



Department of Mathematics, University of Udine

---

# SVM and their use in IR

Lorenzo Zanolin.

[lorenzo.zanolin@spes.uniud.it](mailto:lorenzo.zanolin@spes.uniud.it)

May x, 2023



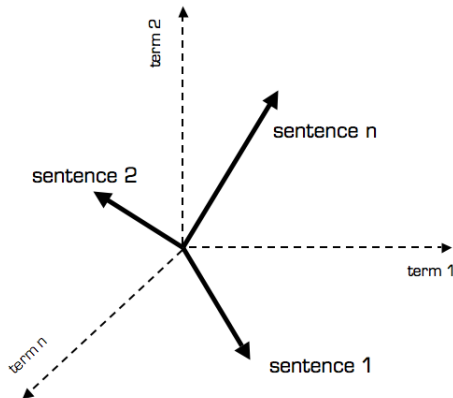
- 1 Linearly separable data
- 2 Non linearly separable data



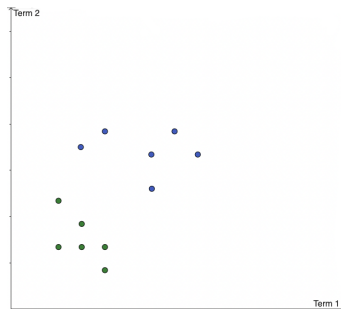
# Vector Space

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.

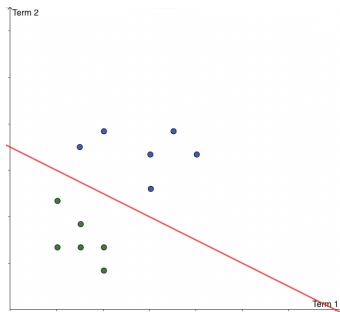
Suppose you have to classify whether a document is *relevant* or not. We can think to use *terms* as features to divide properly the data. We will use a *t-dimensional* vector space to represent our documents.



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  $a + bt_1 + ct_2 = 0$ .

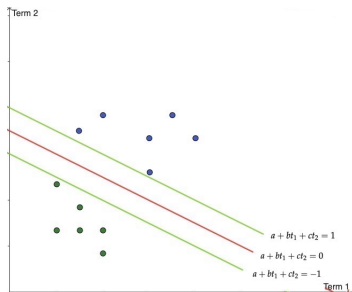
We can introduce two parallel hyperplanes (lines) to the decision boundary, called *support vectors* whose equations are

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

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

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



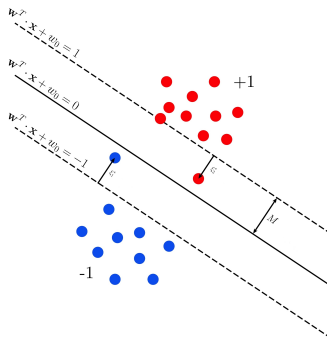


The larger is  $C$  the stricter the classification is, since a larger  $C$  will give more evidence to slack variables.

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

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

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



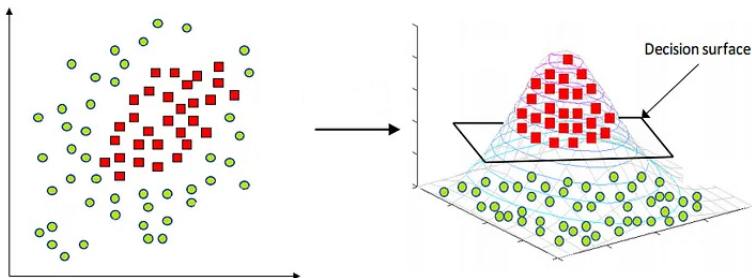


# Non linearly separable data

What if our data is not linearly separable?

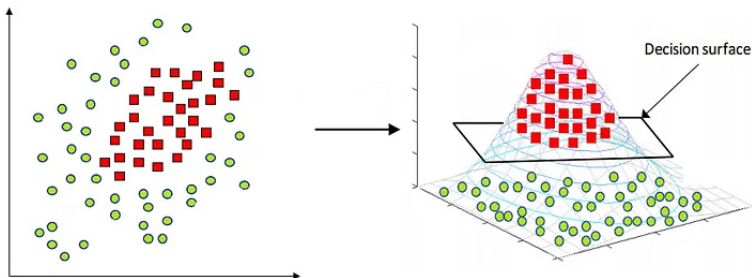
# Non linearly separable data

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...



Here *Kernel Trick* comes handy.



Here *Kernel Trick* comes handy.

Consider the function  $\phi : \mathbb{R}^k \mapsto \mathbb{R}^{k^2}$



Here *Kernel Trick* comes handy.

Consider the function  $\phi : \mathbb{R}^k \mapsto \mathbb{R}^{k^2}$

It allows us to operate in the original feature space without computing the coordinates of the data in a higher dimensional space.