

Support Vector Machine

Prerequisite

- Logistic regression and classification problem.
- Evaluation metrics for classification problems.

Objectives

- Understanding of linear classification and margin classification.
- Understanding the terms like support vectors, soft margin, hyperplane.
- Kernel and kernel trick.

Support Vector Machine

SVM is one of the popular supervised machine learning methods that can be equally used for classification and regression, but SVM is mostly used for classification. The principle of SVM is to find an hyperplane which can classify the training data points into labelled categories. The input of SVM is the training data and uses this training sample point to predict the class of test points.

Consider the following set of points of two classes shown in the graph.

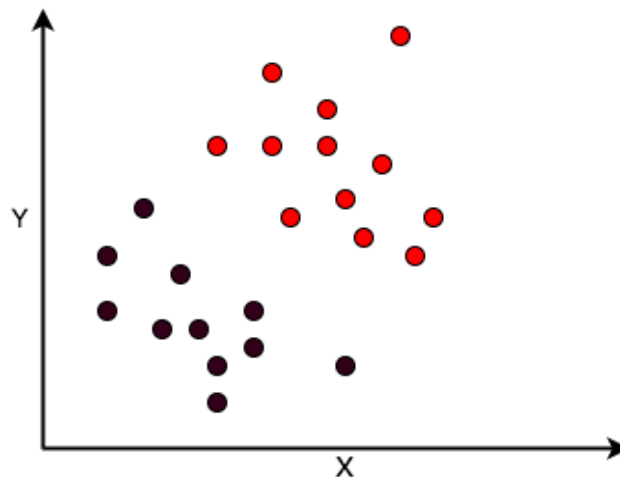


Figure 1 Problem of classification

By looking at the figure we can see that the points can be separated using a hyperplane(line) where + class points are above the line and – class will be below the line. Here we need to remember that there can be many hyperplanes which separate the given points in different ways as shown in figure. Each of

the hyperplanes is valid as it separates the given points successfully. But, here our objective is to find the optimal hyperplane.

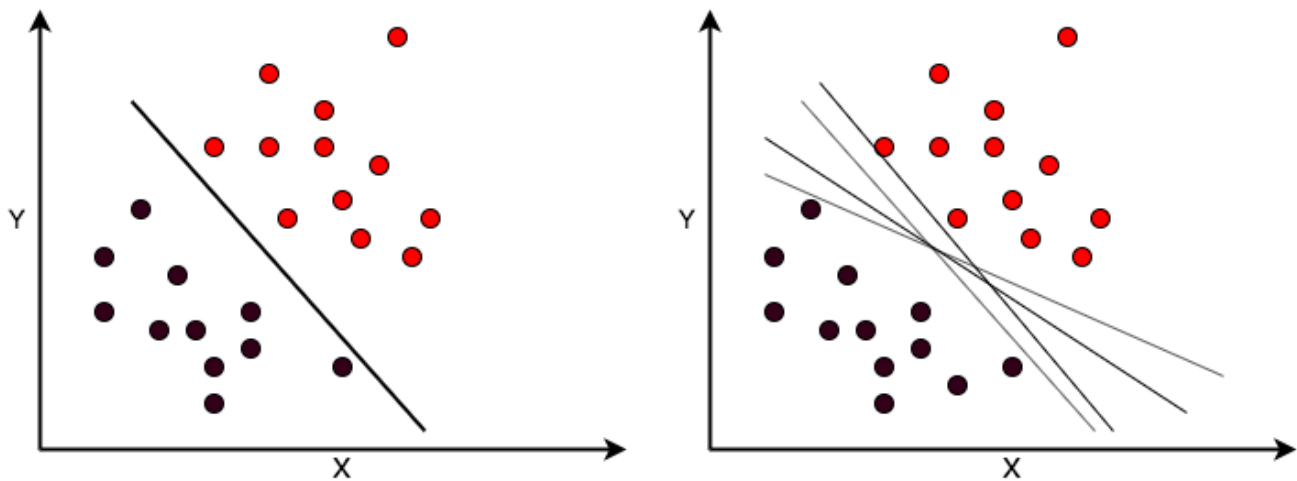


Figure 2 Possible separating hyperplanes

Support vector machine chooses the best hyperplane which is at the maximum distance from the data points from each category. For a given hyperplane, one can compute the distance between the closest data point and hyperplane from both classes. If we double the distance values, we will get a margin.

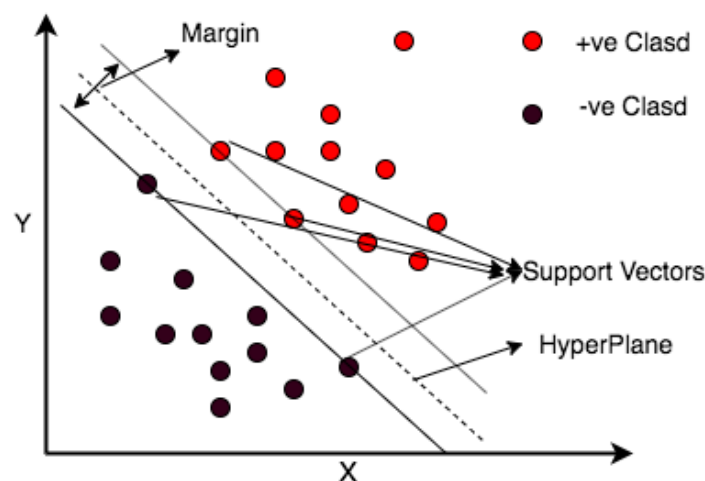


Figure 3 SVM Hyperplane

Margin space is known as no man's land where no data point is present. By looking at the figure we observe that margin width completely depends on how far the points are from the hyperplane. So, the optimal hyperplane is defined by the biggest margin and the objective of SVM is to find a hyperplane with maximum margin from training data. The problem of finding the optimal hyperplane is an optimization problem and can be solved by optimization techniques.

Let us consider the simplest case of linear classification of binary data. we have been given with the data points of two classes known as positive class and negative class. To classify the given data using a linear plane we need to identify a decision boundary of classification on which one side the positive points lie and another side negative point lie as shown in figure. Let anything above the decision boundary have label 1 and anything below have label -1. Mathematically for all data points X_i are subject to $W^T X_i + b > 0$ (H1) will have $y_i = 1$ and data points X_i are subject to $W^T X_i + b < 0$ (H2) will have $y_i = -1$. Here, $W^T X_i + b$ represents the plane (it would be a line in the form of $mx+c$ if we are considering 2-d plane) with slope as W^T and b as intercept. The points on the planes (H1 and H2) are the tips of the Support Vectors.

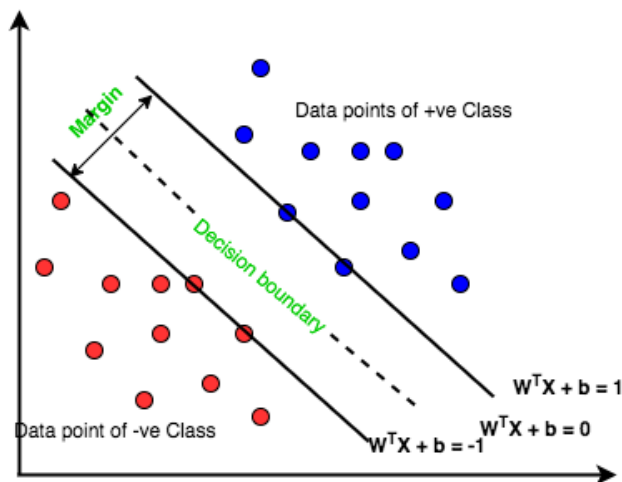


Figure 4 :Example of Binary Classification using SVM

As one can observe from the figure 4 that there is some space between the decision boundary and the nearest points of either class. Thus the decision boundary can be expanded both side and be defined as any point above the boundary $W^T X + b = 1$ will be of one class (label +1) and any point below the boundary $W^T X + b = -1$ will be of another class (label -1).

The margin is defined as the distance between these two boundaries. We want a classifier (linear separator) with as big a margin as possible. In algebra, the distance from a point (x_0, y_0) to a line: $Ax + By + c = 0$ is given by: $|Ax_0 + By_0 + c| / \sqrt{A^2 + B^2}$. So, the distance between H0 (optimal plane/Decision boundary) and H1 is then: $|W^T X + b| / \|W\| = 1 / \|W\|$, So the total distance between H1 and H2 is thus: $2 / \|W\|$. We can prove it as given below:-

As the two boundary lines are parallel to each other so, the distance (margin) can be found by picking an arbitrary point x_1 of line $W^T X + b = -1$ and looking for the closet point x_2 of line $W^T X + b = 1$. We get-

$$X_2 = X_1 + \lambda W \text{ (here, } \lambda \text{ represents the distance between H1 and H2)}$$

Since the closet point X_2 will always be present on perpendicular of line $W^T X + b = -1$.

So, we have

$$\begin{aligned} & W^T X_2 + b = 1 \text{ where } X_2 = X_1 + \lambda W \\ \Rightarrow & W^T (X_1 + \lambda W) + b = 1 \\ \Rightarrow & W^T X_1 + b + \lambda W^T W = 1 \text{ where } W^T X_1 + b = -1 \\ \Rightarrow & -1 + \lambda W^T W = 1 \\ \Rightarrow & \lambda W^T W = 2 \\ \Rightarrow & \lambda = \frac{2}{W^T W} = \frac{2}{\|W\|^2} \end{aligned}$$

So, the distance λ is $\frac{2}{\|W\|^2}$ and it is intuitive that the objective of support vector machine is to maximize the distance or equivalently tries to minimize the distance value $\frac{\|W\|^2}{2}$. The final equation is support vector machine is expressed as:

$$\text{Subject to: } y_i (W^T X_i + b) \geq 1 \text{ (} \forall \text{ data points } X_i \text{)}$$

In order to maximize the margin, we thus need to minimize $\|w\|$. With the condition that there are no data points between H1 and H2.

We can rewrite it as: min f: $\frac{1}{2} \|w\|^2$ with constraint as g: $y_i (W^T X_i) - b = 1$ or $[y_i (W^T X_i) - b] - 1 = 0$. Here, f is a quadratic function hence, minimizing f becomes a constrained optimization problem and it can be solved by the Lagrangian multiplier method.

Soft Margin

There are situations where the given data points are not perfectly linearly separable. (for instance, let the data points are correctly labelled). So, there is a need to define a hyperplane which allows some level of misclassification with data. This situation is defined by soft margin. This is

achieved through introducing slack variables $\epsilon_i \geq 0$ for each X_i and the Support vector machine equation is transformed into:

$$\frac{\|W\|^2}{2} + C \sum_i \epsilon_i$$

Subject to: $y_i(W^T X_i + b) \geq 1 - \epsilon_i$ and $\epsilon_i \geq 0$ (\forall data points X_i)

C is known as a regularization parameter which controls the limit of misclassification for each training example. The optimum value of C is found by cross validation and of multiple tries which gives lowest misclassification rate on testing data. In simple terms parameter C is used for controlling errors due to outliers. Low C value implies allowance of more outliers and high C value implies fewer outliers.

Kernel and Kernel Trick

Data classification can be done linear and nonlinear functions as shown in the following figure. The first figure is a rare case and most real-life situations we face a classification problem which is not linearly separable.

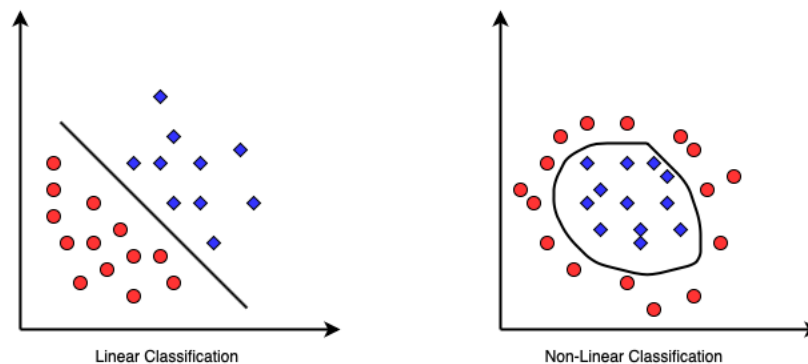


Figure 5: Data Classification

As we know, the support vector machine uses hyper planes to separate data points in the space but with nonlinear data SVM finds it difficult to generate hyperplanes. The solution of this problem is kernel trick.

Kernel function (Kernel trick) maps the original nonlinear separable data into the higher dimensions which helps to find the suitable hyperplane which separates the data points. It can be also understood as the function which maps a low dimensional space into higher dimensional space for making linear separable data points.

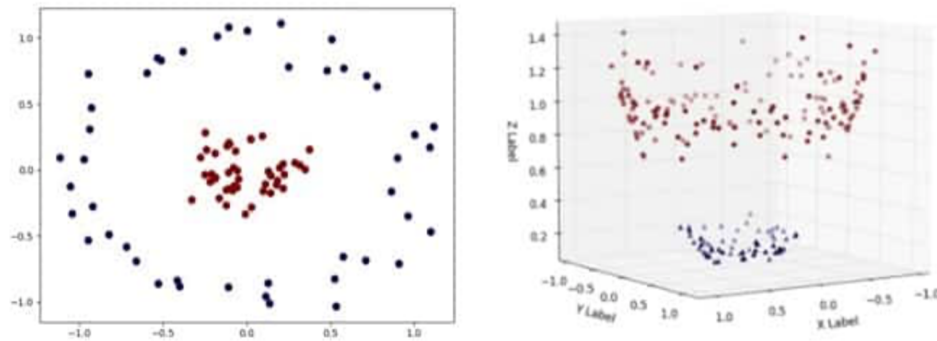


Figure 6: Kernel Trick

As seen from figure 6, that initially the data points are not linearly separable space but after applying kernel trick to it, the data is mapped into higher dimension and become linear separable. Popularly kernel includes three types:

Linear Kernel- It is mostly used kernel specially when data is linearly separable. Data points are separated using a single line or hyperplane. It is widely used in comparison to the other kernels and defined as:

$$K(x_i, x_j) = x_i \cdot x_j$$

Polynomial Kernel- Polynomial kernel is defined as:

$$K(x_i, x_j) = (x_i \cdot x_j + c)^p$$

Where, c is an arbitrary constant and p denotes the degree of polynomial. It is clear that linear kernel is a specialized case of polynomial kernel (c = 0 and p = 1).

Radial basis Kernel- Radial basis kernel is also termed as Gaussian kernel and defined as:

$$K(x_i, x_j) = e^{-\gamma \|x_i - x_j\|^2}$$

Where, γ is hyper parameter controls the variance of model, if γ is large model shows high variance and if γ is small model behaves like linear. $\|x_i - x_j\|$ represent the Euclidean distance between the data point x_i and x_j .

Example-

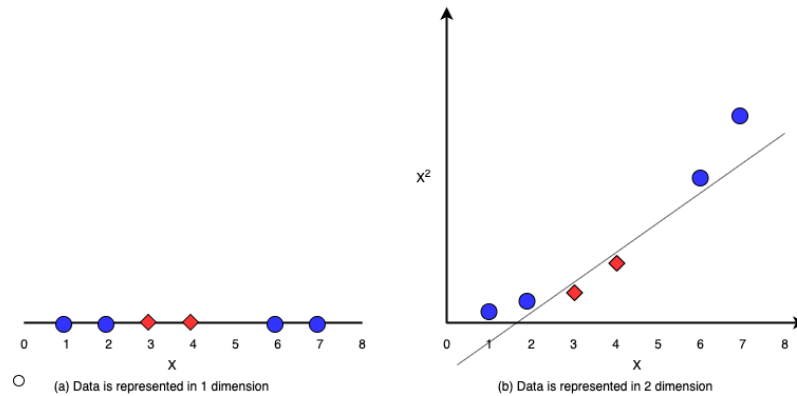
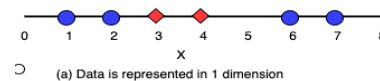


Figure 7 Kernel Trick

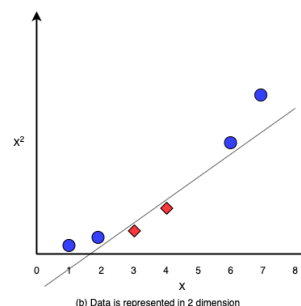
As it can be seen from the figure 7 that initial data is not linearly separable. The given distribution is expressed as:

X	Class
1	Blue
2	Blue
3	Red
4	Red
6	Blue
7	Blue



The given data is in 1 dimension but can be transformed into higher dimension using a kernel trick let a new dimension be formed using squaring the variable x . We will get a two-dimensional distribution of given data which is linearly separable and can be classified using a line.

X	X^2	Class
1	1	Blue
2	4	Blue
3	9	Red
4	16	Red
6	36	Blue
7	49	Blue



Advantages of SVM

- SVM model works well with high dimensional data.
- SVM model equally works well with linear and nonlinear separable data.
- Training model is relatively simple and easy.

Disadvantages of SVM

- Selection of the right kernel and parameters can be computationally expensive.
