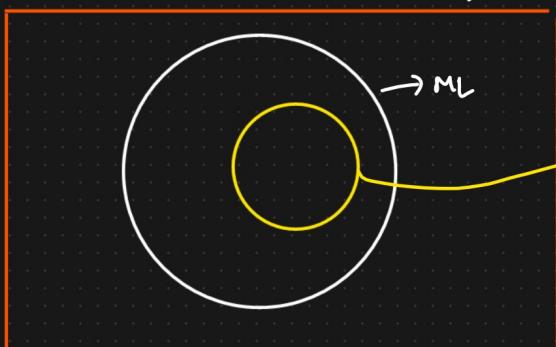


# Deep learning

AI  $\Rightarrow$  Artificial Intelligence



Deep Learning  $\Rightarrow$  Neural Networks  
Multi-layered  
↓  
Mimic the Human Brain

## Deep learning

- ① ANN  $\rightarrow$  Artificial Neural N/w  $\begin{cases} \rightarrow \text{Classification} \\ \rightarrow \text{Regression} \end{cases}$
- ② CNN  $\rightarrow$  Convolutional Neural N/w  $\rightarrow$  I/P : Image, Video  $\rightarrow$  RCNN, MASKED RCNN, frames  
Computer Vision  
Object Detection  
↑
- ③ RNN  $\rightarrow$  Recurrent Neural N/w  $\rightarrow$  NLP  $\rightarrow$  NLP, Time Series  
I/P: Text, Time Series  
Detection, Yolo V3, V6, V7

## FRAMEWORK

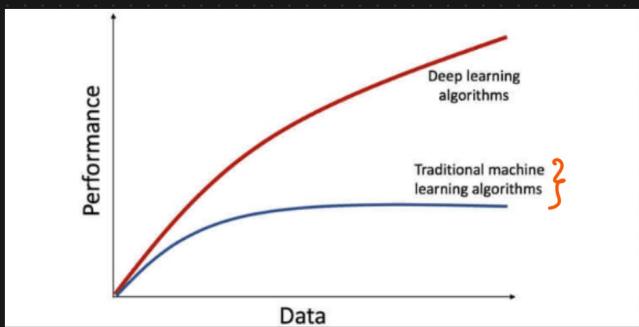
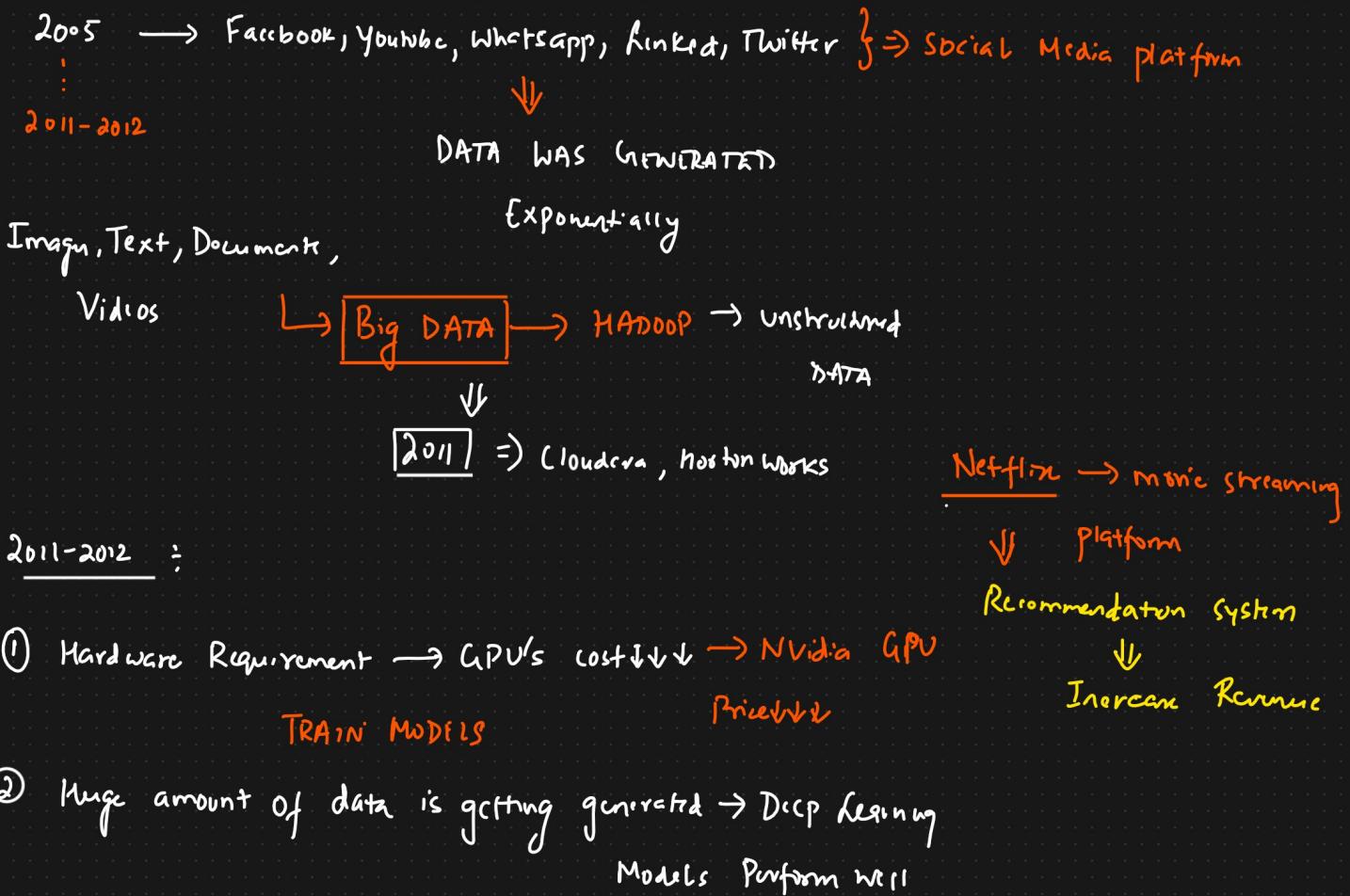
TENSORFLOW

End-to-End Project

{ Word Embedding, LSTM RNN, GRU RNN,

Bidirectional LSTM RNN, Encoder Decoder,  
Transformers, BERT }

## ② Why Deep Learning Is Becoming popular?



③ Deep learning is been used in Many Domains

- ① Medical
  - ② E-commerce
  - ③ Retail
  - ④ Marketing . . . .
- ④ Frameworks Open source

Community size ↑↑

Tensorflow



Google

Pytorch



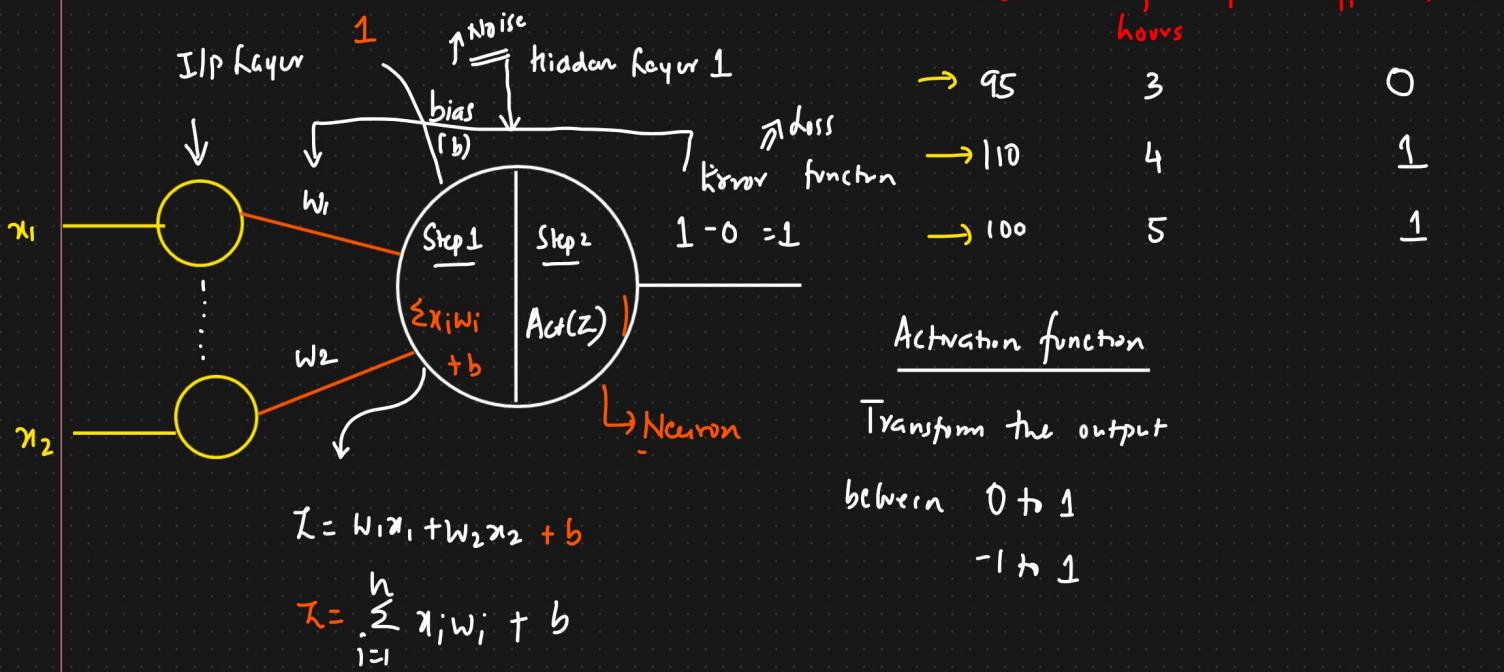
Facebook



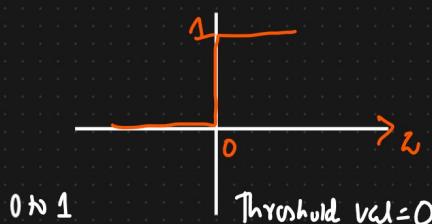
More Research

### ③ Perceptron [Artificial Neuron or Neural Network Unit]

- ① Input layer ✓
  - ② Hidden layer ✓
  - ③ Weights ✓
  - ④ Activation function ✓
- [Single Layered NN]  $\downarrow$
- Binary classifier  $\downarrow$
- DATA SET

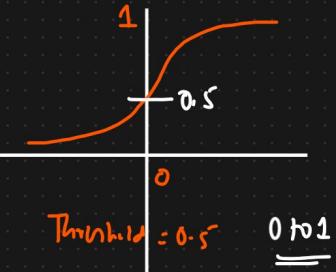


Step Function

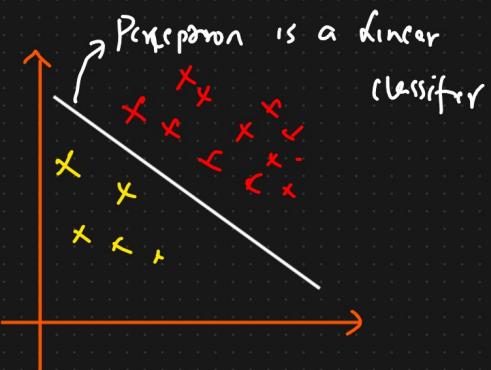


$$= \begin{cases} 0 & z \leq 0 \\ 1 & z > 0 \end{cases}$$

Sigmoid Function



$$= \begin{cases} 1 & z > 0.5 \\ 0 & z \leq 0.5 \end{cases}$$



Step 1

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

$$z = b + w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_n$$

$$y = mx + c$$

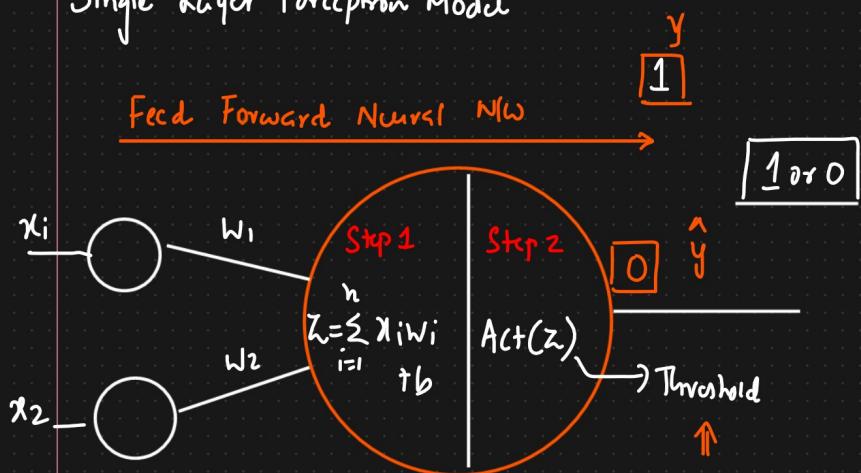
$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n$$

↓  
linear  
problem  
statement

## Perceptron Models

Single Layer Perceptron Model

Feed Forward Neural NW



Multi Layered Perceptron Model

① Forward Propagation

② Backward " "

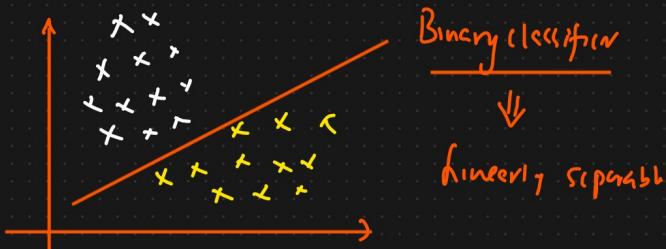
③ Loss functions

④ Activation functions

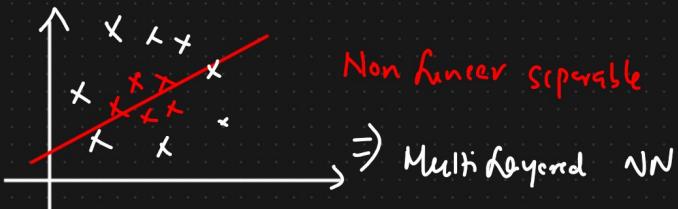
⑤ Optimizers

[ANN]  
↓

Linear Separable problem



Non linearly separable



# 5 Multi layered Perceptron Model [Artificial Neural N/W]



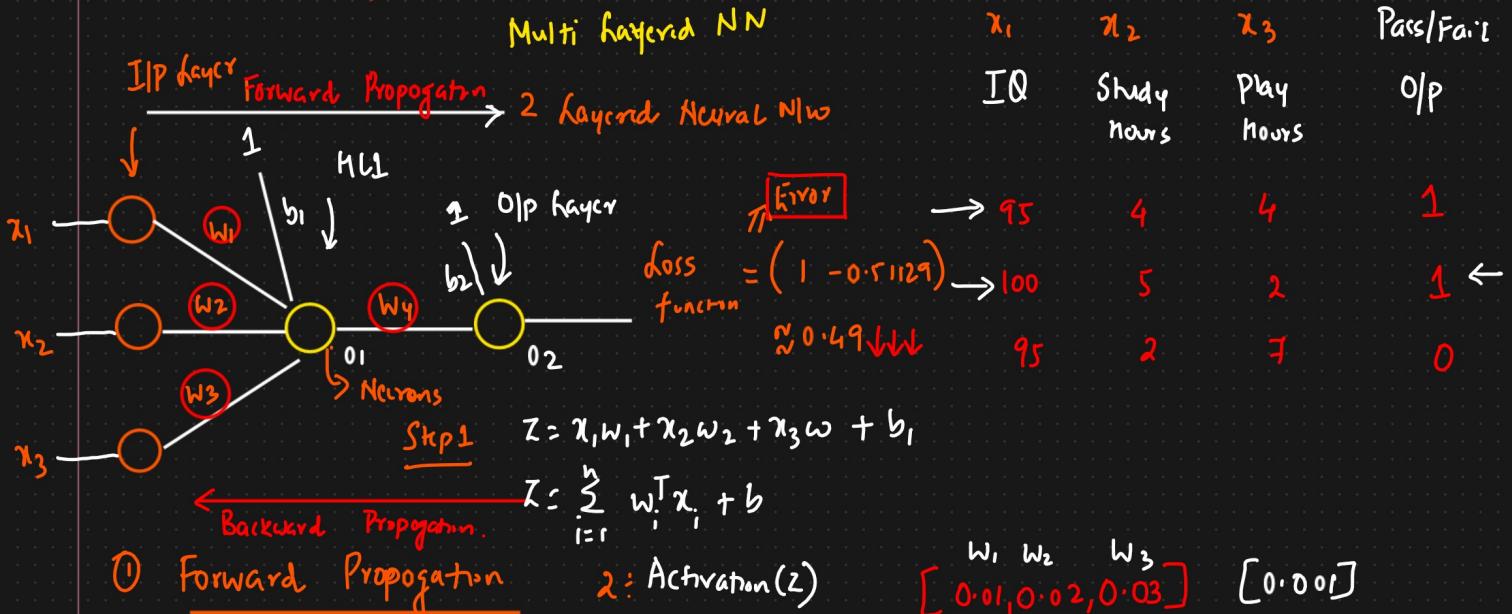
① Forward Propagation

② Backward Propagation → Geoffrey Hinton →

③ Loss function

④ Optimizers ✓

⑤ Activation function.



## Hidden Layer 1

$$\begin{aligned} \text{Step 1: } Z &= 95 \times 0.01 + 4 \times 0.02 + 4 \times 0.03 + 1 \times 0.01 \\ &= 1.151 \end{aligned}$$

$$\text{Sigmoid.} = \frac{1}{1+e^{-z}}$$

$$\text{Step 2: Activation (Z)}$$

$$f(z) = \frac{1}{1+e^{-1.151}} = 0.759$$

$$O_1 = 0.759$$

## Hidden Layer 2

$$w_4 = 0.02$$

$$b_2 = 0.03$$

$$\text{Step 1: } Z = O_1 \times w_4 + b_2$$

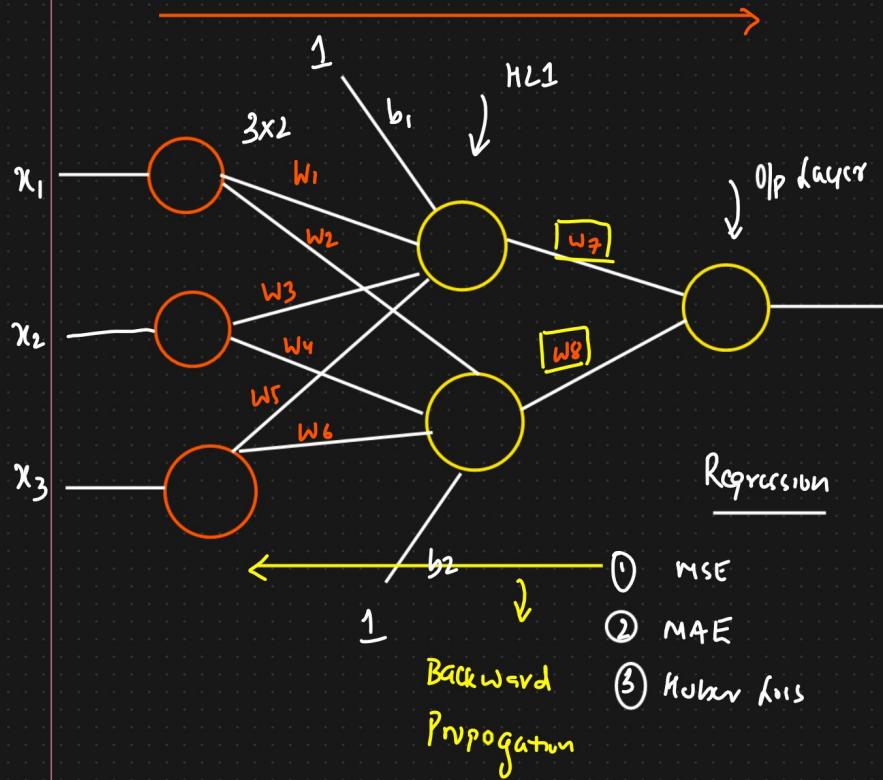
$$= 0.759 \times 0.02 + 0.03$$

$$= 0.04518$$

Step 2 : Activation(z)  $\frac{1}{1+e^{-(0.04518)}} = 0.51129$

$$O_2 = 0.51129 \Rightarrow \hat{y}$$

## ⑥ Back Propagation And Weight Updation Formula



Loss function  
 $(y - \hat{y})^2$

### Regression

- ① MSE
- ② MAE
- ③ Huber Loss

### Classification

- ① Binary Cross Entropy
- ② Categorical Cross Entropy

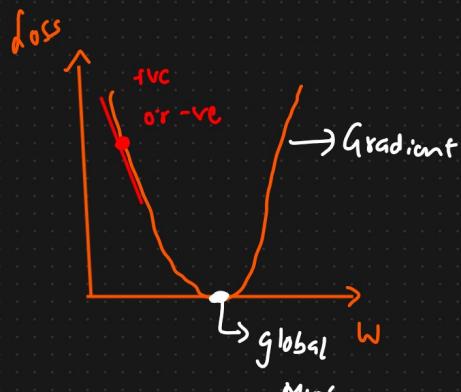
### Weight update Formula

$$w_{7\text{new}} = w_{7\text{old}} - \eta \left[ \frac{\partial h}{\partial w_{7\text{old}}} \right]$$

slope

$$w_{8\text{new}} = w_{8\text{old}} - \eta \left[ \frac{\partial h}{\partial w_{8\text{old}}} \right]$$

Rate



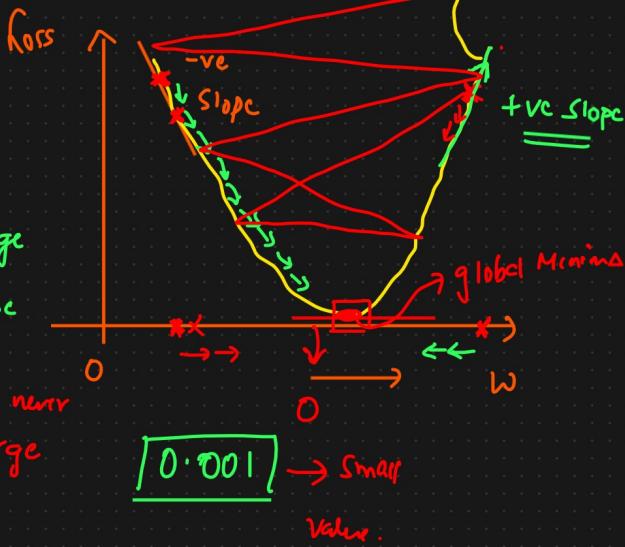
$$w_{\text{new}} = w_{\text{old}} - \eta \left[ \frac{\partial h}{\partial w_{\text{old}}} \right] \Rightarrow \text{Weight Updation Formula.}$$

Gradient Descent

↗ Optimizers

I/P featur.s			Pass/Fail
$x_1$	$x_2$	$x_3$	
IQ	Study hours	Play hours	O/P
> 95	4	4	1
→ 100	5	2	1 ←
95	2	7	0

Optimizers : To reduce the loss value



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

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

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

$\eta$  = large value  
↓  
It may never converge

Learning Rate

$$\boxed{0.001} \rightarrow \text{Small value.}$$

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

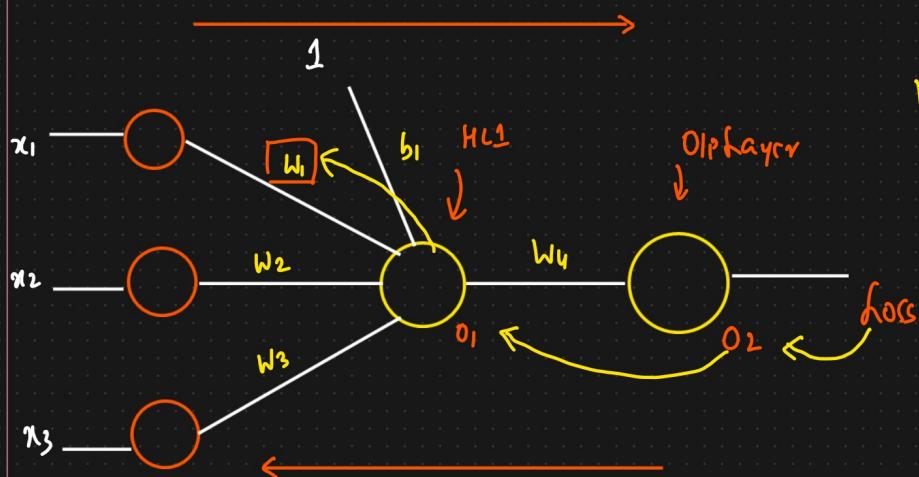
$$= W_{\text{old}} - \eta \quad (+\text{ve})$$

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

When  $w$  reaches Global Minima

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

### ⑦ Chain Rule of Derivative



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

$$W_{4,\text{new}} = W_{4,\text{old}} - \eta \boxed{\frac{\partial h}{\partial W_{4,\text{old}}}}$$

$$\frac{\partial h}{\partial w_{1,old}} = \frac{\partial h}{\partial o_2} * \frac{\partial o_2}{\partial w_{1,old}} \Rightarrow \text{Chain Rule of Derivation}$$

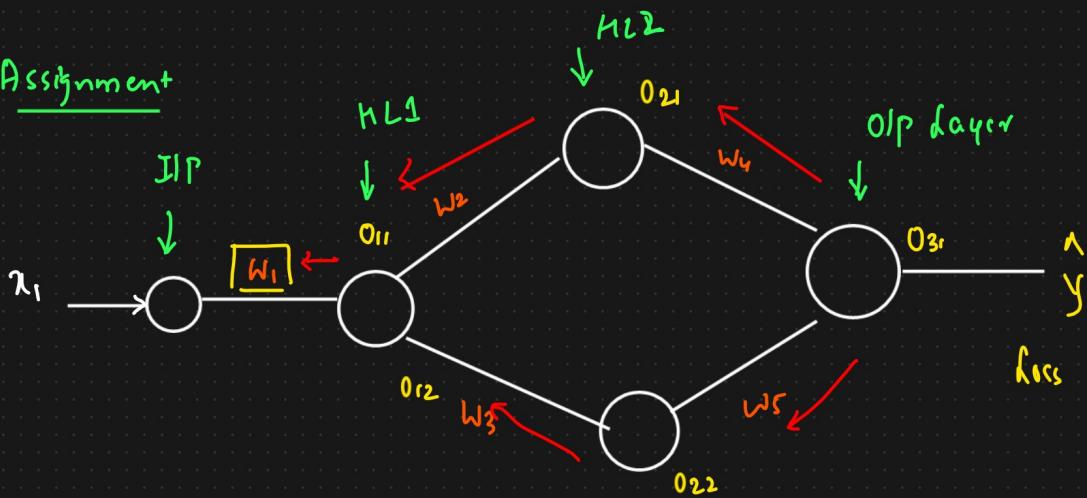
$$w_{1,new} = w_{1,old} - \eta \left[ \frac{\partial h}{\partial w_{1,old}} \right]$$

$$\boxed{\frac{\partial h}{\partial w_{1,old}} = \frac{\partial h}{\partial o_2} * \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_{1,old}}}$$

$w_{2,new}$

$w_{3,new}$

Assignment



$$w_{1,new} = w_{1,old} - \eta \left[ \frac{\partial h}{\partial w_{1,old}} \right]$$

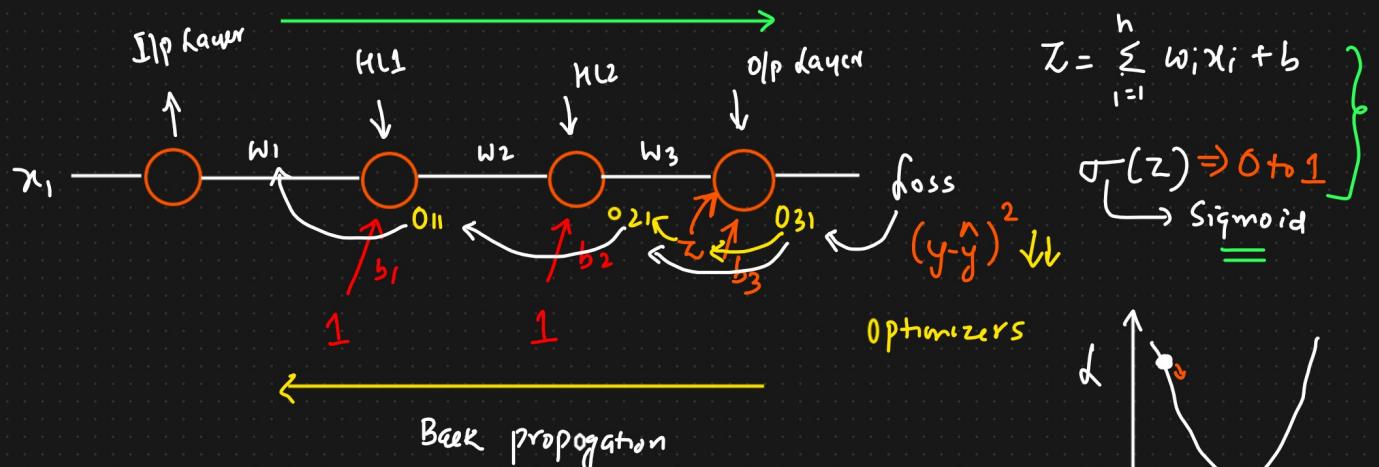


$$\frac{\partial h}{\partial w_{1,old}} = \left[ \frac{\partial h}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial o_{11}} * \frac{\partial o_{11}}{\partial w_{1,old}} \right]$$

+

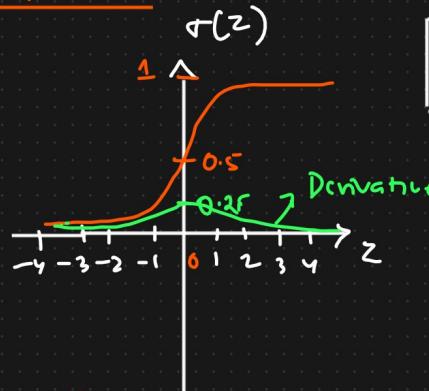
$$\left[ \frac{\partial h}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{22}} * \frac{\partial o_{22}}{\partial o_{12}} * \frac{\partial o_{12}}{\partial w_{1,old}} \right]$$

## ⑧ Vanishing Gradient Problem And Activation functions



### ⑨ Sigmoid Activation function

$$\sigma(z) = \frac{1}{1+e^{-z}}$$



$$0 \leq \sigma(z) \leq 1$$

↓

$$\text{Derivative } (\sigma(z)) \quad 0 \leq \sigma(z) \leq 0.25$$

Vanishing  
Gradient  
Problem.

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

$\Rightarrow \boxed{w_{\text{new}} \approx w_{\text{old}}}$

$$\frac{\partial h}{\partial w_{\text{old}}} = \frac{\partial h}{\partial w} * \boxed{\frac{\partial w}{\partial o_{31}}} * \boxed{\frac{\partial o_{31}}{\partial o_{21}}} * \boxed{\frac{\partial o_{21}}{\partial o_{11}}} * \boxed{\frac{\partial o_{11}}{\partial w_{\text{old}}}}$$

$\Rightarrow \text{Smaller value}$

$$o_{31} = \sigma \left( \underbrace{w_3 * o_{21}}_z + b_3 \right) \quad Z = w_3 * o_{21} + b_3$$

$$o_{31} = \sigma(z)$$

$$\frac{\partial o_{31}}{\partial o_{21}} = \frac{\partial (\sigma(z))}{\partial (z)} * \frac{\partial z}{\partial o_{21}} \quad \left\{ \text{Chain Rule} \right\}$$

$$0 \leq \sigma(z) \leq 0.25 \quad * \quad \frac{\partial((w_3 * o_{21}) + b_3)}{\partial(o_{21})}$$

↓

Derivative of  
Sigmoid

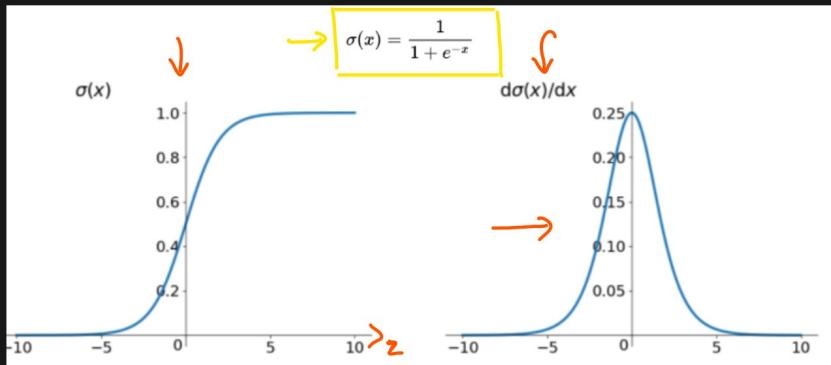
$$\left| \frac{\partial o_{31}}{\partial o_{21}} = 0 \leq \sigma(z) \leq 0.25 \quad * \quad w_3 \right. \quad \left. \begin{matrix} \text{old} \\ \end{matrix} \right\}$$

④ To fix this problem Researchers started exploring other Activation function

① Tanh ② ReLU ③ PReLU ④ Swiss

## Activation Functions

① Sigmoid Activation function  $[0 \text{ to } 1]$   $z = \sum_{i=1}^n w_i^T x + b$

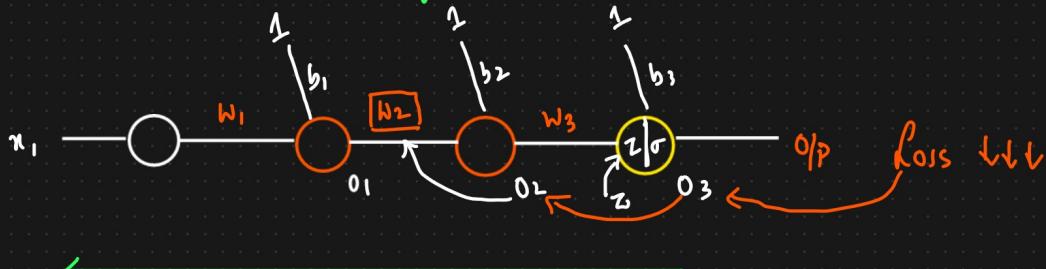


$$\sigma(z) \Rightarrow 0 \text{ to } 1$$

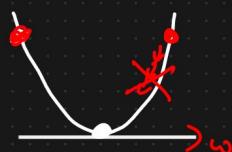
$$\phi(z) \Rightarrow$$

$$\frac{\partial \sigma(z)}{\partial z} = 0 \text{ to } 0.25$$

Forward Propagation



Backward Propagation.



$$w_{2\text{new}} = w_{2\text{old}} - \eta \left[ \frac{\partial L}{\partial w_{2\text{old}}} \right] \quad \left[ 0.001 \right] \quad \left[ 0.0001 \right] \quad \Rightarrow w_{2\text{new}} \approx w_{2\text{old}} \Rightarrow$$

$$\frac{\partial L}{\partial w_{2\text{old}}} = \frac{\partial L}{\partial o_3} * \boxed{\frac{\partial o_3}{\partial o_2}} * \frac{\partial o_2}{\partial w_2}$$

$\downarrow \downarrow \downarrow$

$$0.20 * 0.01 * \det z = (o_2 * w_3) + b_3$$

$$\frac{\partial o_3}{\partial o_2} = \frac{\partial (\sigma(z))}{\partial z} * \frac{\partial z}{\partial o_2}$$

$[0 \text{ to } 1]$

$$= [0 - 0.25] * \frac{\partial [(o_2 * w_3) + b_3]}{\partial o_2}$$

$$= [0 - 0.25] * w_3 \Rightarrow \text{Small value} \Rightarrow$$

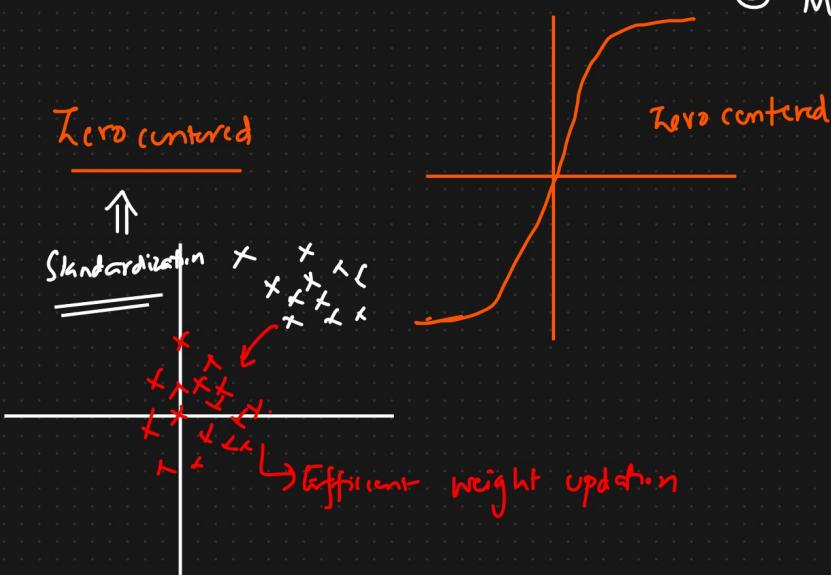
## Advantages

- ① Binary Classification Suitable.
- ② clear prediction i.e. very close 1 or 0

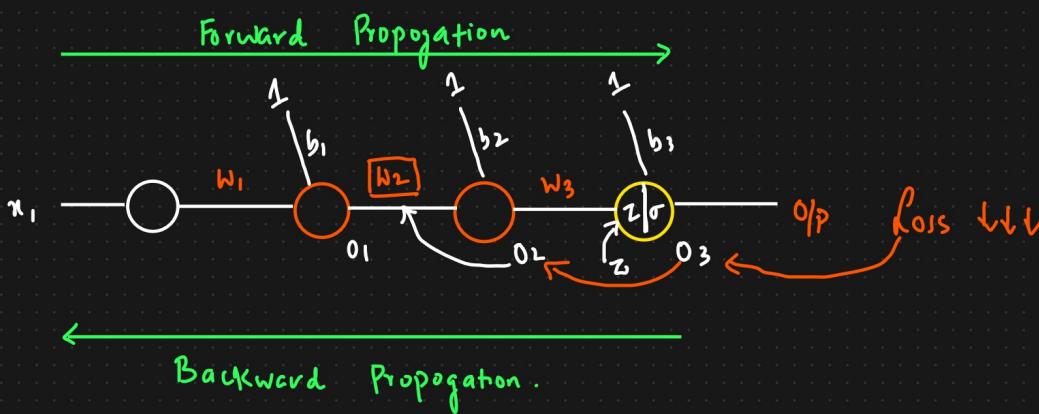
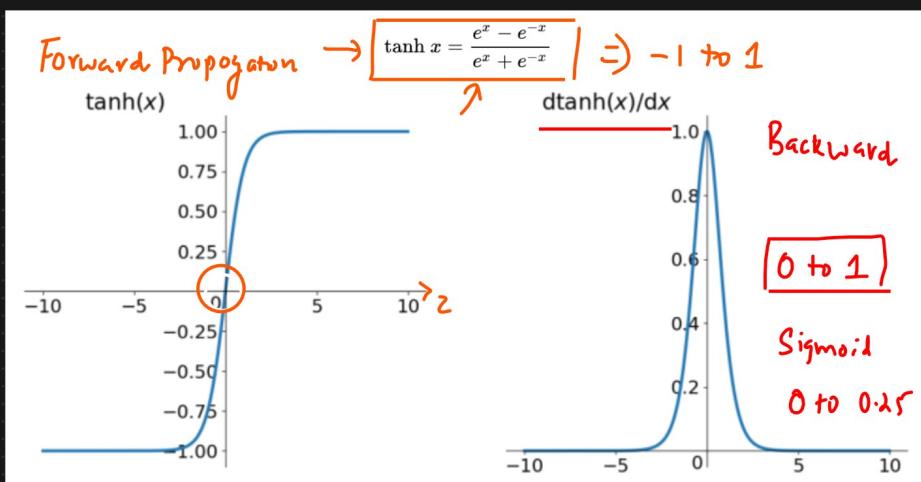
## Disadvantages

- ① Prone to vanishing Gradient Problem.
- ② Function output is not zero centered  $\Rightarrow$  Efficient weight update
- ③ Mathematical operation are relatively time consuming

Zero centered



## ② Tanh Activation Function



$$\frac{\partial h}{\partial w_{2012}} = \frac{\partial h}{\partial o_3} * \boxed{\frac{\partial o_3}{\partial o_2}} * \frac{\partial o_2}{\partial w_2}$$

↓↓↓

$$0 \cdot 20_{11} * 0 \cdot 01 * \text{det } z = (o_2 * w_3) + b_3$$

$$\begin{aligned}\frac{\partial o_3}{\partial o_2} &= \boxed{\frac{\partial (\tanh(z))}{\partial z}} * \frac{\partial z}{\partial o_2} && [0 \text{ to } 1] \\ &= [0 - 1] * \frac{\partial [(o_2 * w_3) + b_3]}{\partial o_2} \\ &= [0 - 1] * w_3 \Rightarrow \text{Small value} \Rightarrow\end{aligned}$$

Advantages

① Zero Centric  $\Rightarrow$  weight update is efficient

Disadvantages

① prone to vanishing gradient problem

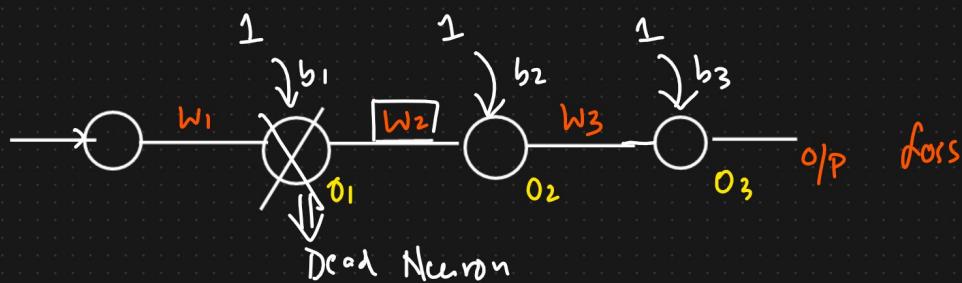
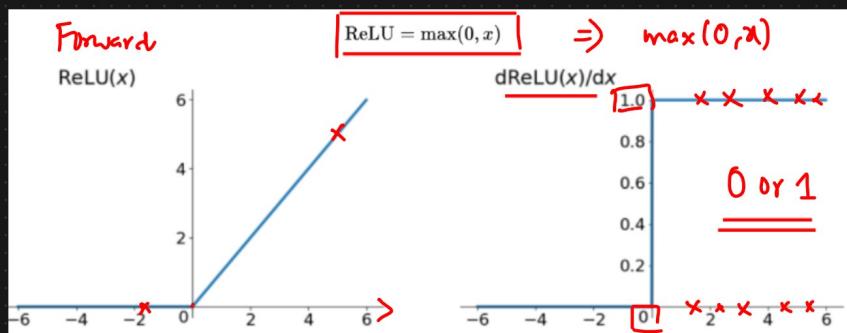
② Train complexity

[Rectified Linear Unit].

③ ReLU Activation Function

Tanh  $\Rightarrow 0 \text{ to } 1$

Sigmoid  $\Rightarrow 0 \text{ to } 0.25$



$$\frac{\partial h}{\partial w_{201d}} = \frac{\partial h}{\partial o_3} * \boxed{\frac{\partial o_3}{\partial o_2}} * \frac{\partial o_2}{\partial w_2}$$

$\downarrow \downarrow \downarrow \downarrow$

$$0.20_{11} * 0.01 * \text{let } z = (o_2 * w_3) + b_3$$

$$\frac{\partial o_3}{\partial o_2} = \boxed{\frac{\partial (\text{relu}(z))}{\partial z}} * \frac{\partial z}{\partial o_2}$$

$[0 \rightarrow 1]$

$$= [0 \text{ or } 1] * \frac{\partial [(o_2 * w_3) + b_3]}{\partial o_2}$$

$$= \begin{bmatrix} \text{-ve} & \text{+ve} \end{bmatrix} * w_3 \Rightarrow \text{Small value} \Rightarrow$$

If ReLU output is  $\boxed{1} \Rightarrow$  Weight updation will happen

If ReLU output is  $\boxed{0} \Rightarrow$  Dead Neuron

$$w_{2\text{new}} = w_{201d} - \eta \boxed{\frac{\partial h}{\partial w_{201d}}} \Rightarrow 0$$

If Derivative of  $\text{ReLU}(z)$  is 0

$$\boxed{w_{2\text{new}} \approx w_{201d}} \Rightarrow \text{Dead Neuron}$$

If  $z = \text{+ve}$   $\frac{\partial \text{ReLU}(z)}{\partial z} = 1$

If  $z = \text{-ve}$   $\frac{\partial \text{ReLU}(z)}{\partial z} = 0$

## Advantages

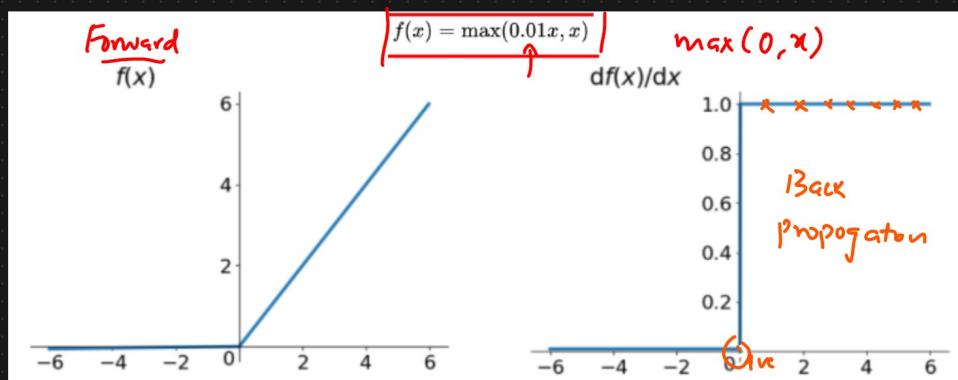
- ① Solving Vanishing Gradient Problem
- ②  $\text{Max}(0, x) \rightarrow$  Calculation is Superfast. The ReLU function has a linear relationship.
- ③ It is much faster than Sigmoid or Tanh.

## Disadvantages

- ① Dead Neuron
- ② ReLU function O/P  $(0, x) \Rightarrow 0$  or zero number  
↓  
It is not zero centric

## ④ Leaky ReLU And Parametric ReLU

$\max(\lambda x, x)$  → hyperparameter  $\lambda = \alpha = 0.01, 0.02, \dots, 0.03$



ReLU → Dead Neuron → Dead ReLU Problem

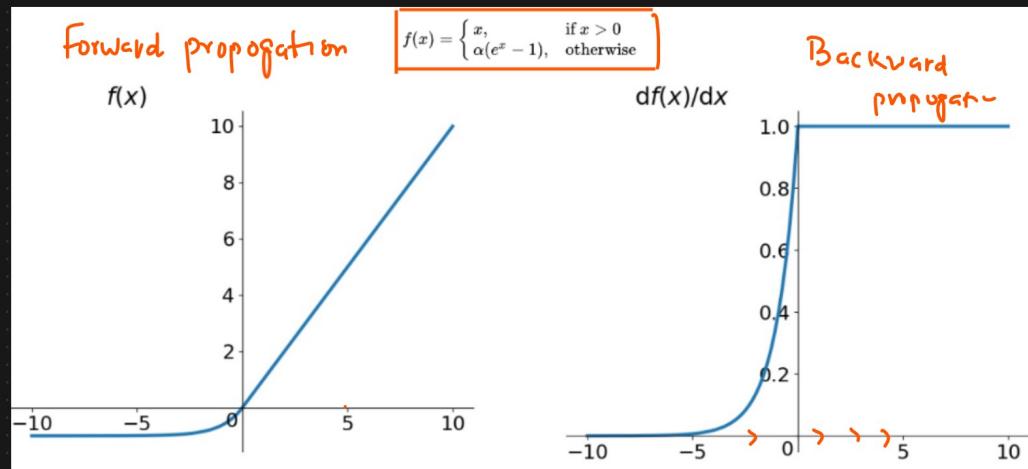
## Advantages

- ① Leaky ReLU has all the advantages of ReLU
- ② It removes the Dead ReLU Problem

## Disadvantage

- ① It is not zero centric

## ⑤ ELU (Exponential Linear Units)



Advantages

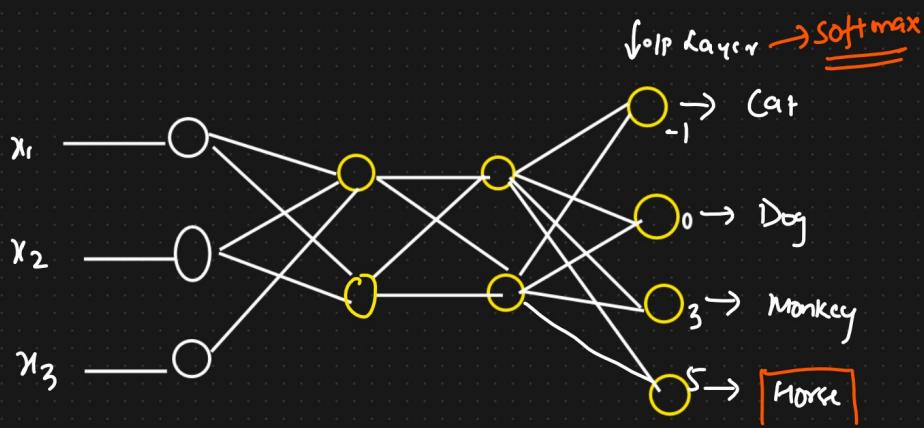
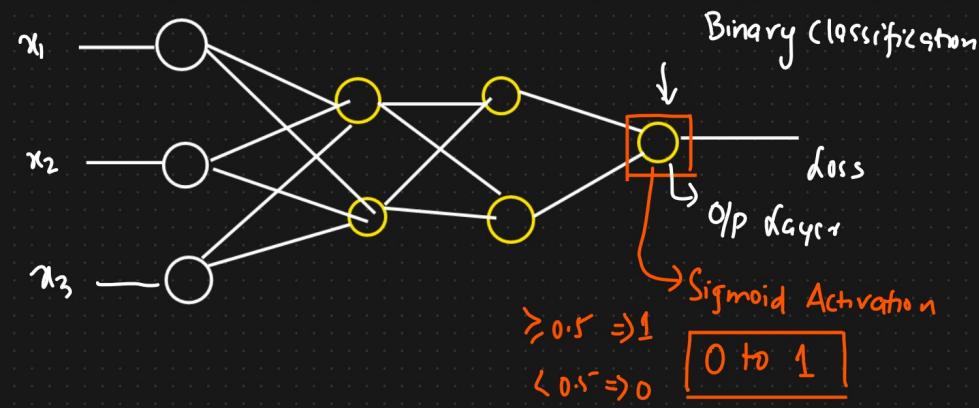
[It is used to solve  
ReLU problems]

Disadvantage

- ① No Dead ReLU Issues
- ② Zero centered

i) Slightly more computationally intensive.

## ⑥ Softmax Activation function [Multiclass classification problem]



$$\text{Softmax} = \frac{e^{y_i}}{\sum_{k=0}^n e^{y_k}}$$

$$y_i = \mathbf{w} \cdot \mathbf{x} + b$$

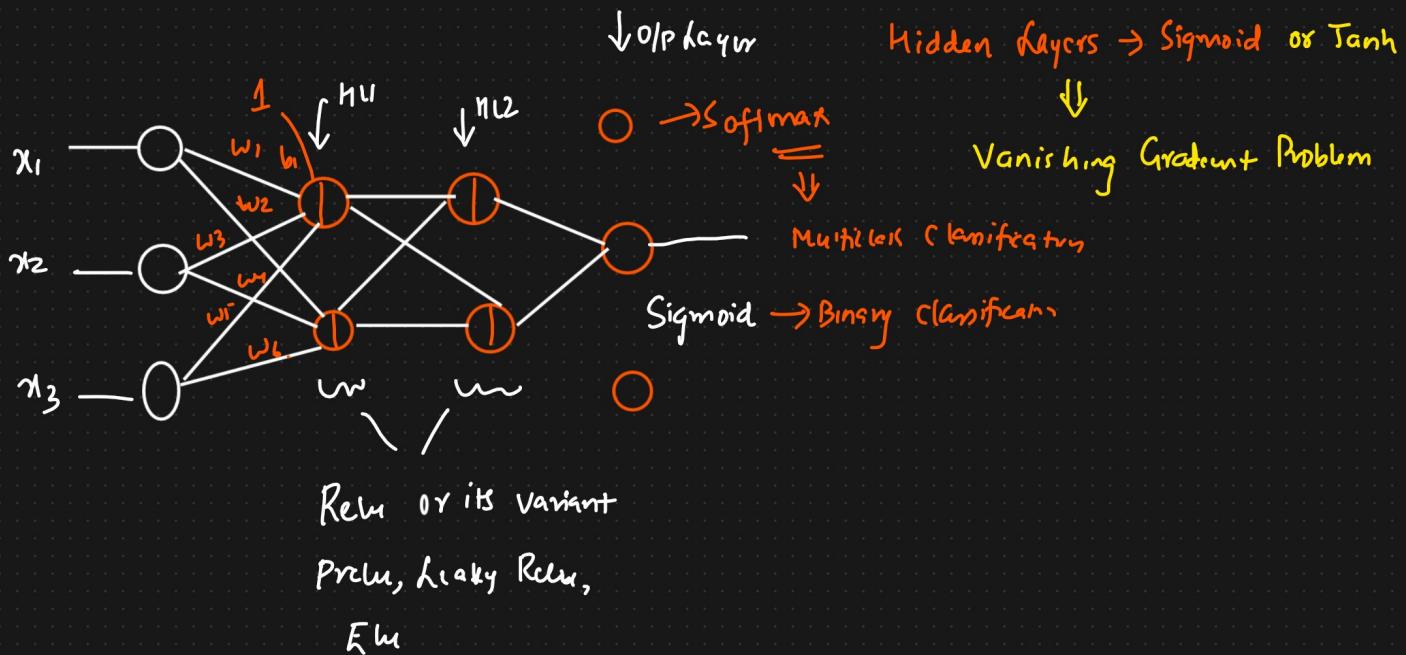
$$\text{Softmax} \Rightarrow \text{Cat} = \frac{e^{-1}}{e^{-1+0+3+5}} = 0.00033 \quad \Pr(\text{Horse}) = \frac{0.1353}{0.00033 + 0.0024 + 0.0183 + 0.1353}$$

$$\text{Dog} = \frac{e^0}{e^{-1+0+3+5}} = 0.0024 \quad \approx 86\%$$

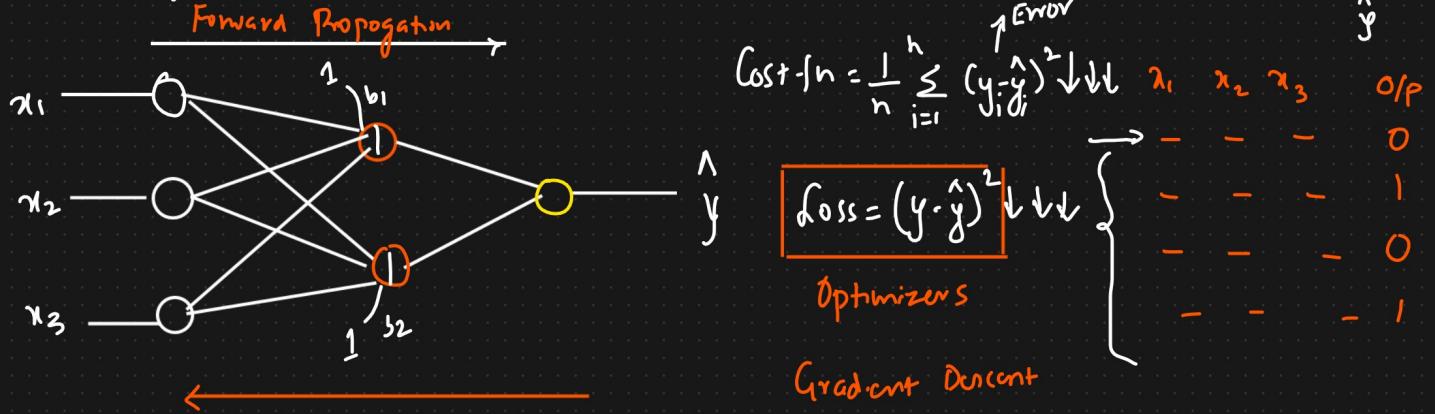
$$\text{Monkey} = \frac{e^3}{e^{-1+0+3+5}} = 0.0183$$

$$\text{Horse} = \frac{e^5}{e^{-1+0+3+5}} = 0.1353$$

## ⑦ Which Activation Function To Use When?



# Loss function And Cost Function

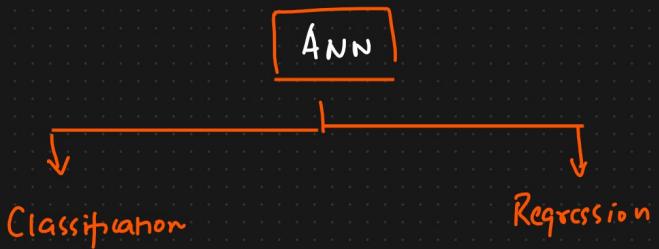


## Loss function

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

## Cost function

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



- ① Mean Squared Error (MSE)
- ② Mean Absolute Error (MAE)
- ③ Huber Loss
- ④ RMSE

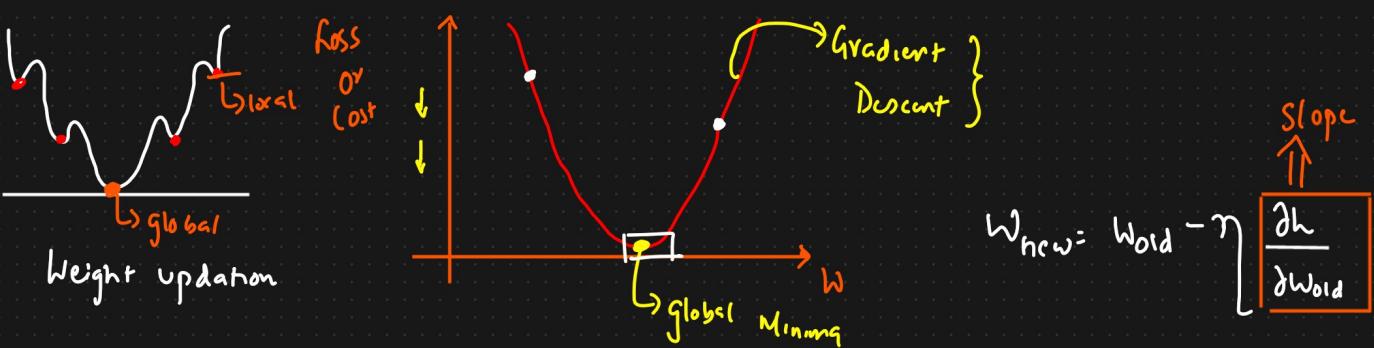
## ① Mean Squared Error (MSE)

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

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

↳ Quadratic Equations



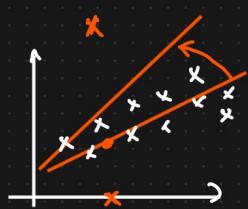


### Advantages

- ① MSE is Differentiable
- ② It has 1 local or global Minima
- ③ It converges faster
- ④ Mean Absolute Error (MAE)

### Disadvantages

- ① Not Robust to outliers → penalizing the error



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

$$\text{Loss fn} = |y - \hat{y}|$$

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

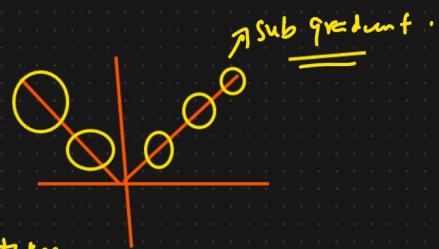
### Advantages

- ① Robust to outliers



### Disadvantages

- ① Convergence usually takes time in MAE



### ③ Huber loss

$$\text{① MSE}$$

$$\text{② MAE}$$

$$\text{Cost fn} = \begin{cases} \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 & \text{if } |y_i - \hat{y}_i| \leq \delta \\ \text{MSL} & \text{if } |y_i - \hat{y}_i| > \delta \end{cases}$$

No outlier ↑ Hypo-parameter

$$\begin{cases} \delta |y - \hat{y}| - \frac{1}{2} \delta^2, & \text{otherwise} \\ \downarrow \\ \text{MAE} \end{cases}$$

#### ④ RMSE (Root Mean Squared Error)

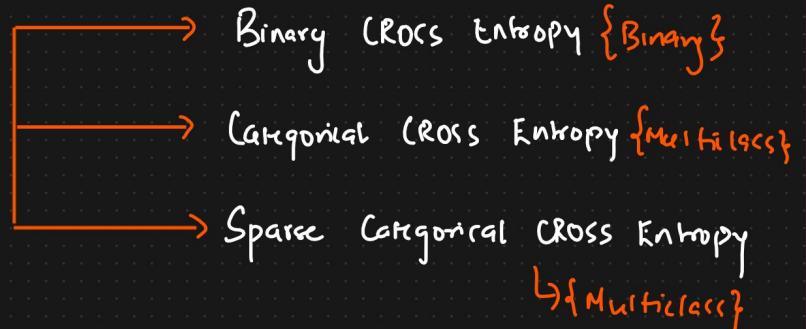
$$\text{Cost function} = \sqrt{\frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{N}}$$

Advantages

Disadvantages

#### ② Loss or Cost function For Classification Problems

Classification  $\rightarrow$  Cross Entropy



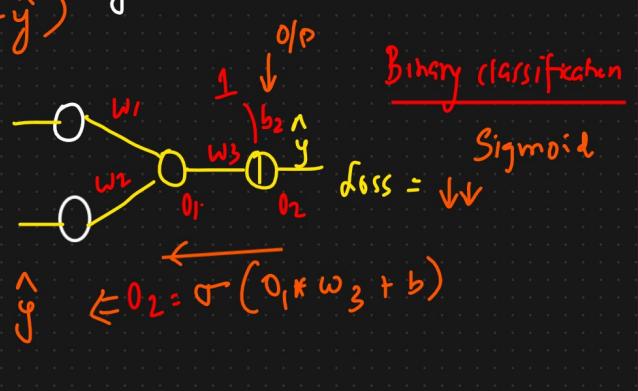
#### ① Binary Cross Entropy

$\downarrow$  Log Loss

$$\text{Loss} = -y * \log(\hat{y}) - (1-y) * \log(1-\hat{y})$$

$y$  = Actual Value

$\hat{y}$  = Predicted Value



$$\text{Loss} = \begin{cases} -\log(1-\hat{y}) & \text{if } y=0 \\ -\log(\hat{y}) & \text{if } y=1 \end{cases}$$

$$\hat{y} = \frac{1}{1+e^{-z}} \Rightarrow \text{Sigmoid Activation function}$$

## ② Categorical Cross Entropy (Multiclass classification)

ONE → One Hot Encoding

$f_1$	$f_2$	$f_3$	O/P	$j=1$ Good	$j=2$ Bad	$j=3$ Neutral	$C = \text{No. of categories}$
→ 2	3	4	Good	[1 0 0]			
→ 5	6	7	Bad	0 [1 0]			
→ 8	9	10	Neutral	0 0 [1]			

$i = 1 \rightarrow n$

$$L(x_i, y_i) = - \sum_{j=1}^C y_{ij} * \ln(\hat{y}_{ij})$$

Actual value  $\leftarrow y_{ij} = [y_{11} \ y_{12} \ y_{13} \ \dots \ y_{1c}]$

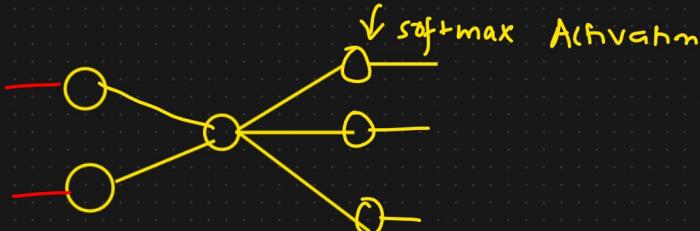
$$[y_{21}, y_{22}, y_{23}, \dots, y_{2c}]$$

⋮

⋮

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

Prediction  $\leftarrow \hat{y}_{ij} \Rightarrow \text{Softmax Activation} = S_{\text{of}}(z) =$



$$S_{\text{of}}(z) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

O/p  $\hat{y}_{ij}$  = Probabilities

$[0.1, 0.2, 0.3, 0.2, 0.2] \leftarrow 1$

Categorical  $\Rightarrow [0.2, 0.3, 0.5]$

(Cross Entropy)

      .

$\downarrow$   
[This also gives the probability of other categories]

### ③ Sparse Categorical Cross Entropy

$[0, 1^{\text{st}}, 2^{\text{nd}}]$   $\rightarrow$  Categories  
 $[0.2, 0.3, 0.5]$

Disadvantage  
① losing info about the probability of other category.

$\boxed{2^{\text{nd}} \text{ Index}} \Rightarrow \underline{\text{O/p}}$

$[0, 1, 2^{\text{nd}}, 3^{\text{rd}}, 4^{\text{th}}]$   
 $[0.2, 0.3, 0.1, 0.2, 0.2]$   $\Rightarrow 1^{\text{st}} \text{ Index}$   
Category  $\Rightarrow \underline{\text{O/p}}$

### ④ Right Combination

Activation applied

Hidden layers

O/p layer

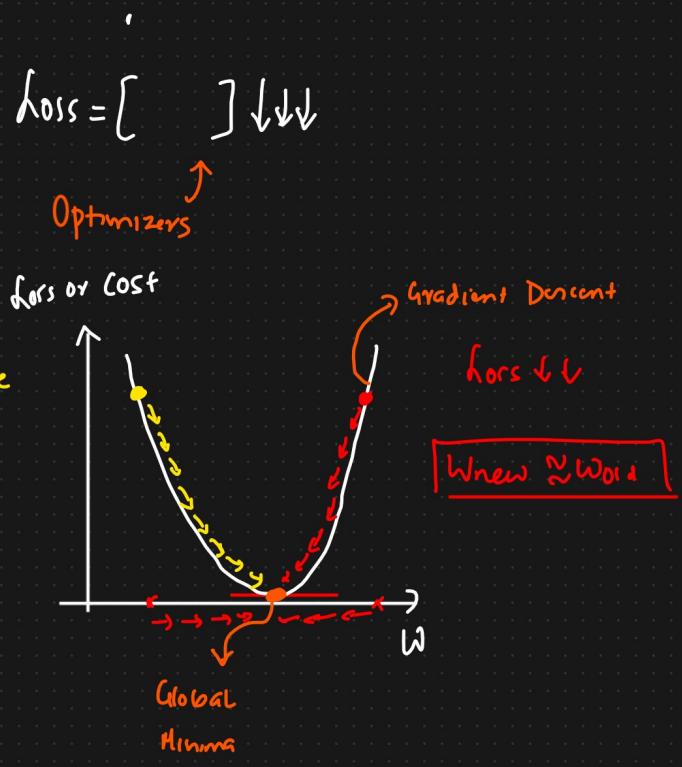
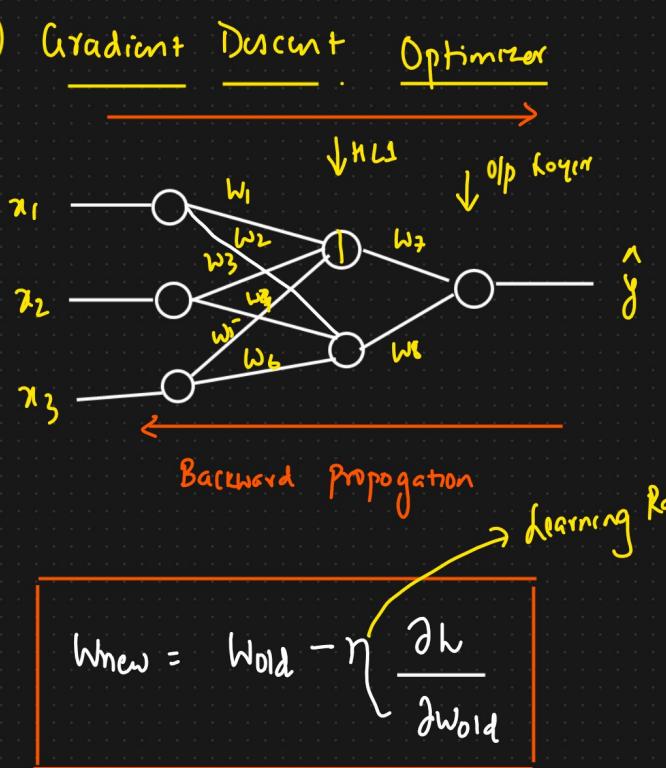
Problem Statement

Loss function

①	ReLU or its variants	Sigmoid	Binary Classification	Binary Cross Entropy
②	ReLU or its variants	Softmax	Multi Class	Categorical or Sparse CE
③	ReLU or its variants	Linear	Regression	MSE, MAE, L1/L2 loss, RMSE

## Optimizers

- ① Gradient Descent
  - ② Stochastic Gradient Descent (SGD)
  - ③ Mini batch SGD
  - ④ SGD With Momentum
  - ⑤ Adagrad and RMSprop
  - ⑥ Adam Optimizers



$$\text{Loss fn} = (\hat{y} - y)^2 \quad \text{Cost fn} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 \quad \text{Data size} = 1000 \text{ Data points}$$

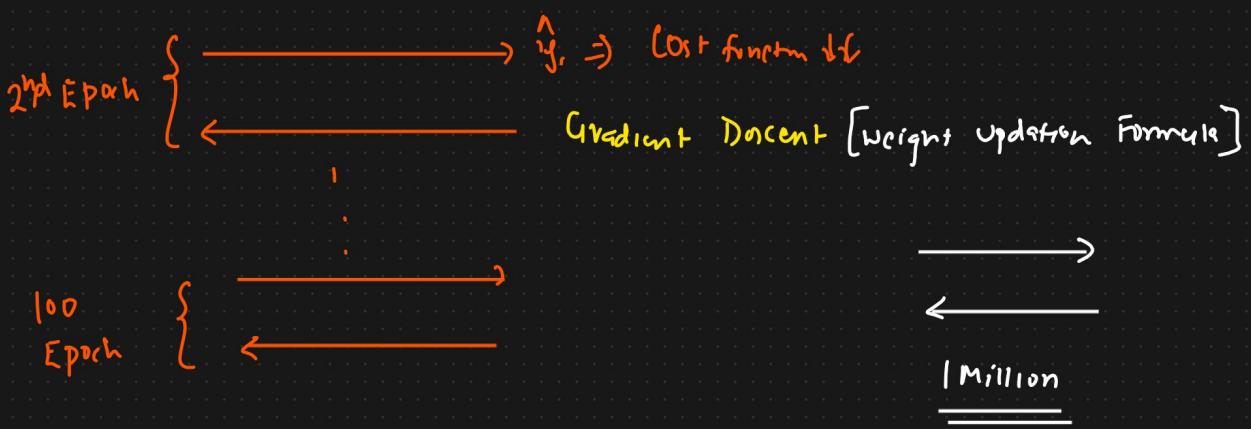
Epochs, Iteration

1 Epoch = 1 Iteration

10 Iteration

Back propagation

## Gradient Descent optimizes [Weight updation]



## Advantages

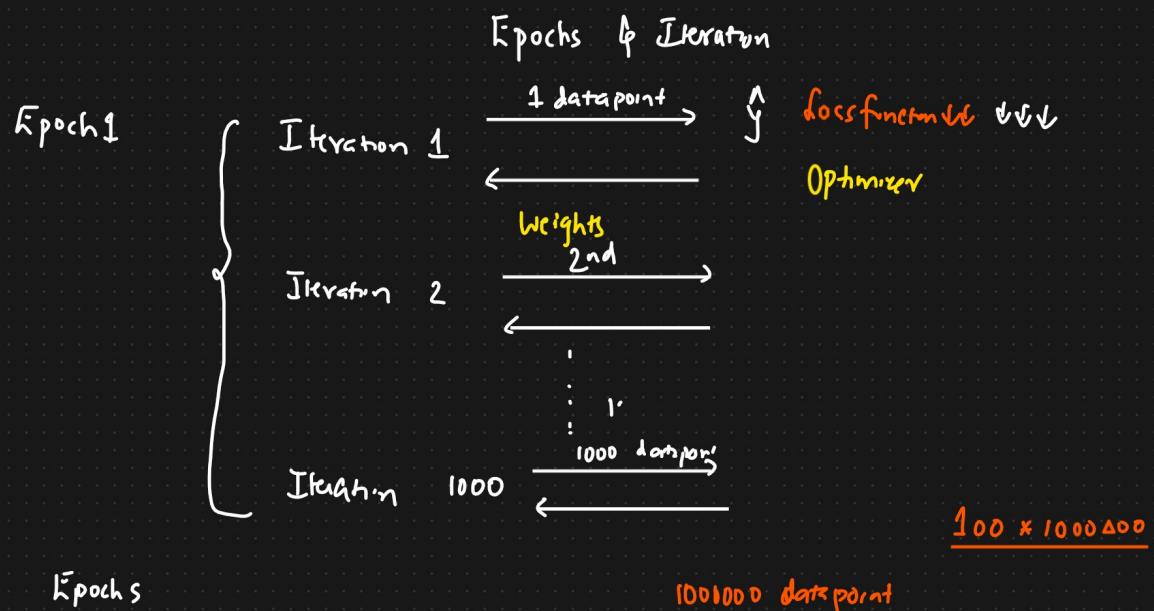
- ① Convergence will happen -

① Huge Resource RAM, GPU  
↓

## Resource Intensive -

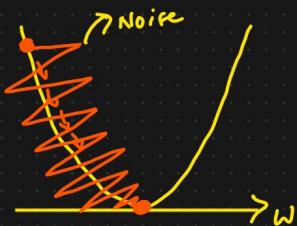
## ② Stochastic Gradient Descent (SGD)

1000 datapoint



## Advantage

- ## Ⓐ Solve Resource Issue



## Disadvantages

- ④ Time Complexity  $\uparrow\uparrow$
  - ④ Convergence will also take more time.
  - ④ Noise gets introduced

$\hat{y}$  Cost

### ③ Mini Batch SGD

Epoch, Iteration, Batch-size

$$\text{No. of iterations} = \frac{100000}{1000} = 100 \text{ iterations}$$

Data points = 100000

batch.size = 1000

MSGD

$$\text{Cost fn} = \sum_{i=1}^{1000} (y_i - \hat{y}_i)^2 \downarrow$$

Epoch 1

Iteration 1

Optimizer  $\Rightarrow$  Mini Batch SGD

change the weights

1000

$\downarrow$

[8gb] 1

Iteration 2

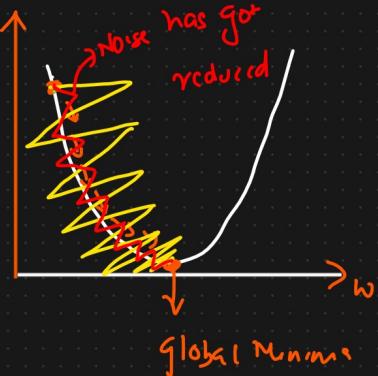
[16gb] 5000

Iteration 3

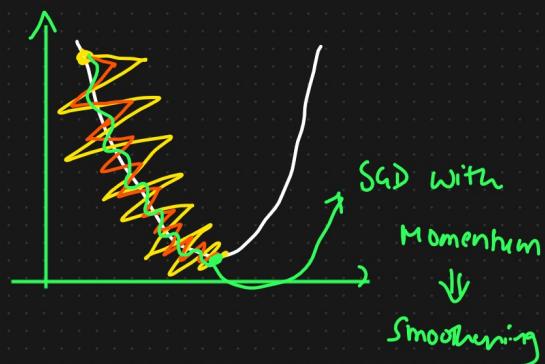
!

Iteration 100

Cost function



Batch Size



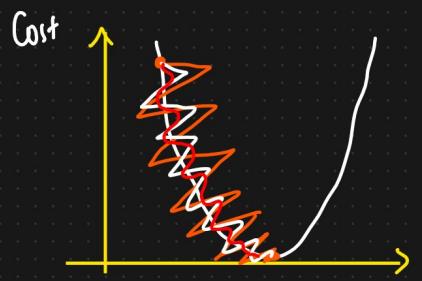
#### Advantages

- ① Convergence speed will increase
- ② Noise will be low when compared to SGD
- ③ Efficient Resource Usage (RAM)

#### Disadvantage

- ① Noise still exists

#### ④ SGD With Momentum



Weight Update formula

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial h}{\partial w_{\text{old}}} \quad \text{Learning Rate}$$

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

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

Exponential Weight Average {Smoothing} }  $\Rightarrow$  ARIMA, SARIMAX

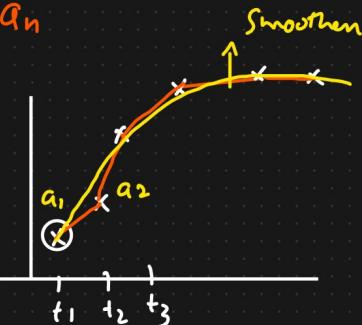
Time =  $t_1, t_2, t_3, t_4, \dots, t_n$  Time Series

Values =  $a_1, a_2, a_3, a_4, \dots, a_n$

$$V_{t_1} = a_1$$

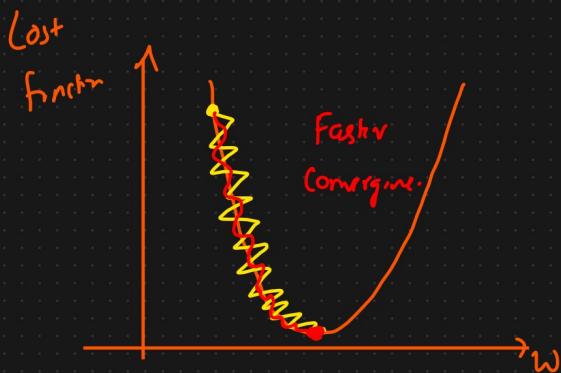
$$V_{t_2} = \boxed{\beta} * V_{t_1} + (1-\beta) * a_2$$

$$\begin{aligned} \beta = 0.95 \\ V_{t_2} = 0.95 * a_1 + (0.05) a_2 \end{aligned}$$



$$V_{t_3} = \beta * V_{t_2} + (1-\beta) * a_3$$

$$= 0.95 \left[ 0.95 * a_1 + (0.05) a_2 \right] + (0.05) * a_3$$



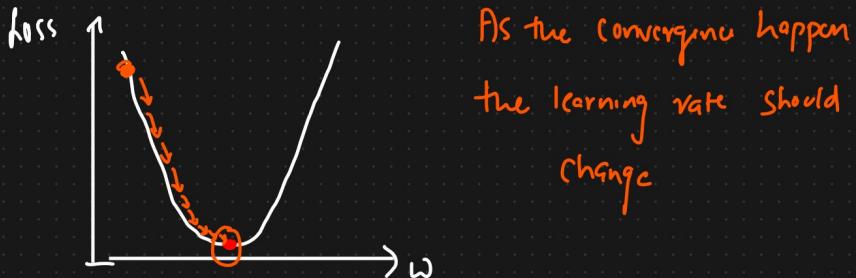
Advantage

- ① Reduces the noise
- ② Quick Convergence

⑤ Adagrad = Adaptive Gradient Descent  $\eta = \underline{\text{fixed}} \Rightarrow \underline{\text{Dynamic learning}}$

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

Learning Rate =  $0.001$



$$w_t = w_{t-1} - \eta' \left[ \frac{\partial h}{\partial w_{t-1}} \right]$$

$0.00001 \approx 0$

$\eta' = \frac{\eta}{\sqrt{d_t + \epsilon}}$   $\eta' \downarrow \text{when } d_t \uparrow$   
 $d_t \uparrow \Downarrow$   
 $d_t = \sum_{i=1}^t \left( \frac{\partial h}{\partial w_i} \right)^2$

$t=1 \quad t=2 \quad \rightarrow \quad t=3 \quad w_t \approx w_{t-1}$

$$\eta = 0.01 \quad \eta = 0.005 \quad \eta = 0.003$$

### Disadvantage

- ①  $\eta' \rightarrow$  Possibility to become a very small value  $\approx 0$

## ⑥ Adadelta And RMSprop

Exponential Weight Average

$$\beta = 0.95 \quad S_{dw_t} = 0$$

$$S_{dw_t} = \beta * S_{dw_{t-1}} + (1-\beta) \left( \frac{\partial h}{\partial w_{t-1}} \right)^2$$

$$\eta' = \frac{\eta}{\sqrt{S_{dw_t} + \epsilon}}$$

Dynamic LR + Smoothing [EWA]

$$w_t = w_{t-1} - \eta' \frac{\partial h}{\partial w_{t-1}}$$

## ⑦ Adam Optimizer

SGD with Momentum + RMSprop [ Dynamic LR + Smoothing ]

$$w_t = w_{t-1} - \eta' V_{dw}$$

$$b_t = b_{t-1} - \eta' V_{db}$$

Weight Updation

Bias Updation

$$\eta' = \frac{\eta}{\sqrt{S_{dw_t} + \epsilon}}$$

EWA

$S_{dw_t} = 0$

$$S_{dw_t} = \beta * S_{dw_{t-1}} + (1-\beta) \left( \frac{\partial h}{\partial w_{t-1}} \right)^2$$

[ EWA ]

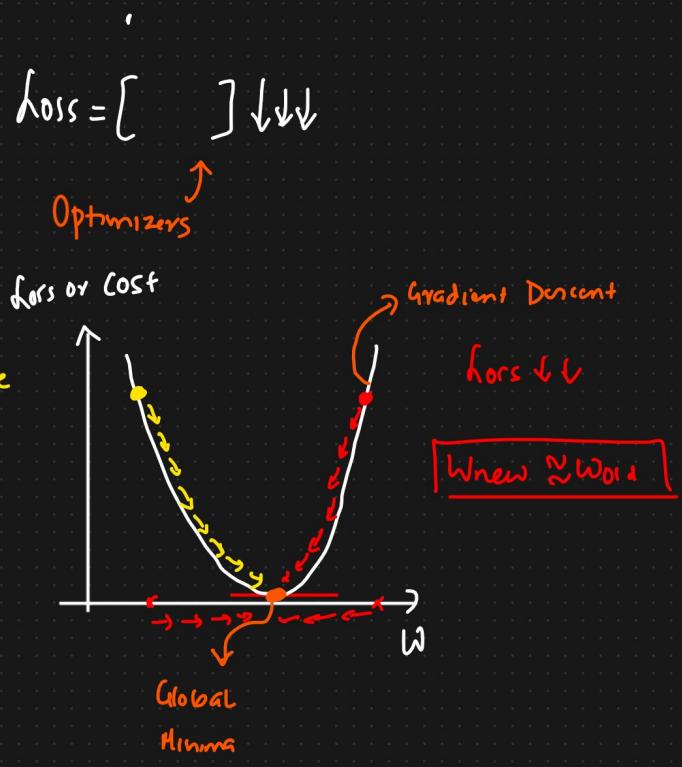
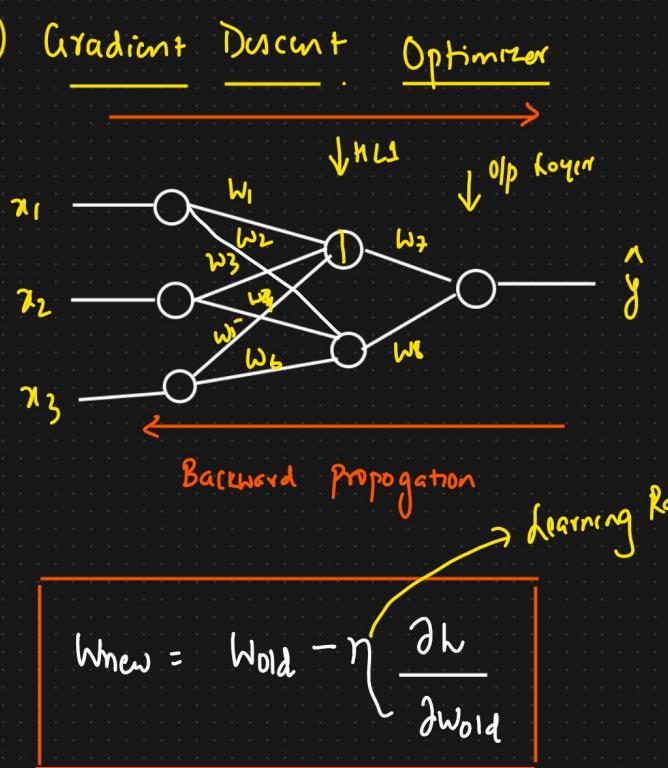
$$V_{dw_t} = \beta * V_{dw_{t-1}} + (1-\beta) \frac{\partial h}{\partial w_{t-1}}$$

$$V_{db_t} = \beta * V_{db_{t-1}} + (1-\beta) \frac{\partial h}{\partial b_{t-1}}$$

$\Rightarrow$  Momentum  
↳ Smoothing

## Optimizers

- ① Gradient Descent
  - ② Stochastic Gradient Descent (SGD)
  - ③ Mini batch SGD
  - ④ SGD With Momentum
  - ⑤ Adagrad and RMSprop
  - ⑥ Adam Optimizers



$$\text{Loss fn} = (\hat{y} - y)^2 \quad \text{Cost fn} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 \quad \text{Dataset} = 1000 \text{ Data Points}$$

Epochs, Iteration

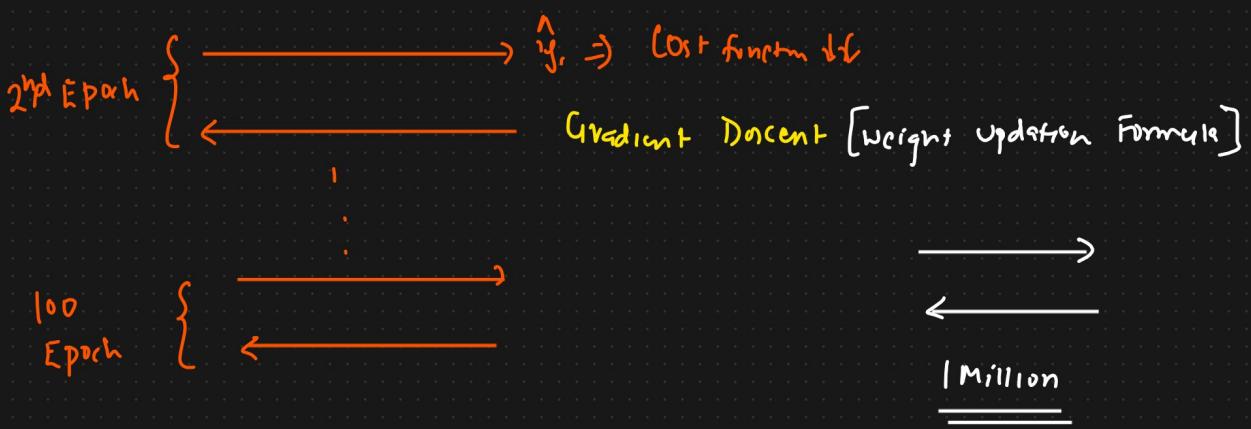
1 Epoch = 1 Iteration

10 Iteration

$J(\theta) = 100$

Final point

## Gradient Descent optimizes [Weight updation]



## Advantages

- ① Convergence will happen -

## Disadvantage

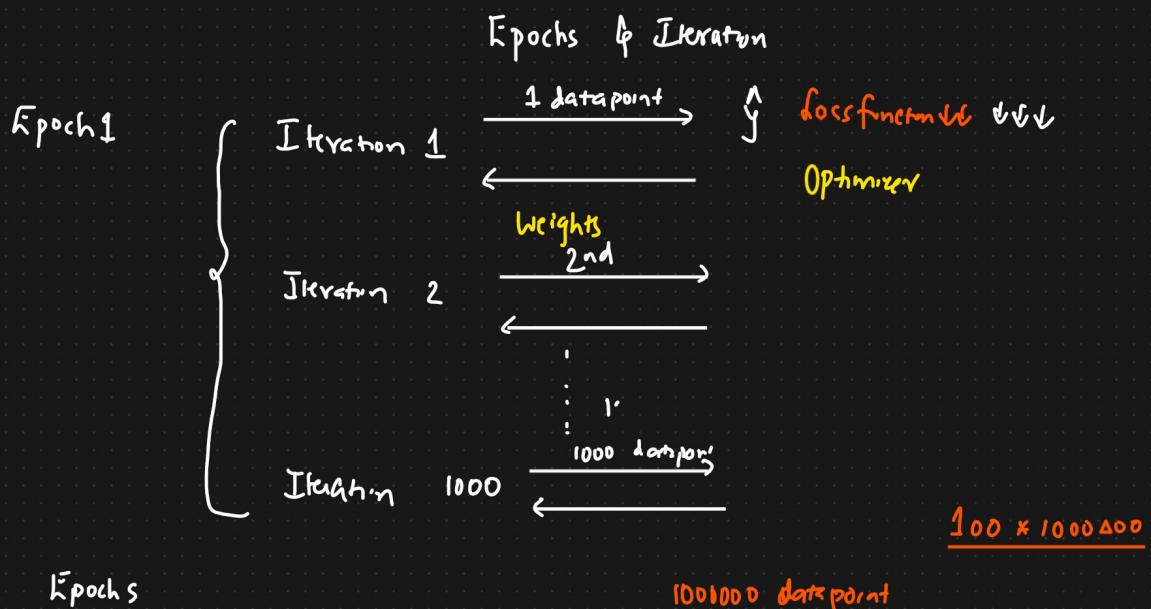
- ① Huge Resource RAM, GPU

11

## Resource Intensive -

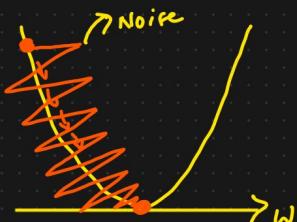
## ② Stochastic Gradient Descent (SGD)

1000 datapoint



## Advantage

- ## Ⓐ Solve Resource Issue



## Disadvantages

- ⑥ Time Complexity  $\uparrow\uparrow$

- ④ Convergence will also take more time.
  - ⑤ Noise gets introduced

$\hat{y}$  Cost

### ③ Mini Batch SGD

Epoch, Iteration, Batch-size

$$\text{No. of iterations} = \frac{100000}{1000} = 100 \text{ iterations}$$

Data points = 100000

batch.size = 1000

MSGD

$$\text{Cost fn} = \sum_{i=1}^{1000} (y_i - \hat{y}_i)^2 \downarrow$$

Epoch 1

Iteration 1

Optimizer  $\Rightarrow$  Mini Batch SGD

change the weights

1000

$\downarrow$

[8gb] 1

Iteration 2

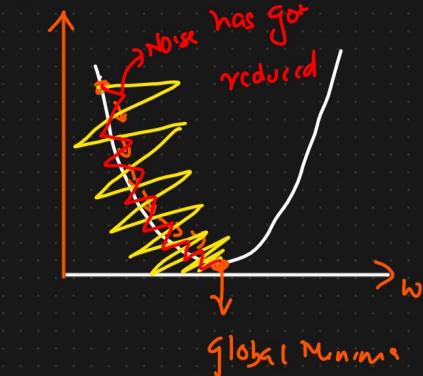
[16gb] 5000

Iteration 3

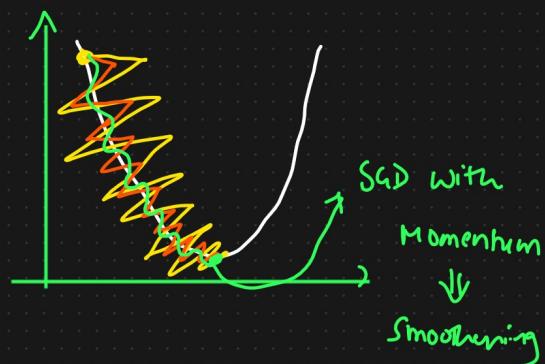
Iteration 100

Iteration 100

Cost function



Batch Size



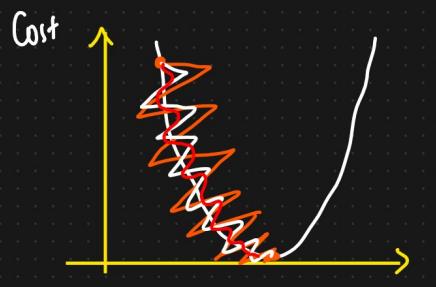
#### Advantages

- ① Convergence speed will increase
- ② Noise will be low when compared to SGD
- ③ Efficient Resource Usage (RAM)

#### Disadvantage

- ① Noise still exists

#### ④ SGD With Momentum



Weight Update formula

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial h}{\partial w_{\text{old}}} \quad \text{Learning Rate}$$

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

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

Exponential Weight Average {Smoothing} }  $\Rightarrow$  ARIMA, SARIMAX

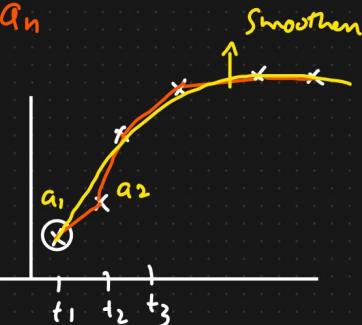
Time =  $t_1, t_2, t_3, t_4, \dots, t_n$  Time Series

Values =  $a_1, a_2, a_3, a_4, \dots, a_n$

$$V_{t_1} = a_1$$

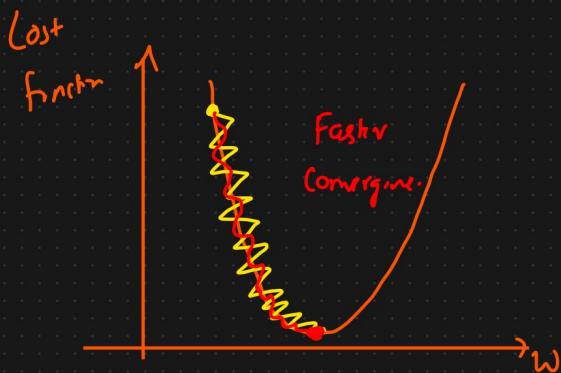
$$V_{t_2} = \boxed{\beta} * V_{t_1} + (1-\beta) * a_2$$

$$\begin{aligned} \beta = 0.95 \\ V_{t_2} = 0.95 * a_1 + (0.05) a_2 \end{aligned}$$



$$V_{t_3} = \beta * V_{t_2} + (1-\beta) * a_3$$

$$= 0.95 \left[ 0.95 * a_1 + (0.05) a_2 \right] + (0.05) * a_3$$



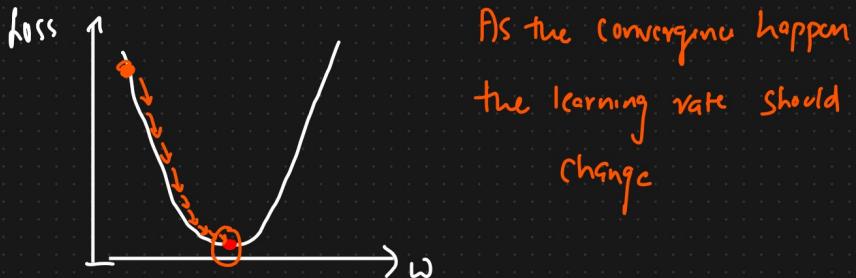
Advantage

- ① Reduces the noise
- ② Quick Convergence

⑤ Adagrad = Adaptive Gradient Descent  $\eta = \underline{\text{fixed}} \Rightarrow \underline{\text{Dynamic learning}}$

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

Learning Rate =  $0.001$



$$w_t = w_{t-1} - \eta' \left[ \frac{\partial h}{\partial w_{t-1}} \right]$$

$0.00001 \approx 0$

$\eta' = \frac{\eta}{\sqrt{d_t + \epsilon}}$   $\eta' \downarrow \text{when } d_t \uparrow$   
 $d_t \uparrow \Downarrow$   
 $d_t = \sum_{i=1}^t \left( \frac{\partial h}{\partial w_i} \right)^2$

$\epsilon \text{ is a small value}$

$t=1 \quad t=2 \quad \rightarrow \quad t=3 \quad w_t \approx w_{t-1}$

$$\eta = 0.01 \quad \eta = 0.005 \quad \eta = 0.003$$

### Disadvantage

- ①  $\eta' \rightarrow$  Possibility to become a very small value  $\approx 0$

## ⑥ Adadelta And RMSprop

Exponential Weight Average

$$\beta = 0.95 \quad S_{dw_t} = 0$$

$$S_{dw_t} = \beta * S_{dw_{t-1}} + (1-\beta) \left( \frac{\partial h}{\partial w_{t-1}} \right)^2$$

$$\eta' = \frac{\eta}{\sqrt{S_{dw_t} + \epsilon}}$$

Dynamic LR + Smoothing [EWA]

$$w_t = w_{t-1} - \eta' \frac{\partial h}{\partial w_{t-1}}$$

## ⑦ Adam Optimizer

SGD with Momentum + RMSprop [ Dynamic LR + Smoothing ]

$$w_t = w_{t-1} - \eta' V_{dw}$$

$$b_t = b_{t-1} - \eta' V_{db}$$

weight  
Bias  
Updation

$$\eta' = \frac{\eta}{\sqrt{S_{dw_t} + \epsilon}}$$

EWA

$S_{dw_t} = 0$

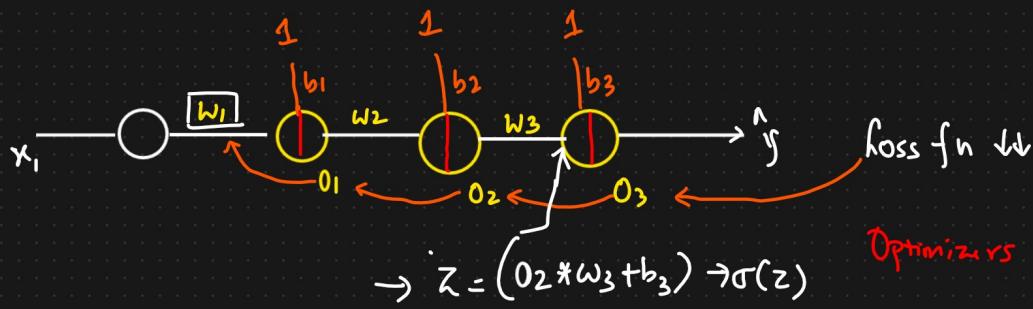
$$S_{dw_t} = \beta * S_{dw_{t-1}} + (1-\beta) \left( \frac{\partial h}{\partial w_{t-1}} \right)^2$$

$$V_{dw_t} = \beta * V_{dw_{t-1}} + (1-\beta) \frac{\partial h}{\partial w_{t-1}}$$

$$V_{db_t} = \beta * V_{db_{t-1}} + (1-\beta) \frac{\partial h}{\partial b_{t-1}}$$

⇒ Momentum  
↳ Smoothing

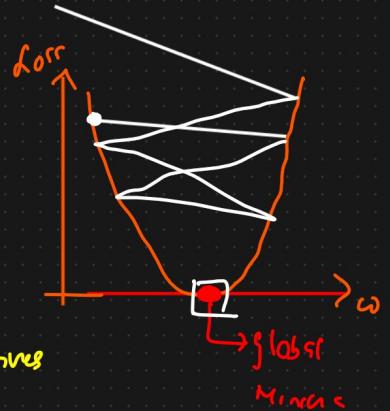
# Exploding Gradient Problem $\Rightarrow$ Weight Initialization Technique



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

$\left\{ \begin{array}{l} w_{\text{new}} >>> w_{\text{old}} \\ w_{\text{new}} <<< w_{\text{old}} \end{array} \right\}$

↗ Chain Rule of Derivatives



$$\frac{\partial h}{\partial w_{\text{old}}} = \frac{\partial h}{\partial o_3} * \left[ \frac{\partial o_3}{\partial o_2} \right] * \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_{\text{old}}}$$

big \* big \* big \* big  $\Rightarrow$  big value

Weight initialization



Very high value

$$\frac{\partial o_3}{\partial o_2} = \left[ \frac{\partial \sigma(z)}{\partial z} \right] * \frac{\partial z}{\partial o_2}$$

$$= [0 - 0.25] * \frac{\partial (o_2 * w_3 + b_3)}{\partial o_2}$$

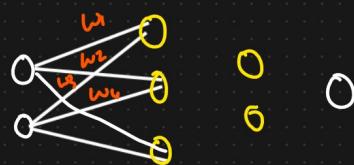
$$= [0 - 0.25] * \frac{w_3}{w_3} = \underline{\underline{5.00 - 1000}}$$

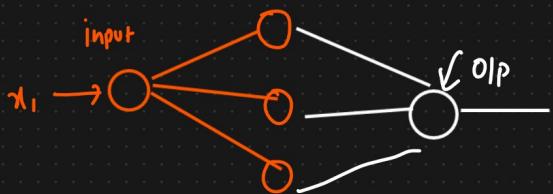
## Weight Initializing Techniques

- ① Uniform Distribution ✓
- ② Xavier / Glorot Initialization ✓
- ③ Kaiming He Initialization ✓

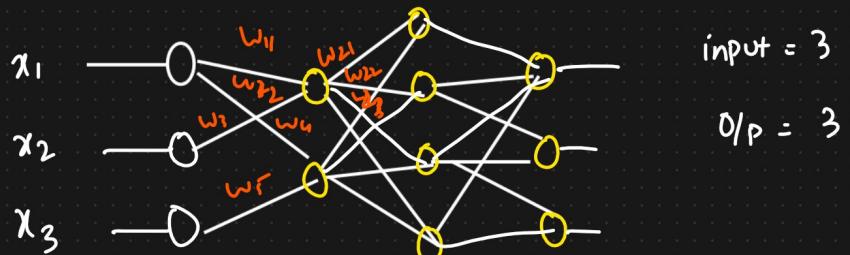
## Key Points

- ① Weights should be small ✓
- ② Weights should not be same ✓
- ③ Weights should have good variance }





input = 1  
Output = 1



input = 3

O/P = 3

## ① Uniform Distribution

$$w_{ij} \sim \text{Uniform Distribution} \left[ \frac{-1}{\sqrt{\text{input}}}, \frac{1}{\sqrt{\text{input}}} \right]$$

$$\left[ \frac{-1}{\sqrt{3}}, \frac{1}{\sqrt{3}} \right]$$

## ② Xavier / Glorot Initialization

Recurrent  $\rightarrow$  Xavier Glorot

### ① Xavier Normal Init

$$w_{ij} \sim N(0, \sigma)$$

$$\sigma = \sqrt{\frac{2}{(\text{input} + \text{output})}}$$

### ② Xavier Uniform

$$w_{ij} \sim \text{Uniform Distribution}$$

$$\left[ \frac{-\sqrt{6}}{\sqrt{\text{input} + \text{output}}}, \frac{\sqrt{6}}{\sqrt{\text{input} + \text{output}}} \right]$$

### (3) Kaiming He Initialization

#### ① He Normal

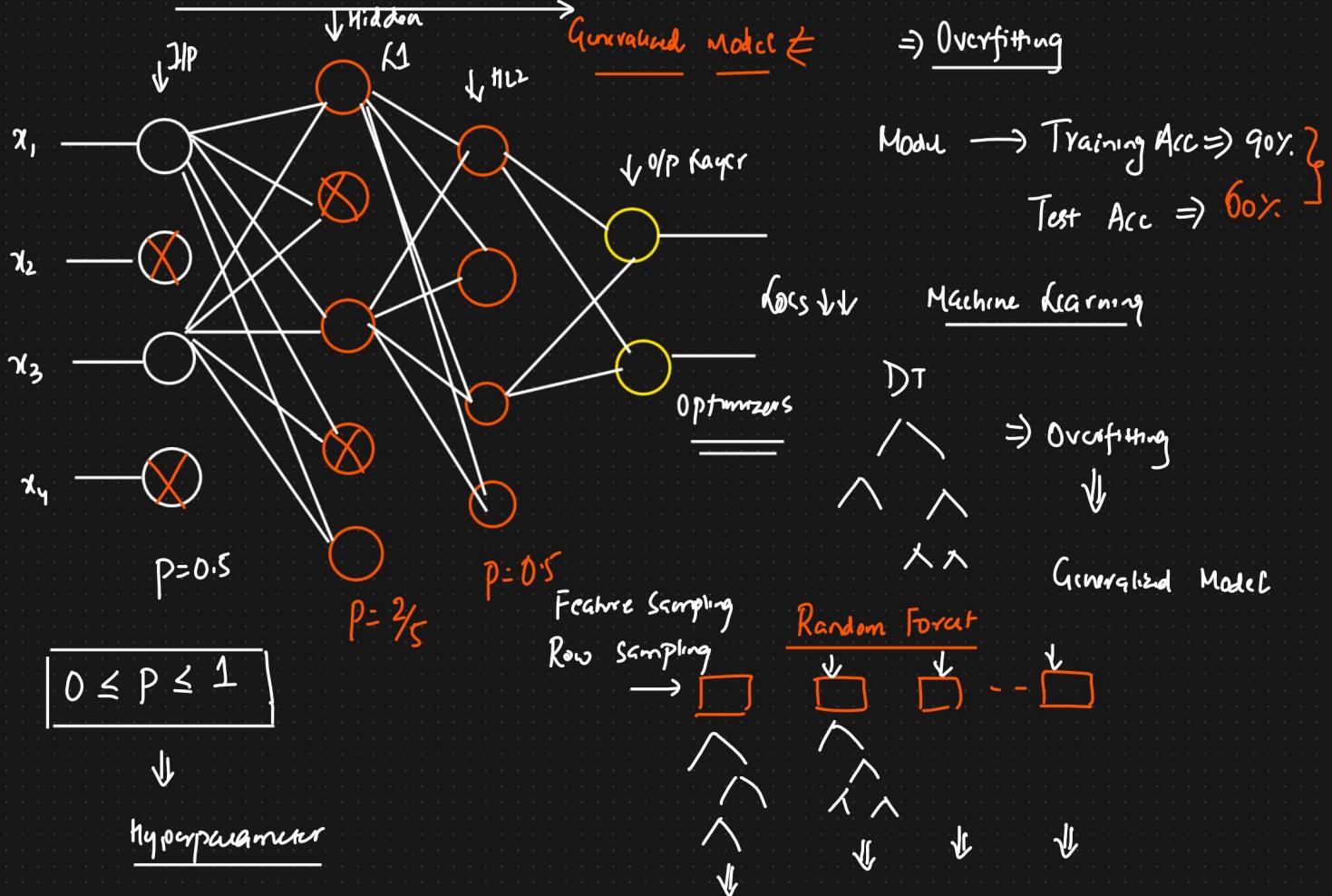
$$w_{ij} \sim N(0, \sigma)$$

$$\sigma = \sqrt{\frac{2}{\text{input}}}$$

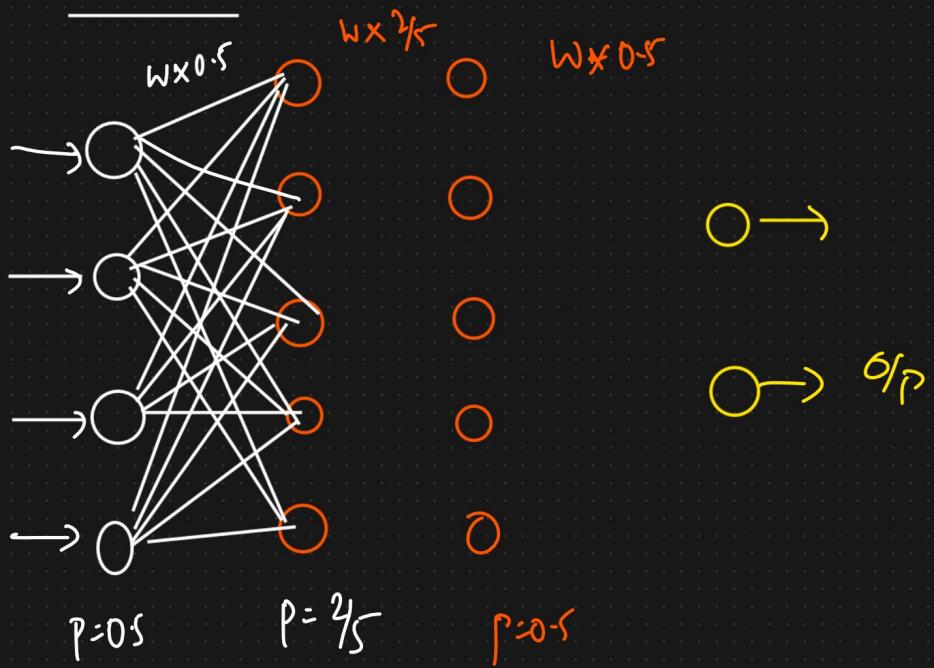
#### ② He uniform

$$w_{ij} \sim \text{Uniform Distribution} \left[ -\sqrt{\frac{6}{\text{input}}}, \sqrt{\frac{6}{\text{input}}} \right]$$

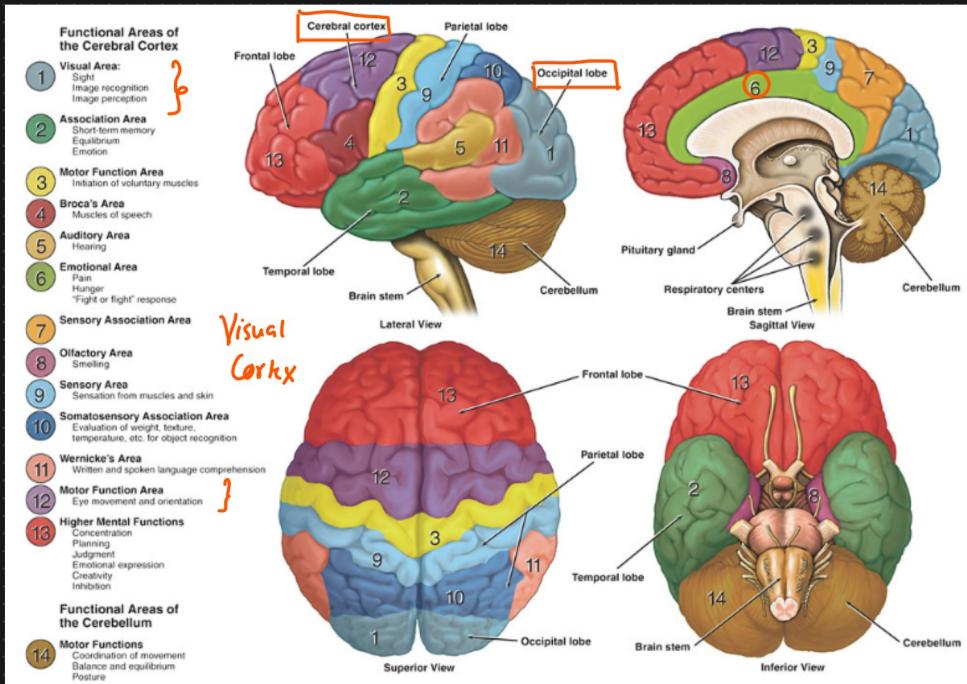
## Drop Out layer



Test data



# Convolutional Neural N/w



<https://www.dana.org>

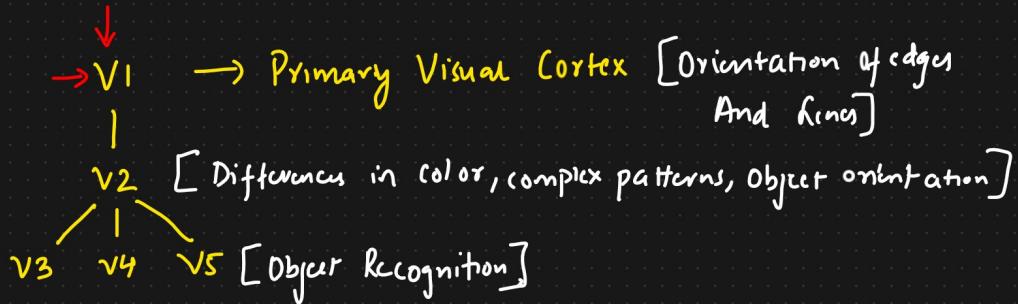
① ANN → Supervised Learning →   
 Classification  
 Regression

Dataset : I/p features O/P

② CNN : I/p ⇒ Images Eg: Image classification,  
 Object Detection, Segmentation

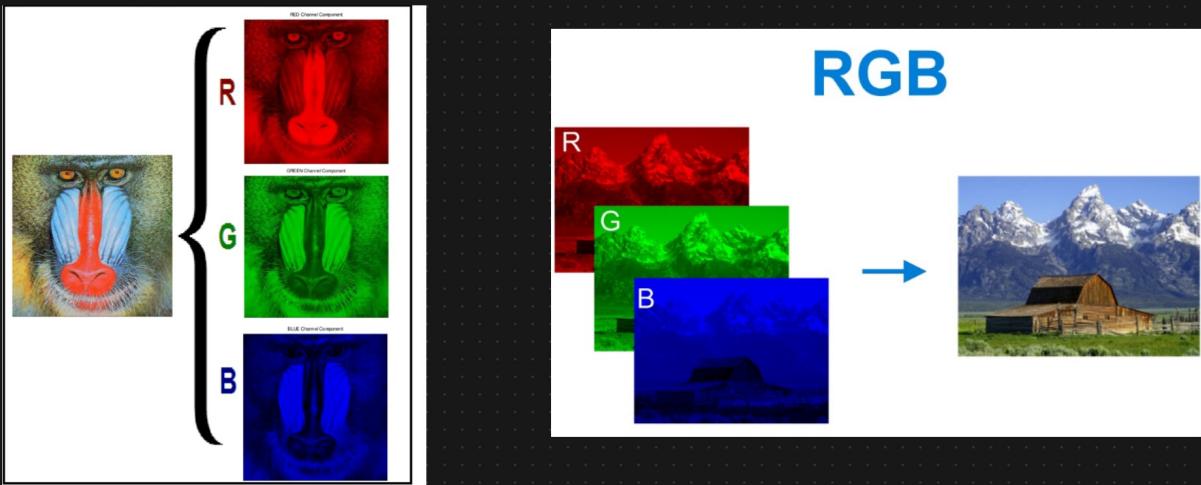
② Cerebral Cortex And Visual Cortex

Visual Cortex (VI-V5) [Region of the brain that receives, integrates and processes visual information relayed from the retinas].

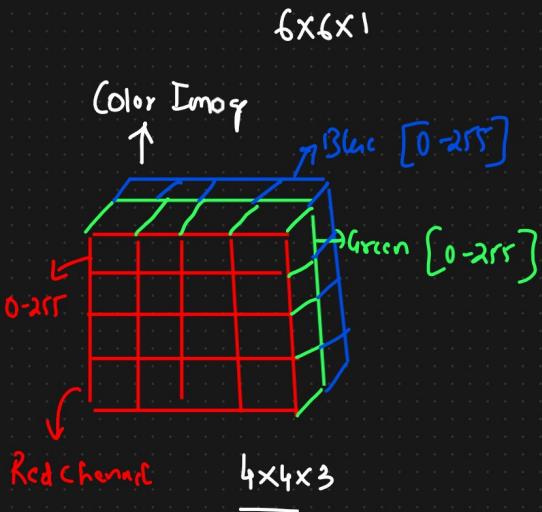
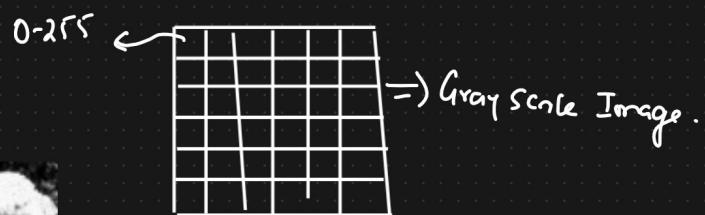
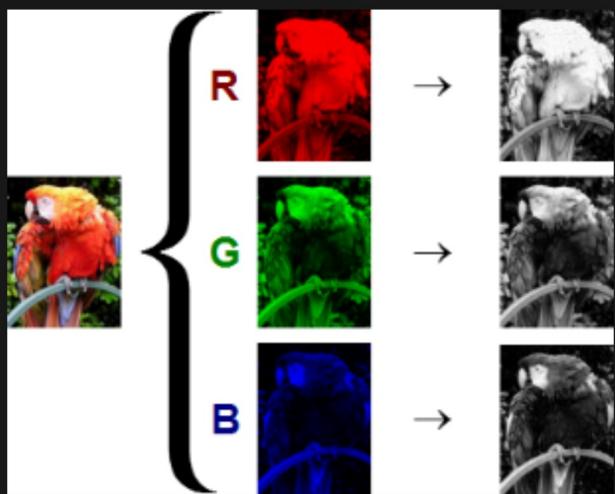


Visualize the Image

### ③ RGB Images And Gray Scale Images



<https://www.researchgate.net/>



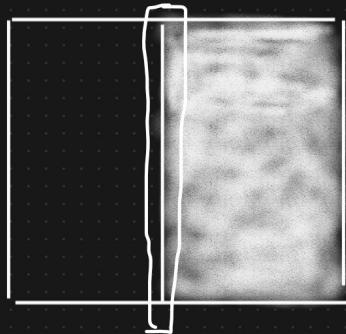
<https://commons.wikimedia.org/>

### ④ Convolution Operation In CNN

→ (0,1)

0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255

→



Convolution operation

Step 1

① Normalize

Divide by 255

+1	+2	-1
0	0	0
-1	-2	-1

$S\text{tride}=1$

0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1

$n=6$

\*

$f=3$

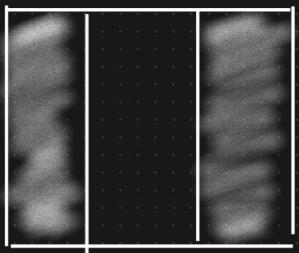
+1	0	-1
+2	0	-2
+1	0	-1

=

0	-4	-4	0
0	-4	-4	0
0	-4	-4	0
0	-4	-4	0

$4 \times 4$

$6 \times 6 \times 1$



filters  
vertical edge filters

$$h - f + 1 = \\ = 6 - 3 + 1 = 4$$

arr	0	0	arr
arr	0	0	arr
arr	0	0	arr
arr	0	0	arr

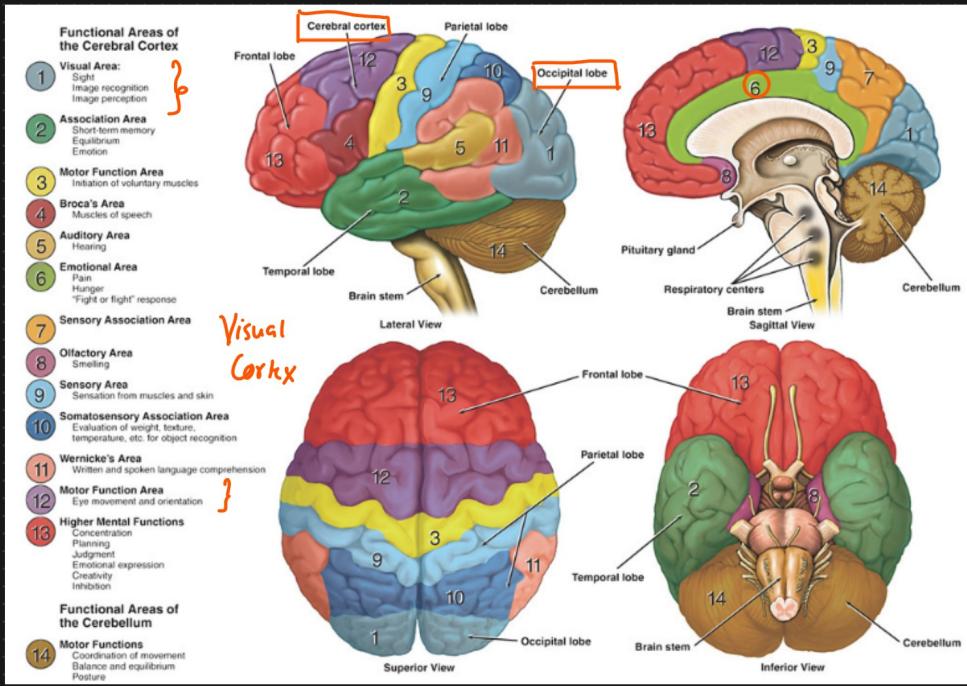
0	0	o		
0	o	o		
0	o	o		
0	o	o		
o	o	o		
0	o	o		
0	o	o		

\*

+1	0	-1
+2	0	-2
+1	0	-1

0	-4	-4	0
0	-4	-4	0
0	-4	-4	0
0	-4	-4	0

# Convolutional Neural N/w



<https://www.dana.org>

① ANN → Supervised Learning →   
 Classification  
 Regression

Dataset : I/p features O/P

② CNN : I/p ⇒ Images Eg: Image classification,  
 Object Detection, Segmentation

② Cerebral Cortex And Visual Cortex

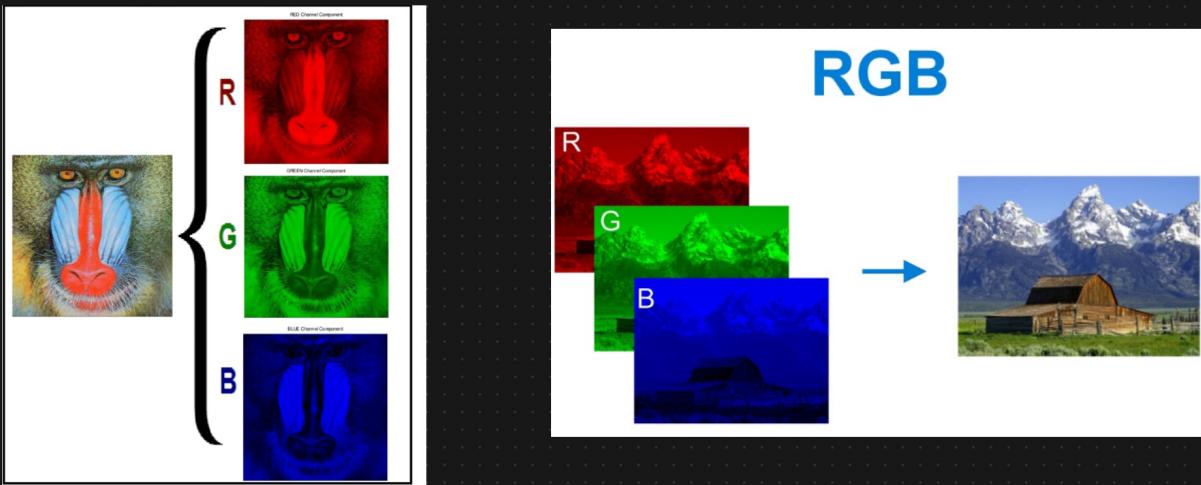
Visual Cortex (VI-V5) [Region of the brain that receives, integrates and processes visual information relayed from the retinas].

↓  
→ VI → Primary Visual Cortex [Orientation of edges And lines]

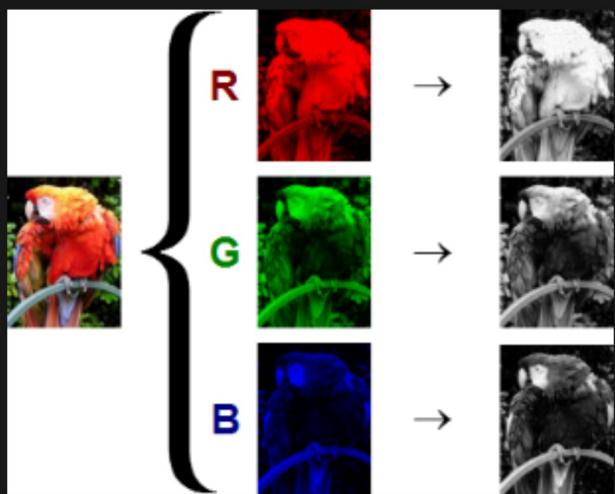
V1  
V2 [Differences in color, complex patterns, object orientation]  
V3 V4 V5 [Object Recognition]

Visualize the Image

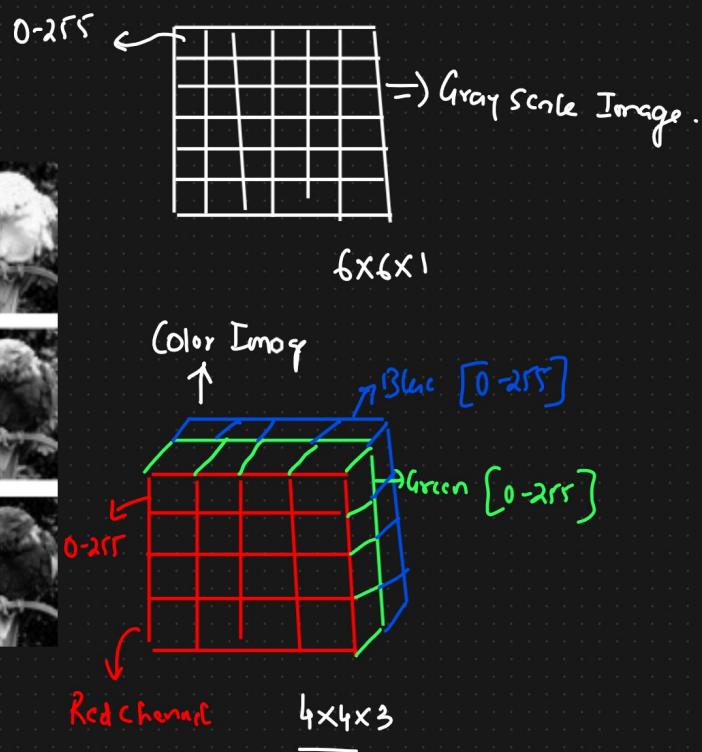
### ③ RGB Images And Gray Scale Images



<https://www.researchgate.net/>



<https://commons.wikimedia.org/>

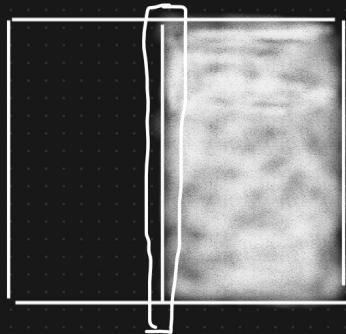


### ④ Convolution Operation In CNN

→ (0,1)

0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255

→



6x6x1

Convolution operation

Step 1

① Normalize

Divide by 255

+1	+2	-1
0	0	0
-1	-2	-1

Stride=1

0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1

$n=6$

\*

$f=3$

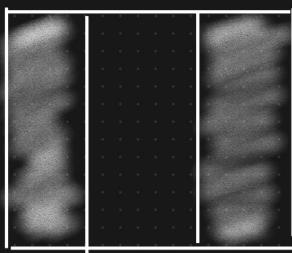
+1	0	-1
+2	0	-2
+1	0	-1

=

0	-4	-4	0
0	-4	-4	0
0	-4	-4	0
0	-4	-4	0

$4 \times 4$

$6 \times 6 \times 1$



filters  
vertical edge filters

$\leftarrow$

arr	0	0	arr
arr	0	0	arr
arr	0	0	arr
arr	0	0	arr

$$\begin{aligned} h - f + 1 &= \\ &= 6 - 3 + 1 = 4 \end{aligned}$$

## ⑤ Padding In CNN

0	0	0	0	1	1	1	1
0	0	0	0	1	1	1	1
0	0	0	0	1	1	1	1
0	0	0	0	1	1	1	1
0	0	0	0	1	1	1	1
0	0	0	0	1	1	1	1
0	0	0	0	1	1	1	1
0							

$6 \times 6$

$n=6$

$8 \times 8$

① Zero padding

② Neighbour padding

$f=3$

+1	0	-1
+2	0	-2
+1	0	-1

0	-4	-4	0	
0	-4	-4	0	
0	-4	-4	0	
0	-4	-4	0	

$6 \times 6$   
=

$n-f+2p+1$

$6 - 3 + 2p + 1 = 6$

$3 + 2p + 1 = 6$

$2p = 6 - 4$

$p = \frac{2}{2} = 1$

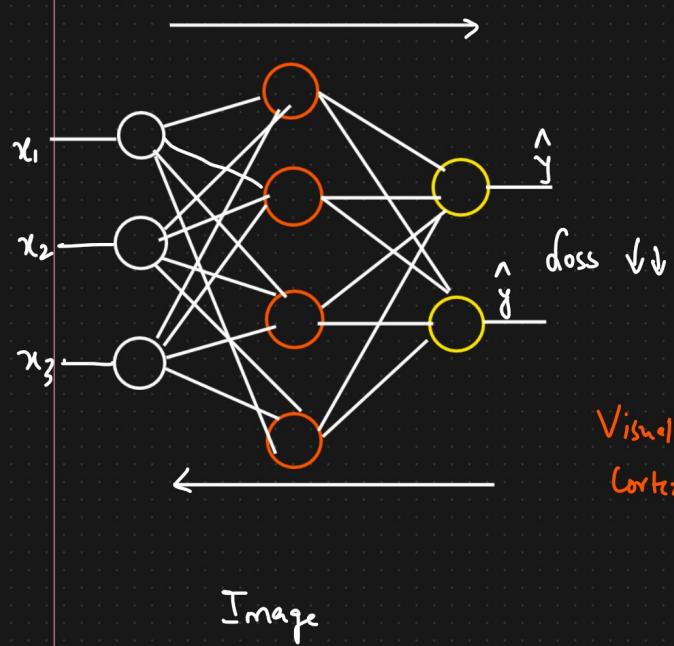
$7 \times 7$

$3 \times 3$

$7 \times 7$

How much padding you need to apply?

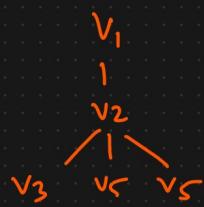
## ⑥ Operation of CNN vs ANN



$$z = w^T x_i + b$$

$$\text{relu}(z)$$

Visual  
Cortex



→ ReLU operation  $\max(0, x)$

0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1

\*

+1	0	-1
+2	0	-2
+1	0	-1


$f_1$

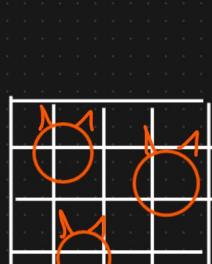
0	-4	-4	0
0	-4	-4	0
0	-4	-4	0
0	-4	-4	0

-	-	-	-
-	-	-	-
-	-	-	-
-	-	-	-

$f_2$   
 $f_3$   
⋮  
 $f_n$

Convolution Layer

## ⑦ Max Pooling, Min Pooling, Mean Pooling



$2 \times 2$

$P=1$   
 $S=1$

Convolution Layer


ReLU

$f_1$   
 $f_2$   
 $f_3$   
 $f_4$

1	
2	
3	

4	6
8	4

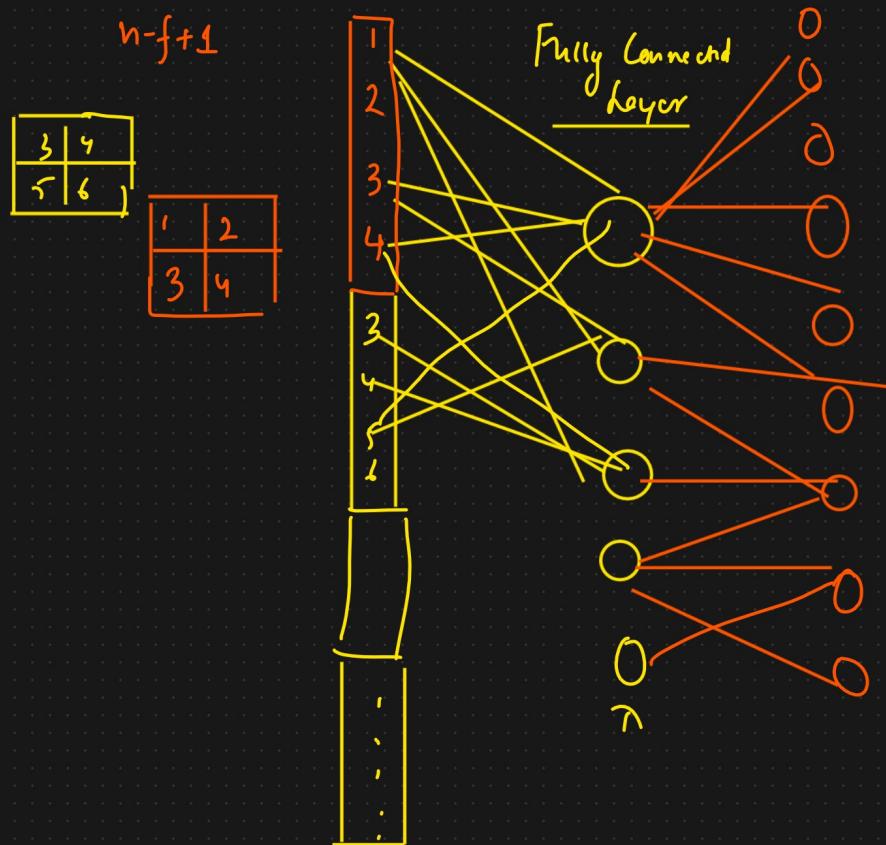
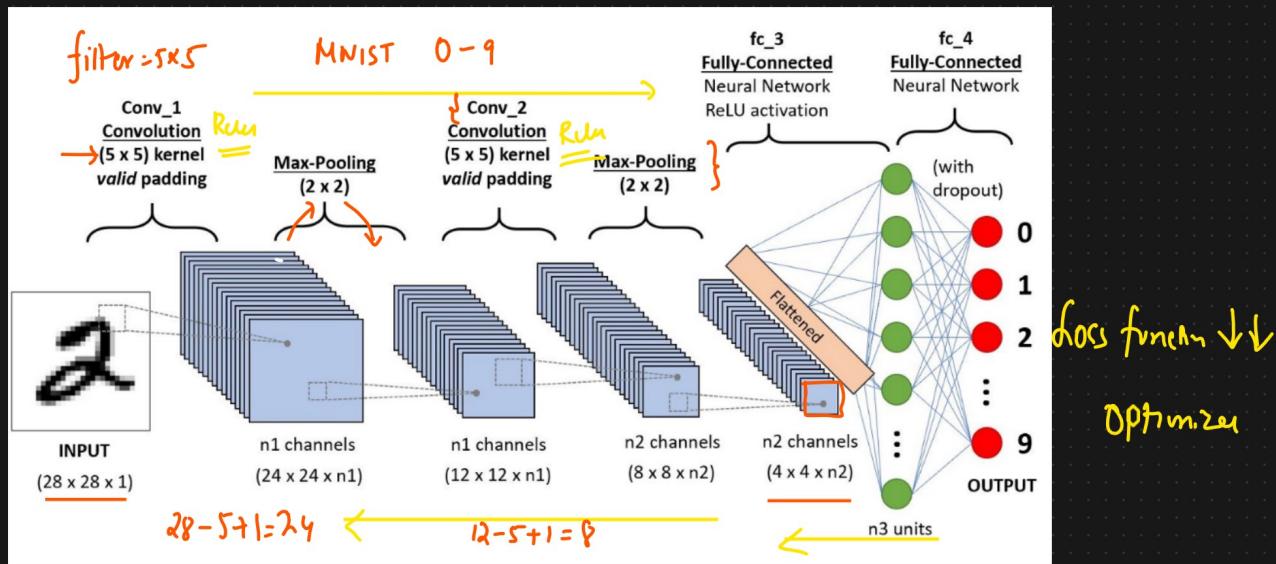
Max Pooling

1	
2	
3	

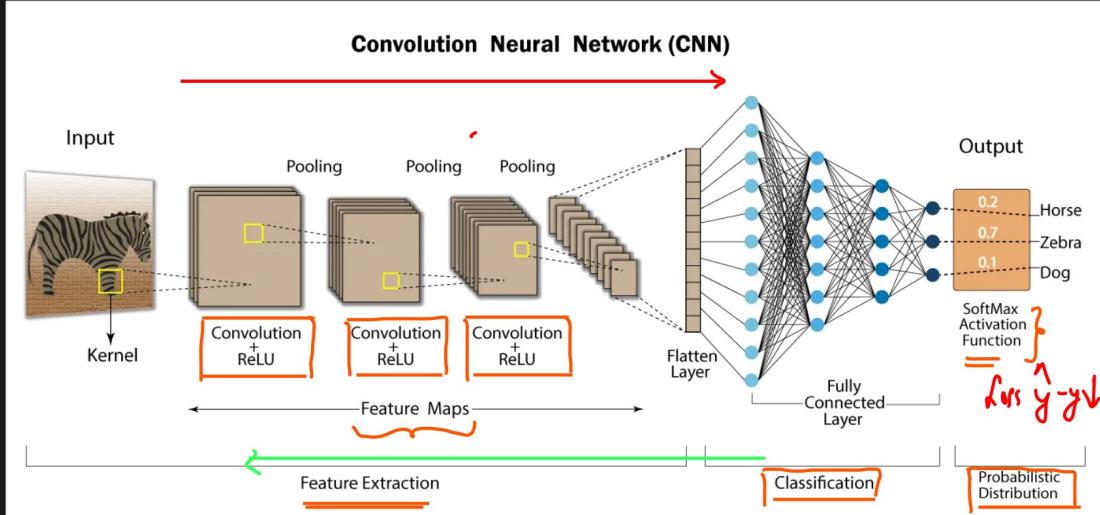
4	6
8	4

## ⑦ Location Invariant

## ⑧ Fully Connected layer In CNN [Flattened Layer]



# ⑨ CNN Complete Example



<https://developersbreach.com/convolution-neural-network-deep-learning/>

0	0	0	
0	0	0	
0	0	0	
0	0	0	
0	0	0	
0	0	0	

\*

+1	0	-1
+2	0	-2
+1	0	-1

0	-4	-4	0
0	-4	-4	0
0	-4	-4	0
0	-4	-4	0