

## DEEP LEARNING

We need to provide a good amount of data to create a model (ML model)

ex: identifying a dog - length of legs, nose size, etc

We need to populate all Features and create the model. Even after this, there is a chance for Failure

→ Rule based system: (Around 1970)

Rules  
+  
Data ] → MODEL → gives output.

→ ML system: (1970 - 2010)

Input data  
+  
Output data ] → MODEL → creates Rules

Need to provide structured / tabular data (rows & columns)

→ DL Model:

- No need to provide structured data
- Finds out Features on itself
- It is a ML model which works on unstructured data and derives the Features on itself

Image → Finds basic primitive Features → uses these Features to form more complex Features

→ / → □ → ... → Final dog picture

many layers are present - SEQUENTIAL LAYER ARCHITECTURE

(Good inputs, good outputs)

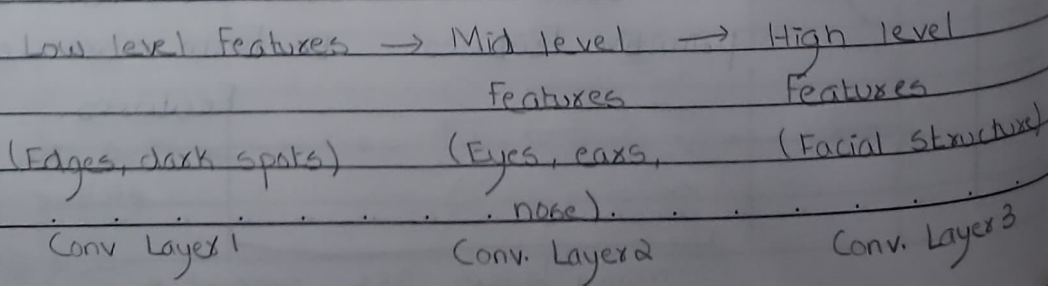
We need to provide the data as accurately as possible - otherwise system gives bad outputs (GIGO - Garbage In Garbage Out)

Learning happens in a sequential manner.

Every layer is responsible for different type of features. More layers means more deeper architecture.

DL model learns the representations from raw data.

- Supervised Learning : labelled data.
- Unsupervised Learning (clustered approach) : Unlabelled  
Find similarity in data and cluster it
- DL model : Supervised learning
- DL is mostly used in unstructured data
- It does not have much applications in structured data (not required).
- DL is mostly designed for unstructured data.  
In case of structured data, the features are already present - does not make sense to use DL. DL model can be used, but not required - other simpler models can be used





Mostly in DL, the first few layers (lower level features) are the same. As we go deeper the features will be different.

More layers / deeper architecture  $\rightarrow$  Better accuracy

$\rightarrow$  Difference b/w DL and ML:

- ① Data dependency : DL requires more data
- ② Hardware dependency : requires specialised hardware like GPUs.
- ③ Training Time : higher for DL.
- ④ Feature selection : Features are automatically extracted by layers.
- ⑤ Interpretability : Very less or not present in DL model.

TENSOR  
FLOW

sklearn helps use numpy functions more efficiently

numpy  $\rightarrow$  sklearn  $\rightarrow$  used in CPUs

numpy  $\rightarrow$  TensorFlow  $\rightarrow$  used in GPUs.

The background of sklearn and TensorFlow is numpy functions (arrays).

TensorFlow Functions are supported by both CPUs and GPUs. Sklearn supports only CPUs.

Mostly in DL, the first few layers (lower level features) are the same. As we go deeper the features will be different.

More layers / deeper architecture  $\rightarrow$  Better accuracy

$\rightarrow$  Difference b/w DL and ML:

- ① Data dependency : DL requires more data
- ② Hardware dependency : requires specialised hardware like GPUs.
- ③ Training Time : higher for DL.
- ④ Feature selection : Features are automatically extracted by layers.
- ⑤ Interpretability : Very less or not present in DL model.

TENSOR  
FLOW

sklearn helps use numpy functions more efficiently

numpy  $\rightarrow$  sklearn  $\rightarrow$  used in CPUs.

numpy  $\rightarrow$  TensorFlow  $\rightarrow$  used in GPUs.

The background of sklearn and TensorFlow is numpy functions (arrays).

TensorFlow functions are supported by both CPUs and GPUs. Sklearn supports only CPUs.

Keras - can run on top of TensorFlow.

TensorFlow consists of tensors (similar to arrays)

Constant: (cannot be changed) and variable tensors

Arithmetic operations - Tensors

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + 20$$

gets converted to

$$\begin{bmatrix} 20 & 20 & 20 \\ 20 & 20 & 20 \end{bmatrix}$$

$$(2, 3, 2) + (2, 3)$$

Need to change dimension of (2, 3)

(can be added  $\begin{matrix} \downarrow \\ (1, 2, 3) \end{matrix}$ )

on either

side)

or

$(2, 3, 1)$

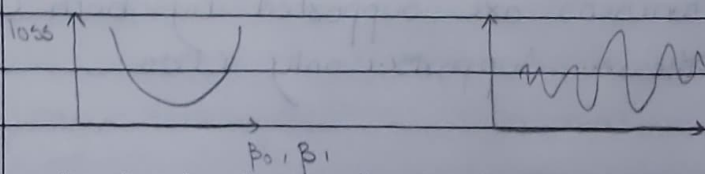
Simple Linear Regression

$$y = \beta_1 x + \beta_0$$

$\beta_0, \beta_1 \rightarrow$  weights

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

Loss Function helps get best fit line - From parabola of curve





For diff values of  $\beta_0, \beta_1$  compute  $\beta_{new}$   
at that  $\beta_{new}$ , check loss.

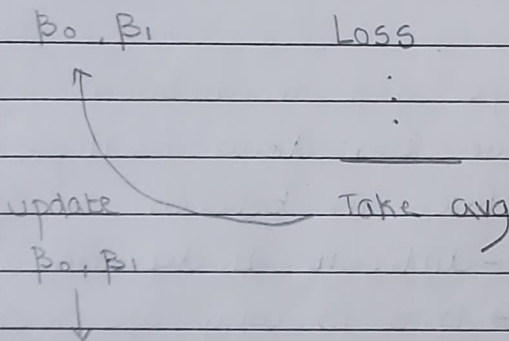
$\frac{\partial \text{Loss}}{\partial \beta_0, \beta_1}$  where slope = 0  $\rightarrow$  that is the global minimum point.

$$\beta_{new} = \beta_{old} - \eta \frac{\partial \text{Loss}}{\partial \beta_{old}}$$

why minus? we need to minimise  $\beta_{new}$

### ① Vanilla gradient descent - batch processing

- slow process
- For ex:



Calculate loss again - slow process

### ② Stochastic gradient descent

In 1 epoch: updating 1000 times

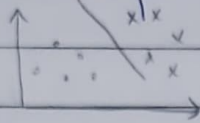
Therefore we apply Vanilla GD in stochastic GD

- mini batch gradient descent (randomly picks batch size)  
Perform stochastic - form batches  $\rightarrow$  take avg

Update  $\beta_0, \beta_1$

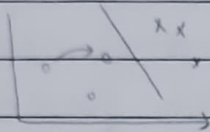
There should not be many batches

Linear Binary Classifier : when you can easily draw a line between two classes / if you can easily differentiate two classes into two parts



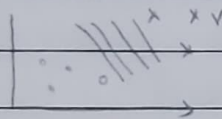
LOGISTIC REGRESSION  $\Rightarrow$

marks	GA	Pass/Fail
20	40	1
20	30	1
10	10	0
10	20	0



suppose one point is marked incorrectly, it shifts to new position  $\rightarrow$  line also shifts

suppose there is a huge gap in between

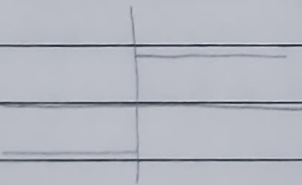


we can draw multiple lines in between

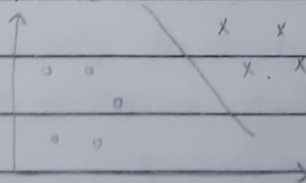
- difficult to determine

$\rightarrow$  this is why we need to use optimisation process.

$\downarrow$   
we get step function



This step function needs to be changed into sigmoid function

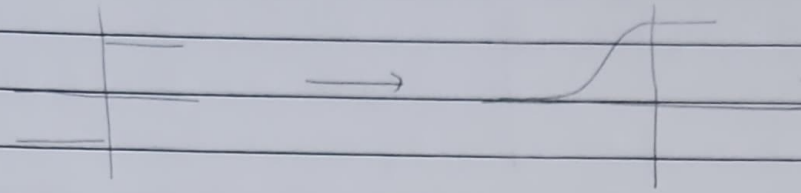


Take a line in between and calculate error

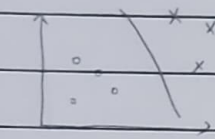
distance b/w point and line

$\rightarrow$  depending upon that taken as

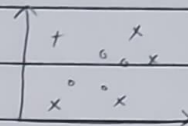
Likelihood is always taken for class 1  
Best line is the one with the best likelihood



Perceptron with Sigmoid Function = Logistic regression

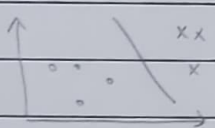


Linear

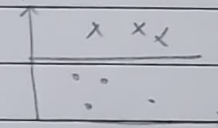


Non linear

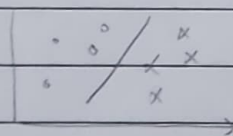
Using multiple perceptrons to get the output  
- Multilayer Perceptron (MLP)



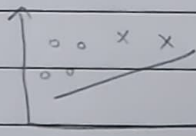
log



The backstory of Logistic regression - Perceptron



P = 0.8



P = 0.6

$$P = 0.8 + 0.6 = 1.4$$

Cannot be 1.4 as  $\leq \text{Prob} = 1$

pass through  
sigmoid function

O/p

