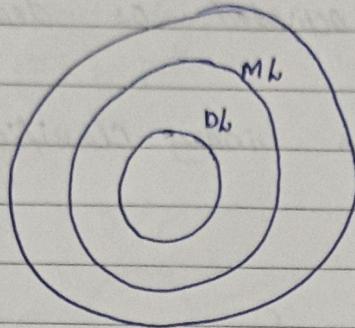


Deep learning

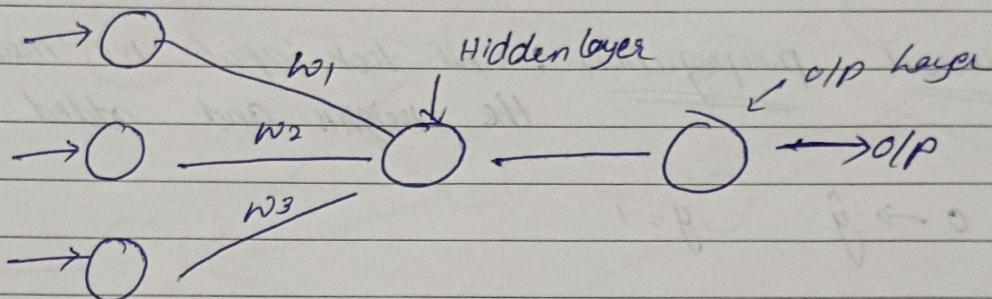


$AI \Rightarrow$ Application which can do its own task without any human intervention

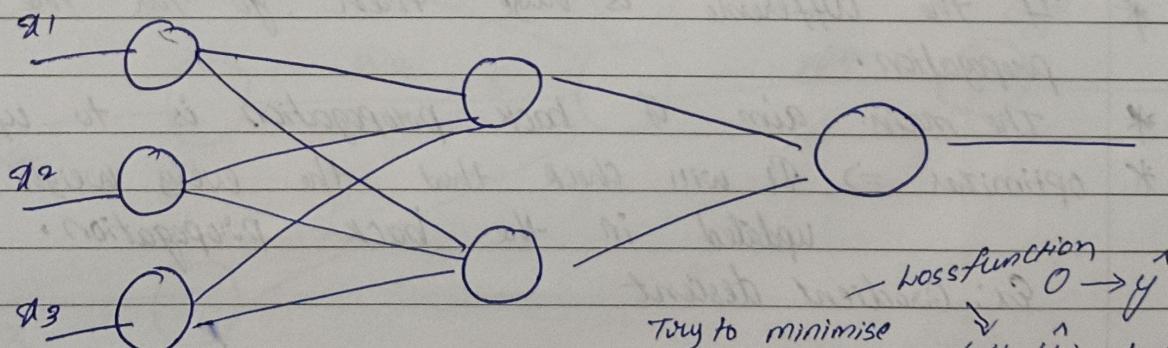
$ML \Rightarrow$ Statistical tool to analyze the data, visualize the data, predictions, forecasting

Perception
IIP layers

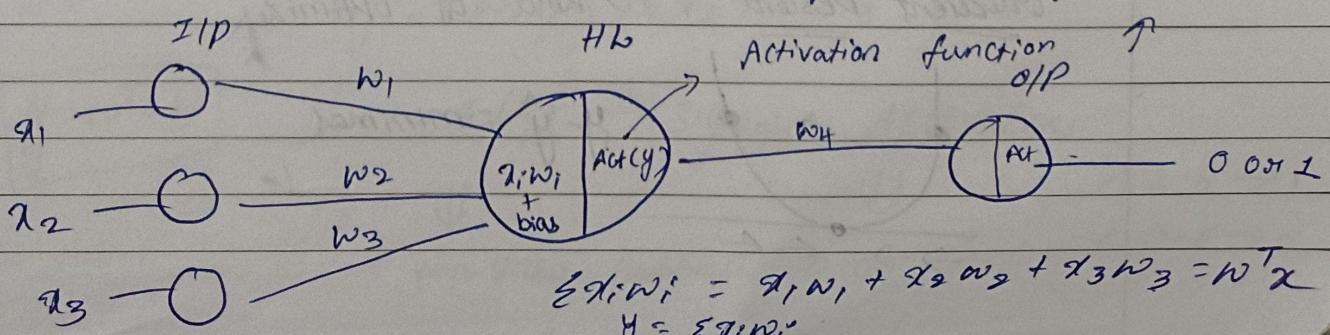
$\{$ Single Layered Neural Network $\}$



$x_i =$	Study (hr)	Play (hr)	Sleep (hr)	Pass/Fail
	7	3	7	1
	2	5	8	0
	4	3	7	1



Try to minimise the difference $(y - \hat{y})^2$; $y = 1$



Weight \rightarrow how much the neuron gets activated or deactivated
 Activation function \rightarrow

Sigmoid $\Rightarrow \frac{1}{1+e^{-y}}$ \Rightarrow used for binary classification

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

↓ {
 0 to 1 }
 0.5 \rightarrow 1
 -0.5 \rightarrow 0

Forward propagation: we took input we multiplied with the weights and added a bias

$$0 \rightarrow \hat{y} \quad y = 1$$

$$(y - \hat{y}) = 1$$

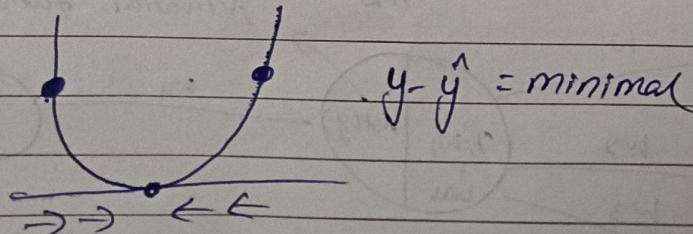
↓ This difference should be near to 0

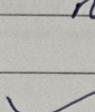
loss function \Rightarrow difference b/w predicted and real value
 we should try to minimize the difference
 and it always nearly to 0

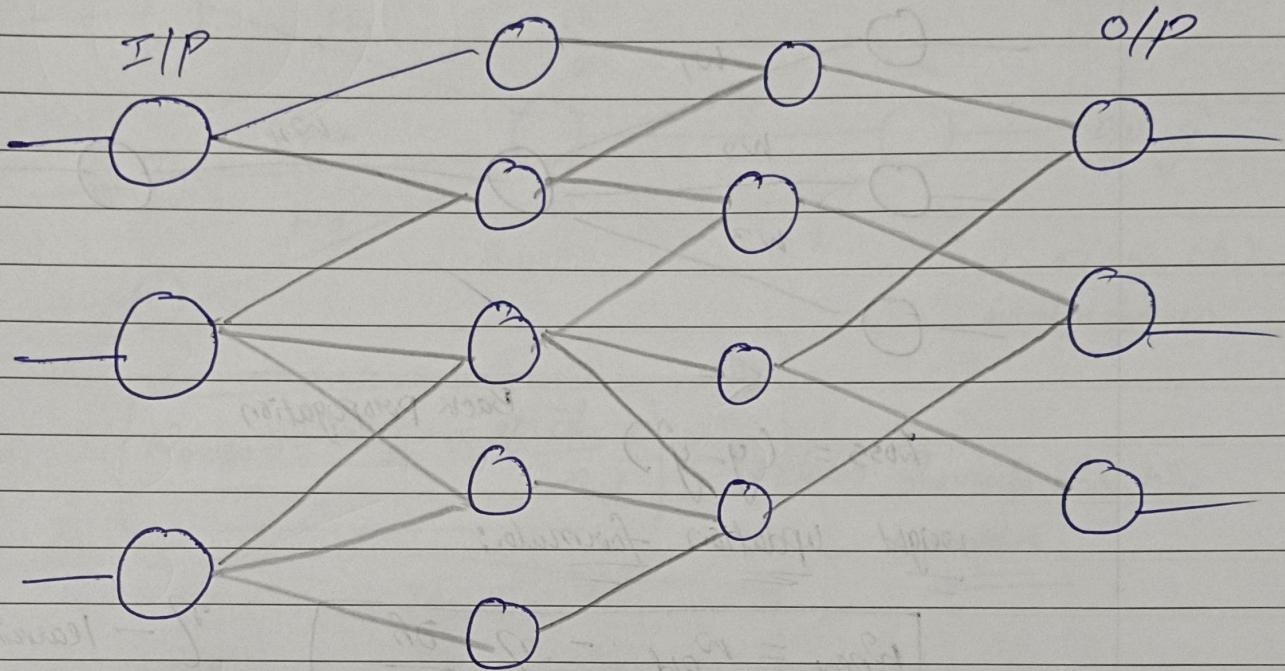
- * If the difference is huge then go for the back propagation.
- * The main aim of back propagation is to update weights
- * Optimizer \Rightarrow It will check that every weight should be updated in the back propagation.

Ex: Gradient descent

Gradient Descent \rightarrow kind of optimizer

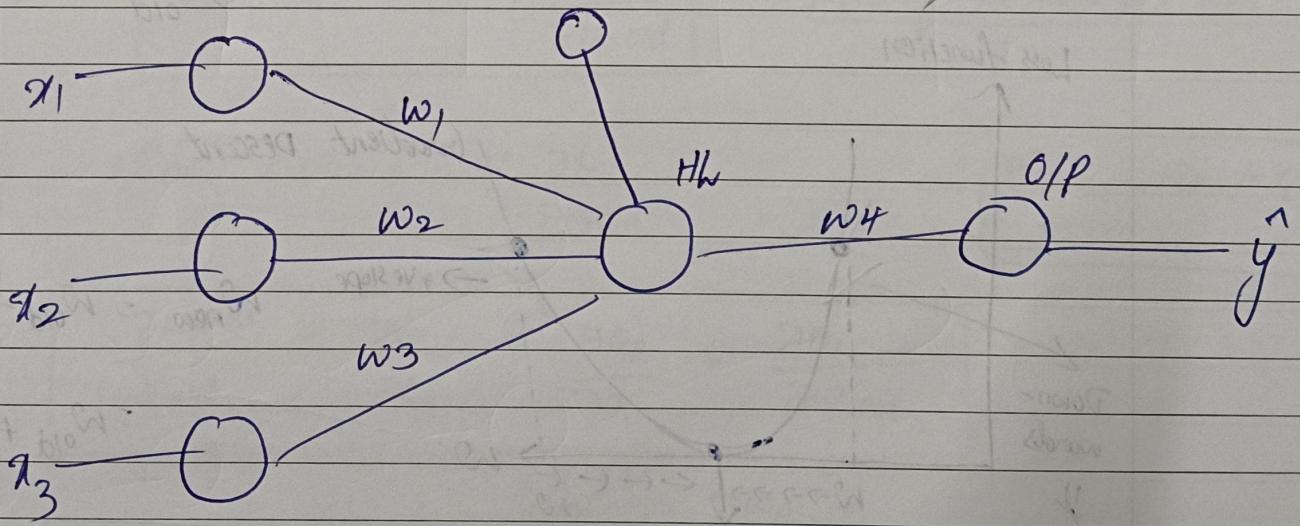


~~Multi layered~~ neural networks = ~~networks~~  Hidden layers



~~DAY - 2~~

Forward propagation bias →



$$y = (w_1 x_1 + w_2 x_2 + w_3 x_3) + \text{bias}$$

$$y = w^T + b$$

$w = \sigma(y) \rightarrow$ Sigmoid Activations



Non-linear problems

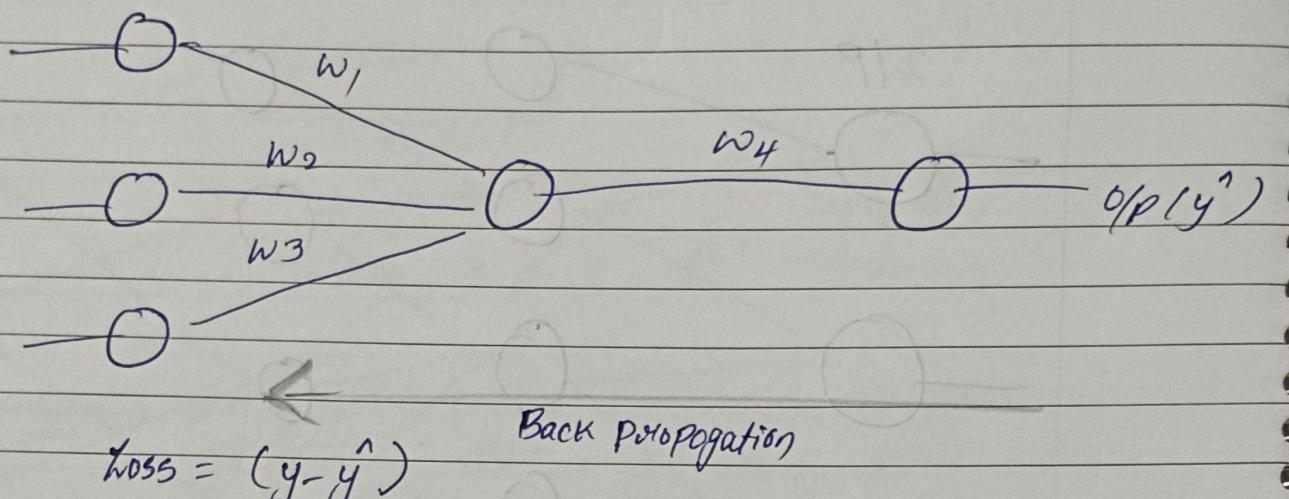


0 or 1

$> 0.5 \rightarrow 1$

$< 0.5 \rightarrow 0$

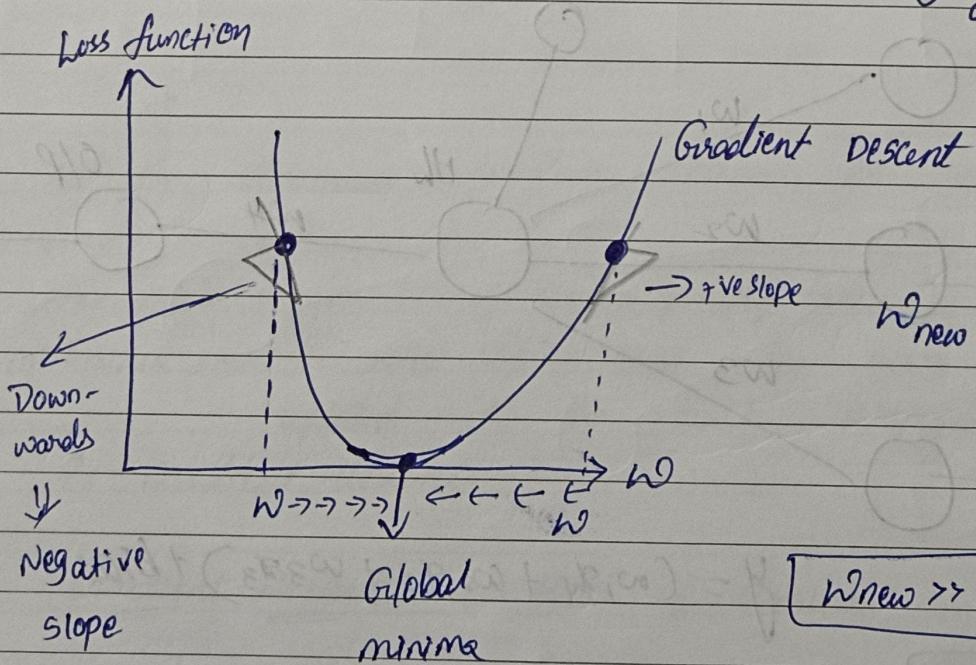
Back propagation



Weight update formula:

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

η - learning rate
 $\frac{\partial h}{\partial w_{\text{old}}}$ - slope



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

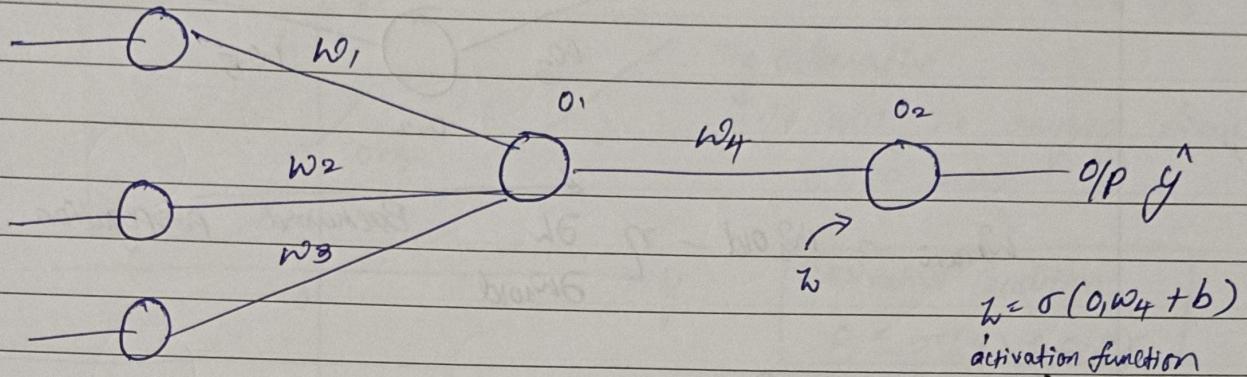
$$= w_{\text{old}} + \eta (\text{ve})$$

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

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

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

Chain Rule of Differentiation



$$w_{4\text{new}} = w_{4\text{old}} - \eta \frac{\partial L}{\partial w_{4\text{old}}}$$

→ updating formula for w_4 weights

$$\frac{\partial L}{\partial w_{4\text{old}}} = \frac{\partial L}{\partial O_2} * \frac{\partial O_2}{\partial w_{4\text{old}}}$$

$$b_{2\text{new}} = b_{2\text{old}} - \eta \frac{\partial L}{\partial b_{2\text{old}}}$$

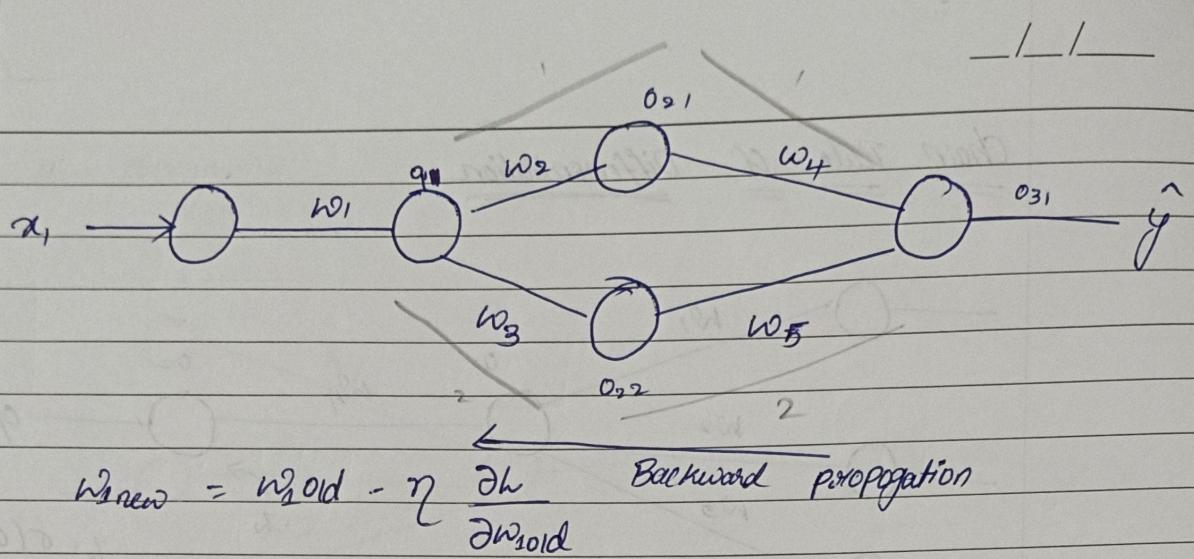
⇒ updating formula for bias

$$w_{1\text{new}} = w_{1\text{old}} - \eta \frac{\partial L}{\partial w_{1\text{old}}}$$

$$w_{2\text{new}} = w_{2\text{old}} - \eta \frac{\partial L}{\partial w_{2\text{old}}}$$

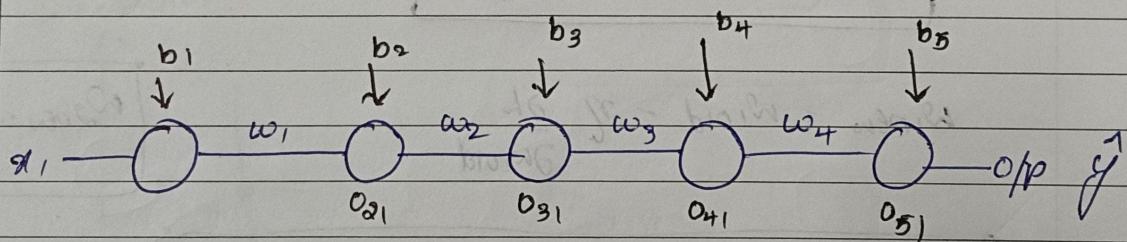
$$\frac{\partial L_{\text{loss}}}{\partial w_{1\text{old}}} = \frac{\partial L}{\partial O_2} * \frac{\partial O_2}{\partial O_1} * \frac{\partial O_1}{\partial w_{1\text{old}}}$$

$$\frac{\partial O_2}{\partial w_4} * \frac{\partial w_4}{\partial O_1}$$



$$\begin{aligned} \frac{\partial L}{\partial w_{1\text{old}}} &= \left[\frac{\partial L}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial o_{11}} * \frac{\partial o_{11}}{\partial w_{1\text{old}}} \right] + \\ &\quad \left[\frac{\partial L}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{22}} * \frac{\partial o_{22}}{\partial o_{11}} * \frac{\partial o_{11}}{\partial w_{1\text{old}}} \right] \end{aligned}$$

③ Vanishing Gradient problem



$$loss = \frac{1}{2} (y - \hat{y})^2 \rightarrow MSE$$

$$w_{1\text{new}} = w_{1\text{old}} - \eta \frac{\partial L}{\partial w_{1\text{new}}}$$

$$\frac{\partial L}{\partial w_{1\text{new}}} = \frac{\partial L}{\partial o_{51}} * \frac{\partial o_{51}}{\partial o_{41}} * \frac{\partial o_{41}}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial w_{1\text{new}}}$$

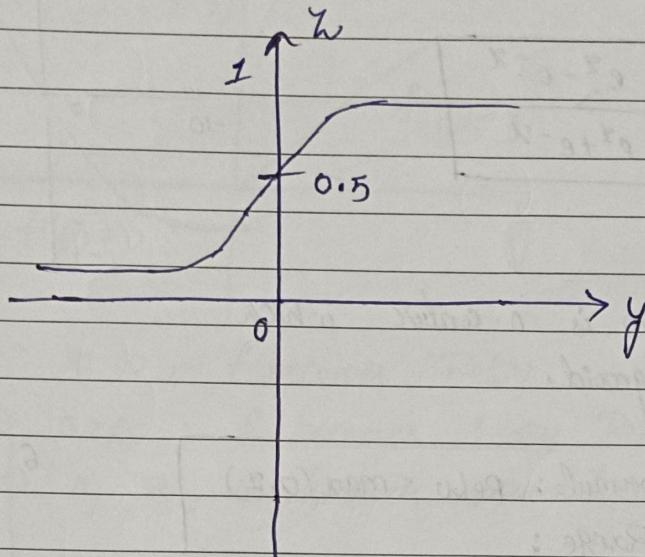
Sigmoid Activation formula

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

$\rightarrow 0 \text{ or } 1$

$y > 0.5 \rightarrow 1$

$y < 0.5 \rightarrow 0$



The derivative

↓ It will be changing from
0 to 0.25

Derivative condition:
 $0 \leq \sigma(y) \leq 0.25$

$$0.51 = \sigma[(0_{41} * w_4) + b]$$

↓
sigmoid Activation

$$w_{\text{new}} = w_{\text{old}} - \eta^{\text{small}} \text{ (small number)}$$

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



Vanishing gradient problem

Activation functions:

①

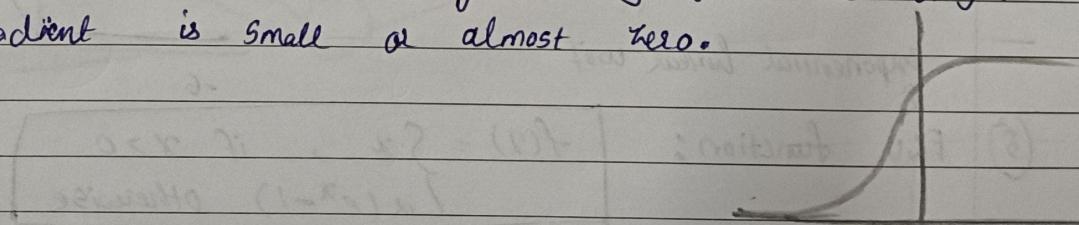
Sigmoid function

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

Range: 0 to 1

⇒ The function output is not centered on 0, which will reduce the efficiency of weight update.

⇒ when the input is slightly away from the origin then the gradient is small or almost zero.

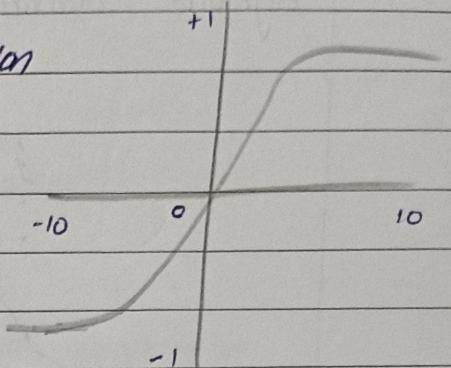


(2) Tanh function :

\Rightarrow Also called hyperbolic tangent function

$$\text{Formula : } \tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

Range : -1 to 1



\Rightarrow The whole function is 0-centric which is better than sigmoid.

(3) ReLU function :

$$\text{formula : ReLU} = \max(0, x)$$

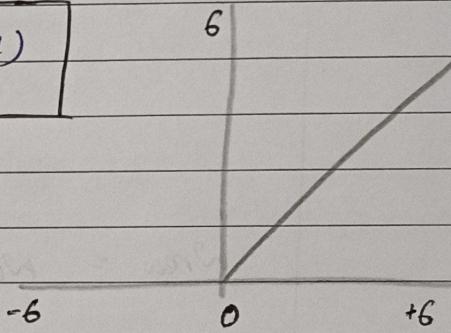
Range :

\Rightarrow Not 0-centric

\Rightarrow Better than sigmoid & tanh

\Rightarrow calculation speed is much faster

\Rightarrow when input is -ve, ReLU completely inactive which means that once a -ve number is entered ReLU will die

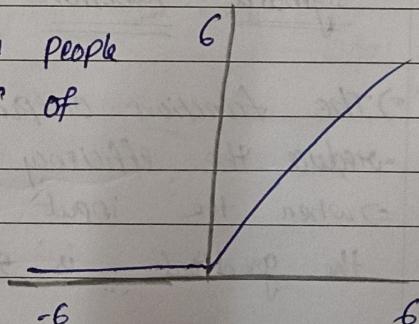


(4) leaky ReLU function :

$$f(x) = \max(0.01x, x)$$

\Rightarrow To solve the Dead ReLU problem people

proposed to set the first half of ReLU 0.01x instead of 0.



exponential linear unit

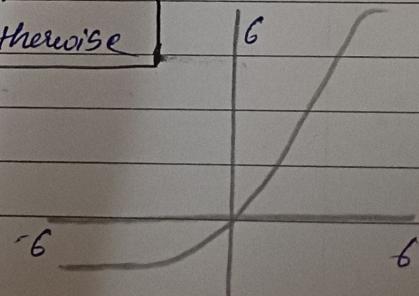
(5) ELU function :

$$f(x) = \begin{cases} x, & \text{if } x > 0 \\ \alpha(e^x - 1), & \text{otherwise} \end{cases}$$

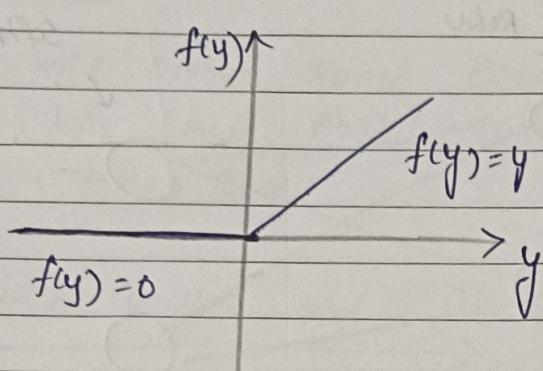
\Rightarrow Slightly more computationally intensive

\Rightarrow No Dead ReLU issues

\Rightarrow The mean of the output is close to 0, zero-centred



(6) PReLU function:



$$f(y_i) = \begin{cases} y_i & \text{if } y_i > 0 \\ a_i y_i & \text{if } y_i \leq 0 \end{cases}$$

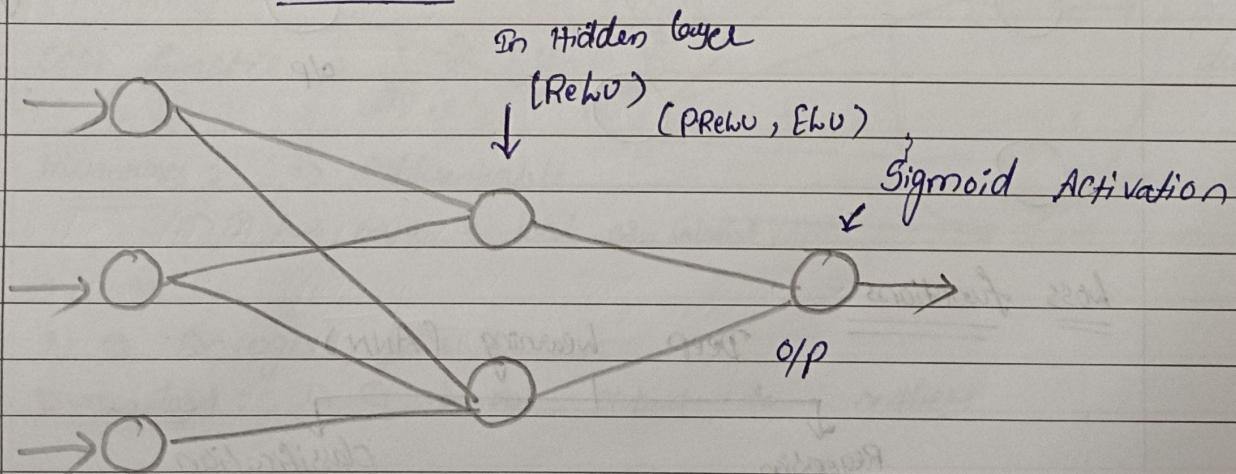
$\Rightarrow a=0$, f becomes ReLU

$\Rightarrow a>0$, f becomes leaky ReLU

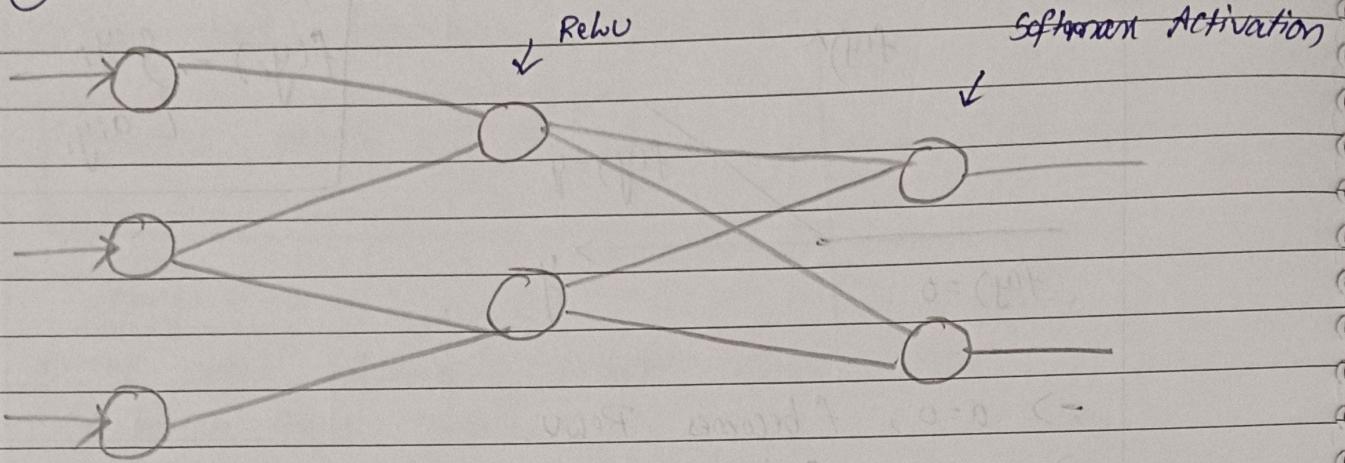
$\Rightarrow a$ is a learnable parameter, f becomes PReLU

(7) Swish (A Self-Gated) function:

Technique which activation function we should use

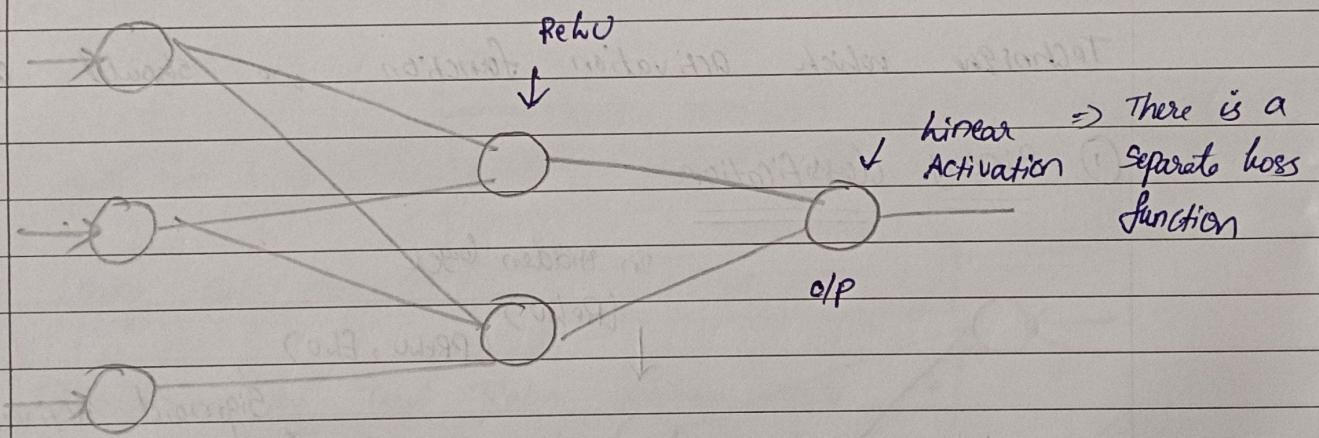
(1) Binary classification

② Multi class classification



* In Hidden layers use ReLU maximum it solves

③ Regression



loss functions

Deep learning (ANN)

Regression			Classification		
Experience	Degree	salary (O/P)	Play	Study	Pass/Fail
10	PhD	-	10	2	F
-	-	-	4	3	F
-	-	-	5	5	P
-	-	-	2	7	P

Regression Problem

Continuous values

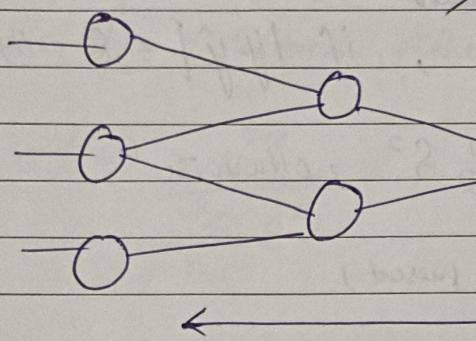
① Regression

(i) MSE (mean squared error)

(ii) MAE (mean absolute error)

(iii) Huber loss

loss function and cost function



Dataset = 100 records

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

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

(Provide batch of data points)

(i) MSE

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

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

Advantages: 1) Differentiable

2) It has only 1 local or Global minima

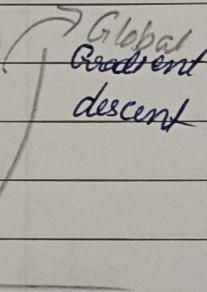
3) It converges faster

Disadvantages: 1) It is not robust to outliers

(ii) MAE

$$\text{Loss function} = \frac{1}{2} |y - \hat{y}|$$

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



Advantages: 1) Robust to outliers

⇒ Subgradient

⇒ Time consuming

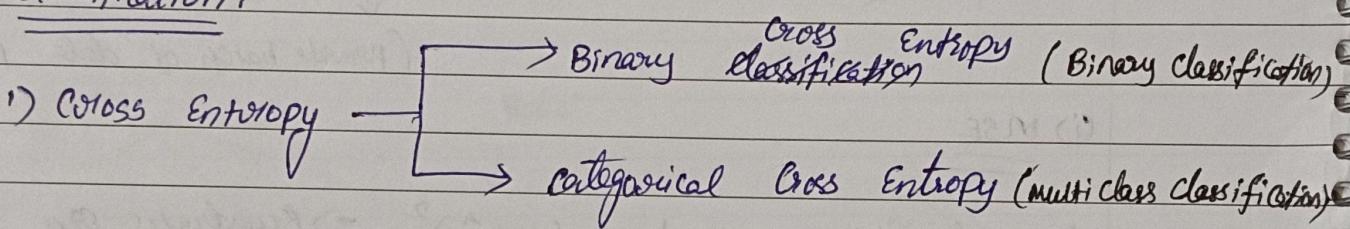
(3) Huber loss : Combination of MSE and MAE

Outliers not present

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

↓
(outliers are present)

Classification:-



(i) Binary Cross Entropy :

$$\text{loss} = -y * \log(\hat{y}) - (1-y) * \log(1-\hat{y}) \quad \text{Logistic Regression}$$

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

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

(ii) Categorical Cross Entropy : (multi-class classification)

f_1	f_2	f_3	O/P	Good	Bad	Neutral
2	3	4	Good	1	0	0
5	6	7	Bad	0	1	0
8	9	10	Neutral	0	0	1

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

i = Row

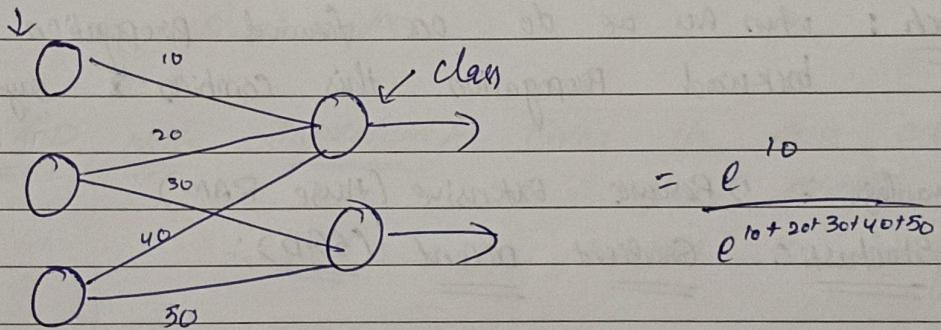
$$y_i = [y_{i1}, y_{i2}, y_{i3}, \dots, y_{ic}] \rightarrow \text{category}$$

j = column

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

\hat{y}_{ij} = Softmax Activation \rightarrow o/p layer

$$\sigma(z) = \frac{e^{z_j}}{\sum_{j=1}^k e^{z_j}}$$



Conclusion:-

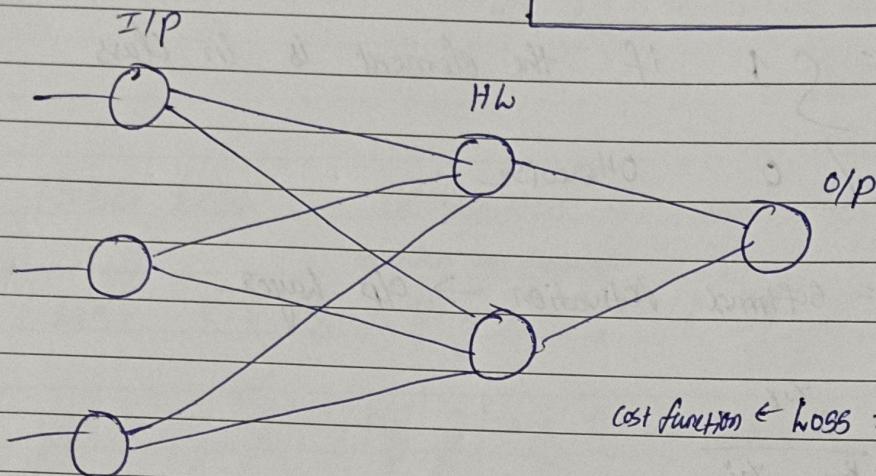
- ① ReLU, Softmax \Rightarrow Multiclass \rightarrow Categorical Cross Entropy
- ② ReLU, Sigmoid \Rightarrow Binary \rightarrow Binary Cross Entropy

ReLU, Linear Activation \rightarrow MSE, MAE, Huber loss

Day-3

① Gradient Descent :-

weight update formula : $w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}}$

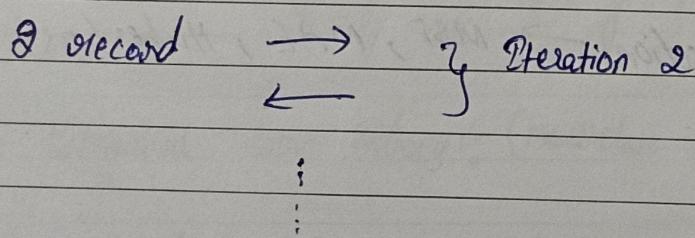
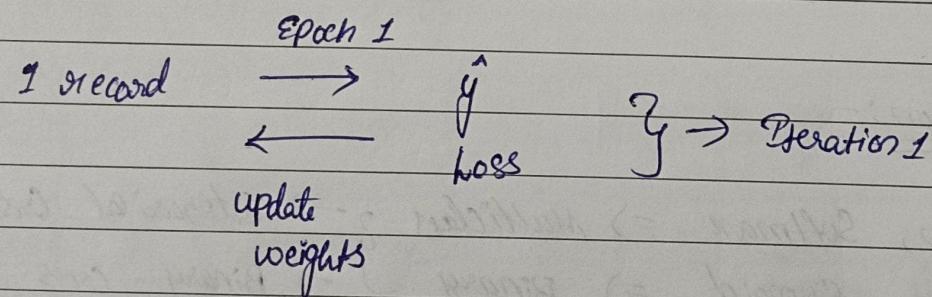


$$\text{cost function} \leftarrow \text{loss} = \frac{1}{2n} \sum_{i=1}^n (y - \hat{y})^2$$

Epoch : whenever we do one forward propagation and one backward propagation this combined says one epoch.

Disadvantage : 1) Resource extensive (huge RAM)

② Stochastic Gradient descent (SGD) :-



Disadvantage :-

1) Convergence will be very slow

2 million Iteration

2) Time complexity will also be high

1/1

③ Mini batch SGD :-

1000000 batch size = 1000

epoch - 1

1000 →

Iteration 1

i) Advantages :-

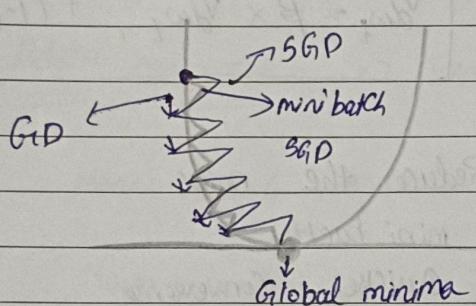
- (i) Resource intensive reduced
- (ii) Convergence will be better
- (iii) Time complexity will be improved

→

Iteration 2

↓

Iteration 1000



SGD having highest noise

④ SGD with Momentum :- Reduce the noise and comes to the global minima smoothen the journey.

Exponential Weighted Average

↓

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

Time series

↓

ARIMA, ARMA

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

Exponential weighted average

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

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

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

Forecasting

$$v_{t_1} = a_1$$

$\beta \Rightarrow$ Hyperparameter

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

↓ in which value we should focus more

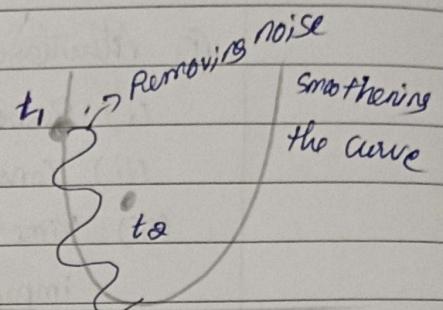
$$\beta = 0 \text{ to } 1$$

↓

0.95

$$v_{t_2} = (0.95) * v_{t_1} + (0.05) * a_2$$

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



Exponential weighted Average

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

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

⇒ { Reduce the noise
mini batch
faster convergence }

③ Adagrad → Adaptive Gradient Descent

η = learning rate
helps do

maintain
the speed
of convergence

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

η = fixed \rightarrow adaptive

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

$$\eta' = \frac{\eta}{\sqrt{\alpha_t + \epsilon}}$$

ϵ = small number (denominator never become zero)

$$\alpha_t = \sum_{i=1}^t \left(\frac{\partial h}{\partial w_i} \right)^2$$

t = current timestamp

η' decreases as we are reaching the global minima

AdaDelta and RMSprop

$$\eta' = \frac{\eta}{\sqrt{sdw} + \epsilon}$$

$$sdw = 0$$

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

$$\beta = 0.95$$

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

Adam optimizer \Rightarrow Best optimizer

Momentum + RMSprop (Adaptive learning Rate)

$$vdw = 0 \quad vdb = 0 \quad sdw = 0 \quad sdb = 0$$

$$\text{For weight} \Rightarrow w_t = w_{t-1} - \eta' vdw$$

$$\text{For bias} \Rightarrow b_t = b_{t-1} - \eta' vdb$$

$$\eta' = \frac{\eta}{\sqrt{sdw} + \epsilon}$$

$$vdw_t = \beta \times vdw_{t-1} + (1-\beta) \frac{\partial L}{\partial w_{t-1}}$$

① Smoothening

② learning rate adaptive

$$vdb_t = \beta \times vdb_{t-1} + (1-\beta) \frac{\partial L}{\partial b_{t-1}}$$

ANN (Artificial Neural Networks)

① For which all algorithms feature scaling is required?

ANN \Rightarrow Any algorithm that relates to some distance based problems

LR \Rightarrow Gradient Descent

KNN

KMeans

Tensorflow \Rightarrow popular library \Rightarrow Google
Tensorflow - wrapper Keras > 2.0 \Rightarrow They both got integrated

Dense \Rightarrow we can create I/P, O/P, H/W
Sequential =)
Dropout =)

Black Box Model vs White Box Model

Random Forest \rightarrow Black Box model ANN, xgboost
Decision Tree \rightarrow White Box Model linear regression

CNN (Convolutional Neural Networks) \Rightarrow Image, video frame

CNN

Visual cortex

\Rightarrow convolution

$v_1 \rightarrow$ moving things



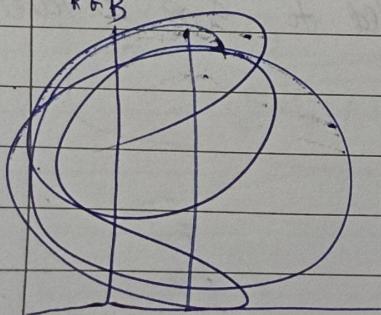
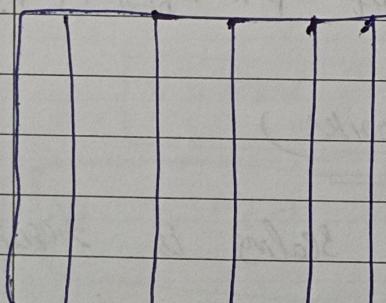
Images =)

$v_2 \rightarrow$ Animals (Cat)

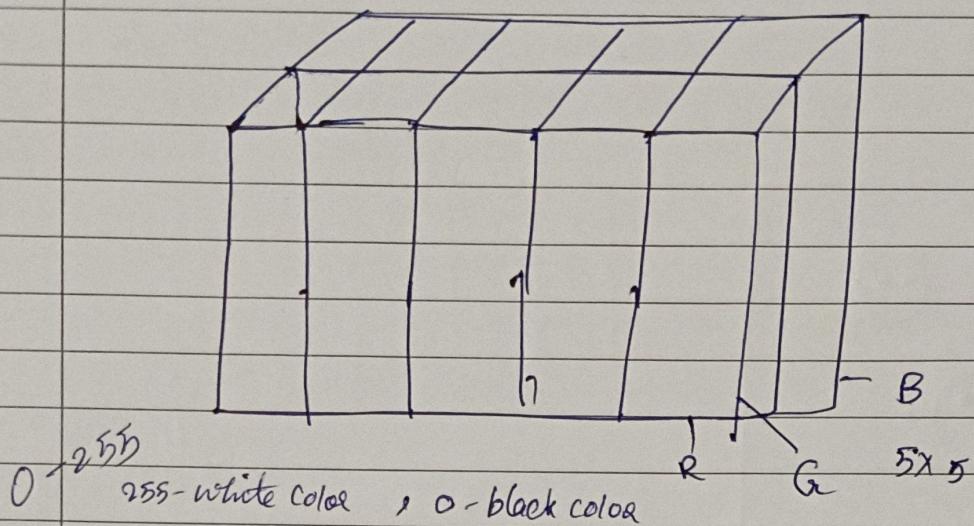
Black & white



$v_3 \rightarrow$ Map the environment



RGB



Convolution

Image

filter / kernel

o/p

$$Stride = 1$$

min max

Gx G

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

3x3

Horizontal
edge filter

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

4×4

Channing

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

filter/kernel

padding => Build the Compound Around the image

$$N - fH = 6 - 3H = 3H = 4$$

$$\text{After Padding} = 8 - 3 + 1 = 6$$

11

padding filter

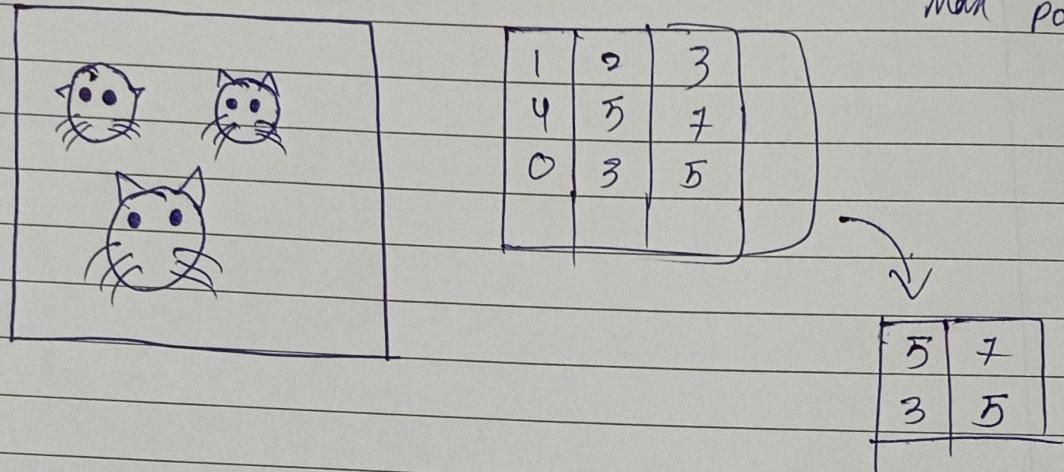
$$n=6, p=1, f=3 \Rightarrow n \geq p - f + 1 = 6 \geq 2 - 3 + 1 = 6$$

Max Pooling

Avg Pooling

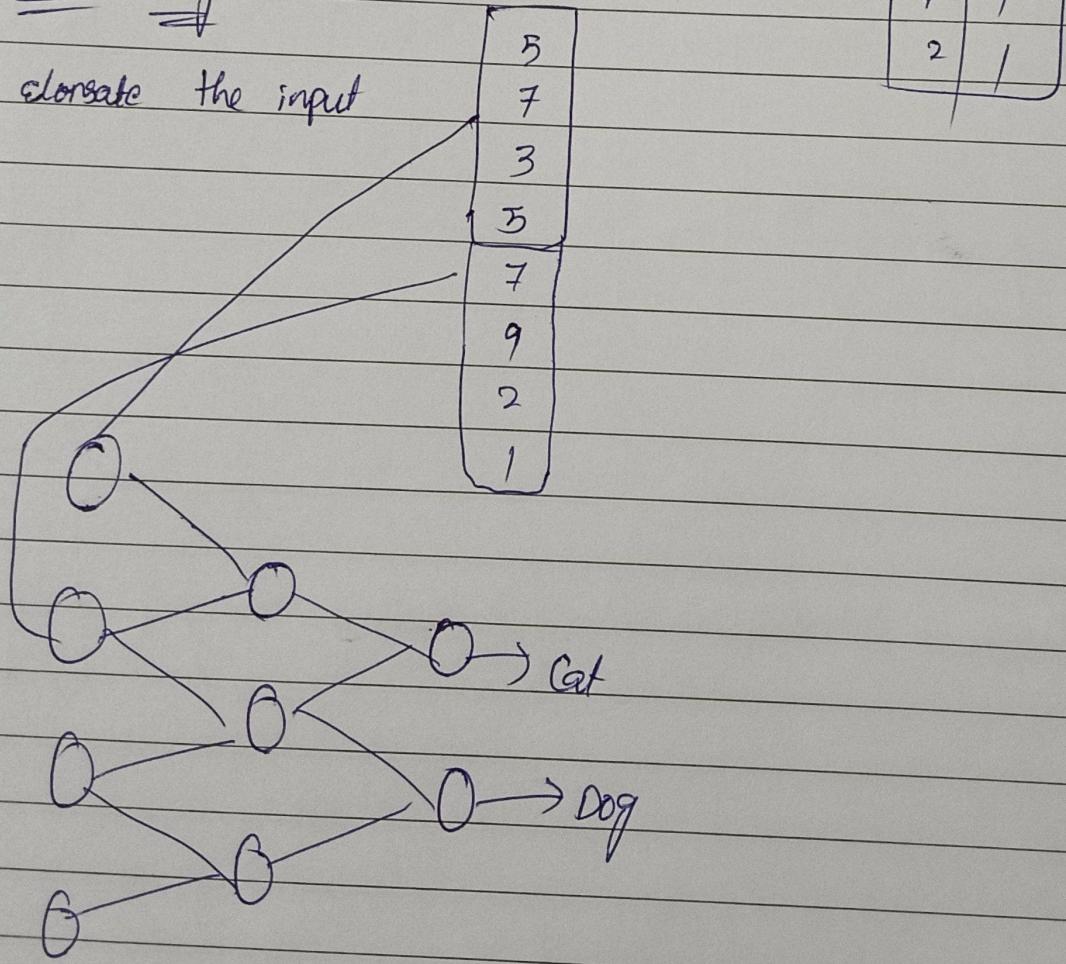
Min Pooling

Max Pooling



Flattening Layer

Flatten the input



— / —

