

* ARTIFICIAL NEURAL NETWORKS

INTRODUCTION TO ANN:

Backpropagation

→ Human brain has millions of

neurons. In Biological term, these are the smallest unit/organ of human brain which performs sensory task (i.e. all the sense organs of our body are stimulated by the neurons).

→ In Biological term $\xrightarrow{\text{called as}}$ Neurons
In formal term (i.e. in AI) \rightarrow Nodes }

→ Ex: When we touch a hot bowl, then we feel immediately take out our hand. Because our skin senses that, bowl is hot & pass this input to the neurons. Then neurons ~~at~~ process the input & generates some actions (i.e., we take out our hand). All these sensing, processing & generating actions are performed within a fraction of seconds. This is how neurons work.

→ For the same way, In AI,
natural functions are done
artificially.

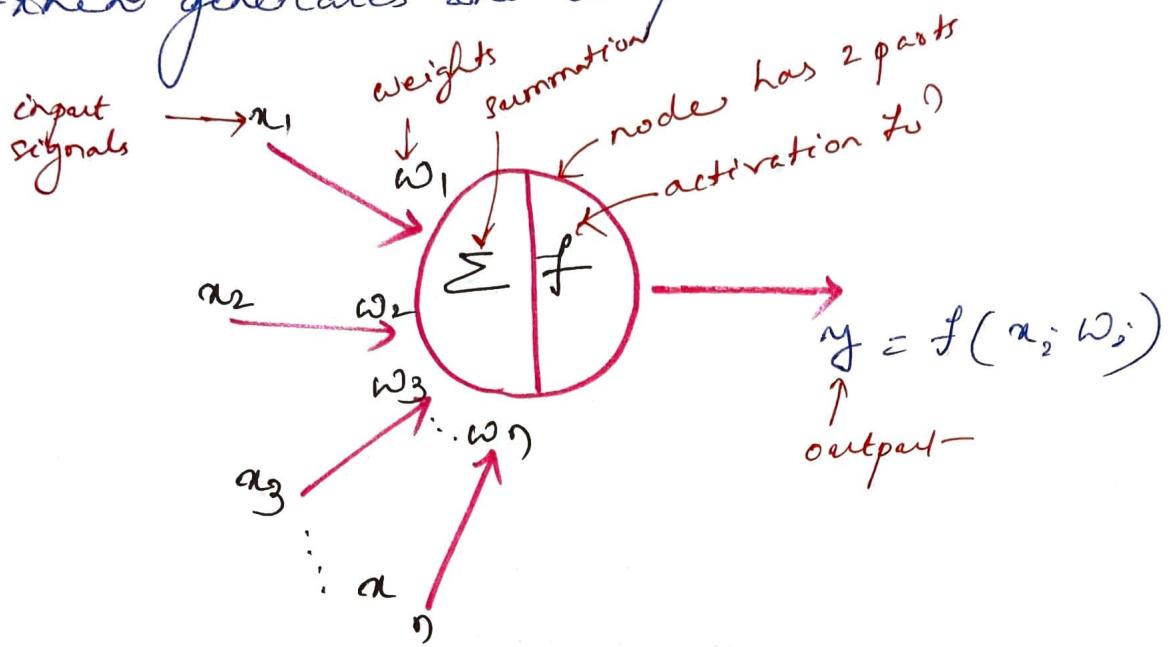
→ def :-

An ANN is usually a computational NW, based on biological neural networks that construct the structure of the human brain.

Similar to a human brain has neurons interconnected to each other, ANNs also have neurons that are linked to each other in various layers of the networks. These neurons are called nodes.

→ In the same way nodes can AI also works.

{ Nodes take the input,
using some f_o? it processes,
then generates the output.



A node has 2 parts

1) Summation

→ Calculates the weighted sum of all the inputs.

$$\sum_{i=1}^n x_i w_i = x_1 w_1 + x_2 w_2 + x_3 w_3 + \dots + x_n w_n$$

→ Once weighted sum is calculated

↓
it is sent to the activation function.

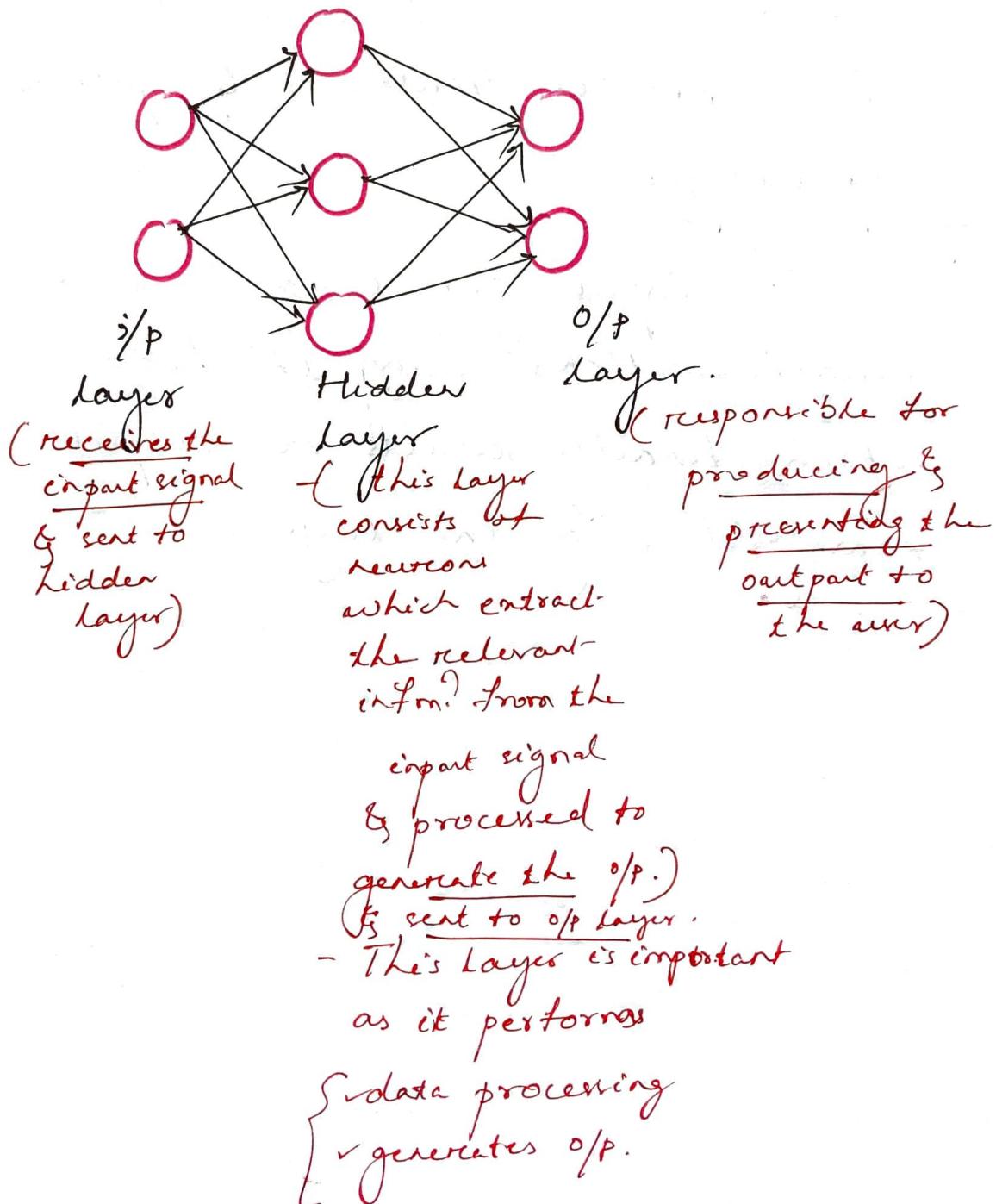
2) Activation Function

→ Generates the output based on input given
 $y = f(x_i; w_i)$

REPRESENTATION OF ANN:-

→ ANNs are divided into 3 parts

- { 1. Input layer
- 2. Hidden layer $\xrightarrow{\text{can be}}$ 1 layer (single layer ANN)
- 3. Output layer. $\xrightarrow{\text{or}}$ Σ Layer (multilayer ANN)



This is how ANN works.

Input layer :-

It accepts inputs in several different formats provided by the programme.

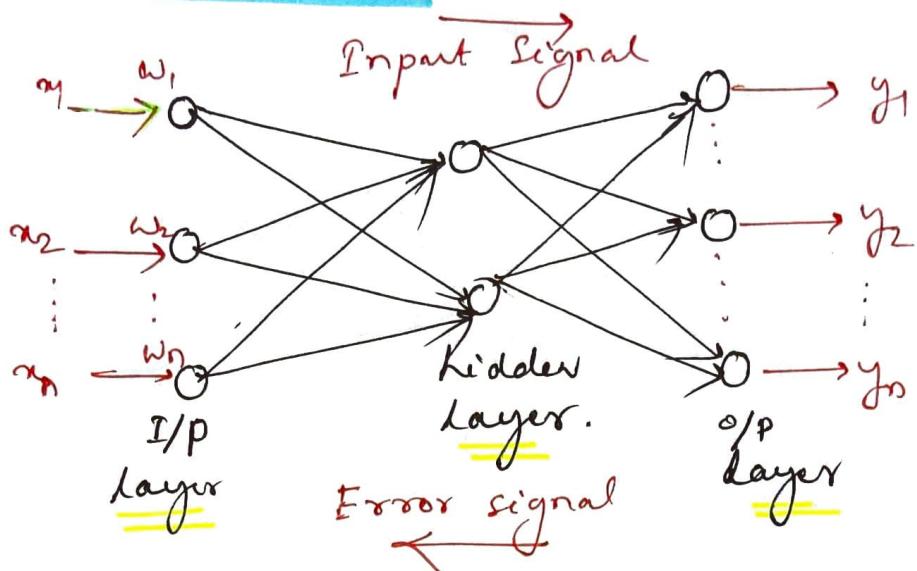
Hidden layer :-

- It is present bet. input & output layers.
- It performs all the calculations to find hidden features & patterns.

Output layer :-

The input goes through a series of transformations using the hidden layer, which finally results in output - that is conveyed using this layer.

HOW ANNs WORK?



- ANN can be best represented as a weighted graph directed graph, where artificial neurons form the node.
- The association between the neurons outputs & neuron inputs can be viewed as the directed edges with weights.
- The ANN receives the input signal from the external source in the form of a pattern & image in the form of a vector. These inputs are then mathematically assigned by the notations (x_n) for every no. of inputs like x_1, x_2, \dots, x_n .
- Then, each of the input is multiplied by its corresponding weights.

SS(b)

(Here, weights → the strength of the interconnection between neurons encode the ANN)

All the weighted inputs are summarized inside the computing unit.

→ If the weighted sum is = 0, then bias is added to make the output non-zero (i.e., something is added to scale-up to the system's response).

(Bias has the same input, weight = 1)

Hence the total of weighted inputs can be in the range of (0 to +∞)

Then the total of weighted input is passed through the activation fu.?

→ The activation fu.? refers to the set of transfer fu.? used to achieve the desired output.

⇒ This is how ANN works.

APPROPRIATE PROBLEMS FOR ANN

i.e., Where it works ?

Situations where ANN algm. can apply;

→ Instances have many attribute value pairs.

→ Target function has discrete values, continuous values or combination of both.

→ Training examples with errors or missing values.

: machine is trained with many attribute value pairs.

→ long training times are acceptable.

→ fast evaluation of the target function learnt. *→ testing is faster*

→ Ability for humans to understand the target function learnt by machine is not important. (*one, what machine has learnt difficult to understand by human in ANN*)

ADVANTAGES OF ANNs

→ parallel processing capability.

→ Storing data on the entire network.

→ Capability to work with incomplete knowledge.

→ Having fault tolerance.

DISADVANTAGES OF ANNs

- Assurance of proper network structure
- Hardware dependence.
- Difficulty of showing the issue to the network.
- The duration of the network is unknown.

* [SIGMOID ACTIVATION FUNCTION PROBLEM] (1)

Question:-

Obtain the output of the neuron Y using

- i) bipolar binary sigmoidal &
- ii) bipolar sigmoidal activation function.

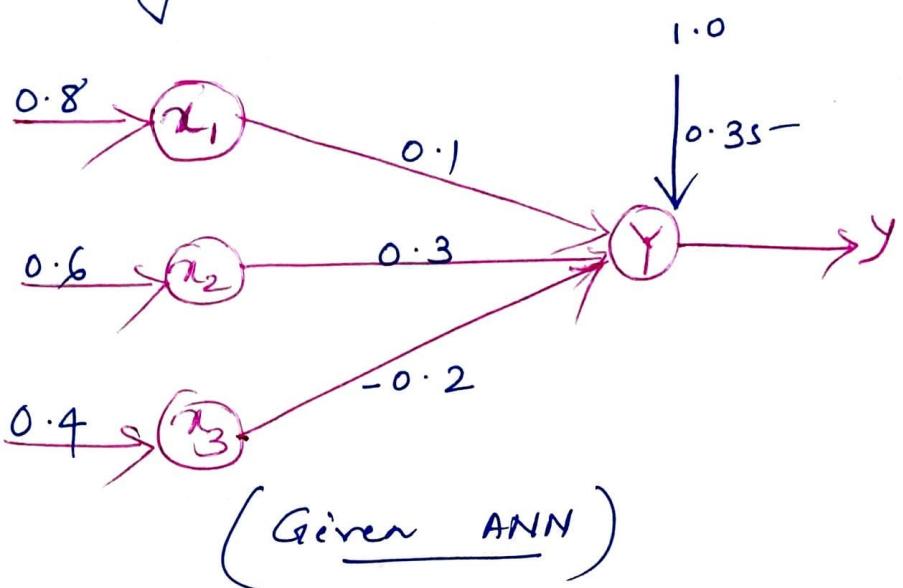
Solution:-

→ In this ANN,

2 layers are given :-

Input layer

Output layer.



→ Input layer is having 3 neurons. (x_1, x_2, x_3). Output layer is having 1 neuron. (Y)

Each input neuron is connected to output neuron with some weight.

$$\overrightarrow{x_1 Y} = 0.1$$

$$\overrightarrow{x_2 Y} = 0.3$$

$$\overrightarrow{x_3 Y} = -0.2$$

3 inputs to 3 neurons are $0.8, 0.6, 0.4$.

0.35 is the bias weight

→ There are 2 types of sigmoid functions:-

- 1) Binary Sigmoid activation function.

$$Y = f(Y_{in}) = \frac{1}{1 + e^{-Y_{in}}}$$

(Where, $Y_{in} \rightarrow$ total input to Y)

ii) Bipolar Sigmoid Activation Function. ②

$$y = f(y_{in}) = \frac{2}{1 + e^{-y_{in}}} - 1$$

→ Calculation of net input to the output neuron. (y_{in})

$$y_{in} = b + \sum_{i=1}^3 x_i w_i$$

(here $n = 3$ = no. of input neurons.
 $b = \text{bias} = 0.35$)

$$\begin{aligned} &= b + x_1 w_1 + x_2 w_2 + x_3 w_3 \\ &= 0.35 + (0.8 \times 0.1) + (0.6 \times 0.3) + (0.4 \times 0.2) \\ &= 0.35 + 0.08 + 0.18 - 0.08 \\ &= 0.53 \end{aligned}$$

→ Calculation of Binary Sigmoid Activation Function (y)

$$y = f(y_{in})$$

$$= \frac{1}{1 + e^{-y_{in}}} = \frac{1}{1 + e^{-0.53}} = 0.625$$

→ Calculation of Bipolar Sigmoid Activation Function (y)

$$y = f(y_{in})$$

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

$$= \frac{2}{1 + e^{-0.53}} - 1 = 0.259$$

→ **Perceptron** is a single layer neural network & a multilayer perceptron is called Neural Network.

→ Perceptron is a linear classifier (binary)

→ It is used in supervised learning.

→ It helps to classify the given input data.

→ The perceptron consists of 4 parts:-

- input values or one input layer.
- weights or Bias
- Net Sum
- Activation f(x) (transformation / step f(x))

→ Perceptron works on 3 steps:-

a) All the inputs are multiplied with their weights w .

b) Add all the multiplied values & call them as weighted sum.

c) Apply that weighted sum to the correct activation function.

→ Perceptron are the basic unit to build ANN.

✓ It takes real valued i/p, calculates linear combination of these inputs & generates the output.

$$\begin{cases} \text{Output} = 1 & \text{if result} > \text{threshold} \\ & \\ & = -1 \text{ otherwise.} \end{cases}$$

Hence the word linear combination mean,

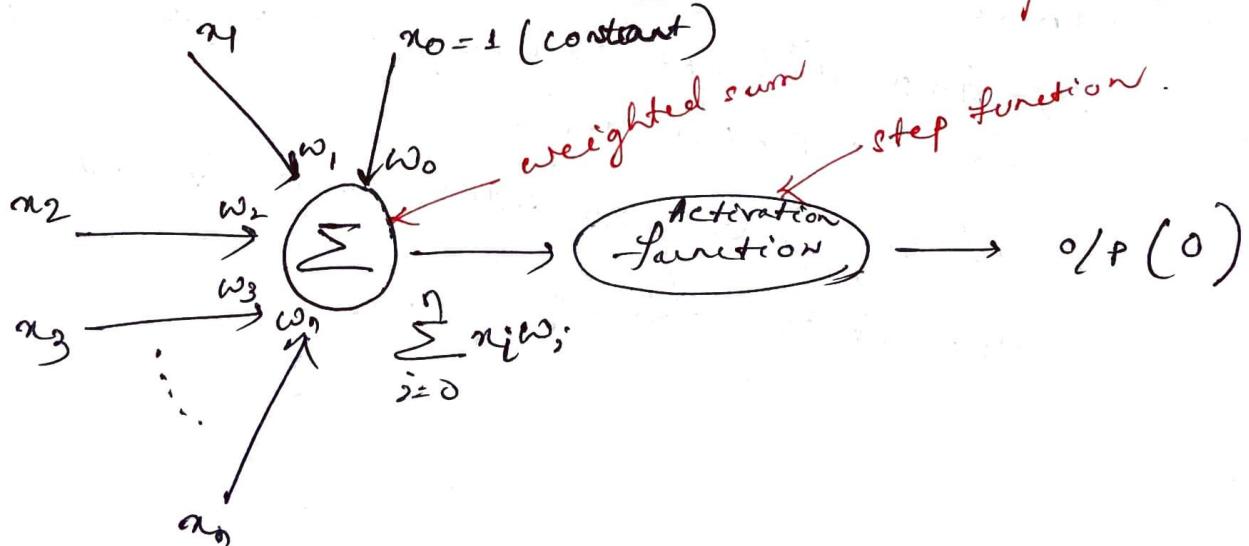
$$o(x_1, x_2, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n > 0 \\ -1 & \text{otherwise} \end{cases}$$

↑
threshold value

(Where, $w_0 \rightarrow$ initial weight -
 $w_1, w_2, \dots, w_n \rightarrow$ corresponding
weights of inputs.)

→ Receptron Rule,

$x_0 = 1$ (initial output) \rightarrow constant -
 $x_1, x_2, \dots, x_n \rightarrow$ inputs).



* This rule says that,
Each node has 2 parts.

(8)

↳ summation
 ↳ activation funⁿ.

- * Here summation adds all the inputs along with their corresponding weights.

Here, x_0 (critical input) = 1.

$x_1, x_2, \dots, x_n \rightarrow$ inputs.

$w_1, w_2, \dots, w_n \rightarrow$ weight.

$$\sum_{i=0}^n w_i x_i = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n.$$

$$= \begin{cases} 1 & / \\ -1 & \end{cases}$$

$$\sum_{i=0}^n w_i x_i > 0 \quad \sum_{i=0}^n w_i x_i < 0$$

- * Then this value is sent to the activation function which ~~processes~~ is responsible for processing the input & generates the output.
(English alphabet 0 not zero)

- * Here, $\hat{o} \rightarrow$ actual output (generated by the activation funⁿ)
 $t \rightarrow$ target output (expected output by us)

- * $\hat{o} \neq t \Rightarrow$ actual output = target output -

\Rightarrow then weights are fined

i.e., we need to consider critical $w_0, w_1, w_2, w_n \dots$ &

for generating the output. (S-9)

But, if actual output \neq target output
 \Rightarrow then we have to change the weights.

How to change weights?

Theoretically, take some random weights.

Later, keep on applying iterations

i.e., find $\sum_{i=0}^n w_i x_i$

and check if ($o = t$)

If $o \neq t$, then stop the process.

If $o \neq t$, then (apply formula to charge weight)

$$w_i \leftarrow w_i + \Delta w_i$$

$$\Delta w_i = n(t - o)x_i$$

add this small change in weight to make ($o = t$)

(where, n = learning rate (i.e., take small values like 0.1, 0.2, 0.3, ...)
 t = target o/p
 o = actual o/p)

x_i = i/p associated with w_i)

Hence $x_1 \rightarrow w_1$ is associated with
 $x_2 \rightarrow w_2$
 \vdots
 $x_n \rightarrow w_n$.

$$\text{i.e., } \Delta w_1 = n(t - o)x_1$$

$$\Delta w_2 = n(t - o)x_2$$

$$\Delta w_n = n(t - o)x_n$$

like this we need to find.

(60)

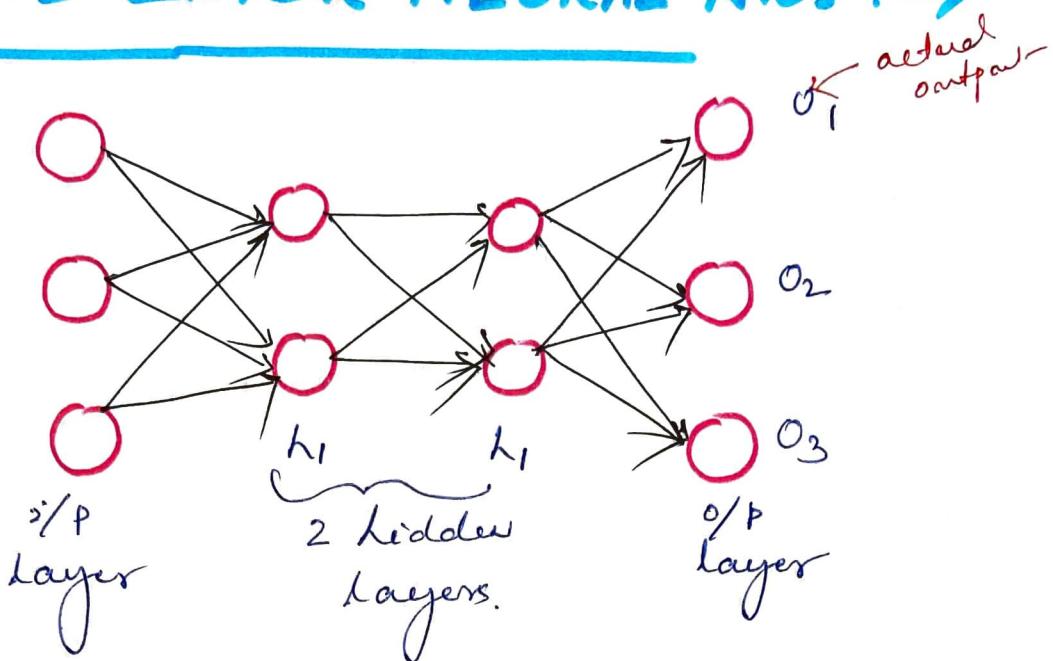
like this, find all the weights &
 find the $\sum_{i=0}^j w_i x_i$ for several iterations

~~↓~~ until $o = t$. i.e. actual output = target output.

Once it satisfies, then stop the process and
~~then~~ ~~train the~~ finalize the weights and
 train the perceptron.

⇒ This is how perceptron training
 goes on.

★ MULTI LAYER NEURAL NNs :→



- The single layer neural network will work only for linearly separable data. It will not work for non linearly separable data. Multi layer neural NNs solves this problem.
- In multi layer neural network, there are 2 hidden layers and each node in a connected single layer is connected to each and every node of the next layer and so on.
- Here O₁, O₂, O₃ are the actual outputs and are compared against the target outputs t₁, t₂, t₃ which we are supposed to get. If both doesn't match, then we have to go back & modify the weights.

(62)

which is called as Back-propagation.

⇒ This is all how multi-layer neural network works.

* BACK PROPAGATION : → (Backward propagation error)

- Back propagation (backward propagation) is an important mathematical tool for improving the accuracy of predictions in data-mining & machine learning.
- Artificial Neural Networks use backpropag. as a learning algorithm to compute a gradient descent w.r.t weights. Desired outputs are compared to achieved systems outputs, (Here desired outputs → target outputs system achieved outputs → actual outputs) and then systems are tuned by adjusting connection weights to narrow the difference between the two as much as possible.

- This algorithm is called as Back-propagation because the weights are updated backwards, from outputs towards inputs.

→ So, when error occurs⁽⁶³⁾, we go in backward direction.

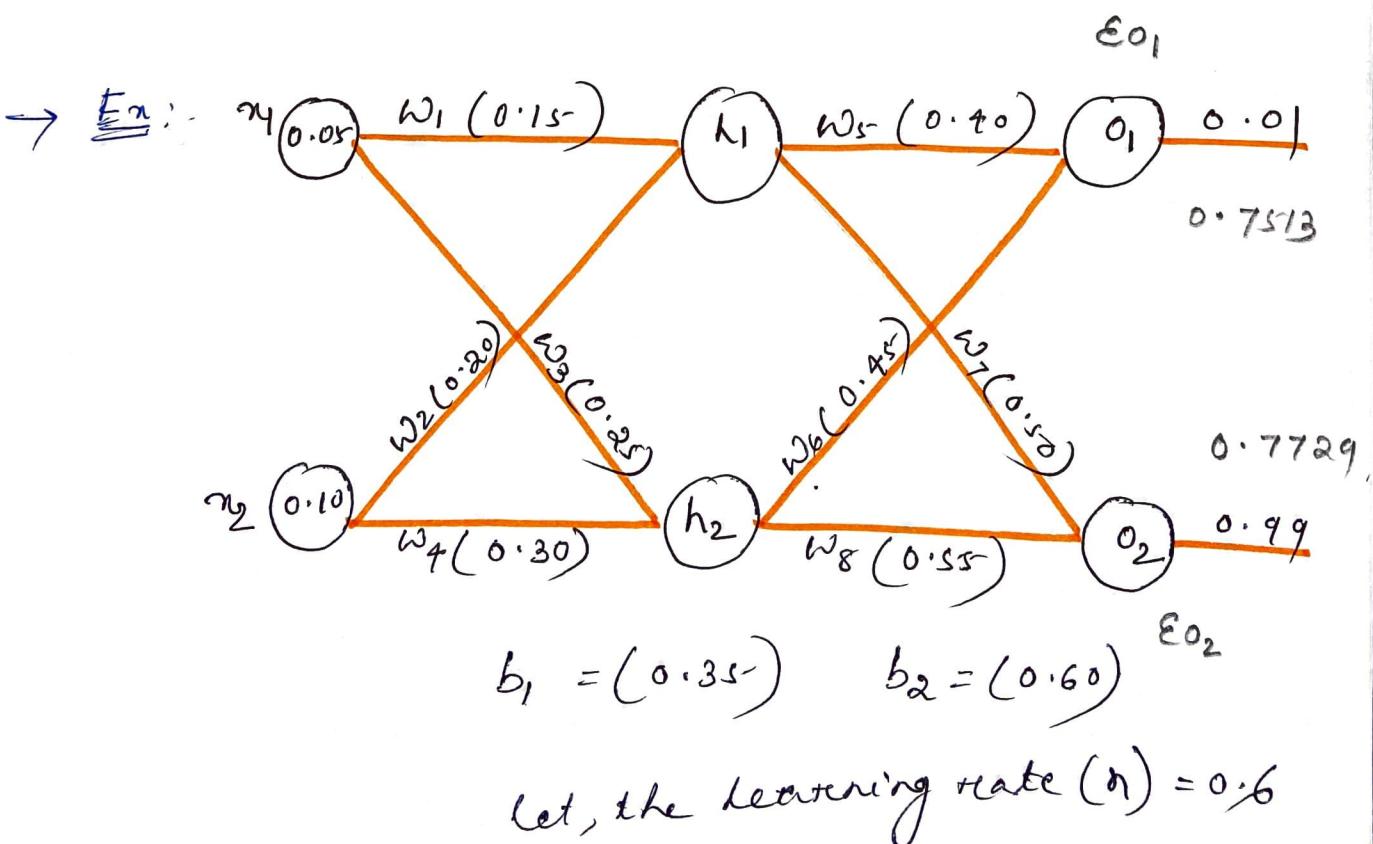
i.e., output → hidden → input layer.

→ Here the propagation errors will be calculated in 3 parts.

part 1: Calculate forward propagation error.
i.e., from $\text{input} \rightarrow \text{hidden} \rightarrow \text{output}$ layer.

part 2: Calculate backward propagation error.
i.e. from output layer \rightarrow hidden layer.

part 3: Calculating backward propagation error
i.e. from hidden \rightarrow input layer.



(64)

Here, $x_1, x_2 \rightarrow$ are the inputs $(0.05, 0.10)$
 $h_1, h_2 \rightarrow$ hidden layer.

$o_1, o_2 \rightarrow 0.01, 0.99$ (target-
outputs are
 $w_1, w_2 \dots w_p \rightarrow$ corresponds given).
weights of all
edges,

$b_1 \rightarrow$ is the bias factor from input to
hidden layer. (0.35) .

$b_2 \rightarrow$ is the bias factor from hidden
to output layer. (0.60)

orange lines \rightarrow edges.

We have to findout the actual outputs,
if it matches with the target outputs,
then it is said to be no error. If it does not
match, then we have to change the weights
is using back propagation algorithm.

part 1: Calculate forward propagation
error (input \rightarrow hidden \rightarrow output layer)

1) Calculate h_1 (in and out)

$$\begin{aligned}
h_1(in) &= x_1 w_1 + x_2 w_2 + b_1 \\
&= (0.05 \times 0.15) + (0.10 \times 0.20) + 0.35 - \\
&= 0.377
\end{aligned}$$

(65)

$$h_1(\text{out}) = \frac{1}{1 + e^{-h_1(\text{in})}}$$

$$= \frac{1}{1 + e^{-(0.377)}}$$

$$= 0.5932$$

2) Calculate $h_2(\text{in}$ and $\text{out})$

$$h_2(\text{in}) = w_3 w_3 + w_4 w_4 + b_1$$

$$= (0.05 \times 0.25) + (0.10 \times 0.30) + 0.35$$

$$= 0.0125 + 0.03 + 0.35$$

$$= 0.3925$$

$$h_2(\text{out}) = \frac{1}{1 + e^{-h_2(\text{in})}}$$

$$= \frac{1}{1 + e^{-(0.3925)}}$$

$$= 0.5968$$

3) Calculate $o_1(\text{in}$ and $\text{out})$

$$o_1(\text{in}) = h_1(\text{out}) \times w_5 + h_2(\text{out}) \times w_6 + b_2$$

$$= (0.593 \times 0.40) + (0.596 \times 0.45) + 0.60$$

$$= 1.105$$

$$o_1(\text{out}) = \frac{1}{1 + e^{-o_1(\text{in})}}$$

$$= \frac{1}{1 + e^{-(1.105)}} = 0.7513$$

+) Calculate O_2 (in and out) (66)

$$O_2(\text{in}) = h_1(\text{out}) \times w_1 + h_2(\text{out}) \times w_2 * b_2$$

$$= (0.5932 \times 0.50) + (0.5968 \times 0.58)$$

$$+ 0.60$$

$$= 1.22484$$

$$O_2(\text{out}) = \frac{1}{1 + e^{-O_2(\text{in})}}$$

$$= \frac{1}{1 + e^{-(1.22484)}}$$

$$= 0.7729$$

e) Calculate E_{Total}

$$E_{\text{Total}} = \sum \frac{1}{2} (\text{target} - \text{actual})^2$$

$$= E_{O_1} + E_{O_2}$$

$$= \frac{1}{2} (0.01 - 0.7573)^2 + \frac{1}{2} (0.99 - 0.7729)^2$$

$$= \frac{1}{2} (-0.7413)^2 + \frac{1}{2} (0.2171)^2$$

$$= 0.274 + 0.0235$$

$$= 0.29837 \text{ (approximately)}$$

(here $E_{\text{Total}} \rightarrow$ Total error
 $E_{O_1} \rightarrow$ error at output O_1 ,
 $E_{O_2} \rightarrow$ " " " " O_2)

part 2

: Calculate backward propagation error (output layer \rightarrow hidden layer)

Now, here we have to adjust the weights which are between output layer & hidden layer.

i.e., w_5 , w_6 , w_7 and w_8 .

first let us adjust w_5 ,

$$w_5^* = w_5 - \eta \frac{\partial E_{\text{Total}}}{\partial w_5}$$

(Hence, $\eta \rightarrow$ learning rate (0.6)^{let}
 $w_5 \rightarrow$ weight is given.
 $\partial \rightarrow$ partial differentiation)

$$\frac{\partial E_{\text{Total}}}{\partial w_5} = \frac{\partial E_{\text{Total}}}{\partial \text{out}_{0_1}} \times \frac{\partial \text{out}_{0_1}}{\partial \text{net}_{0_1}} \times \frac{\partial \text{net}_{0_1}}{\partial w_5}$$

(apply chain rule)

$$\begin{aligned} \frac{\partial E_{\text{total}}}{\partial \text{out}_{0_1}} &= \text{out}_{0_1} - \text{target}_{0_1} \\ &= 0.751365 - 0.01 \\ &= 0.7413565 \end{aligned}$$

$$\begin{aligned} \frac{\partial \text{out}_{0_1}}{\partial \text{net}_{0_1}} &= \text{out}_{0_1}(1 - \text{out}_{0_1}) \\ &= 0.751365(1 - 0.751365) \\ &= 0.186815602 \end{aligned}$$

(68)

$$\frac{\partial \text{Net O}_1}{\partial w_5} = \text{out}_1 \\ = 0.59326992$$

Now, $\frac{\partial E_{\text{Total}}}{\partial w_5} = 0.7913565 \times 0.186815602$
 $\times 0.59326992$
 $= 0.08216704$

Hence, $w_5^* = w_5 - n \frac{\partial E_{\text{Total}}}{\partial w_5}$
 $= 0.4 - (0.6) \times 0.08216704$
 $= 0.350699776$

Here,

w_5^* → modified weight at w_5
 w_5 → actual weight
 n → learning rate (0.6) ut)

B.

Similarly we have to find the modified weights for w_6^* , w_7^* and w_8^* .

After finding the modified weight, repeat part 1 to check the output to check the output. Repeat this process several times until the actual output matches the target output (or threshold value) i.e. almost ^{nearly} equal to zero.

(69)

part 3 :- Calculating backward propagation
of error. (Hidden \rightarrow input layer)

Here we have to modify the weights
 w_1, w_2, w_3 and w_4 .

First let us adjust w_1 .

$$w_1^* = w_1 - \eta \frac{\partial E_{\text{Total}}}{\partial w_1}$$

$$\frac{\partial E_{\text{Total}}}{\partial w_1} = \frac{\partial E_{\text{Total}}}{\partial \text{out}(h_1)} \times \frac{\partial \text{out}(h_1)}{\partial \text{net}(h_1)} \times \frac{\partial \text{net}(h_1)}{\partial w_1}$$

(apply chain rule)

$$\frac{\partial E_{\text{Total}}}{\partial \text{out}(h_1)} = \frac{\partial E_{O_1}}{\partial \text{out}(h_1)} + \frac{\partial E_{O_2}}{\partial \text{out}(h_1)}$$

$$\frac{\partial E_{O_1}}{\partial \text{net}_{O_1}} \times \frac{\partial \text{net}_{O_1}}{\partial \text{out}(h_1)}$$

$$\frac{\partial E_{O_1}}{\partial \text{out}_{O_1}} \times \frac{\partial \text{out}_{O_1}}{\partial \text{net}_{O_1}}$$

$$\frac{\partial E_{O_2}}{\partial \text{net}_{O_2}} \times \frac{\partial \text{net}_{O_2}}{\partial \text{out}(h_1)}$$

$$\frac{\partial E_{O_2}}{\partial \text{out}_{O_2}} \times \frac{\partial \text{out}_{O_2}}{\partial \text{net}_{O_2}}$$

$$\frac{\partial E_{O_2}}{\partial \text{out}_{O_2}} = (\text{out}_{O_2} - \text{target}_{O_2})$$

$$= 0.772928465 - 0.99$$

$$= -0.217071535$$

$$\begin{aligned} \frac{\partial \text{out}_{O_2}}{\partial \text{net}_{O_2}} &= \text{out}_{O_2} (1 - \text{out}_{O_2}) \\ &= 0.7729 (1 - 0.7729) \\ &= 0.175510052 \end{aligned}$$

(70)

Now,

$$\checkmark \frac{\partial E_{O_2}}{\partial \text{net } O_2} = \frac{\partial E_{O_2}}{\partial \text{out } O_2} \times \frac{\partial \text{out } O_2}{\partial \text{net } O_2}$$

$$= -0.2170715^{\circ} 35^{\circ} X$$

$$0.175 \leftarrow 1005^{\circ} 2.$$

$$= \underline{\quad}$$

Similarly, find

$$\checkmark \frac{\partial E_{O_1}}{\partial \text{net } O_1} = \frac{\partial E_{O_1}}{\partial \text{out } O_1} \times \frac{\partial \text{out } O_1}{\partial \text{net } O_1}$$

$$\downarrow \qquad \qquad \qquad \nearrow \text{out } O_1 (1 - \frac{\text{out } O_1}{O_1})$$

$$(\text{out}_2 - \text{target } O_1)$$

$$= \text{Find the product.}$$

Now,

$$\frac{\partial \text{net } O_1}{\partial \text{out } L_1} = \text{on } O_1 \text{ from } L_1 \Rightarrow w_s -$$

$$= 0.40$$

$$\frac{\partial \text{net } O_2}{\partial \text{out } L_1} = \text{on } O_2 \text{ from } L_1 \Rightarrow w_f$$

$$= 0.50$$

Now, $(0.13849856 \times 0.40) + (-0.0380982 \times 0.50)$

$$\checkmark \frac{\partial E_{\text{Total}}}{\partial \text{out } L_1} = 0.055399425 + (-0.019049119)$$

$$= 0.036350306$$

(71)

$$\checkmark \frac{\partial o_{\text{auth}_1}}{\partial w_{\text{eth}_1}} = o_{\text{auth}_1} (1 - o_{\text{auth}_1}) \\ = 0.241300709$$

$$\checkmark \frac{\partial o_{\text{eth}_1}}{\partial w_1} = \frac{\partial (w_1 x_1 + w_2 x_2 + b_1)}{\partial w_1} \\ = x_1 = 0.05 -$$

Now,

$$\frac{\partial E_{\text{Total}}}{\partial w_1} = 0.036350306 \times 0.241300709 \times 0.05 \\ = 0.00438568$$

$$\text{Now, } w_1^* = w_1 - \eta \frac{\partial E_{\text{Total}}}{\partial w_1}$$

$$= (0.15) - ((0.6) \times 0.00438568) \\ = 0.1497368592$$

Similarly we have to calculate

$$w_2^*, w_3^*, w_4^*.$$

In the previous part 2, we have calculated

$$w_5^*, w_6^*, w_7^* \text{ & } w_8^*$$

Using all these updated weights, proceed part 1 i.e., calculate the error again & again until the error is being reduced almost nearly equal to zero or below the target/ threshold value.



SUPPORT VECTOR MACHINES (SVM)

Presented by
Ms. Bidyutlata Sahoo
Assistant Professor
CSE (AI & ML)

OBJECTIVES

The objective of SVM algorithm is ***to find a hyperplane in an N-dimensional space that distinctly classifies the data points.***

The dimension of the hyperplane depends upon the number of features.

OUTCOMES

SVM or Support Vector Machine is **a linear model for classification and regression problems.** It **can solve linear and non-linear problems** and work well for many practical problems.

SUPPORT VECTOR MACHINES (SVM)

- Support Vector Machine or SVM is one of the most popular ***Supervised Learning algorithms***, which is ***used for Classification as well as Regression problems.***
- However, primarily, it is used for Classification problems in Machine Learning.
- The goal of the SVM algorithm is to create the best line or decision boundary that can segregate n-dimensional space into classes so that we can easily put the new data point in the correct category in the future. **This best decision boundary is called a hyperplane.**

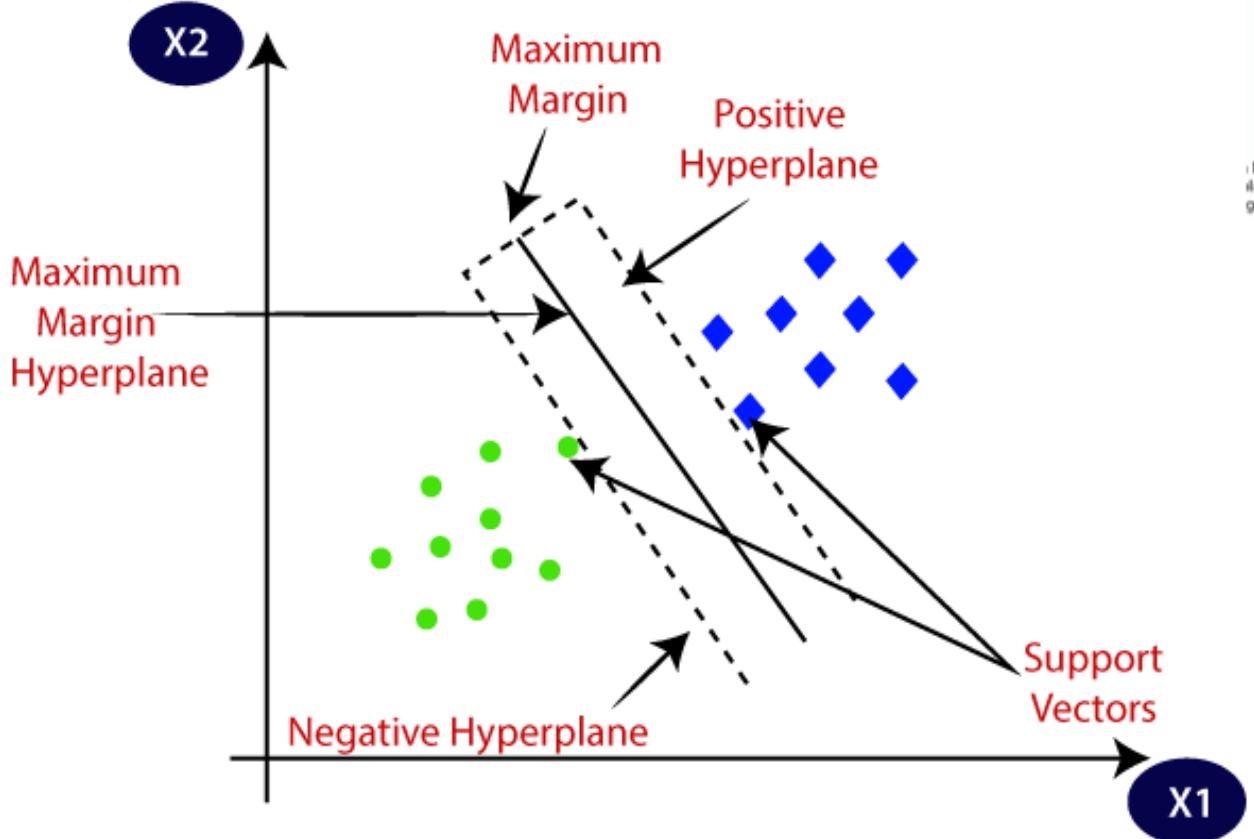
SUPPORT VECTOR MACHINES (SVM)

SVM chooses the ***extreme points/vectors*** that help in creating the hyperplane. These extreme cases are called as ***support vectors***, and hence algorithm is termed as Support Vector Machine.

SUPPORT VECTOR MACHINES (SVM)

Diagram

Consider the diagram in which there are two different categories that are classified using a decision boundary or hyperplane:



SUPPORT VECTOR MACHINES (SVM)

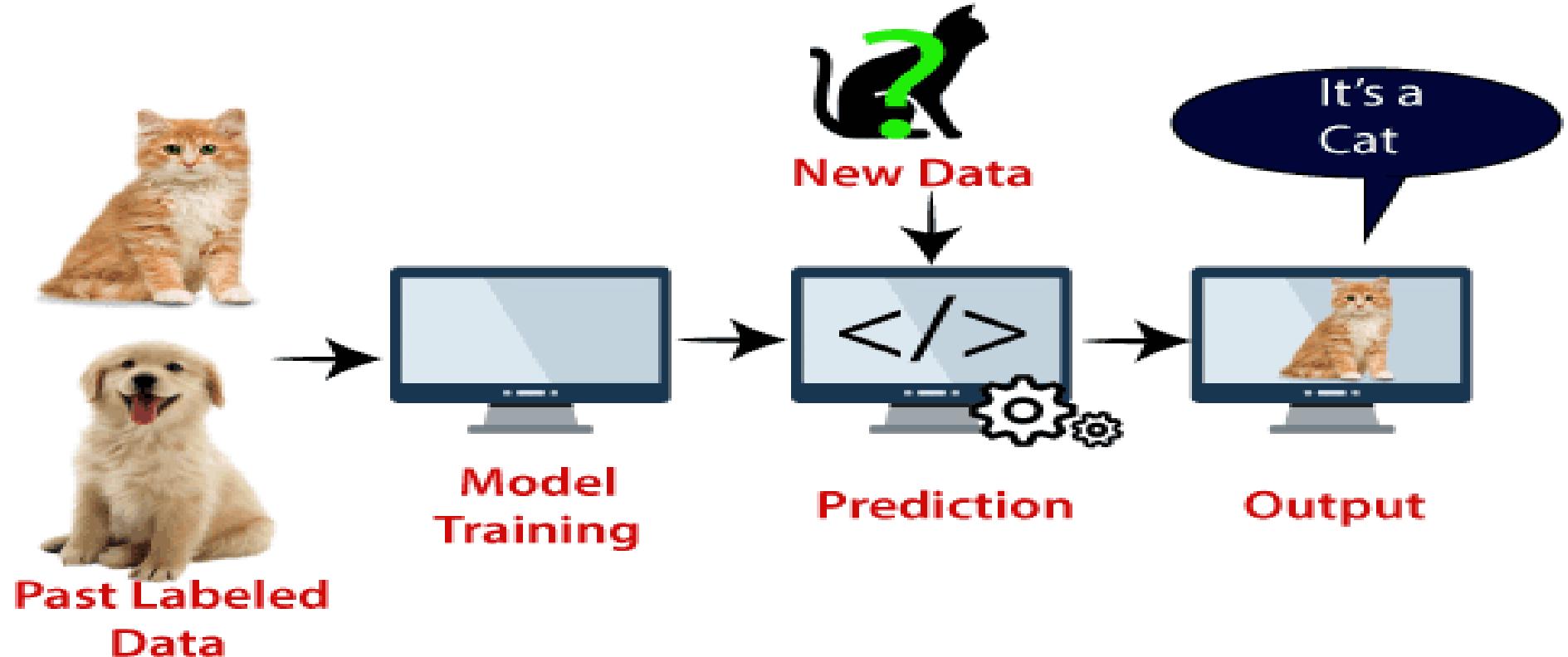
Example

Suppose we see a strange cat that also has some features of dogs, so if we want a model that can accurately identify whether it is a cat or dog, so such a model can be created by using the SVM algorithm. We will first train our model with lots of images of cats and dogs so that it can learn about different features of cats and dogs, and then we test it with this strange creature. So as support vector creates a decision boundary between these two data (cat and dog) and choose extreme cases (support vectors), it will see the extreme case of cat and dog. On the basis of the support vectors, it will classify it as a cat. Consider the below diagram:

SUPPORT VECTOR MACHINES (SVM)

Example

Consider the below diagram:

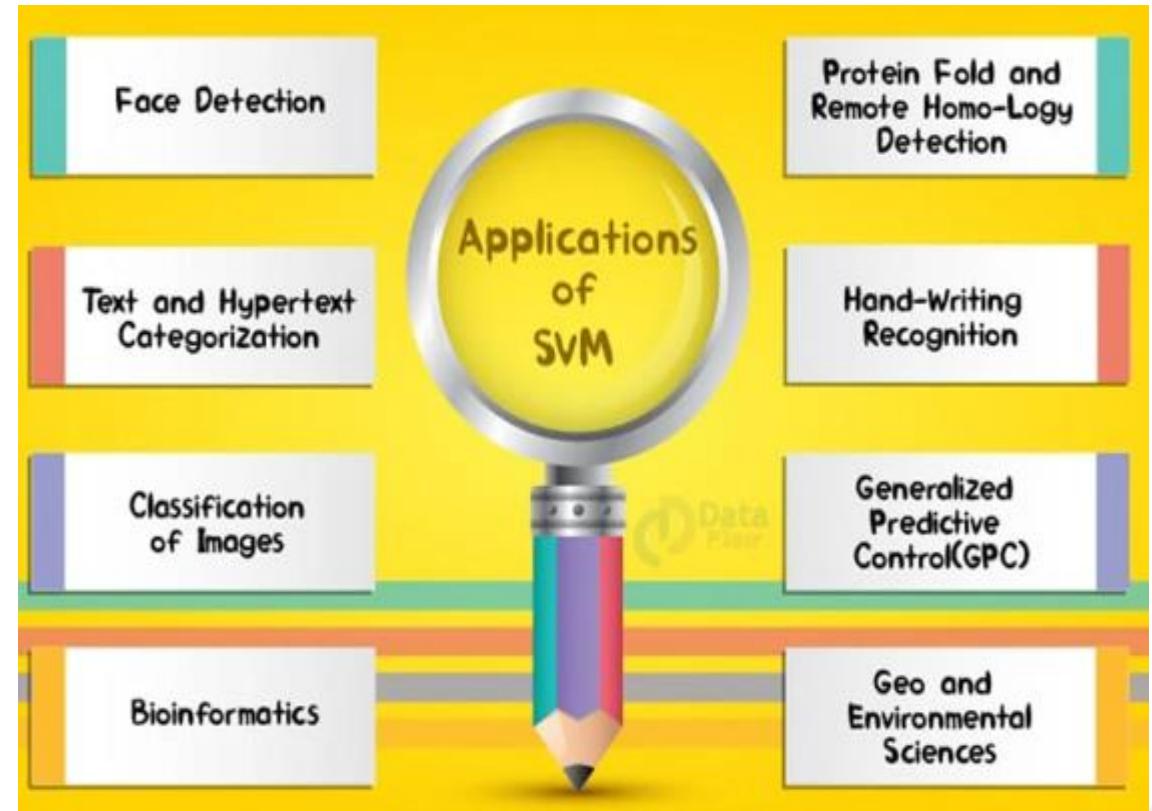


SUPPORT VECTOR MACHINES (SVM)

Applications

Some common applications of SVM are-

- Face Detection*
- Text and Hypertext Categorization*
- Classification of Images*
- Bioinformatics*
- Protein fold and remote homology detection*
- Handwriting recognition*
- Generalized Predictive Control*



SUPPORT VECTOR MACHINES (SVM)

Advantages of SVM

- SVM offers **best classification performance (accuracy) on the training data.**
- SVM renders **more efficiency for correct classification of the future data.**
- It **does not make any strong assumptions on data.**
- SVM **robust to outliers.**
- SVM **are conceptually easy to understand.**
- **The main strength of SVM is that they work well even when the number of SVM features is much larger than the number of instances.**

SUPPORT VECTOR MACHINES (SVM)

Disadvantages of SVM

- SVM is **not suitable for large data sets.**
- It **does not perform very well when the data has more noise** i.e. target classes are overlapping.
- In cases where the number of features for each data point exceeds the number of training data samples, the **SVM will underperform.**



LINEAR AND NON-LINEAR SVM

Presented by:
Ms. Bidyutlata Sahoo
Assistant Professor
CSE (AI & ML)

OBJECTIVES

The objective of SVM algorithm is ***to find a hyperplane in an N-dimensional space that distinctly classifies the data points.***

The dimension of the hyperplane depends upon the number of features.

OUTCOMES

SVM or Support Vector Machine is **a linear model for classification and regression problems.** It **can solve linear and non-linear problems** and work well for many practical problems.

LINEAR AND NON-LINEAR SVM

Types of SVMs

SVM can be of ***two types:***

- Linear SVM***
- Non-Linear SVM***

LINEAR AND NON-LINEAR SVM

Types of SVMs

Linear SVM:

Linear SVM is used for linearly separable data, which means *if a dataset can be classified into two classes by using a single straight line*, then such data is termed as linearly separable data, and classifier is used called as **Linear SVM classifier**. (i.e, the training data/classifiers are separated by a hyperplane.)

LINEAR AND NON-LINEAR SVM

Types of SVMs

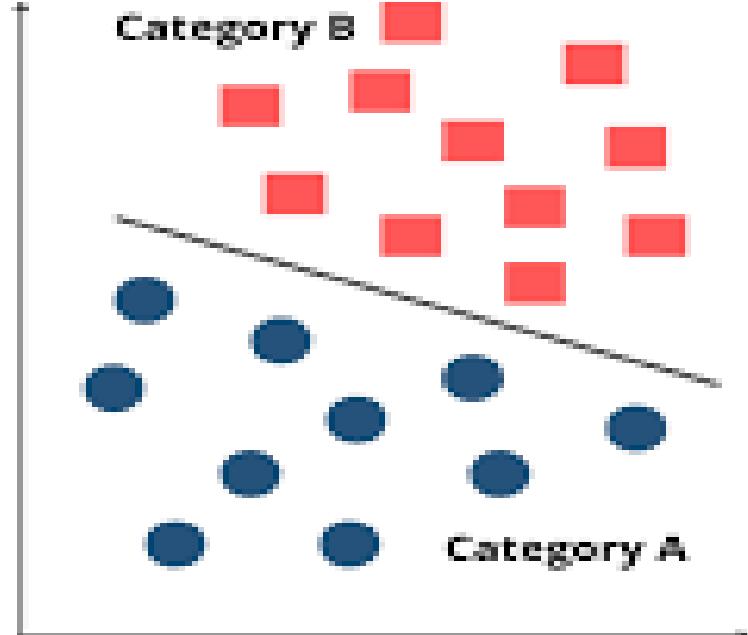
Non-linear SVM:

Non-Linear SVM is used for non-linearly separated data, which means ***if a dataset cannot be classified by using a straight line***, then such data is termed as non-linear data and classifier used is called as ***Non-linear SVM classifier***. (i.e, it is not possible to separate the training data using a hyperplane)

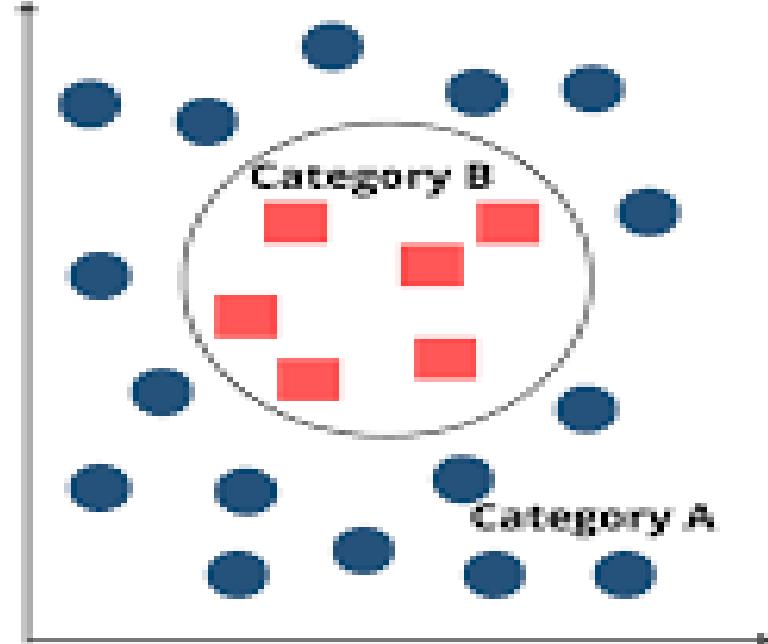
LINEAR AND NON-LINEAR SVM

Types of SVMs

Linear SVM



Non Linear SVM



LINEAR AND NON-LINEAR SVM

Hyperplane in SVM

- There can be multiple lines/decision boundaries to segregate the classes in n-dimensional space, but we need to find out the best decision boundary that helps to classify the data points. ***This best boundary is known as the hyperplane of SVM.***
- The dimensions of the ***hyperplane depend on the features present in the dataset,*** which ***means if there are 2 features (as shown in image), then hyperplane will be a straight line. And if there are 3 features, then hyperplane will be a 2-dimension plane.***

LINEAR AND NON-LINEAR SVM

Hyperplane in SVM

[NOTE: We always create a hyperplane that has a maximum margin, which means the maximum distance between the data points.]

LINEAR AND NON-LINEAR SVM

Support Vectors in SVM

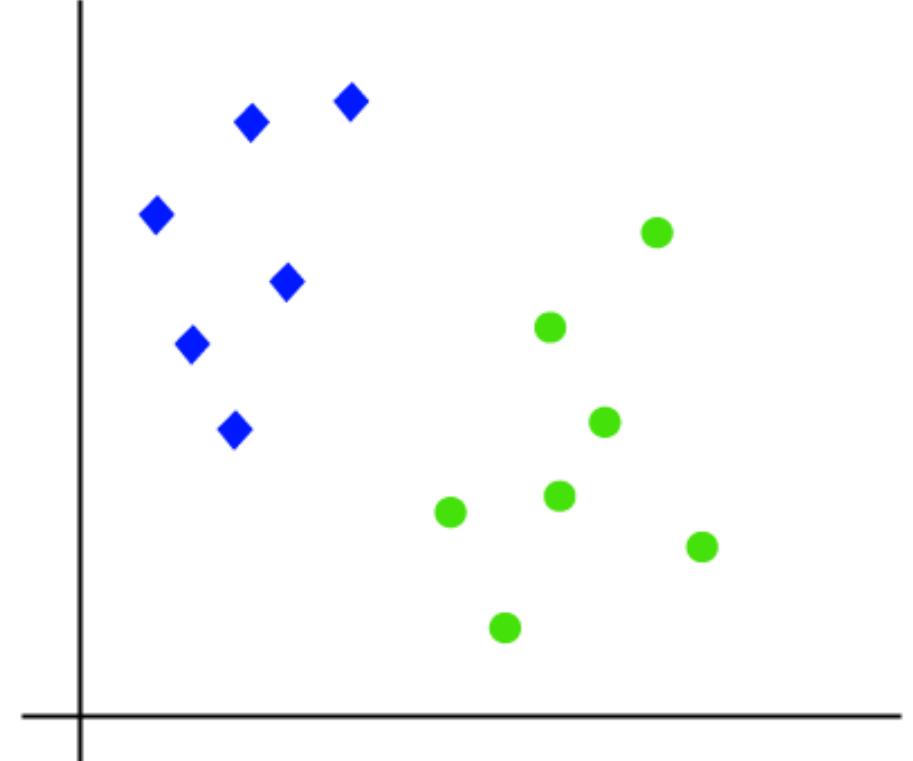
The **data points or vectors that are the closest to the hyperplane and which affect the position of the hyperplane** are termed as **Support Vector**. Since these vectors support the hyperplane, hence called a Support vector.

LINEAR AND NON-LINEAR SVM

Working of Linear SVM with Example

The working of the SVM algorithm can be understood by using an example.

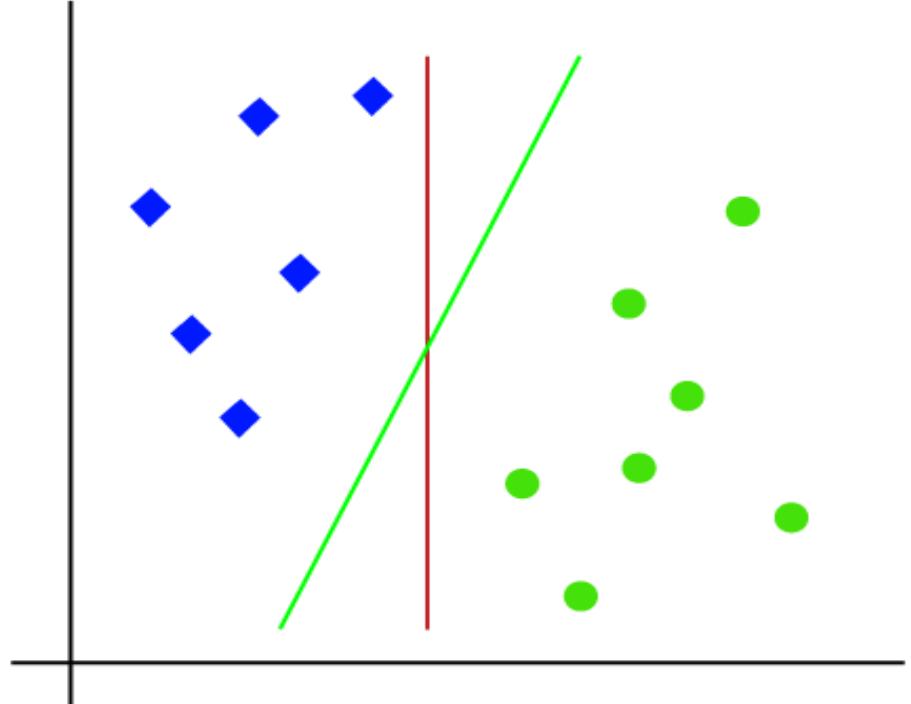
Suppose we have a dataset that has two tags (green and blue), and the dataset has two features x_1 and x_2 . We want a classifier that can classify the pair (x_1, x_2) of coordinates in either green or blue. Consider the given image:



LINEAR AND NON-LINEAR SVM

Working of Linear SVM with Example

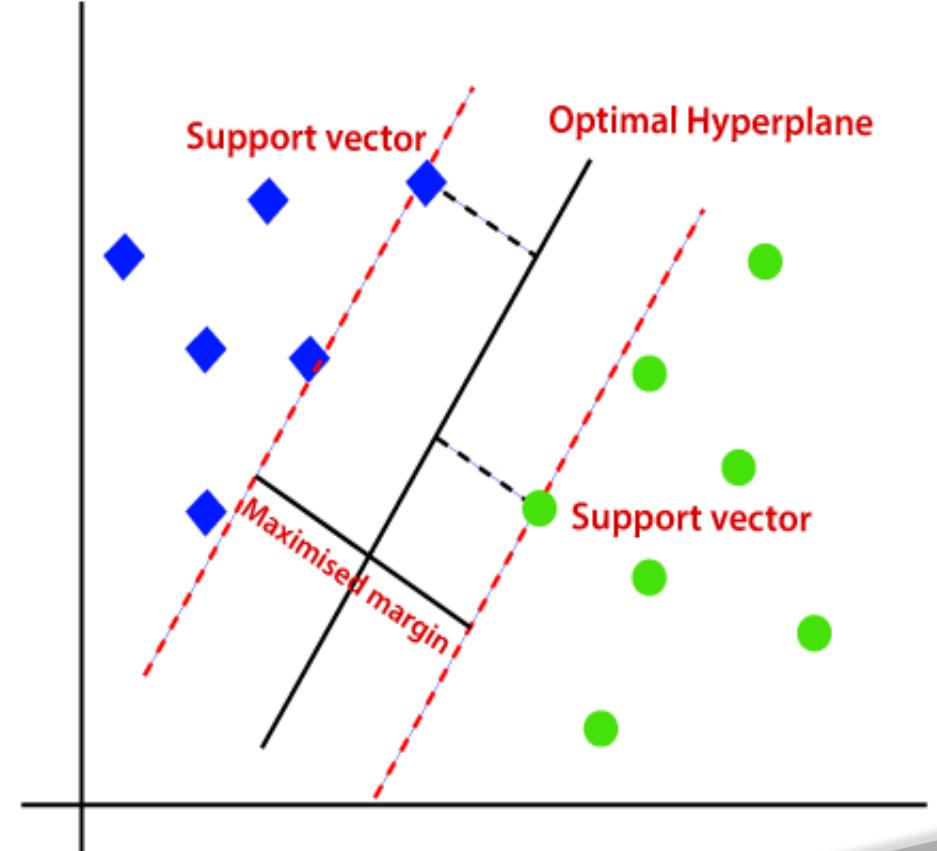
So as it is a 2-d space, by just using a straight line, we can easily separate these two classes. But there can be multiple lines that can separate these classes. Consider the given image:



LINEAR AND NON-LINEAR SVM

Working of Linear SVM with Example

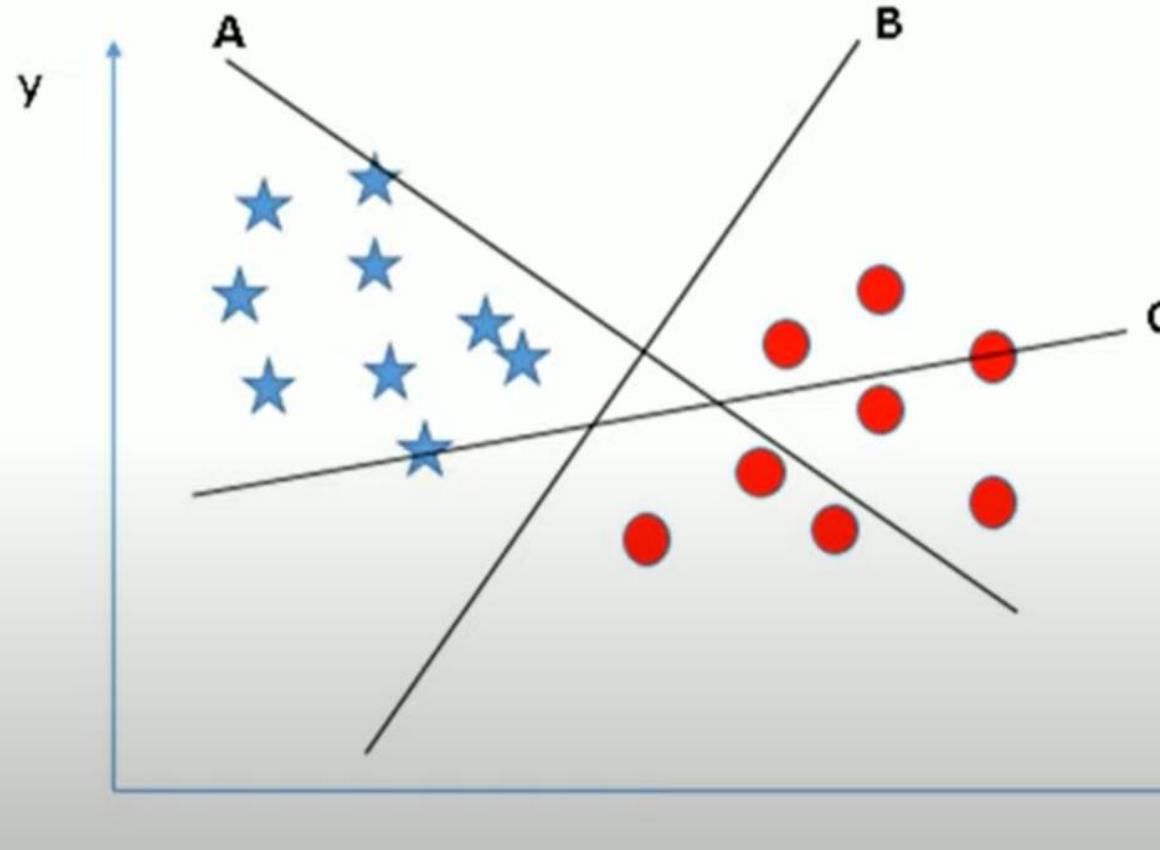
Hence, the **SVM algorithm helps to find the best line or decision boundary**; this best boundary or region is called as a **hyperplane**. SVM algorithm finds the closest point of the lines from both the classes. These points are called support vectors. The distance between the vectors and the hyperplane is called as **margin**. And the goal of SVM is to maximize this margin. The **hyperplane with maximum margin is called the optimal hyperplane**.



LINEAR AND NON-LINEAR SVM

Working of Linear SVM with more Examples – Scenario-1

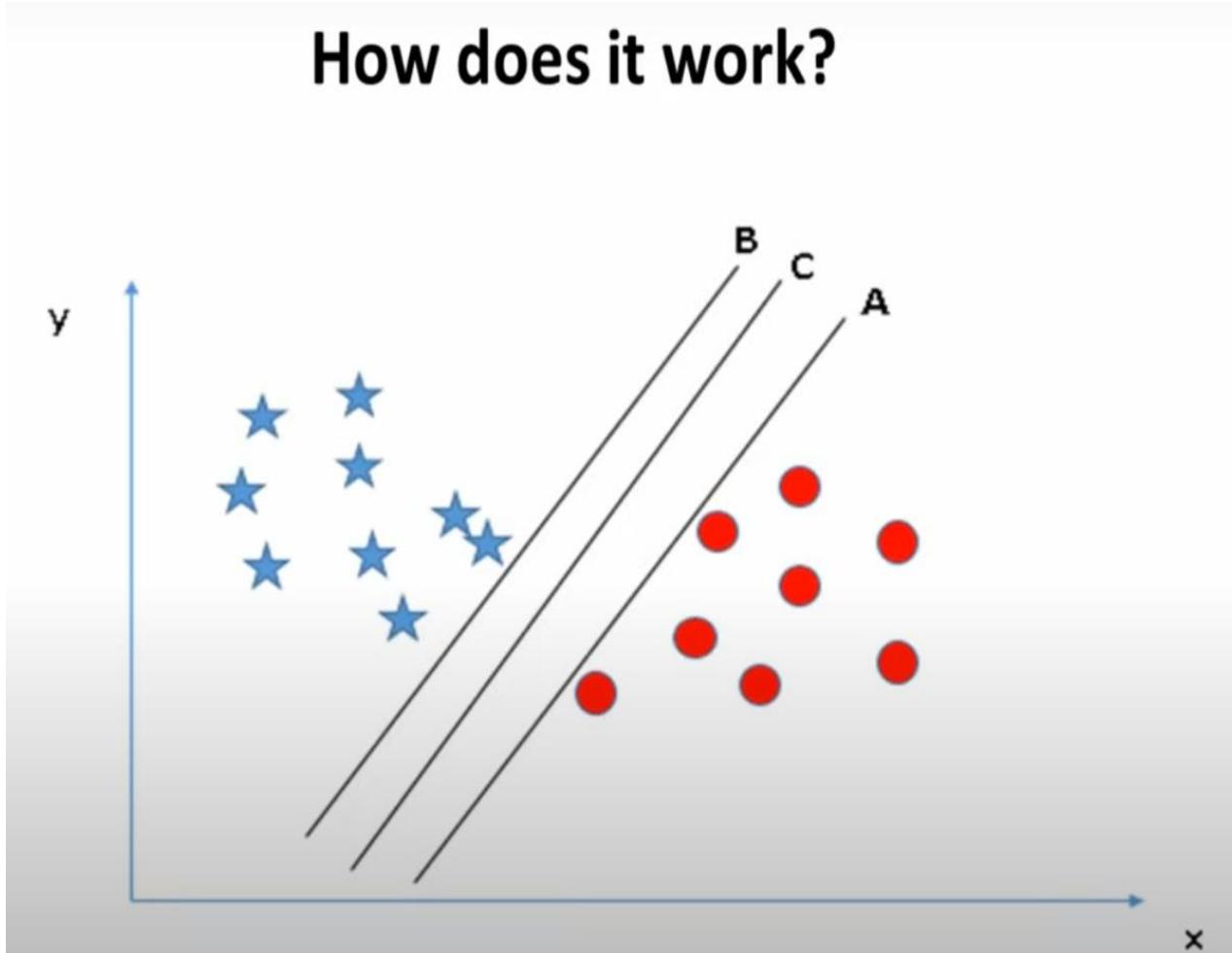
How does it work?



LINEAR AND NON-LINEAR SVM

Working of Linear SVM with more Examples – Scenario-2

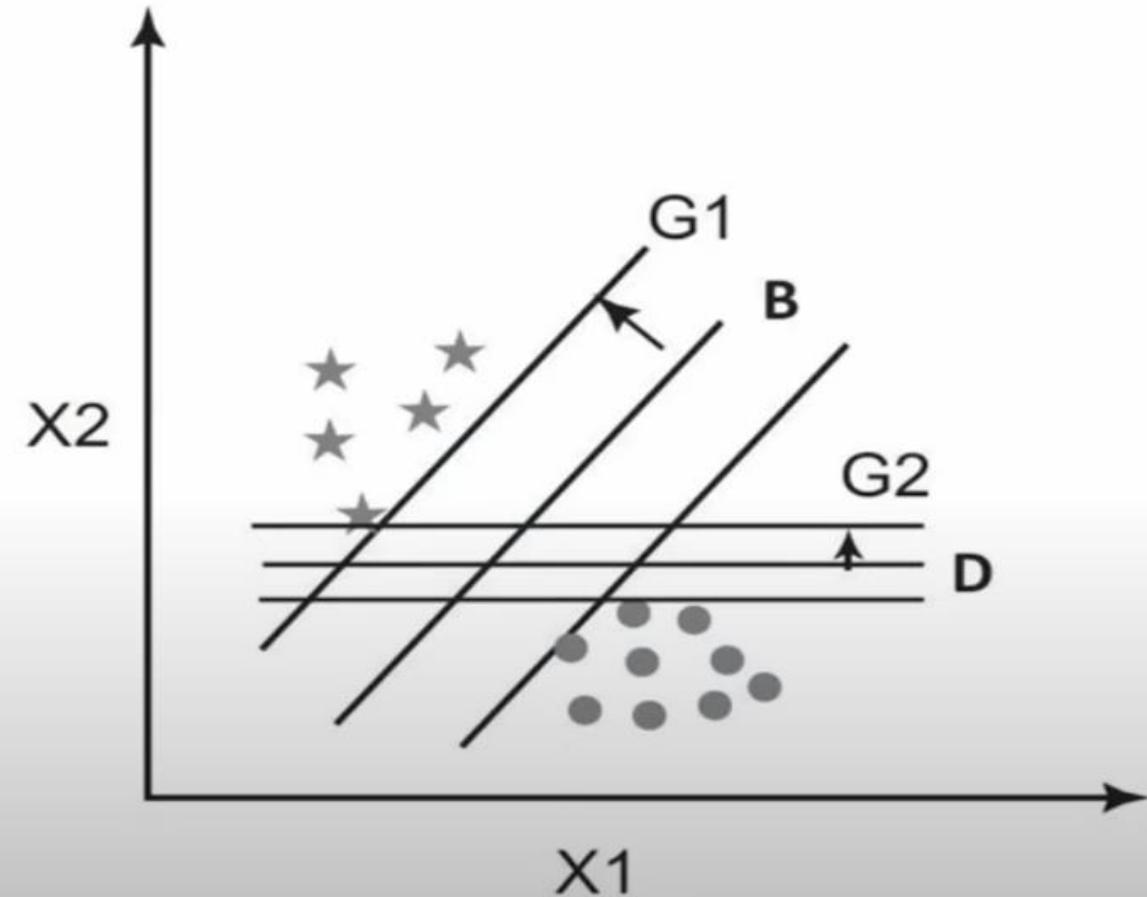
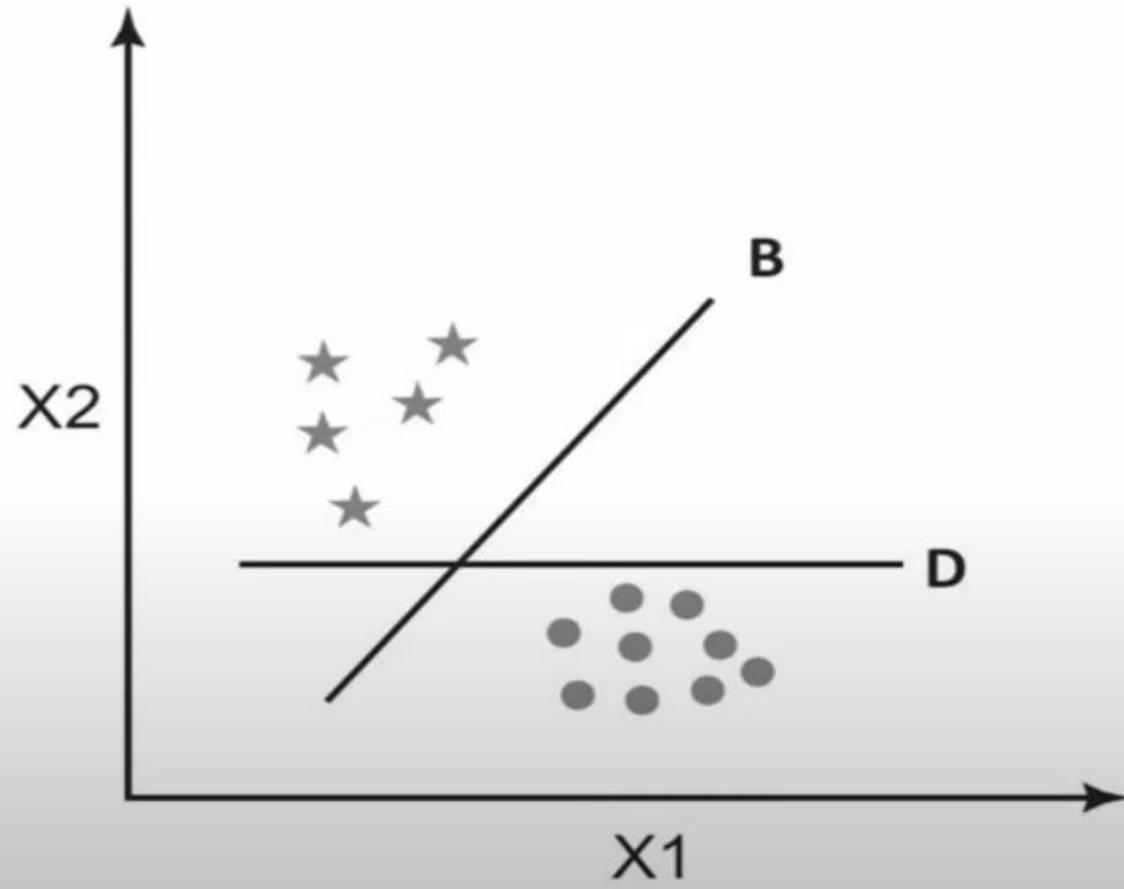
How does it work?



LINEAR AND NON-LINEAR SVM

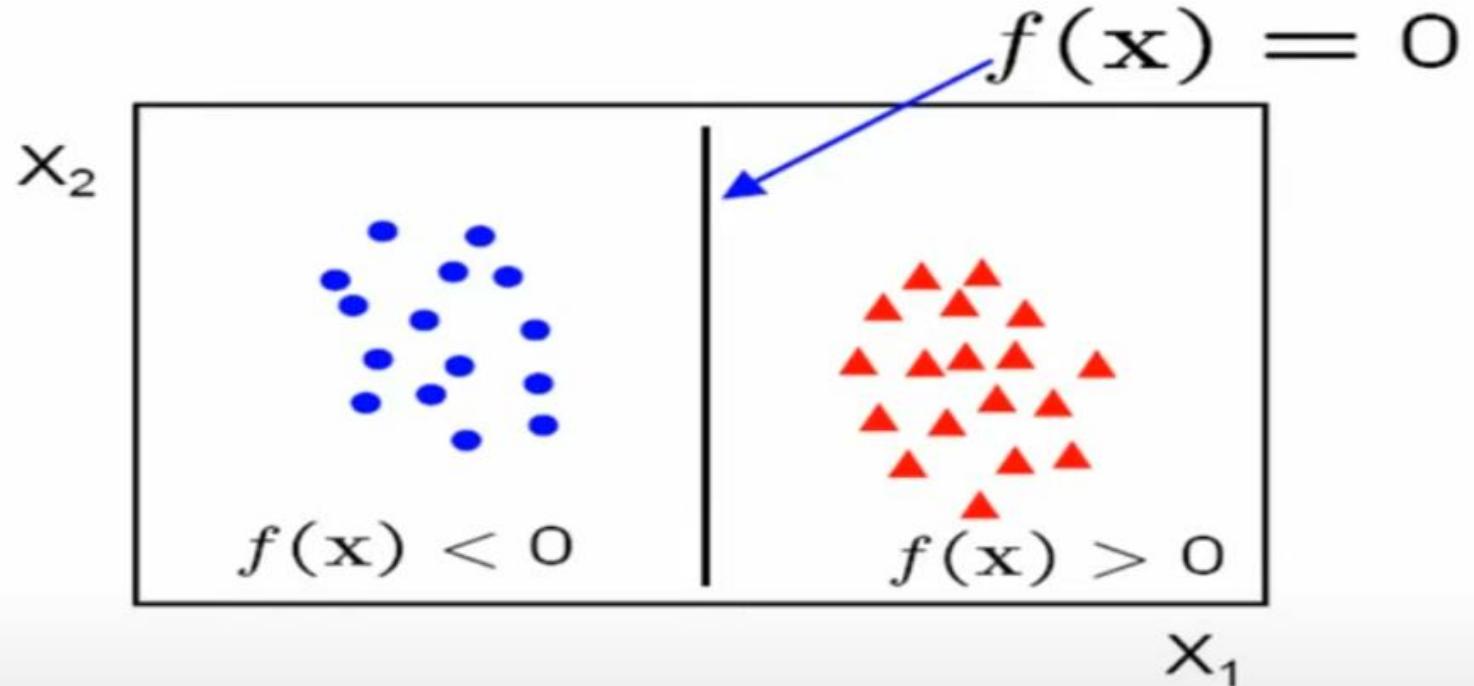
Working of Linear SVM with more Examples – Scenario-3

How does it work?



LINEAR AND NON-LINEAR SVM

Working of Linear SVM: (Mathematical Representation)

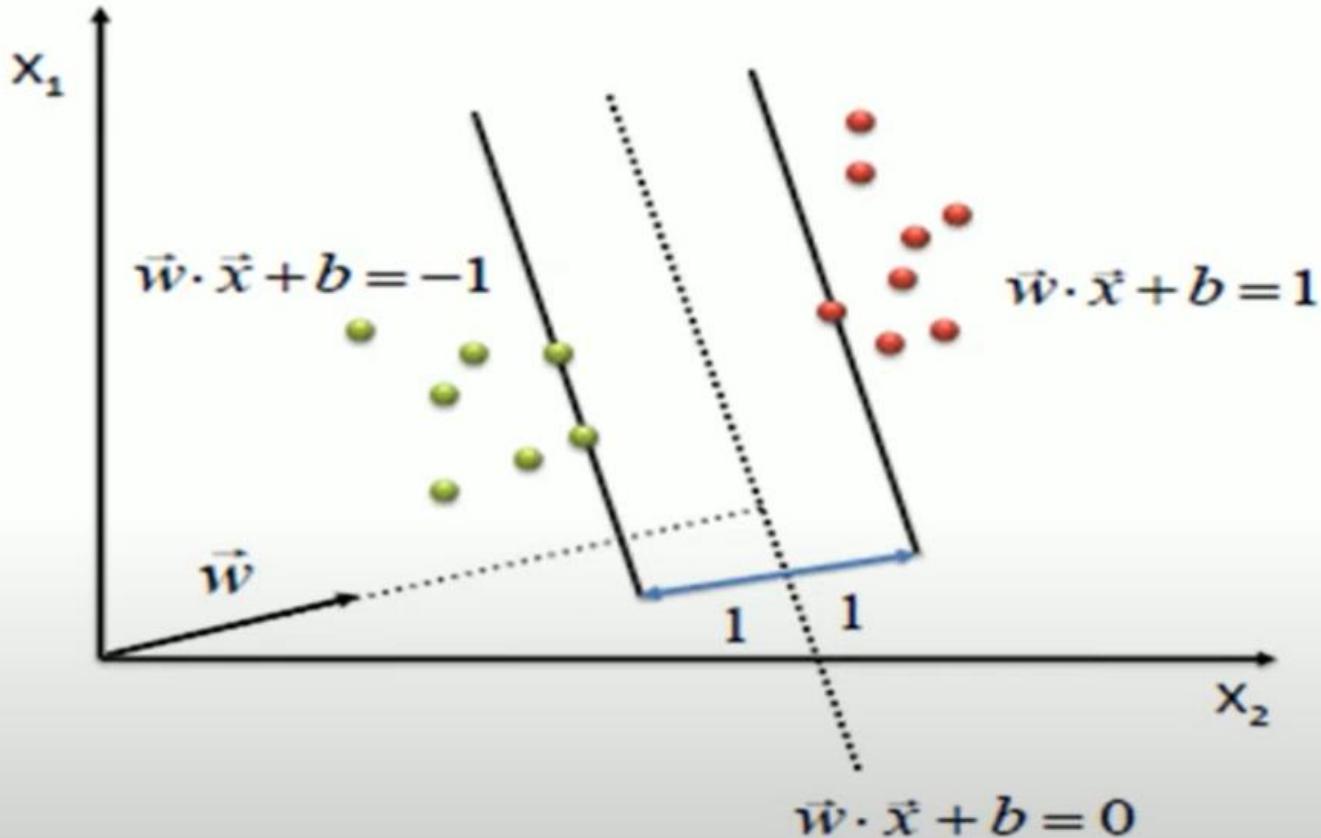


- $f(\mathbf{x}) = \mathbf{W} \cdot \mathbf{X} + b$
- \mathbf{W} is the normal to the line, \mathbf{X} is input vector and b the bias
- \mathbf{W} is known as the weight vector

LINEAR AND NON-LINEAR SVM

Working of Linear SVM: (mathematical representation):

SVM Model



$$\max \frac{2}{\|\vec{w}\|}$$

s.t.

$$(\vec{w} \cdot \vec{x} + b) \geq 1, \forall x \text{ of class 1}$$

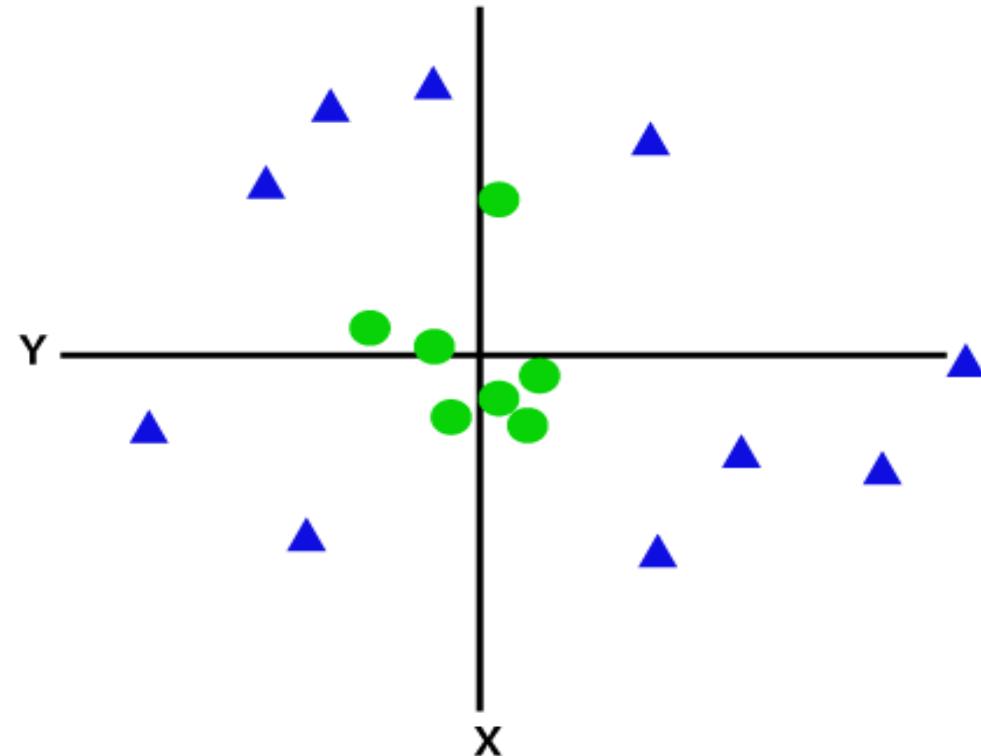
$$(\vec{w} \cdot \vec{x} + b) \leq -1, \forall x \text{ of class 2}$$

LINEAR AND NON-LINEAR SVM

Working of Non-Linear SVM with Example

If data is linearly arranged, then we can separate it by using a straight line, but **for non-linear data, we cannot draw a single straight line.**

Consider the given image:

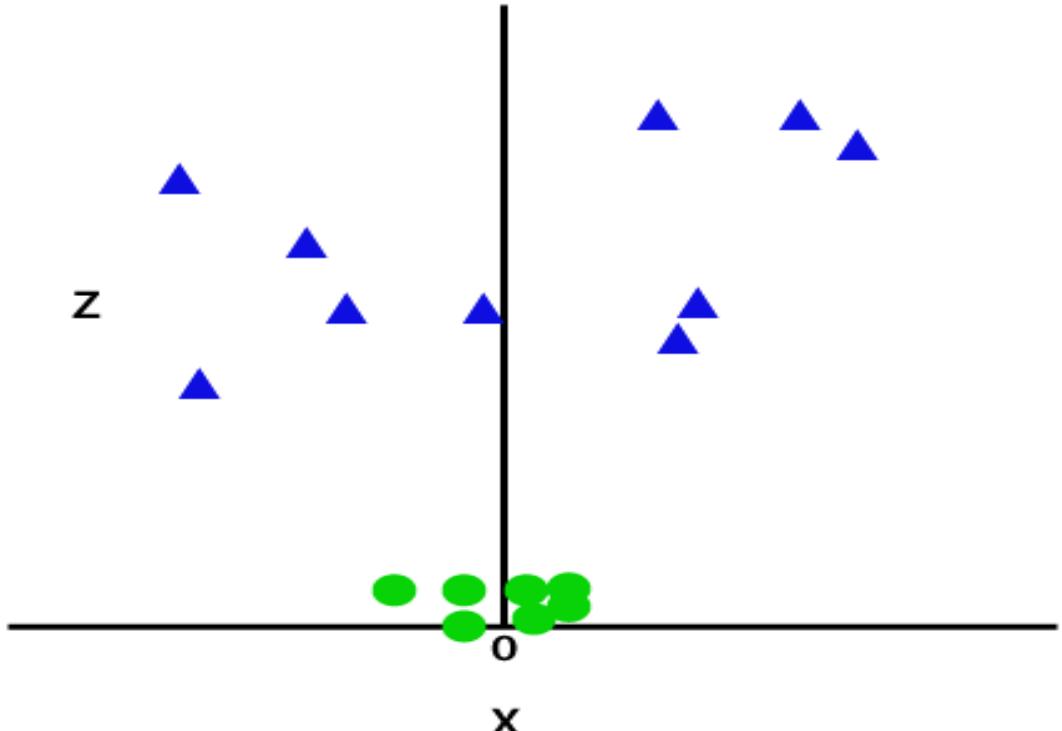


LINEAR AND NON-LINEAR SVM

Working of Non-Linear SVM with Example

So, to separate these data points, we need to add one more dimension. For linear data, we have used two dimensions x and y, so **for non-linear data, we will add a third-dimension z**. It can be calculated as:

$$z=x^2 + y^2$$



LINEAR AND NON-LINEAR SVM

Working of Non-Linear SVM with Example

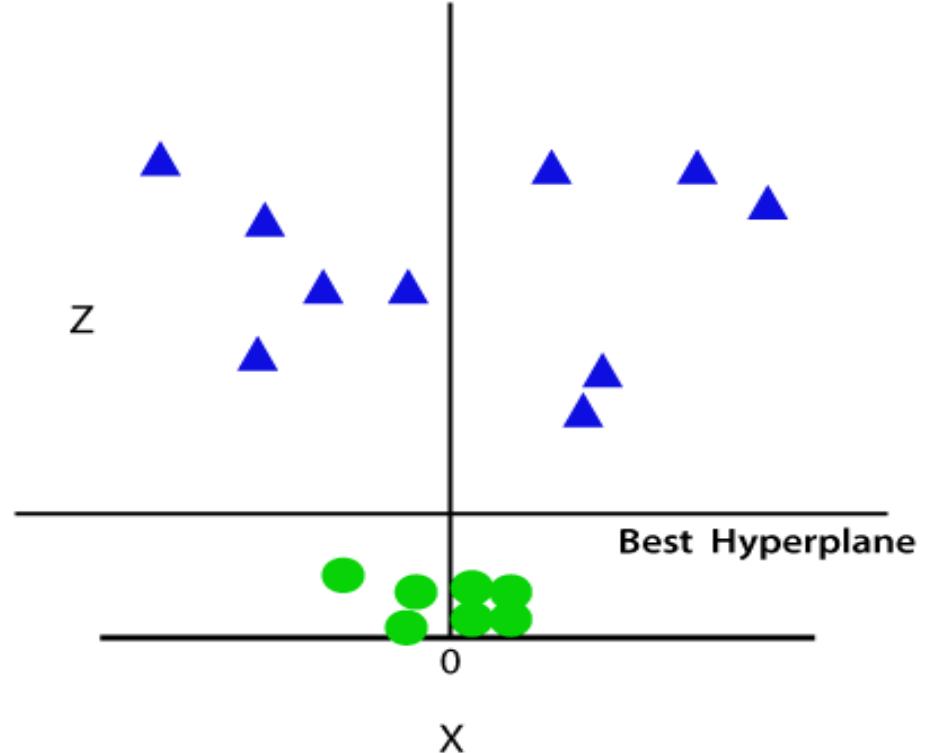
In the above plot, points to consider:

- **All values of z would be positive always** because z is the squared sum of both x and y.
- In the original plot, green circles appear close to the origin of x and y axes, leading to lower value of z and blue triangles relatively away from the origin result to higher value of z. In SVM, it is easy to have a linear hyperplane between these two classes. But another question arises that, should we need to add this feature manually to have a hyperplane.
- No, SVM has a technique called the **KERNEL TRICK. (SVM KERNEL).**
- **SVM Kernel are the functions which takes low dimensional input space and transform it to a higher dimensional space** i.e. it converts non separable problem to separable problem, these functions are called **kernels**. It is mostly useful in non-linear separation problem.

LINEAR AND NON-LINEAR SVM

Working of Non-Linear SVM with Example

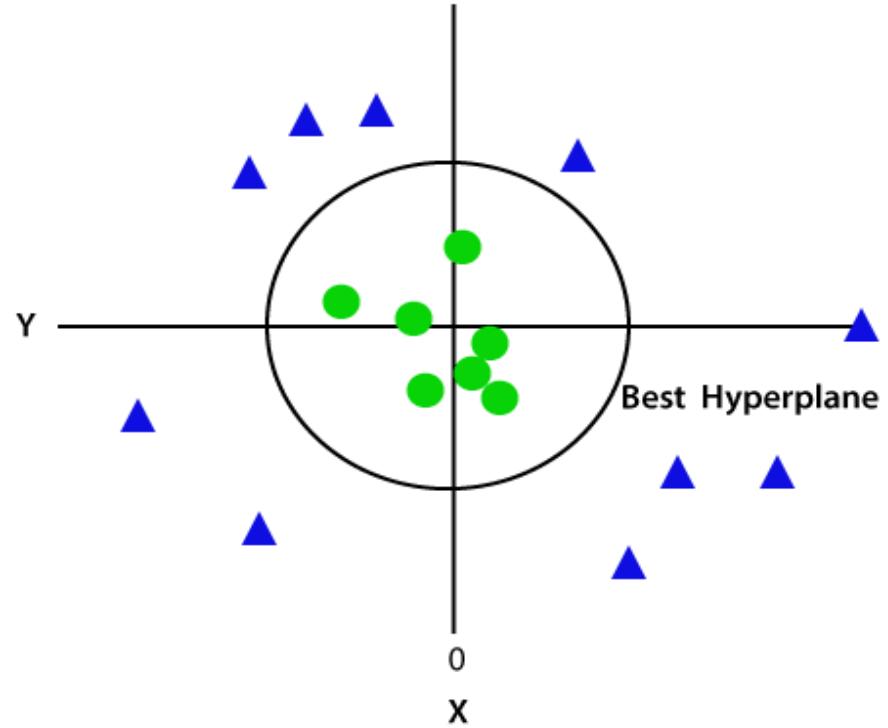
So now, SVM will divide the datasets into classes in the following way. Consider the given image:



LINEAR AND NON-LINEAR SVM

Working of Non-Linear SVM with Example

Since we are in 3-d Space, hence it is looking like a plane parallel to the x-axis. If we convert it in 2d space with $z=1$, then it will become as: (i.e. when we look at the hyperplane in original input space it looks like a circle)



Hence, we get a circumference of radius 1 in case of non-linear data.

LINEAR AND NON-LINEAR SVM

Data Science Libraries to implement SVM

- SciKit Learn
- PyML
- SVM^{struct} Python
- LIBSVM

K-NEAREST NEIGHBORS (KNN)

KNN is a reasonably simple classification technique that identifies the class in which a sample belongs by measuring its similarity with other nearby points.

It is a powerful *technique for identifying the class of an unknown sample point*.

K-nearest neighbors (KNN) algorithm is a *type of supervised ML algorithm* which can be *used for both classification as well as regression predictive problems. However, it is mainly used for classification predictive problems in industry*.

The following two properties would define KNN well –

- **Lazy learning algorithm** – KNN is a lazy learning algorithm because it does *not have a specialized training phase and uses all the data for training while classification*.
- **Non-parametric learning algorithm** – KNN is also a non-parametric learning algorithm because it *doesn't assume anything* about the underlying data.

KNN ALGORITHM

The steps for the KNN Algorithm in Machine Learning are as follows:

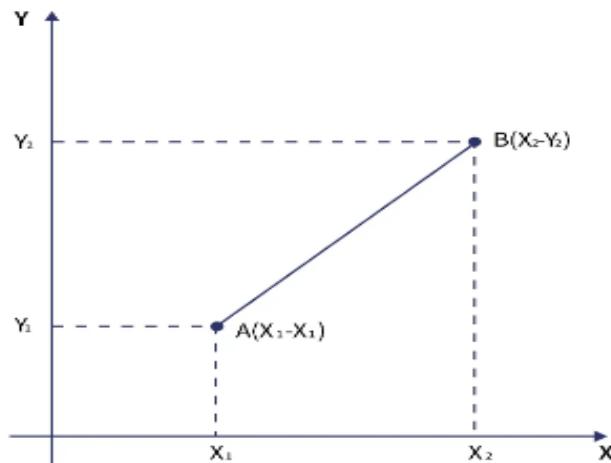
- **Step - 1 :**
Select the number K of the neighbors
- **Step - 2 :**
Calculate the Euclidean distance of each point from the target point.
- **Step - 3 :**
Take the K nearest neighbors per the calculated Euclidean distance.
- **Step - 4 :**
Among these k neighbors, count the number of the data points in each category.
- **Step - 5 :**
Assign the new data points to that category for which the number of neighbors is maximum.

EXAMPLE

Suppose we have a new data point and need to put it in the required category. Consider the below image :

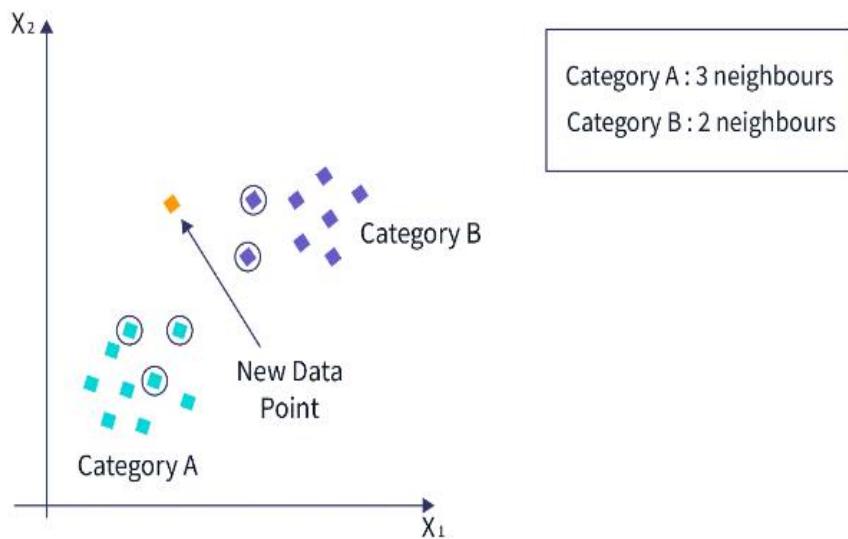


- First, we will choose the number of neighbors as $k=5$.
- Next, we will calculate the Euclidean distance between the data points. The Euclidean distance is the distance between two points, which we have already studied in geometry. It can be calculated as :



$$\text{Euclidean Distance } A_1 \text{ and } B_2 = \sqrt{(X_2 - X_1)^2 + (Y_2 - Y_1)^2}$$

- By calculating the Euclidean distance, we got the nearest neighbors, as three nearest neighbors in category A and two nearest neighbors in category B. Consider the below image :



New data point will fall in category A.

HOW TO DETERMINE THE K VALUE?

- The value of k is very crucial in the KNN algorithm to define the number of neighbors in the algorithm.
- The value of k in the k-nearest neighbors (k-NN) algorithm should be chosen based on the input data. If the input data has more outliers or noise, a higher value of k would be better.
- It is **recommended to choose an odd value for k** to avoid ties in classification.
- **Cross-validation** methods can help in selecting the best k value for the given dataset.

DISTANCE METRICS USED IN KNN

Euclidean Distance $d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$

Manhattan Distance $d(x, y) = \sum_{i=1}^n |x_i - y_i|$

Minkowski Distance $d(x, y) = \left(\sum_{i=1}^n (x_i - y_i)^p \right)^{\frac{1}{p}}$

APPLICATIONS OF KNN

- KNN can be **used in plagiarism checks** to check if two documents are semantically identical.
- ML in **healthcare uses KNN to diagnose diseases based on subtle symptoms**.
- Netflix uses KNN in its **recommendation systems**.
- KNN is used in **face recognition software**.

- **Political scientists and psychologists use** KNN to classify people based on their behavior and response to stimuli.
- **Banking System**

KNN can be used in banking system to predict whether an individual is fit for loan approval? Does that individual have the characteristics similar to the defaulters one?

- **Calculating Credit Ratings**

KNN algorithms can be used to find an individual's credit rating by comparing with the persons having similar traits.

- **Politics**

With the help of KNN algorithms, we can classify a potential voter into various classes like "Will Vote", "Will not Vote", "Will Vote to Party 'Congress'", "Will Vote to Party 'BJP'".

ADVANTAGES OF KNN

- It is very **simple algorithm to understand and interpret**.
- It is very **useful for nonlinear data because there is no assumption** about data in this algorithm.
- It is a **versatile algorithm** as we can use it for classification as well as regression.
- It has **relatively high accuracy** but there are much better supervised learning models than KNN.

DISADVANTAGES OF KNN

- It is **computationally a bit expensive algorithm** because it stores all the training data.
- **High memory storage required** as compared to other supervised learning algorithms.
- **Prediction is slow** in case of big N.
- It is very **sensitive to the scale of data as well as irrelevant features**.

PROBLEM: (Prediction of Sugar of Diabetic Patient given BMI and Age)

Apply K nearest neighbor classifier to predict the diabetic patient with the given features BMI, Age. If the training examples are,

BMI	Age	Sugar
33.6	50	1
26.6	30	0
23.4	40	0
43.1	67	0
35.3	23	1
35.9	67	1
36.7	45	1
25.7	46	0
23.3	29	0
31	56	1

Assume K=3,

Test Example BMI=43.6, Age=40, Sugar=?

Solution:

The given training dataset has 10 instances with two features BMI (Body Mass Index) and Age. Sugar is the target label. The target label has two possibilities 0 and 1. 0 means the diabetic patient has no sugar and 1 means the diabetic patient has sugar.

Given the dataset and new test instance, we need to find the distance from the new test instance to every training example. Here we use the Euclidean distance formula to find the distance.

$$Distance = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

In the next table, you can see the calculated distance from text example to training instances.

BMI	Age	Sugar	Formula	Distance
33.6	50	1	$\sqrt{(43.6-33.6)^2+(40-50)^2}$	14.14
26.6	30	O	$\sqrt{(43.6-26.6)^2+(40-30)^2}$	19.72
23.4	40	O	$\sqrt{(43.6-23.4)^2+(40-40)^2}$	20.20
43.1	67	O	$\sqrt{(43.6-43.1)^2+(40-67)^2}$	27.00
35.3	23	1	$\sqrt{(43.6-35.3)^2+(40-23)^2}$	18.92
35.9	67	1	$\sqrt{(43.6-35.9)^2+(40-67)^2}$	28.08
36.7	45	1	$\sqrt{(43.6-36.7)^2+(40-45)^2}$	8.52
25.7	46	O	$\sqrt{(43.6-25.7)^2+(40-46)^2}$	18.88
23.3	29	O	$\sqrt{(43.6-23.3)^2+(40-29)^2}$	23.09
31	56	1	$\sqrt{(43.6-31)^2+(40-56)^2}$	20.37

Once you calculate the distance, the next step is to find the nearest neighbors based on the value of k. In this case, the value of k is 3. Hence, we need to find 3 nearest neighbors.

BMI	Age	Sugar	Distance	Rank
33.6	50	1	14.14	2
26.6	30	O	19.72	
23.4	40	O	20.20	
43.1	67	O	27.00	
35.3	23	1	18.92	
35.9	67	1	28.08	
36.7	45	1	8.52	1
25.7	46	O	18.88	3
23.3	29	O	23.09	

31	56	1	20.37	
----	----	---	-------	--

Now, we need to **apply the majority voting technique** to decide the resulting label for the new example. Here the 1st and 2nd nearest neighbors have target label 1 and the 3rd nearest neighbor has target label 0. Target label 1 has the majority. Hence the new example is classified as 1, That is the diabetic patient has **Sugar**.

Test Example BMI=43.6, Age=40, Sugar=1