

Support Vector Machines - Passerini

Mattia Carolo - @Carolino96

Support Vector Machines

Support Vector Machine, SVM from now on, are linear classifiers that separate using a **large margin classifier** which solution depends only on a small subset of the training examples called **support vector**. It's very important to note that it has a sound generalization theory (not to study) and they can be easily extended to non linear separation retaining the separation properties thanks to *kernel machines*.

Maximum margin Classifier & Hard Margin SVM

Let's try to formalize the margin. We already know that $yf(x)$ is the confidence on the correct prediction, if negative the prediction is wrong otherwise if positive correct and the value is the confidence on the prediction. Now suppose we have a classifier that correctly separates with no training errors. If this is the case the minimum value among the training examples is called *confidence margin* and it's written like

$$\rho = \min_{(\mathbf{x}, y) \in D} yf(\mathbf{x})$$

Since it depends on w we can compute the distance from the minimal distance to our classifier and it's called **geometric margin** which is formalized like

$$\frac{\rho}{\|\mathbf{w}\|} = \min_{(\mathbf{x}, y) \in D} \frac{yf(\mathbf{x})}{\|\mathbf{w}\|}$$

Ideally we want to maximize the last formula in order to get w in order to maximize the margin. However if we put in an optimization problem we have actually one degree of freedom that is being removed. Suppose we have a solution where

$$\mathbf{w}^T \mathbf{x} + w_0 = 0$$

now if we want to characterize further the plane we can, for example, multiply the terms with an $\alpha \neq 0$ and still we will return to a formula that look like

before since we can incorporate the α in our formalization. This is because there is an infinite number of equivalent formulation for the same hyperplane even with different parameters.

We can counter this problem through the introduction of the *canonical hyperplane* in which we set the constraint that ρ must be equal to a number given a priori (in our case we take 1) in order to get:

$$\rho = \min_{(\mathbf{x}, y) \in D} yf(\mathbf{x}) = 1$$

and it's geometric margin will be $\frac{\rho}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}$

the numerical value in the geometric margin must match

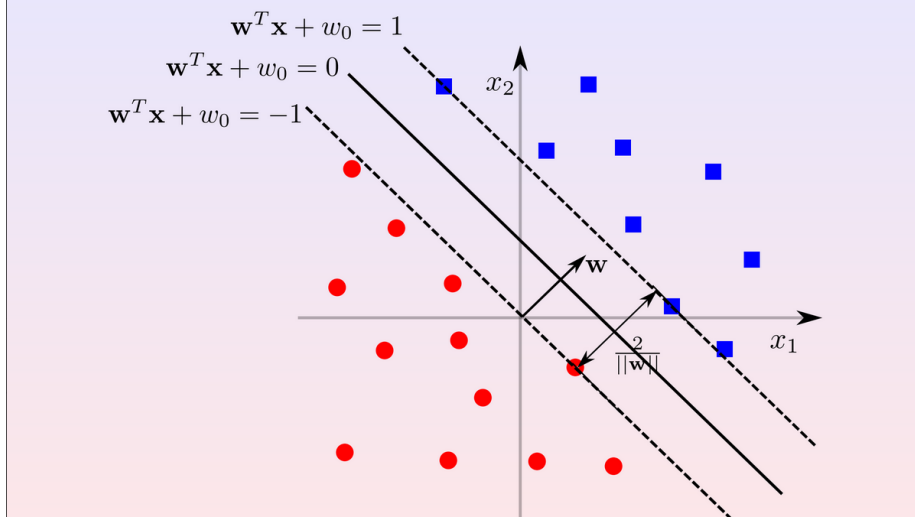


Figure 1: MLClassifier.png

As we can see from the image above the two dotted lines are the two canonical hyperplanes with their ρ set to 1 so summing their respective geometric margin we get that the total geometric margin is equal to $\frac{2}{\|\mathbf{w}\|}$.

We can take this and convert it to an optimization problem.

First of all we want to maximize the margin so $\frac{2}{\|\mathbf{w}\|}$ and we want to do it by enforcing all examples to stay on the correct part of the hyperplane for both canonical ones. Formalized will be

$$\max \frac{2}{\|\mathbf{w}\|} \text{ s.t. } \forall x_i : y_i = 1 \Rightarrow \mathbf{w}^T x_i + w_0 \geq 1 \quad \& \quad \forall x_i : y_i = -1 \Rightarrow \mathbf{w}^T x_i + w_0 \leq -1$$

and the term to maximize can be inverted in order to get

$$\min \frac{\|\mathbf{w}\|}{2} = \frac{\sqrt{w^T w}}{2}$$

which is not a quadratic function but it's monotonic so if we found a maximum it will be the same even squared so we can minimize doing $\min \frac{\|\mathbf{w}\|^2}{2}$ and to summarize the constraints we can just say $yf(x) \geq 1$ since it's our confidence

Margin Error Bound (just a citation not study material)

$$\text{Margin Error Bound: } \nu + \sqrt{\frac{c}{m} \left(\frac{R^2 \wedge^2}{\rho^2} \ln^2 m + \ln\left(\frac{1}{\delta}\right) \right)}$$

The probability of test error so depends on:

- ν is number of margin errors (samples that are outside the confidence margin, correctly classified samples with low confidence)
- m training example in the $\sqrt{\frac{\ln^2 m}{m}}$ so the result goes down if m goes up
- R is the radius of the space containing all the samples
- larger the margin ρ , the smaller test error (so we want the margin $\frac{2}{\|\mathbf{w}\|}$ to be large)
if ρ is fixed to 1, maximizing margin corresponds to minimizing $\|\mathbf{w}\|$
- c is a constant
it makes an upper bound of the generalization error (?)

The name **hard margin** is because we require all examples to be at confidence margin at least one.

Learning Problem

The learning problem is formalized like $\min \frac{\|\mathbf{w}\|^2}{2}$ with linear constraints in w $y_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1, \forall (\mathbf{x}_i, y_i) \in D$. Still this is a quadratic optimization problem which means that is convex and it has only one global optimum. Problem now is that we need to minimize respect to the constraints and one way to do this is the **KKT approach**

Karush-Kuhn-Tucker (KKT) approach

With this approach basically we turn a *constrained problem* into an *unconstrained* one with the same solution. To do this suppose we have $f(z)$ to minimize with some constraints like $g_i(z) \geq 0 \forall i$. Now how can we get rid of the constraints? To do so we introduce a non negative variable called **Lagrange multiplier** noted with $\alpha_i \geq 0$ for each constraint and we rewrite the optimization problem as a **Lagrangian**:

$$\min_z \max_{\alpha \geq 0} f(z) - \sum_i \alpha_i g_i(z)$$

If we find an optimum of this lagrangian called z^* it's still an optimum for the original constrained problem. That's because suppose we find a solution called z' than:

- if at least one constraint is not satisfied ($\exists i \mid g_i(z') < 0$), maximizing over α_i leads to an infinite value;
- if all constraints are satisfied, maximizing over α sets all elements in the sum to zero so that z' is a solution for $\min_z f(z)$.

Applying the approach to our learning problem we will get that

$$\begin{aligned} & \min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|^2 \\ & \text{subject to:} \\ & y_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1 \\ & \forall (\mathbf{x}_i, y_i) \in D \end{aligned}$$

where :

- z will be our \mathbf{w} so $f(z)$ will be $\frac{1}{2} \|\mathbf{w}\|^2$
- the constraint will be $g_i(z) \geq 0$ so we need to turn $y_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1$ into the g_i which will become $y_i(\mathbf{w}^T \mathbf{x}_i + w_0) - 1 \geq 0$

By substituting the new obtained terms we get

$$L(\mathbf{w}, w_0, \alpha) = \frac{\|\mathbf{w}\|^2}{2} - \sum_{i=1}^m \alpha_i (y_i(\mathbf{w}^T \mathbf{x}_i + w_0) - 1)$$

$$m = |D|$$

This lagrangian should be minimized with respect to \mathbf{w} and w_0 and maximized with respect to α_i . The solution is called saddle point and it's located where there is the solution for both problems

Going further with the calculus we get that

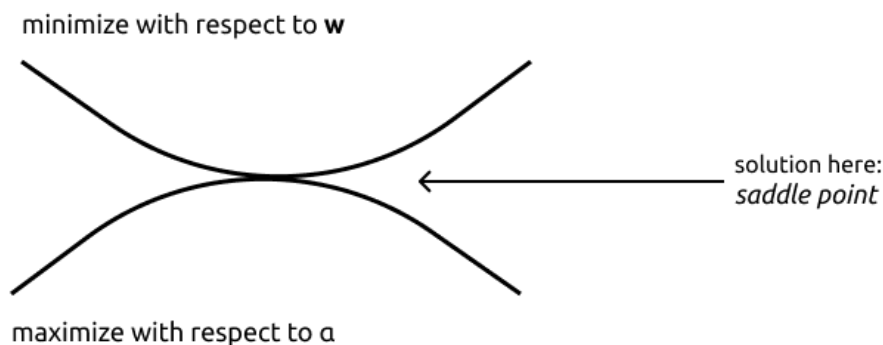


Figure 2: MLClassifier.png

$$L(\mathbf{w}, w_0, \alpha) = \frac{\|\mathbf{w}\|^2}{2} - \sum_{i=1}^m \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + w_0) - 1)$$

which is our lagrangian. Now we want to minimize with respect to \mathbf{w}, w_0 and maximize with respect to α . We are interested in computing the gradient respect to our **primal variables** \mathbf{w} and w_0

$$\text{So we take } \nabla_{\mathbf{w}} L = \nabla_{\mathbf{w}} \frac{\mathbf{w}^T \mathbf{w}}{2} - \nabla_{\mathbf{w}} \sum_i \alpha_i y_i \mathbf{w}^T \mathbf{x} = \frac{\mathbf{w}}{1} - \sum_i \alpha_i y_i \mathbf{x}$$

and now we set $\mathbf{w} - \sum_i \alpha_i y_i \mathbf{x} = 0$ getting $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}$, but this is not the solution because \mathbf{w} is defined in terms of α

Then we can also take the derivative in order to get

$$\frac{\delta L}{\delta w_0} = \frac{\delta(-\sum_i \alpha_i y_i w_0)}{\delta w_0} = -\sum_i \alpha_i y_i = 0 \rightarrow \sum_i \alpha_i y_i = 0$$

CERCASI ANIMA PIA CHE TRASCRIVI LA FORMULA LEZ 18.6 DAL MIN 46:00 O DALLE SLIDE

After all the steps we get a function which is only a function of the dual variables (the alphas) without our primal variable which is like

$$-\frac{1}{2} \sum_i \sum_j \alpha_i y_i \alpha_j y_j x_i^T x_j + \sum_i \alpha_i$$

which needs to be maximized respect to α with constraints of

$$\begin{aligned}\alpha_i &\geq 0 \quad \forall i \\ \sum_i \alpha_i y_i &= 0\end{aligned}$$

But still this result is a quadratic optimization problem respect to alpha. In all of this we can see that the beforementioned **primal variables** are missing and replaced with a new pair of variables which will be called **dual variables**(the alphas). This new type of formalization is called **dual formulation**

The result is that $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ can be written both in form of the primal and of the dual because we know that \mathbf{w} is equal to $\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i$

Decision function

When we did the gradient with respect to \mathbf{w} previously we got $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$.

Now if we plug it into our $f(\mathbf{x})$ we get that

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = \sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + w_0$$

The decision $f(\mathbf{x})$ (defined as **decision function**) on \mathbf{x} is basically taken as a linear combination of dot products between training points and \mathbf{x} , so if \mathbf{x}_i is similar to \mathbf{x} it will have a high dot product because here the dot product works kind of a similarity between the points. Plus the weights of the combination are $\alpha_i y_i$ where large α_i implies large contribution toward class y_i

KKT conditions

To understand wheter a training examples contributes or not to our decision function can be found by applying KKT. The formulation remains the same so:

$$L(\mathbf{w}, w_0, \alpha) = \frac{\|\mathbf{w}\|^2}{2} - \sum_{i=1}^m \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + w_0) - 1)$$

In the optimal solution $\alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + w_0) - 1)$ should be = 0, either:

- $\alpha_i = 0$, so the example \mathbf{x}_i does not contribute to the final solution
- if $\alpha_i > 0$ than $y_i (\mathbf{w}^T \mathbf{x}_i + w_0) = 1$, so the confidence for the example should be 1

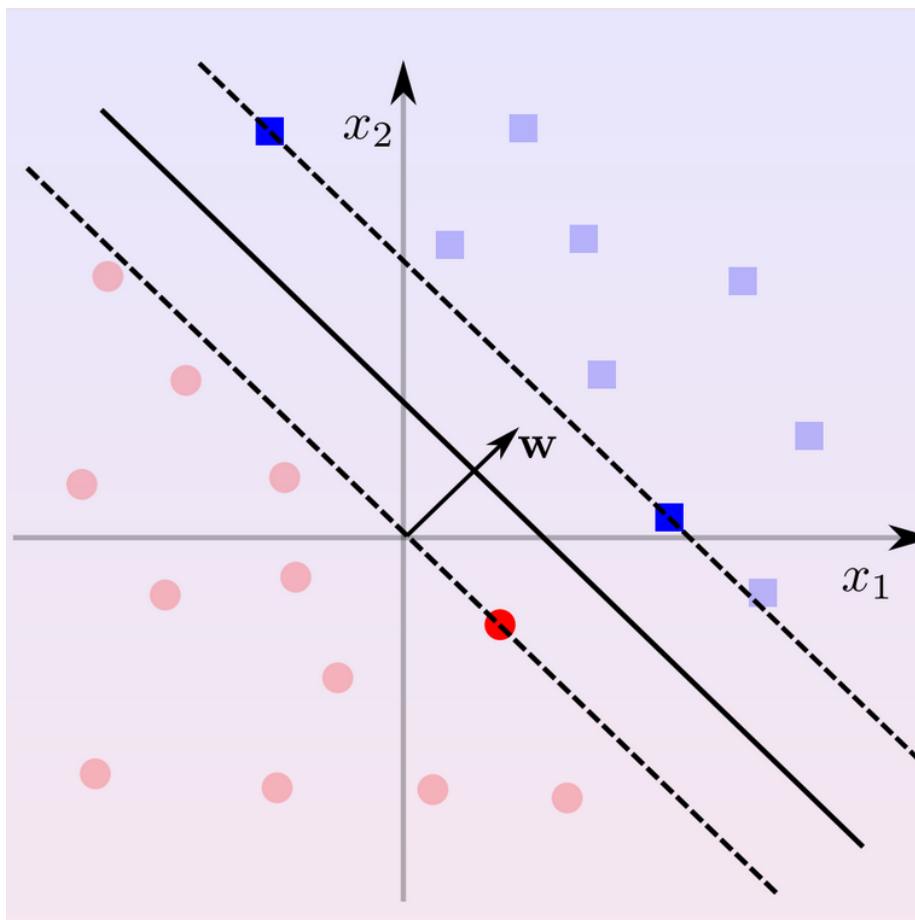


Figure 3: SVM.png

Graphically we will achieve something like this:

The *support vectors* are where $\alpha_1 > 0$ that lays in the hyperplanes where *confidence* $f(\mathbf{x}) = 1$