

UNIT IV CLASSIFICATION AND CLUSTERING 10

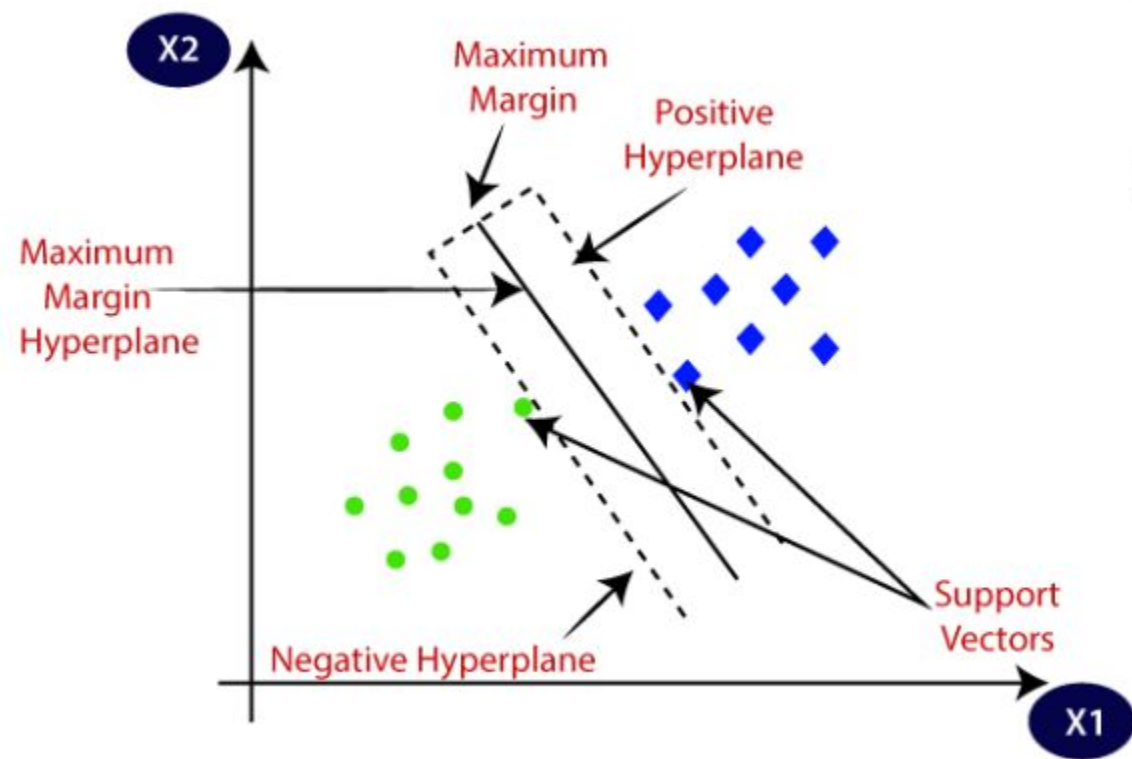
Decision Tree Induction - Bayesian Classification – Rule Based Classification – Classification by Back propagation – Support Vector Machines – Associative Classification – Lazy Learners – Other Classification Methods – Clustering techniques – , Partitioning methods- k-means- Hierarchical Methods – distance based agglomerative and divisible clustering, Density-Based Methods – expectation maximization -Grid Based Methods – Model-Based Clustering Methods – Constraint – Based Cluster Analysis – Outlier Analysis

Support Vector Machines

- 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**.

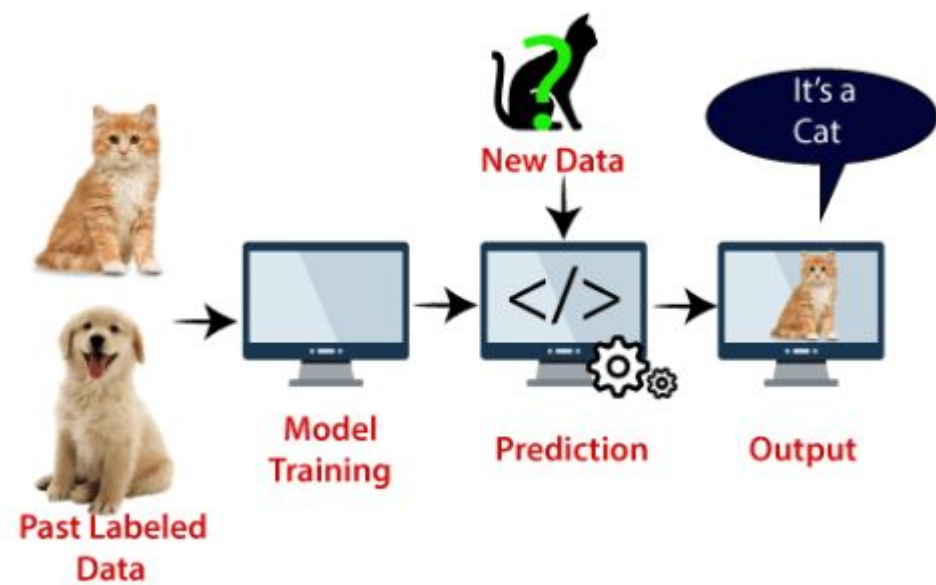
- The **goal of the SVM algorithm is to create the best line or decision boundary that can segregate n-dimensional space into classes** so that we can easily put the new data point in the correct category in the future. This best decision boundary is called a hyperplane.
- **SVM chooses the extreme points/vectors that help in creating the hyperplane.** These extreme cases are called as support vectors, and hence algorithm is termed as Support Vector Machine. Consider the below diagram in which there are two different categories that are classified using a decision boundary or hyperplane:

-



Example:

- SVM can be understood with the **example that we have used in the KNN classifier.**
- 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



- SVM algorithm can be used for **Face detection, image classification, text categorization**, etc.

Types of SVM

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

Hyperplane and Support Vectors in the SVM algorithm:

- **Hyperplane:** 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**.
- We always create a hyperplane that has a maximum margin, which means the maximum distance between the data points.

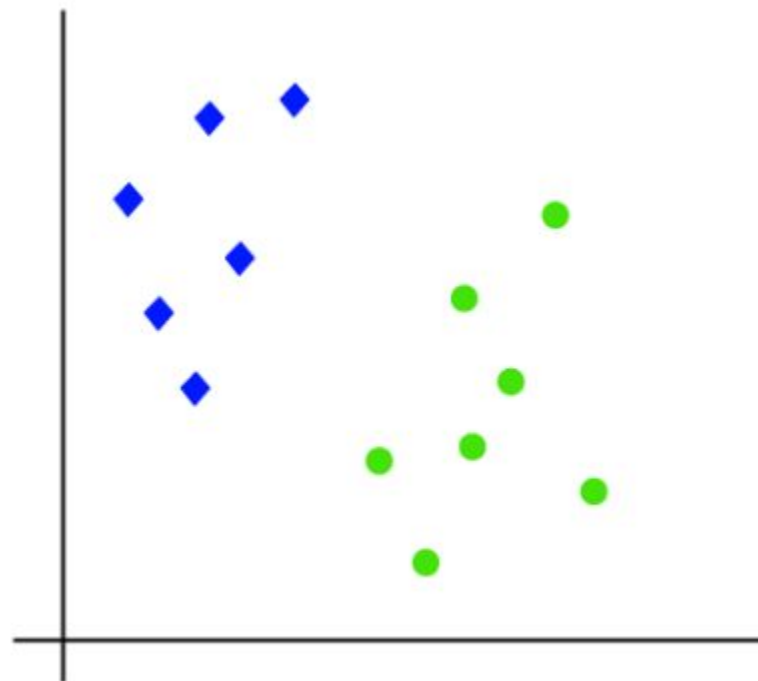
- **Support Vectors:**

- 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.

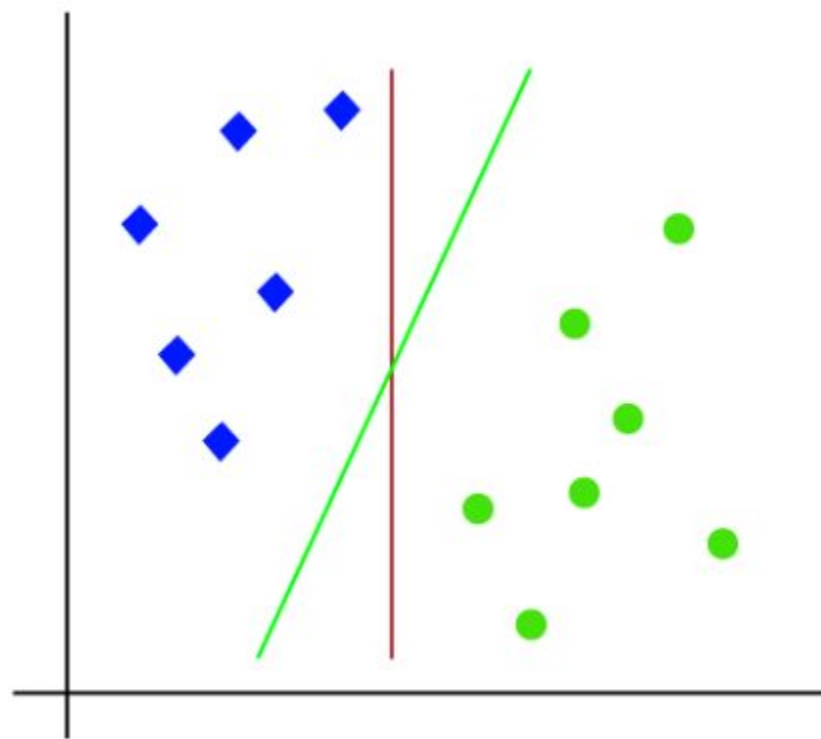
How does SVM works?

- **Linear SVM:**

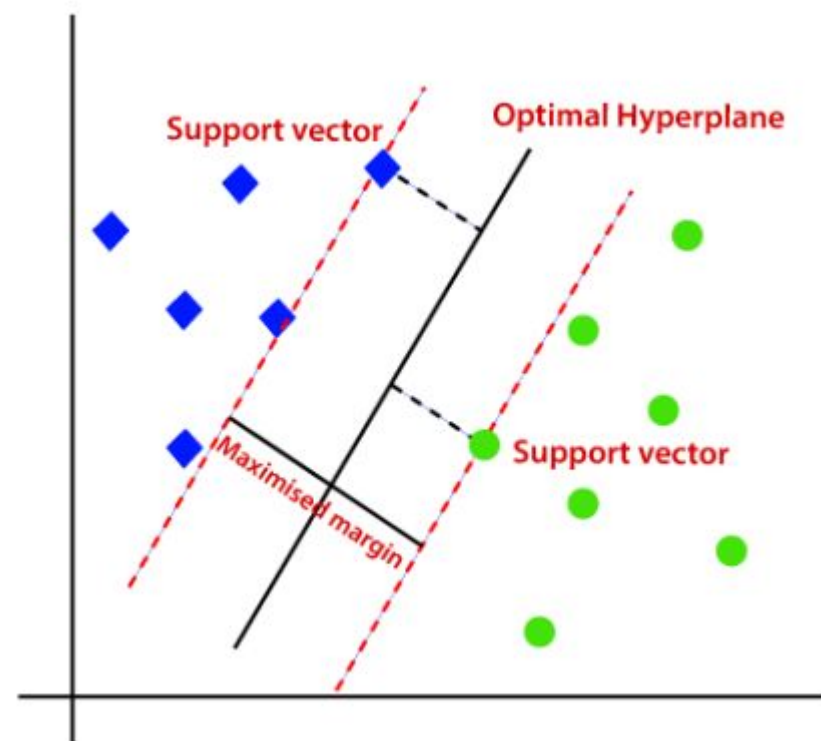
- The working of the SVM algorithm can be understood by using an example.
- Suppose we have a dataset that has two tags (green and blue), and the dataset has two features x_1 and x_2 . We want a classifier that can classify the pair(x_1 , x_2) of coordinates in either green or blue. Consider the below image:
 -



- So as it is 2-d space so 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 below image:



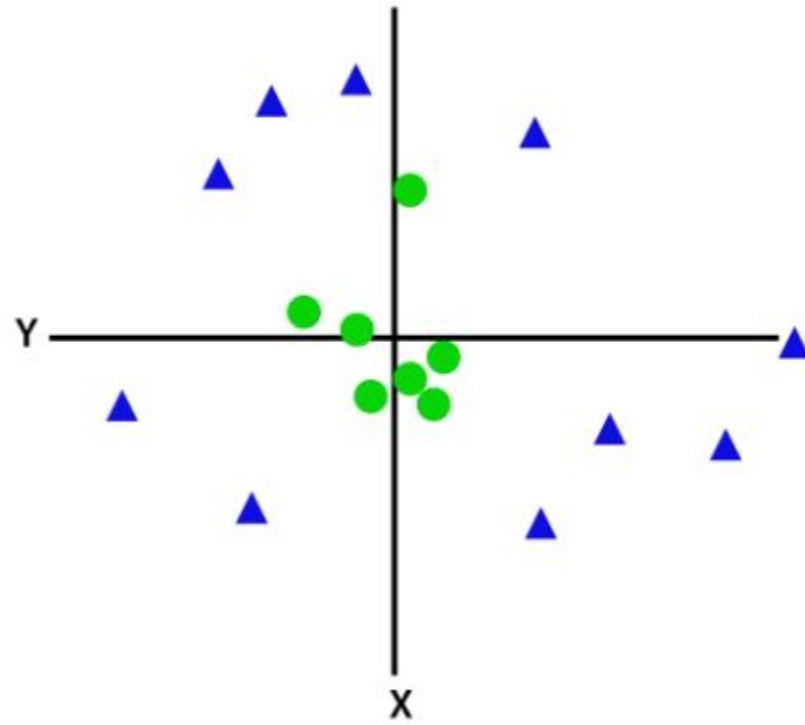
- Hence, the SVM algorithm helps to find the best line or decision boundary; this best boundary or region is called as a **hyperplane**. SVM algorithm finds the closest point of the lines from both the classes. These points are called support vectors. The distance between the vectors and the hyperplane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



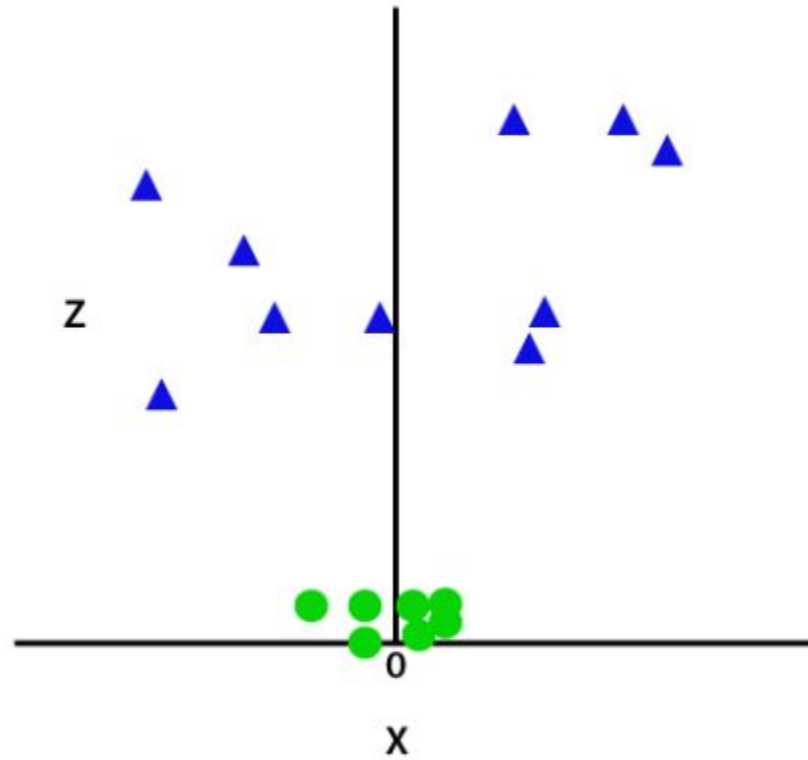
Non-Linear SVM:

- 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 below image:

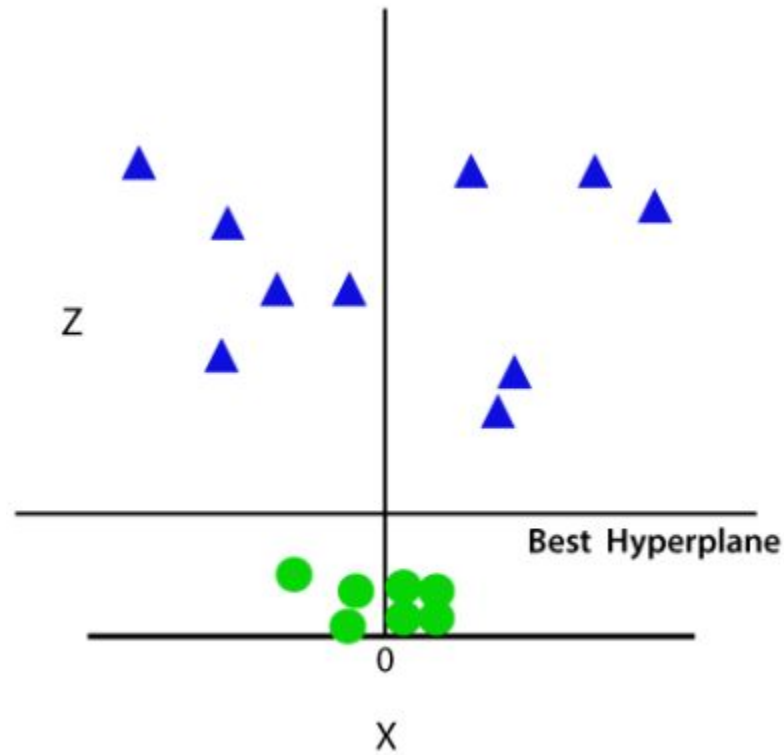
$$z = x^2 + y^2$$



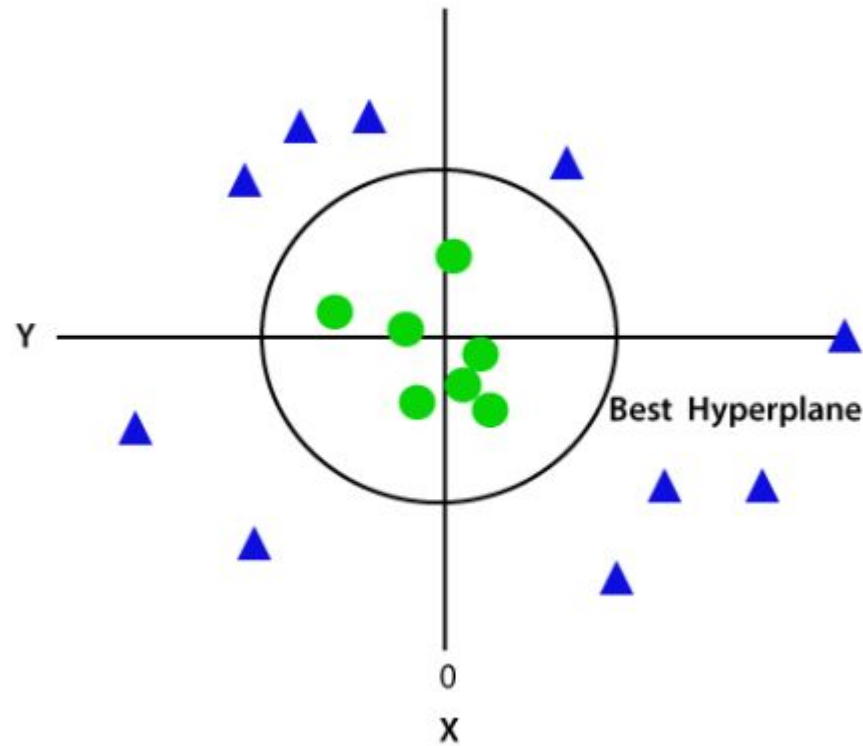
By adding the third dimension, the sample space will become as below image:



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



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:



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