

Algorithmes de Plans Coupants, Schémas de Compressions et Apprentissage Actif

Ugo Louche et Liva Ralaivola *

Qarma, LIF - CNRS, Aix-Marseille Université.

April 30, 2015

Abstract

Les méthodes de plans coupants (*Cutting Planes methods*) sont des algorithmes d’optimisations bien connus. Nous montrerons que ces méthodes peuvent être directement utilisées pour l’apprentissage automatique, et non simplement comme un outil d’optimisation. Plus particulièrement, elles permettent l’apprentissage de classifieurs parcimonieux disposants de propriétés de compressions avantageuses. De plus, nous montrerons que de nouveaux algorithmes d’apprentissages actifs peuvent être facilement dérivés à partir de ces méthodes. Plus précisément, nous décrirons un processus générique permettant d’adapter une large gamme d’algorithmes d’apprentissages passifs pour l’apprentissage actif. La pertinence de notre approche sera appuyée par des simulations numériques dans le cadre de problèmes d’apprentissages actifs comme pas-

Mots-clef: Apprentissage Actif; Schéma de Compression; Classification; Méthode de Plans Coupants; Centre de Gravité

1 Introduction

We show that localization methods based on cutting planes provide a natural framework to derive machine learning algorithms for classification, both in the supervised learning framework and the active learning framework. Our claim is that cutting plane algorithms, beyond their optimization purposes, embed features that are beneficial for generalization purposes. In particular a) under mild conditions, they may provide compression scheme with a compression rate that is

directly related to their aim at rapidly finding a solution of the localization problem and b) the pivotal step of such algorithms, namely, the querying step, may be slightly twisted so as to be active-learning friendly.

In the present paper, we show that existing learning algorithms might be revisited from the cutting planes point of view. Not only might the active learning SVM procedure of Tong and Koller [TK02] be reinterpreted as an algorithm falling under the framework we describe but so are the Bayes Point Machines [RHC01], for which we will propose an active learning version of it.

The problems we are interested in are linear classification problems. Given a training sample $D \doteq \{(x_n, y_n)\}_{n \in [N]}$, with $x_n \in \mathcal{X} \doteq \mathbb{R}^d$, $y_n \in \mathcal{Y} \doteq \{-1, +1\}$, and $[N] \doteq \{1, \dots, N\}$, we are looking for a classification vector $w \in \mathcal{X}$ that is an element of the version space

$$\mathcal{W}_0(D) \doteq \{w \in \mathcal{X} : y_n \langle w, x_n \rangle \geq 0, n \in [N]\}, \quad (1)$$

of D , i.e. the set of vectors w from \mathcal{X} such that the corresponding linear predictors

$$f_w(x) \doteq \text{sign}(\langle w, x \rangle) \quad (2)$$

make no mistake on the training set D . In order to render the exposition clearer, we make the assumption that the training data are linearly separable so that $\mathcal{W}_0(D)$ is not empty. The case where $\mathcal{W}_0(D) = \emptyset$ can be tackled with usual machine learning techniques — e.g. the “ λ -trick” and/or kernels [FS99] [RHC01].

Also, for the sake of brevity, we may use \mathcal{W}_0 instead of $\mathcal{W}_0(D)$ and thus drop the explicit dependence on D .

With the relevant notation at hand, the problem we are interested in may be stated as:

$$\text{find } w \in \mathcal{W}_0, \quad (3)$$

*Email: firstname.lastname@lif.univ-mrs.fr

which might be simply rewritten as the problem of solving a set of linear inequalities

$$\text{find } w \text{ s.t. } \begin{cases} w \in \mathcal{X} \\ y_n \langle w, x_n \rangle \geq 0, n \in [N]. \end{cases} \quad (4)$$

There is a variety of methods in the optimization literature from as back as the 50’s that are available to solve such problems. Among them, we may mention (*over*-)*relaxation* based methods [Gof80, MS54], simplex-based algorithms and, of course, the Perceptron algorithm and its numerous variants [Blo62, Nov62, Ros58]. Localization methods based on *cutting planes*, or, in short, cutting planes algorithms, are well-studied algorithms, well-known to be very efficient to solve such problems. We will show that, when used to solve (4), i) they naturally provide compression scheme algorithms [FW95], and thus, *learning* algorithms that embed features designed to ensure good generalization properties and ii) they also set the ground for the development of new active learning algorithms.

1.1 Related Works

Cutting-plane methods provide a family of optimization procedures that have received some interest from the machine learning community [JFY09, FS09, TVSL10]. However, they have mainly been considered as optimization methods to solve problems such as those posed by support vector machines or, more generally, regularized risk functionals. The more profound connection of these methods with *learning* algorithms, that is, procedures that are designed in a way to ensure generalization ability to the predictor they build (e.g. the Perceptron algorithm) has less been studied; this is one of the peculiarities of the present paper to discuss this feature—to some extent, the work of [TVSL10], which pinpoints how statistical regularization is beneficial for the stabilization of cutting-plane methods, skims over this connection. Within the vast literature of active learning (see, e.g. [Set12]), we may single out a few contributions our work is closely related to; they share the common feature of focusing on/exploiting the geometry of the version space. The query strategies proposed by [GK10] and [GSSS13] are based on multiple estimations of the volume of the (potential) version space, which, when added together might be computationally expensive. In comparison, in the active learning strategy we derive from the general cutting-plane approach, we compute our queries from an approximated center of gravity of the version space, which is computationally equivalent to a single volume estimation. The work of [BBZ07], who propose a margin-

based query strategy provide theoretical justifications of such strategies and gives insights on the foundations the work of [TK02] hinges on. Our contribution is to show how the cutting planes literature and its accompanying worst-case convergence analyzes may give rise to theoretically supported query strategies that do not have to hinge on margin-based arguments. To some extent, our work has connections with *uncertainty-based active learning* (see, e.g. [LG94]) which advocates to query the points whose class is the most uncertain; our approach may be re-interpreted as a theoretically motivated uncertainty measure based on the volume reduction of the version space.

1.2 Outline

The paper is structured as follows. Section 2 provides some background to cutting planes methods and their possible application to learning. Section 3 further explores the connections between cutting planes and learning algorithms and then provides a way to turn cutting planes methods into an active learning algorithms. Section 4 reports empirical results for algorithms derived from our argumentation on the relevance of cutting plane methods to machine learning.

2 Background

In this section, we first recall the general form of a cutting plane algorithm to solve a localization problem. We then specialize this algorithm to the case where the convex space into which we want to find a point is the version space associated to training set D . Finally, in order for the reader to get a taste on how cutting planes algorithms give rise to *learning algorithms*, i.e. algorithms that embed features, namely, they define *compression schemes* with targeted small compression size, that are beneficial for generalization.

2.1 Vanilla Localization Algorithm with Cutting Planes

In order to solve a problem like

$$\text{find } w \in \mathcal{C},$$

for \mathcal{C} some closed convex set, a localization algorithm based on cutting planes works as follows (see also the synthetic depiction in Algorithm 1) [Kel60]. The algorithm maintains and iteratively refines (i.e. reduces) a closed convex set \mathcal{C}^t that is known to contain \mathcal{C} . From \mathcal{C}^t a *query point* is computed —there are several ways to

compute such query points; we will mention some when specializing localization methods to the specific problem of finding a point in the version space later on—which leads to two possible options: either a) w^t is in \mathcal{C} and the tackled problem is solved or b) $w^t \notin \mathcal{C}$. In the latter case, a so-called *cutting plane oracle* is queried with w^t upon which it returns the parameters (a_t, b_t) of the hyperplane $\{z : \langle a_t, z \rangle = b_t\}$ such that this hyperplane separates w_t from \mathcal{C} , i.e., $\forall w \in \mathcal{C}, \langle a_t, w \rangle > b_t$ and $\langle a_t, w_t \rangle \leq b_t$. The hyperplane is used to reduce \mathcal{C}^t into $\mathcal{C}^t \cap \{w : \langle a_t, w \rangle > b_t\}$ (which still contains \mathcal{C}). For the specific problem (4) of finding a point in the version space, the cutting planes rendered by the oracle will be such that $b_t = 0$.

Algorithm 1 Classical Cutting Plane Algorithm for the localization of $w \in \mathcal{C}$.

Ensure: $w \in \mathcal{C}$

- 1: compute \mathcal{C}^0 , such that $\mathcal{C}^0 \supset \mathcal{C}$ and \mathcal{C}^0 is convex and closed.
- 2: $t \leftarrow 0$
- 3: **repeat**
- 4: Compute *query point* w^t in \mathcal{C}^t
- 5: Ask the *cutting plane oracle* whether $w^t \in \mathcal{C}$
- 6: **if** $w^t \notin \mathcal{C}$ **then**
- 7: Receive a cutting plane (a_t, b_t)
- 8: $\mathcal{C}^{t+1} \leftarrow \mathcal{C}^t \cap \{x : \langle a_t, x \rangle > b_t\}$
- 9: $t \leftarrow t + 1$
- 10: **end if**
- 11: **until** $w^t \in \mathcal{C}$
- 12: **return** w^t

Algorithm 2 The Cutting Plane approach instantiated to the problem of finding a point from the version space of D .

Ensure: w solution of Problem (7)

- 1: $\mathcal{C}^0 \leftarrow \mathcal{B}$
- 2: $t \leftarrow 0$
- 3: **repeat**
- 4: $w^t \leftarrow \text{QUERY}(\mathcal{C}^t) \triangleright$ Compute *query point* w^t in \mathcal{C}^t
- 5: **if** $w^t \notin \mathcal{W}$ **then**
- 6: $n_t \leftarrow \text{PICK}(\mathcal{C}^t, w^t) \triangleright$ pick a cutting plane index
- 7: $\mathcal{C}^{t+1} \leftarrow \mathcal{C}^t \cap \{z : y_{n_t} \langle z, x_{n_t} \rangle > 0\}$
- 8: $t \leftarrow t + 1$
- 9: **end if**
- 10: **until** $w^t \in \mathcal{W}$
- 11: **return** w^t

2.2 Cutting Planes to Localize a Point in the Version Space

Note that problem (4) is scale-insensitive: if $w \in \mathcal{W}_0$, then $\lambda w \in \mathcal{W}_0$ as well for any $\lambda > 0$. In order to get rid of this degree of freedom *and* to make the use of cutting planes algorithms possible (they require the sets \mathcal{C}^t to be bounded), we will restrict ourselves to finding a solution vector w^* both in \mathcal{W}_0 and in the unit ball

$$\mathcal{B} \doteq \{w \in \mathcal{X} : \|w\| \leq 1\}. \quad (5)$$

In other words, we will be looking for w^* in the constrained version space

$$\mathcal{W} \doteq \mathcal{W}_0 \cap \mathcal{B}, \quad (6)$$

and the problem we face is therefore:

$$\text{find } w \text{ such that } \begin{cases} w \in \mathcal{B} \\ y_n \langle w, x_n \rangle \geq 0, n \in [N] \end{cases} \quad (7)$$

In the case of Problem (7), the localization algorithm described earlier translates into the one given in Algorithm 2. The following changes might be observed when comparing with Algorithm 1: \mathcal{C}^0 is now initialized to \mathcal{B} , the unit ball, and the cutting planes are picked among the hyperplanes —i.e. the points of D — defining the version space.

2.3 Query Point Generation

In both Algorithm 1 and Algorithm 2, the strategy to compute a query point is left unspecified. There actually exist many ways to compute such query points, but they all aim at a query point which calls for a cutting plane that will divide the current enclosing convex set \mathcal{C}^t in the most stringent way. It turns out that such guarantee might be expected when the query point is as close as possible to the ‘center’ of \mathcal{C}^t , so that the volume of \mathcal{C}^t is reduced with a positive factor —just as in the well-known bisection method, where the factor is 1/2. The center of \mathcal{C}^t is not defined in a unique way, but for the most popular query methods, it may refer to: a) the center of gravity of \mathcal{C}^t , b) the center of the largest ball inscribed in \mathcal{C}^t , which is called the *Chebyshev center* or c) the analytic center, which we will not discuss further (the interested reader may refer to [Nes95] for further details). We may mention three things regarding the center of gravity: i) it is NP-hard¹ to exactly compute the center of gravity of a convex set in an arbitrary n -dimensional space even

¹To be precise, it is actually #P-hard.

though some practical approximation algorithms exist; ii) it is the query point that comes with the best guarantees in terms of convergence speed of the cutting plane method [DSP10]; iii) the center of gravity of a polytope is precisely the point that is looked for in the case of the theoretically founded Bayes Point Machines of [RHC01].

3 Results

This section is devoted to some algorithmic results that can be obtained when analyzing the behavior of cutting-plane methods for the localization of a point in the version space.

3.1 Cutting Planes Provide Sample Compression Schemes

Let $\mathcal{D} \doteq \bigcup_{n=1}^{\infty} (\mathcal{X} \times \mathcal{Y})^n$ be the set of all finite training samples made of pairs from $\mathcal{X} \times \mathcal{Y}$. In short, sample compression schemes [FW95] are learning algorithms $\mathcal{A} : \mathcal{D} \rightarrow \mathcal{Y}^{\mathcal{X}}$ that are associated with a compression function $\mathcal{S} : \mathcal{D} \rightarrow \mathcal{D}$ so that, given any training sample D , we have $\mathcal{A}(D) = \mathcal{A}(\mathcal{S}(D))$. Sample compression schemes are especially interesting when the size $|\mathcal{S}(D)|$ of the compression set $\mathcal{S}(D)$ is small. Indeed, generalization guarantees that come with these procedures say that the generalization error of $f_D \doteq \mathcal{A}(D)$ is, with high probability (over the random draw of training set D according to an unknown and fix distribution) bounded from above by something like

$$\frac{1}{N - |\mathcal{S}(D)|} \sum_{n=1}^N \mathbb{I}[f_D(x_n) \neq y_n] + \mathcal{O}\left(\sqrt{\frac{1}{N - |\mathcal{S}(D)|}}\right) \quad (8)$$

(see [FW95, GHST05] for a precise statement of the bound). Among the most well-known learning compression schemes, we find the Perceptron and the Support Vector Machines.

We claim that Algorithm 2, which finds a point in the version space using cutting planes, may be a compression scheme.

Proposition 1. *If QUERY(\mathcal{C}^t) (line 4, Algorithm 2) and PICK(\mathcal{C}^t, w^t) (line 6) are both deterministic then Algorithm 2 is a sample compression scheme.*

Proof. If the compression set is made of the training examples that define the cutting planes, this result is a direct consequence of the structure of Algorithm 2. A proof by induction that essentially hinges on the fact that, at each iteration t , the next query point is deterministically computed from \mathcal{C}^t (only) gives the result. \square

Algorithm 3 Top : A Perceptron-based localization algorithm for the case of problem (7). Bottom : The slightly modified perceptron algorithm for compression scheme.

Ensure: Problem (7)

```

1:  $\mathcal{C}^0 \leftarrow \mathcal{B}$ 
2:  $t \leftarrow 1, w^0 \leftarrow 0, \tilde{w}^0 \leftarrow 0$ 
3: repeat
4:    $\tilde{w}^t \leftarrow \text{PERCEPTRON}(\tilde{w}^{t-1}, x_{n_0}, \dots, x_{n_t})$ 
5:    $w^t \leftarrow \tilde{w}^t / \|\tilde{w}^t\|_2$ 
6:   if  $w^t \notin \mathcal{W}$  then
7:     Pick a cutting plane index  $n_t$ 
8:      $\mathcal{C}^{t+1} \leftarrow \mathcal{C}^t \cap \{z : y_{n_t} \langle z, x_{n_t} \rangle \geq 0\}$ 
9:      $t \leftarrow t + 1$ 
10:  end if
11: until  $w^t \in \mathcal{W}$ 
12: return  $w^t$ 
13:
14: function PERCEPTRON( $w^{start}, x_{n_0}, \dots, x_{n_N}$ )
15:    $t \leftarrow 0$ 
16:    $w^0 \leftarrow w^{start}$ 
17:   while  $\exists n_i : \langle w^t, x_{n_i} \rangle < 0$  do
18:      $w^{t+1} \leftarrow w^t + x_{n_i}$ 
19:      $t \leftarrow t + 1$ 
20:   end while
21:   return  $w^t$ 
22: end function

```

A few observations can be made. First, the learning algorithm obtained with the assumptions of Proposition 1 is a *process sample compression scheme*, that is, even if we interrupt the learning before convergence has occurred, running the algorithm on the partial compression scheme obtained so far gives exactly the same predictor. Second, it is obviously an aim to have fast convergence of the localization procedure, where fast convergence means few iterations of the cutting-plane procedure. This directly translates into the idea of finding a point in the version space that is expressed as a combination as few vectors as possible, which, by (8), is very beneficial for generalization purposes. Later, we will see that there are settings for cutting-plane methods that come with guarantees on the number of iterations, and therefore on $|\mathcal{S}(D)|$, to reach convergence.

3.2 Perceptron-based Localization Algorithm

One of the simplest ways to compute a query point w^t for Algorithm 2 is to run Rosenblatt's Perceptron algorithm [Ros58] at each step and query the normalized

solution $w^t = \tilde{w}^t / \|\tilde{w}^t\|_2$. Intuitively, we may expect \tilde{w}^{t+1} to be ‘close’ to \tilde{w}^t because \mathcal{C}^{t+1} is essentially the intersection of \mathcal{C}^t with a cutting plane and much of the geometry of \mathcal{C}^t might be preserved. According to this intuition, \tilde{w}^t should be a good starting point for the Perceptron algorithm to be run and to have it output \tilde{w}_{t+1} . Algorithm 3 implements that idea, and reuses the last query point as an initialization vector for the Perceptron to compute the next query point. Additionally, note that for Algorithm 3 to match Proposition 1 a little technicality is needed: we require that datapoints are selected in the lexicographical order² when multiple choices are possible (e.g. line 7 and 17). It turns out this simple querying procedure enjoys the same convergence rate than a regular Perceptron, with the added empirically observed benefit of providing stronger compression (see Section 4 for empirical results).

Proposition 2. *Consider Problem (7) and let γ be the radius of the largest inscribed sphere in \mathcal{W} . Define M the number of Perceptron updates performed by the Perceptron-based Localization Algorithm 3 (i.e. M is the number of times line 18 of PERCEPTRON() of Algorithm 3 is executed). Then the following holds: $M \leq 1/\gamma^2$.*

Proof. We recall that the usual definition of the margin of D is $\min_{x \in D} \langle w^*, x \rangle$ and note that γ is related to it since $\forall n \in [N], \langle w^*, x_n \rangle / \|x_n\|_2 \geq \gamma$. Let $\mathcal{S} \doteq \{a_1, \dots, a_M\}$ be the sequence of points used to perform Perceptron updates across a complete execution of Algorithm 3. Thus, \mathcal{S} is a sequence from D (with possible duplicates) and w^* achieves a margin at least γ with all points in \mathcal{S} . From [Blo62, Nov62] we know that the number M of Perceptron updates on any arbitrary sequence linearly separable with margin γ is no more than $1/\gamma^2$. Since we use w^t as a starting point to compute w^{t+1} , the execution of the cutting-plane algorithm is tied to the execution of the Perceptron algorithm on \mathcal{S} . Therefore, there is less than $1/\gamma^2$ Perceptron updates during the execution of the algorithm. Alternatively, $|\mathcal{S}| \leq 1/\gamma^2$ since all points in \mathcal{S} correspond to a Perceptron update, thus a mistake. \square

On a side note, the same argument can be applied to obtain similar results with most Perceptron-like learning procedures (see for instance [LZH⁺02, CDK⁺06]).

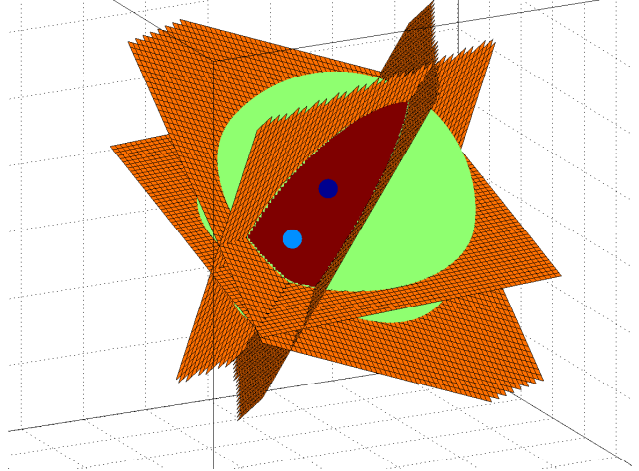


Figure 1: An Example of version space where the Chebyshev Center (light blue) is a bad approximation of the gravity center (dark blue).

3.3 Center of Gravity and Approximations

The question of computing a query point w^t is of central importance in cutting-plane localization algorithms. As we have seen, a simple Perceptron can already yield interesting computational results for that matter. A more assiduous analysis of this question can be conducted by looking at the volume reduction $\text{Vol}(\mathcal{C}^{t+1})/\text{Vol}(\mathcal{C}^t)$ of \mathcal{C}^t from one iteration to the next. The notion of center of gravity is going to be pivotal to this end.

Definition 1 (Center of Gravity). Let \mathcal{C} be a closed set in \mathbb{R}^n . The *center of gravity* (CG) $\mathbf{cg}(\mathcal{C})$ of \mathcal{C} is defined by as $\mathbf{cg}(\mathcal{C}) \doteq \int_{\mathcal{C}} z dz / \int_{\mathcal{C}} dz$.

The center of gravity is deeply tied to the volume of \mathcal{C}^t and plays a central role in devising cutting-plane algorithms for which the volume reduction $\text{Vol}(\mathcal{C}^{t+1})/\text{Vol}(\mathcal{C}^t)$ is the largest. Theorem 1 reports one of the most fundamental property of the center of gravity (see [Gru60, New65, Lev65, BV08])

Theorem 1 (Partition of Convex bodies). *Let $\mathcal{C} \in \mathbb{R}^d$ a convex body of center of gravity $\mathbf{cg}(\mathcal{C})$ and h a hyperplane such that $\mathbf{cg}(\mathcal{C}) \in h$. Thus, h divide \mathcal{C} in two subsets \mathcal{C}_1 and \mathcal{C}_2 and the following relations hold for $i = 1, 2$: $\text{Vol}(\mathcal{C}_i) \geq e^{-1} \text{Vol}(\mathcal{C})$*

The center of gravity method proposed by [New65, Lev65] consists in querying $w^t = \mathbf{cg}(\mathcal{C}^t)$ and typically

²This is an arbitrary choice and any total order over \mathbb{R}^d can be used instead

have a very fast convergence rate as the version space is almost halved at each step. More precisely, a direct consequence of Theorem 1 is that the volume of \mathcal{C}^t is bounded by $\text{Vol}(\mathcal{C}^t) \leq (1 - 1/e)^t \text{Vol}(\mathcal{C}^0)$. However, computing the center of gravity is hard, making the center of gravity method impractical. Instead, one has to consider structural or numerical approximations to the center of gravity.

Definition 2 (Chebyshev’s Center). Let \mathcal{C} a set in \mathbb{R}^n . Chebyshev’s center (CC) of \mathcal{C} , $\text{cc}(\mathcal{C})$ is the center of the largest inscribed ball in \mathcal{C} :

$$\text{cc}(\mathcal{C}) = \arg \min_{\hat{z}} \max_z \|z - \hat{z}\|_2^2.$$

Chebyshev’s center is used as a computationally efficient approximation of the center of gravity for cutting-plane algorithms since the late 70’s [EM75] (see, e.g. [BV04] for a linear formulation of the problem). Unfortunately, the interesting property of Theorem 1 does not carry over with Chebyshev’s center. One problem in machine learning related to Chebyshev’s center is the extensively studied Support Vector Machine (SVM) [Vap95] defined as :

$$\min_w \frac{1}{2} \|w\|_2^2 \quad \text{s.t.} \quad \begin{cases} w \in \mathcal{X} \\ y_n \langle w, x_n \rangle \geq 1, n \in [N]. \end{cases} \quad (9)$$

A notable property of the SVM is that its solution w_{SVM} is closely related to the center of the largest inscribed ball in \mathcal{W} and is an approximation of the center of gravity [RHC01]. Indeed, w_{SVM} is actually a rescaled Chebyshev’s center [TK02] [RHC01].

On the other hand, numerical approximations aim at finding a point that is in the close neighborhood of the center of gravity. One of the contributions of this paper is to give a generalized version of Theorem 1 for approximations of the center of gravity, thus laying a theoretical justification for these methods.

Theorem 2 (Generalized Partition of Convex bodies). Let \mathcal{C} be a closed convex body in \mathbb{R}^d and $\text{cg}(\mathcal{C})$ its center of gravity. Let h_x a hyperplane of normal vector x , $\|x\|_2 = 1$ and define the upper (resp. lower) partition \mathcal{C}^+ (resp. \mathcal{C}^-) of \mathcal{C} by h_x as

$$\begin{aligned} \mathcal{C}^+ &\doteq \mathcal{C} \cap \{w \in \mathbb{R}^d : \langle x, w \rangle \geq 0\} \\ \mathcal{C}^- &\doteq \mathcal{C} \cap \{w \in \mathbb{R}^d : \langle x, w \rangle < 0\}. \end{aligned}$$

The following holds true: if $\text{cg}(\mathcal{C}) + \Lambda x \in \mathcal{C}^+$ then

$$\text{Vol}(\mathcal{C}^+) / \text{Vol}(\mathcal{C}) \geq e^{-1} (1 - \lambda)^d,$$

where

$$\Lambda = \lambda \Theta_d \frac{\text{Vol}(\mathcal{C}) H_{\mathcal{C}^+}}{R^d H_{\mathcal{C}^-}},$$

with $\lambda \in \mathbb{R}$ an arbitrary real, Θ_d a constant depending only on d , R the radius of the $(d - 1)$ -dimensional ball B of volume $\text{Vol}[B] \doteq \text{Vol}[\mathcal{C} \cap \{w \in \mathbb{R}^d : \langle x, w \rangle = 0\}]$ and $H_{\mathcal{C}^+} = \max_{a \in \mathcal{C}^+} a^T x$ (resp. $H_{\mathcal{C}^-} = \min_{a \in \mathcal{C}^-} a^T x$)

Proof. The proof is a (non-trivial) extension of Grunbaum’s one for Theorem 1 [Gru60]. Due to space restriction, we cannot expose it here in full and refer the interested reader to <http://pageperso.lif.univ-mrs.fr/~ugo.louche/paper/activeCPSuppl.pdf>. \square

Theorem 2 extends Theorem 1 to the situation when an approximation of the center of gravity is considered; it reduces to Theorem 1 when applied to the very center of gravity. This is to the best of our knowledge the first result of this kind and this is a result that is of its own interest, which may benefit to many fields of computer science. Here, the purpose of Theorem 2 is essentially to validate the use of approximations of the center of gravity $\text{cg}(\mathcal{C})$ in the procedures at hand, which is inevitable due to the complexity of exactly finding this point. We will more precisely use it in two occasions: a) for center-of-gravity-based compression scheme methods and b) in the active learning setting (see below).

3.4 Active Learning with Cutting Planes

An interesting situation of learning is that of active learning when the algorithm is presented with unlabelled data and it has to query for the labels of the training points that carry the most information to build a relevant decision boundary. Given a volume \mathcal{C} inside which a good classifier w^* for the classification task at hand is known to lie, the amount of information carried by a labeled training point (x, y) (where y has been queried) might be for instance measured by how (x, y) can be used to identify within \mathcal{C} an (hopefully small) volume $\mathcal{C}' \subseteq \mathcal{C}$ where w^* lives. Termed otherwise, the amount of information provided by (x, y) might be measured as the volume reduction induced by the knowledge of (x, y) : this is exactly the type of information cutting-plane methods build upon. We take advantage of this philosophy shared by active learning methods and cutting-plane algorithms to argue it is easy to transform a cutting-plane algorithm into an active learning method. Based on the idea of maximum volume reduction, the question to address is simply that of identifying a training pattern x in D such that, independently of the label it might receive, is guaranteed to define a cutting hyperplane of equation

Algorithm 4 Top: a generic cutting-plane active learning procedure; w^t is computed as the ‘center’ of \mathcal{C}^t —center may refer to the center of gravity of the Chebyshev center. Bottom: a possible implementation of `QUERY()`: sampling strategies are given in, e.g., [RHC01, LV06, KN12].

```

1:  $\mathcal{C}^0 \leftarrow \mathcal{B}$ 
2:  $t \leftarrow 0$ 
3: repeat
4:    $w^t \leftarrow \text{center}(\mathcal{C}^t)$ 
5:    $x_{n_t}, y_{n_t} \leftarrow \text{QUERY}(\mathcal{C}^t, D)$ 
6:   if  $y_{n_t} \langle w^t, x_{n_t} \rangle < 0$  then
7:      $\mathcal{C}^{t+1} \leftarrow \mathcal{C}^t \cap \{z : y_{n_t} \langle z, x_{n_t} \rangle \geq 0\}$ 
8:      $t \leftarrow t + 1$ 
9:   end if
10: until  $\mathcal{C}^t$  is small enough
11: return  $w^t$ 
12:
13: function QUERY( $\mathcal{C}, D$ )
14:   Sample  $M$  points  $s_1, \dots, s_M$  from  $\mathcal{C}$ 
15:    $\mathbf{g} \leftarrow \sum_{k=1}^M s_k / M$ 
16:    $x \leftarrow \arg \min_{x_i \in D} \langle \mathbf{g}, x_i \rangle$ 
17:    $y \leftarrow \text{get label from an expert}$ 
18:   return  $x, y$ 
19: end function

```

$\langle x, w \rangle = 0$ that intersects the current convex \mathcal{C} in a controlled way. To do so, a typical good query point is one that is as close as possible to the ‘center’ of \mathcal{C} , where center may have the few meanings discussed above (cf. center of gravity, Chebyshev’s center). The algorithm given in Table 4 is a generic active learning algorithm that is based on the classical cutting-plane approach. Making active learning algorithms from cutting-plane methods is a route that has been taken by [TK02], even though the connection with cutting-plane algorithms was not clearly identified.

Being able to approximate the center of gravity of a convex polytope is pivotal for the design of active learning strategies. It is interesting to note that in the recent years, methods have been devised to uniformly sample from the version space such as the *Hit-and-Run* algorithm of [LV06] or a billiard algorithm of [Ruj97]. More recently, the *Dikin Walk* algorithm of [KN12] provided a strongly polynomial algorithm for approximate uniform sampling over the version space while the *Expectation Propagation* method of [Min13] gave a Bayesian interpretation of billiard algorithms. Notably, these methods have been successfully used with cutting planes for active Boosted Learning [TSCD11]. Another practical approach we should mention is the

one proposed in [RHC01] that consists in repeatedly running a Perceptron over a permutation of the training set: in the active learning setting, the number of labeled points available is just too low to produce interesting approximation of the center of gravity with this method.

A by-product of our active learning procedure is that we now solve a Bayes Point Machine (BPM) problem [RHC01] at each step t by finding the center of gravity of the current convex body \mathcal{C}^t . Therefore, we can turn our active learning procedure into a full active learning algorithm—that we dub **Active-BPM**—for free by using the center of gravity for classification. Note that this is one of many possible instantiations of our procedure, which is nonetheless of interest as it is the BPM-counterpart the **Active-SVM** algorithm of Tong and Koller [TK02].

In conclusion, Theorem 2 provides a general guideline to systematically query the training point that comes with the best volume reduction guarantees. This is a theoretically sound and viable strategy for active learning that comes with a theoretical bound on the induced volume reduction, the lack of which was an essential limit of the Chebyshev’s center-based method of [TK02].

4 Numerical Simulations

Here, we present some empirical simulations based on the algorithms described throughout this paper in both passive and active learning settings.

4.1 Synthetic Data and Perceptron-based Localization Algorithm

We generate a toy dataset of 1,000 2-dimensional datapoints. Each point is uniformly drawn on a 20-by-20 square centered at the origin. We label this dataset according to a classifier w^* uniformly drawn over the unit circle. In order to have only positive labels, negative examples are reflected through the origin. We then enforce a minimal margin γ by pruning examples x_i for which $\langle w^*, x_i \rangle < \gamma$. This last modification allows us to have some control over the size of the version space \mathcal{W} . The downside of this is that we no longer have exactly 1,000 datapoints (though during our experiments we noted that the size of the dataset stays mostly the same for reasonable margin values).

For these experiments, we use the Perceptron-based Localization algorithm (Algorithm 3). We implement it with three different oracle strategies for selecting cut-

ting planes. The first strategy (which we call *Largest Error*) picks the cutting plane with the lowest margin. The second one (*Smallest Error*) picks the cutting plane with the highest negative margin, that is to say points that are incorrectly classified but close to the decision boundary. Finally, the third one (*Random Error*) simply picks a cutting plane with negative margin at random. It should also be noted that our instantiation of the Perceptron algorithm picks the update vector that realizes the lowest margin for its internal update—line (18) of PERCEPTRON() in Algorithm 3. This is mostly an arbitrary choice and we only mention it for the sake of reproducibility.

The first experiment consists in a single run over a dataset of margin $\gamma = 0.1$. We monitor both the number of cutting planes generated and the number of internal Perceptron updates for each cutting plane. The presented results are averaged over 1,000 runs.

The left pane of Figure 2 supports the soundness of our approach in the case of a compression scheme with no more than 6 cutting planes for the best strategy (Largest Error). Additionally, we can observe a sharp decrease after the third cutting plane with this strategy and 80% of the time, only 4 cutting planes are required to model the dataset. In contrast, the right-hand side of Figure 2 reveals a trade-off between the number of cutting planes used and the number of internal updates for each cutting plane. We observe a smooth shift across our three strategies with Smallest Error putting the emphasis on small number of internal updates. In all respect, the Random Error strategy acts as a middle ground between the two other extreme approaches.

For the second experiment the margin (i.e. the volume of \mathcal{W}) is variable with values between 0.01 and 0.3. We also monitor the total number of internal updates rather than the *per cutting plane* value for the three strategies and a regular Perceptron Algorithm³. Remind that this value is bounded from Proposition 2. This bound also holds for the regular Perceptron.

The previously observed behavioral shift across the three strategies is confirmed by Figure 3. Additionally, some relative robustness is observed with respect to γ , especially when the emphasis is put on querying a small number of cutting planes. It is interesting to note that the Random Strategy makes nearly as few updates as Smallest error while still querying a—relatively—low number of cutting planes. Finally, all three strategies are making slightly less updates than the regular Perceptron. To conclude, note that the theoretical bound

³More precisely, we use the exact same Perceptron than the one used for the internal loop but ran on the full dataset

of Proposition 2 is far too big to be plotted on the plot on the left of Figure 2.

4.2 Active Learning on Real Data

We illustrate our method for active learning on text classification data. For easy comparison, we follow an experimental procedure similar as the one in [TK02]. Namely, we use the *Reuters-21578 —ModApte* variation— and Newsgroups datasets⁴. The Reuters dataset is composed of 8,293 documents represented in TF-IDF form for 18,933 words. The dataset spans 65 topics such as *Earn*, *Coffee* or *Cocoa* and is split in 5,946 training examples and 2,347 test examples. On the other hand, the Newsgroups dataset accounts for 18,846 documents of 26,214 features splitted in 20 topics. Half of this dataset is uniformly picked for training while the rest is kept for testing purposes. On both datasets we train a “one-versus-all” classifier for each class. We start by creating a pool of unlabeled training examples sampled from the training set. Then we run Algorithm 4. We use two variations of the QUERY() function: one based on the Chebyshev center (note that this is equivalent to the Active-SVM of [TK02]), and the other based on an approximation of the center of gravity from Minka’s Expectation Propagation method [Min13]. This last approach corresponds to the Active-BPM algorithm and has, to the best of our knowledge, never been used before. It is a direct application of Active Learning algorithms with Cutting planes method to the Bayes Point Machine. For both methods, we use two pools of different sizes (500 and 1,000 examples). For initialization reasons, each pool comes with two already labeled vectors.⁵ All the computations are done with a linear kernel and the presented results are class-wise accuracy measurements on the test examples over the 10 most represented classes. The values reported here are an average of these measures over 25 runs. We complement these two datasets with Gunnar Raetsch’s Banana dataset. The Banana dataset is a widely used dataset of 2-dimensional points split into two classes from which we extract 400 training and 4900 test examples. Due to its small size, the whole training set is used for the pool of unlabeled example. The computations are realized with an RBF kernel of parameter $\sigma = 0.5$ and

⁴Available at <http://www.cad.zju.edu.cn/home/dengcai/Data/TextData.html>

⁵SVM and CC are computed with libSVM: <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>. BPM and CG are computed from Minka’s own implementation of EP for BPM in matlab: <http://research.microsoft.com/en-us/um/people/minka/papers/ep/bpm/>

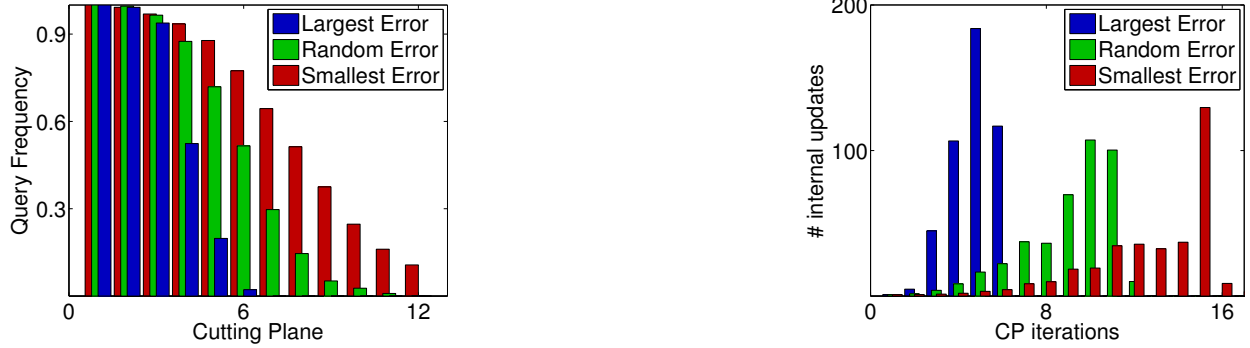


Figure 2: *Left* : for each value i the bar represents the empirical probability (over 1,000 runs) to query at least i cutting planes. *Right*: each bar represents the number of internal Perceptron updates computed after each Cutting Plane loop.

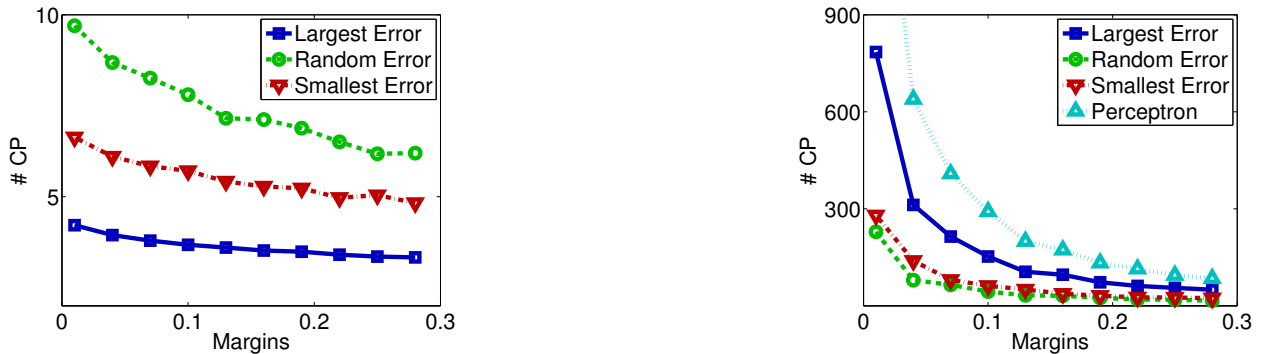


Figure 3: *Left*: The average number of cutting planes used for each strategy with respect to the value of γ . *Right*: the total number of internal updates with respect to γ . The fourth plot corresponds to a regular Perceptron

presented results are averaged over 50 runs.

Figure 4 graphically depicts the behavior of the so-called **Active-SVM** [TK02] and the **Active-BPM** algorithms on each dataset. Namely, in both algorithms, the queries are selected according to their distance to the “centroid” of \mathcal{C} , which, in turn, serves as classifier. The difference between these two algorithms lies in that **Active-SVM** uses the Chebyshev center and **Active-BPM** the center of gravity for centroid. In Figure 4, data are represented by circles of squares whether they correspond to results achieved by **Active-SVM** or **Active-BPM**. Additionally, for the Reuters and Newsgroups datasets, dashed plots correspond to the pool of 500 examples while dotted plots relate to the pool of 1000 examples. The error bounds on the third plot (Banana) correspond to the usual standard deviation. Each plot represents the accuracy of those algorithms with respect to the number of queries made. We can see that **Active-BPM** systematically outperforms **Active-SVM** and increases its accuracy faster for all datasets, already attaining an accuracy of 0.9

after roughly 10 queries for both Reuters and Newsgroups datasets. Both algorithm seem to stabilize after 30 queries, with the **Active-BPM** being slightly more accurate than its SVM counterpart. For the Banana dataset, the accuracy increase in the first queries is a lot smoother, with an accuracy for **Active-BPM** of roughly 0.8 after 20 queries. Both algorithms seem to have converged after 60 queries. Comparatively, not only does **Active-BPM** clearly dominate its SVM counterpart but it is also more stable as evidenced by the error bars which become negligible past the 60th query.

5 Conclusion and Future Directions

In this paper, we have shown that deep connections exist between Localization methods and Learning algorithms. Both fields have extensively characterized and studied similar concepts over the past years, sometime independently. On the other hand, complementary re-

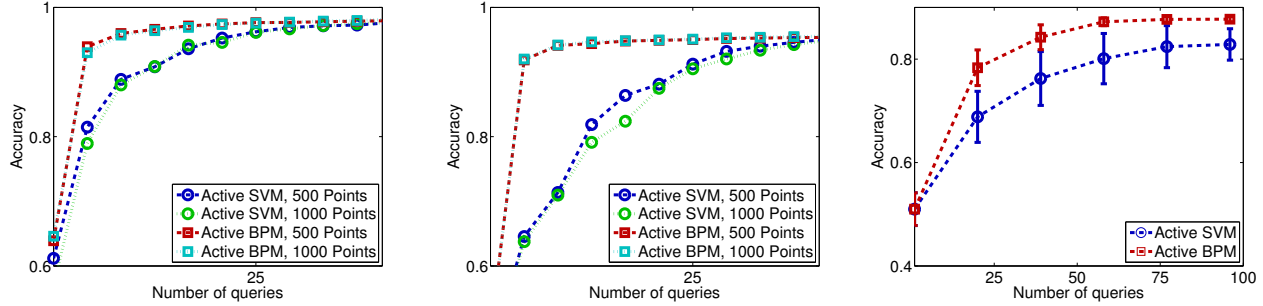


Figure 4: Accuracy on the Reuters (left) and Newsgroups (middle) datasets for **Active-SVM** and **Active-BPM** for pools of 500 and 1000 examples. Left: accuracy with error bars on the Banana dataset (Gunnar Raetsch) for **Active-SVM** and **Active-BPM**.

sults have been found in each community. A notable example is the absence of a kernel approach in the Cutting Planes literature while center of gravity methods were mostly unknown in machine learning until Herbrich’s BPM [RHC01]. We may also mention that the Cutting planes’ equivalent of the famous SVM [Vap95] appears as soon as the 70’s in [EM75]. This work is a testimony on how it is possible to derive new learning algorithms, both efficient and theoretically funded, by reformulating Cutting Planes approach for the learning paradigm. Besides the cutting plane-related flavor of the present work, it should be restated that Theorem 2 has a value that goes beyond the scope of this paper. A field that may be impacted by this result is obviously that of computational geometry where most of the results about the computation of centers of gravity come from; nonetheless, it should be noted that more closely related works could also benefit from our result. For instance, if we consider the active learning methods whose query steps rely on explicit exploration of all the possible query/label combinations (see, e.g. [RM01]), then Theorem 2 provides a tool to devise natural and theoretically sound heuristics to effectively locate the most informative query points, or, in other words, those that may lead to the smallest expected error.

Among all the possible extensions of this work, one we are particularly interested in is to study how these results may carry over to the multiclass setting and provide proper multiclass active algorithms based on, for example, Crammer’s Ultraconservative Additive Algorithms [CS03].

References

[BBZ07] M. F. Balcan, A. Broder, and T. Zhang. Margin based active learning. In *COLT*, 2007.

[Blo62] H. Block. The perceptron: a model for brain functioning. *Reviews of Modern Physics*, 1962.

[BV04] S. P. Boyd and L. Vandenberghe. *Convex optimization*. 2004.

[BV08] S. Boyd and L. Vandenberghe. Localization and cutting-plane methods. 2008.

[CDK⁺06] K. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz, and Y. Singer. Online passive-aggressive algorithms. *JMLR*, 2006.

[CS03] K. Crammer and Y. Singer. Ultraconservative online algorithms for multiclass problems. *JMLR*, 2003.

[DSP10] F. Dabbene, P. S. Shcherbakov, and B. T. Polyak. A randomized cutting plane method with probabilistic geometric convergence. *SIAM Journal on Optimization*, 2010.

[EM75] J. Elzinga and T. G. Moore. A central cutting plane algorithm for the convex programming problem. *Mathematical Programming*, 1975.

[FS99] Y. Freund and R. E. Schapire. Large margin classification using the perceptron algorithm. *Machine learning*, 1999.

[FS09] V. Franc and S. Sonnenburg. Optimized cutting plane algorithm for large-scale risk minimization. *JMLR*, 2009.

[FW95] S. Floyd and M.K. Warmuth. Sample compression, learnability, and the vapnik-chervonenkis dimension. *Machine Learning*, 1995.

[GHST05] T. Graepel, R. Herbrich, and J. Shawe-Taylor. PAC-Bayesian Compression Bounds on the Prediction Error of Learning Algorithms for Classification. *Machine Learning*, 2005.

[GK10] D. Golovin and A. Krause. Adaptive submodularity: A new approach to active learning and stochastic optimization. *CoRR*, 2010.

[Gof80] J.-L. Goffin. The relaxation method for solving systems of linear inequalities. *Mathematical Operations Research*, 1980.

- [Gru60] B. Grunbaum. Partitions of mass-distributions and of convex bodies by hyperplanes. *Pacific Journal of Mathematics*, 1960.
- [GSSS13] A. Gonen, S. Sabato, and S. Shalev-Shwartz. Efficient active learning of halfspaces: an aggressive approach. *JMLR*, 2013.
- [JFY09] T. Joachims, T. Finley, and C.-N. Yu. Cutting-plane training of structural svms. *Machine Learning*, 2009.
- [Kel60] J. E. Kelley. The cutting plane method for solving convex programs. *SIAM*, 1960.
- [KN12] R. Kannan and H. Narayanan. Random walks on polytopes and an affine interior point method for linear programming. *Mathematics of Operations Research*, 2012.
- [Lev65] A. Levin. On an algorithm for the minimization of convex functions. *Soviet Math. Doklady*, 1965.
- [LG94] David D Lewis and William A Gale. A sequential algorithm for training text classifiers. In *Proceedings of the 17th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 3–12. Springer-Verlag New York, Inc., 1994.
- [LV06] L. Lovász and S. Vempala. Hit-and-run from a corner. *SIAM Journal on Computing*, 2006.
- [LZH⁺02] Y. Li, H. Zaragoza, R. Herbrich, J. Shawe-Taylor, and J. Kandola. The perceptron algorithm with uneven margins. 2002.
- [Min13] T. P. Minka. Expectation propagation for approximate bayesian inference. *CoRR*, 2013.
- [MS54] T. S. Motzkin and I. J. Schoenberg. The relaxation method for linear inequalities. *Canadian Journal of Mathematics*, 1954.
- [Nes95] Y. Nesterov. Cutting plane algorithms from analytic centers: efficiency estimates. *Mathematical Programming*, 1995.
- [New65] D. J. Newman. Location of the maximum on unimodal surfaces. *J. ACM*, 1965.
- [Nov62] A.B.J. Novikoff. On convergence proofs on perceptrons. In *Proceedings of the Symposium on the Mathematical Theory of Automata*, 1962.
- [RHC01] T. Graepel R. Herbrich and C. Campbell. Bayes point machines. *JMLR*, 2001.
- [RM01] Nicholas Roy and Andrew McCallum. Toward optimal active learning through sampling estimation of error reduction. In *ICML*, pages 441–448. Morgan Kaufmann Publishers Inc., 2001.
- [Ros58] F. Rosenblatt. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 1958.
- [Ruj97] P. Ruján. Playing billiards in version space. *Neural Computation*, 1997.
- [Set12] B. Settles. Active learning. In *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 2012.
- [TK02] S. Tong and D. Koller. Svm active learning with applications to text classification. *JMLR*, 2002.
- [TSCD11] K. Trapeznikov, V. Saligrama, D. A. Castañón, and A. David. Active boosted learning (act-boost). In *AISTATS*, 2011.
- [TVSL10] C. H. Teo, S. V. N. Vishwanathan, A. J. Smola, and Q. V. Le. Bundle methods for regularized risk minimization. *JMLR*, 2010.
- [Vap95] V. N. Vapnik. *The Nature of Statistical Learning Theory*. 1995.