

Support vector machines

Course of Machine Learning
Master Degree in Computer Science
University of Rome “Tor Vergata”
a.a. 2023-2024

Giorgio Gambosi

Idea

The binary classification problem is approached in a direct way, that is:

We try and find a plane that separates the classes in feature space (indeed, a “best” plane, according to a reasonable characteristic)

If this is not possible, we get creative in two ways:

- We soften what we mean by “separates”, and
- We enrich and enlarge the feature space so that separation is (more) possible

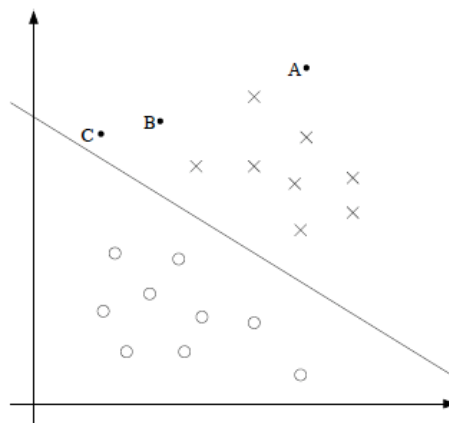
Most classification algorithms output not just a class label but a score. For example, logistic regression returns the class probability

$$p(y = 1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0) \in [0, 1]$$

If the class probability is > 0.5 , the model outputs class 1. The score is an estimate of confidence; it also represents how far we are from the decision boundary.

Intuitively, we want to select boundaries with high **margin**. This means that we are as confident as possible for every point and we are as far as possible from the decision boundary.

For example, in the following case,



A can be assigned to \mathcal{C}_1 with greater confidence than B and even greater confidence than C .

Binary classifiers

Consider a binary classifier which, for any element \mathbf{x} , returns a value $y \in \{-1, 1\}$, where we assume that \mathbf{x} is assigned to \mathcal{C}_0 if $y = -1$ and to \mathcal{C}_1 if $y = 1$.

Moreover, we consider linear classifier such as

$$h(\mathbf{x}) = g(\mathbf{w}^T \mathbf{x} + w_0)$$

where $g(z) = 1$ if $z \geq 0$ and $g(z) = -1$ if $z < 0$. The prediction on the class of \mathbf{x} is then provided by deriving a value in $\{-1, 1\}$ just as in the case of a perceptron, that is with no estimation of the probabilities $p(\mathcal{C}_i|\mathbf{x})$ that \mathbf{x} belongs to each class.

How can we define the concept of margin more formally?

For any training set item (\mathbf{x}_i, t_i) , the **functional margin** of $\bar{\mathbf{w}} = \mathbf{w} \cup \{w_0\}$ wrt such item is defined as

$$\bar{\gamma}_i = t_i(\mathbf{w}^T \mathbf{x}_i + w_0)$$

Observe that the resulting prediction is correct iff $\bar{\gamma}_i > 0$. Moreover, larger values of $\bar{\gamma}_i$ denote greater confidence on the prediction.

Given a training set $\mathcal{T} = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_n, t_n)\}$ the functional margin of $\bar{\mathbf{w}}$ wrt \mathcal{T} is the minimum functional margin for all items in \mathcal{T}

$$\bar{\gamma} = \min_i \bar{\gamma}_i$$

However, we have a problem.

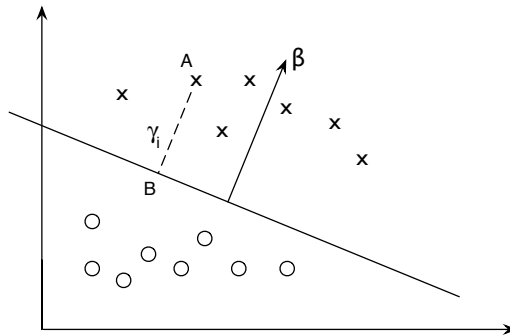
- If we rescale the parameters \mathbf{w}, w_0 by a scalar $\alpha > 0$, we get new parameters $\alpha\mathbf{w}, \alpha w_0$.
- Using $\alpha\mathbf{w}, \alpha w_0$ doesn't change the classification of points.
- However, the margin $(\alpha\mathbf{w}^T \mathbf{x}_i + \alpha w_0) = \alpha(\mathbf{w}^T \mathbf{x}_i + w_0)$ is now scaled by α !

It doesn't make sense that the same classification boundary can have different margins when we rescale it.

We define the **geometric margin** γ_i with respect to a training example (\mathbf{x}_i, t_i) as

$$\gamma_i = t_i \left(\frac{\mathbf{w}^T \mathbf{x}_i}{\|\mathbf{w}\|} + \frac{w_0}{\|\mathbf{w}\|} \right) = \frac{\bar{\gamma}_i}{\|\mathbf{w}\|}$$

We call it geometric because γ_i equals the (signed) distance between \mathbf{x}_i and the hyperplane defined by \mathbf{w}, w_0 , that is as the length of the line segment from \mathbf{x}_i to its projection on the boundary hyperplane



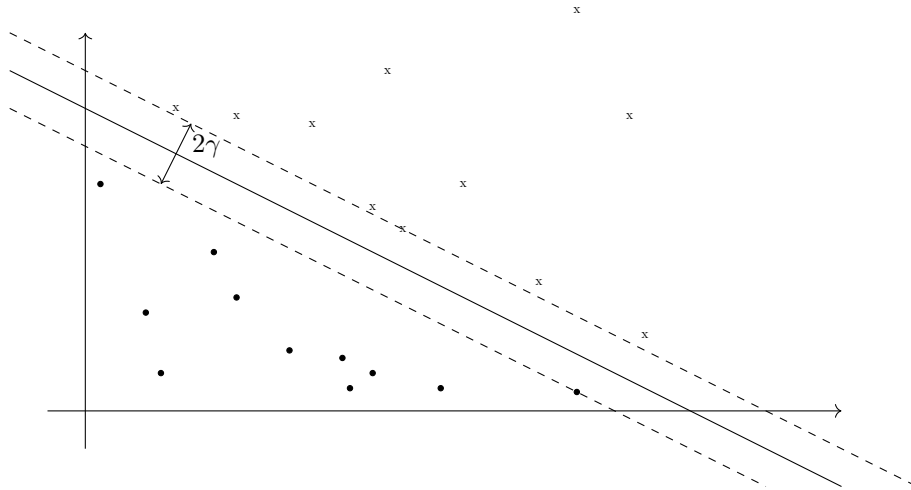
Differently from $\bar{\gamma}_i$, the geometric margin γ_i is invariant wrt parameter scaling. In fact, by substituting $\alpha\mathbf{w}$ to \mathbf{w} and αw_0 to w_0 , we get

$$\begin{aligned} \bar{\gamma}_i &= t_i(\alpha\mathbf{w}^T \mathbf{x}_i + \alpha w_0) = \alpha t_i(\mathbf{w}^T \mathbf{x}_i + w_0) \\ \gamma_i &= t_i \left(\frac{\alpha\mathbf{w}^T \mathbf{x}_i}{\|\alpha\mathbf{w}\|} + \frac{\alpha w_0}{\|\alpha\mathbf{w}\|} \right) = t_i \left(\frac{\mathbf{w}^T \mathbf{x}_i}{\|\mathbf{w}\|} + \frac{w_0}{\|\mathbf{w}\|} \right) \end{aligned}$$

- The geometric margin wrt the training set $\mathcal{T} = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_n, t_n)\}$ is then defined as the smallest geometric margin for all items (\mathbf{x}_i, t_i)

$$\gamma = \min_i \gamma_i$$

- a useful interpretation of γ is as half the width of the largest strip, centered on the hyperplane $\mathbf{w}^T \mathbf{x} + w_0 = 0$, containing none of the points $\mathbf{x}_1, \dots, \mathbf{x}_n$
- the hyperplanes on the boundary of such strip, each at distance γ from the hyperplane and passing (at least one of them) through some point \mathbf{x}_i are said **maximum margin hyperplanes**.



Optimal margin classifiers

Given a training set \mathcal{T} , we wish to find the hyperplanes which separates the two classes (if one does exist) and has maximum γ : by making the distance between the hyperplanes and the set of points corresponding to elements as large as possible, the confidence on the provided classification increases.

Assume classes are linearly separable in the training set: hence, there exists a hyperplane (an infinity of them, indeed) separating elements in C_1 from elements in C_2 . In order to find the one among those hyperplanes which maximizes γ , we have to solve the following optimization problem

$$\begin{aligned} & \max_{\bar{\mathbf{w}}, w_0} \gamma \\ & \text{where } \gamma_i = \frac{t_i}{\|\bar{\mathbf{w}}\|} (\bar{\mathbf{w}}^T \mathbf{x}_i + w_0) \geq \gamma \quad i = 1, \dots, n \end{aligned}$$

which maximises the smallest margin over the (\mathbf{x}_i, t_i) . It guarantees each point has margin at least γ .

This corresponds to,

$$\begin{aligned} & \max_{\bar{\mathbf{w}}, w_0} \gamma \\ & \text{where } t_i (\bar{\mathbf{w}}^T \mathbf{x}_i + w_0) \geq \gamma \|\bar{\mathbf{w}}\| \quad i = 1, \dots, n \end{aligned}$$

As observed, if all parameters are scaled by any constant α , all geometric margins γ_i between elements and hyperplane are unchanged: we may then exploit this freedom to introduce the constraint

$$\gamma = \min_i t_i (\bar{\mathbf{w}}^T \mathbf{x}_i + w_0) = 1$$

This can be obtained by assuming $\|w\| = \frac{1}{\gamma}$, which corresponds to considering a scale where the maximum margin has width 2. This results, for each element \mathbf{x}_i, t_i , into a constraint

$$\gamma_i = t_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1$$

An element (point) is said **active** if the equality holds, that is if

$$t_i(\mathbf{w}^T \mathbf{x}_i + w_0) = 1$$

and **inactive** if this does not hold. Observe that, by definition, there must exist at least one active point.

For any item (\mathbf{x}, t) , different cases correspond to the value of the signed linear combination $t(\mathbf{w}^T \mathbf{x} + w_0)$:

1. $t(\mathbf{w}^T \mathbf{x} + w_0) > 1 \Rightarrow \mathbf{x}$ is in the region corresponding to its class, outside the margin strip
2. $t(\mathbf{w}^T \mathbf{x} + w_0) = 1 \Rightarrow \mathbf{x}$ is in the region corresponding to its class, on the maximum margin hyperplane
3. $0 < t(\mathbf{w}^T \mathbf{x} + w_0) < 1 \Rightarrow \mathbf{x}$ is in the region corresponding to its class, inside the margin strip
4. $t(\mathbf{w}^T \mathbf{x} + w_0) = 0 \Rightarrow \mathbf{x}$ is on the separating hyperplane
5. $-1 < t(\mathbf{w}^T \mathbf{x} + w_0) < 0 \Rightarrow \mathbf{x}$ is in the region corresponding to the other class, inside the margin strip
6. $t(\mathbf{w}^T \mathbf{x} + w_0) = -1 \Rightarrow \mathbf{x}$ is in the region corresponding to the other class, on the maximum margin hyperplane
7. $t(\mathbf{w}^T \mathbf{x} + w_0) < -1 \Rightarrow \mathbf{x}$ is in the region corresponding to the other class, outside the margin strip

The optimization problem, is then transformed into

$$\begin{aligned} \max_{\mathbf{w}, w_0} \gamma &= \|\mathbf{w}\|^{-1} \\ \text{where } t_i(\mathbf{w}^T \mathbf{x}_i + w_0) &\geq 1 \quad i = 1, \dots, n \end{aligned}$$

Maximizing $\|\mathbf{w}\|^{-1}$ is equivalent to minimizing $\|\mathbf{w}\|^2$: hence we may formulate the problem as

$$\begin{aligned} \min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{where } t_i(\mathbf{w}^T \mathbf{x}_i + w_0) &\geq 1 \quad i = 1, \dots, n \end{aligned}$$

This is a **convex quadratic** optimization problem. The function to be minimized is in fact convex and the set of points satisfying the constraint is a convex polyhedron (intersection of half-spaces).

Duality

From optimization theory it derives that, given the problem structure (linear constraints + convexity):

- there exists a **dual formulation** of the problem
- the optimum of the dual problem is the same of the original (**primal**) problem

Lagrangian

Consider the optimization problem

$$\min_{\mathbf{x} \in \Omega} f(\mathbf{x})$$

where Ω is the feasible region, defined by the constraints

$$g_i(\mathbf{x}) \leq 0 \quad i = 1, \dots, k$$

where $f(\mathbf{x}), g_i(\mathbf{x})$ are convex functions and Ω is a convex set.

Define the **Lagrangian**

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \sum_{i=1}^k \lambda_i g_i(\mathbf{x})$$

Consider the maximum wrt to λ

$$\begin{aligned} \max_{\lambda} L(\mathbf{x}, \lambda) &= f(\mathbf{x}) + \max_{\lambda} \sum_{i=1}^k \lambda_i g_i(\mathbf{x}) \\ \lambda_i &\geq 0 \quad i = 1, \dots, k \end{aligned}$$

which is a function of \mathbf{x} .

- if \mathbf{x} is a feasible solution, $g_i(\mathbf{x}) \leq 0$ for all i and the maximum is obtained for $\lambda_i = 0$: as a consequence,

$$\max_{\lambda: \lambda_i \geq 0} L(\mathbf{x}, \lambda) = f(\mathbf{x})$$

- if \mathbf{x} is an infeasible solution, then $g_i(\mathbf{x}) > 0$ for some i and the maximum is unbounded, as λ_i can be arbitrarily great

As a consequence, the maximum is equal to $f(\mathbf{x})$ if \mathbf{x} is feasible, while it is unbounded if \mathbf{x} is not feasible.

This results, in the case that an optimum \mathbf{x}^* exists, into

$$\min_{\mathbf{x} \in \Omega} f(\mathbf{x}) = \min_{\mathbf{x} \in \Omega} \max_{\lambda: \lambda_i \geq 0} L(\mathbf{x}, \lambda)$$

In general, the **weak duality property** holds

$$\max_{\lambda: \lambda_i \geq 0} \min_{\mathbf{x} \in \Omega} L(\mathbf{x}, \lambda) \leq \min_{\mathbf{x} \in \Omega} \max_{\lambda: \lambda_i \geq 0} L(\mathbf{x}, \lambda) = \min_{\mathbf{x} \in \Omega} f(\mathbf{x})$$

where $\max_{\lambda: \lambda_i \geq 0} \min_{\mathbf{x} \in \Omega} L(\mathbf{x}, \lambda)$ is the **dual** problem of $\min_{\mathbf{x} \in \Omega} f(\mathbf{x})$

Moreover, in the case of convex optimization (our case here) the strong duality property holds

$$\max_{\lambda: \lambda_i \geq 0} \min_{\mathbf{x} \in \Omega} L(\mathbf{x}, \lambda) = \min_{\mathbf{x} \in \Omega} \max_{\lambda: \lambda_i \geq 0} L(\mathbf{x}, \lambda) = \min_{\mathbf{x} \in \Omega} f(\mathbf{x})$$

Karush-Kuhn-Tucker conditions

The following necessary and sufficient conditions hold for the existence of an optimum $(\mathbf{x}^*, \lambda^*)$, and can be used to simplify the dual problem definition.

$$\begin{aligned} \frac{\partial L(\mathbf{x}, \lambda)}{\partial \mathbf{x}} \Big|_{\mathbf{x}^*, \lambda^*} &= \mathbf{0} \\ \frac{\partial L(\mathbf{x}, \lambda)}{\partial \lambda_i} \Big|_{\mathbf{x}^*, \lambda^*} &= g_i(\mathbf{x}^*) \geq 0 & i = 1, \dots, k \\ \lambda_i^* &\geq 0 & i = 1, \dots, k \\ \lambda_i^* g_i(\mathbf{x}^*) &= 0 & i = 1, \dots, k \end{aligned}$$

In order for the optimum to be a minimum, the second order condition must hold that the Hessian H_x evaluated at \mathbf{x}^* must be positive definite.

Note: the last condition (**complementary slackness**) states that a Lagrangian multiplier λ_i^* can be non-zero only if $g_i(\mathbf{x}^*) = 0$, that is of \mathbf{x}^* is "at the limit" for the constraint $g_i(\mathbf{x})$. In this case, the constraint is said **active**.

Application to SVM

In our case,

- $f(\mathbf{x})$ corresponds to $\frac{1}{2} \|\mathbf{w}\|^2$
- $g_i(x)$ corresponds to $t_i(\mathbf{w}^T \mathbf{x}_i + w_0) - 1 \geq 0$
- Ω is the intersection of a set of hyperplanes, that is a polyhedron, hence convex.

The corresponding Lagrangian is

$$L(\bar{\mathbf{w}}, \lambda) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \lambda_i (t_i(\mathbf{w}^T \mathbf{x}_i + w_0) - 1) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \lambda_i t_i \mathbf{w}^T \mathbf{x}_i - w_0 \sum_{i=1}^n \lambda_i t_i + \sum_{i=1}^n \lambda_i$$

and the dual problem (with same optimum) is

$$\begin{aligned} \max_{\lambda} \min_{\bar{\mathbf{w}}} L(\bar{\mathbf{w}}, \lambda) \\ \lambda_i \geq 0 \quad i = 1, \dots, k \end{aligned}$$

Applying the KKT conditions

In order to state the dual problem as a function of λ , let us derive the values of the coefficients $\bar{\mathbf{w}}$ which minimize $L(\bar{\mathbf{w}}, \lambda)$.

Such values satisfy the following conditions:

$$\begin{aligned} \frac{\partial L(\bar{\mathbf{w}}, \lambda)}{\partial \mathbf{w}} \Big|_{\bar{\mathbf{w}}^*} &= \mathbf{w}^* - \sum_{i=1}^n \lambda_i t_i \mathbf{x}_i = 0 \\ \frac{\partial L(\bar{\mathbf{w}}, \lambda)}{\partial w_0} \Big|_{\bar{\mathbf{w}}^*} &= \sum_{i=1}^n \lambda_i t_i = 0 \end{aligned}$$

The resulting KKT conditions are

$$\begin{aligned}
\mathbf{w}^* &= \sum_{i=1}^n \lambda_i t_i \mathbf{x}_i \\
\sum_{i=1}^n \lambda_i t_i &= 0 \\
t_i(\mathbf{w}^* \cdot \mathbf{x}_i + w_0^*) - 1 &\geq 0 & i = 1, \dots, n \\
\lambda_i &\geq 0 & i = 1, \dots, n \\
\lambda_i (t_i(\mathbf{w}^* \cdot \mathbf{x}_i + w_0^*) - 1) &= 0 & i = 1, \dots, n
\end{aligned}$$

By substituting the values of the coefficients $\bar{\mathbf{w}}^*$ according to the above equations, we may state that at the optimum $\bar{\mathbf{w}}^*$ the dual problem asks to maximize

$$\begin{aligned}
L(\lambda) &= \frac{1}{2} \|\mathbf{w}^*\|^2 - \sum_{i=1}^n \lambda_i t_i \left(\sum_{j=1}^n \lambda_j t_j (\mathbf{x}_j \cdot \mathbf{x}_i) \right) - w_0^* \sum_{i=1}^n \lambda_i t_i + \sum_{i=1}^n \lambda_i \\
&= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j (\mathbf{x}_i \cdot \mathbf{x}_j) \sum_{i=1}^n - \sum_{j=1}^n \lambda_i t_i \lambda_j t_j (\mathbf{x}_i \cdot \mathbf{x}_j) + \sum_{i=1}^n \lambda_i \\
&= \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j (\mathbf{x}_i \cdot \mathbf{x}_j)
\end{aligned}$$

with the constraints on λ

$$\begin{aligned}
\sum_{i=1}^n \lambda_i t_i &= 0 \\
\lambda_i &\geq 0 & i = 1, \dots, n
\end{aligned}$$

The remaining two KKT conditions can be proved to be always verified.

Dual SVM problem

We have modified the definition of the dual problem by applying the KKT conditions to drop the occurrences of coefficients $\bar{\mathbf{w}}$ from $L(\bar{\mathbf{w}}, \lambda)$.

The new problem has the same optimum of the original **primal**, where the KKT conditions will indeed hold, connecting the values of the optimal solutions of the two problems

$$\begin{aligned}
\max_{\lambda} L(\lambda) &= \max_{\lambda} \left(\sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j t_i t_j (\mathbf{x}_i \cdot \mathbf{x}_j) \right) \\
\lambda_i &\geq 0 & i = \dots, n \\
\sum_{i=1}^n \lambda_i t_i &= 0
\end{aligned}$$

All the considerations above clearly hold if we assume a set of base functions ϕ is applied, thus resulting into the dual problem

$$\begin{aligned} \max_{\lambda} L(\lambda) &= \max_{\lambda} \left(\sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j t_i t_j (\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)) \right) \\ \lambda_i &\geq 0 \quad i = \dots, n \\ \sum_{i=1}^n \lambda_i t_i &= 0 \end{aligned}$$

By defining the **kernel function**

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

the dual problem's formulation can be given as

$$\begin{aligned} \max_{\lambda} L(\lambda) &= \max_{\lambda} \left(\sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j t_i t_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \right) \\ \lambda_i &\geq 0 \quad i = 1, \dots, n \\ \sum_{i=1}^n \lambda_i t_i &= 0 \end{aligned}$$

Passing from primal to dual

Disadvantage The number variables increases from m to n (in particular, if $\mathbf{x} = \mathbf{x}$, from d to n).

Advantage The number of variables to be considered, which are relevant for classification, turns out to be quite smaller than n .

Deriving coefficients

By solving the dual problem, the optimal values of Lagrangian multipliers λ^* are obtained.

The optimal values of parameters \mathbf{w}^* are then derived through the relations

$$w_i^* = \sum_{j=1}^n \lambda_j^* t_j \phi_i(\mathbf{x}_j) \quad i = 1, \dots, m$$

The value of w_0^* can be obtained by observing that, for any support vector \mathbf{x}_k (characterized by the condition $\lambda_k \geq 0$), it must be

$$\begin{aligned} 1 &= t_k (\phi(\mathbf{x}_k)^T \mathbf{w}^* + w_0^*) = t_k \left(\sum_{j=1}^n \lambda_j^* t_j \phi(\mathbf{x}_j)^T \phi(\mathbf{x}_k) + w_0^* \right) \\ &= t_k \left(\sum_{j=1}^n \lambda_j^* t_j \kappa(\mathbf{x}_j, \mathbf{x}_k) + w_0^* \right) = t_k \left(\sum_{j \in \mathcal{S}} \lambda_j^* t_j \kappa(\mathbf{x}_j, \mathbf{x}_k) + w_0^* \right) \end{aligned}$$

where \mathcal{S} is the set of indices of support vectors.

As a consequence, since $t_k = \pm 1$, in order to have a unitary product it must be

$$t_k = \sum_{j \in \mathcal{S}} \lambda_j^* t_j \kappa(\mathbf{x}_j, \mathbf{x}_k) + w_0^*$$

and

$$w_0^* = t_k - \sum_{j \in \mathcal{S}} \lambda_j^* t_j \kappa(\mathbf{x}_j, \mathbf{x}_k)$$

A more precise solution can be obtained as the mean value obtained considering all support vectors

$$w_0^* = \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left(t_i - \sum_{j \in \mathcal{S}} \lambda_j^* t_j \kappa(\mathbf{x}_j, \mathbf{x}_i) \right)$$

Classification through SVM

A new element \mathbf{x} can be classified, given a set of base functions ϕ or a kernel function κ , by checking the sign of

$$y(\mathbf{x}) = \sum_{i=1}^m w_i^* \phi_i(\mathbf{x}) + w_0^* = \sum_{j=1}^n \lambda_j^* t_j \kappa(\mathbf{x}_j, \mathbf{x}) + w_0^*$$

As noticed, if \mathbf{x}_i is not a support vector, then it must be $\lambda_i^* = 0$. Thus, the above sum can be written as

$$y(\mathbf{x}) = \sum_{j \in \mathcal{S}} \lambda_j^* t_j \kappa(\mathbf{x}_j, \mathbf{x}) + w_0^*$$

The classification performed through the dual formulation, using the kernel function, does not take into account all training set items, but only support vectors, usually a quite small subset of the training set.

Non separability in the training set

- The approach described before, when applied to non linearly separable sets, does not provide acceptable solutions: it is in fact impossible to satisfy all constraints

$$t_i(\mathbf{w}^T \phi(\mathbf{x}_i) + w_0) \geq 1 \quad i = 1, \dots, n$$

- These constraints must then be relaxed in order to allow them to not hold, at the cost of some increase in the objective function to be minimized
- A **slack variable** ξ_i is introduced for each constraint, to provide a measure of how much the constraint is not verified

This can be formalized as

$$\begin{aligned} \min_{\mathbf{w}, w_0, \boldsymbol{\xi}} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \xi_i \\ & t_i(\mathbf{w}^T \phi(\mathbf{x}_i) + w_0) \geq 1 - \xi_i \quad i = 1, \dots, n \\ & \xi_i \geq 0 \quad i = 1, \dots, n \end{aligned}$$

where $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n)$

- By introducing suitable multipliers, the following Lagrangian can be obtained

$$\begin{aligned} L(\mathbf{w}, w_0, \boldsymbol{\xi}, \boldsymbol{\lambda}, \boldsymbol{\alpha}) &= \\ &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \lambda_i (y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + w_0) - 1 + \xi_i) - \sum_{i=1}^n \alpha_i \xi_i \\ &= \frac{1}{2} \sum_{i=1}^n w_i^2 + \sum_{i=1}^n (C - \alpha_i) \xi_i - \sum_{i=1}^n \lambda_i (t_i (\sum_{j=1}^m w_j \phi_j(\mathbf{x}_i)) + w_0) - 1 + \xi_i \\ &= \frac{1}{2} \sum_{i=1}^n w_i^2 + \sum_{i=1}^n (C - \alpha_i - \lambda_i) \xi_i - \sum_{i=1}^n \sum_{j=1}^m \lambda_i t_i w_j \phi_j(\mathbf{x}_i) + w_0 \sum_{i=1}^n \lambda_i t_i + \sum_{i=1}^n \lambda_i \end{aligned}$$

where $\alpha_i \geq 0$ and $\lambda_i \geq 0$, for $i = 1 \dots, n$.

The Karush-Kuhn-Tucker conditions are now:

$$\begin{array}{ll}
\frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}, w_0, \boldsymbol{\xi}, \lambda, \boldsymbol{\alpha}) = \mathbf{0} & \text{null gradient} \\
\frac{\partial}{\partial w_0} L(\mathbf{w}, w_0, \boldsymbol{\xi}, \lambda, \boldsymbol{\alpha}) = 0 & \text{null gradient} \\
\frac{\partial}{\partial \boldsymbol{\xi}} L(\mathbf{w}, w_0, \boldsymbol{\xi}, \lambda, \boldsymbol{\alpha}) = \mathbf{0} & \text{null gradient} \\
t_i(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) + w_0) - 1 + \xi_i \geq 0 & i = 1, \dots, n & \text{constraints} \\
\xi_i \geq 0 & i = 1, \dots, n & \text{constraints} \\
\lambda_i \geq 0 & i = 1, \dots, n & \text{multipliers} \\
\alpha_i \geq 0 & i = 1, \dots, n & \text{multipliers} \\
\lambda_i (t_i(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) + w_0) - 1 + \xi_i) = 0 & i = 1, \dots, n & \text{complementary slackness} \\
\alpha_i \xi_i = 0 & i = 1, \dots, n & \text{complementary slackness}
\end{array}$$

Deriving a dual formulation

From the null gradient conditions wrt w_i, b, ξ_j it derives

$$\begin{aligned}
w_i &= \sum_{j=1}^n \lambda_j t_j \phi_i(\mathbf{x}_j) & i = 1, \dots, m \\
0 &= \sum_{i=1}^n \lambda_i t_i \\
\lambda_i &= C - \alpha_i \leq C & i = 1, \dots, n
\end{aligned}$$

By plugging the above relations into $L(\bar{\mathbf{w}}, \boldsymbol{\xi}, \lambda, \boldsymbol{\alpha})$, the dual problem results

$$\begin{aligned}
\max_{\lambda} \tilde{L}(\lambda) &= \max_{\lambda} \left(\sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j t_i t_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \right) \\
0 &\leq \lambda_i \leq C & i = 1, \dots, n \\
\sum_{i=1}^n \lambda_i y_i &= 0
\end{aligned}$$

Observe that the only difference wrt the linearly separable case is given by constraints $0 \leq \lambda_i$ transformed into $0 \leq \lambda_i \leq C$

Classification

From the optimal solution λ^* of the dual problem, the coefficients \mathbf{w}^* and w_0^* can be derived just as done in the linearly separable case.

A new element \mathbf{x} can then be classified, again, through the sign of

$$y(\mathbf{x}) = \sum_{i=1}^m w_i^* \phi_i(\mathbf{x}) + w_0^*$$

or, equivalently, of

$$y(\mathbf{x}) = \sum_{i \in \mathcal{S}} \lambda_j^* t_j \kappa(\mathbf{x}_i, \mathbf{x}_j) + w_0^*$$

Extensions

The approach can be extended to

- More than 2 classes (multiclass classification): solve one vs all binary classification problem for all classes
- Real-valued outputs (support vector regression)

Computational issues

- Training time of the standard SVM is $O(n^3)$ (solving QP)
 - Can be prohibitive for large datasets
- Lots of research has gone into speeding up the SVMs
 - Many approximate QP solvers are used to speed up SVMs
 - Gradient descent faster and with possibility of limiting the computation time

SVM and gradient descent

Recall the formalization of the problem in the general case

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \xi_i \\ & t_i(\mathbf{w}^T \phi(\mathbf{x}_i) + w_0) \geq 1 - \xi_i \quad i = 1, \dots, n \\ & \xi_i \geq 0 \quad i = 1, \dots, n \end{aligned}$$

Given \mathbf{w} , w_0 , the slack variable ξ_i is minimized as

$$\xi_i = \begin{cases} 0 & t_i(\mathbf{w}^T \phi(\mathbf{x}_i) + w_0) \geq 1 \\ 1 - t_i(\mathbf{w}^T \phi(\mathbf{x}_i) + w_0) & \text{otherwise} \end{cases}$$

The optimal value of ξ_i corresponds to the **hinge loss** of the corresponding item

$$L_H(\mathbf{w}, w_0, \mathbf{x}_i, t_i) = \max(0, 1 - t_i(\mathbf{w}^T \phi(\mathbf{x}_i) + w_0))$$

We may then define the cost function to be minimized as

$$\begin{aligned} C(\mathbf{w}) &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n L_H(\mathbf{w}, w_0, \mathbf{x}_i, t_i) \\ &\propto \sum_{i=1}^n L_H(\mathbf{w}, w_0, \mathbf{x}_i, t_i) + \frac{1}{2C} \|\mathbf{w}\|^2 \end{aligned}$$

That is, SVM correspond to hinge loss with ridge regularization

Since hinge loss is not differentiable at $x = 1$, as discussed above, subgradient descent can be applied to iteratively find the optimal solution, with

$$\frac{\partial L_H}{\partial w_i} = w_i - \sum_{\mathbf{x}_k \in L} t_k \phi_i(\mathbf{x}_k)$$

where $\mathbf{x}_k \in L$ iff $t_k(\mathbf{w}^T \phi(\mathbf{x}_k) + w_0) < 1$.

The resulting iteration is

$$w_i^{(r+1)} = w_i^{(r)} - \eta w_i^{(r)} + \eta \sum_{\mathbf{x}_k \in L} t_k \phi_i(\mathbf{x}_k) = (1 - \eta) w_i^{(r)} + \alpha \sum_{\mathbf{x}_k \in L} t_k \phi_i(\mathbf{x}_k)$$

SVM and SGD

In stochastic gradient descent, single items are considered at each iteration. This results in the following update rule, assuming \mathbf{x}_k is the element considered at the current step

$$\begin{aligned} w_i^{(r+1)} &= (1 - \eta)w_i^{(r)} + \alpha\phi_i(\mathbf{x}_k) && \text{if } t_k(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_k) + w_0) < 1 \\ w_i^{(r+1)} &= w_i^{(r)} && \text{otherwise} \end{aligned}$$

Kernel methods motivation

- Often we want to **capture nonlinear patterns** in the data
 - Nonlinear Regression: Input-output relationship may not be linear
 - Nonlinear Classification: Classes may not be separable by a linear boundary
- Linear models (e.g., linear regression, linear SVM) are not just rich enough
- Kernels: Make linear models work in nonlinear settings
 - By mapping data to higher dimensions where it exhibits linear patterns
 - Apply the linear model in the new input space
 - Mapping changing the feature representation
- Note: Such mappings can be expensive to compute in general
 - Kernels give such mappings for (almost) free
 - * In most cases, the mappings need not be even computed
 - * .. using the Kernel Trick!

Kernels: Formally Defined

- Recall: Each kernel k has an associated basis function $\boldsymbol{\phi}$
- $\boldsymbol{\phi}$ takes input $x \in \mathcal{X}$ (input space) and maps it to \mathcal{F} (feature space)
- Kernel $\kappa(\mathbf{x}_1, \mathbf{x}_2)$ takes two inputs and gives their similarity in \mathcal{F} space

$$\begin{aligned} \boldsymbol{\phi} : \mathcal{X} &\mapsto \mathcal{F} \\ \kappa : \mathcal{X} \times \mathcal{X} &\mapsto \mathbb{R} \qquad \kappa(\mathbf{x}_1, \mathbf{x}_2) = \boldsymbol{\phi}(\mathbf{x}_1)^T \boldsymbol{\phi}(\mathbf{x}_2) \end{aligned}$$

- \mathcal{F} needs to be a vector space with a dot product defined on it (Hilbert space)
- Can just any function be used as a kernel function?
 - No. It must satisfy a suitable condition

Verifying a given function is a kernel

A necessary and sufficient condition for a function $\kappa : \mathbb{R}^n \times \mathbb{R}^n \mapsto \mathbb{R}^n$ to be a kernel is that, for all sets $(\mathbf{x}_1, \dots, \mathbf{x}_n)$, the Gram matrix \mathbf{K} such that $k_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$ is semidefinite positive, that is

$$\mathbf{v}^T \mathbf{K} \mathbf{v} \geq 0$$

for all vectors^v.

This is equivalent to the condition that all eigenvalues of \mathbf{K} are non negative.

Constructing kernel functions

Example:

Let $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^2$: $\kappa(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1 \cdot \mathbf{x}_2)^2$ is a valid kernel function?

This can be verified by observing that

$$\begin{aligned}\kappa(\mathbf{x}_1, \mathbf{x}_2) &= (x_{11}x_{21} + x_{12}x_{22})^2 \\ &= x_{11}^2x_{21}^2 + x_{12}^2x_{22}^2 + 2x_{11}x_{12}x_{21}x_{22} \\ &= (x_{11}^2, x_{12}^2, x_{11}x_{12}, x_{11}x_{12}) \cdot (x_{21}^2, x_{22}^2, x_{21}x_{22}, x_{21}x_{22}) \\ &= \boldsymbol{\phi}(\mathbf{x}_1) \cdot \boldsymbol{\phi}(\mathbf{x}_2)\end{aligned}$$

and by defining the base functions as $\mathbf{x} = (x_1^2, x_2^2, x_1x_2, x_1x_2)^T$.

- In general, if $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$ then $\kappa(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1 \cdot \mathbf{x}_2)^2 = \boldsymbol{\phi}(\mathbf{x}_1)^T \boldsymbol{\phi}(\mathbf{x}_2)$, where

$$\mathbf{x} = (x_1^2, \dots, x_d^2, x_1x_2, \dots, x_1x_d, x_2x_1, \dots, x_dx_{d-1})^T$$

- the d -dimensional input space is mapped onto a space with dimension $m = d^2$
- observe that computing $\kappa(\mathbf{x}_1, \mathbf{x}_2)$ requires time $O(d)$, while deriving it from $\boldsymbol{\phi}(\mathbf{x}_1)^T \boldsymbol{\phi}(\mathbf{x}_2)$ requires $O(d^2)$ steps

The function $\kappa(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1 \cdot \mathbf{x}_2 + c)^2$ is a kernel function. In fact,

$$\begin{aligned}\kappa(\mathbf{x}_1, \mathbf{x}_2) &= (\mathbf{x}_1 \cdot \mathbf{x}_2 + c)^2 \\ &= \sum_{i=1}^n \sum_{j=1}^n x_{1i}x_{1j}x_{2i}x_{2j} + \sum_{i=1}^n (\sqrt{2c}x_{1i})(\sqrt{2c}x_{2i}) + c^2 \\ &= \boldsymbol{\phi}(\mathbf{x}_1)^T \boldsymbol{\phi}(\mathbf{x}_2)\end{aligned}$$

for

$$\mathbf{x} = (x_1^2, \dots, x_d^2, x_1x_2, \dots, x_1x_d, x_2x_1, \dots, x_dx_{d-1}, \sqrt{2c}x_1, \dots, \sqrt{2c}x_d, c)^T$$

This implies a mapping from a d -dimensional to a $(d+1)^2$ -dimensional space.

Function $\kappa(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1 \cdot \mathbf{x}_2 + c)^t$ is a kernel function corresponding to a mapping from a d -dimensional space to a space of dimension

$$m = \sum_{i=0}^t d^i = \frac{d^{t+1} - 1}{d - 1}$$

corresponding to all products $x_{i_1}x_{i_2} \dots x_{i_l}$ with $0 \leq l \leq t$.

Observe that, even if the space has dimension $O(d^t)$, evaluating the kernel function requires just time $O(d)$.

Techniques for constructing kernel functions

Given kernel functions $\kappa_1(\mathbf{x}_1, \mathbf{x}_2)$, $\kappa_2(\mathbf{x}_1, \mathbf{x}_2)$, the function $\kappa(\mathbf{x}_1, \mathbf{x}_2)$ is a kernel in all the following cases

- $\kappa(\mathbf{x}_1, \mathbf{x}_2) = e^{\kappa_1(\mathbf{x}_1, \mathbf{x}_2)}$
- $\kappa(\mathbf{x}_1, \mathbf{x}_2) = \kappa_1(\mathbf{x}_1, \mathbf{x}_2) + \kappa_2(\mathbf{x}_1, \mathbf{x}_2)$
- $\kappa(\mathbf{x}_1, \mathbf{x}_2) = \kappa_1(\mathbf{x}_1, \mathbf{x}_2)\kappa_2(\mathbf{x}_1, \mathbf{x}_2)$
- $\kappa(\mathbf{x}_1, \mathbf{x}_2) = c\kappa_1(\mathbf{x}_1, \mathbf{x}_2)$, for any $c > 0$

- $\kappa(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^T \mathbf{A} \mathbf{x}_2$, with \mathbf{A} positive definite
- $\kappa(\mathbf{x}_1, \mathbf{x}_2) = f(\mathbf{x}_1) \kappa_1(\mathbf{x}_1, \mathbf{x}_2) g(\mathbf{x}_2)$, for any $f, g : \mathbb{R}^n \mapsto \mathbb{R}$
- $\kappa(\mathbf{x}_1, \mathbf{x}_2) = p(\kappa_1(\mathbf{x}_1, \mathbf{x}_2))$, for any polynomial $p : \mathbb{R}^q \mapsto \mathbb{R}$ with non-negative coefficients
- $\kappa(\mathbf{x}_1, \mathbf{x}_2) = \kappa_3(\boldsymbol{\phi}(\mathbf{x}_1), \boldsymbol{\phi}(\mathbf{x}_2))$, for any vector $\boldsymbol{\phi}$ of m functions $\phi_i : \mathbb{R}^n \mapsto \mathbb{R}$ and for any kernel function $\kappa_3(\mathbf{x}_1, \mathbf{x}_2)$ in \mathbb{R}^m

$\kappa(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1 \cdot \mathbf{x}_2 + c)^d$ is a kernel function. In fact,

1. $\mathbf{x}_1 \cdot \mathbf{x}_2 = \mathbf{x}_1^T \mathbf{x}_2$ is a kernel function corresponding to the base functions $\boldsymbol{\phi} = (\phi_1, \dots, \phi_n)$, with $\phi_i(\mathbf{x}) = \mathbf{x}$
2. c is a kernel function corresponding to the base functions $\boldsymbol{\phi} = (\phi_1, \dots, \phi_n)$, with $\phi_i(\mathbf{x}) = \frac{\sqrt{c}}{n}$
3. $\mathbf{x}_1 \cdot \mathbf{x}_2 + c$ is a kernel function since it is the sum of two kernel functions
4. $(\mathbf{x}_1 \cdot \mathbf{x}_2 + c)^d$ is a kernel function since it is a polynomial with non negative coefficients (in particular $p(z) = z^d$) of a kernel function

$$\kappa(\mathbf{x}_1, \mathbf{x}_2) = e^{-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|^2}{2\sigma^2}}$$

is a kernel function. In fact,

1. since $\|\mathbf{x}_1 - \mathbf{x}_2\|^2 = \mathbf{x}_1^T \mathbf{x}_1 + \mathbf{x}_2^T \mathbf{x}_2 - 2\mathbf{x}_1^T \mathbf{x}_2$, it results

$$\kappa(\mathbf{x}_1, \mathbf{x}_2) = e^{-\frac{\mathbf{x}_1^T \mathbf{x}_1}{2\sigma^2}} e^{-\frac{\mathbf{x}_2^T \mathbf{x}_2}{2\sigma^2}} e^{\frac{\mathbf{x}_1^T \mathbf{x}_2}{\sigma^2}}$$

2. $\mathbf{x}_1^T \mathbf{x}_2$ is a kernel function (see above)
3. then, $\frac{\mathbf{x}_1^T \mathbf{x}_2}{\sigma^2}$ is a kernel function, being the product of a kernel function with a constant $c = \frac{1}{\sigma^2}$
4. $e^{\frac{\mathbf{x}_1^T \mathbf{x}_2}{\sigma^2}}$ is the exponential of a kernel function, and as a consequence a kernel function itself
5. $e^{-\frac{\mathbf{x}_1^T \mathbf{x}_1}{\sigma^2}} e^{-\frac{\mathbf{x}_2^T \mathbf{x}_2}{2\sigma^2}} e^{\frac{\mathbf{x}_1^T \mathbf{x}_2}{\sigma^2}}$ is a kernel function, being the product of a kernel function with two functions $f(\mathbf{x}_1) = e^{-\frac{\mathbf{x}_1^T \mathbf{x}_1}{\sigma^2}}$ and $g(\mathbf{x}_2) = e^{-\frac{\mathbf{x}_2^T \mathbf{x}_2}{2\sigma^2}}$

Relevant kernel functions

1. Polynomial kernel

$$\kappa(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1 \cdot \mathbf{x}_2 + 1)^d$$

2. Sigmoidal kernel

$$\kappa(\mathbf{x}_1, \mathbf{x}_2) = \tanh(c_1 \mathbf{x}_1 \cdot \mathbf{x}_2 + c_2)$$

3. Gaussian kernel

$$\kappa(\mathbf{x}_1, \mathbf{x}_2) = \exp\left(-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|^2}{2\sigma^2}\right)$$

where $\sigma \in \mathbb{R}$

Observe that a gaussian kernel can be derived also starting from a non linear kernel function $\kappa(\mathbf{x}_1, \mathbf{x}_2)$ instead of $\mathbf{x}_1^T \mathbf{x}_2$.