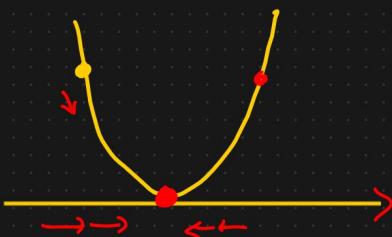
② CNN

③ RNN

④ Object Detection

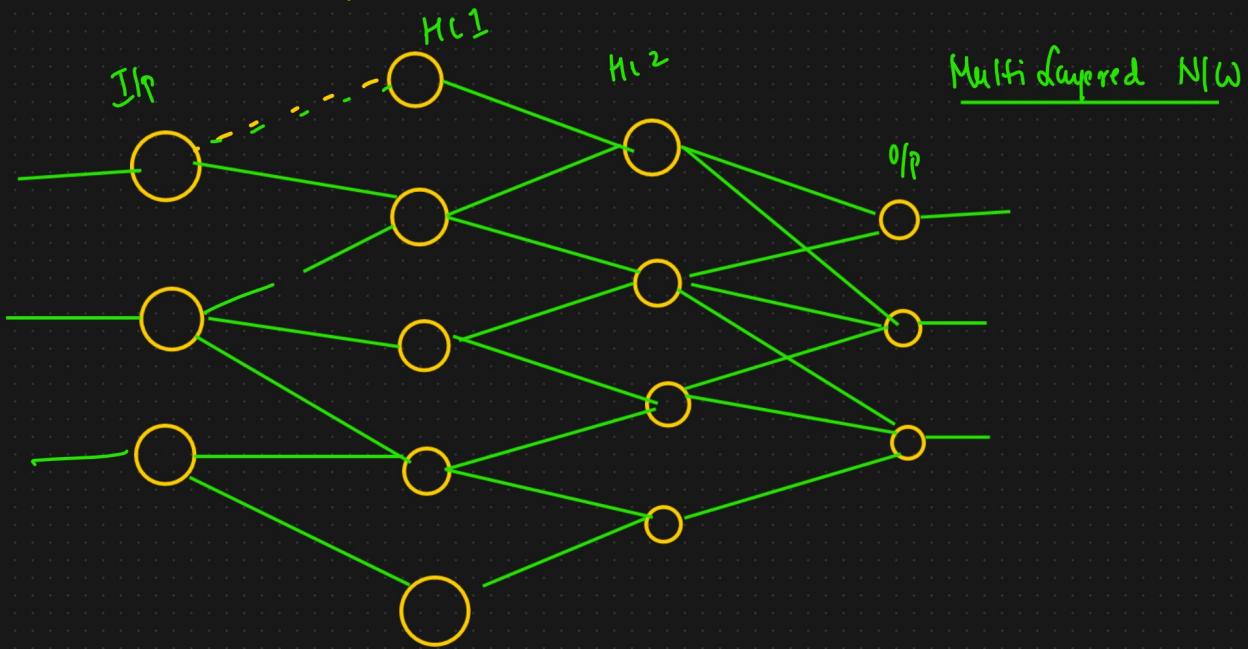
GRADIENT DESCENT  $\Rightarrow$  Kind of Optimizers

$$\hat{y} - y \Rightarrow \text{minimal}$$



## Conclusion

- ① I/P layers
  - ② Weights
  - ③ Bias Day 2
  - ④ Activation function  $\times$
  - ⑤ Loss function  $\times \{ \hat{y} - y \}$
  - ⑥ Optimizers
  - ⑦ Update the weight
- Forward Propagation      (300s)
- Scribble Ink
- Day 2      Day 3      Day 3
- Backward Propagation
- ANN
- Chain Rule of Differentiation



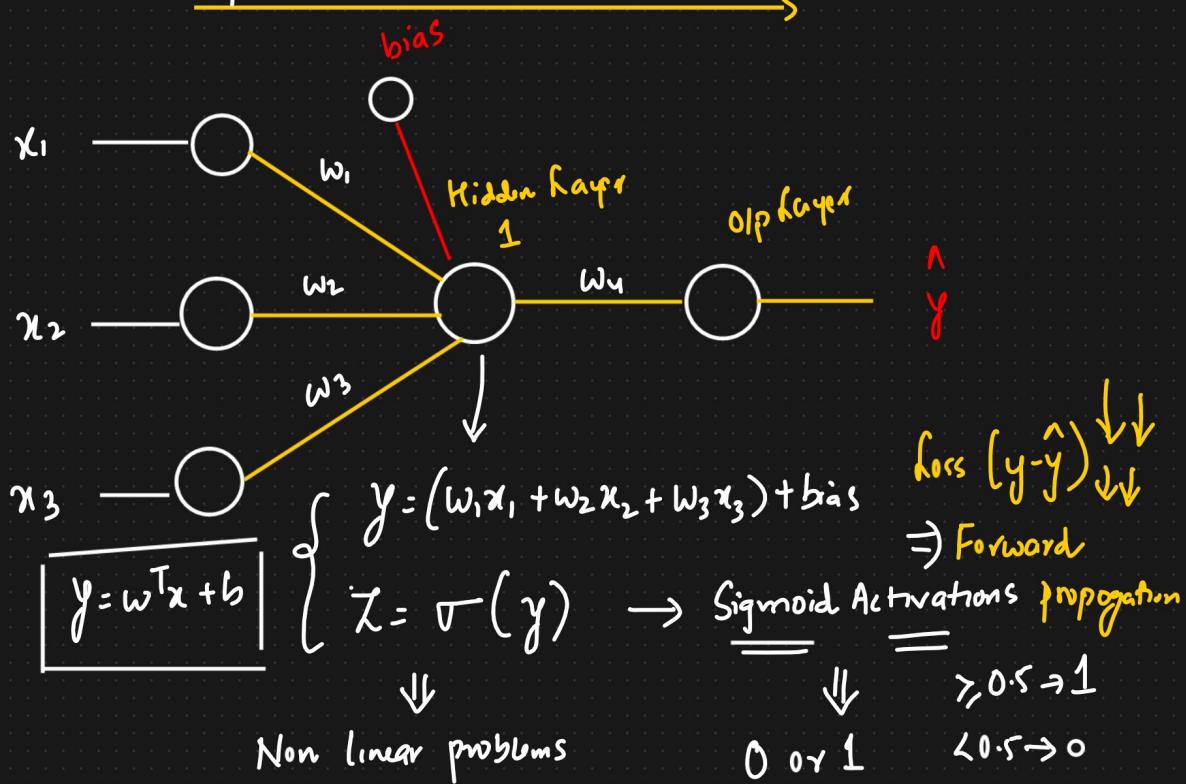
- ① ANN
  - ② CNN
  - ③ RNN
- NLP ~ 7 days

# Day 2 - Deep Learning

## Agenda

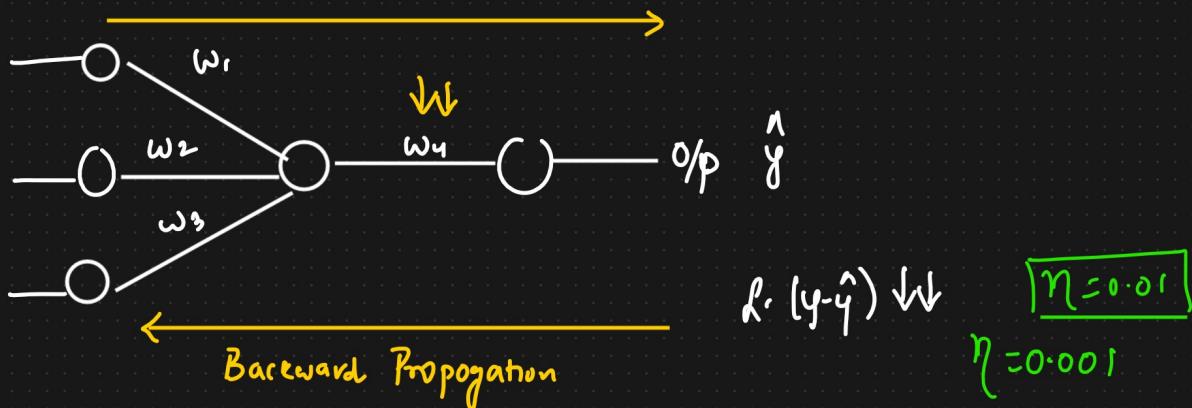
- ① Forward Propagation ✓
- ② Chain Rule of Differentiation ✓
- ③ Vanishing Gradient Problem ✓
- ④ Loss functions ✓

## Activation functions



## ② Backpropagation

- ① Weight update formula
- ② Chain Rule of Differentiation ✓



## ① Weight updation formula

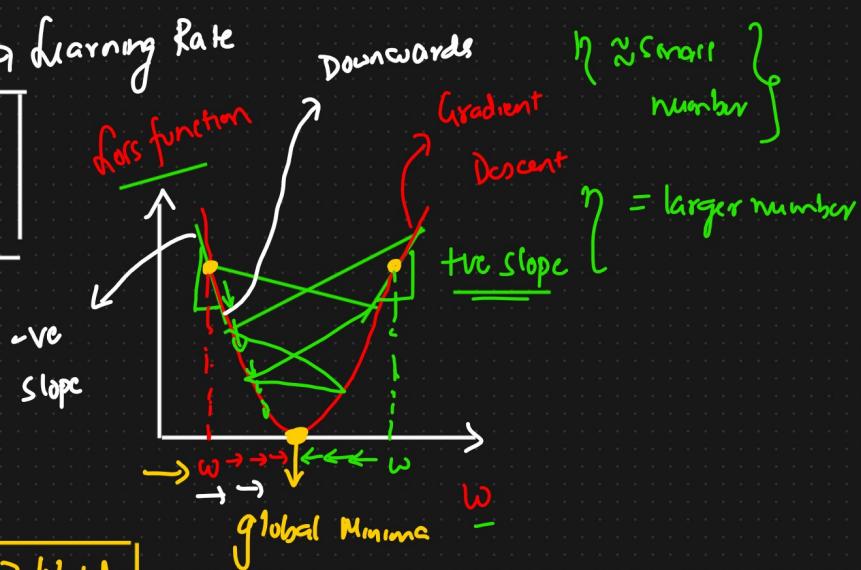
$$W_{\text{new}} = W_{\text{old}} - \eta \left[ \frac{\partial h}{\partial W_{\text{old}}} \right]$$

↑ Slope

$$\frac{\partial h}{\partial W_{\text{old}}} = \boxed{-\text{ve Slope}}$$

$$W_{\text{new}} = W_{\text{old}} - \eta (-\text{ve}) \quad \boxed{W_{\text{new}} > W_{\text{old}}}$$

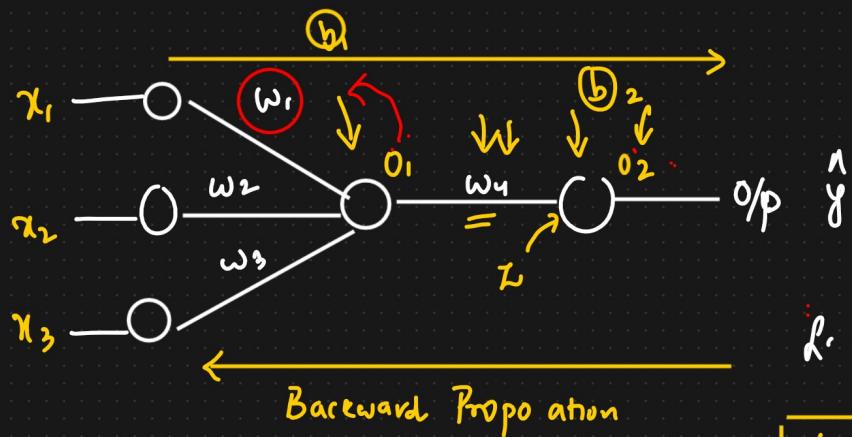
$$= W_{\text{old}} + \eta (\text{ve})$$



$$W_{\text{new}} = W_{\text{old}} - \eta (\text{+ve})$$

$$\boxed{W_{\text{new}} < W_{\text{old}}}$$

## ② Chain Rule of Differentiation



$$W_{\text{new}} = W_{\text{old}} - \eta \left[ \frac{\partial L}{\partial W_{\text{old}}} \right]$$

$$L = \sigma(o_1 w_4 + b)$$

$$L = (y - \hat{y})$$

$$W_4_{\text{new}} = W_4_{\text{old}} - \eta \left[ \frac{\partial L}{\partial W_4_{\text{old}}} \right]$$

$$b_2_{\text{new}} = b_2_{\text{old}} - \eta \left[ \frac{\partial L}{\partial b_2_{\text{old}}} \right]$$

{ Chain Rule of  
Derivative }

$$\frac{\partial h}{\partial W_{\text{old}}} = \frac{\partial L}{\partial o_2} * \frac{\partial o_2}{\partial W_{\text{old}}}$$

$$\omega_{1, \text{new}} = \omega_{1, \text{old}} - \eta$$

$$\frac{\partial h}{\partial \omega_{1, \text{old}}}$$

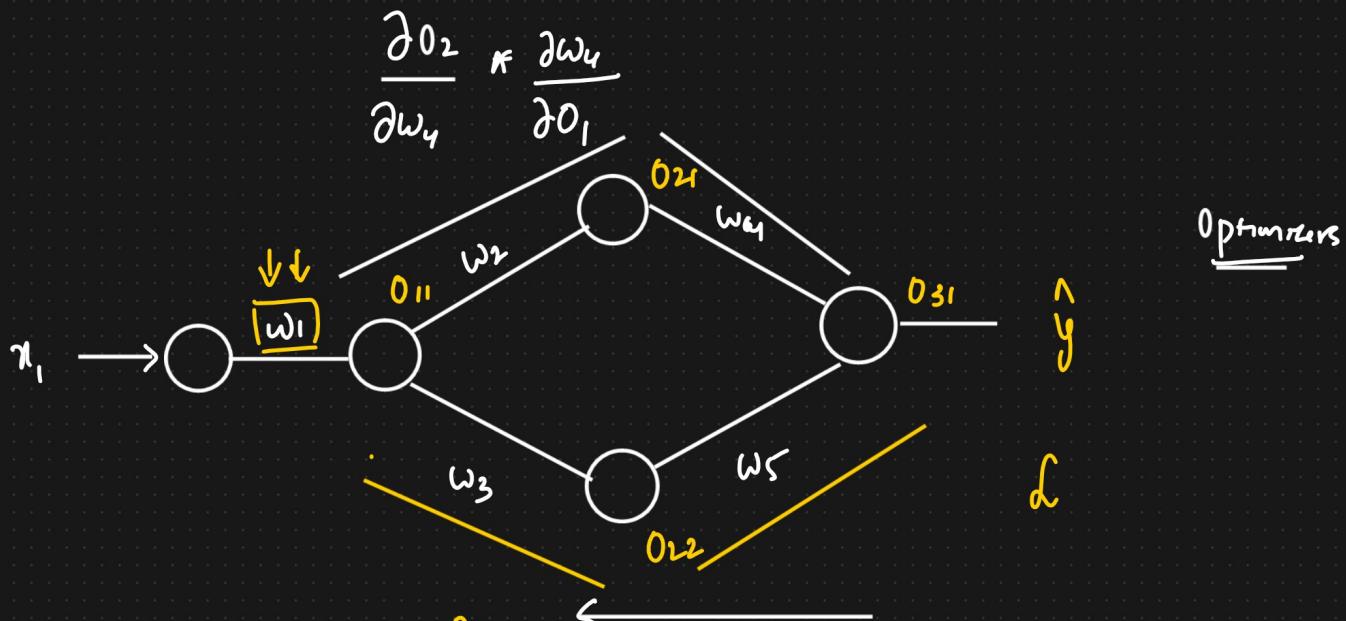
$$\omega_{2, \text{new}} = \omega_{2, \text{old}} - \eta$$

$$\frac{\partial h}{\partial \omega_{2, \text{old}}}$$

loss

$$\frac{\partial h}{\partial \omega_{1, \text{old}}} = \frac{\partial h}{\partial o_2} * \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial \omega_{1, \text{old}}}$$

↓



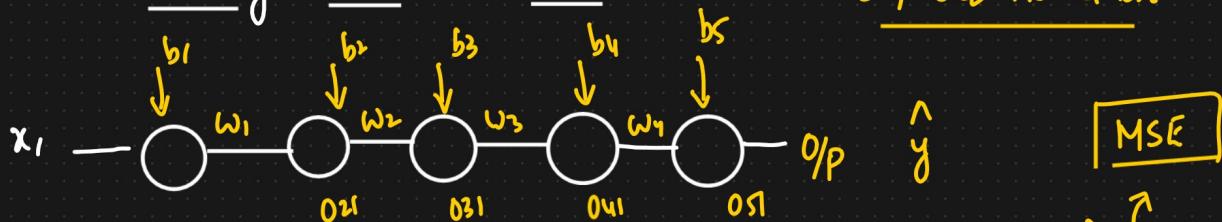
$$\omega_{1, \text{new}} = \omega_{1, \text{old}} - \eta \frac{\partial h}{\partial \omega_{1, \text{old}}}$$

↗ Chain Rule of Derivatives

$$\frac{\partial h}{\partial \omega_{1, \text{old}}} = \left[ \frac{\partial h}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial o_{11}} * \frac{\partial o_{11}}{\partial \omega_{1, \text{old}}} \right]$$

$$+ \left[ \frac{\partial h}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{22}} * \frac{\partial o_{22}}{\partial o_{11}} * \frac{\partial o_{11}}{\partial \omega_{1, \text{old}}} \right]$$

### ③ Vanishing Gradient Problem



$$w_{i,\text{new}} = w_{i,\text{old}} - \eta \left[ \frac{\partial L}{\partial w_{i,\text{new}}} \right]$$

$$o_{51} = \sigma \left[ (o_{41} * w_4) + b \right]$$

↓  
Sigmoid Activation

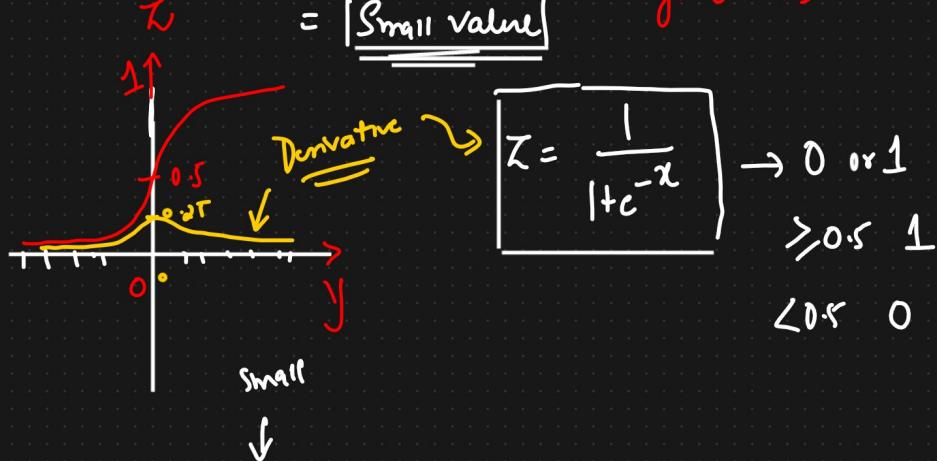
$$\frac{\partial L}{\partial w_{i,\text{new}}} = \frac{\partial L}{\partial o_{51}} * \frac{\partial o_{51}}{\partial o_{41}} * \frac{\partial o_{41}}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial w_i}$$

$$= 0.25 * 0.15 * 0.10 * 0.05 * 0.02$$

$\pi$  = Small value       $y = w^T x + b$

Derivative:

$$0 \leq \sigma(y) \leq 0.25$$



$$w_{i,\text{new}} = w_{i,\text{old}} - \eta \left( \text{Small number} \right)$$

Another Activation  
function.

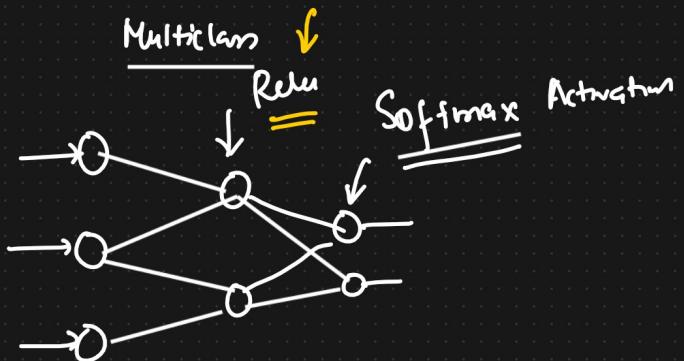
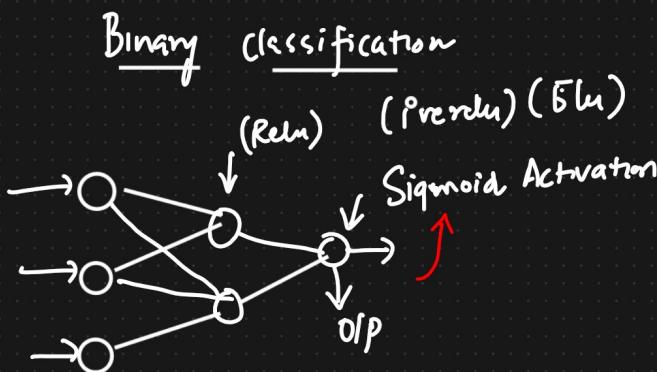
$w_{i,\text{new}} \approx w_{i,\text{old}}$

$\Rightarrow$  Vanishing gradient Problem

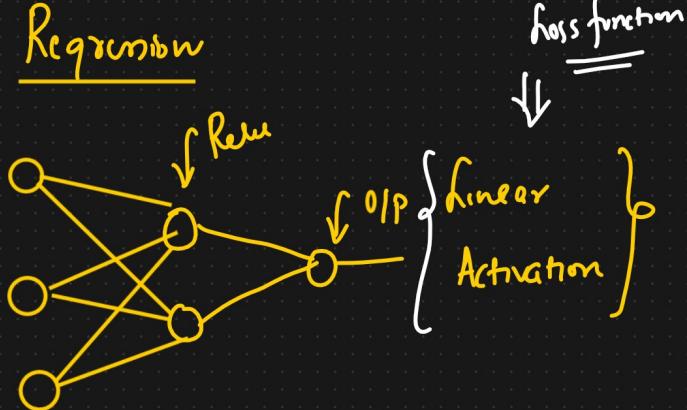
No change in weights

- ① Sigmoid  
 ② Tanh  
 ③ ReLU  
 ④ Leaky ReLU  
 ⑤ PReLU

# Technique Which Activation fn we should Use



## Regression



④

## loss functions

### Deep Learning (ANN)



#### Regression

Exp	Degree	Salary	O/P
10	Phd	—	
—	—	—	
—	—	—	

Regression Problem

#### Classification

Dataset	Play	Study	Pas/Fail
10	2		Fail
4	3		Fail
5	5		Maybe
2	7		Pas

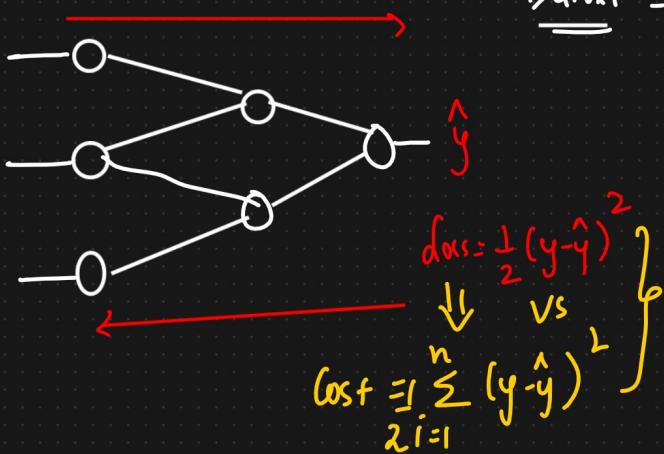
# ① Regression

- ① MSE {Mean Squared Error}
- ② MAE {Mean Absolute Error}
- ③ Huber Loss

Loss function AND

Cost Function

10 records  
↓  
Dataset = 100 records



## ① Mean Squared Error (MSE)

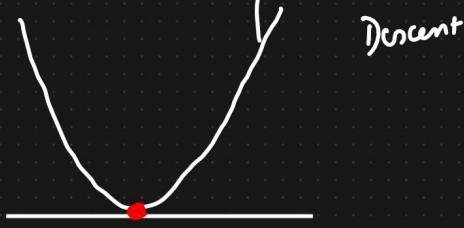
Loss function =  $\frac{1}{2} (y - \hat{y})^2$  Cost function =  $\frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$

↓ Error

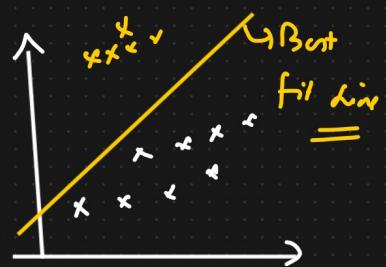
Quadratic Equation

$$(a-b)^2 = a^2 - 2ab + b^2$$

$$a x^2 + b x + c$$



Penalizing Error

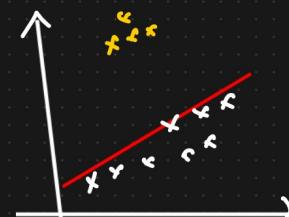


### Advantages

- ① Differentiable
- ② It has only 1 local or Global Minima.
- ③ It converges faster

### Disadvantage

- ① Not Robust to outliers



## ② Mean Absolute Error

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

$$\text{Cost fn} = \frac{1}{2} \sum_{i=1}^n |y_i - \hat{y}_i|$$

① Robust to outliers  $\Rightarrow$  Time consuming ↑↑

$\Rightarrow$  Subgradient



③ Huber loss

① MSE

② MAE

outliers are not  
present

Hyperparameter

$$\text{loss} = \begin{cases} \frac{1}{2} (y - \hat{y})^2 & \text{if } |y - \hat{y}| \leq \delta \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2, & \text{otherwise} \end{cases}$$



③ Classification

CROSS ENTROPY  $\rightarrow$  Binary Cross Entropy  $\rightarrow$  Binary Classification

CROSS ENTROPY  $\rightarrow$  Categorical Cross Entropy  $\rightarrow$  Multiclass Classification.

① Binary Cross Entropy

$$\text{loss} = -y * \log(\hat{y}) - (1-y) * \log(1-\hat{y}) \Rightarrow \text{logistic regression}$$



$$\text{loss} = \begin{cases} -\log(1-\hat{y}) & \text{if } y=0 \\ -\log(\hat{y}) & \text{if } y=1 \end{cases} \quad \text{Binary Classification}$$

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

## ② Categorical Cross Entropy {Multi-class Classification Problem}

$f_1$	$f_2$	$f_3$	O/P	Good	Bad	Neutral	$i = \text{Row}$
1	3	4	Good	[ 1 ]	0	0 ]	$j = \text{Column}$ }
2	6	7	Bad	0	1	0	
5	9	10	Neutral	0	0	1	
8							

$$d(x_i, y_i) = - \sum_{j=1}^c \boxed{y_{ij}} \neq \ln(\hat{y_{ij}})$$

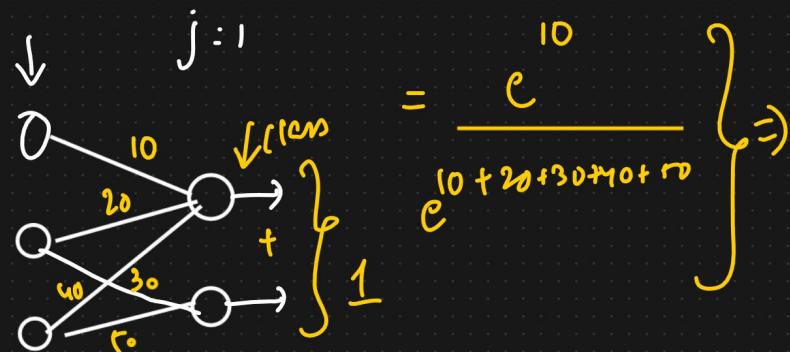
$$y_i = [y_{i1}, y_{i2}, y_{i3}, \dots, y_{ic}]$$

$$y_{ij} = \begin{cases} 1 & \text{if the element is in class.} \\ 0 & \text{Otherwise} \end{cases}$$

$$\hat{y}_{ij} = \text{Softmax Activation} \xrightarrow{\text{O/p Layer}} =$$

$$f(z) = \frac{e^{z_i}}{\sum_j e^{z_j}}$$

} Activation } Softmax }



## Conclusions

ReLU, Softmax  $\Rightarrow$  MultiClass  $\left\{ \begin{array}{l} \text{Categorical Cross} \\ \text{Entropy} \end{array} \right.$

ReLU, Sigmoid  $\Rightarrow$  Binary  $\rightarrow$  Binary Cross Entropy.

## Linear Regression

ReLU, Linear Activation  $\rightarrow$  MSE, MAE, Huber Loss

## Day 3 - Deep Learning.

## Agenda

## ① Optimizers

- i) Gradient Descent
  - 2) SGD (Stochastic Gradient Descent)
  - 3) Mini Batch SGD
  - 4) SGD with Momentum
  - 5) Adagrad
  - 6) RMSProp
  - 7) Adam Optimizer

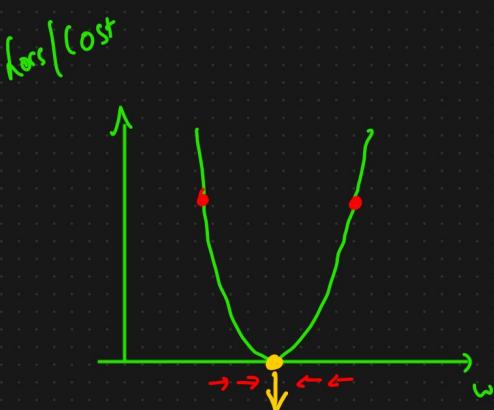
## Batch, Epochs, Iterations

ANN

## ① GRADIENT DESCENT $\xrightarrow{\text{Optimieren}}$

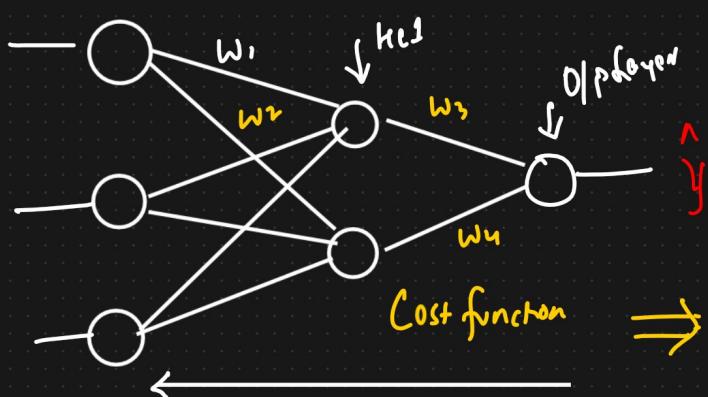
## Weight Updation Formula

$$W_{nw} = W_{old} - \eta \frac{\partial h}{\partial W_{old}}$$



## Global Minima.

↓ IPhone



$$\hat{y} \xrightarrow{\text{MSE}} \text{Optimizers}$$

$$\downarrow \quad \downarrow$$

$$w$$

$$\Rightarrow \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \{$$

A diagram illustrating an epoch. At the top, the word "Epoch" is written in yellow, with a horizontal yellow arrow pointing to the right above it. Below this, a vertical yellow bracket spans the length of the arrow. At the bottom, a yellow arrow points back to the left, ending at the bottom of the bracket. Above the bracket, the number "1000000" is written in black, indicating the number of steps or iterations within an epoch.



## ④ SGD With Momentum

{ Exponential Weighted Average }

$$w_{\text{new}} = w_{\text{old}} - \eta \left[ \frac{\partial L}{\partial w_{\text{old}}} \right]$$

$$b_{\text{new}} = b_{\text{old}} - \eta \left[ \frac{\partial L}{\partial b_{\text{old}}} \right]$$

↓  
Time Series

↓  
ARIMA, ARMA,

$$w_t = w_{t-1} - \eta \left[ \frac{\partial L}{\partial w_t} \right]$$

Exponential Weighted Average

{ Forecasting }

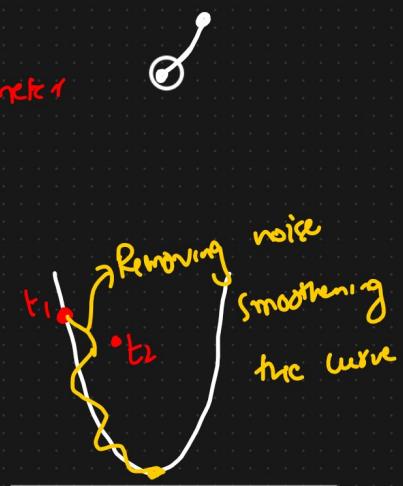
$$t_1 \quad t_2 \quad t_3 \quad t_4 \quad \dots \quad t_n$$

$$a_1 \quad a_2 \quad a_3 \quad a_4 \quad \dots \quad a_n$$

$\beta \Rightarrow$  Hyper parameter

$$\beta = 0 \text{ to } 1$$

$$0.95$$



$$V_{t_1} = a_1$$

$$V_{t_2} = \beta \times V_{t_1} + (1 - \beta) \times a_2$$

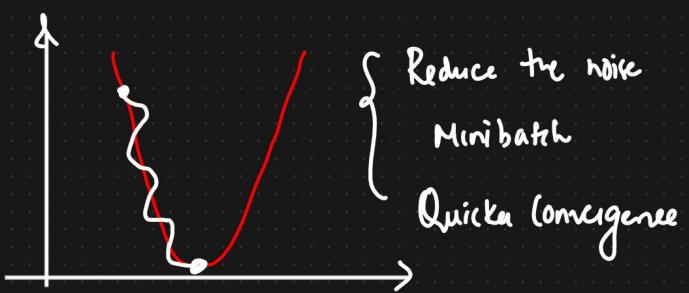
$$= (0.95) \times V_{t_1} + (0.05) \times a_2$$

Exponential Weighted Avg

↑

$$w_t = w_{t-1} - \eta V_{dw}$$

$$V_{dw} = \beta \times V_{dw_{t-1}} + (1 - \beta) \times \frac{\partial L}{\partial w_{t-1}}$$

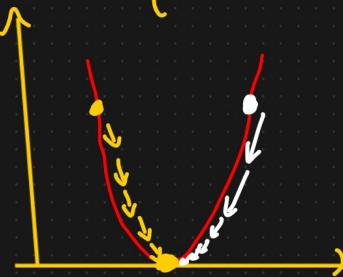


## Recap

- ① Gradient Descent
- ② SGD
- ③ Mini batch SGD
- ④ SGD with Momentum
- ⑤ Adagrad  $\rightarrow$  Adaptive Gradient Descent

$\{\text{fixed}\} \Rightarrow \text{optimizer}$

$\uparrow$   
 $\eta = \text{Learning Rate}$



Global Minima

$\downarrow$   
Minima

$\downarrow$

$$w_t = w_{t-1} - \eta \frac{dl}{\partial w_{t-1}}$$

$$w_t = w_{t-1} - \eta \left[ \frac{\partial l}{\partial w_{t-1}} \right]$$

$$w_t \approx w_{t-1}$$

$$\eta = \eta \leq 0.01$$

$$\Rightarrow \sqrt{d_t + \epsilon}$$

$$d_t = \sum_{i=1}^t \left( \frac{\partial l}{\partial w_t} \right)^2$$

Huge  
number

$t=1 \quad t=2 \quad t=3$

$$\eta = 0.01 \quad \eta = 0.005 \quad \eta = 0.002$$

- ⑥ Adadelta and Rmsprop

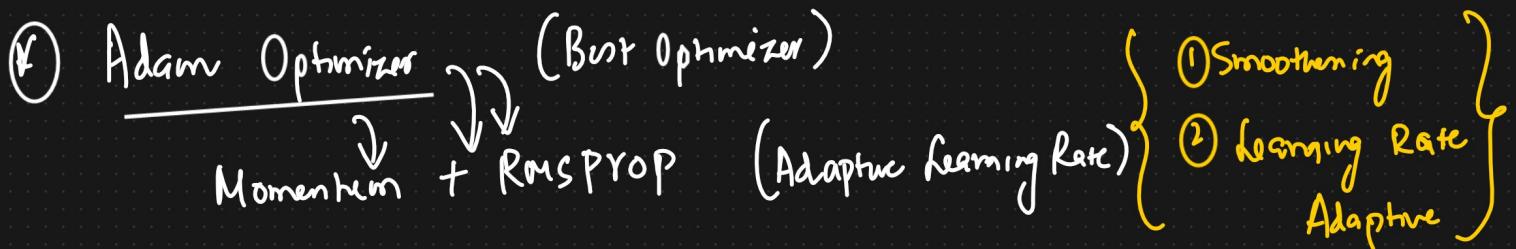
$\downarrow$   
Exponential Weighted Average

$$\eta^t = \frac{\eta}{\sqrt{S_{dw} + \epsilon}}$$

$$S_{dw}^t = \beta S_{dw}^t + (1-\beta) \left( \frac{\partial L}{\partial w_{t-1}} \right)^2$$

$$\beta = 0.95$$

$$S_{dw}^t = (0.95) S_{dw}^t + (0.05) \left( \frac{\partial L}{\partial w_{t-1}} \right)^2$$



$$V_{dw,0} = 0 \quad S_{dw} = 0 \quad S_{db} = 0$$

↳

$$w_t = w_{t-1} - \eta^t V_{dw}$$

$$b_t = b_{t-1} - \eta^t V_{db}$$

$$\eta^t = \frac{\eta}{\sqrt{S_{dw} + \epsilon}}$$

$$V_{dw}^t = \beta \times V_{dw}^t + (1-\beta) \frac{\partial L}{\partial w_{t-1}}$$

$$V_{db}^t = \beta \times V_{db}^t + (1-\beta) \frac{\partial L}{\partial b_{t-1}}$$