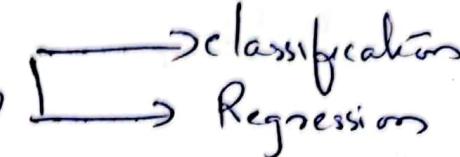
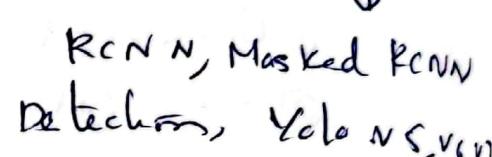
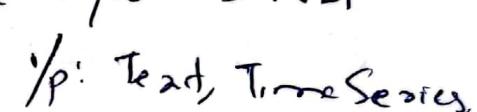
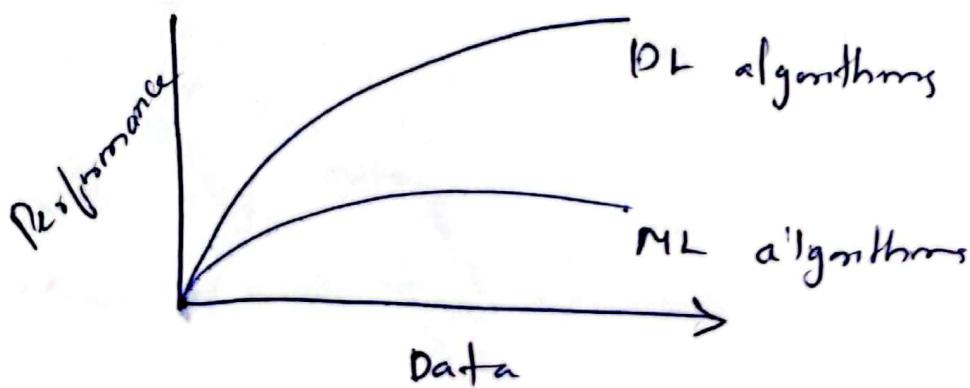


DEEP LEARNING FUNDAMENTALS

Deep Learning

- 1) ANN → Artificial Neural N/w 
 - classification
 - Regression
- 2) CNN → Convolutional Neural N/w → I/p → Images/Video frames

- 3) RNN → Recurrent Neural N/w → NLP

 - ↓
 - { Word Embedding, LSTM RNN,
GRU RNN, Bidirectional
LSTM RNN, Encoders,
Decoders, Transformers,
BERT }

Framework : TensorFlow



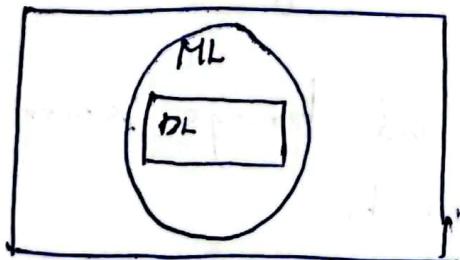
Deep Learning

Day 1

Agenda:

- 1) Deep learning: Perception. {AI v/s ML v/s DL v/s DS}
- 2) Forward Propagation
- 3) Backward propagation
- 4) Loss function
- 5) Activation function
- 6) Optimizers

AI : Application which can do its own task without any human intervention



- eg:
- 1) Netflix App
 - 2) Self driving Cars
 - 3) Amazon Applications
 - 4) Sofia

ML: provides stats tools to analyze the data, visualize the data, predictions, forecasting - clustering.

DL: Mimic Human brain

The final goal of Data scientist is to create an AI application

Q) Why DL becoming so popular?

1. In 2005, → ORCUT

After that, → Facebook, Instagram, WhatsApp,
LinkedIn, Twitter.

Because of this data started to get generated exponentially.

2. Concept of Bigdata was very high in demand.

3. Company had huge amount of data.

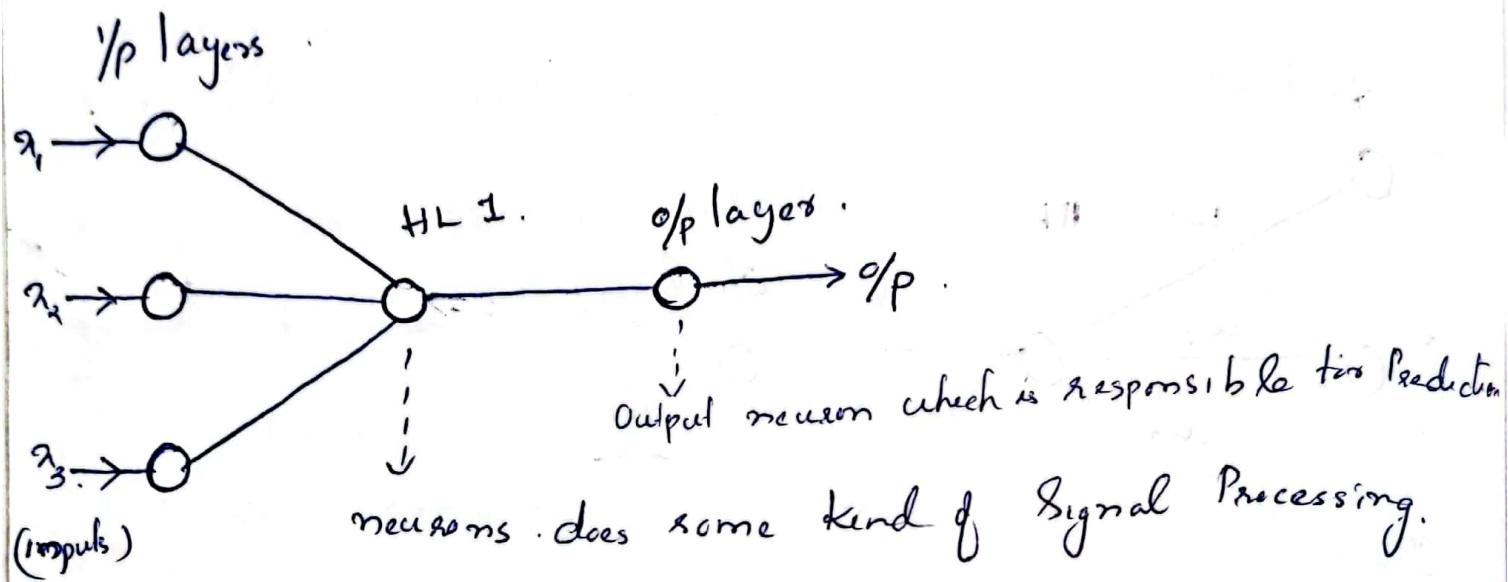
↓
use this data and brings seamless effect on their product.

↓

AI became so popular.

Q) Hardware Advancements (NVIDIA) → GPU's →
(Training the model)

Perceptrons : Single layer and Multilayer Neural N/w

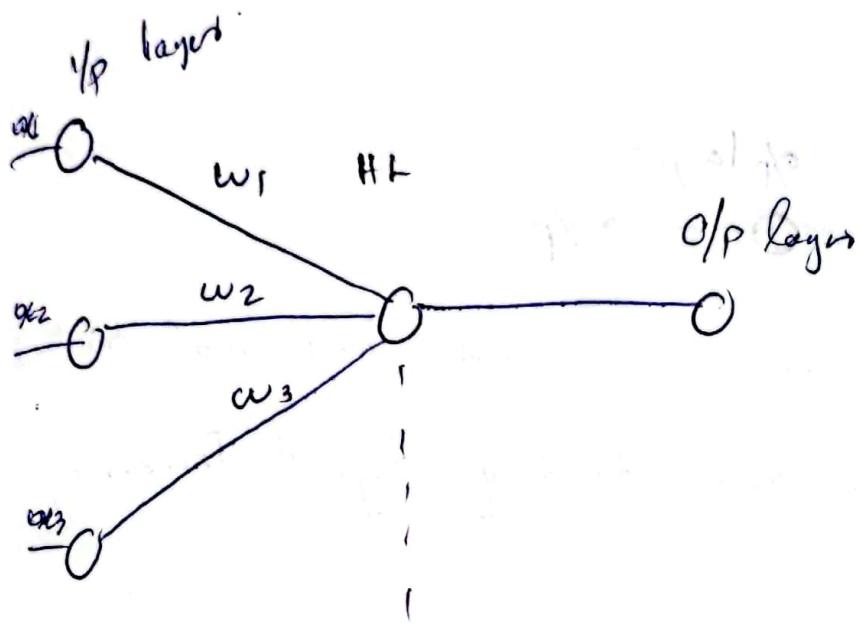


Dataset f

	x_1	x_2	x_3	Pass/Fail
Study	7	3	7	1
Play	2	5	8	0
Sleep	4	3	7	1

Whenever we create a neural n/w from scratch we need to train them based on the o/p.

When the signal is getting passed from the up layers, what is happening?



As soon as the signal goes to the H.L., assigning different weights.

$$\textcircled{1} \quad \sum_i x_i w_i = x_1 w_1 + x_2 w_2 + x_3 w_3 \\ = \mathbf{w}^T \mathbf{x}$$

\textcircled{2} Pass the above function to an activation function

Imp
Weights will actually help the neurons to act whether it should be activated or what level it should be activated

For weights, we initialize it as zero
then $\sum z_i w_i$ will be zero.

So Inorder to avoid that we introduce a parameter called bias.

as if by mistake we initializes the weight as zero, bias will have some value so that some value it will not equal to zero so we can continue the training and later can update the weights.

(Bias initialization will automatically happen by neural n/u) bias will required in every HL.

Then in step 3,
we pass an activation function on top of y where $y = \sum z_i w_i$.

Activation function

It is a function that calculates the o/p of a node based on the input it receives.

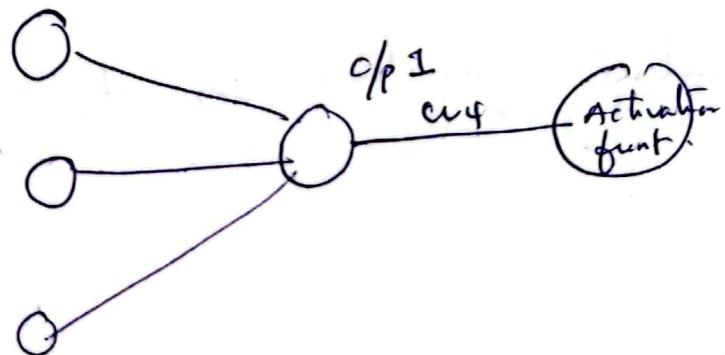
Sigmoid Activation function (used for binary classification)

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

Output of the equation will be zero or one.
Or it gives values b/w 0 and 1.

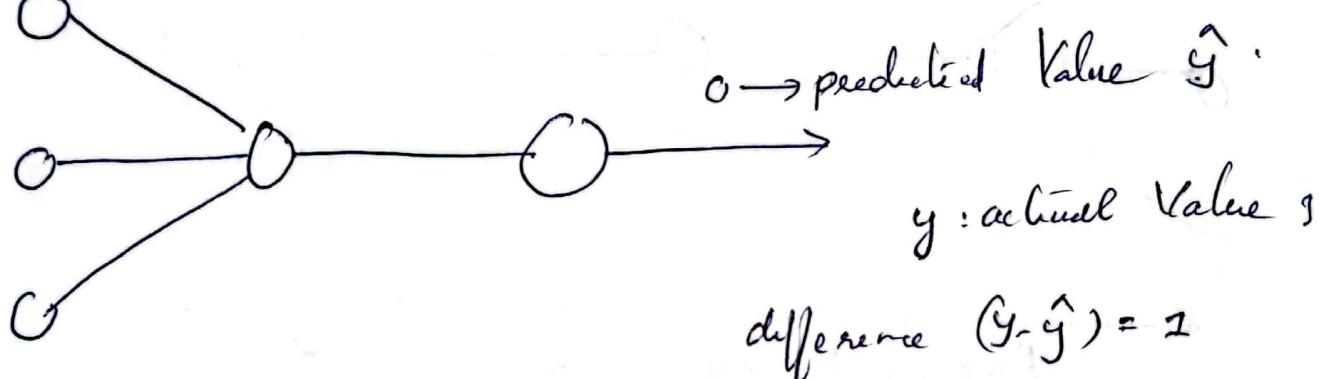
If ≥ 0.5 considered as 1
 < 0.5 considered as 0.

Then we go the next layers, another set of weights will be assigned and next activation function will be applied and finally will be able to get the o/p.



We basically took the i/p multiplied by weights, added a bias and activated.

So the entire process from the beginning till the end is called Forward Propagation.



The difference should be ~~st~~ near to zero

Loss function

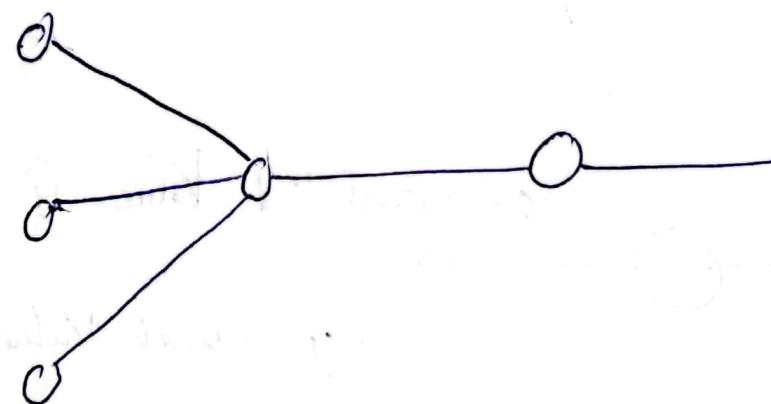
It finds out / measures the discrepancy b/w the predicted outcomes and the actual targets. The primary goal of a loss func is to guide the optimization process towards reducing the error.

Near to zero.

then we do an important process known as backward propagation.

The main aim of backward propagation is to update the weights.

forward Propagation



Backward propagation

An Optimizer.

forward Propagation

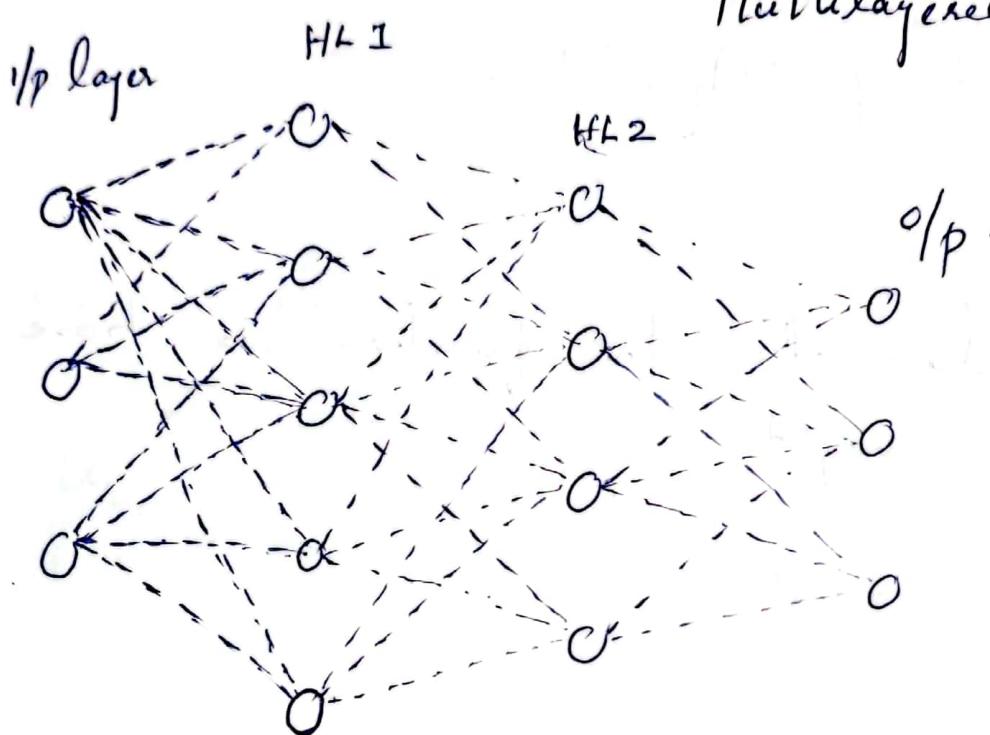
Input layer \rightarrow weights \rightarrow Bias \rightarrow Activation function

update the weights \leftarrow optimizers \leftarrow loss function
 $(\hat{y} - y)^2 / 2$

backward propagation

④ ANN, ENN, RNN, object detection will works on forward and backward propagation.

Multilayered Neural N/w



Neural N/w are for mimicing the human brain

Day 2

Agenda:

1. Forward Propagation
2. Chain rule of derivatives
3. Vanishing Gradient Problem
4. Loss functions

Backward propagation

1. Weight update formula
2. Chain rule of Differentiation

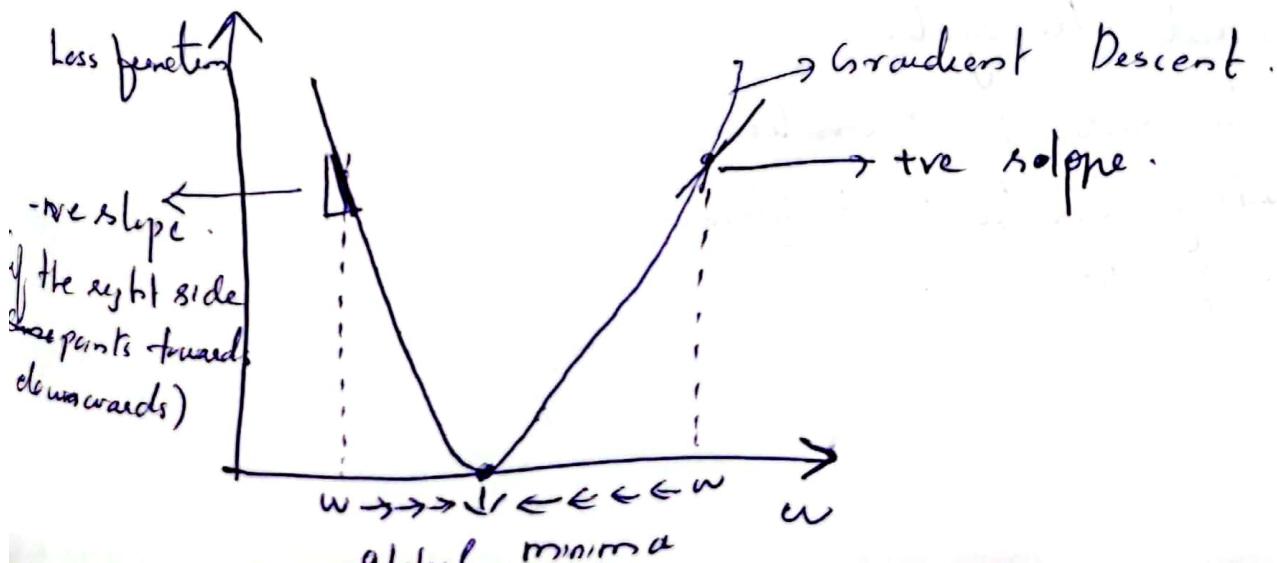
In order to reduce the loss function we have to update the weights.

Weight Update formula

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

$\eta \rightarrow$ learning rate.

$\frac{\partial L}{\partial w_{\text{old}}}$ \rightarrow slope.



$$\begin{aligned} w_{\text{new}} &= w_{\text{old}} - \eta (-\text{ve}) \\ &= w_{\text{old}} + \eta (+\text{ve}) \end{aligned}$$

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

(+ve case)

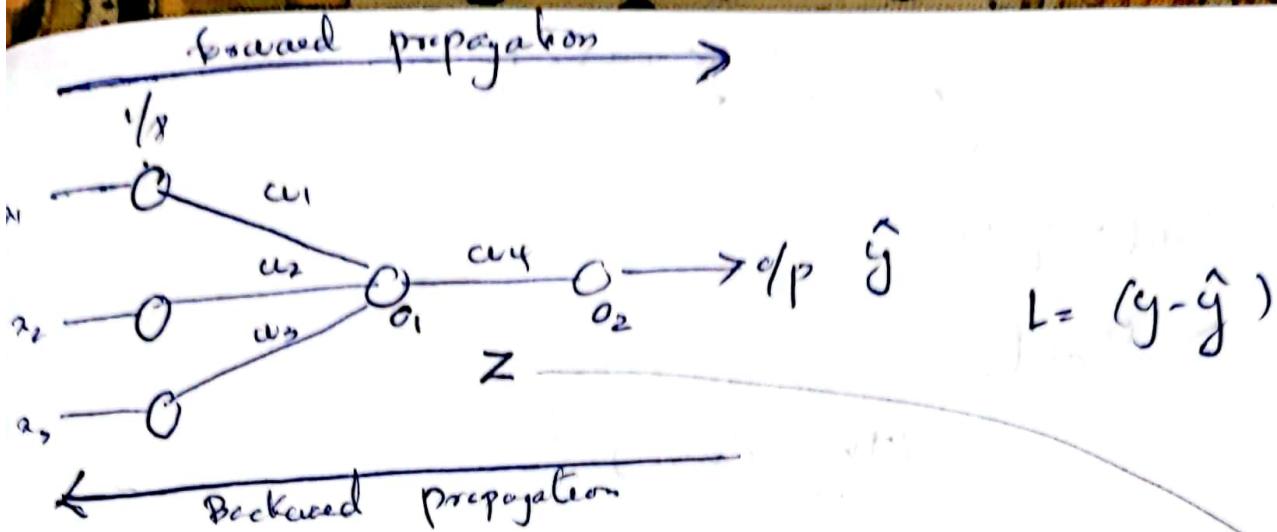
$$w_{\text{new}} = w_{\text{old}} - \eta (+\text{ve})$$

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

Here the main aim is to come to the global minima.

- Learning rate (η) should be a small number, so that we are slowly converging to the global minima.
- If η = large number, there may ~~not~~ be chance that it will never reach global minima.
- Usually $\boxed{\eta = 0.001}$ or $\boxed{\eta = 0.01}$

Chain Rule of differentiation



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

Loss is dependant on o_2 chain rule of derivative.

$$\text{So, } \frac{\partial L}{\partial w_{4\text{old}}} = \frac{\partial L}{\partial o_2} \times \frac{\partial o_2}{\partial w_4}$$

$$z = \sigma(o_1 w_4 + \text{bias})$$

it is nothing but activation function applied on $\underbrace{o_1 w_4 + \text{bias}}$

$$b_2^{\text{new}} = b_2^{\text{old}} - \eta \frac{\partial L}{\partial b_2^{\text{old}}}.$$

bias updation formula.

$$\omega_{1, \text{new}} = \omega_{1, \text{old}} - \eta \frac{\partial L}{\partial \omega_{1, \text{old}}} \rightarrow \text{for updating } \omega_1$$

$\boxed{\text{Loss}}$ dependant on $\boxed{o_2}$ dependant on $\frac{\boxed{o_1}}{\text{eq}}$ dependant
on $\boxed{\omega_1}$ dependant on ω_4

$$\frac{\partial L}{\partial \omega_{1, \text{old}}} = \frac{\partial L}{\partial o_2} \times \frac{\partial o_2}{\partial o_1} \times \frac{\partial o_1}{\partial \omega_{1, \text{old}}}$$

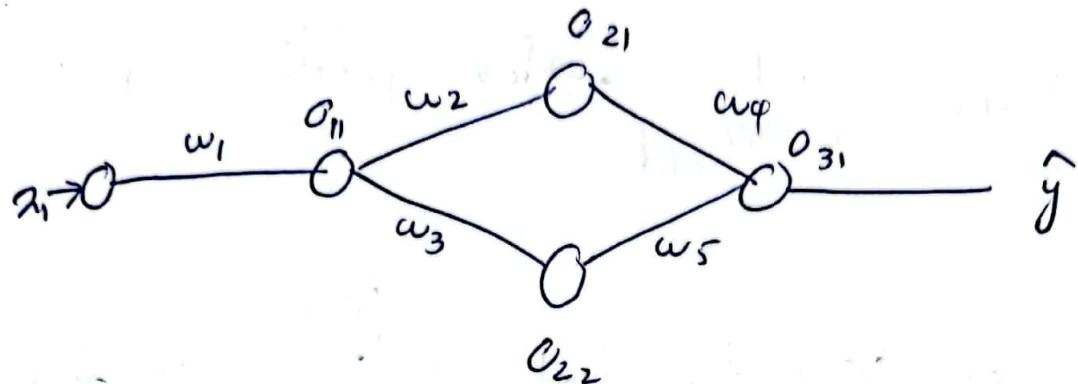
$$\frac{\partial L}{\partial \omega_{2, \text{old}}} = \frac{\partial L}{\partial o_2} \times \frac{\partial o_2}{\partial o_1} \times \frac{\partial o_1}{\partial \omega_{2, \text{old}}}$$

both can also be written as

$$\frac{\partial L}{\partial \omega_{1, \text{old}}} = \frac{\partial L}{\partial o_2} \times \frac{\partial o_2}{\partial \omega_4} \times \frac{\partial \omega_4}{\partial o_1} \times \frac{\partial o_1}{\partial \omega_{1, \text{old}}}$$

and

$$\frac{\partial L}{\partial \omega_{2, \text{old}}} = \frac{\partial L}{\partial o_2} \times \frac{\partial o_2}{\partial \omega_4} \times \frac{\partial \omega_4}{\partial o_1} \times \frac{\partial o_1}{\partial \omega_{2, \text{old}}}$$



$$w_{i,\text{new}} = w_{i,\text{old}} - \gamma \frac{\partial L}{\partial w_{i,\text{old}}}.$$

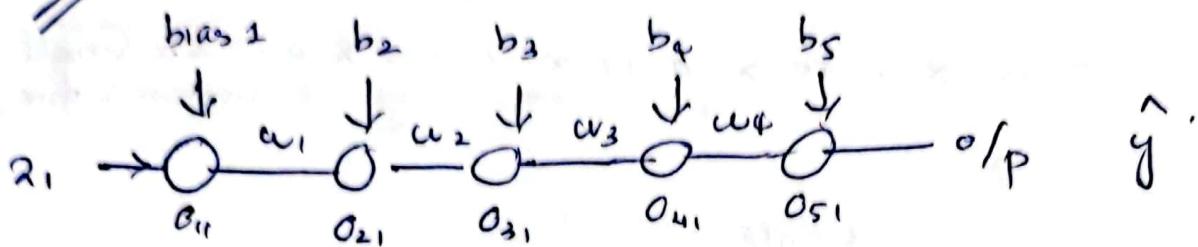
$$\frac{\partial L}{\partial w_{1,\text{old}}} = \left[\frac{\partial L}{\partial o_{31}} \times \frac{\partial o_{31}}{\partial w_4} \times \frac{\partial w_4}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial w_2} \times \frac{\partial w_2}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w_{1,\text{old}}} \right] + \left[\frac{\partial L}{\partial o_{31}} \times \frac{\partial o_{31}}{\partial w_5} \times \frac{\partial w_5}{\partial o_{22}} \times \frac{\partial o_{22}}{\partial w_3} \times \frac{\partial w_3}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w_{1,\text{old}}} \right]$$

can write as

$$\left[\frac{\partial L}{\partial o_{31}} \times \frac{\partial o_{31}}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w_{1,\text{old}}} \right] + \left[\frac{\partial L}{\partial o_{31}} \times \frac{\partial o_{31}}{\partial o_{22}} \times \frac{\partial o_{22}}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w_{1,\text{old}}} \right] = \frac{\partial L}{\partial w_{1,\text{old}}}.$$

V.Imp

Vanishing Gradient Problem

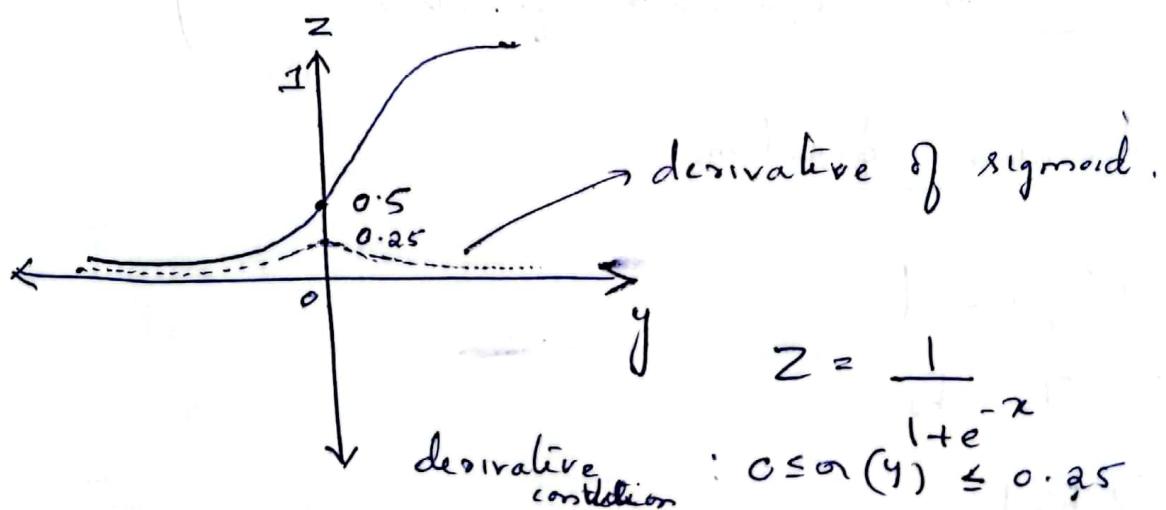


$$\text{Loss} = \frac{1}{2} (y - \hat{y})^2 : \text{MSE}$$

Sigmoid activation functions

- a) Properties of Sigmoid activation functions?
- b) Loss function MSE?

Sigmoid Activation functions



The derivative curve of sigmoid function will be ranging between 0 and 0.25

Because of this,

Since derivative is b/w 0 and 0.25.

for eg: $0.25 \times 0.15 \times 0.10 \times 0.05 \times 0.02 =$ Small value.

Then,

$$W_{\text{new}} = W_{\text{old}} - \eta \downarrow \times (\text{small number})$$

$\Rightarrow W_{\text{new}} \approx W_{\text{old}} \Rightarrow$ Vanishing Gradient Problem.

No change in weights.

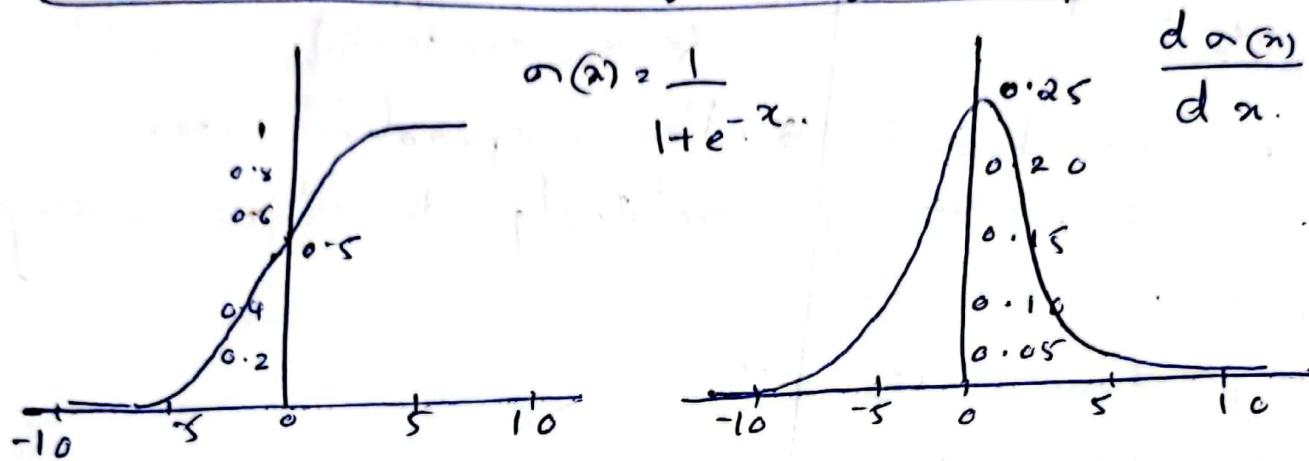
How do we basically solve this? }
Use another activation function }

Activation functions

- 1) Sigmoid
- 2) Tanh
- 3) ReLU
- 4) Leaky ReLU
- 5) Pre ReLU

Commonly Used Activation functions

Sigmoid functions / Logistic functions



- It is the most frequently used activation function in the beginning of Deep learning. It is a smoothing function that is easy to derive.
- In the Sigmoid function, we can see that its o/p is in the open interval (0, 1). It can be thought of as the firing rate of a neuron. In the middle where the slope is relatively large, it is the sensitive area of the neuron. On the sides where the slope is very gentle, it is the neuron's inhibitory area.

Advantages

- 1. Smooth gradients, preventing "jumps" in o/p values.
- 2. Output values bound b/w 0 and 1, normalizing the o/p of each neuron.
- 3. Clear predictions i.e., very close to 1 or 0.

Disadvantages

- 1. Prone to gradient Vanishing
- 2. Function o/p is not zero centered
- 3. Back propagation operation are relatively time consuming.

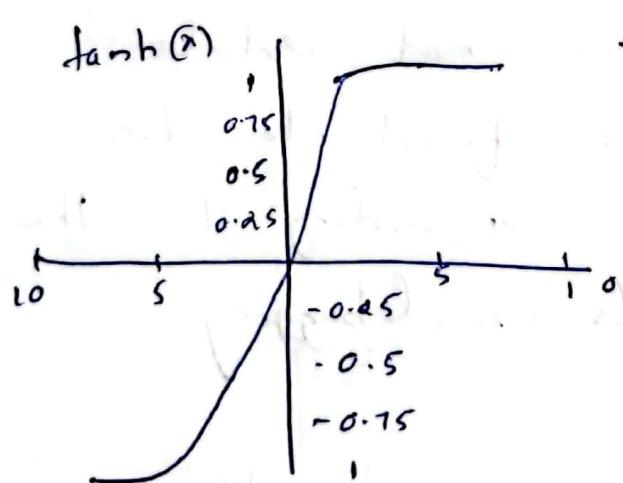
Defects

i) When the i/p is slightly away from the coordinate origin, the gradient of the function becomes very small, almost zero. In the process of neural n/w backpropagation, we all use the chain rule of differential to calculate the differential of each weight w. When the back propagation passes through the sigmoid function the differential on this chain is very small. Moreover, it may pass through many sigmoid functions which will eventually

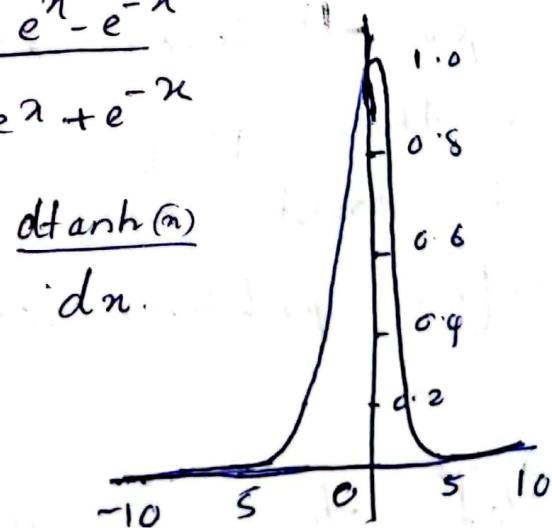
cause the weight to do have little effect on the loss function, which is not conducive to the optimization of the weight. This problem is called gradient saturation or gradient desorption.

- 2) The function o/p is not centered on 0, which will reduce the efficiency of weight update.
- 3) The Sigmoid function performs exponential operations, which is slower for computers.

2. tanh function (Hyperbolic Tangent function)



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



Tanh is a hyperbolic tangent function. The curves of tanh function and sigmoid are relatively similar.

Range $-1 \rightarrow 1$.

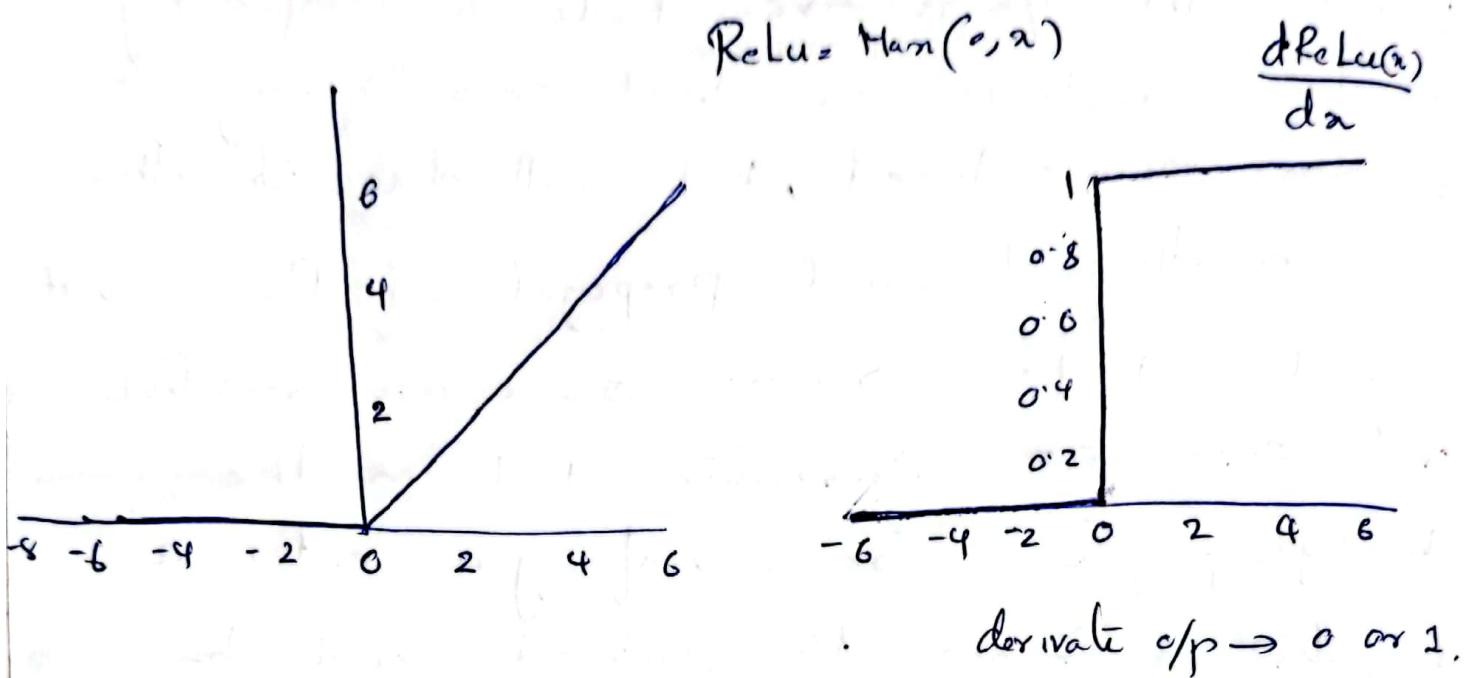
Derivative range $0 - 1$.

first of all when the i/p is large or small, the o/p is almost smooth and the gradient is small, which is not conducive to weight update. the difference is the o/p interval, the o/p of interval of tanh is $-1 - 1$ and the whole function is decreasing much better than sigmoid.

In general binary classification Problems, the tanh function is used for the hidden layers and sigmoid function is used for the o/p layers. However, these are not stable, and the specific activation function to be used must be analysed according to the specific problems, as depends on debugging.

Yet there will be chances of Vanishing Gradient Descent.

3. ReLU function/Rectified Linear unit



→ The ReLU funct is actually a funct that takes the Max value. Note that this is not fully interval - derivable, but can take sub-gradient as shown in the above fig ,

→ Comparing to Sigmoid & tanh it has the following advantages.

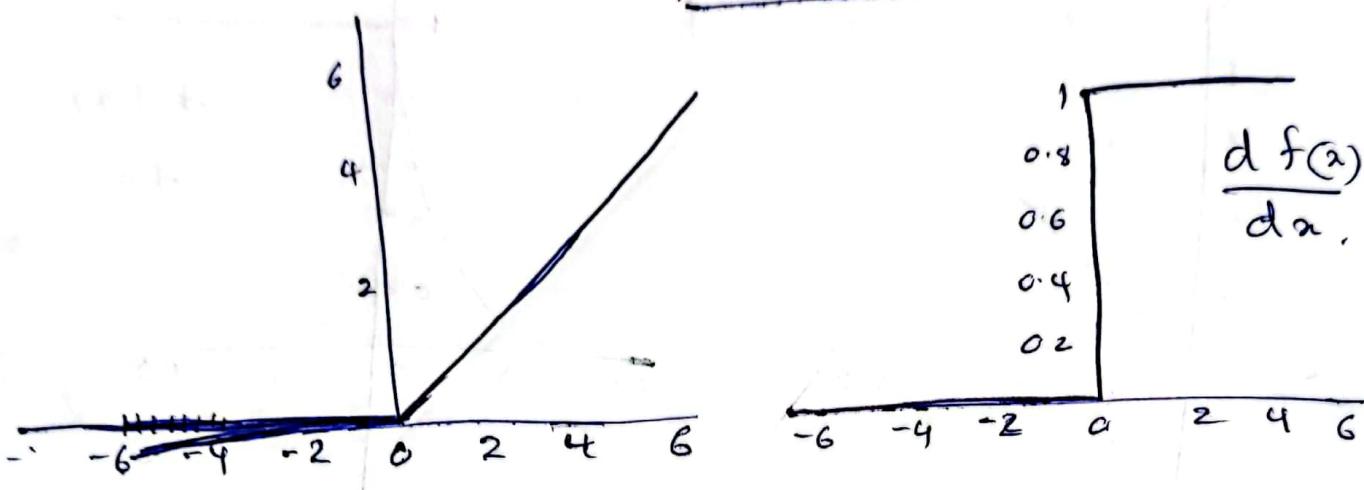
1. When the o/p is positive, there is no gradient saturation problem
2. The calculation Speed is much faster. The ReLU funct has only a linear Relationship. Whether it is forward or backward, it is much faster than Sigmoid and tanh

Disadvantages

- ① When the i/p is -ve, ReLu is completely inactive, which means that once a negative number is rendered, ReLu will die. In this way, in the forward propagation process, it is not a problem. Some areas are sensitive and some are insensitive. But in the backpropagation process, if you enter a negative number, the gradient will be completely zero, which has the same problem as the sigmoid function and tanh funct.
- ② We find that the o/p of the ReLu funct is either 0 or a positive number, which means that the ReLu function is not a 0-centered funct.

4. Leaky ReLU functions

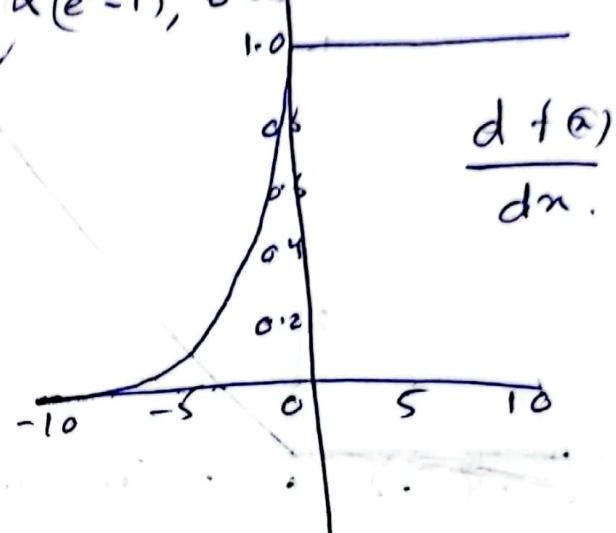
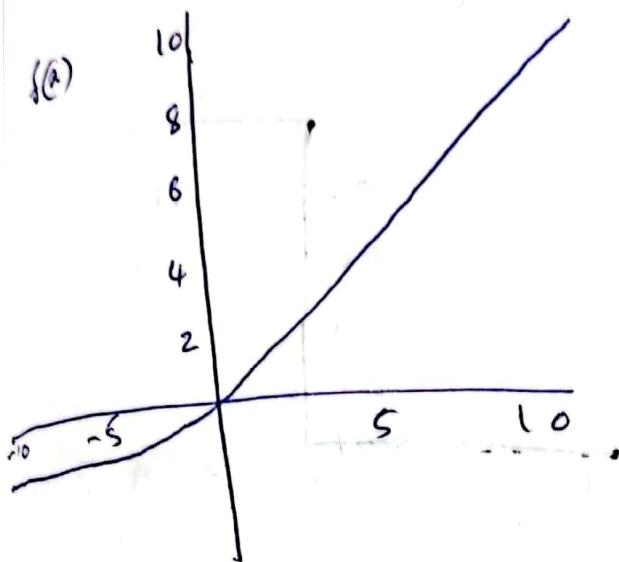
$$f(x) = \max(0.01x, x)$$



- In order to solve the Dead ReLU Problem, people proposed to set the first half of ReLU $0.01x$ instead of 0. Another intuitive idea is a parameter-based method, Parameteric ReLU: $f(x) = \max(\alpha x, x)$, where α can be learned from back propagation.
- In theory Leaky ReLU has all the advantages of ReLU, plus there will be no problems with Dead ReLU, but in actual operations it has not been fully proved that Leaky ReLU is always better than ReLU.

5. ELU (Exponential Linear Unit) function

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



ELU is proposed to solve the problem of ReLU. Obviously, ELU has all the advantages of ReLU.

- No Dead ReLU Issues

- The mean of the c/p is close to 0, zero-centered.

The small problem is that it is slightly more computationally sensitive. Similar to Leaky ReLU, although theoretically better than the ReLU, there is currently no good evidence in practice that ELU is always better than ReLU.

6. Softmax

$$S(x_j) = \frac{e^{x_j}}{\sum_{k=1}^K e^{x_k}} \quad ; j = 1, 2, \dots, K.$$

- For arbitrary real vector of length K , Softmax can compress it into a real vector of length K with a value in the range $(0, 1)$ and the sum of the elements in the vector is 1.

- It also has many applications in Multiclass classification and neural n/w. Softmax is different from the normal max function; the max function only outputs the largest value and Softmax ensures that smaller values have a smaller probability and will be discarded directly. It is a 'max' that is soft.
- The denominator of the Softmax function combines all factors of the original o/p value, which means that the different Probabilities obtained by the Softmax function are related.

to each other in the case of binary classification. for Sigmoid.

there are,

$$P(y=1/x) = \frac{1}{1+e^{-\theta^T x}}$$

$$P(y=0/x) = 1 - P(y=1/x) = \frac{e^{-\theta^T x}}{1+e^{-\theta^T x}}$$

for Softmax with $K=2$, there are,

$$P(y=1/x) = \frac{e^{\theta_1^T x}}{e^{\theta_0^T x} + e^{\theta_1^T x}} = \frac{1}{1 + (e^{\theta_0^T} - e^{\theta_1^T})x} = \frac{1}{1 + e^{-\beta x}}$$

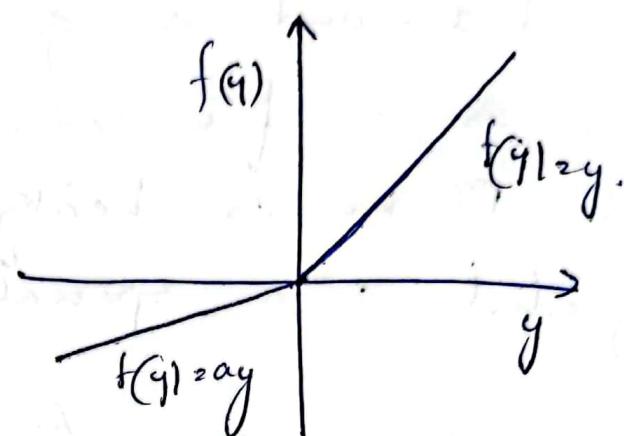
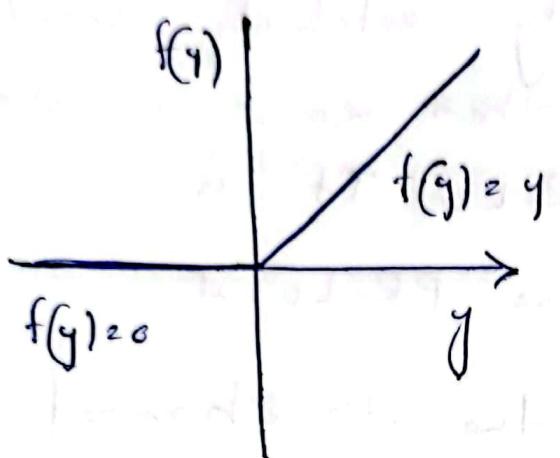
$$\begin{aligned} P(y=0/x) &= \frac{e^{\theta_0^T x}}{e^{\theta_0^T x} + e^{\theta_1^T x}} = \frac{e^{(\theta_0^T - \theta_1^T)x}}{1 + e^{(\theta_0^T - \theta_1^T)x}} \\ &= \frac{e^{-\beta x}}{1 + e^{-\beta x}} \end{aligned}$$

Among them, it can be seen that

$$\beta = -(\theta_0^T - \theta_1^T)$$

In the case of binary classification, Softmax is degraded to Sigmoid.

7. PReLU (Parametric ReLU)



ReLU V/S PReLU

For PReLU, the coefficient of the negative Part is not constant and is adaptively learned.

PReLU is also an improved version of ReLU. In the negative region, PReLU has a small slope, which can also avoid the problem of ReLU death. Compared to ELU PReLU is a linear operation in the negative Region. Although the slope is small it does not tend to 0, which is a certain advantage.

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

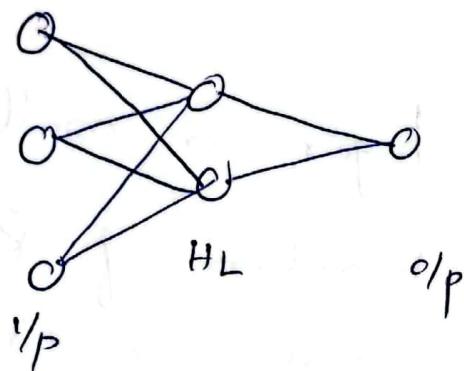
We look at the formula of PReLU. The parameter α is generally not a number like 0 and 1, and is generally relatively small such as a few zeroes. When $\alpha = 0.01$ we call PReLU as Leaky ReLU. It is regarded as a special case PReLU.

Above y_i is any if on the i th channel a_i is the negative slope which is a learnable parameter.

$\left\{ \begin{array}{l} \text{if } a_i = 0, f \text{ becomes ReLU} \\ \text{if } a_i > 0, f \text{ becomes Leaky ReLU.} \\ \text{if } a_i \text{ is a learnable parameter, } f \text{ becomes PReLU} \end{array} \right\}$

Imp Techniques on which activation functions are we should use?

If Suppose you have binary classification,
You have a neural n/w

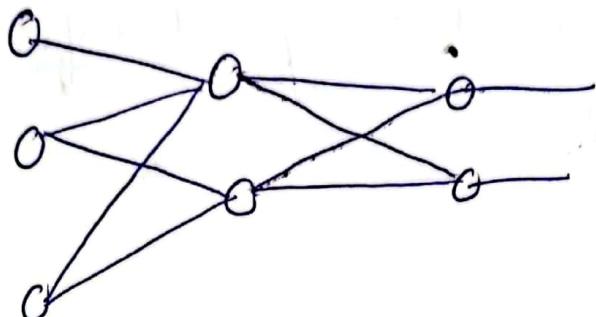


- In the hidden layers always try to use ReLu Activation funt.
- In o/p layers, Sigmoid activation functions.

Nb: If using ReLu, convergence is not happening then change ReLu to PReLU or ELU

Imp But the o/p must have Sigmoid for binary classifications.

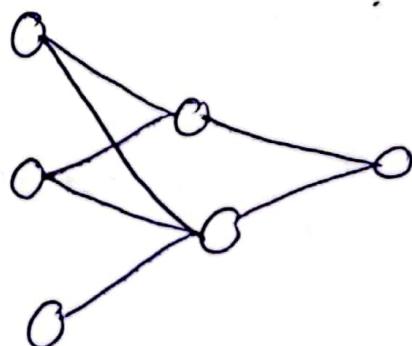
Suppose You have multiclass classification problems,



- Use Softmax for o/p layers in multiclass.
- (or) Use ReLU or the combinations of ReLU, PReLU, ELU if the convergence is not happening.

Suggested : ReLU

In Case of Regression Problems,



- In the hidden layers use any variant of ReLU.
- In the o/p layer use Linear Activation function

- There will be separate loss function for this.

Loss functions

For Regression

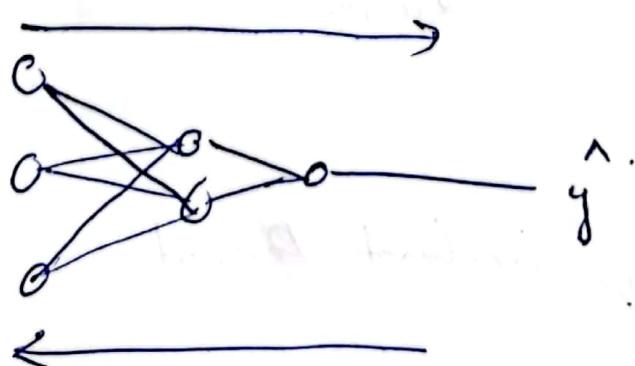
In Deep learning specifically - ANN, we try to solve two types of problems classification and regression.

In Regression Problems,

1. MSE (Mean Squared Error)
2. MAE (Mean Absolute Error)
3. Huber loss.
4. RMSE

Epochs, Batch size?

Loss function and Cost function



Suppose I have 100 records, If I Pass

each records separately and do the forward and backward propagation it takes time. So instead I pass a batch of records for eg: 10.

difference b/w Loss func & Cost func

Loss function \rightarrow Specifically for single record
Cost function \rightarrow Specifically for batch of records.

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

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

I. Mean Squared Error (MSE)

Loss func

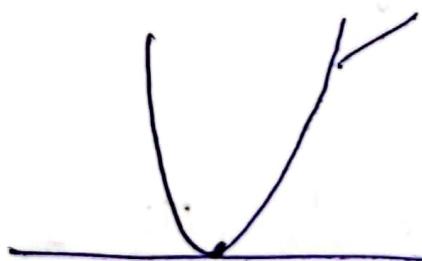
$$\frac{1}{2n}(y - \hat{y})^2$$

Cost func

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

\downarrow
Quadratic equation.

Gradient Descent.



Advantages

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



if I have some outliers.



The reason it changes because we are penalizing the error. Since we are squaring the error and it will have major shift.

2. Mean Absolute Errors (MAE)

Loss function

$$\frac{1}{2} \sum |y - \hat{y}|$$

Cost function

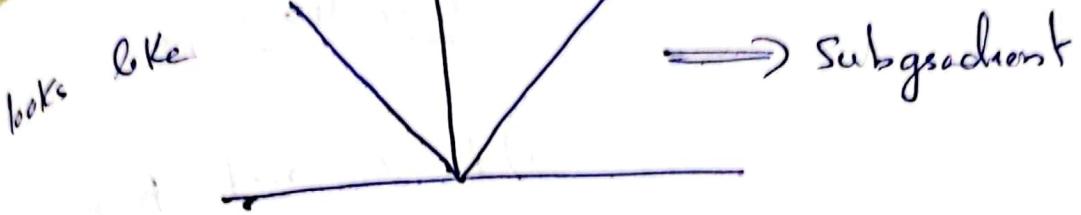
$$\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Advantages

1. Robust to outliers.

Disadvantages

1. Time consuming
 $> MSE$.



3) Huber loss.

It is a Combination of MSE and MSF.

Cost function

~~Loss function~~

(outliers are not present)

$$\frac{1}{2}(y - \hat{y})^2 \text{ if } |y - \hat{y}| \leq \delta$$

$$\delta |y - \hat{y}| - \frac{1}{2} \delta^2, \text{ o.w.}$$



Hyperparameters.

Cost functions

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

$$\delta |y - \hat{y}| - \frac{1}{2} \delta^2, \text{ o.w.}$$

4. Root Mean Squared Error (RMSE).

Loss function

$$\sqrt{\frac{(y_i - \hat{y}_i)^2}{n}}$$

Cost function

$$\sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n}}$$

Advantages

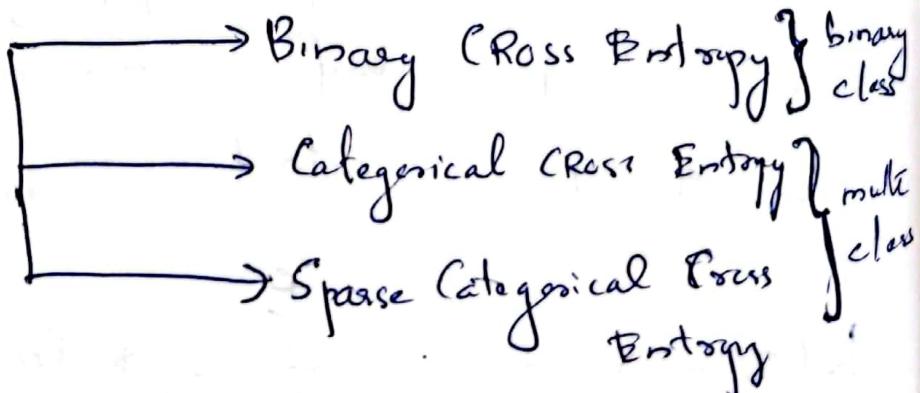
- Penalize the large errors
- Optimizes According to MSE
- Differentiable
- Respects zero values.

Disadvantages

- Sensitive to outliers
- Does not consider error relative to actual value
- Not Robust to outliers

For Classification

Cross ENTROPY



Binary Cross Entropy

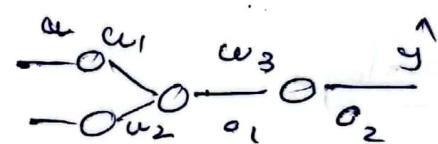
$$\text{Loss} = -y * \log(\hat{y}) - (1-y) * \log(1-\hat{y}) \Rightarrow \text{Logistic Regression}$$

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

How do we calculate \hat{y} ?

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

sigmoid activation function in last layer.



$$\hat{y} \leftarrow o_2 = o_1 * w_3 + \text{bias}$$

Categorical Cross Entropy (Multiclass classification Problem)

f_1	f_2	f_3	$\%P$		Good	Bad	Neutral
2	3	4	Good	one hot encoding	1	0	0
5	6	7	Bad		e	1	0
8	9	10	Neutral		0	0	1

$$h(x_i, y_i) = - \sum_{j=1}^c y_{ij} * \ln(\hat{y}_{ij}) \quad i = \text{row} \quad j = \text{column}$$

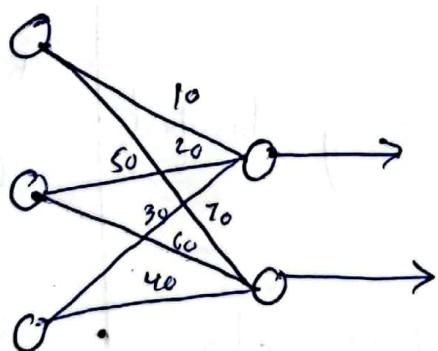
$$y_i = [y_{i1}, y_{i2}, y_{i3}, \dots, y_{ic}]$$

c: no. of categories.

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

\hat{y}_{ij} = Softmax Activation function
 (It is applied on o/p layer of multiclass classification Problem)

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



$$\text{for } z_1 = 10 \quad \frac{e^{10}}{e^{10+20+30+40+50+70}}$$

will be probability.

Total will be 1.

3. Sparse Categorical Cross Entropy

0 1st 2nd → categories.
 [0.2, 0.3, 0.5
 ↓
 2nd Indon → o/p.

$[0, 0.3, 0.1, 0.2, 0.2]$ → 1st Index.
 category → o/p.

Disadvantage

Loose info about the probability of other category

* Right Combinations.

Hidds layers	O/p layers	Problem Statement	Loss function
ReLU or z^2 's variant	Sigmoid	Binary Classification	Binary Cross Entropy
ReLU or z^2 's variant	Sofmax	Multi class	Categorical or Sparse CE
ReLU or z^2 's variant	Linear	Regression	MSE, MAE, Huber Loss, R MSE.

The sparse categorical cross entropy loss function computes the cross-entropy loss b/w

the integer labels and the predicted Probabilities.
It is suitable for scenarios where there are
two or more label classes, and the labels are
provided as integers rather than one-hot
representations.

Usage

This loss function is particularly useful when
working with classification tasks that have
multiple classes and whose labels are represented
as integers. It expects integer labels for each
sample and compute the CE loss based on
these integer labels.

Difference b/w Categorical CE and Sparse Categorical CE

Categorical CE	Sparse Categorical CE
<ul style="list-style-type: none">Requires Label to be one-hot encoded, where each label is represented as a binary vector with a single "1" indicating the class	<ul style="list-style-type: none">Utilizes integer-encoded labels where each label is represented by an integer value corresponding to the class.

Ideal for Scenarios where each sample belongs to only one class, making it suitable for mutually exclusive classes.

② Suited for problems where sample can belong to multiple classes or when dealing with soft probabilities.

3 Mathematical formulation

the loss function calculates the cross entropy b/w the true one-hot encoded labels and the predicted probabilities

the loss function computes cross entropy with integer targets, optimizing based on the difference b/w Predicted Probabilities and integer labels

TensorFlow → Google
PyTorch → Facebook.

Keras, PyTorch.

Day ³

1) Optimizers

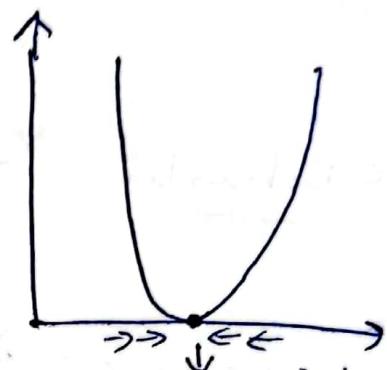
- ① Gradient Descent
- ② SGD (Stochastic Gradient Descent)
- ③ Mini Batch SGD
- ④ SGD with momentum
- ⑤ Adagrad - Adaptive Gradient
- ⑥ RMS Prop
- ⑦ Adam optimizers

- 2) Batch
- 3) Epochs
- 4) Iterations
- 5) ANN practical.

~~GRT~~ Optimizers

1. Gradient Descent

Loss/cost

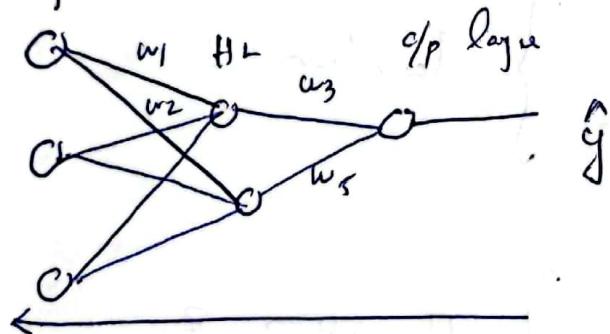


Weight Update Formula.

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}} \quad \xrightarrow{\text{Learning rate}}$$

forward prop. →

y_p layer



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

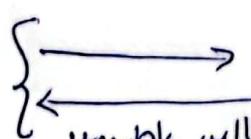
Epochs

A combination of forward and backward propagation

Epochs, Iterations

Dataset = 1000

1 Epoch



1 Epoch = 1 iteration

weights will get updated.

$\hat{y}_i \Rightarrow$ cost function

100 iterations,
batch size = 10



Advantages

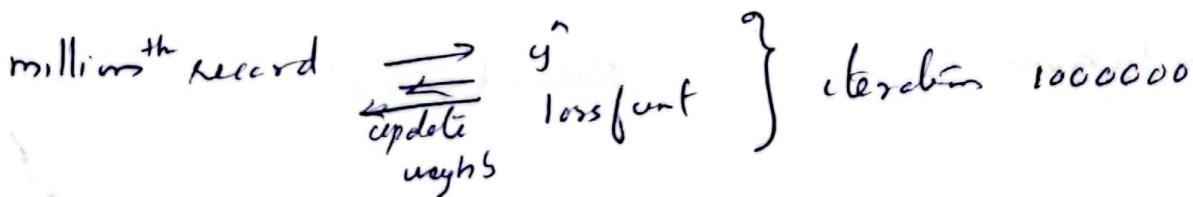
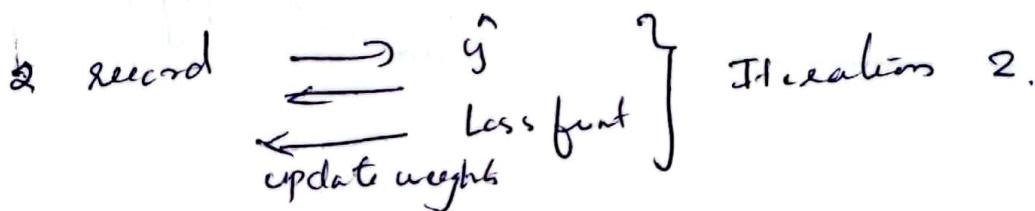
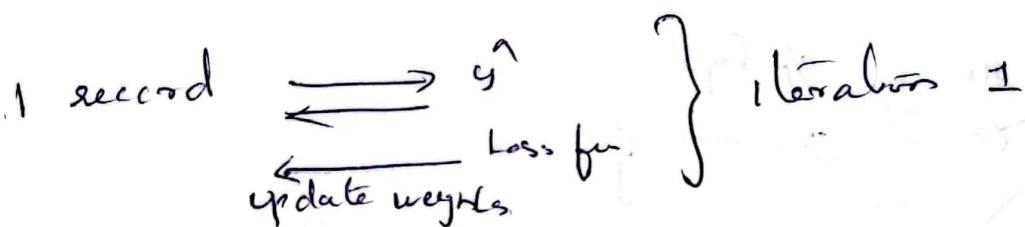
Convergence will happen.

Disadvantages

- ① Resource Extensive & required huge RAM to process millions records.

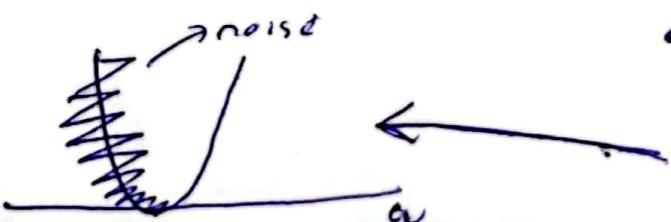
2. Stochastic Gradient Descent

1000000s of records.



Disadvantage

- Solve Resource Issue.

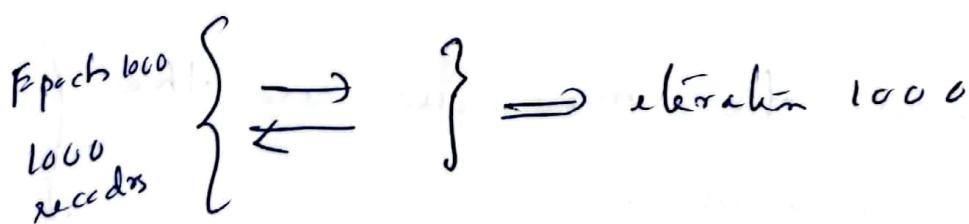
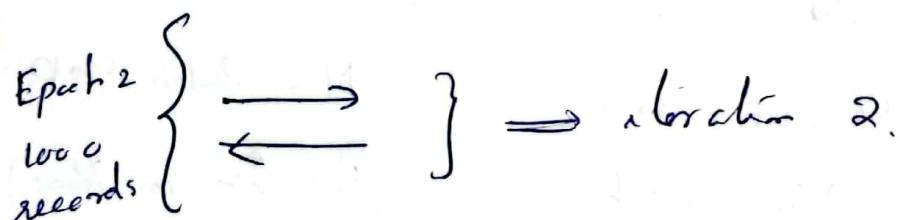
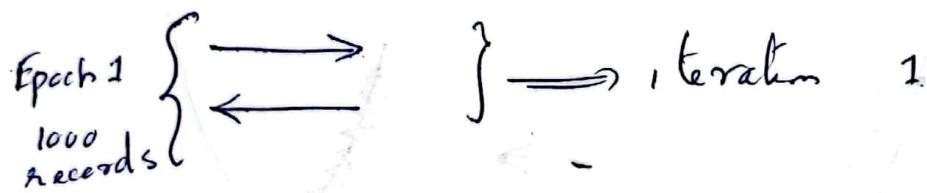


- Time Complexity ↑↑
- Convergence will also take more time
- Noise gets introduced.

3. Mini Batch SGD

Instead of passing single records, we pass batch of records.

If we have 1000000 of records, and batch size = 1000. Then for every iteration,



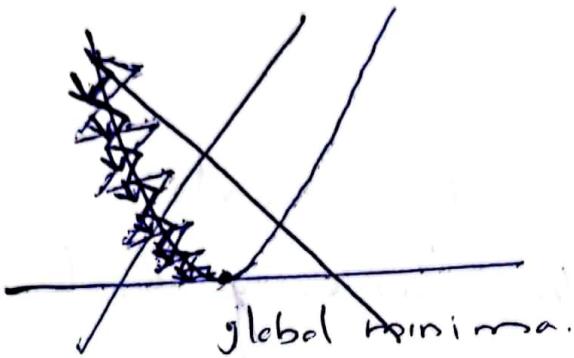
Advantages.

1. Convergence Speed will increase
2. Noise will be less when compared to SGD
3. Efficient Resource Usage (RAM)

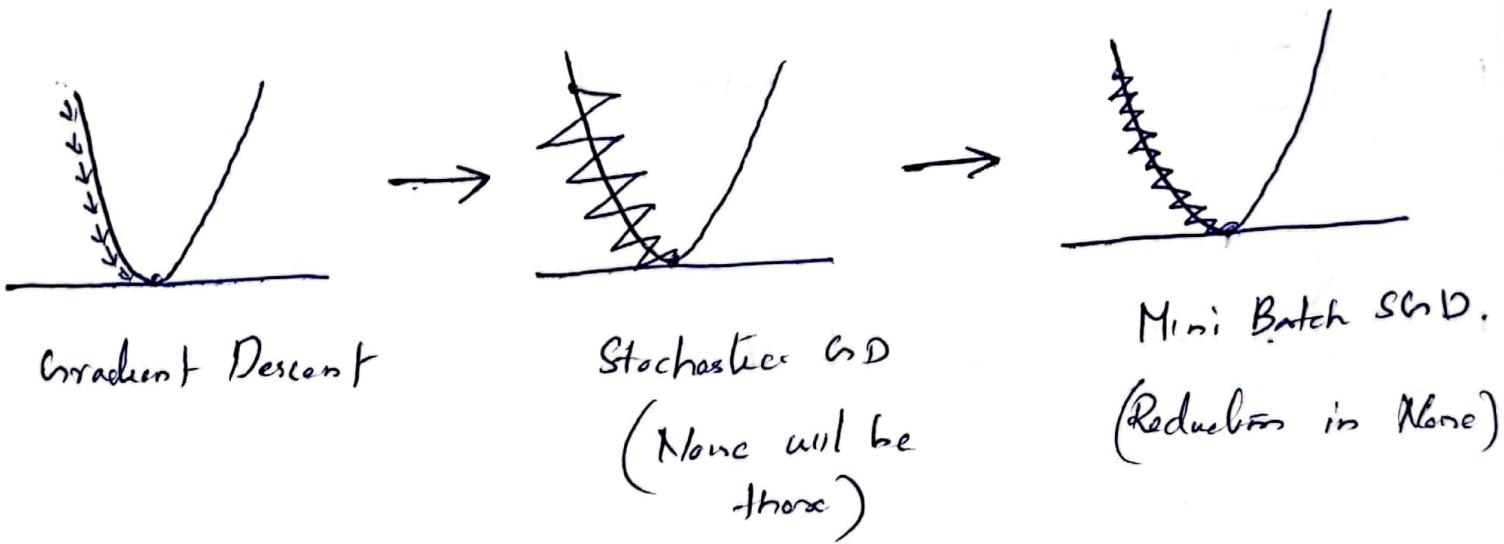
Disadvantages

1. Noise still exists

4. Time Complexity will improve.



Main aim is to reach the global minima.



{ So how do we remove the noise? }

{ In order to remove the noise we use the concept called Momentum }

Min Batch 4. SGD with Momentum.

Momentum will smoothen the journey and reduce the noise.

Exponential Weighted Average.



Time Series



ARIMA, ARMA, ...

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

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

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

Exponential Weighted Average :

time $t_1 t_2 t_3 t_4 \dots t_n$

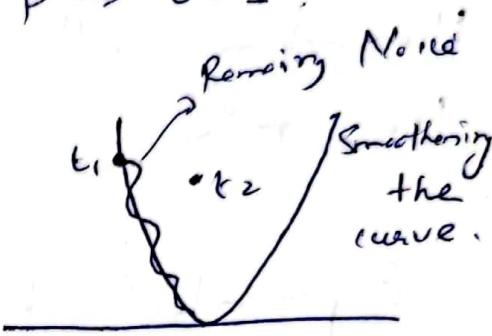
Value $a_1 a_2 a_3 a_4 \dots a_n$

Value at time $t_1 = a_1$

$\beta \rightarrow$ hyperparameter

$$V_{t_2} = \beta * V_{t_1} + (1 - \beta) * a_2 \quad \beta \rightarrow 0 - 1$$

$$\boxed{\beta = 0.95}$$



$$V_{t_2} = 0.95 * a_1 + 0.05 * a_2$$

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

$$= 0.95 (0.95 * a_1 + 0.05 * a_2) + 0.05 * a_3$$

Where do we apply this? (Exponential Weighted Average)

We apply this in the derivative of loss cost derivative of weights.

Exp. Weighted Avg

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

$$w_t = w_{t-1} - \eta Vdw,$$

$$\therefore Vdw = \beta \times V_{dw,t-1} + (1 - \beta) * \frac{\partial L}{\partial w_{t-1}}$$

Advantages .

- Reduce the noise
- Working on Mini Batch
- Quicker convergence

Interview
an

If there is so much noise in the Gradient Descent . what should we do?

ans: Try to bring the smoothening factor in the optimizer.

Imp
5. Adagrad → Adaptive Gradient Descent

$\eta \rightarrow$ helps us to maintain the speed of the convergence: learning rate

? How can we change the value of the learning rate while we move towards the global minima.

In previous optimizers the learning rate is fixed.
So here instead of fixing the learning rate, we need to make it adaptive as we move towards the global minima.

$$w_t = w_{t-1} - \eta \frac{dL}{dw_{t-1}}$$



$$w_t = w_{t-1} - \eta' \frac{dL}{dw_{t-1}}$$

$$\eta' = \frac{\eta}{\sqrt{\alpha_t + \epsilon}} \rightarrow \text{learning rate}$$

$\sqrt{\alpha_t + \epsilon} \rightarrow$ small number to avoid the divide by zero condition i.e., the denominator never becomes zero.

$$\alpha_t = \sum_{i=1}^t \left(\frac{\partial L}{\partial w_i} \right)^2 \quad t \rightarrow \begin{matrix} \text{current time stamp,} \\ (\text{the weights at that particular time stamp}) \end{matrix}$$

Main aim to decrease the learning rate as we move to global minima.

As time goes,

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

because α_t is the summation from $i=1$ to $i=t$.

As α_t value \uparrow , $\eta' \downarrow$

Then there is possibility that α_t value will be huge.

Because of this its $\eta' = \frac{\eta}{\sqrt{d_t + \epsilon}}$ there will be

negligible change

6. AdaDelta and

Disadvantage

$\eta' \rightarrow$ Possibility to become
a very small value ≈ 0

So we need to ensure that η' should reach
a huge value.

6. AdaDelta and RMSProp

$$\eta' = \frac{\eta}{\sqrt{s_{dw} + \epsilon}}$$

{We are going to apply exponential}
 $s_{dw} = 0 \downarrow$ weighted avg

$$s_{dw_t} = \beta s_{dw_{t-1}} + (1-\beta) \left(\frac{\partial L}{\partial w_t} \right)^2$$

(Dynamic LR + Smoothing EWA)

if $\beta = 0.95$

$$s_{dw_t} = 0.95 \times s_{dw_{t-1}} + 0.05 \left(\frac{\partial L}{\partial w_t} \right)^2$$

Here we can control the increase of s_{dw} .

Hence there will be slow decrease in η' and will be able to reach the global minima.

Whenever we combine the exponential weighted avg and the adaptive learning rate we get,

[1. Adam Optimizer]

[Best Optimizer]

Momentum + RMS Prop (Adaptive Learning Rate)

We have V_{dw} , s_{dw} ,

$$\downarrow$$
$$V_{db} \quad s_{db}$$

initially $V_{db} = 0$, $s_{db} = 0$.

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

weight updation

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

bias updation

$$\eta' = \frac{\eta}{\sqrt{s_{dw}} + \epsilon}$$

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

$$V_{d_{wt}} = \beta * V_{d_{wt-1}} + (1-\beta) \frac{\partial L}{\partial w_{t-1}}$$

\Rightarrow Momentum Smoothening

$$V_{d_{bt}} = \beta * V_{d_{bt-1}} + (1-\beta) \frac{\partial L}{\partial b_{t-1}}$$

$$\eta = \frac{\eta}{\sqrt{S_{d_{wt}} + \epsilon}}$$

$$S_{d_{wt}} = 0.$$

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

- It solves problem of Smoothening
- Make sure that learning rate becomes adaptive

Adams Optimizer

- 1) It is an optimizer which combines the best properties of Adaline and RMSProp.
- 2) It is based on adaptive estimation of first order and second order moments, making it computationally efficient and memory efficient.

3) Adam Adjusts learning rates individually for each parameter leading to more efficient optimization, especially in complex models with many parameters.

Day 4 -

1. ANN Practical Implementation
2. Early Stopping
3. Black box Models vs White Box Models
4. CNN Introduction.

In Case of ANN

Interview Question

For which all algorithms feature scaling is required?

ans: ANN? Linear? Log? Decision Tree? RF?
✓ ✓ ✓ Not necessary Not necessary

Necessary

XGBoost Not necessary

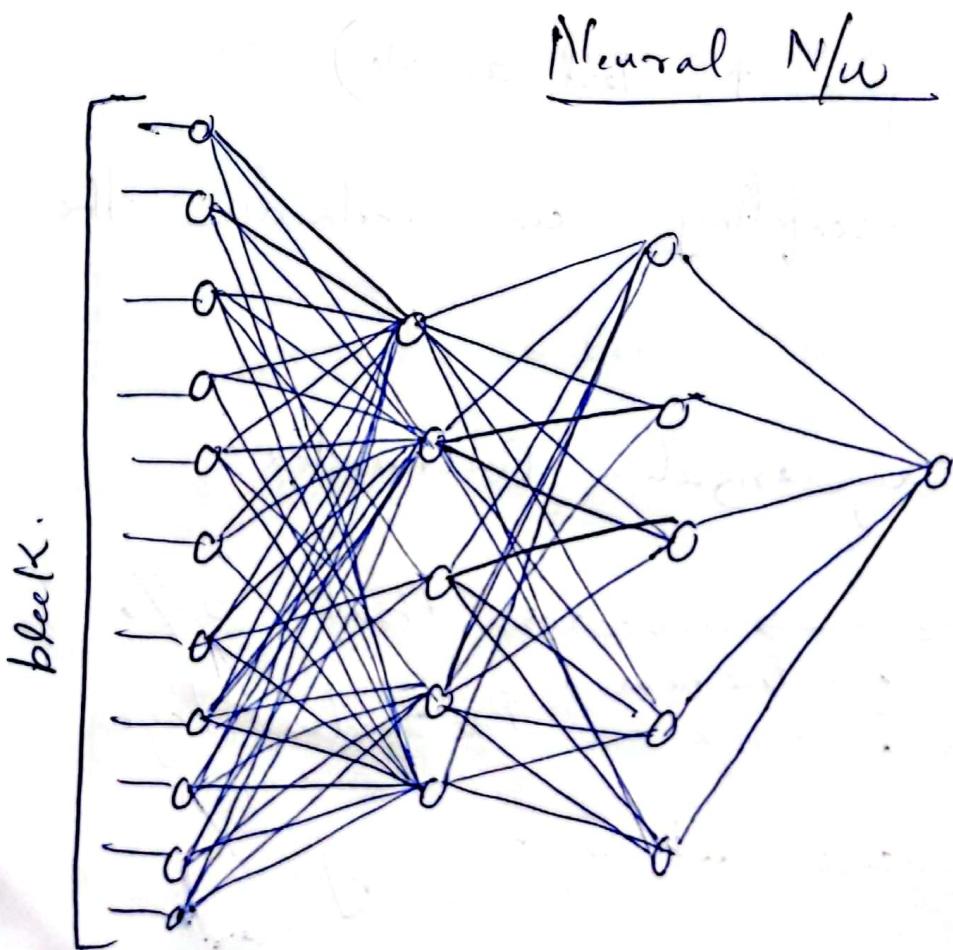
K-Means ✓

KNN ✓

⇒ Anything related to Distance based, we need to do feature scaling.

⇒ Wherever gradient Descent / optimizers are used, scaling will be required.

1. Sequential
2. Dense
3. Activation



if I take this entire neural net together,
as a block \rightarrow Sequential.

And can do forward and backward propagation.

\Rightarrow Whenever we want to create neurons, we
use [dense]

With the help of dense \rightarrow we will be
able to create HLs, i/p layers, o/p layers.

\Rightarrow Dropout

Sometimes the entire neural net leads to
overfitting (Training acc \uparrow , Test acc \downarrow)

To reduce the overfitting, we introduce the
dropout layers.

It is like a regularization parameter.

If dropout = 0.3 \rightarrow indicates that 30%
of data present ^{of entire neurons} inside a particular layer,
will get deactivated while training.

Q) How many numbers of epochs should we basically stop?

While I was faced with overfitting, used the
 $\text{dropout} = 0.2$ to reduce the overfitting in each
HL.

Black Box Model vs White Model

Random Forest \rightarrow Black Box Model

Decision Tree \rightarrow Decision for White Box Model

ANN \rightarrow Black Box Model

Xgboost \rightarrow Black Box Model

Linear Regression \rightarrow White box Model.

In Some Algorithms, the internal working and all are difficult to monitor.

People ~~with~~ come up with solutions to monitor each model, how the model is performing and details about each I/P feature → That is above explainable AI comes up.

White Box Model	Black Box Model
<ul style="list-style-type: none">They are transparent and interpretable, allowing users to understand how predictions are made.Suitable for scenarios requiring high accountability and trust, like risk assessment, robotics, and scientific research.Utilizes models with higher transparency and interpretability such as linear models or decision tree.	<ul style="list-style-type: none">Black box models prioritize innovation and accuracy but lack transparency in how predictions are generated.Commonly used for complex scenarios with deep, non-linear interactions b/w data.Employ techniques like deep learning, boosting models, and random forests.

Pros

Pros

- Rising in popularity due to increasing mistrust of

- High predictive accuracy, making them indispensable for tasks like speech

Opaque AI Systems.

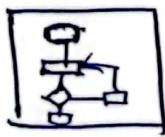
- Provides reliable predictions and is suitable for scenarios where understanding the model's decision making process is crucial.

recognition and fraud detection.
• Effective at handling large datasets with complex relationships.

Cons

- May not produce groundbreaking results or immoral ideas.
- Less suitable for modelling complex relationships, potentially leading to lower accuracy.
- Low algorithm interpretability, making it challenging to explain decisions to stakeholders.
- Results may not be perfectly reproducible due to the complexity of the model.

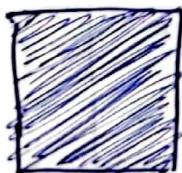
I/P



o/p.

White box Testing

I/P



o/p.

Black Box testing

Day 5

Convolutional Neural N/w.

- Images / Videos or Video frames
- object detections

Agenda

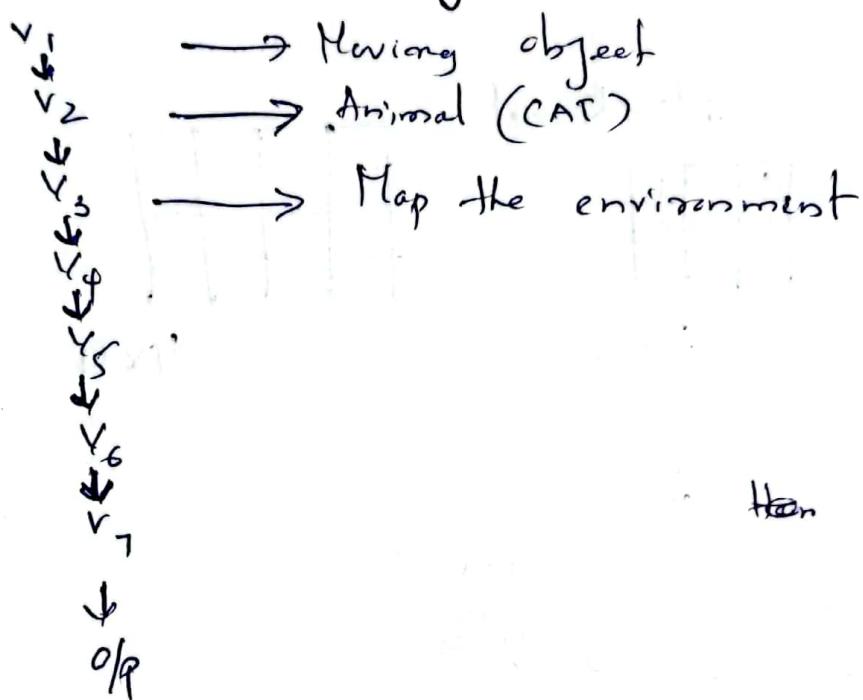
- 1) CNN v/s Humans Brains
- 2) Convolution Operations.
 - Convolution
 - Padding
 - Stride
 - Filters (Kernels)
- 3) Max-Pooling
- 4) Flattening
- 5) Practical.



Major Part of Brain \rightarrow Cerebral Cortex.

\downarrow
visual Cortex (Part)

It has many layers



Similarly as case of CNN.

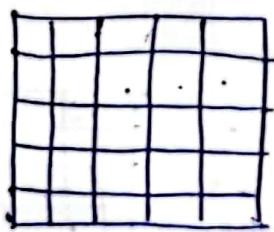
CNN.

1. Convolution:

Images?

B & W, RGB.

\downarrow
1 channel



5x5

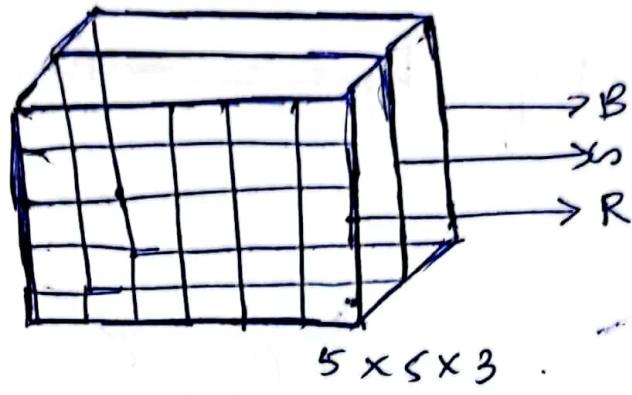
0.255

0 - black
255 - white

Every Pixel may have different values based on composition.

Not

RGB,



Convolution

Image.

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

↓
 6×6

filter/Kernel.

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

↓
 3×3

1	1	1	1
1	1	1	1
1	1	1	1
1	1	1	1

4×4

minimum scaling,
converting all the values
of the pixels to 0, 1

horizontal edge filter

We have 6×6 image.

Passed it through ~~a~~ a 3×3 filter specifically called
horizontal edge filter.

1)

0	0	0
0	0	0
0	0	0

 \rightarrow

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

 \rightarrow

$$\begin{array}{l} 0+0+0 \\ 0+0+0 \\ 0+0+0 \end{array}$$
 \rightarrow

0	1	1
0	1	1
0	1	1

2)

0	0	1
0	0	1
0	0	1

 \rightarrow

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

 \rightarrow

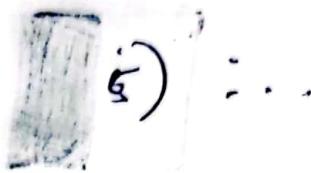
$$\begin{array}{l} 0+0+1 \\ 0+0+0 \\ 0+0-1 \end{array}$$
 \rightarrow

0	0	1
0	0	1
0	0	1

Q8. Now

3) $\begin{array}{|c|c|c|} \hline c & 1 & 1 \\ \hline c & 1 & 1 \\ \hline c & 1 & 1 \\ \hline \end{array}$

4) $\begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array}$



Let us consider

$$\begin{array}{|c|c|c|c|} \hline 0 & 0 & 0 & 0 \\ \hline -4 & -4 & -4 & -4 \\ \hline -4 & -4 & -4 & -4 \\ \hline 0 & 0 & 0 & 0 \\ \hline \end{array}$$

as the output.

If the filter is

$$\begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 2 & 0 & -2 \\ \hline 1 & 0 & -1 \\ \hline \end{array}$$

Vertical edge filter/kernel.

If we are reverting back the feature scaling we get the values ranging b/w 0-255

it becomes

$$\begin{array}{|c|c|c|c|} \hline 255 & 255 & 255 & 255 \\ \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 255 & 255 & 255 & 255 \\ \hline \end{array}$$

The highest value will be converted to 255 (white) and lowest value will be " " 0 (black.)



Similarly in the case of human brain,
there a nervous layer we inserted a layer.

Here horizontal edge filter / vertical edge filter / other filters:

⇒ Filters: with the help of these filters we are trying to extract the information from the image.

~~Imp~~ Whenever we pass a $\boxed{6 \times 6}$ image into a $\boxed{3 \times 3}$ filter we get a $\boxed{4 \times 4}$ o/p.
How this happens?

Let say size of the image $\approx n = 6$ and
filter size = 3.

$$\text{o/p size} = \boxed{n - f + 1} = 6 - 3 + 1 = 4$$

? How do you define the ~~filter~~ ?

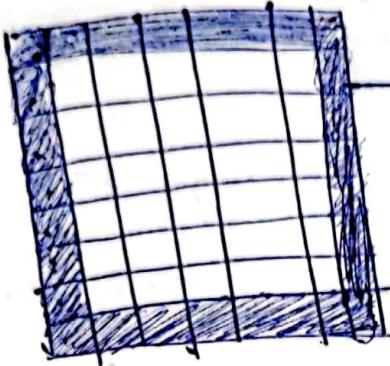
? While we are passing a 6×6 image into a 3×3 filter we get 4×4 . That means we are losing some kind of information. So in order to prevent that what can we actually do?

ans: Solution is Padding.

Padding

Padding is building a compound around the image.

It is a technique used to preserve the spatial dimensions of the i/p data after convolution operations. It involves adding extra elements, such as zeros, to the i/p data before applying any convolutional filter. This helps prevent information loss, especially from the edge of the image, and maintains the spatial relationship and characteristic of the i/p data.



Zero Padding → fill the shaded cells with zero.
→ fill the cells with nearest value.

1. Same Padding
2. Valid Padding
3. Causal Padding

The updated formula = $\lceil \frac{m+2p-f+1}{s} \rceil$

$m \rightarrow$ size of the image

$p \rightarrow$ No of layers of Padding

$f \rightarrow$ filter size

$$m=6 \quad p=1 \quad f=3$$

$$6+2-3+1 = \underline{\underline{6}}$$

In CNN we have to make sure that we update the filters based on the o/p image. So with the help of back propagation we need to update the filters.

Where do we specifically use the activation function?

After the convolution operation, whenever we get the output, on top of each value we apply a ReLu activation function ($\max(0, x)$).

? But why do we apply ReLu

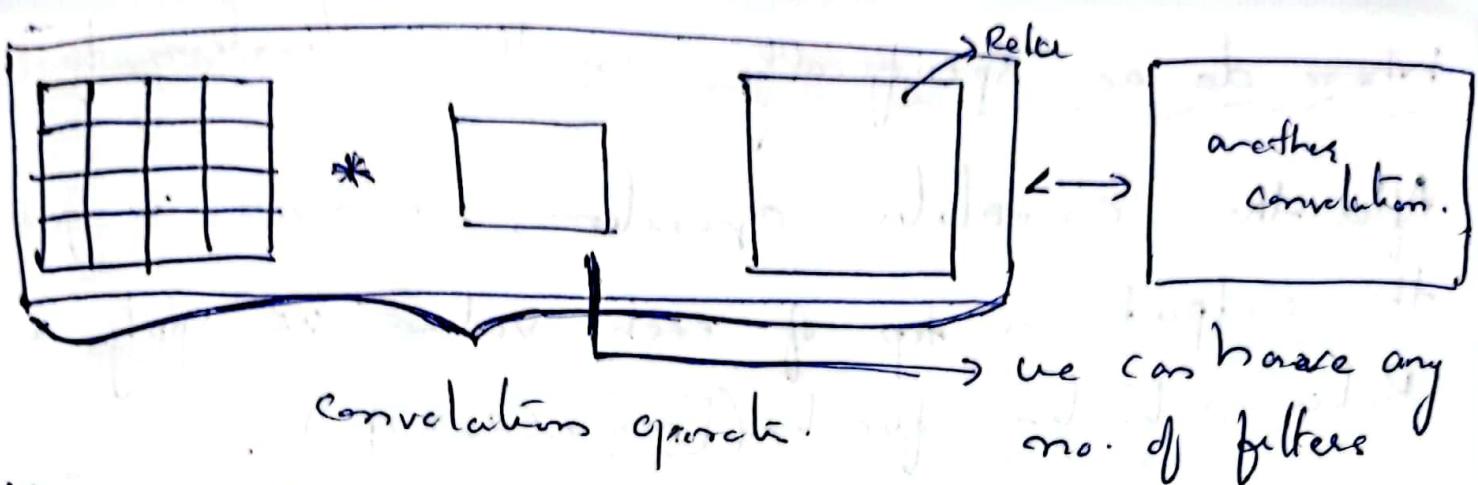
During the back propagation the derivative can be found out. Since we can find out the derivative we can do the back propagation and hence update the filter.

Stride

It is the number of pixels by which a filter moves across the input image during the convolution operations. It is a crucial parameter that determines how many units the filter shifts at each step, either horizontally, vertically, or both.

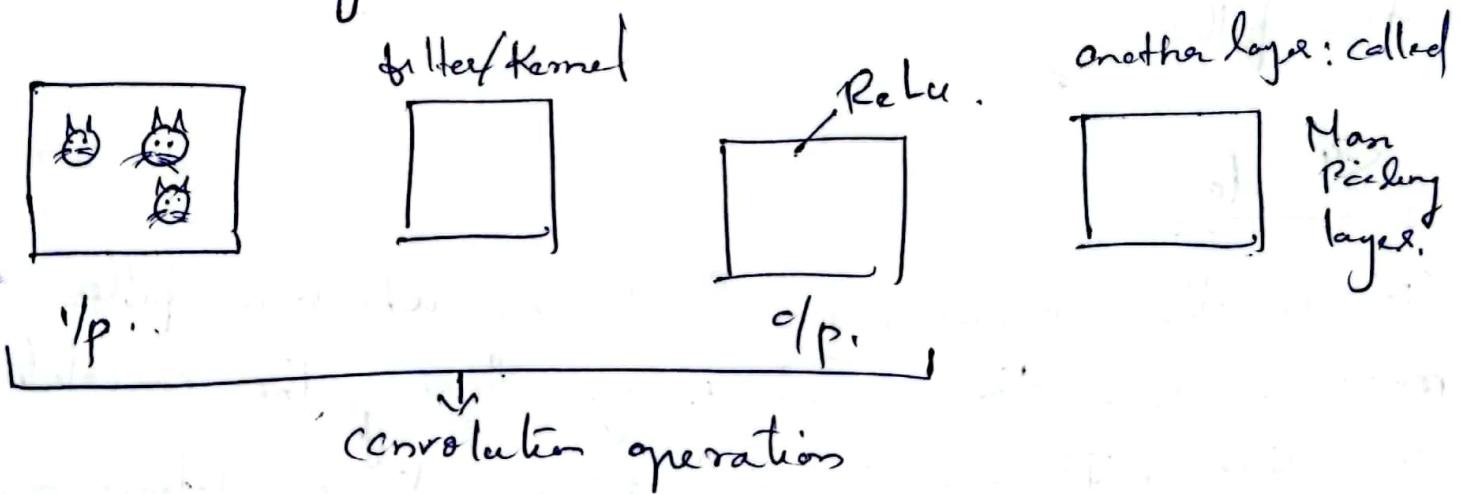
if $\text{Stride} \neq 1$, then - the updated formula

$$\frac{n+2P-F+1}{S}$$



Main aim of convolution operations is we need to learn ~~filters~~ and update based on the I/P images

Max-pooling



Location invariant

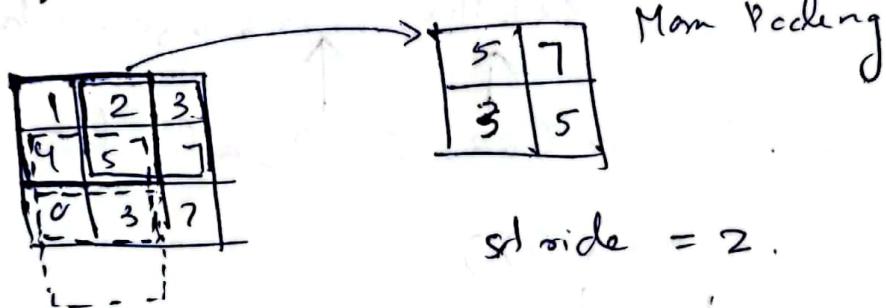
Let's say there are multiple ~~images~~ objects in the images, the cnns till it goes ahead handles the next neural n/w it should be able to extract the more and more clear info.

So we use max pooling.

Type of Max Pooling

1. Avg Pooling
2. Min Pooling
3. Max Pooling

We are going to place the max Pooling as top of the o/p, and the o/p cell only be focusing on the biggest number



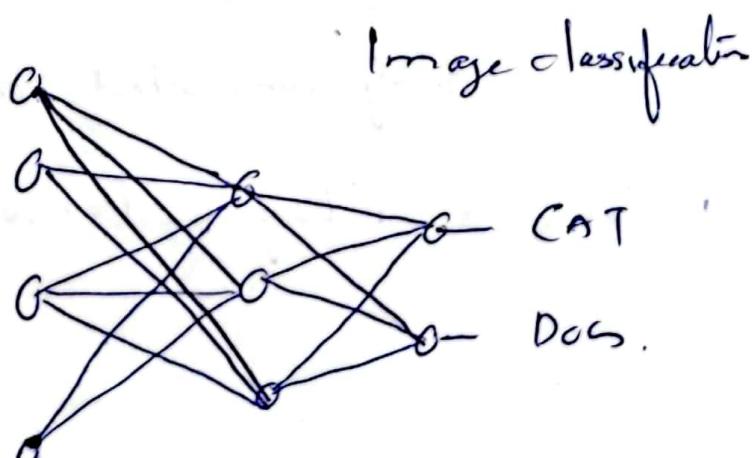
Based on different Problem statement, we can use min pooling, max pooling, avg pooling.

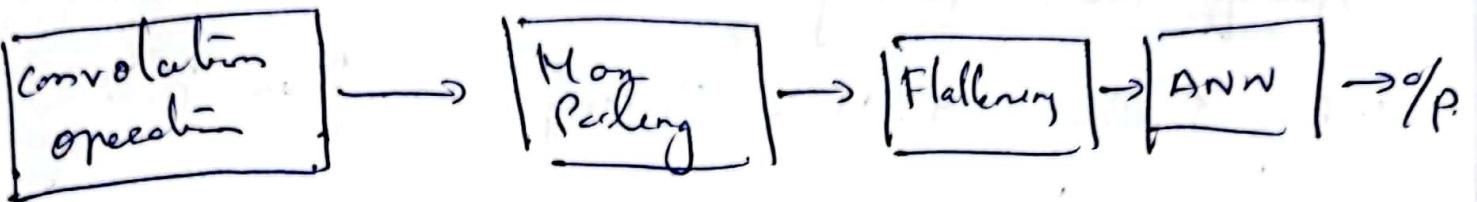
Flattening layer (like ANN Dense layer)

The pooling o/p will be elongated in this stage.

5
7
3
5
7
9
2
1

this will
be i/p
to the
dense
layer.





Max pooling → refers to extracting more information from the images.

32 filters 3×3 size each.

`model.add(layers.Conv2D(32, (3, 3), activation = 'relu'))`
`input_shape = (32, 32, 3))`

\downarrow \downarrow
 32×32 3 channels
 (1/p image size)

? How do you decide the number of filters?

~~Flatton layer~~:

`model.add(layers.Flatten())`

~~ANN~~: `model.add(layers.Dense(64, activation = 'relu'))`

fully connected neural net with 64 dense neurons

~~O/P layer~~:

`model.add(layers.Dense(10))`

Total number of o/p