

## Assignment

**B ANIRUDH SRINIVASAN**  
COE17B019

## 1. Support Vector Machine (SVM)

### 1.1 Introduction

SVM is a method for the classification of both linear and nonlinear data. In a nutshell, an SVM is an algorithm that works as follows. It uses a nonlinear mapping to transform the original training data into a higher dimension. Within this new dimension, it searches for the linear optimal separating hyperplane (i.e., a “decision boundary” separating the tuples of one class from another). With an appropriate nonlinear mapping to a sufficiently high dimension, data from two classes can always be separated by a hyperplane. The SVM finds this hyperplane using support vectors (“essential” training tuples) and margins (defined by the support vectors). Although the training time of even the fastest SVMs can be extremely slow, they are highly accurate, owing to their ability to model complex nonlinear decision boundaries. They are much less prone to overfitting than other methods. The support vectors found also provide a compact description of the learned model. SVMs can be used for numeric prediction as well as classification. They have been applied to a number of areas, including handwritten digit recognition, object recognition, and speaker identification, as well as benchmark time-series prediction tests.

#### Intuition and Math behind SVM

For linearly separable data, a hyper plane divides the hyperspace in such a way that, points to one side of hyperplane is class 1 and the points on the other side of the hyperplane is class 2.

In mathematical sense, if

$$w^t x + b = 0 \quad (1)$$

is the equation of the hyperplane, then,  $w^t x + b > 0$ , all  $x$  satisfying this equation, will be of one class and the points satisfying the eq

$w^t x + b < 0$

will be another class.

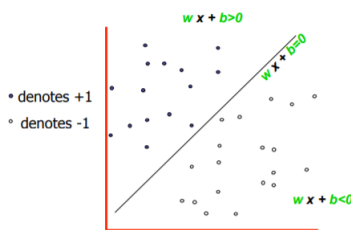


Figure 1: Illustration of how the hyperplane separates points.

We can easily separate the datapoints with many hyperplanes, For example,

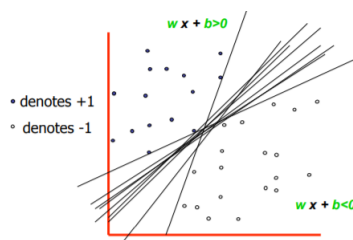


Figure 2: This problem has infinitely many solutions.

But only one of those hyperplane classify the data optimally, and we use svm to identify them.

The optimal plane is the one in which the distance between the plane and the class points are maximum. Moreover, the minimum distance between points in class A and the hyperplane and points in class B and the hyperplane must be same, i.e. if the minimum distance of points class A with the hyperplane is greater than that of minimum distance of points in classe B with the hyper plane, then the hyperplane wouldnt be optimum as we tend to classify points more towards class A than class B. For example, consider this one dimensional data set,

Table 1: Example dataset

class	x
class A	1
class A	2
class A	3
class B	5
class B	6
class B	7

The optimum plane (or point) in this case would be  $x = 4$ , which is equally spaced between both the classes.

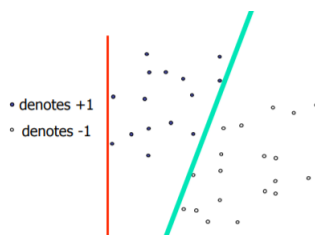


Figure 3: Optimal hyper plane

Support vectors are those planes which run parallel to the decision boundary and pass through the points of both the classes, which have minimum distance with the hyperplane.

The equation of the support vectors is,

$$w^t x + b = \pm 1 \quad (2)$$

The optimal decision boundary is the one in which the margin (distance between the support vectors) is maximum.

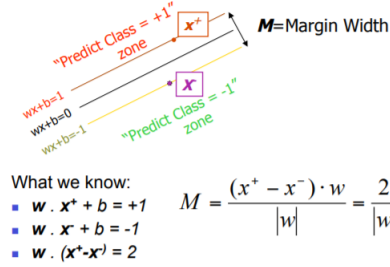


Figure 4: Optimal hyper conditions

From above figure, we can tell that,

$$w^t x + b \geq 1 \text{ if } x \in \text{class A} \quad (3)$$

$$w^t x + b \leq -1 \text{ if } x \in \text{class B} \quad (4)$$

We can create a label such that, points in class B have label -1 and point in class A has label 1.

Therefore, the above equation becomes,

$$w^t x_i + b \geq 1 \text{ if } y_i \in \text{class A} \quad (5)$$

$$w^t x_i + b \leq -1 \text{ if } y_i \in \text{class B} \quad (6)$$

We can combine them into one equation,

$$y_i(w^t x_i + b) \geq 1 \quad (7)$$

Moreover, we know that margin width,

$$M = \frac{2}{|w|} \quad (8)$$

Maximising M, is same as minimising  $|w|$ .

Thus, we minimise,

$$|w| \quad (9)$$

$$\text{and } \mathcal{N}(y_i, x_i) : y_i(w^t x_i + b \geq 1) \quad (10)$$

This minimisation problem can be solved by associating Lagrange multiplier with each constraint of this problem.

### Cases of SVM

There are three cases of SVM. They are:

- **Hard Margin**

Here, we try to find a decision boundary that classifies all training points correctly.

Extremely useful in the case of linearly separable, but performs poorly in case of non linearly separable and is susceptible to noises.

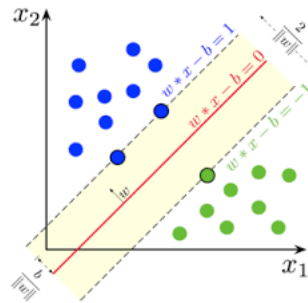


Figure 5: Example of Hard Margin.

- **Soft Margin**

Here, we try to find a decision boundary that classifies most of the points correctly.

We allow some points to be misclassified. The amount of error allowed is set by us.

This is useful to classify noisy data and some non-linearly separable data.

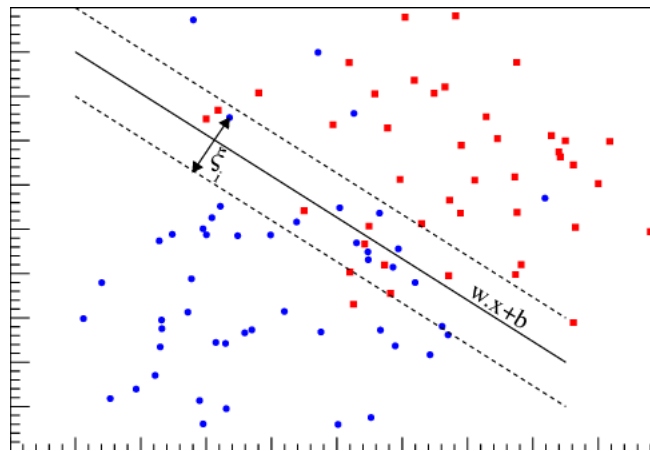


Figure 6: Example of Soft Margin.

- **Kernal**

Here, we transform the points to a higher dimension (where they are linearly separable) and find a linear hyperplane that classifies the points accurately in that dimension and we plot the contours of that hyperplane on our lower dimension to get a decision boundary.

This is extremely useful to classify non-linearly seperable data.

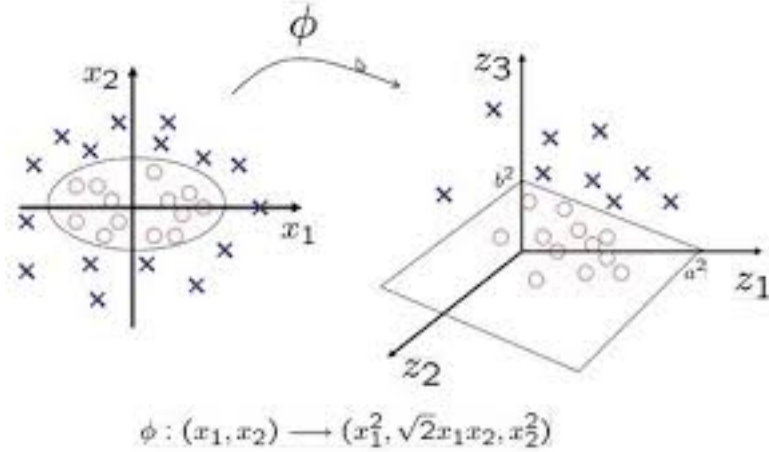


Figure 7: Example of the use of kernel function to transform the input to higher dimension, where they are linearly separable.

## 1.2 Classifying using SVM with example

Example Dataset

Table 2: Example dataset

class	x	y
class 1 (w1)	2	2
class 1 (w1)	-1	2
class 1 (w1)	1	3
class 1 (w1)	-1	-1
class 1 (w1)	0.5	0.5
class 2 (w2)	-1	-3
class 2 (w2)	0	-1
class 2 (w2)	1	-2
class 2 (w2)	-1	-2
class 2 (w2)	0	-2

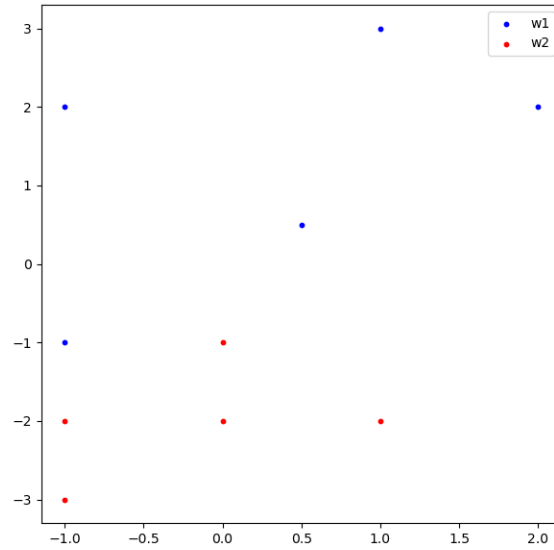


Figure 8: Visualisation of the dataset

### Case 1 : Hard Margin

**Step 1:** Decide the initial weights of the plane.

The Equation of the plane that separates the classes is :

$$w^t x + b = 0 \quad (11)$$

where,

- $w$  is the normal vector of the hyperplane that separates the data. It is represented as a matrix. Hence  $w^t$  means transpose of  $w$  matrix.
- $x$  is a variable vector.
- $b$  is the bias.

The  $z$  matrix would be a matrix where,  $z_i = 1$  if  $x_i$  in class 1  
 $z_i = -1$  if  $x_i$  in class 2

**Step 2:** Finding the Lagrange Multipliers

We need to find the lagrange multipliers for our hyperplane, which inturn would help us to identify the equation of the hyperplane.

$$L(w, b, \alpha) = \frac{\|w\|^2}{2} - \sum_{k=1}^n \alpha_k (1 - z_k (w^t x_k + b)) \quad (12)$$

where,

- $w^t$  is the transpose of weight matrix, and  $\|w\|$  is the magnitude of the normal vector of the hyper plane.
- $b$  is the bias.
- $\alpha$  is a matrix of all the lagrange multipliers. There are totally  $n$  lagrange multiplier, where  $n$  is the number of datapoints.
- $z_k$  is the class value of a data point  $x_k$ . Note that  $z_k$  takes either 1 or -1 and not 1 or 0.

After finding few relation between  $\alpha$  and weight matrix and bias, the above equation reduces to,

$$L(\alpha) = \sum_{k=1}^n \alpha_k - \frac{1}{2} \sum_{k,j}^n \alpha_k \alpha_j z_k z_j x_j^t x_k \quad (13)$$

We should, find the minimum value of  $L(\alpha)$ , which satisfy the following constraint,

- (a)  $\sum_{k=1}^n z_k \alpha_k = 0$
- (b)  $\alpha_k \geq 0, k = 1, 2, \dots, n$

For our example, we got,

**Step 3:** Finding Weight Vector and Bias

The weight vector is given by,

$$w = \sum_{k=1}^n \alpha_k z_k x_k \quad (14)$$

We know that,

$$z_i (w^t x_i + b) = 1$$

is the equation of the support planes.

And these support planes would pass through the corner points(that are to the other class end) of each class.



Table 3:  $\alpha$  values

$\alpha$	Value
$\alpha_1$	$1.26183502 \times 10^{-1}$
$\alpha_2$	$1.03087027 \times 10^{-14}$
$\alpha_3$	$2.58278642 \times 10^{-15}$
$\alpha_4$	$3.45312250 \times 10^0$
$\alpha_5$	$4.21384434 \times 10^1$
$\alpha_6$	$1.64516838 \times 10^{-15}$
$\alpha_7$	$3.01092476 \times 10^0$
$\alpha_8$	$1.82009356 \times 10^{-14}$
$\alpha_9$	$9.89765678 \times 10^1$
$\alpha_{10}$	$8.76371458 \times 10^{-15}$

Therefore, we need to find the points.

$\alpha_i$  would be negligible for all non - other points/ points through which the support vector wont pass through.

Therefore, with this we can locate the corner point and take avgerage to find avgerage b value.

In our example we get,

$w = [-2.0002976 \ 2.00039283]$  and  $b = 1.0001476102615467$ .

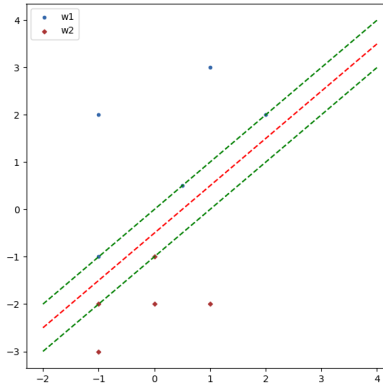


Figure 9: The decision boundary for the example dataset. The red line is the actual decision boundary with green line as the support vectors.

**Case 2 : Soft Margin**

The steps are same as that of Hard Margin, we only change the equation to allow error.

$$w^t x + b = 1 - \varepsilon \quad (15)$$

$$L(w, b, \varepsilon, \alpha, \gamma) = \frac{\|w\|^2}{2} + C * \sum_{k=1}^n \varepsilon_k - \sum_{k=1}^n \alpha_k (z_k (w^t x_k + b) - 1 + \varepsilon_k) - \sum_{i=k}^n \gamma_k \varepsilon_k \quad (16)$$

where,

- $\varepsilon$  is the error in our classifier.
- $\alpha$  and  $\gamma$  are two lagrange multiplier.
- $C$  is a regularisation parameter, which tells us how much error to be present. If  $C = 0$ , we get hard margin, and if  $C$  is a large value, then we allow a lot of miscalculation as the error term in eq(31) increases and dominates the equation.

This may seem difficult to solve, but they give the same equation, but differ in constraint.

$$L(\alpha) = \sum_{k=1}^n \alpha_k - \frac{1}{2} \sum_{k,j}^n \alpha_k \alpha_j z_k z_j x_j^t x_k \quad (17)$$

We should, find the minimum value of  $L(\alpha)$ , which satisfy the following constraint,

- (a)  $\sum_{k=1}^n z_k \alpha_k = 0$
- (b)  $C \geq \alpha_k \geq 0, k = 1, 2, \dots, n$

The weight vector is given by,

$$w = \sum_{k=1}^n \alpha_k z_k x_k \quad (18)$$

The error or slack is given by,

$$\varepsilon_i = \frac{\alpha_i}{C} \quad (19)$$

Similarly, corresponding  $\alpha$  of non corner point will be negligible. With this we find the corner points.

Same way we substitute the corner points in eq(32) to find bias .

**Case 3 : Kernal method**

Here we transform the input data to a higher dimension using kernel functions.

$$y = \Phi(x) \quad (20)$$

where,

- $y$  is the transformed data.
- $x$  is the input data points.
- $\phi(.)$  is a transformation or kernel function.

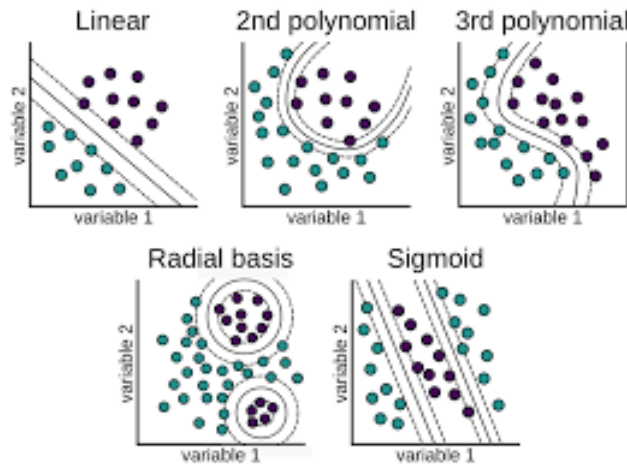


Figure 10: Different kernel functions

After transforming the input data, follow hard margin algorithm to get the equation of plane in the higher dimension.

To make the computation at higher dimension easier, we use a special property of kernel functions. i.e. the dot product of the function is conserved throughout the dimensions.

The contours of that plane on the input dimension would be our decision tree.

### 1.3 Pseudo-Code SVM

#### Hard Margin

```

1 Begin
2 Initialize all weights to 0;
3 for each training tuple X in D {
4    $z_i = \text{classlabel}(X)$ 
5 }
6  $n = \text{len}(X)$ 
7  $L(\alpha) = \sum_{k=1}^n \alpha_k - \frac{1}{2} \sum_{k,j} \alpha_k \alpha_j z_k z_j x_j^t x_k$ 
8  $\alpha = \alpha : L(\alpha)$  is maximum,
9    $\alpha_i \geq 0$ , and
10   $\sum_{i=1}^n \alpha_i z_i = 0$ 
11
12 Weights =  $\sum_{i=1}^n z_i \alpha_i x_i$ 
13
14 for i in  $\text{len}(\alpha)$  {
15   if ( $\alpha_i \geq 0.001$ )
16      $b = \frac{1}{z_i} - w^t x_i$ 
17     add b to the set bias
18 }
19 bias = bias.mean()

```

#### Soft Margin SVM

```

1 Begin
2 Initialize all weights to 0;
3 for each training tuple x in X {
4    $z_i = \text{classlabel}(x)$ 
5 }
6 initialise the desired value of c.
7  $n = \text{len}(X)$ 
8  $L(\alpha) = \sum_{k=1}^n \alpha_k - \frac{1}{2} \sum_{k,j} \alpha_k \alpha_j z_k z_j x_j^t x_k$ 
9  $\alpha = \alpha : L(\alpha)$  is maximum,
10   $c \geq \alpha_i \geq 0$ , and
11   $\sum_{i=1}^n \alpha_i z_i = 0$ 
12
13 Weights =  $\sum_{i=1}^n z_i \alpha_i x_i$ 
14 for i in  $\text{len}(X)$  {
15    $\varepsilon_i = \frac{1}{\alpha_i} \{c\}$ 
16 }
17 for i in  $\text{len}(X)$  {
18   if ( $\alpha_i \geq 0.001$ )
19      $b = \frac{1 - \varepsilon_i}{z_i} - w^t x_i$ 
20     add b to the set bias
21 }
22 bias = bias.mean()

```

## Kernal trick SVM

```

1 Begin
2 Initialize all weights to 0;
3 for each training tuple x in X {
4    $z_i = \text{classlabel}(x)$ 
5 }
6  $Y = \psi(X)$  //transforming the input into higher dimensions.
7
8 weight, bias = HardMargin(Y)
```

### 1.4 Additional Notes

#### More about Kernal functions

Consider the following example. A 3-D input vector  $X = (x_1, x_2, x_3)$  is mapped into a 6-D space,  $Z$ , using the mappings  $\phi_1(X) = x_1$ ,  $\phi_2(X) = x_2$ ,  $\phi_3(X) = x_3$ ,  $\phi_4(X) = (x_1)^2$ ,  $\phi_5(X) = x_1.x_2$ , and  $\phi_6(X) = x_1.x_3$ . A decision hyperplane in the new space is

$$d(Z) = W^t Z + b \quad (21)$$

where,  $W$  and  $Z$  are vectors. This is linear. We solve for  $W$  and  $b$  and then substitute back so that the linear decision hyperplane in the new ( $Z$ ) space corresponds to a nonlinear second-order polynomial in the original 3-D input space:

But there are some problems. First, how do we choose the nonlinear mapping to a higher dimensional space? Second, the computation involved will be costly. Refer to

Given the test tuple, we have to compute its dot product with every one of the support vectors. In training, we have to compute a similar dot product several times in order to find the MMH. This is especially expensive. Hence, the dot product computation required is very heavy and costly.

Luckily, we can use another math trick. It so happens that in solving the quadratic optimization problem of the linear SVM (i.e., when searching for a linear SVM in the new higher dimensional space), the training tuples appear only in the form of dot products,  $\phi(X_i) \cdot \phi(X_j)$ , where  $\phi(X)$  is simply the nonlinear mapping function applied to transform the training tuples. Instead of computing the dot product on the transformed data tuples, it turns out that it is mathematically equivalent to instead apply a kernel function,  $K(X_i, X_j)$ , to the original input data. That is,

$$K(X_i, X_j) = \phi(X_i) \cdot \phi(X_j) \quad (22)$$

In other words, everywhere that  $\phi(X_i) \cdot \phi(X_j)$  appears in the training algorithm, we can replace it with  $K(X_i, X_j)$ . In this way, all calculations are made in the original input space, which is of potentially much lower dimensionality! We can safely avoid the mapping—it turns out that we don't even have to know what the mapping is! We will talk more later about what kinds of functions can be used as kernel functions for this problem.

After applying this trick, we can then proceed to find a maximal separating hyperplane. kernel functions that could be used to replace the dot product scenario just described,

$$\textit{Polynomialkernelofdegreeh} : K(X_i, X_j) = (X_i \Delta X_j + 1)^h \quad (23)$$

$$\textit{Gaussianradialbasisfunctionkernel} : K(X_i, X_j) = e^{-\frac{\|kX_iX_jk\|^2}{2\sigma^2}} \quad (24)$$

$$\textit{Sigmoidkernel} : K(X_i, X_j) = \tanh(\kappa X_i X_j - \delta) \quad (25)$$