# Modelli Lineari (Generalizzati) e SVM

## Corso di AA, anno 2017/18, Padova

Fabio Aiolli
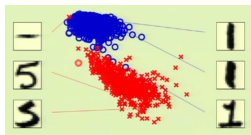
13/15 Novembre 2017

# Outline

- Linear methods for classification and regression
- Non-linear transformations
- Optimal hyperplane and Support Vector Machine (SVM)
- SVM for non linearly-separable data

# Linear Models

- One of the most important type of models in ML
- A linear model is in the form $f_{\mathbf{w},b}(\mathbf{x}) = \sum_{i=1}^{m} w_i x_i + b = \mathbf{w} \cdot \mathbf{x} + b$
- Which can also be written as $f_{\mathbf{w}}(\mathbf{x}) = \sum_{i=0}^{m} w_i x_i = \mathbf{w} \cdot \mathbf{x}$ where $x_0 = 1$ is a an ad-hoc artificial feature (coordinate)
- For classification, the sign is returned, that is
  $h(\mathbf{x}) = \text{sign}(f_{\mathbf{w}}(\mathbf{x})) \in \{-1, +1\}$
- For regression, the original function can be taken, that is
  $h(\mathbf{x}) = f_{\mathbf{w}}(\mathbf{x}) \in \mathbb{R}$

# Real Data



- $16 \times 16$ grey-level images
- Standard representation: raw input $\mathbf{x} = (x_1, \ldots, x_{256})$
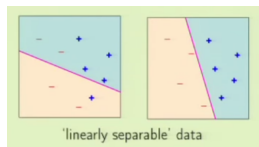- You can also use other representations (e.g. intensity, symmetry)

# The Perceptron Algorithm

$$h_{\mathbf{w}}(\mathbf{x}) = \text{sign}(\sum_{i=0}^{n} w_i x_i) = \text{sign}(\mathbf{w} \cdot \mathbf{x})$$

1. Given a training set

$$\{(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_n, y_n)\}$$

2. Pick a misclassified point $\mathbf{x}_i$ ($\text{sign}(\mathbf{w} \cdot \mathbf{x}_i) \neq y_i$)

3. Update the weight vector $\mathbf{w} \leftarrow \mathbf{w} + y_i \mathbf{x}_i$

4. Repeat from step 2 until all points are correctly classified



'linearly separable' data

If data are linearly separable, then the perceptron algorithm always converges to a valid solution!
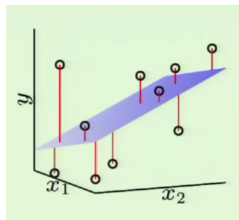
# Linear models for classification

- Perceptron
- Constrained linear problem
- Least squares for classification
- Logistic regression models
- SVM
- ...

# (Multi-variate) Linear Regression

- Example: Given a set of characteristics for a user: age, annual salary, years in residence, years in job, current debits, etc. Predict the credit line, that is the amount of credit that can be granted to a customer.

- Given TRAIN $= \{(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_n, y_n)\}$, in linear regression we look for a hypothesis $h_\mathbf{w}$ (a linear space) which minimizes the mean squared error on the training set, that is $\frac{1}{n} \sum_{i=1}^{n} (h_\mathbf{w}(\mathbf{x}_i) - y_i)^2$

# Solving the Linear Regression problem

$$
\begin{aligned}
E(\mathbf{w}) &= \frac{1}{n}\sum_{i=1}^{n}(h_\mathbf{w}(\mathbf{x}_i) - y_i)^2 \\
&= \frac{1}{n}\sum_{i=1}^{n}(\mathbf{w}\cdot\mathbf{x}_i - y_i)^2 \\
&= \frac{1}{n}||\mathbf{X}\mathbf{w} - \mathbf{y}||^2
\end{aligned}
$$

where
$$
\mathbf{X} = \begin{pmatrix} \ldots \mathbf{x}_1^\top \ldots \\ \ldots \mathbf{x}_2^\top \ldots \\ \vdots \\ \ldots \mathbf{x}_n^\top \ldots \end{pmatrix}, \qquad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}
$$

$$\min_{\mathbf{w}} E(\mathbf{w}) \equiv \frac{1}{n} ||\mathbf{X}\mathbf{w} - \mathbf{y}||^2$$

$$\nabla E(\mathbf{w}) = \frac{2}{n}\mathbf{X}^{\top}(\mathbf{X}\mathbf{w} - \mathbf{y}) = \mathbf{0}$$
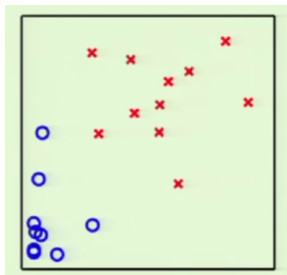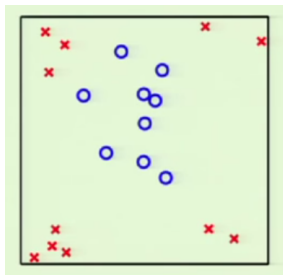
$$\mathbf{X}^{\top}\mathbf{X}\mathbf{w} = \mathbf{X}^{\top}\mathbf{y}$$

$$\mathbf{w} = \mathbf{X}'\mathbf{y}, \text{ where } \mathbf{X}' = (\mathbf{X}^{\top}\mathbf{X})^{-1}\mathbf{X}^{\top}$$
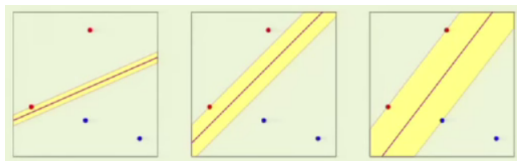
The matrix $\mathbf{X}'$ is the pseudo-inverse of $\mathbf{X}$

$$(x_1, x_2) \xrightarrow{\Phi} (x_1^2, x_2^2)$$



Generalized Linear Models: In general, any non-linear transformation $x \xrightarrow{\Phi} z$ (a.k.a. basis function) can be applied to the data. An hyperplane, hence a linear model, in the transformed space will correspond to a non-linear decision surface in the original space!

# Linear separability



- Consider the hypothesis space of hyperplanes
- Take a set of linearly separable points
- We have different separating hyperplanes fitting data
- Which is the best?

Two questions:

1. Why the widest possible margin (or optimal) hyperplane is better?
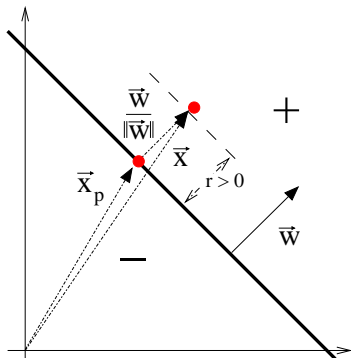2. To which **w**, $b$ this corresponds?

# Margin of a hyperplane

Given the hyperplane $\mathbf{w} \cdot \mathbf{x} + b = 0$, the "distance" of a point $\mathbf{x}$ from the hyperplane can be expressed by the *algebraic* measure $g(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$.

We can write $\mathbf{x} = \mathbf{x}_p + r\frac{\mathbf{w}}{||\mathbf{w}||}$
where:

- $\mathbf{x}_p$ is the normal projection of $\mathbf{x}$ onto the hyperplane
- $r$ is the desired algebraic distance ($r > 0$ if $\mathbf{x}$ is on the positive side of the hyperplane, otherwise $r < 0$)

# Margin

$$g(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b \quad \text{and} \quad \mathbf{x} = \mathbf{x}_p + r\frac{\mathbf{w}}{||\mathbf{w}||},$$

Two facts:

- Note that $g(\mathbf{x}_p) = 0$ (because $\mathbf{x}_p$ in on the optimal hyperplane)
- Since the absolute distance from any nearest positive example is the same as the absolute distance from any nearest negative example, then we can consider hypotheses $\mathbf{w}, b$ such that $g(\mathbf{x}) = 1$ when $\mathbf{x}$ is in the (positive side) margin hyperplane and $g(\mathbf{x}) = -1$ when $\mathbf{x}$ is in the (negative side) margin hyperplane.

Take $\mathbf{x}_k$ in the positive magin hyperplane, then

$$g(\mathbf{x}_k) = \underbrace{\mathbf{w} \cdot \mathbf{x}_p + b}_{=0} + r\frac{\mathbf{w} \cdot \mathbf{w}}{||\mathbf{w}||} = r||\mathbf{w}|| \Rightarrow r = \frac{g(\mathbf{x}_k)}{||\mathbf{w}||} = \frac{1}{||\mathbf{w}||}$$

and hence the margin will be $\rho = \frac{2}{||\mathbf{w}||}$.

# Support Vector Machines: basic idea

- Can we apply the Structural Risk Minimization principle to hyperplanes?
- We have seen that a hyperplane in $\mathbb{R}^m$ has $VC = m + 1$.
- In fact, if we add further constraints on the hyperplanes we can do better!

# Margin: Link with SRM

**Theorem** *Let R denote the diameter of the smallest ball containing all the input points. The set of optimal hyperplanes described by the equation $\mathbf{w} \cdot \mathbf{x} + b = 0$ has a VC-dimension $VC_{opt}$ bounded from above as*

$$VC_{opt} \leq \min\{\lceil \frac{R^2}{\rho^2} \rceil, m\} + 1$$

*where $\rho = \frac{2}{||\mathbf{w}||}$ and m is the dimensionality of the input space.*

Thus, if we consider the hypothesis spaces

$$\mathcal{H}_k = \{\mathbf{w} \cdot \mathbf{x} + b \mid ||\mathbf{w}||^2 \leq c_k\} \text{ where } c_1 < c_2 < c_3 < \dots$$

and linearly separable data, then the empirical error of $\mathcal{H}_k$ is 0 for each $k$ and the bound on the true risk can be minimized by *maximizing the margin of separation* (i.e. minimizing the weight norm).

# Separable Case: Quadratic optimization

If we have $n$ linearly separable examples $\{(\mathbf{x}_i, y_i)\}_1^n$, it is possible to find the optimal hyperplane solving the following constrained quadratic optimization problem:

$$\min_{\mathbf{w},b} \frac{1}{2}||\mathbf{w}||^2$$
$$\text{subject to: } \forall i \in \{1,\ldots,n\} : y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$$

- This is a (convex) constrained quadratic problem. Thus it guarantees a unique solution!
- Many QP algorithms exist with polynomial complexity to find the solution of this quadratic problem

# Solving the optimization problem

- The problem above, called primal problem, can be solved more easily using the dual formulation.
- In the dual problem Lagrange multipliers $\alpha_i \geq 0$ are associated with every constraint in the primal problem (one for each example).

- The dual formulation is:

$$\max_\alpha \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$
$$\text{subject to: } \forall i \in \{1, \ldots, n\} : \alpha_i \geq 0 \text{ e } \sum_{i=1}^n y_i \alpha_i = 0.$$

- At the solution most of the $\alpha_i$'s are zeros. Examples associated with non zero multipliers are called support vectors.

## SVM solution

The primal solution turns out to be:

$$
\begin{aligned}
\mathbf{w} &= \sum_{i=1}^{n} y_i \alpha_i \mathbf{x}_i \\
b &= y_k - \mathbf{w} \cdot \mathbf{x}_k \quad \text{for any } \mathbf{x}_k \text{ such that } \alpha_k > 0
\end{aligned}
$$

and hence:

$$
h(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b) = \text{sign}(\sum_{i=1}^{n} y_i \alpha_i (\mathbf{x}_i \cdot \mathbf{x}) + b)
$$

that is, the decision function only depends on dot products between the point and other points in the training set (the support vectors).

## SVM for the Non-separable case

If the examples are NOT linearly separable we have to allow that some constraints are violated. This can be done by

- introducing *slack* variables $\xi_i \geq 0, \ i = 1, \dots, n$ , one for each constraint:

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$$

- modifying the cost function so to penalize slack variables which are not 0:

$$\frac{1}{2}||\mathbf{w}||^2 + C \sum_{i=1}^{n} \xi_i$$

where $C$ (regularization parameter) is a positive constant controlling the tradeoff between the complexity of the hypothesis space and the number of margin errors.
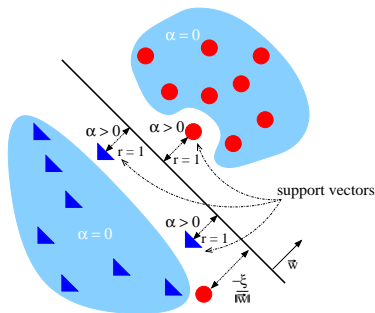
# SVM for the Non-separable Case

The dual of this new formulation is very similar to the previous one:

$$\max_\alpha \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

subject to: $\forall i \in \{1, \ldots, n\} : 0 \leq \alpha_i \leq C$ e $\sum_{i=1}^n y_i \alpha_i = 0$.

The main difference is due to the fact that the dual variables are upper bounded by $C$. The value for $b$ is obtained similarly to the separable case (with some minor differences...).

- The parameter C can be seen as a way to control overfitting
- As C becomes larger it is unattractive to not respect the data at the cost of reducing the geometric margin
- When C is small, larger margin is possible at the cost of increasing errors in training data
- Interestingly, the SVM solution is in the same form as in the hard margin case!

Nevertheless, this formulation is not always satisfactory because of the limited separation capability of a hyperplane.
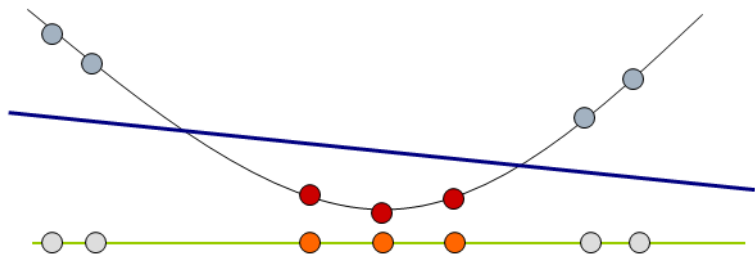
How can we separate these data?

# Another approach

Projecting them into a higher dimensional space

$$\mathbf{x} \rightarrow \varphi(\mathbf{x})$$

When the examples are not linearly separable, an alternative approach based on the following two steps can be used

1. the input vectors (input space) are projected into a larger space (feature space);

2. the optimal hyperplane in feature space is computed (using the formulation with slack variables)

# Nonseparable case

Step 1 is justified by Cover's theorem on separability, which states that non-linearly separable patterns may be transformed into a new feature space where the patterns are linearly separable with high probability, provided that the transformation is nonlinear, and that the dimensionality of the feature space is high enough.

Step 2 is trivially justified by the fact that the optimal hyperplane (in the feature space) minimizes the VC-dimension.

## Nonseparable case

We can assume that any of the new feature space coordinate is generated by a nonlinear function $\varphi_j(\cdot)$. Thus, we can consider $M$ functions $\varphi_j(\mathbf{x})$ with $j = 1, \ldots, M$. A generic vector $\mathbf{x}$ is thus mapped into the $M$-dimensional vector

$$\varphi(\mathbf{x}) = [\varphi_1(\mathbf{x}), \ldots, \varphi_M(\mathbf{x})]$$

Step 2 asks to find the optimal hyperplane into the $M$-dimensional feature space. A hyperplane into the feature space is defined as

$$\sum_{j=1}^{M} w_j \varphi_j(\mathbf{x}) + b = 0$$

or, equivalently

$$\sum_{j=0}^{M} w_j \varphi_j(\mathbf{x}) = \mathbf{w} \cdot \varphi(\mathbf{x}) = 0$$

where we added a coordinate $\varphi_0(\mathbf{x}) = 1$ and $w_0 = b$.

## Nonseparable case

By

$$\mathbf{w} = \sum_{k=1}^{n} y_k \alpha_k \vec{\varphi}(\mathbf{x}_k)$$

the equation defining the hyperplane becomes

$$\sum_{k=1}^{n} y_k \alpha_k \varphi(\mathbf{x}_k) \cdot \varphi(\mathbf{x}) = 0$$

where $\varphi(\mathbf{x}_k) \cdot \varphi(\mathbf{x})$ represents the dot product (in feature space) between vectors induced by the $k$-th training instance and the input $\mathbf{x}$.

## Kernel functions

Now, what we need is a function $K(\cdot, \cdot)$ (called *kernel function*) such that

$$K(\mathbf{x}_k, \mathbf{x}) = \varphi(\mathbf{x}_k) \cdot \varphi(\mathbf{x}) = \sum_{j=0}^{M} \varphi_j(\mathbf{x}_k) \varphi_j(\mathbf{x}) = K(\mathbf{x}, \mathbf{x}_k) \text{ (symmetric function)}$$

If we get such a function, we could compute the decision function in feature space WITHOUT explicitly representing the vectors into the feature space:

$$\sum_{k=1}^{n} y_k \alpha_k K(\mathbf{x}_k, \mathbf{x})$$

Functions with this property do actually exist, if some conditions are satisfied... namely, Mercer's conditions.

# Kernel functions

In general, a kernel function satisfying Mercer's conditions represents a dot-product between vectors generated by some (non-linear) transformation.

Note that we do not need to know such a transformation!!

Examples of popular kernel functions:

- linear kernel, $K(\mathbf{x}, \mathbf{x}') = \mathbf{x} \cdot \mathbf{x}'$
- polynomial kernel of degree $p$, $K(\mathbf{x}, \mathbf{x}') = (\mathbf{x} \cdot \mathbf{x}' + u)^p$
- radial-basis function (RBF) kernel, $K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$

# Formulation with Kernel

The introduction of a kernel does not modify the problem formulation:

$$\max_\alpha \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j)$$
subject to: $\forall i \in \{1, \ldots, n\} : 0 \leq \alpha_i \leq C$ and $\sum_{i=1}^n y_i \alpha_i = 0$.

where the needed kernel values are computed over all pairs of vectors ($K(\mathbf{x}_i, \mathbf{x}_j)$, with $i, j = 1, \ldots, n$) and arranged into a matrix $\boldsymbol{K} \in \mathbb{R}^{n \times n}$ (symmetric and positive definite) known as kernel matrix or gram matrix.

# Formulation with Kernel

E.g., if we use a polynomial kernel with degree $p = 3$ we obtain $K_{i,j} = (\mathbf{x}_i \cdot \mathbf{x}_j + 1)^3$ and a new instance $\mathbf{x}$ is classified by the following discriminant function

$$h(\mathbf{x}) = \text{sign}(\sum_{\mathbf{x}_k \in SV} y_k \alpha_k^* K(\mathbf{x}_k, \mathbf{x})) = \text{sign}(\sum_{\mathbf{x}_k \in SV} y_k \alpha_k^* (\mathbf{x}_k \cdot \mathbf{x} + 1)^3)$$
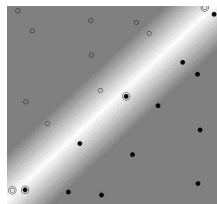
where $SV$ is the set of support vectors and $\alpha_k^*$ are the optimal values for the support vectors (the remaining dual variable are $0 \Rightarrow$ corresponding vectors do not contribute to the sum).
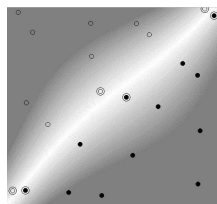
# Formulation with Kernel

Using this approach, we can use a nonlinear transformation $\varphi(\cdot)$ IMPLICITLY, in fact what we need is not the explicit representation of vectors in feature space, but their dot product into the feature space. This can be directly computed in the input space via the kernel function.
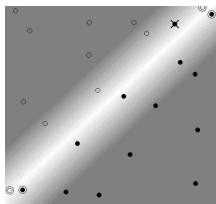
Examples of decision surfaces generated IN THE INPUT SPACE with or without kernel (polynomial with degree 3) both for the separable and nonseparable case:
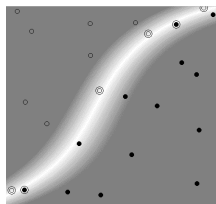


separable          degree 3 poly          not separable          degree 3 poly

## Example: polynomial kernel

Given two vetors **x** and **z** and the following mapping $\varphi()$

$$\mathbf{x} = (x_1, x_2); \varphi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)$$
$$\mathbf{z} = (z_1, z_2); \varphi(\mathbf{z}) = (z_1^2, z_2^2, \sqrt{2}z_1z_2)$$

A dot product between $\varphi(\mathbf{x})$ and $\varphi(\mathbf{z})$ corresponds to evaluate the function $K_2(\mathbf{x}, \mathbf{z}) = \langle x, z \rangle^2$
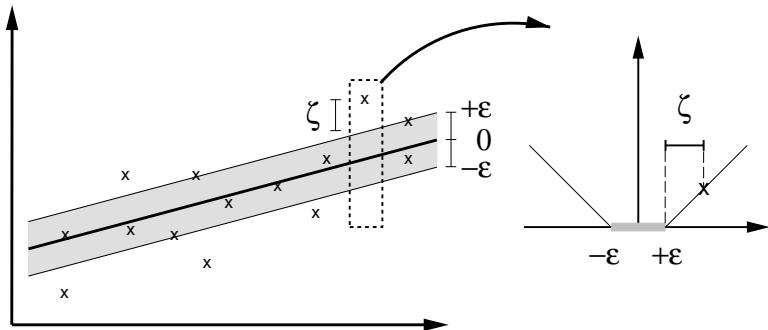
$$\langle \varphi(\mathbf{x}), \varphi(\mathbf{z}) \rangle = \langle (x_1^2, x_2^2, \sqrt{2}x_1x_2), (z_1^2, z_2^2, \sqrt{2}z_1z_2) \rangle =$$
$$= x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 z_1 x_2 z_2 = \langle x, z \rangle^2 = (x_1z_1 + x_2z_2)^2 = K_2(\mathbf{x}, \mathbf{z})$$

$K_2()$ is faster to evaluate than $\langle \varphi(\mathbf{x}), \varphi(\mathbf{z}) \rangle$!

# Regression with SVM: Basic idea

When considering a regression problem, the idea is to define an $\epsilon$-tube: predictions which differ from the desired value for more that $\epsilon$ in absolute error are linearly penalized, otherwise they are not considered errors.

# Regression with SVM

The $\epsilon$-tube idea leads to the following formulation

$$\min_{\mathbf{w},b,\xi,\xi^*} \frac{1}{2}||\mathbf{w}||^2 + C \sum_{i=1}^{n}(\xi_i + \xi_i^*)$$
subject to:
$$\forall i \in \{1, \ldots, n\}$$
$$y_i - \mathbf{w} \cdot \mathbf{x}_i - b \leq \epsilon + \xi_i$$
$$\mathbf{w} \cdot \mathbf{x}_i + b - y_i \leq \epsilon + \xi_i^*$$
$$\xi_i, \xi_i^* \geq 0$$

which has the following dual formulation...

$$\max_{\alpha,\alpha^*} -\epsilon \sum_{i=1}^{n}(\alpha_i + \alpha_i^*) + \sum_{i=1}^{n} y_i(\alpha_i - \alpha_i^*)+$$
$$-\frac{1}{2}\sum_{i,j=1}^{n}(\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*)K(\mathbf{x}_i, \mathbf{x}_j)$$
subject to:
$$\sum_{i=1}^{n}(\alpha_i - \alpha_i^*) = 0$$
$$\alpha_i, \alpha_i^* \in [0, C]$$

## Recap

- Motivare dal punto di vista teorico il metodo SVM (Support Vector Machine)
- Discutere i metodi per il trattamento di dati non linearmente separabili con le SVM
- Nel contesto di metodi kernel (p.e. SVM), descrivere il concetto di kernel e il suo legame stretto con la rappresentazione dei dati
- Relazioni tra Perceptron e SVM: è possibile vedere il Perceptron come metodo kernel?
- Relazioni tra Reti Neurali e SVM: dov'è il "kernel" nelle Reti Neurali? Quali sono le differenze in termini di funzione obiettivo tra i due metodi?