

Learning with Kernels

Proseminar Data Mining

Matthieu Bulté

Fakultät für Mathematik

Technische Universität München

Email: matthieu.bulte@tum.de

Abstract—Text der Kurzfassung

Index Terms—support vector machines, machine learning, statistics

I. INTRODUCTION

todo

II. FIRST STEPS: HARD MARGIN CLASSIFIER

The hard margin classifier learns the parameters $\mathbf{w} \in \mathbb{R}^n$ and $b \in \mathbb{R}$ of a hyperplane $\mathcal{H} \subset \mathbb{R}^n$, separating two classes of observations \mathbf{x} of dimension n . To train it, the algorithm takes a set of m observations \mathbf{x}_i together with a vector of labels $y_i = \pm 1$ for each of the observations:

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)$$

The learned parameters are then used to define a decision function f , assigning to points of \mathbb{R}^n a label ± 1 corresponding to the side of the hyperplane on which the points stands:

$$f(\mathbf{x}) = \text{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle + b) \quad (1)$$

A. Margin maximization

One particularity about the learning algorithm used for Support Vector Machines is the loss function the algorithm tries to optimize. Other learning algorithms typically chose a loss function based on the empirical risk, defined as following for a function f and a training set (\mathbf{x}, \mathbf{y}) of size m :

$$R_{\text{emp}}[f] = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} |f(\mathbf{x}_i) - y_i| \quad (2)$$

The margin maximization algorithm instead, searches within the space of hyperplanes properly classifying the training set, for the hyperplane with the largest distance from the training points to the hyperplane. This can be translated in the following constrained optimization problem:

NOTE: before jumping into equations, maybe explain the reason why we want to maximize the margin. argument using the noise around sample.

$$\begin{aligned} & \underset{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}}{\text{maximize}} & M &:= \frac{1}{\|\mathbf{w}\|} \min_i y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \\ & \text{subject to} & & y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq M \text{ for all } i = 1 \dots m \end{aligned} \quad (3)$$

Assuming that a solution to this optimization problem exists, it still has the issue that although the solution hyperplane

is unique, there is an infinity of parameter vectors $\|\mathbf{w}\|$ resulting in the solution hyperplane, only differing in their length. Uniqueness of \mathbf{w} can be ensured by adding to its length the following constraint:

$$M \|\mathbf{w}\| = 1$$

We can thus reformulate the optimization problem into the following quadratic optimization problem which is known to have a numerical solution:

$$\begin{aligned} & \underset{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}}{\text{minimize}} & & \|\mathbf{w}\| \\ & \text{subject to} & & y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 \text{ for all } i = 1 \dots m \end{aligned} \quad (4)$$

B. Support vectors

Although the previous formulation of the problem helped us to better understand the ideas of support vector machines, the resulting quadratic problem will become impractical when later introducing kernels methods. We will introduce in this section an equivalent formulation of the optimization problem (4) **can latex handle references instead of hard coding?** that will not only solve computational issues, but will also give us more insights about the resulting solution.

In order to get to these benefits, we introduce the Lagrangian L together with Lagrange multipliers $\alpha_i \geq 0$:

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^m \alpha_i [y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1] \quad (5)$$

This is called the dual formulation of our original optimization problem, called the primal problem, in which \mathbf{w} and b are called primal variables, and α the dual variable. It can be shown **ref?** that the primal optimization problem is equivalent to finding a saddle point of L , minimizing with respect to \mathbf{w} and b , while maximizing it with respect to α .

Because the solution of this dual problem is a saddle point, thus an extremum with respect to all variables, following equalities must hold:

$$\frac{\partial}{\partial b} L(\mathbf{w}, b, \alpha) = 0, \quad \frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}, b, \alpha) = 0 \quad (6)$$

leading to

$$\sum_{i=1}^m \alpha_i y_i = 0 \quad (7)$$

and

$$\begin{aligned} \mathbf{w} - \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i &= 0 \\ \Rightarrow \mathbf{w} &= \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \end{aligned} \quad (8)$$

Further, by replacing (7) and (8) into the Lagrangian, we obtain an optimization problem free of primal variables:

$$\begin{aligned} \underset{\alpha \in \mathbb{R}^n}{\text{maximize}} \quad & W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\ \text{subject to} \quad & \alpha_i \geq 0 \text{ for all } i = 1 \dots m \\ & \sum_{i=1}^m \alpha_i y_i = 0 \end{aligned} \quad (9)$$

Finally, the Karush-Kuhn-Tucker theorem states that the solution α satisfies the following equality for all $i = 1 \dots m$

$$\alpha_i [y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1] = 0 \quad (10)$$

The importance of this equality is twofold. First, it allows us to compute b from an observation i for which $\alpha_i \neq 0$. Second and most importantly, this equality implies that every observations not lying on the margin must have a vanishing Lagrangian coefficient. These observations are called *support vectors*.

III. NON LINEARLY SEPARABLE

A. Noisy training data

B. Non linear target boundary

IV. KERNEL METHODS

A. The kernel trick

B. Some useful kernels

V. ZUSAMMENFASSUNG UND AUSBLICK

blabla