



Department of Mathematics, University of Udine

SVM and their use in IR

Lorenzo Zanolin.

lorenzo.zanolin@spes.uniud.it

May x, 2023



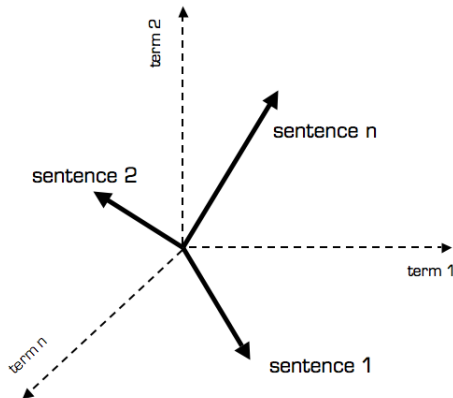
- 1 SVM over linearly separable data
- 2 SVM over non linearly separable data
- 3 SVM use in IR



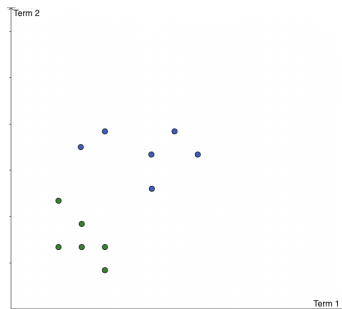
Vector Space

Suppose you have to classify whether a document is *relevant* or not. We can think to use *terms* as features to divide properly the data. We will use a *t-dimensional* vector space to represent our documents.

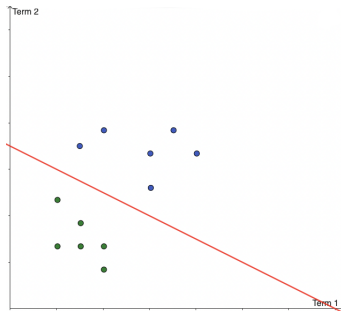
Suppose you have to classify whether a document is *relevant* or not. We can think to use *terms* as features to divide properly the data. We will use a *t-dimensional* vector space to represent our documents.



For simplicity, we can reason using a 2D Space.
Data can be separated using a *decision boundary*, which is an *hyperplane*.



For simplicity, we can reason using a 2D Space.
Data can be separated using a *decision boundary*, which is an *hyperplane*.





Support Vectors

In the previous example the *decision boundary* is a line, represented by the equation $a + bt_1 + ct_2 = 0$.

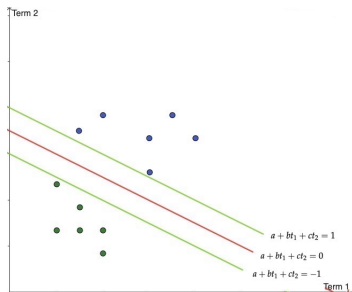
We can introduce two parallel hyperplanes (lines) to the decision boundary, called *support vectors* whose equations are

$$\begin{cases} a + bt_1 + ct_2 = 1 \\ a + bt_1 + ct_2 = -1 \end{cases} \quad (1)$$

In the previous example the *decision boundary* is a line, represented by the equation $a + bt_1 + ct_2 = 0$.

We can introduce two parallel hyperplanes (lines) to the decision boundary, called *support vectors* whose equations are

$$\begin{cases} a + bt_1 + ct_2 = 1 \\ a + bt_1 + ct_2 = -1 \end{cases} \quad (1)$$





For convenience, we can rewrite equations using *vectorization* notation.

$$\begin{cases} b^T t + a = 0 \\ b^T t + a = \pm 1 \end{cases} \quad (2)$$

For convenience, we can rewrite equations using *vectorization* notation.

$$\begin{cases} b^T t + a = 0 \\ b^T t + a = \pm 1 \end{cases} \quad (2)$$

Intuitively, we can define the *margin* of the two support vectors as the distance between them.



For convenience, we can rewrite equations using *vectorization* notation.

$$\begin{cases} b^T t + a = 0 \\ b^T t + a = \pm 1 \end{cases} \quad (2)$$

Intuitively, we can define the *margin* of the two support vectors as the distance between them.

We can consider two points x_1, x_2 that lie respectively on the two support vectors, their distance is $\lambda \|b\| = \frac{2}{\sqrt{b^T b}}$.

SVM are used to find the *maximum margin linear classifier*, thus we want to maximize the margin.

Remembering we want to classify documents, our goal is to find specific b s.t. given a document x belonging to class y the decision boundary behave the following:

$$\begin{cases} b^T x + a \geq 1 & \text{if } y = 1 \\ b^T x + a \leq -1 & \text{if } y = -1 \end{cases} \quad (3)$$

Remembering we want to classify documents, our goal is to find specific b s.t. given a document x belonging to class y the decision boundary behave the following:

$$\begin{cases} b^T x + a \geq 1 & \text{if } y = 1 \\ b^T x + a \leq -1 & \text{if } y = -1 \end{cases} \quad (3)$$

We can define the *cost function* as a system of equations.

$$\begin{cases} \min_{b,a} \frac{\sqrt{b^T b}}{2} \\ \text{subject to } y_i(b^T x_i + a) \geq 1 \quad \forall x_i \end{cases} \quad (4)$$



In real world problem it is not likely to get an exactly separate line dividing the data within the space. It would be better for the smooth boundary to ignore few data points than be curved or go in loops, around the outliers.



In real world problem it is not likely to get an exactly separate line dividing the data within the space. It would be better for the smooth boundary to ignore few data points than be curved or go in loops, around the outliers.

So, we will use *slack variables* to introduce a penalty for each misclassified point.



In real world problem it is not likely to get an exactly separate line dividing the data within the space. It would be better for the smooth boundary to ignore few data points than be curved or go in loops, around the outliers.

So, we will use *slack variables* to introduce a penalty for each misclassified point.

The new cost function will be

$$\begin{cases} \min_{b,a} \frac{\sqrt{b^T b}}{2} + C \sum_i \xi_i \\ \text{subject to} & y_i(b^T x_i + a) \geq 1 - \xi_i \quad \text{and} \quad \xi_i > 0 \quad \forall x_i \end{cases} \quad (5)$$

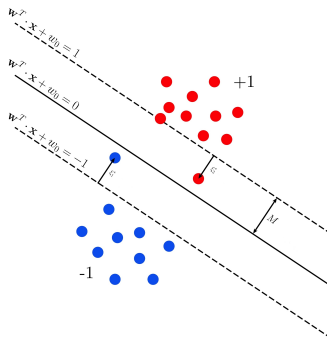


The larger is C the stricter the classification is, since a larger C will give more evidence to slack variables.

$$\begin{cases} \min_{b,a} \frac{\sqrt{b^T b}}{2} + C \sum_i \xi_i \\ \text{subject to} & y_i(b^T x_i + a) \geq 1 - \xi_i \quad \text{and } \xi_i > 0 \quad \forall x_i \end{cases} \quad (6)$$

The larger is C the stricter the classification is, since a larger C will give more evidence to slack variables.

$$\begin{cases} \min_{b,a} \frac{\sqrt{b^T b}}{2} + C \sum_i \xi_i \\ \text{subject to } y_i(b^T x_i + a) \geq 1 - \xi_i \quad \text{and } \xi_i > 0 \quad \forall x_i \end{cases} \quad (6)$$



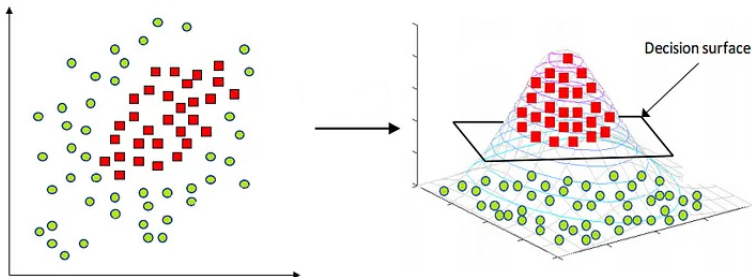


Non linearly separable data

What if our data is not linearly separable?

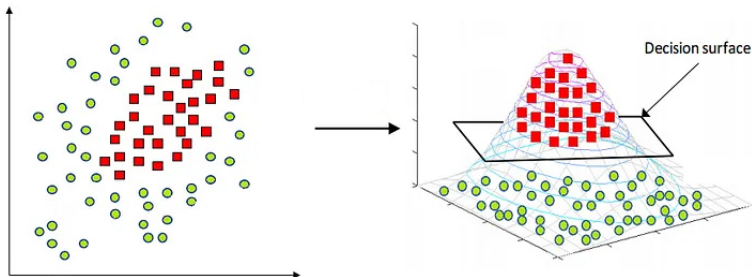
Non linearly separable data

What if our data is not linearly separable? It should be better to reason in a bigger dimensional space.



Non linearly separable data

What if our data is not linearly separable? It should be better to reason in a bigger dimensional space.



But augmenting dimensions costs a lot...



Consider the function $\phi : \mathbb{R}^3 \mapsto \mathbb{R}^{10}$ used to map points in a new vector space. Calculating the *similarity* $\phi(x_i)^T \phi(x_j)$ between each point may be intractable.



Kernel Trick

Consider the function $\phi : \mathbb{R}^3 \mapsto \mathbb{R}^{10}$ used to map points in a new vector space. Calculating the *similarity* $\phi(x_i)^T \phi(x_j)$ between each point may be intractable.

Here *Kernel Trick* comes handy.



Consider the function $\phi : \mathbb{R}^3 \mapsto \mathbb{R}^{10}$ used to map points in a new vector space. Calculating the *similarity* $\phi(x_i)^T \phi(x_j)$ between each point may be intractable.

Here *Kernel Trick* comes handy.

It consists of a simple linear algebra reformulation,

$$\phi(x_i)^T \phi(x_j) = K(x_i, x_j) = (1 + x_i^T x_j)^2. \quad (7)$$



Consider the function $\phi : \mathbb{R}^3 \mapsto \mathbb{R}^{10}$ used to map points in a new vector space. Calculating the *similarity* $\phi(x_i)^T \phi(x_j)$ between each point may be intractable.

Here *Kernel Trick* comes handy.

It consists of a simple linear algebra reformulation,

$$\phi(x_i)^T \phi(x_j) = K(x_i, x_j) = (1 + x_i^T x_j)^2. \quad (7)$$

Instead of doing the complex computations in the 10-dimensional space, we reach the same result within the 3-dimensional space by calculating the dot product.



Gaussian Kernel

There are lots of Kernel functions, an example is the *Gaussian*.



Gaussian Kernel

There are lots of Kernel functions, an example is the *Gaussian*.

The idea behind is to use all n training points as *landmarks*.

Then we can calculate the similarities of each point and all the landmarks.



There are lots of Kernel functions, an example is the *Gaussian*.

The idea behind is to use all n training points as *landmarks*.

Then we can calculate the similarities of each point and all the landmarks.

$$\forall l_i \text{ similarity}(x, l^{(i)}) = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right) \quad (8)$$

There are lots of Kernel functions, an example is the *Gaussian*.

The idea behind is to use all n training points as *landmarks*.

Then we can calculate the similarities of each point and all the landmarks.

$$\forall l_i \text{ similarity}(x, l^{(i)}) = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right) \quad (8)$$

The result will be a mapping for each point in a n -dimensional space. Finally, we can identify an hyperplane which can divide correctly the two classes of points.

An hybrid technique used in the IR field will be presented; it uses two components:

- **K-Means:** model for *unsupervised classification*;
- **SVM:** model for *supervised classification*.

An hybrid technique used in the IR field will be presented; it uses two components:

- **K-Means:** model for *unsupervised classification*;
- **SVM:** model for *supervised classification*.

K-means algorithm is one of the most used clustering algorithms. It consists of partitioning unlabeled objects into k classes, where k is no predefined.

A brief explanation will follow.



K-means

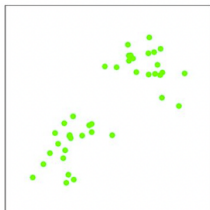
The aim is to create clusters that contain *similar documents*.

Given a training set $x^{(1)}, x^{(2)}, \dots, x^{(m)}$, the algorithm used is the following:

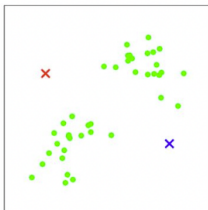
The aim is to create clusters that contain *similar documents*.

Given a training set $x^{(1)}, x^{(2)}, \dots, x^{(m)}$, the algorithm used is the following:

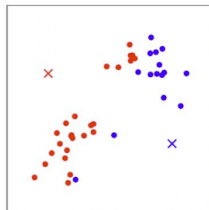
```
1:  $K \leftarrow \text{randomValue}$                                 ▷ multiple K will be tried
2: procedure K-MEANS( $K$ )
3:   initialize cluster centroids  $\mu_1, \dots, \mu_k$  randomly.
4:   repeat until convergence{
5:     for  $i \leftarrow 1$  to  $m$  do
6:        $c^{(i)} = \arg \min_j \|x^{(i)} - \mu_j\|$            ▷ Closest centroid to  $x^{(i)}$ 
7:     end for
8:     for  $k \leftarrow 1$  to  $K$  do
9:        $\mu_j = \frac{\sum_{i=1}^m \mathbf{1}\{c^{(i)}=j\} x^{(i)}}{\sum_{i=1}^m \mathbf{1}\{c^{(i)}=j\}}$        ▷ New centroid of cluster  $k$ 
10:    end for
11:  }
12: end procedure
```



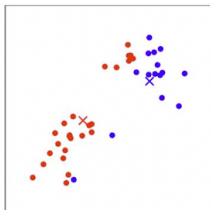
(a)



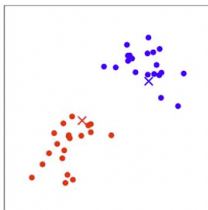
(b)



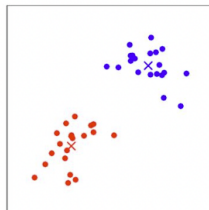
(c)



(d)



(e)



(f)