

# Support Vector Machine

1. Introduction of SVM
2. What is Hyperplane, Support Vector and Margin
3. How to Select Optimal Hyperplane
4. Soft Margin and Hard Margin SVM
5. Regulation in SVM
6. What is Kernel
7. Types of Kernel
8. Tuning Parameter: Regulation, gamma

Support vector machines (**SVMs**) are **supervised learning models with associated learning algorithms** that analyze data used for classification and regression analysis. It can easily handle multiple continuous and categorical variables. SVM constructs a hyperplane in multidimensional space to separate different classes. SVM generates optimal hyperplane in an iterative manner, which is used to minimize an error. The core idea of SVM is to find a maximum marginal hyperplane (MMH) that best divides the dataset into classes.

SVM approach has some advantages compared to others classifiers. They are robust, accurate and very effective even in cases where the number of training samples is small. SVM technique also shows greater ability to generalize and greater likelihood of generating good classifiers.

Wikipedia tells us that SVMs can be used to do two things: classification or regression.

1. SVM is used for classification
2. SVR (Support Vector Regression) is used for regression

Goal of SVM:

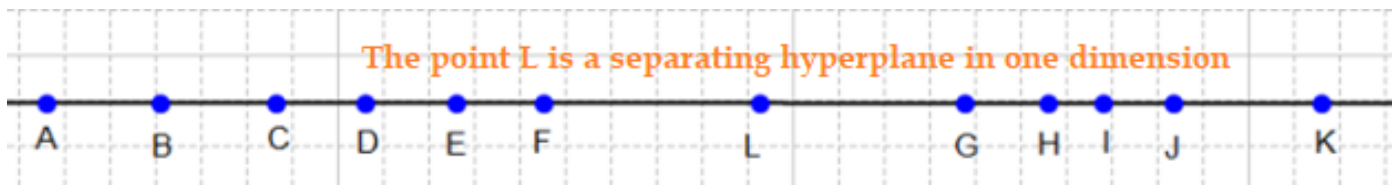
*The goal of a support vector machine is to find the optimal separating hyperplane which maximizes the margin of the training data.*

# Hyperplane

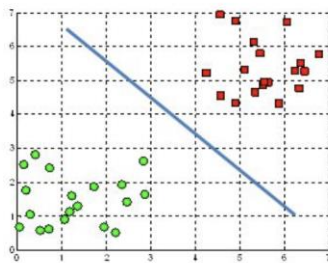
A hyperplane is a decision plane which separates between a set of objects having different class memberships.

A hyperplane is a generalization of a plane.

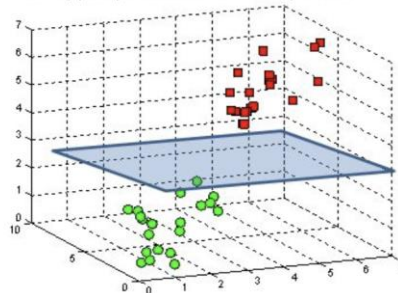
- in one dimension, a hyperplane is called a point
- in two dimensions, it is a line
- in three dimensions, it is a plane
- in more dimensions you can call it a hyperplane



A hyperplane in  $\mathbb{R}^2$  is a line



A hyperplane in  $\mathbb{R}^3$  is a plane

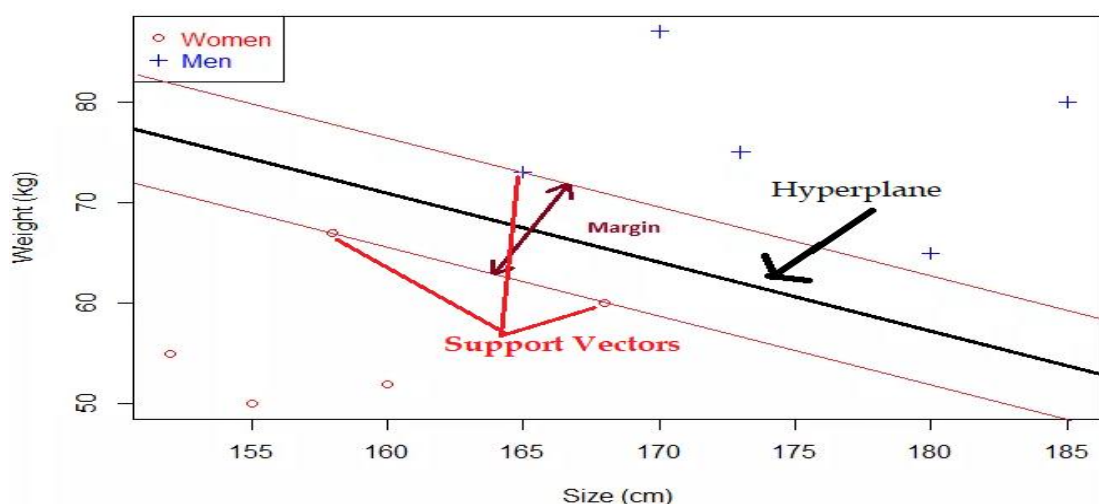


## Support Vectors

Support vectors are the data points, which are closest to the hyperplane. These points will define the separating line better by calculating margins.

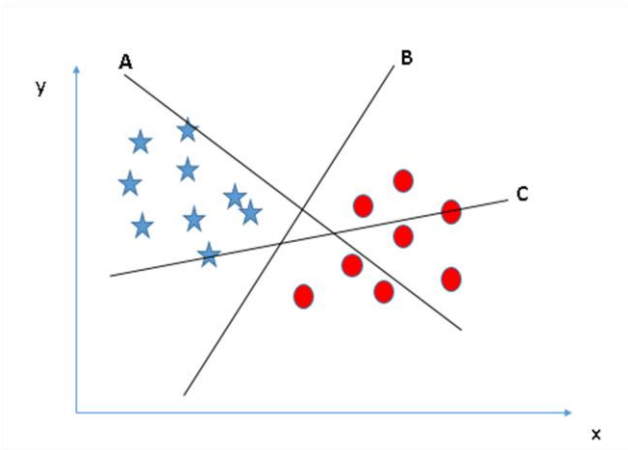
## Margin

A margin is a gap between the two lines on the closest class points. This is calculated as the perpendicular distance from the line to support vectors or closest points. If the margin is larger in between the classes, then it is considered a good margin, a smaller margin is a bad margin.

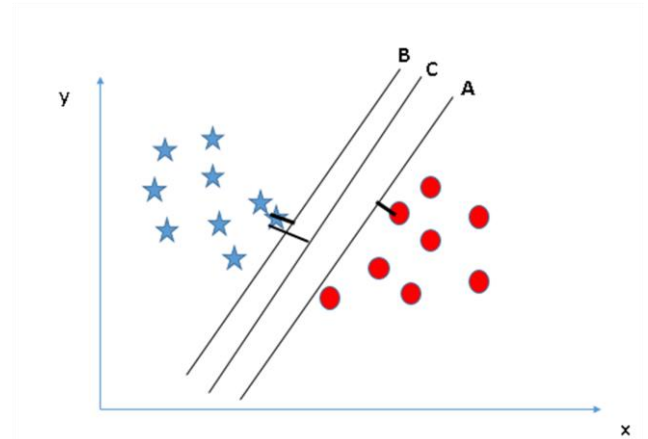


## How can we identify the right hyper-plane?

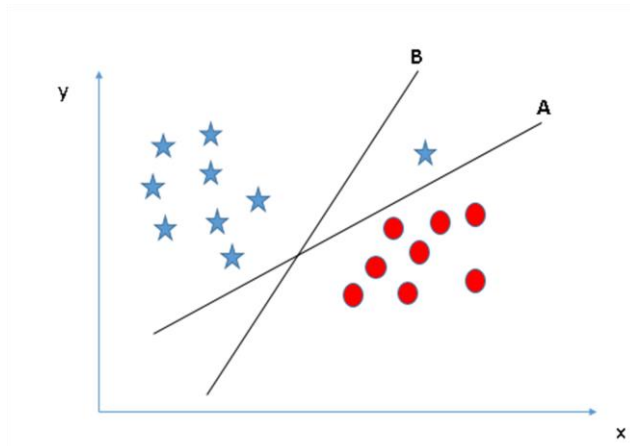
1. The two data points should be separated by maximum distance.
2. These two data points closer to the margin
3. Generate hyperplanes which segregates the classes in the best way



Select the hyper-plane which segregates the two classes better". In this scenario, hyper-plane "B" has excellently performed this job:- Rule 3



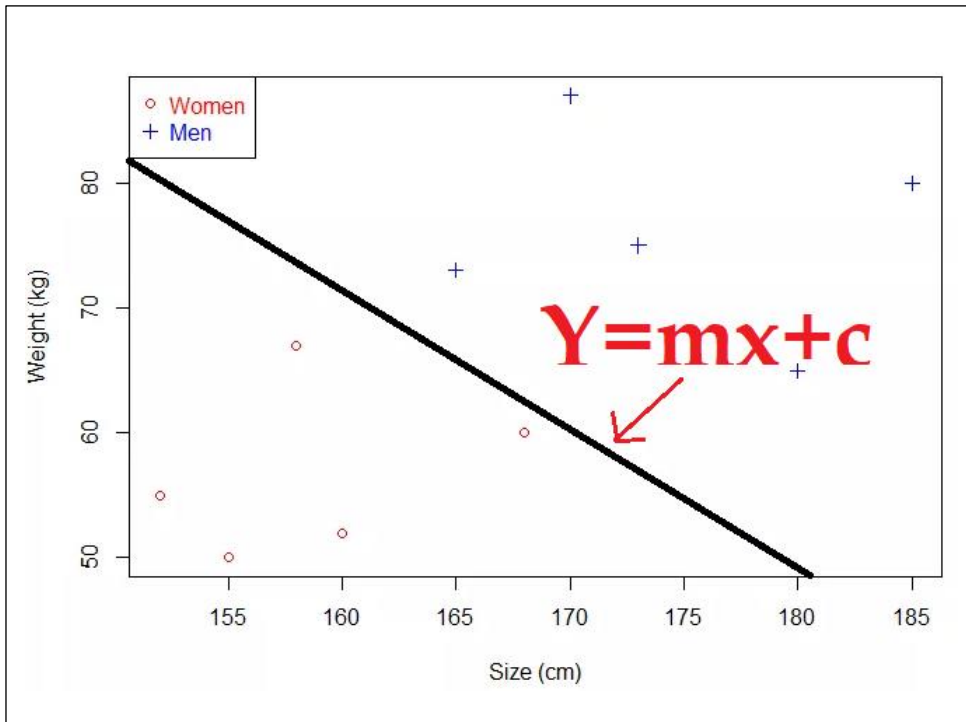
We can see that the margin for hyper-plane C is high as compared to both A and B. Hence, we name the right hyper-plane as C:- Rule 1,2



SVM selects the hyper-plane which classifies the classes accurately prior to maximizing margin. Here, hyper-plane B has a classification error and A has classified all correctly. Therefore, the right hyper-plane is A.

# How to Select Optimal Hyperplane?

## Understanding the equation of the hyperplane



Equation of line:

$$y = mx + b$$

$$y - mx - b = 0 \dots\dots\dots(\text{Eq 1})$$

Given two vector W and X in 3-dimension

$$W = \begin{bmatrix} -b \\ -m \\ 1 \end{bmatrix} \quad \text{and} \quad X = \begin{bmatrix} 1 \\ x \\ y \end{bmatrix}$$

$$\begin{aligned} W^T X &= [-b - m \ 1] \begin{bmatrix} 1 \\ x \\ y \end{bmatrix} \\ &= (-b \times 1) + (-m \times x) + (1 \times y) \\ &= y - mx - b \dots\dots\dots(\text{Eq 2}) \end{aligned}$$

Substitute Eq 2 in Eq 1 then equation of Hyperplane is,

$$W^T X = 0$$

## Why do we use the hyperplane equation $W^T X$ Instead of $y = mx + b$ ?

1. It is easier to work in more than two dimensions with this notation
2. The vector  $w$  will always be normal to the hyperplane

## Convert 3-dimension to 2-dimension

Given,  $W = [-m \ 1]$  and  $X = [x \ y]$

$$W \cdot X = (-m \times x) + (1 \times y)$$

$$= -mx + y$$

$$= W^T X + b$$

$$W \cdot X - b = W^T X$$

Equation of Hyperplane in 2-Dimension:

$$W \cdot X - b = 0$$

## How can we find the biggest margin?

1. Select two hyperplanes which separate the data with no points between them
2. Maximize their distance (the margin)

We can select two others hyperplanes  $H_1$  and  $H_2$  which also separate the data and have the following equations:

$$W \cdot X - b = \delta$$

$$W \cdot X - b = -\delta$$

If  $\delta = 1$  then,

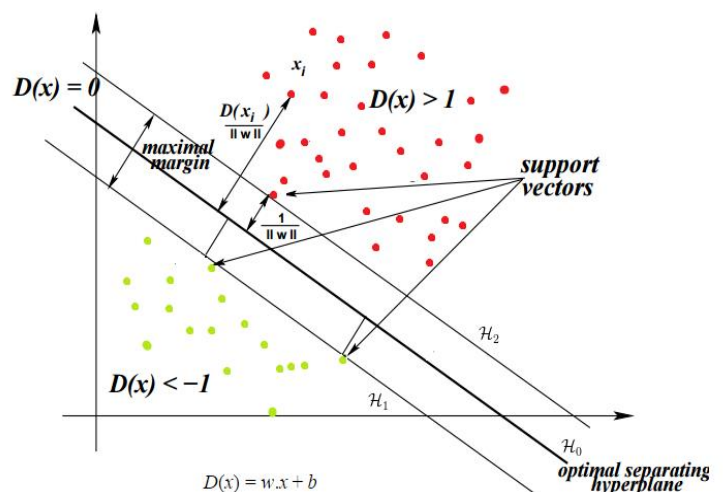
$$W \cdot X - b = 1$$

$$W \cdot X - b = -1$$

For each vector  $x_i$  either:

$$W \cdot X_i - b \geq 1 \quad \text{for } X_i \text{ having class 1}$$

$$W \cdot X_i - b \leq -1 \quad \text{for } X_i \text{ having class } -1$$



## Combining the constraints

Both the constraints stated above can be combined into a single constraint.

### Constraint 1:

For  $x_i$  having the class -1,  $(w \cdot x - b) \leq -1$

Multiplying both sides by  $y_i$  (which is always -1 for this equation)

$y_i(w \cdot x - b) \geq y_i(-1)$  which implies  $y_i(w \cdot x - b) \geq 1$  for  $x_i$  having the class -1.

### Constraint 2:

$y_i(w \cdot x - b) \geq 1$  for  $x_i$  having the class 1

Combining both the above equations, we get

$$y_i(w \cdot x - b) \geq 1 \text{ for } 1 \leq i \leq n$$

This leads to a unique constraint instead of two which are mathematically equivalent. The combined new constraint also has the same effect, i.e., no points between the two hyperplanes.

## Hard margin SVM

*If  $y_i(w \cdot x - b) \geq 1$*

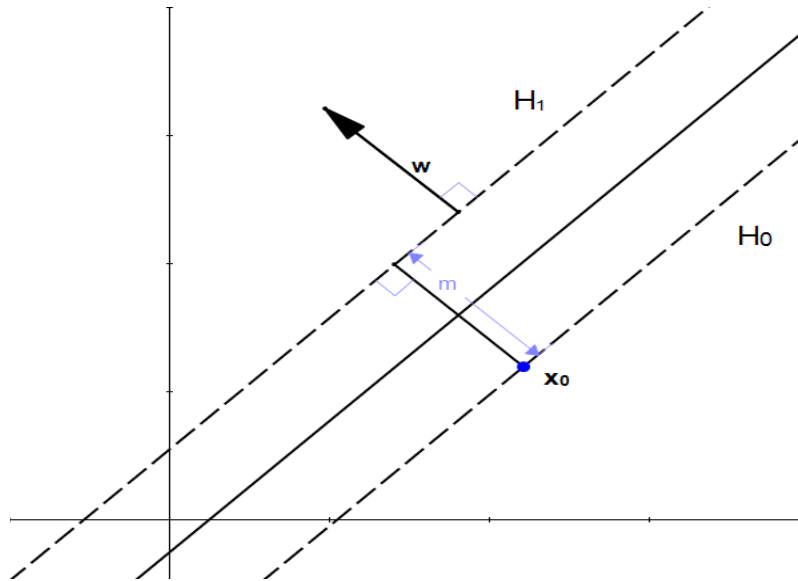
*then  $x_i$  is correctly classified*

*else:*

*$x_i$  is incorrectly classified.*

we can see that if the points are linearly separable then only our hyperplane is able to distinguish between them and if any outlier is introduced then it is not able to separate them. So, this type of SVM is called as **hard margin SVM**.

## Step 2: Maximize the distance between the two hyperplanes



Let:

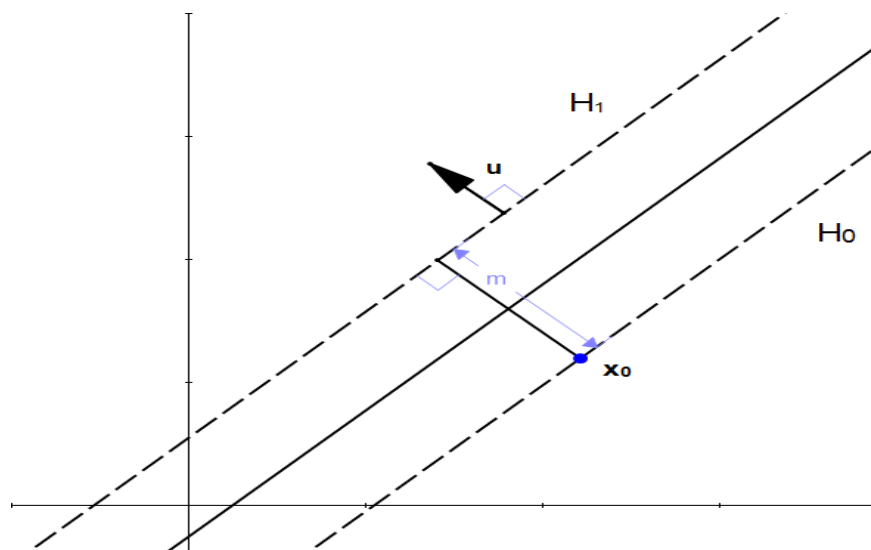
- $H_0$  be the hyperplane having the equation  $(w \cdot x - b) = -1$
- $H_1$  be the hyperplane having the equation  $(w \cdot x - b) = 1$
- $x_0$  be a point in the hyperplane  $H_0$ .

We will call  $m$  the perpendicular distance from  $x_0$  to the hyperplane  $H_1$ . By definition,  $m$  is what we are used to call **the margin**.

As  $x_0$  is in  $H_0$ ,  $m$  is the distance between hyperplanes  $H_0$  and  $H_1$ .

we already know a vector perpendicular to  $H_1$ , that is  $w$  (because  $H_1 = (w \cdot x - b) = 1$ )

Let's define  $u = w / \|w\|$  the unit vector of  $w$ . As it is a unit vector  $\|u\|=1$  and it has the same direction as  $w$  so it is also perpendicular to the hyperplane.

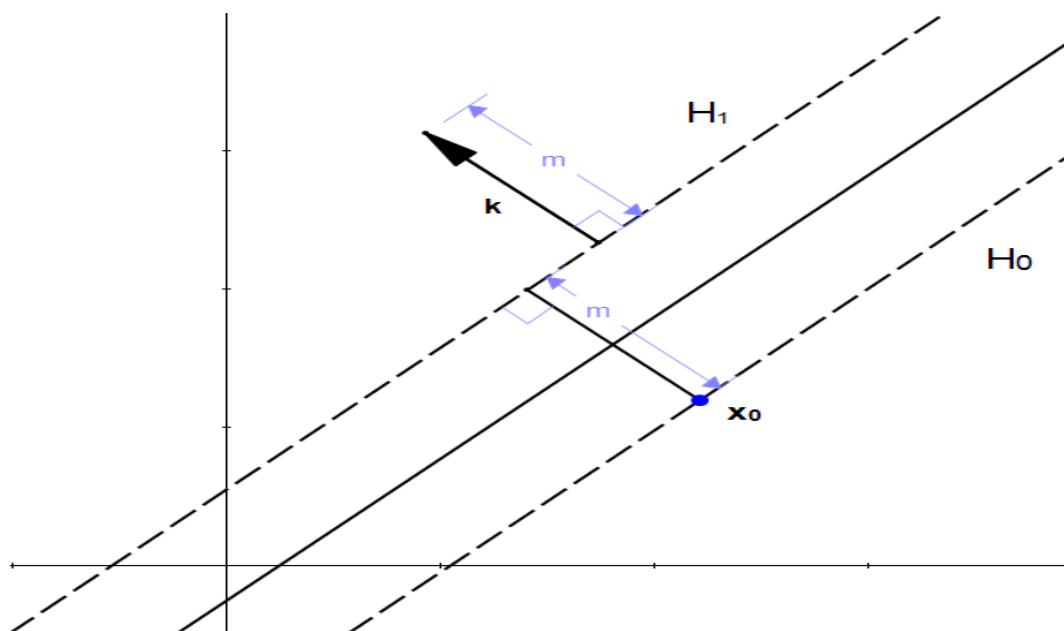


If we multiply  $u$  by  $m$  we get the vector  $k=mu$  and :

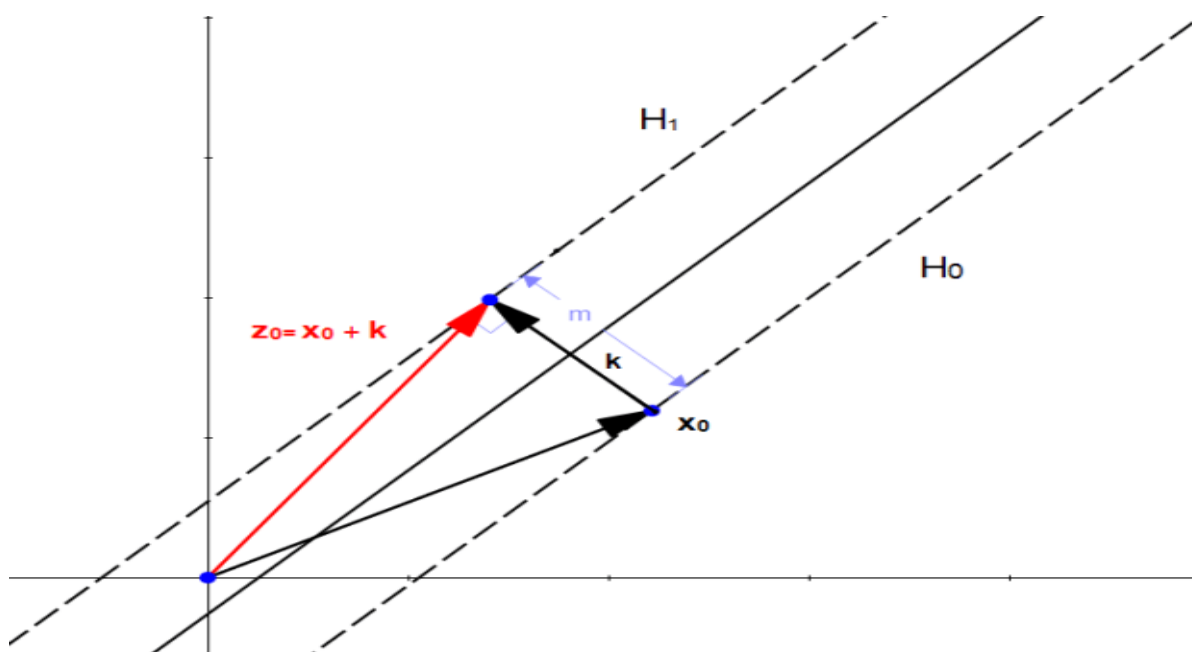
1.  $\|k\|=m$
2.  $k$  is perpendicular to  $H_1$ (because it has the same direction as  $u$ )

From these properties we can see that  $k$  is a vector of length  $m$  perpendicular *to*  $H_1$

$$K = mu = m \frac{W}{\|W\|}$$



If we start from the point  $X_0$  and add  $k$  we find that the point  $z_0 = X_0 + K$  is in the hyperplane  $H_1$





The fact that  $z_0$  is in  $H_1$  means that

$$W \cdot Z_0 - b = 1$$

$$W \cdot (x_0 + K) - b = 1$$

$$W \cdot \left( x_0 + m \frac{W}{\|W\|} \right) - b = 1$$

$$W \cdot x_0 + m \frac{\|W\|^2}{\|W\|} - b = 1$$

$$W \cdot x_0 - b = 1 - m\|W\|$$

As  $x_0$  is in  $H_0$  then  $w \cdot x_0 - b = -1$

$$-1 = -1 - m\|W\|$$

$$m = \frac{2}{\|W\|}$$

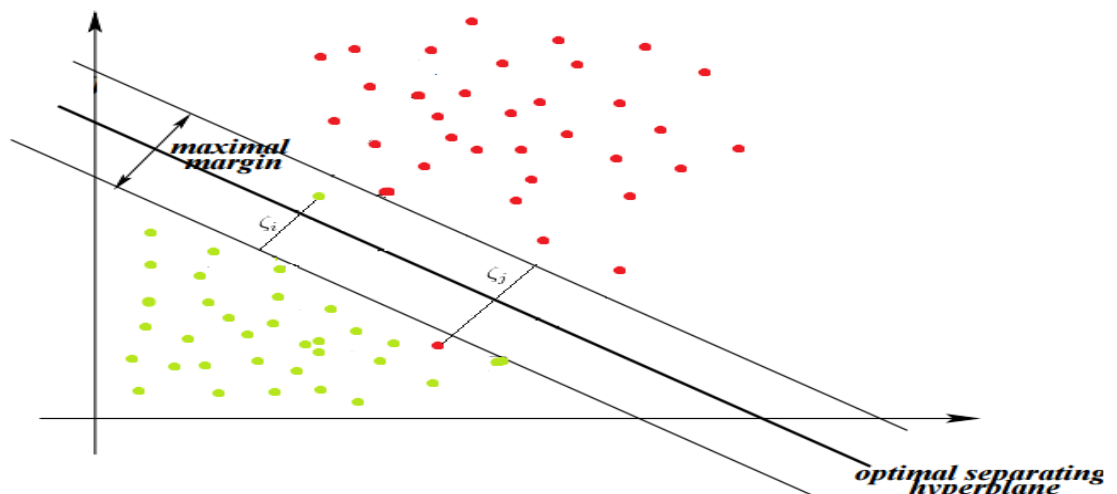
Our goal is to maximize the margin. Among all possible hyperplanes meeting the constraints, we will choose the hyperplane with the smallest  $\|w\|$  because it is the one which will have the biggest margin.

The only variable in this formula is  $w$ , which is indirectly proportional to  $m$ , hence to maximize the margin we will have to minimize  $\|w\|$ . This leads to the following optimization problem:

$$\text{Minimize in } (w, b) \frac{\|w\|^2}{2} \text{ subject to } y_i(w \cdot x_i - b) \geq 1 \text{ for any } i = 1, \dots, n$$

The above is the case when our data is linearly separable. There are many cases where the data cannot be perfectly classified through linear separation. In such cases, Support Vector Machine looks for the hyperplane that maximizes the margin and minimizes the misclassifications.

For this, we introduce the slack variable,  $\epsilon_i$  which allows some objects to fall off the margin but it penalizes them.



In this scenario, the algorithm tries to maintain the slack variable to zero while maximizing the margin. However, it minimizes the sum of distances of the misclassification from the margin hyperplanes and not the number of misclassifications.

Constraints now changes to  $\mathbf{y}_i (\mathbf{w} \cdot \mathbf{x} - \mathbf{b}) \geq 1 - \epsilon_i$  for  $1 \leq i < n, \epsilon_i \geq 0$  and the optimization problem changes to

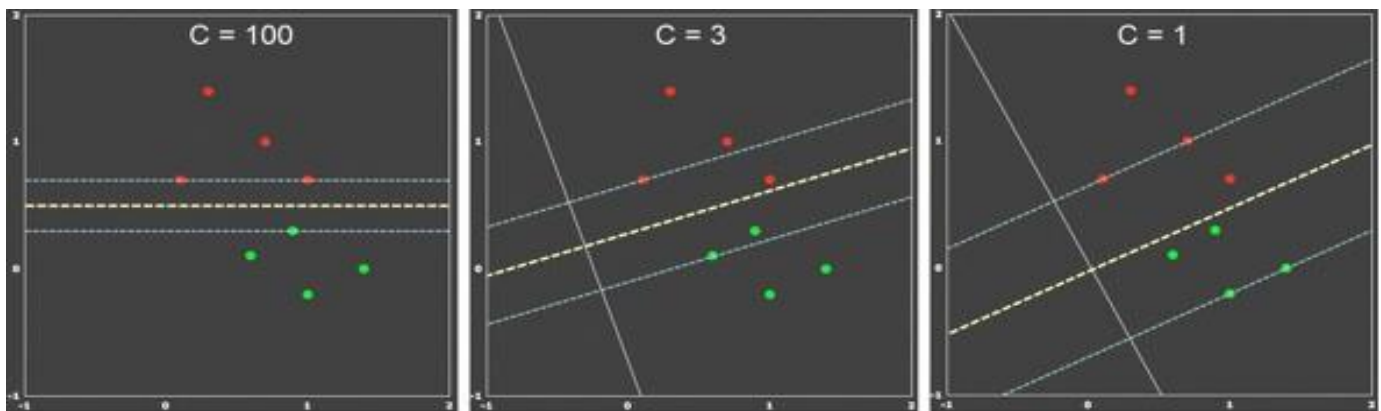
$$\text{Minimize in } (w, b) \frac{\|w\|^2}{2} + C \sum_i \epsilon_i \text{ subject to } y_i (w \cdot x - b) \geq 1 - \epsilon_i \text{ for any } i = 1 \dots n$$

Here, the parameter  $C$  is the **regularization** parameter that controls the trade-off between the slack variable penalty (misclassifications) and width of the margin.

1. Small  $C$  makes the constraints easy to ignore which leads to a large margin.
2. Large  $C$  allows the constraints hard to be ignored which leads to a small margin.
3. For  $C=\infty$ , all the constraints are enforced.

Large Value of parameter  $C \Rightarrow$  small margin

Small Value of parameter  $C \Rightarrow$  Large margin



## Soft margin SVM:

We basically consider that the data is linearly separable and this might not be the case in real life scenario. We need an update so that our function may skip few outliers and be able to classify almost linearly separable points. For this reason, we introduce a new *Slack variable* ( $\epsilon$ ) which is called  $X_i$ .

if we introduce  $\epsilon$  it into our previous equation we can rewrite it as

$$\mathbf{y}_i (\mathbf{w} \cdot \mathbf{x} - \mathbf{b}) \geq 1 - \epsilon_i$$

if  $\epsilon_i = 0$ ,

the points can be considered as correctly classified.

else:

$\epsilon_i > 0$ , Incorrectly classified points.

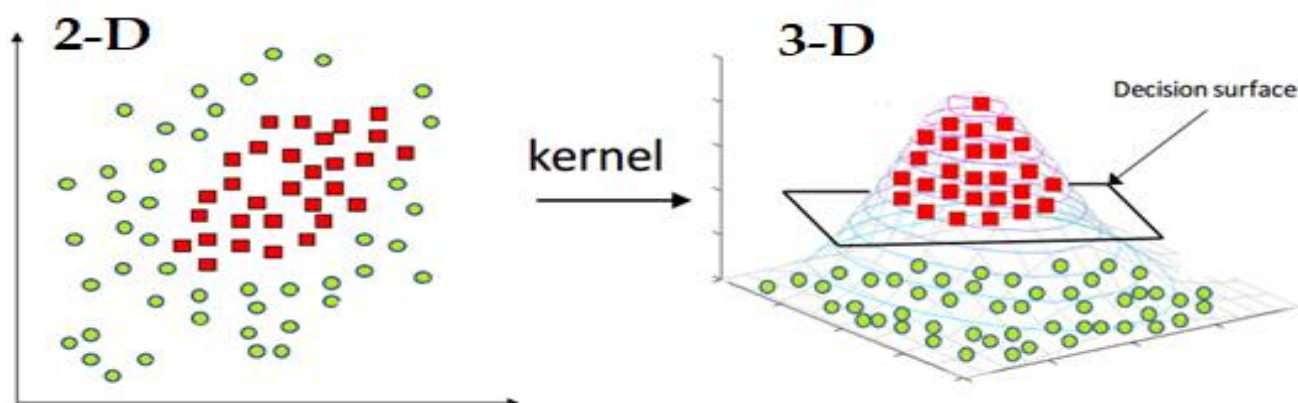
so, if  $\epsilon_i > 0$  it means that  $X_i$  (variables) lies in incorrect dimension, thus we can think of  $\epsilon_i$  as an error term associated with  $X_i$  (variable). The average error can be given as;

$$\frac{1}{n} \sum_i^n \epsilon_i$$

## Kernel:

The easiest way to separate two classes of data is a line in case of 2D data and a plane in case of 3D data. But it is not always possible to use lines or planes and one requires a nonlinear region to separate these classes. Support Vector Machines handle such situations by using a kernel function which maps the data to a different space where a linear hyperplane can be used to separate classes. This is known as the **kernel trick** where the kernel function transforms the data into the higher dimensional feature space so that a linear separation is possible.

The kernel takes a low-dimensional input space and transforms it into a higher dimensional space. In other words, we can say that it converts non-separable problems to separable problems by adding more dimension to it. It is most useful in non-linear separation problems. Kernel trick helps you to build a more accurate classifier.



1. **Linear Kernel:** A linear kernel can be used as normal dot product any two given observations. The product between two vectors is the sum of the multiplication of each pair of input values.

$$K(X, X_i) = \text{Sum}(X * X_i)$$

The kernel defines the similarity or a distance measure between new data and the support vectors. The dot product is the similarity measure used for linear SVM or a linear kernel because the distance is a linear combination of the inputs.

2. **Polynomial Kernel:** A polynomial kernel is a more generalized form of the linear kernel. The polynomial kernel can distinguish curved or nonlinear input space.

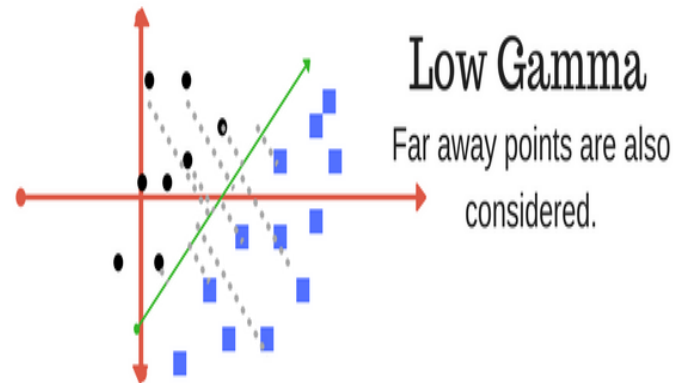
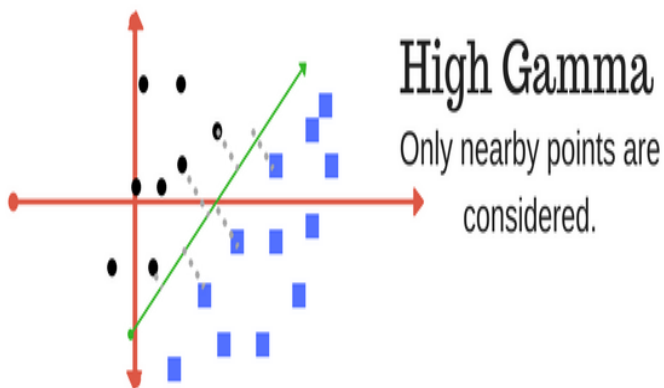
$$K(X, X_i) = 1 + \text{Sum}(X * X_i)^d$$

Where  $d$  is the degree of the polynomial.  $d=1$  is similar to the linear transformation. The degree needs to be manually specified in the learning algorithm.

3. **Gaussian Radial Basis Function Kernel (RBF):** The Radial basis function kernel is a popular kernel function commonly used in support vector machine classification. RBF can map an input space in infinite dimensional space.

$$K(X, X_i) = \exp(-\gamma \cdot \text{Sum}(X - X_i^2))$$

Here gamma ( $\gamma$ ) is a parameter, which ranges from 0 to 1. A higher value of gamma will perfectly fit the training dataset, which causes over-fitting. Gamma=0.1 is considered to be a good default value. The value of gamma needs to be manually specified in the learning algorithm.



4. **Sigmoid Kernel:** Sigmoid kernel is also known as hyperbolic tangent (sigmoid) kernel and multilayer perceptron (MLP) kernel. The origin of this kernel is from theory of neural networks and has been found to be performed well in practice as well. There are two adjustable parameters of sigmoid kernel, *slope of alpha* and the *intercept parameter C*. Scoring function of sigmoid kernel is as follow:

$$K(X, X_i) = \tanh(\alpha x x^T + c)$$

### Tuning Parameters:

1. **C:** Inverse of the strength of regularization.

**Behaviour:** As the value of 'c' increases the model gets overfits.  
As the value of 'c' decreases the model underfits.

2.  **$\gamma$ :** Gamma (used only for RBF kernel)

**Behaviour:** As the value of ' $\gamma$ ' increases the model gets overfits.  
As the value of ' $\gamma$ ' decreases the model underfits.

## References:

1. <https://towardsdatascience.com/support-vector-machines-svm-c9ef22815589>
2. <https://www.svm-tutorial.com/2014/11/svm-understanding-math-part-1/>
3. <https://medium.com/@pushkarmandot/what-is-the-significance-of-c-value-in-support-vector-machine-28224e852c5a>
4. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4260380/>
5. <https://www.intechopen.com/books/advances-in-character-recognition/svm-classifiers-concepts-and-applications-to-character-recognition>
6. <https://medium.com/machine-learning-101/chapter-2-svm-support-vector-machine-theory-f0812effc72>