

$$\textcircled{1} \sqrt{\underset{2}{(6-4)}^2 + \underset{4}{(8-3)}^2} = \sqrt{29} = 5.38$$

$$\textcircled{2} \sqrt{\underset{1}{(6-6)}^2 + \underset{1}{(8-7)}^2} = \textcircled{1} - \text{Pass}$$

$$\textcircled{3} \sqrt{\underset{1}{(6-7)}^2 + \underset{0}{(8-8)}^2} = \textcircled{1} - \text{Pass} \quad \text{see the nearest value}$$

$$\textcircled{4} \sqrt{\underset{2}{(6-5)}^2 + \underset{3}{(8-5)}^2} = \sqrt{10} = 3.16$$

$$\textcircled{5} \sqrt{(6-8)^2 + (8-8)^2} = \textcircled{2} - \text{Pass}$$

for given  $x \rightarrow$   
consider 3 value

As  $K=3$  we have to

$\Rightarrow P \quad P \quad P = 3$  [based on nearest neighbour]  
fail  $\textcircled{0}$

$$\boxed{3 > 0}$$

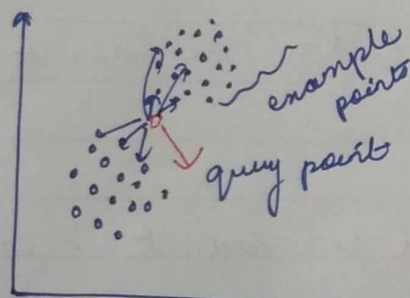
So  $x$  is declare to be Pass

\*  $KNN \rightarrow$  It is a supervised learning algorithm in which you have some  $\sim$  data points or data vector which is separated into different no of several categories or classes and it try to predict the classification of a new sample from a particular population set

⇒ A lazy algo<sup>n</sup> → (It tries to only memorize the process. It does not learn any process)

(means it does not take its own decision)  
suppose you call a KNN algorithm to go there  
it will go there you will call it to come here  
it will come to that point

⇒ It classifies new points based on a similarity measure like - euclidean distance



Assume  $N = 30$  datapoints

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

Suppose we choose any point and  
we find the neighbourhood of that point  
It finds out the distance between query point and  
example point with help of euclidean distance  
or Minkowski distance

\* -  $K$  must be odd eg (3, 5, 7)

→  $K$  must not be a multiple of classes

classification if  $K$  is even then there will be tie in the  
of that particular label

eg

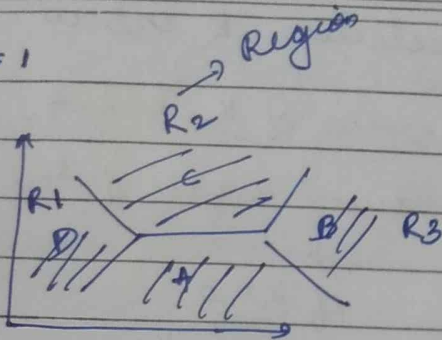
if

$C = 7$

$K \neq 14$



\* if  $K=1$



Then the particular data points will be assigned certain region into the partition

Basically it partitions the space into  $n$  no. of regions. This is called vector partition space.

It means if there is any data element falling in  $R_1$  will be assigned to region  $D$  class label. similarly  $R_2 \rightarrow \underline{C}$  class label

Algo<sup>n</sup>

- S1: load data from csv or xls file
- S2: Initialize  $K$  hyperparameters [assign nearest neighbour to it] (any odd no)
- S3: for each sample in training data.
  - S31: calculate distance between query point & the current point
  - S32: Add the distance & the index of the example to an ordered collection
- S4: Sort  $\rightarrow$  ordered collection of distance & indexes from small to large. (ascending order)

Teacher's Signature \_\_\_\_\_

S<sub>5</sub> - Pick first  $K$  entries from sorted collection

S<sub>6</sub> - Get the labels of selected  $K$  entries

if regression  $\rightarrow$  return mean of  $K$  labels

if classification  $\rightarrow$  return mode of  $K$  labels

$\downarrow$   
if  $\underline{K=3}$   
 $\underline{\text{mode}=1}$

### Advantages of KNN algo<sup>n</sup>

- $\rightarrow$  It is simple to implement
- $\rightarrow$  It is robust to the noisy training data
- $\rightarrow$  It can be more effective if the training data is large

### Disadvantages of KNN algo<sup>n</sup>

- $\rightarrow$  Always needs to determine the value of  $K$  which may be complex some time.
- $\rightarrow$  The computation cost is high because of calculating the distance between the data points for all the training sample

\* How to select the value of  $K$  in the K-NN Algo<sup>n</sup>?

Below are some points to remember while selecting the value of  $K$  in the K-NN algo<sup>n</sup>.



- There is no particular way to determine the best value for "K" so we need to try some values to find the best out of them. The most preferred value for K is 5.
- A very low value for K such as  $K=1$  or  $K=2$ , can be noisy and lead to the effects of outlier in the model.
- Large values for K are good, but it may find some difficulties.

### \* Support Vector Machine

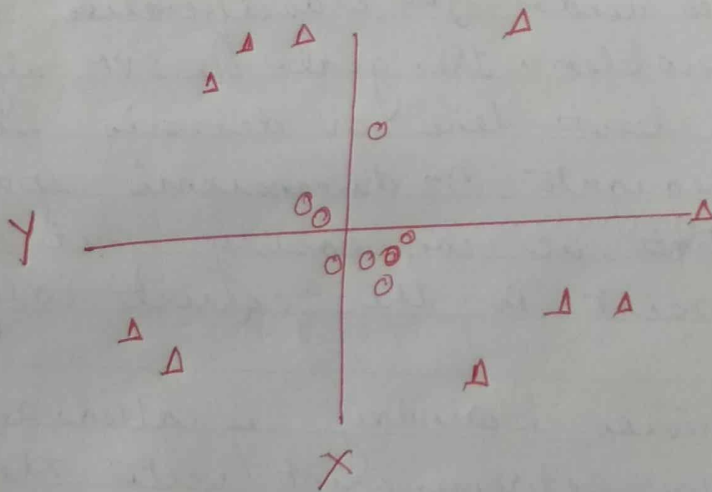
It is one of the most popular supervised learning algo<sup>n</sup> which is used for classification as well as Regression problem. The goal of SVM algo<sup>n</sup> is to create the best line or decision boundary that can segregate  $n$ -dimensional space into classes so that we can easily put the new data point in the correct category in future.

- ⇒ The best decision boundary is called a hyperplane.
- ⇒ SVM chooses the extreme point/vector the help is creating the hyperplane. These extreme cases are called as support vector. Hence the algo<sup>n</sup> is termed as SVM.

## Types of SVM

① 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 & classifier is used called as linear SVM classifier.

② Non-linear SVM - Non-linear SVM is used for non-linearly separable 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.

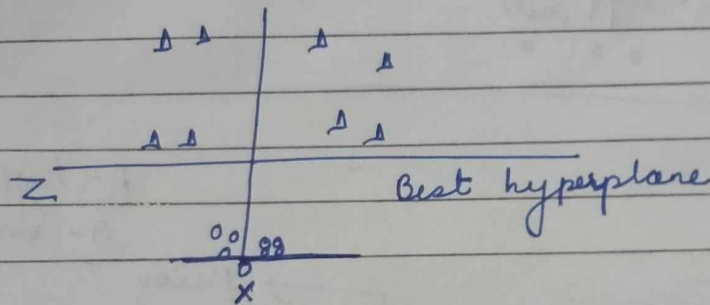


\* if data is linearly arranged, then we can separate by using a straight line. but for non-linear data, we cannot draw a single straight line

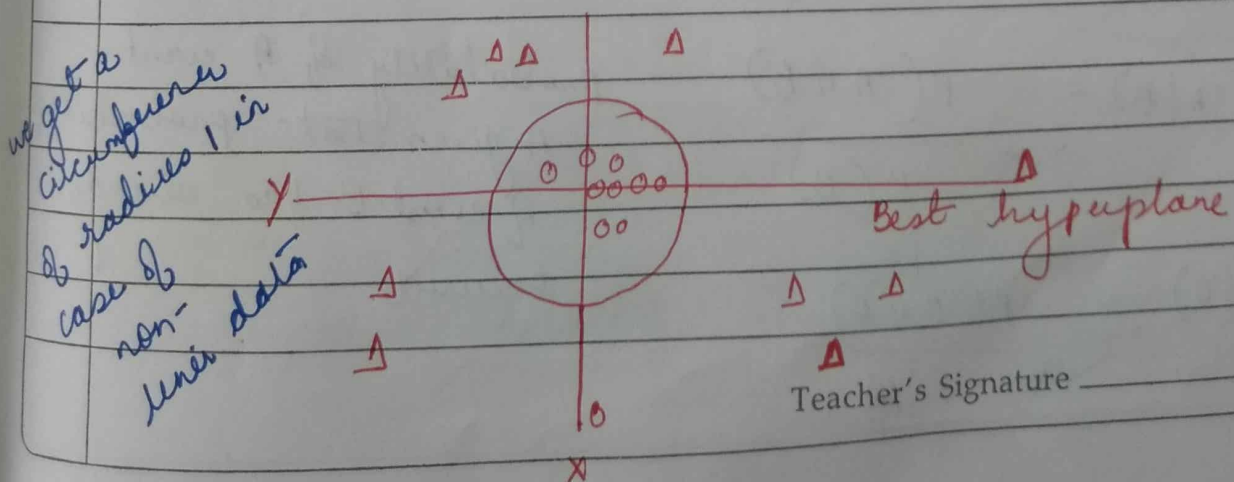
→ So to separate these data points, we need to add one more dimension

$$x, y, z$$

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

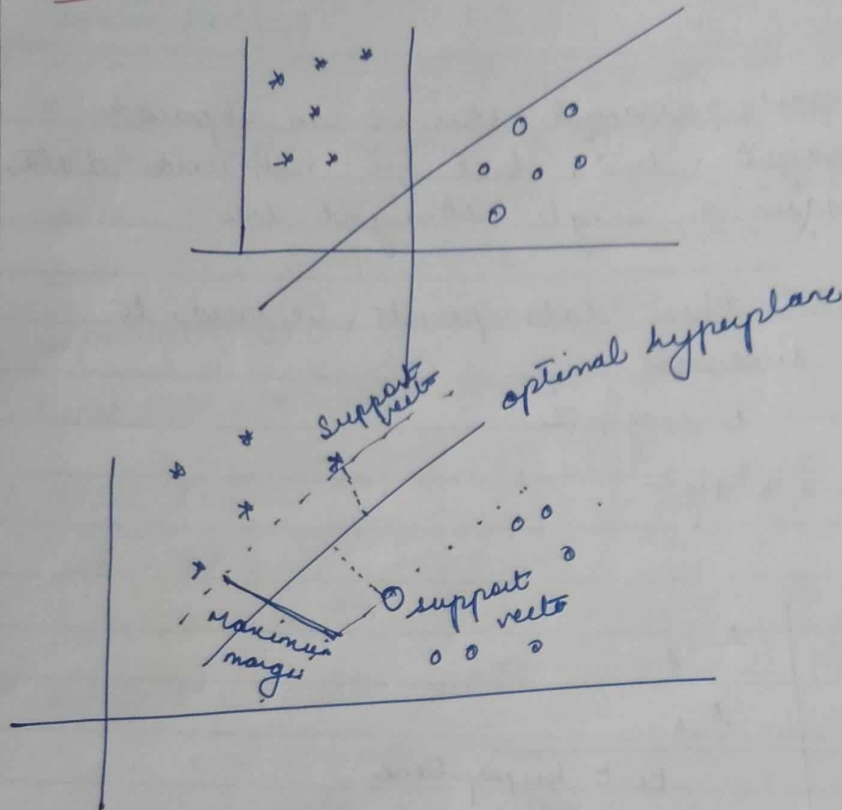


Since we are in 3-d space, hence it is looking like a plane parallel to the  $x$ -axis. If we convert it in 2d space with  $z=1$ , then it will become a





Linear SVM



Bayes Theorem

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

← Posterior  
 ← likelihood  
 ← marginal  
 ← Prior

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(B|A) = \frac{P(B \cap A)}{P(A)}$$

probability of A event  
A gives that probability  
of event B has already  
occurred

A - hypothesis  
B - evidence / data