

The main aim of Deep learning is to mimic human Brain.

Why Deep Learning is becoming popular?

① 2005 → FB, Insta, WhatsApp, LinkedIn

Data → Exponentially ↑↑↑

2008 → {Big Data} → Efficiently

2013 → Company had huge amount of Data
↓ {AI → popular}

Seamless Experience for Products
↓
(Company will generate Revenue, Better Decision)

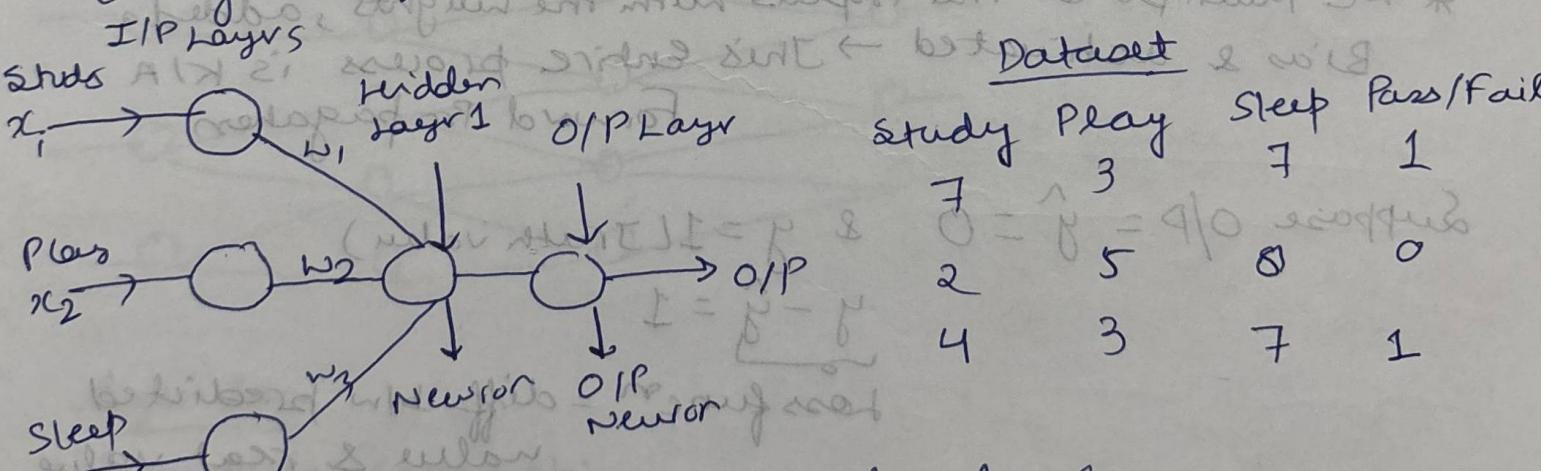
② Hardware advancement (Nvidia)

↪ GPUs → Training the model

Gpus cost is decreasing

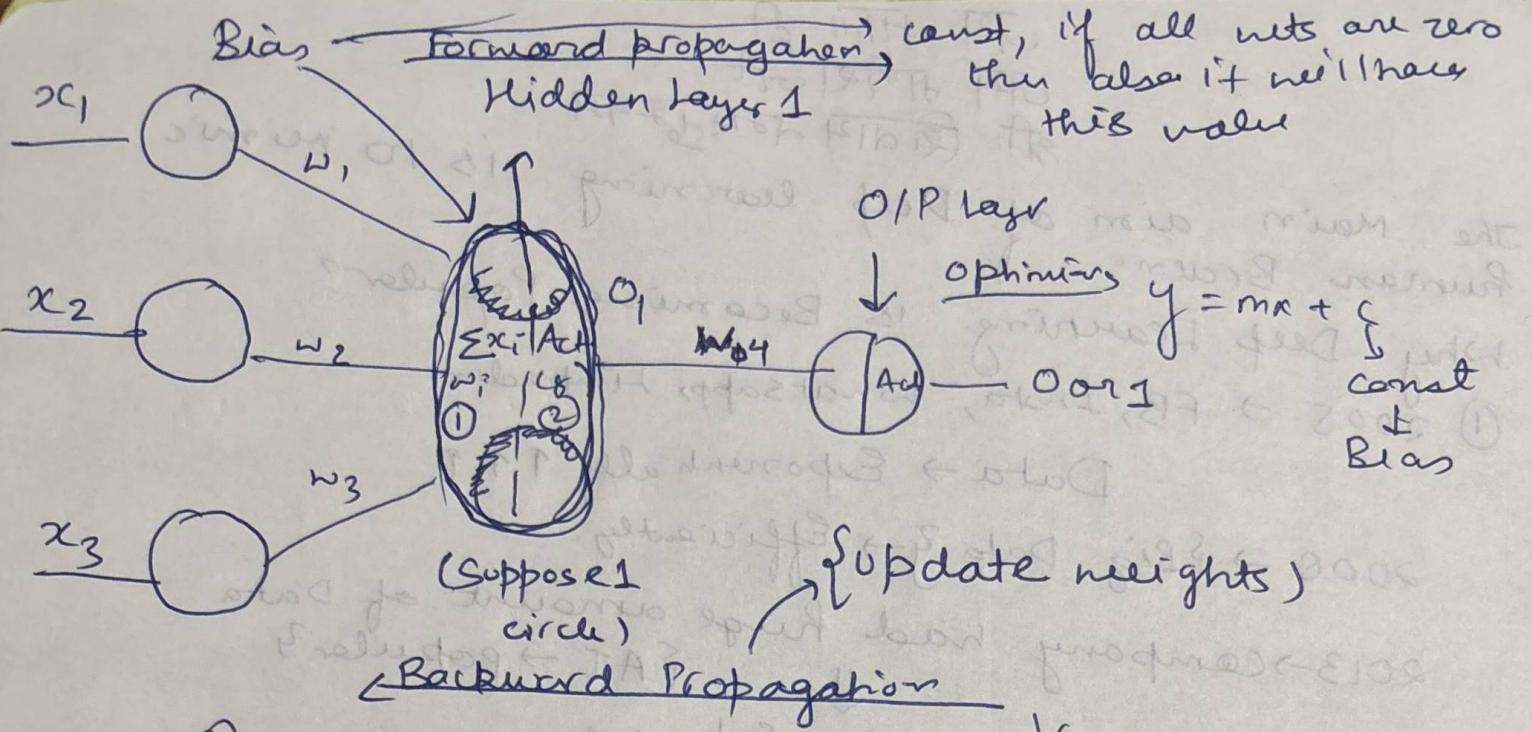
* Perceptron:

{Single Layered Neural Network}



* There can be more than 1 hidden layer

w_1, w_2, w_3 → weights basically tell how much a neuron should get activated



Backward Propagation

$$\textcircled{1} \quad \sum x_i w_i = x_1 w_1 + x_2 w_2 + x_3 w_3 + b \quad y = \beta_0 + \beta x$$

$$w^T x \quad y = \beta^T x$$

$$b \rightarrow \text{bias}$$

② Activation func.

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

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

* we multiplied the inputs with the weights, added a bias & activated \rightarrow This entire process is KIA
Forward Propagation

Suppose O/P = $\hat{y} = 0$ & $y = 1$ (Truth value)

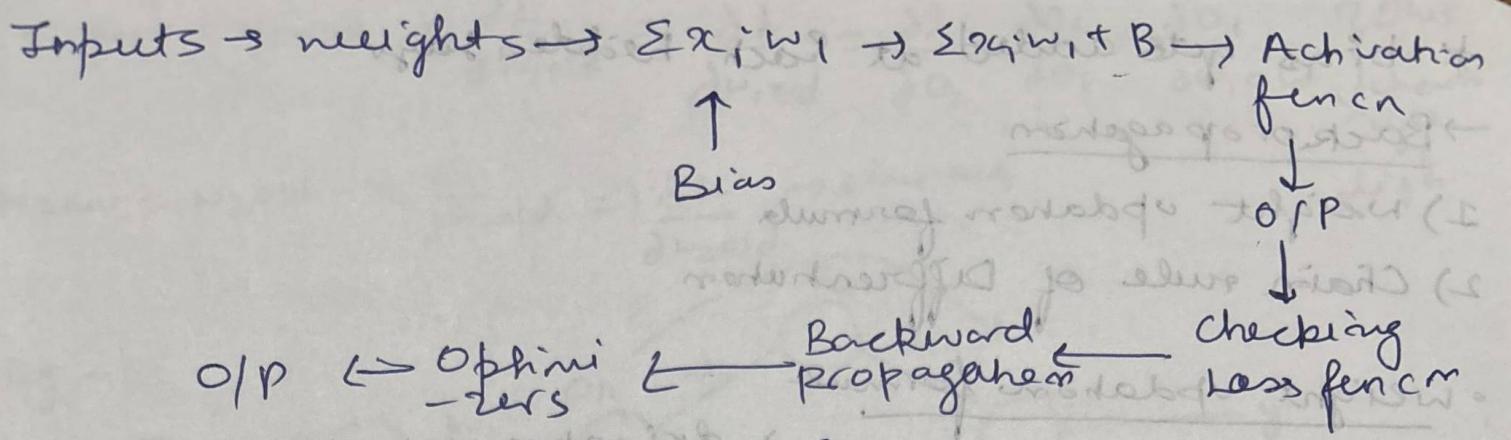
$$y - \hat{y} = 1$$

Loss function = diff b/w predicted value & real value

The main aim should be to minimize difference and make it near to zero

→ The main aim of Backward propagation is to update weights & hence reducing loss fun.

→ Optimizers - main aim is to make sure that each & every weight will be updated while we are Back - Propagation

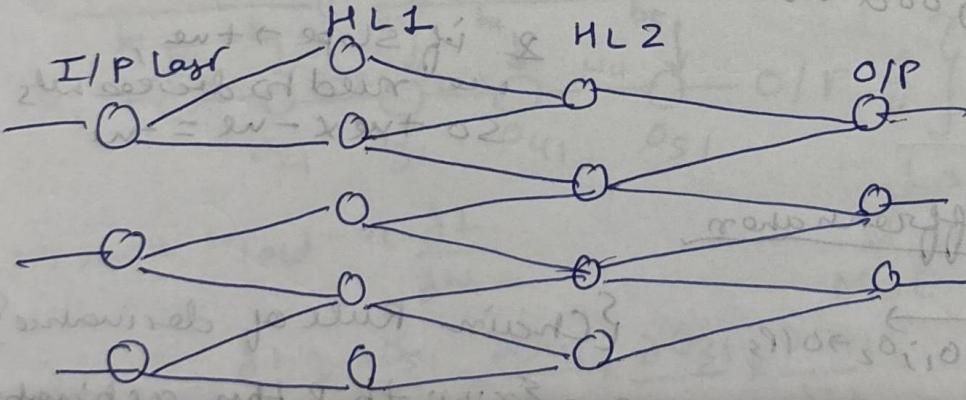


\rightarrow Gradient descent is an eg. of optimiser

- ① I/P layers
- ② Weights
- ③ Bias
- ④ Activation function
- ⑤ Loss function $\{y - y'\}^2$
- ⑥ Optimizers
- ⑦ Update the weight

} forward propagation

multi-layer Neural Network



Agendas -

- 1) Forward Propagation
- 2) Chain rule of Derivative
- 3) Vanishing Gradient Problem
- 4) Loss function

* bias gets added at each hidden layer
 → Backpropagation

- 1) weight update formula
- 2) Chain rule of differentiation

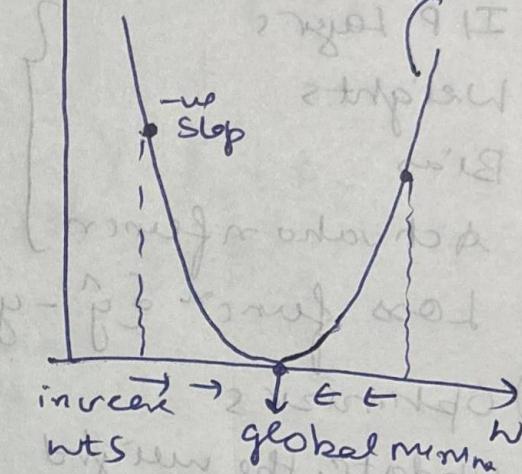
- weight update formula

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}} \quad \begin{matrix} \rightarrow \text{derivative of Loss} \\ \downarrow \end{matrix}$$

$$\frac{\partial h}{\partial w_{\text{old}}} \quad \begin{matrix} \rightarrow \text{Slope} \\ \downarrow \end{matrix}$$

Loss function

Gradient descent

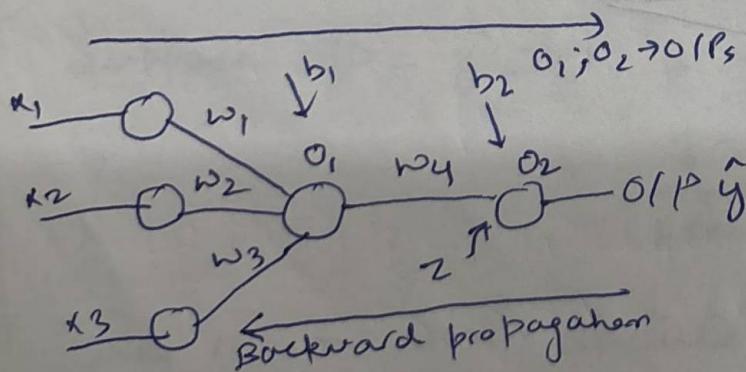


* Learning rate would be better if it is a small number as it will slowly converge to global minima, if it would be larger it would take larger steps & hence might never reach global minima.

$$\eta = 0.001 \text{ is a good value}$$

So $x - w$ must be true
 2 if slope $\rightarrow +ve$, we need to decrease wts
 So $+ve x - ve = -ve$

- Chain rule of differentiation



For w_4 ,

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

{Chain Rule of derivative}

$o_1 \rightarrow \sum_i w_i + b$ & then activation function applied then $o_1/P = o_1$

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

$z = \text{activation function}$

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

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

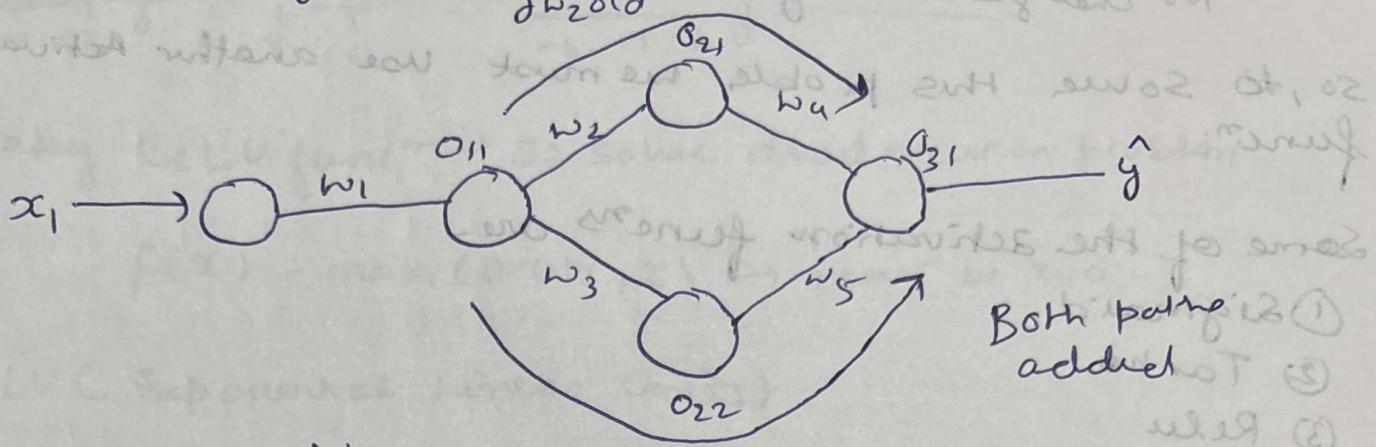
$$w_{1,\text{new}} = w_{1,\text{old}} - \eta \frac{\partial L}{\partial w_{1,\text{old}}} \Rightarrow \frac{\partial L}{\partial w_{1,\text{old}}} = \frac{\partial L}{\partial o_2} \times \frac{\partial o_2}{\partial o_1} \times \frac{\partial o_1}{\partial w_{1,\text{old}}}$$

Thus
it's
chain
rule

$$w_{2,\text{new}} = w_{2,\text{old}} - \eta \frac{\partial L}{\partial w_{2,\text{old}}} \Rightarrow \frac{\partial L}{\partial w_{2,\text{old}}} \rightarrow \text{similar}$$

diff

eg -

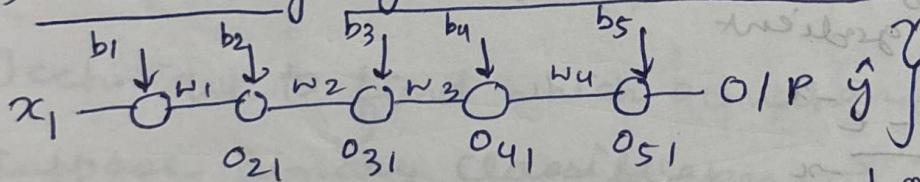


$$w_{1,\text{new}} = w_{1,\text{old}} - \eta \frac{\partial L}{\partial w_{1,\text{old}}}$$

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

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

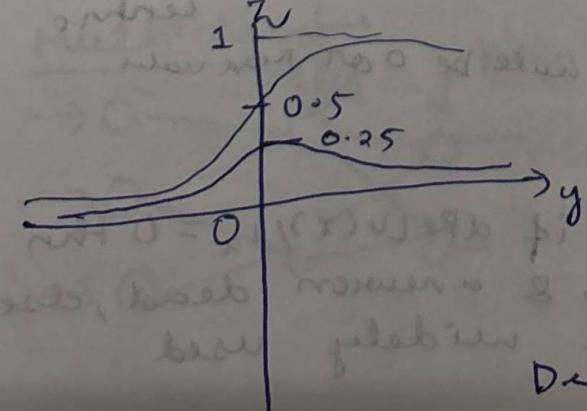
→ Vanishing Gradient Problem



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

$$\frac{\partial L}{\partial w_{1,\text{new}}} = \frac{\partial L}{\partial o_{51}} \times \frac{\partial o_{51}}{\partial o_{41}} \times \frac{\partial o_{41}}{\partial o_{31}} \times \frac{\partial o_{31}}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial w_1}$$

graph of Sigmoid



eg $\rightarrow o_{51} = \sigma([o_{41} \times w_4]) + b$
so, the values multiplied in
chain will give
very small values

$$2 = \frac{1}{1 + e^{-x}}$$

The derivative of Sigmoid weight
is $0 < \sigma'(y) < 0.25$

$$\text{Derivatives } 0 \leq \sigma(y) \leq 0.25$$

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

$w_{\text{new}} \approx w_{\text{old}} \Rightarrow$ vanishing gradient problem
 ↓
 No change in weights

So, to solve this problem we must use another activation funcⁿ.

Some of the activation funcs are-

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

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

→ smooth gradient, prevents jumbies

→ funcⁿ having zero centred curve update weights easily

→ time complexity is high

→ prone to vanishing gradient

• Tanh $\Rightarrow \tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}}$
 [-1 to 1]

derivative (0 to 1)

• still prone to vanishing gradient in case of Deep neural network

• zero centred

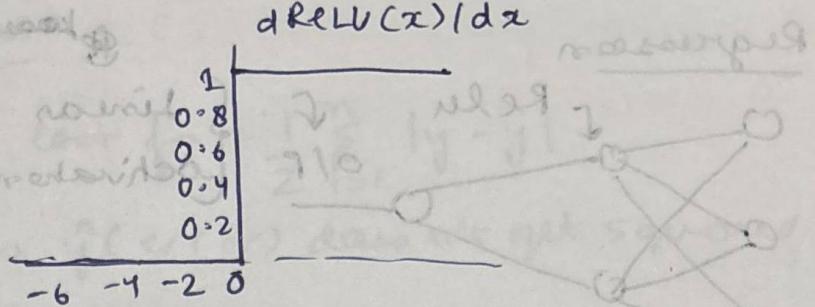
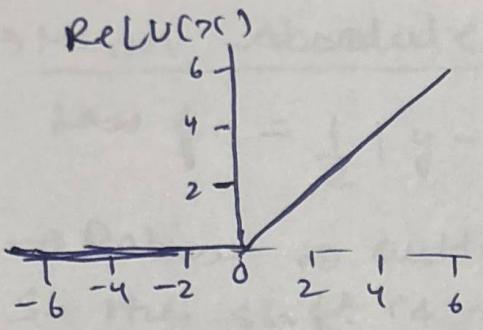
• hyperbolic tanh funcⁿ

• ReLU (Rectified Linear Unit) → quickest, non-zero centric

$$\text{ReLU} = \max(0, x) \rightarrow \text{either it will be } 0 \text{ or Max value}$$

$$d\text{ReLU}(x)/dx = \text{either } 0 \text{ or } 1$$

Only one disadvantage is that if $d\text{ReLU}(x)/dx = 0$ then it will make again $w_{\text{old}} = w_{\text{new}}$ & a neuron dead, else it is the best & the most widely used activation funcⁿ.

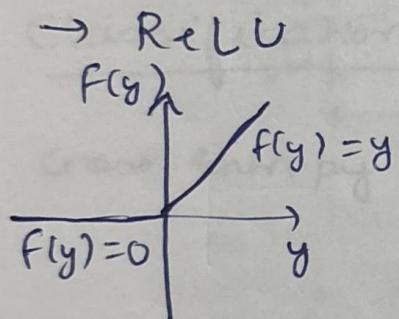


- Leaky ReLU function → to solve dead neuron problem

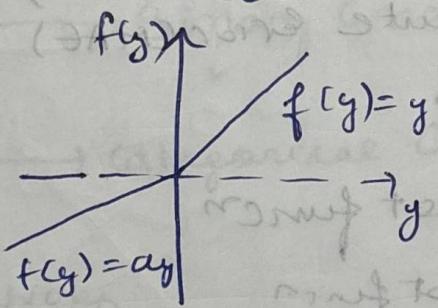
$$f(x) = \max(0.01x, x) \rightarrow \text{never be zero}$$

- ELU (Exponential Linear Units)

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



v/s (→ PReLU)

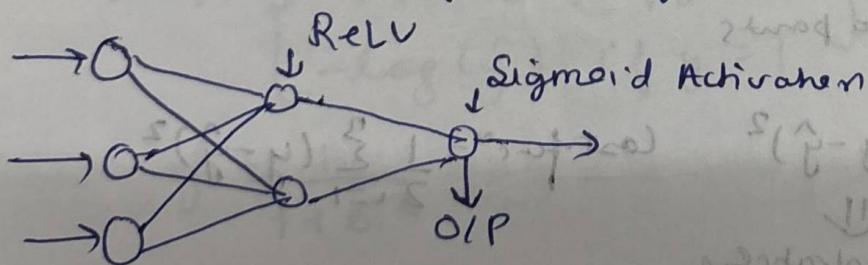


Technique to find which activation funcn to used?

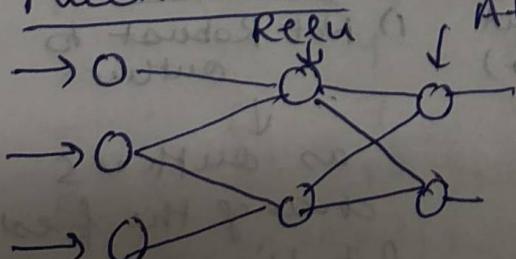
Suppose Binary Classification →

In these at O/P, Sigmoid Activation must be used

& ReLU at Hidden layers or PReLU or ELU



Multiclass →

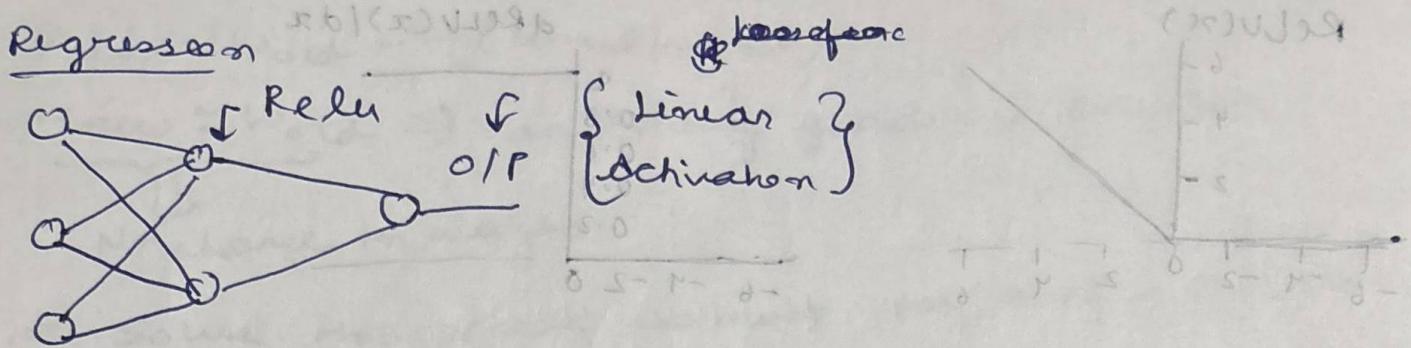


At O/P

use Softmax Activation

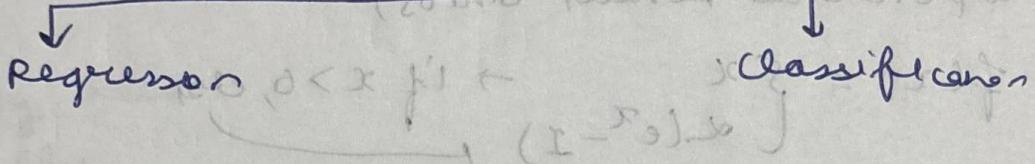
level 1 program to convert into integer value

loop continues



(*) Loss funcs

Deep Learning (Artificial Neural Network)



For Regression

- Mean Squared Error (MSE)
- Mean Absolute Error (MAE)
- Huber Loss

→ Diff b/w Loss & cost func

Loss func

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

cost func

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

→ it takes a single data point

→ it takes a set of data points

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

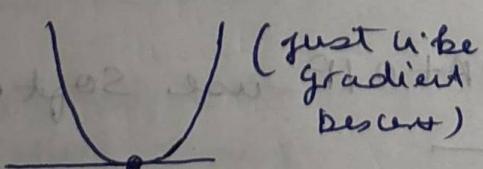
$$\text{cost func} = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

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

Quadratic Eqn
↓ graph

Advantages of Quadratic Eqn:

- ① Differentiable
- ② It has only 1 local or global minima
- ③ converges faster



Disadvantages

- ① Not Robust to outliers
as outliers change the fit line as error $y - \hat{y}$ gets squared

→ Mean Absolute Error -

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

→ Robust to outliers as \hat{y} (error) doesn't get squared

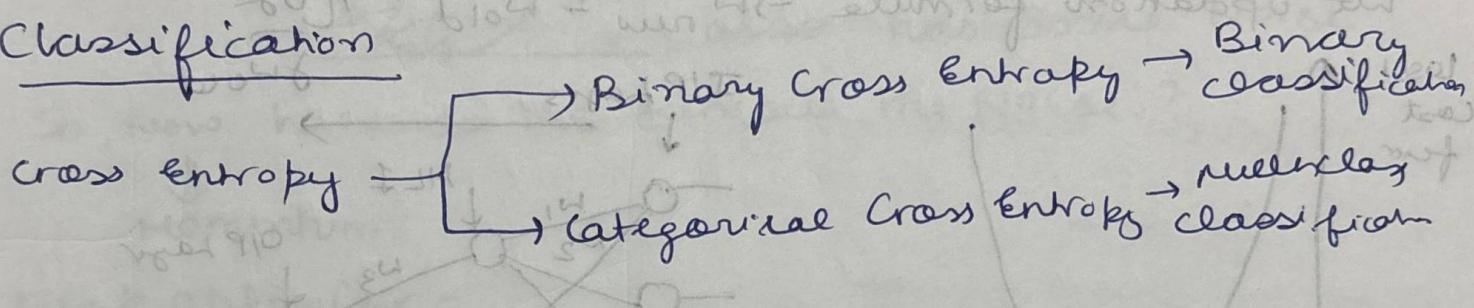
So the shift is minor

→ Time consuming

→ Huber Loss → combination of MSE & MAE & is specifically used when outliers are not there

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

Classification



i) Binary Cross Entropy

$$\text{Loss} = -y \times \log(\hat{y}) - (1-y) \times \log(1-\hat{y}) \Rightarrow \text{Loss is higher if } y \neq \hat{y}$$

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

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

ii) Categorical Cross Entropy -

	f_1	f_2	f_3	O/P	Good	Bad	Neutral
1					0	0	1
2		3	4	good	0	0	1
5	6	7		Bad	0	0	1
8	g	h	l	Neutral	0	0	1

$$L(x_i, y_i) = - \sum_{j=1}^C y_{ij} \ln(\hat{y}_{ij}) \quad y_i = [y_{i1}, y_{i2}, \dots, y_{iC}]$$

$$\hat{y}_{ij} = \text{Softmax Activation} \quad y_{ij} = \begin{cases} 1 & \text{if Element is in class} \\ 0 & \text{else} \end{cases}$$

$$\text{softmax func} \Rightarrow \sigma(z) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

Conclusion

ReLU, Softmax \Rightarrow Multi-class \rightarrow Categorical Cross Entropy

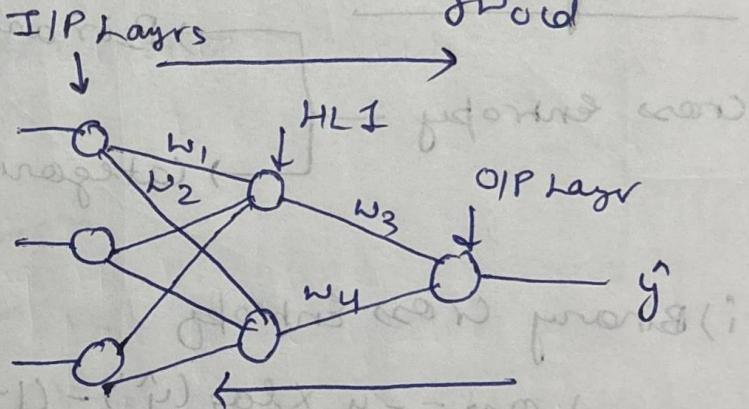
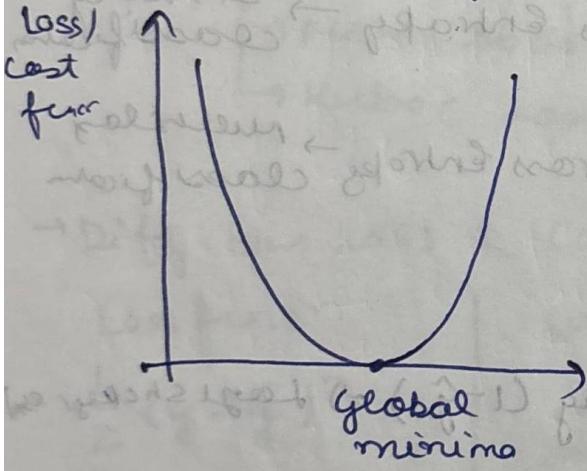
ReLU, Sigmoid \Rightarrow Binary \rightarrow Binary Cross Entropy

Linear Regression - MSE, MAE, Huber Loss

Optimizers -

① Gradient Descent:

$$\text{wt update formula} \rightarrow w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial}{\partial w_{\text{old}}}$$



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

MSE

Epoch \rightarrow 1 forward propagation & 1 backward propagation forms 1 Epoch

Generally we pass million Records i.e. $n = 1 \text{ million}$

Disadvantage of Gradient Descent -
i) Resource Extensive { Huge RAM }

1 million Record
in each Epoch

② Stochastic Gradient Descent \rightarrow 1 record in 1 Epoch

* RAM lit. but convergence will be very slow.

1 record $\rightarrow \hat{y} \rightarrow$ iteration 1

update w.

* High Time Complexity

but this will have million epochs

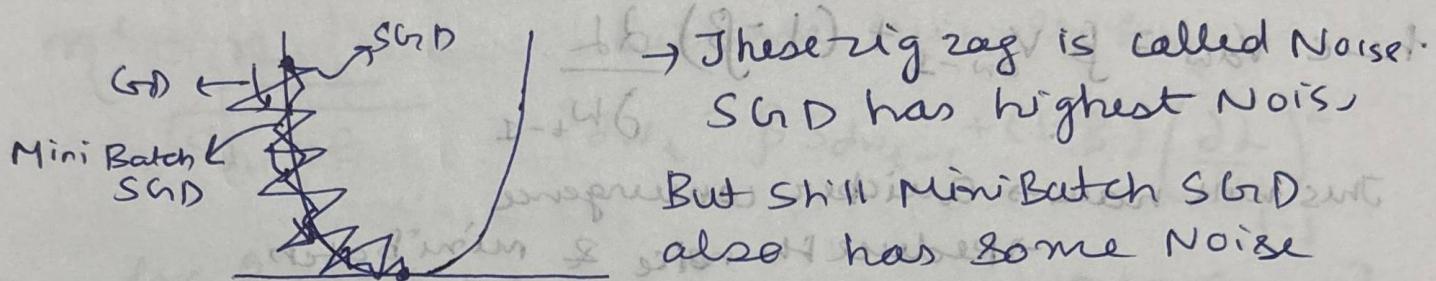
③ Mini Batch SGD (Stochastic Gradient Descent)

→ we will set batch size. Suppose \rightarrow batch-size = 1000

Epoch - 1 $\xrightarrow{\text{1000}} \text{Iteration 1}$

so for million records it is 1000 iterations only

- * Resource Intensive
- * Convergence will be better
- * Time complexity will Improve



So, now how do we remove the noise?

\downarrow
Momentum

Next Optimizer -

④ Mini Batch SGD with Momentum -

→ Exponential Weighted Average } Also used in Time Series
 " Moving " }
 " $b_{\text{new}} = b_{\text{old}} - \eta \frac{dL}{db_{\text{old}}}$

As we knew, $w_{\text{new}} = w_{\text{old}} - \eta \frac{dL}{dw_{\text{old}}}$

$t \rightarrow \text{current time}$
 $t-1 \rightarrow \text{previous time}$

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

time $t_1, t_2, t_3, t_4, \dots, t_n$
 value $a_1, a_2, a_3, a_4, \dots, a_n$

$\beta \rightarrow \text{hyperparameter}$

$$v_{t_1} = a_1$$

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

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

$$= 0.95 v_{t_1} + 0.05 a_2$$

$V_{t_2} = 0.95 \times V_{t_1} + 0.05 a_2$

taking $\beta = 0.95$ we are giving more imp. to V_{t_1}
we are removing noise & smoothening curve

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

So, Exponential Weighted Avg: ~~monotonic~~ smooth *
 $(w_t = w_{t-1} - \eta V_{dw})$ ~~used old bias expression~~ *
~~weights new pixel frame~~ ~~smooth~~ *

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

This will have \rightarrow Quicker convergence

\rightarrow Reduce the noise & minibatch

⑤ Adagrad \rightarrow Adaptive Gradient Descent

$$\eta = \text{fixed}$$

* It is slow & fixed

\rightarrow can't we do something that initial learning rate becomes fast when far from global minima,
& when comes closer then becomes slow?

\rightarrow for this is Adagrad

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

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

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

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

so α_t will always increase

$\therefore \eta'$ will decrease when we reach global minima

small number to avoid divide by zero

(for cases $\alpha_t = 0$)

Sometimes η_t can be huge number which may grow in units. So, we will use another optimiser to solve this problem.

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

huge \leftrightarrow no

$\Rightarrow \eta' \approx \eta_{t-1}$

6) AdaDelta & RMSprop -

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

Suppose initially $S_{dw} = 0$

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

$$\text{for } \beta = 0.95; S_{dw_t} = 0.95 S_{dw_{t-1}} + 0.05 \left(\frac{\partial L}{\partial w_{t-1}} \right)^2$$

→ concept of Exponential weighted average is used

- Solves Smoothering problem
- Assures that Learning Rate becomes Adaptive

7) Adam Optimiser

Momentum + RMSprop → Best Optimizer

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

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

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

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

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

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

ANN Practical Implementation

Churn - Modeling.csv → Dataset

```
!pip install tensorflow-gpu
```

```
import tensorflow as tf
```

```
print(tf.__version__)
```

```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
import pandas as pd
```

```
dataset = pd.read_csv('Churn - Modeling.csv')
```

```
dataset.head()
```

```
X = dataset.iloc[:, 3:13] → 3 to 12 columns
```

```
Y = dataset.iloc[:, 13]
```

```
X.head
```

```
geo = pd.get_dummies(X['geography'], drop_first=True)
```

```
gender = pd.get_dummies(X['Gender'], drop_first=True)
```

```
X = X.drop(['geography', 'Gender'], axis=1)
```

```
X = pd.concat([X, geo, gender], axis=1)
```

```
from sklearn.model_selection import train_test_split
```

```
X_train, X_test, y_train, y_test = train_test_split
```

```
(X, y, test_size=0.2,  
random_state=0)
```

```
# ANN {Feature Scaling}
```

→ For which all algorithms feature scaling is required.

→ ANN, LR, Log. R, k-NN, k-means etc. for distance based or gradient descent

feature scaling
is required

```

from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)

```

Part 2 Now let's create ANN

Tensorflow => open sourced by Google

PyTorch => by facebook

In version < 2.0 Tensorflow -> keras
 In version ≥ 2.0 Tensorflow + keras -> wrapper + keras

```

from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense
from tensorflow.keras.layers import LeakyReLU,
from tensorflow.keras.layers import PReLU, ELU, ReLU
from tensorflow.keras.layers import Dropout

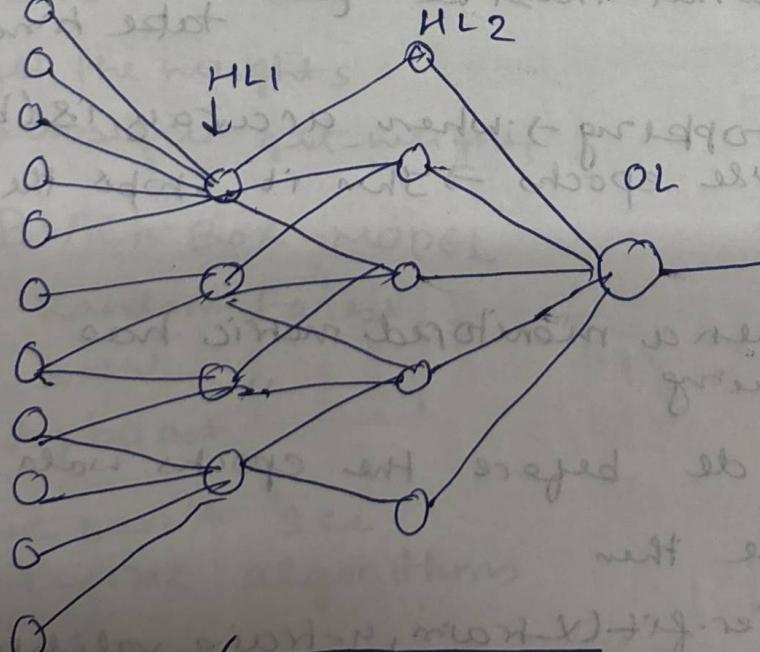
```

① Sequential - Neural block which can do FP & BP
 ② Dense - we are able to create Neurons & its layers

③ Activation - To use activation function

④ Dropout

IL



Sometimes, Neural network may lead to overfitting i.e. my training accuracy is high but test accuracy goes down, so we use Dropout layer - which acts as Regularization parameter, which deactivates neurons, & drops out layers.

```
# Let's initialize ANN  
classifier = Sequential()  
  
# Adding the input layer  
classifier.add(Dense(units=11, activation='relu'))  
  
# adding 1st HL  
classifier.add(Dense(units=7, activation='relu'))  
  
# 3rd HL  
classifier.add(Dense(units=6, activation='relu'))  
  
# Adding O/P Layer  
classifier.add(Dense(1, activation='sigmoid'))  
  
classifier.compile(optimizer='adam', loss='binary_crossentropy',  
metrics=['accuracy'])  
  
import tensorflow  
tensorflow.keras.optimizers.Adam(learning_rate=0.01)
```

```
model_history = classifier.fit(X_train, y_train,  
(taking only 67% of data) validation_split=0.33,  
batch_size=10, epochs=1000)  
↓  
this will
```

So how to know that what must be the no. of epochs?
← take time
we must use Early Stopping → when accuracy is (becoming approx constant while epochs → then it stops the epochs at that time.

```
# Stop training when a monitored metric has stopped improving  
we must paste code before the epochs code.
```

Early Stopping - code then

```
model_history = classifier.fit(X_train, y_train, validation_split=0.33,  
batch_size=10, epochs=1000, callbacks=[EarlyStopping])
```

model.history.history.keys()

tells what all factors we have focussed on.

Summarize history for accuracy

```

plt.plot(model.history.history['accuracy'])
plt.plot(model.history.history['val_accuracy'])
plt.title('model accuracy')
plt.ylabel('accuracy')
plt.xlabel('epoch')
plt.legend(['train', 'test'], loc='upper left')
plt.show()

```

Making predictions & Evaluate Model

Predicting the test set results

```

y_pred = classifier.predict(x_test)
y_pred = (y_pred >= 0.5)

```

make confusion matrix

```

from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
cm

```

Calculate accuracy

```

from sklearn.metrics import accuracy_score
score = accuracy_score(y_pred, y_test)
score

```

get the weights

```

classifier.get_weights()

```

→ BLACK BOX MODEL

- Random Forest
- ANN
- Xg boost

- we can't see internal algorithms
- here we can & internal algorithm

* To look into the Black box Models, nowadays Explainable AI is being discussed.

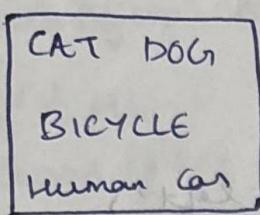
CAT PAPER
BLACK

Convolution Neural Network (CNN)

ANN \rightarrow Theoretical, practical

CNN \rightarrow Image, video frames

CNN VS Human Brain



CEREBRAL CORTEX

VISUAL CORTEX

many layers like,

$V_1 \rightarrow$ Meaning objects } features

$V_2 \rightarrow$ Animals }

$V_3 \rightarrow$ Map Environment }

$V_4 \rightarrow$ Colors }

$V_5 \rightarrow$ Shapes }

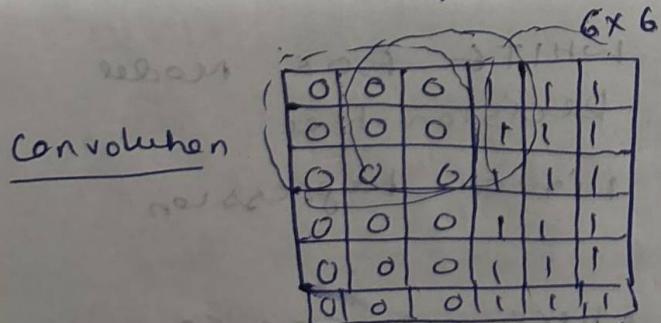
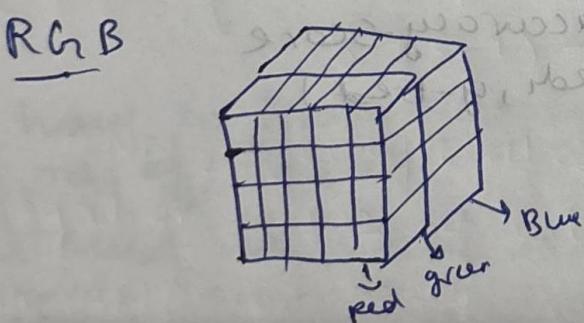
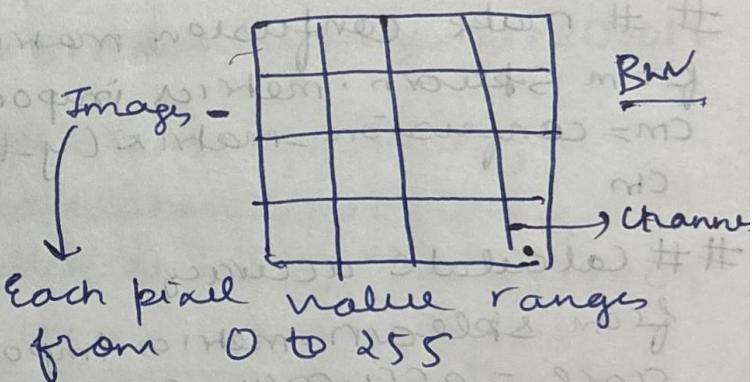
$V_6 \rightarrow$ Lines }

$V_7 \rightarrow$ O/P }

CNN is divided into following steps -

i) Convolution

Image \rightarrow Black White
 \rightarrow RGB



If we want to range each pixel value b/w 0 to 1

each pixel value must divide each pixel value by 255 \rightarrow This is called as min / Max scaling

passed way

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

Suppose this is horizontal edge filter

O/P \rightarrow 4x4

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

Suppose that filter is placed at 1st quarter of image.
 → we will multiply each cell value to other & then add it. $(0+0+0+0+0) = 0$
 & then 1 step to right (suppose stride = 1) if,
 Stride = 2 then 2 steps
 ② $0+0+1+0+0+0+0+0 = 1 = 0$
 ③ $0+2+1+0+0+0 = 2 = 0$
 ④ $\rightarrow 0$

O/P =	0	0	0	0
	0	0	0	0
	0	0	0	0
	0	0	0	0
	0	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

This operation is termed as convolution.

Suppose O/P = $\begin{bmatrix} 0, -4, -4, 0 \\ 0, -4, -4, 0 \\ 0, -4, -4, 0 \\ 0, -4, -4, 0 \end{bmatrix}$ & the filter applied is vertical edge filter

After min max scaling highest value = 255
 lowest value = 0

255, 0, 0, 255
255, 0, 0, 255
255, 0, 0, 255
255, 0, 0, 255

255 → white color

& 0 → black color

White	Black	White
White	Black	White

vertical
edge filter

→ Similarly many types of filters → which can take out specific information from a specific image.

6 × 6 image → 3 × 3 filter

$$n=6$$

$$f=3$$

→ O/P =	1	2	3	4	5	6	7	8	9
---------	---	---	---	---	---	---	---	---	---

$$\boxed{O/P = n - f + 1}$$

In case of $6 \times 6 \rightarrow 3 \times 3$ filter $\rightarrow O/P 4 \times 4$
 So as we see our image size is reducing
 & this should not happen that means we are losing
 some kind of information - In order to prevent
 that we can do padding i.e. building a
 specific compound around the image
 process of applying a layer or
 the top of 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

Now we can fill that
 by anything may be
 zero or nearest value
 ↓
 zero padding

Now $n=8$ so $n-3+1=6$

So updated formula

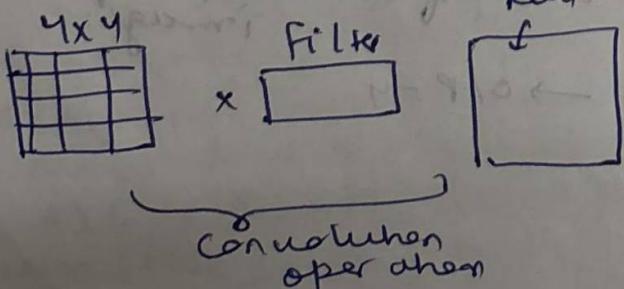
$$n + 2p - f + 1$$

$228 = \text{euler total}$
 $6 = \text{euler total}$

$$6 + 2(1) - 3 + 1 = 6$$

In ANN, we update weights, here in back propagation we update filters based on IP image
 → On O/P on each value we apply ReLU Activation function

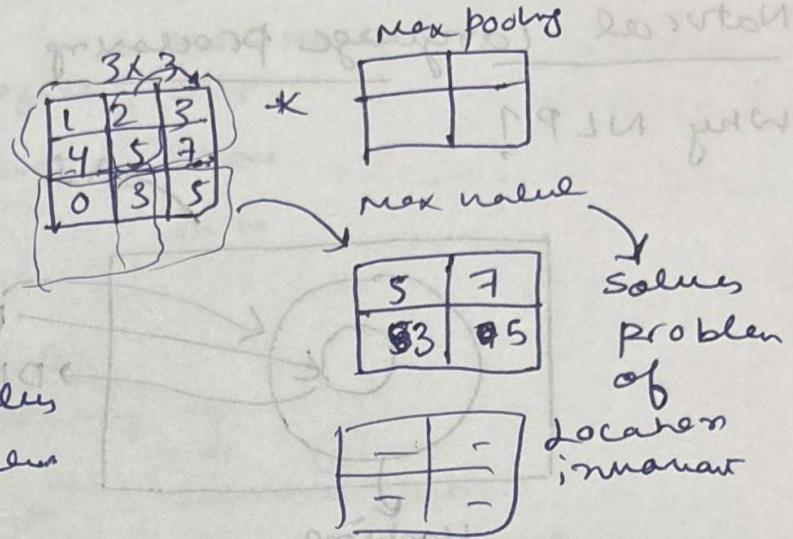
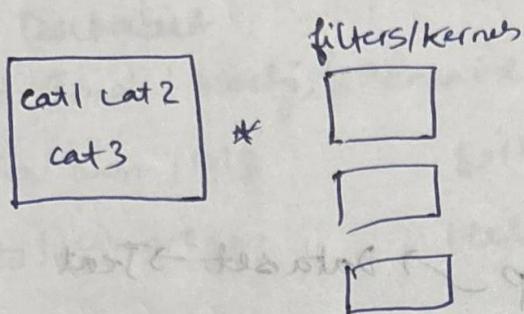
Stride also plays an important role so, updated formula $\Rightarrow n + 2p - f + 1$



$$(f-f) \times r = 910$$

② Max pooling -
pooling is of 3 types

Aug pool way
min pooling
Max pooling



→ Flatenning layer -

ANN Dense layer

