



# On the Existence of Linear Weak Learners and Applications to Boosting\*

SHIE MANNOR

RON MEIR

*Department of Electrical Engineering, Technion, Haifa 32000, Israel*

shie@tx.technion.ac.il

rmeir@ee.technion.ac.il

**Editors:** Yoshua Bengio and Dale Schuurmans

**Abstract.** We consider the existence of a linear weak learner for boosting algorithms. A weak learner for binary classification problems is required to achieve a weighted empirical error on the training set which is bounded from above by  $1/2 - \gamma$ ,  $\gamma > 0$ , for any distribution on the data set. Moreover, in order that the weak learner be useful in terms of generalization,  $\gamma$  must be sufficiently far from zero. While the existence of weak learners is essential to the success of boosting algorithms, a proof of their existence based on a geometric point of view has been hitherto lacking. In this work we show that under certain natural conditions on the data set, a linear classifier is indeed a weak learner. Our results can be directly applied to generalization error bounds for boosting, leading to closed-form bounds. We also provide a procedure for dynamically determining the number of boosting iterations required to achieve low generalization error. The bounds established in this work are based on the theory of geometric discrepancy.

**Keywords:** boosting, weak learner, geometric discrepancy

## 1. Introduction

One of the most exciting developments in the fields of Machine Learning and Pattern Recognition in recent years has been the development of the boosting approach to learning (Schapire, 1990; Freund & Schapire, 1996a). Boosting, similarly to other ensemble based methods, such as Bagging (Breiman, 1996) and Mixture of Experts (Jordan & Jacobs, 1994), has been shown to be effective in the construction of successful classifiers. In addition to many impressive empirical demonstrations of the utility of the procedure, there has also been recently a great deal of theoretical work providing guaranteed performance bounds (e.g., Schapire et al., 1998). It turns out that a key ingredient in the success of the boosting approach is its ability to yield classifiers that achieve a large margin, implying that the decision boundary defining the classifier, while separating the training points, is able to retain as large a distance as possible from them. The notion of margin also plays an essential role in the theory of support vector machines (Vapnik, 1998; Cristianini & Shawe-Taylor, 2000).

\*Support from the Ollendorff center of the Department of Electrical Engineering at the Technion is acknowledged.

In order to derive useful generalization bounds for boosting-type algorithms, it has turned out to be essential to be able to guarantee the existence of a so-called *weak learner*. A weak learner for binary classification problems is one for which the weighted empirical error is guaranteed to be smaller than  $1/2 - \gamma$ ,  $\gamma > 0$ , for *any* distribution on the data. In boosting, the final composite classifier is formed by a convex combination of base classifiers. In order for the procedure to be effective, these base classifiers, also known as weak hypotheses or weak learners, are required to obey the so-called weak-learning assumption, defined in Section 2 below. While many base classifiers have been used in practice, it has not been clear up to this point under what conditions weak learners actually exist. For example, we show in Section 2, that a commonly used learner, the so-called decision stump, cannot constitute a useful weak learner. This situation is somewhat disturbing, as the existence of a weak learner is essential to the theory (and practice) of boosting. Some work providing conditions for the existence of weak learners has been carried out in Freund (1995), Freund and Schapire (1996b), and Breiman (1998). The main distinction with our work is our concentration on purely geometric arguments and real-valued inputs.

In this work we prove that linear classifiers are indeed weak learners. Clearly, any other system based on linear classifiers, such as decision trees with oblique splits and neural networks (both with limited complexity) are also weak learners. Moreover, the proof technique used suggests a randomized algorithm, which achieves the desired result (see the comments at the end of Section 4.2). An important message of this paper, following the discussion in Section 2, is that weak learning in itself is not sufficient to guarantee that boosting performs well. This prompts us to introduce the notion of an *effective* weak learner (see Definition 3), which ensures low error of the combined classifier formed in boosting. In order to establish the latter result, a simple assumption needs to be made concerning the data. As we show in Section 4, some regularity assumption is necessary, as no *effective* weak learner can exist for arbitrary data sets.

The proof method employed in this paper is based on the theory of geometric discrepancy, a sub-field of combinatorial geometry, which deals with irregularities of distributions. The reader is referred to the excellent text-book by Matoušek (1999) for a general introduction to this field.

The remainder of the paper is organized as follows. Section 2 recalls the boosting framework and motivates the importance of finding weak classifiers. In Section 3 we introduce some basic definitions and key results from the field of geometric discrepancy. The main result regarding the existence of a weak linear learner is presented in Section 4, and the application to error bounds for boosting is provided as a corollary of the main result. We close in Section 5 by drawing some conclusions and mentioning several open questions. Some of the more technical proofs have been relegated to appendices. It should be commented that no attempt has been made to optimize the constants appearing in the bounds. Finally, we note that we use the term ‘classifier’ and ‘hypothesis’ interchangeably throughout this work.

## 2. Boosting and weak learners

Boosting is a general meta-procedure for constructing a strong classifier from a set of weak classifiers (Schapire, 1990). Each weak classifier is constructed based on a re-weighted

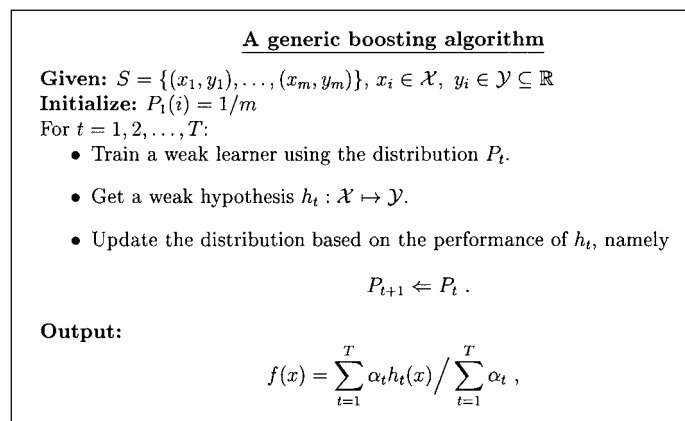


Figure 1. A generalized boosting algorithm. The update of  $P_t$  and the calculation of  $\alpha_t$  are algorithm-specific. Some examples can be found in Freund and Schapire (1996), Friedman, Hastie, and Tibshirani (2000), Schapire and Singer (1999), and Meir, El-Yaniv, and Ben-David (2000).

version of the data set. Although there exist many versions of boosting to-date (e.g., Freund & Schapire, 1996a; Friedman, Hastie, & Tibshirani, 2000; Mason et al., 2000; Schapire & Singer, 1999), most of them fall into the general framework depicted in figure 1, adapted from Schapire and Singer (1999). Note that the number of cycles  $T$  appearing in the algorithm is often determined by assessing the error on an independent validation set (or using cross-validation). In Section 4.2 we provide some theoretical guidelines for the selection of  $T$ .

We begin our discussion by recalling the notion of the Vapnik-Chervonenkis (VC) dimension.

*Definition 1.* Let  $\mathcal{H}$  be a class of functions from  $\mathcal{X}$  to  $\{-1, +1\}$ , and let  $X$  be a set of  $m$  points from  $\mathcal{X}$ . We say that  $\mathcal{H}$  shatters  $X$  if  $\mathcal{H}$  can compute all  $2^m$  dichotomies on  $X$ . The VC-dimension of  $\mathcal{H}$  is the size of the largest shattered subset of  $\mathcal{X}$ .

Since the final classifier constructed within the boosting framework is a convex combination of weak classifiers, it is of interest to investigate error bounds for convex combinations of classifiers. In this context we recall the elegant results derived by Schapire et al. (1998). Let  $\mathcal{H}$  be a class of binary classifiers of VC-dimension  $d_H$ , and denote by  $\text{co}(\mathcal{H})$  the convex hull of  $\mathcal{H}$ ,

$$\text{co}(\mathcal{H}) = \left\{ f : f(x) = \sum_i \alpha_i h_i(x), h_i \in \mathcal{H}, \alpha_i \geq 0, \sum_i \alpha_i = 1 \right\} .$$

Given a sample  $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$ ,  $x_i \in \mathcal{X}$  and  $y_i \in \{-1, +1\}$ , of  $m$  examples drawn independently at random from a probability distribution  $D$  over  $\mathcal{X} \times \{-1, +1\}$ , Schapire et al. (1998) showed that for  $m > d_H$ , with probability at least  $1 - \delta$ , for every

$f \in \text{co}(\mathcal{H})$  and  $\theta > 0$ ,

$$\mathbf{P}_D[Yf(X) \leq 0] \leq \mathbf{P}_S[Yf(X) \leq \theta] + O\left(\frac{1}{\sqrt{m}} \left(\frac{d_{\mathcal{H}}(\log(m/d_{\mathcal{H}}))^2}{\theta^2} + \log \frac{1}{\delta}\right)^{1/2}\right), \quad (1)$$

where the *margin-error*  $\mathbf{P}_S[Yf(X) \leq \theta]$  denotes the fraction of training points for which  $y_i f(x_i) \leq \theta$ , namely

$$\mathbf{P}_S[Yf(X) \leq \theta] = \frac{1}{m} \sum_{i=1}^m I(y_i f(x_i) \leq \theta),$$

where  $I(E)$  is the indicator function for the event  $E$ . Note that the term  $\mathbf{P}_D[Yf(X) \leq 0]$  is simply the probability of misclassification of the classifier  $h(x) = \text{sgn}(f(x))$ .

It is useful to compare (1) to the classic VC bounds (Vapnik & Chervonenkis, 1971), which *do not* take into account the fact that  $f$  is real-valued and that it is a convex combination of functions from  $\mathcal{H}$ . Let  $d_{\text{co}(\mathcal{H})}$  denote the VC-dimension of  $\text{co}(\mathcal{H})$ . Then one finds (e.g., Anthony & Bartlett, 1999) that with probability at least  $1 - \delta$ , for every  $f \in \text{co}(\mathcal{H})$  and  $m > d_{\text{co}(\mathcal{H})}$ ,

$$\mathbf{P}_D[Yf(X) \leq 0] \leq \mathbf{P}_S[Yf(X) \leq 0] + O\left(\frac{1}{\sqrt{m}} (d_{\text{co}(\mathcal{H})} \log(m/d_{\text{co}(\mathcal{H})}) + \log(1/\delta))^{1/2}\right). \quad (2)$$

Note that the first term in (2) is simply the fraction of misclassified data points. Comparing (1) and (2) we note two major differences. The former bound contains the extra parameter  $\theta$ , which allows one to fine tune the bound in order to achieve better performance. However, since  $\mathbf{P}_S[Yf(X) \leq 0] \leq \mathbf{P}_S[Yf(X) \leq \theta]$  the first term in (2) is always *smaller* than the corresponding term in (1). However, observe that the second term in (1) may be significantly smaller than the corresponding term in (2), since it contains the VC-dimension of  $\mathcal{H}$  rather than the VC-dimension of  $\text{co}(\mathcal{H})$ , which may be *significantly* larger (e.g., Section 14.4 in Anthony & Bartlett, 1999). It should be clear that there is a trade-off between the two terms appearing on the r.h.s. (1). While the first term is monotonically increasing with  $\theta$ , the second term decreases monotonically. This situation is very much reminiscent of the idea of structural risk minimization, suggested by Vapnik (1982) in the context of hierarchical hypothesis classes.

In order to obtain a useful bound, one needs to be able to estimate the first term in (1). In order to do so, we formally introduce the notion of a *weak hypothesis*.

*Definition 2.* Let  $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$ ,  $x_i \in \mathcal{X}$ ,  $y_i \in \{-1, +1\}$ , be a set of  $m$  examples. We say that a hypothesis class  $\mathcal{H}$   $\gamma$ -satisfies the weak hypothesis assumption for  $S$  if, for any probability distribution  $P$  on  $S$ , there exists  $h \in \mathcal{H}$  such that the weighted empirical error of  $h$  is smaller than  $1/2 - \gamma$ , namely

$$\sum_{i=1}^m P_i I(h(x_i) \neq y_i) \leq \frac{1}{2} - \gamma. \quad (3)$$

There are two important facts to note concerning this definition. First, it is required to hold for *any* distribution  $P$ , and second, one only demands that the error be smaller than  $1/2 - \gamma$ . Since an error of  $1/2$  is always achievable,<sup>1</sup> this is a very weak assumption if  $\gamma$  is small.

Schapire et al. (1998) have provided a bound on the empirical error of the composite classifier formed by the AdaBoost algorithm (Freund & Schapire, 1996a) run for  $T$  steps. In particular, denote the error of each weak learner by  $\epsilon_t$ ,  $t = 1, \dots, T$ , then Theorem 5 in Schapire et al. (1998) shows that

$$\begin{aligned} \mathbf{P}_S[Yf(X) \leq \theta] &\leq \prod_{t=1}^T \sqrt{4\epsilon_t^{1-\theta}(1-\epsilon_t)^{1+\theta}} \\ &\leq ((1-2\gamma)^{1-\theta}(1+2\gamma)^{1+\theta})^{T/2} \quad (\gamma > \theta), \end{aligned} \quad (4)$$

where the second inequality holds if each weak classifier attains an error smaller than  $1/2 - \gamma$ ,  $\gamma > \theta$ . This bound decreases to zero exponentially fast if  $\gamma > \theta$ . In other words, if a sufficiently large value of  $\gamma$  can be guaranteed, then the first term in (1) converges to zero. However, if  $\gamma$  (and thus  $\theta$ ) behaves like  $m^{-\beta}$  for some  $\beta > 0$ , the rate of convergence in the second term in (1) will increase, possibly leading to worse bounds than those available by using standard VC results (2). What is needed then is a characterization of conditions under which the achievable  $\theta$  does not decrease too rapidly with  $m$ . For example, if  $\gamma = \Omega(1/\log m)$  this will only contribute a logarithmic factor to the complexity penalty, and not change the rate which will remain  $O(1/\sqrt{m})$ , up to logarithmic factors. On the other hand, we have shown in Mannor and Meir (2001) that a value of  $\gamma$  of order  $1/m$  can always be achieved; however, such a small value is useless from the point of view of the bound (1). While the requirement that  $\gamma > \theta$  is stronger than we actually need, finding more refined conditions will be left to future work.

We introduce the notion of an *effective* weak learner, which characterizes a weak learner for which  $\gamma$  is sufficiently large to guarantee that the second term in (1) converges to zero as the sample size increases. In the following definition we use the notation  $f(m) = \omega(g(m))$  to indicate that  $f(m)$  becomes arbitrarily large relative to  $g(m)$  when  $m$  approaches infinity.

*Definition 3.* A weak hypothesis class  $\mathcal{H}$  is effective for a set of points  $S$  if for any  $P$  there exists  $h \in \mathcal{H}$  such that (3) is satisfied with  $\gamma > \omega(\log m/\sqrt{m})$ .

In this work we consider conditions under which an effective weak learner exists for one of the simplest possible classifiers, namely the linear classifier. A particularly interesting situation occurs when the data can be split into a small number of ‘simple’ regions. For the simplest case we now provide an intuitive demonstration without resorting to complex tools. Let  $P_+/P_-$  denote the total weight of the positive/negative examples, respectively. We will show in Lemma 4.1 that it suffices to consider the case  $P_+ = P_- = 1/2$ . Consider the problem of a convex set  $B \subset \mathbb{R}^d$  that contains all the positively labeled data points, such that all the negatively labeled points are contained in  $\mathbb{R}^d \setminus B^\eta$ , where  $B^\eta$  is the  $\eta$ -expansion of  $B$ —see figure 2. In this case, it is easy to show that an effective linear weak learner

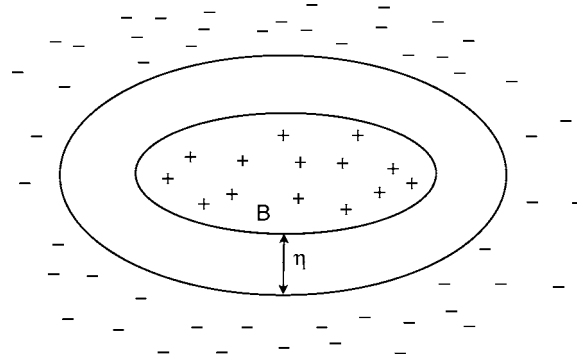


Figure 2. A single convex set containing the positively labeled examples separated from the negative examples by a gap of  $\eta$ .

exists. Moreover, the value of  $\gamma$  is independent of the number of points  $m$ . To see this, construct a convex polytope containing  $B$  (and thus all the positively labeled examples) and contained in  $B^\eta$ . Consider now the set of linear classifiers  $\{\xi_1, \dots, \xi_q\}$  defined by the  $q$  faces of the polytope, such that  $\xi_i = +1$  on the halfspace containing the polytope and  $-1$  otherwise. Clearly each of these classifiers provides the correct classification for the positively labeled points, erring on some of the negatively labeled points. Consider the (non-linear) classifier  $\xi$  defined by  $\min(\xi_1, \dots, \xi_q)$ , which by construction of the polytope has zero error. Since  $\xi$  classifies all the negative examples correctly, and since there are  $q$  linear classifiers  $\{\xi_1, \dots, \xi_q\}$ , there exists a linear classifier (say  $\xi_1$ ) for which the weighted success ratio on the negative examples  $\sum_{i=1}^m P_i I(y_i = \xi_1(\mathbf{x}_i)) I(y_i = -1) \geq P_-/q = 1/2q$ . Thus, the weighted error of  $\xi_1$  on the negative examples is smaller than  $1/2 - 1/2q$ . Since  $\xi_1$  classifies all the positive examples correctly, the total weighted misclassification error of  $\xi_1$  is smaller than  $1/2 - 1/2q$ , which is independent of the number of points  $m$ . Note that while the current example can be addressed using elementary means, we do not see how to extend it to more complex situations, which require the machinery of geometric discrepancy reviewed in Section 3.

### 3. Background from geometric discrepancy

We introduce the key definitions and results which will enable us to show in Section 4 that an *effective* weak learner exists. The results of this section essentially summarize the work of Alexander (1990, 1991, 1994), setting the nomenclature and drawing the connections to classification and the notion of weak learners. We start with several definitions. Let  $h$  denote a hyperplane

$$h = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{w} \cdot \mathbf{x} + b = 0\}.$$

The corresponding (open) half-spaces  $h^+$  and  $h^-$  are given, respectively, by  $h^+ = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{w} \cdot \mathbf{x} + b > 0\}$  and  $h^- = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{w} \cdot \mathbf{x} + b < 0\}$ . We denote by  $H$  a generic half-space.

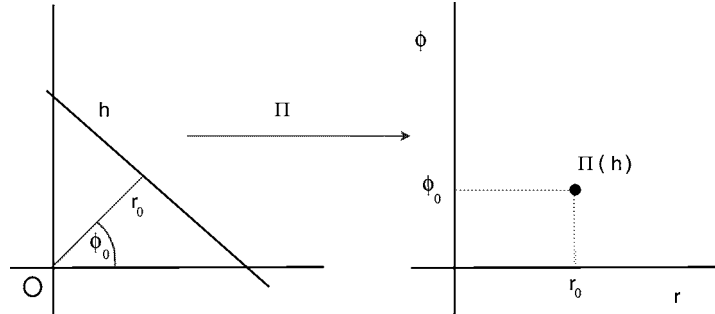


Figure 3. Mapping from hyperplanes to spherical coordinates (based on Matoušek, 1999).

The linear classifier associated with the hyperplane  $h$  will be denoted by  $\xi_h$ , namely

$$\xi_h(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x} + b).$$

An essential ingredient of the proof of our main result is the notion of a *motion invariant measure* on the space of hyperplanes in  $\mathbb{R}^d$ . We follow Section 6.4 in Matoušek (1999) in describing this measure. First, we ignore the set of hyperplanes passing through the origin, which has measure zero. Second, note that any hyperplane  $h$  is *uniquely* determined by  $r$ , the distance from the origin, and  $\Omega$ , the  $d$ -dimensional angle describing the inclination of the shortest line from the origin to the hyperplane (see figure 3 for the case  $d = 2$ , where  $\Omega$  is simply the polar angle  $\phi$ ,  $0 \leq \phi \leq 2\pi$ ). Let  $\Pi$  be the bijective map from the set of hyperplanes to the  $d$ -dimensional Euclidean space  $(r, \Omega)$ . If  $L$  is a set of hyperplanes, we define the measure  $\mu(L)$  as the Lebesgue measure of the set  $\Pi(L)$  in the  $(r, \Omega)$  space, where  $r$  and  $\Omega$  are interpreted as cartesian coordinates. For example, if the set  $L$  is supported in the unit ball in  $\mathbb{R}^2$ , then  $\Pi(L) = \{(r, \phi), 0 \leq r \leq 1, 0 \leq \phi \leq 2\pi\}$ . It is easy to see that the measure  $\mu$  is *motion-invariant*, i.e., if  $L'$  arises from  $L$  by rigid motion (i.e., without changing the relative positions of the lines) then  $\mu(L') = \mu(L)$ . In fact, it turns out that, up to scaling, this is the only motion-invariant measure on hyperplanes (see Section 6.4 in Matoušek, 1999).

Given a set of points  $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$  in  $\mathbb{R}^d$ , and a signed atomic measure  $\nu$ , the geometric discrepancy  $D_s(\nu)$  is given by

$$D_s(\nu) \triangleq \sup_H \{\nu(H)\} = \sup_H \left\{ \sum_{\mathbf{x}_i \in X \cap H} \nu(\mathbf{x}_i) \right\},$$

where the supremum is over all halfspaces.

Intuitively, the discrepancy relates to the maximal measure of points in  $X$ , which belong to a half-space. Consider, for example, the case of an even number of points, such that  $v_i = 1, i = 1, 2, \dots, m/2$ , and  $v_i = -1$  for  $i = m/2 + 1, \dots, m$ , where  $v_i = \nu(\mathbf{x}_i)$ . Then,  $D_s(\nu)$  is the maximal difference between the sizes of positive and negative sets belonging

to any halfspace. In the sequel we will be interested in finding a *lower* bound on  $D_s(\nu)$ . The main motivation for the introduction of  $D_s(\nu)$  is its relationship to the weighted empirical error defined in (3). Let  $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $y_i \in \{-1, +1\}$ , be a labeled sample of  $m$  points weighted according to  $\{P_1, \dots, P_m\}$ ,  $\sum_{i=1}^m P_i = 1$ . Set

$$v_i \triangleq y_i P_i \quad (i = 1, 2, \dots, m), \quad (5)$$

a signed measure over the points  $X$ . We will observe in Section 4 that when  $\sum_i v_i = 0$ ,

$$\inf_{\xi} \left\{ \sum_{i=1}^m P_i I(\xi(\mathbf{x}_i) \neq y_i) \right\} = \frac{1}{2} - D_s(\nu) \quad (v_i = y_i P_i) \quad (6)$$

where the infimum is over all linear classifiers  $\xi$ .

Furthermore, it follows from Lemma 4.1 that if a lower bound on  $D_s(\nu)$  is available for any measure  $\nu$  such that  $\sum_{i=1}^m v_i = 0$ , then a lower bound can be obtained for any signed measure  $\nu$ . For this purpose we define a set of measures obeying the condition,

$$\Psi(\mathbb{R}^d) = \{ \nu : \nu \text{ is a signed measure such that } |\nu|(\mathbb{R}^d) < \infty \text{ and } \nu(\mathbb{R}^d) = 0 \}.$$

Let  $K$  be the support of  $\nu$ . Then for any  $\nu \in \Psi(\mathbb{R}^d)$ ,  $\nu(K) = 0$ . In this work  $K$  will be a finite set,  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ . Furthermore, since  $P$  is a normalized probability distribution, clearly we must also demand that  $\sum_{i=1}^m |v_i| = 1$ . Thus, we make the following assumption throughout this section.

*Assumption 3.1.* The signed measure  $\nu$  obeys the two conditions

$$\sum_{i=1}^m |v_i| = 1 \quad \text{and} \quad \sum_{i=1}^m v_i = 0.$$

A crucial quantity used in the sequel is the function

$$I(\nu) \triangleq \sum_{i,j} \|\mathbf{x}_i - \mathbf{x}_j\|_2 v_i v_j \quad \left( \sum_i v_i = 0 \right),$$

which is a weighted measure of dispersion of the points  $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ . This function turns out to be very useful for our purposes since it is connected to the discrepancy, through the following result.

**Lemma 3.1** (Alexander, 1975, Proposition 3.1). *Let  $\nu$  belong to  $\Psi(\mathbb{R}^d)$ , and denote by  $\mu$  the unique motion invariant measure over hyperplanes. Then*

$$-I(\nu) = \int v(h^+)^2 d\mu(h) \leq D_s(\nu)^2 \mu(\mathbf{H}) \quad (7)$$

where  $\mathbf{H}$  is the set of planes that cut the convex hull of the support of  $\nu$ .

Two important observations follow from (7).

1. The function  $I(v)$  is non-positive. This is in itself a surprising and non-trivial result, following from the classic work of Schoenberg (1937) on isometric embeddings in Hilbert space. The reader is referred to Section 6.7 in Matoušek (1999) for a more modern proof of this fact using the notion of positive definite functions.
2. The resulting inequality

$$D_s(v) \geq \sqrt{\frac{-I(v)}{\mu(\mathbf{H})}} \tag{8}$$

immediately yields a lower bound on  $D_s(v)$  if a lower bound for  $-I(v)$  and an *upper* bound for  $\mu(\mathbf{H})$  are available. An exact value for  $\mu(\mathbf{H})$  is provided by Alexander for hyper-cubes and hyper-spheres (Alexander, 1994), and can be immediately used in our case. For example, from Lemma 3 in Alexander (1994) we find that when  $\mathbf{H}$  is a ball of radius  $r$  in  $\mathbb{R}^d$  (namely, the points  $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$  are contained in a ball of radius  $r$ ),

$$\mu(\mathbf{H}) = 2r(d-1) \frac{O_{d-1}}{O_{d-2}} = 2\pi r(d-1) \frac{\Gamma(\frac{d-1}{2})}{\Gamma(\frac{d}{2})},$$

where  $O_k = 2\pi^{(k+1)/2} (\Gamma((k+1)/2))^{-1}$  is the volume of the unit ball. Note that for large  $d$ ,  $\mu(\mathbf{H})$  behaves like  $1/\sqrt{d}$ . We further comment that Alexander’s Lemma discusses the unit ball. However, the extension to any ball is straightforward using results from Section 13.6 in Santaló (1976). Without loss of generality, we may always assume that the set  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$  is contained in a ball of radius  $\max_{i,j} \|\mathbf{x}_i - \mathbf{x}_j\|/2$ .

One final ingredient needed in order to prove our main result is the construction of an auxiliary measure in  $\mathbb{R}^{d+1}$ . The main idea here, due to Alexander (1990), is to view  $\mathbb{R}^d$  as a hyper-plane in  $\mathbb{R}^{d+1}$ , and construct a measure in  $\mathbb{R}^{d+1}$  for which a lower bound on  $D_s(v)$  can be established. We briefly describe the construction and properties of the new measure, referring the reader to Alexander (1990, 1991, 1994) for the full details. Let  $\Phi$  be an atomic measure over  $\mathbb{R}$ , concentrated on the finite set  $R = \{r_1, r_2, \dots, r_n\}$ . For  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ , consider the set  $X \times R \subset \mathbb{R}^{d+1}$ , and define the convolution measure  $v \star \Phi$  over  $X \times R$  by

$$(v \star \Phi)(\mathbf{x}_i, r_k) \triangleq v(\mathbf{x}_i)\Phi(r_k) = v_i\phi_k.$$

From Corollary 3 in Alexander (1991) we conclude that

$$-I(v \star \Phi) \leq -\|\Phi\|_1^2 I(v),$$

where  $\|\Phi\|_1 = \sum_{k=1}^{\ell} |\phi_k|$ . Finally, let  $\mathbf{q}_{ik} = (\mathbf{x}_i^T, r_k)^T \in \mathbb{R}^{d+1}$  be a point in  $X \times R$ . From Lemma 9 in Alexander (1990) we have that

$$-I(v \star \Phi) = -I(\Phi) \sum_i v_i^2 - \sum_{i \neq j} \sum J_{\Phi}(\mathbf{x}_i, \mathbf{x}_j) v_i v_j, \tag{9}$$

where

$$I(\Phi) = \sum_{k=1}^n \sum_{l=1}^n |r_k - r_l| \phi_k \phi_l, \quad (10)$$

$$J_\Phi(\mathbf{x}_i, \mathbf{x}_j) \triangleq \sum_{k=1}^n \sum_{l=1}^n \|\mathbf{q}_{ik} - \mathbf{q}_{jl}\|_2 \phi_k \phi_l. \quad (11)$$

It is important to observe that since  $\|\mathbf{q}_{ik} - \mathbf{q}_{jl}\|_2^2 = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 + |r_k - r_l|^2$ , the function  $J_\Phi(\mathbf{x}_i, \mathbf{x}_j)$  depends on  $\mathbf{x}_i$  and  $\mathbf{x}_j$  only through  $\rho_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$ . In Section 4 we use (9) and (10) as a starting point for our proof. We quote two essential features of the function  $J_\Phi(\mathbf{x}_i, \mathbf{x}_j)$ .

**Lemma 3.2** (Alexander, 1991, Theorem 4 and Lemma 7). *Let  $\Phi$  be an atomic measure concentrated on the  $n$  points  $\{r_1, \dots, r_n\}$  in  $\mathbb{R}$ , and such that  $\sum_{i=1}^n \Phi(r_i) = 0$ . Assume further that  $\rho = \|\mathbf{x} - \mathbf{x}'\|_2$  exceeds  $\delta$ , the diameter of the support of  $\Phi$ . Then*

$$J_\Phi(\mathbf{x}, \mathbf{x}') = \rho \sum_{k=1}^{\infty} c_k I^{2k}(\Phi) \rho^{-2k},$$

where  $c_k$  is defined by  $(1 + x^2)^{1/2} = \sum_{k=0}^{\infty} c_k x^{2k}$ , and

$$I^\ell(\Phi) \triangleq \sum_{i=1}^n \sum_{j=1}^n |r_i - r_j|^\ell \Phi(r_i) \Phi(r_j). \quad (12)$$

Moreover,

$$I^\ell(\Phi) \leq \|\Phi\|_1^2 \delta^\ell. \quad (13)$$

**Lemma 3.3** (Alexander, 1991, Theorem 6). *Let  $\Phi$  be an atomic measure concentrated on the  $n$  points  $\{r_1, \dots, r_n\}$  in  $\mathbb{R}$ , and such that  $\sum_{i=1}^n \Phi(r_i) = 0$ . Then  $-J_\Phi(\mathbf{x}, \mathbf{x}')$  is a strictly decreasing positive function of  $\rho = \|\mathbf{x} - \mathbf{x}'\|_2$ .*

#### 4. On the existence of linear weak learners

With the basic mathematical tools of Section 3 in hand, we proceed now to the proof of the existence of an effective weak learner based on hyperplanar decision boundaries. A few words are in order concerning the construction of linear classifiers for arbitrary data sets, which are, of course, not linearly separable in general. When the points are linearly separable, it is known that there are effective (i.e., polynomial time) procedures for determining the separating hyperplane. For example, some algorithms for linear programming achieve the desired goal. When the points are *not* linearly separable, it has been proved that the problem of finding the linear classifier with the minimal number of errors is NP-hard (Johnson & Preparata, 1978); see also Section 24.2 in Anthony and Bartlett (1999). Recently, it

was shown in Bartlett and Ben-David (1999) that the problem is NP-hard even if only an approximate solution is required. Note, however, that our problem is somewhat different, as we are interested only in finding, for any distribution on the data points, a ‘weak’ solution, i.e., one that is always slightly better than the trivial error of  $1/2$ . Our main focus here is on the *existence* of such a linear weak learner. We proceed now to the mathematical analysis of the problem.

Let

$$\xi(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x} + b),$$

be the classification according to the hyper-plane  $h$  defined by  $h = \{\mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = 0\}$ .

Let  $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $y_i \in \{-1, +1\}$ , be a sample of  $m$  points. We need to establish an upper bound on the error

$$\epsilon(P, \xi) \triangleq \sum_{i=1}^m P_i I(y_i \neq \xi(\mathbf{x}_i)) \quad (14)$$

where  $P_i \geq 0$  and  $\sum_i P_i = 1$ . Observe that  $\epsilon(P, \xi)$  may also be expressed as  $\sum_{i=1}^m P_i [1 - y_i \xi(\mathbf{x}_i)]/2$ .

First, we wish to show that unless some conditions are specified, there can be no effective linear classifiers. To show this consider the following example, which demonstrates that, for any classifier  $\xi$ , there exists a set of points such that  $\epsilon(P, \xi) \geq 1/2 - \gamma$ , where  $\gamma \leq O(1/m)$ .

*Example 1.* Consider the points  $\mathbf{x}_1, \dots, \mathbf{x}_{2m}$  on the boundary of a circle, and set  $v_i = (-1)^i (2m)^{-1}$ . Clearly  $\sum_i v_i = 0$  and  $\sum_i |v_i| = 1$ . For example, in figure 4 there are 8 grey (positive) points and 8 white (negative) points alternating on the circumference of a circle. From the symmetry of the problem, it is clear that any line cutting the circle (and not passing through any of the points) is such that there is a difference of at most 1 between the number of points from either color on either side of the line. Clearly the error incurred by *any* hyperplane  $h$  is at least  $1/2 - 1/m$ , and cannot fulfill the condition required from an effective weak learner (recall Definition 3).

A more general proof of the non-existence of effective linear weak learners follows from Matoušek (1995). Let  $U$  denote the uniform distribution over a set of  $m$  points, namely  $P_i = 1/m$ ,  $i = 1, 2, \dots, m$ . Then Matoušek (1995) showed that for *any* set of points there exists a dichotomy such that

$$\epsilon(U, h) \geq \frac{1}{2} - \frac{c}{m^{1/2+1/2d}}, \quad (15)$$

for any linear classifier  $h$ , where the positive constant  $c$  does not depend on  $m$ . Note that in our case, the dichotomy is predefined. Our goal is to find conditions on the points which guarantee the existence of an effective linear weak learner.

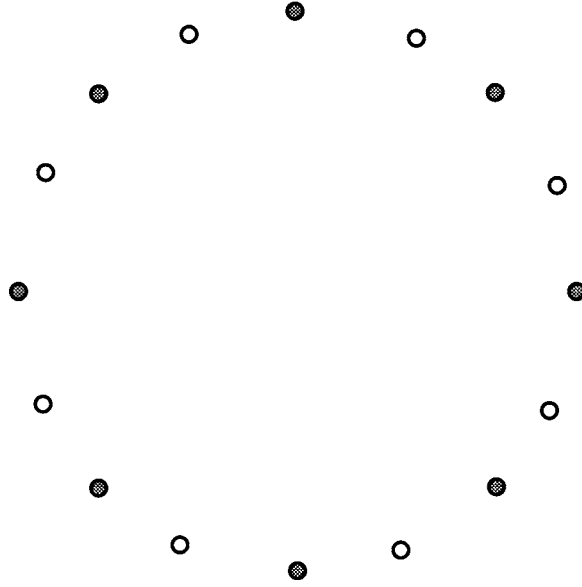


Figure 4. Sixteen alternating points on a sphere. No line can achieve an error lower than  $1/2 - 1/8 = 3/8$ .

Recall the auxiliary signed measure  $\nu_i = P_i y_i$  defined in (5). For any classifier  $\xi$  let

$$\tilde{D}(\nu, \xi) = \sum_{i=1}^m \nu_i \xi(\mathbf{x}_i) \quad (\nu_i = y_i P_i),$$

and note that

$$\epsilon(P, \xi) = \frac{1 - \tilde{D}(\nu, \xi)}{2} \quad (\nu_i = y_i P_i).$$

Let  $\mathcal{I}^\pm$  be the subsets of indices for which  $y_i = \pm 1$ . A probability distribution is said to be *symmetric*, with respect to a sample  $S$ , if

$$\sum_{i \in \mathcal{I}^+} P_i = \sum_{i \in \mathcal{I}^-} P_i, \tag{16}$$

i.e.,  $P$  assigns equal weights to positive and negative examples. We show in Lemma 4.1 that it suffices for our purposes to consider symmetric distributions. Observe that the symmetry and normalization conditions on  $P$  are equivalent to the conditions  $\sum_{i=1}^m \nu_i = 0$  and  $\sum_{i=1}^m |\nu_i| = 1$ , which constitute Assumption 3.1.

For any halfspace  $h^+ = \{\mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b > 0\}$  and points  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ , let

$$D(v, h^+) \triangleq v(h^+) = \sum_{\mathbf{x}_i \in X \cap h^+} v_i = \sum_{\mathbf{x}_i \in X} \left( \frac{1 + \xi_h(\mathbf{x}_i)}{2} \right) v_i,$$

where  $v_i = v(\mathbf{x}_i)$  and  $\xi_h(\mathbf{x}) = 2I(\mathbf{x} \in h^+) - 1$ .

Under the assumption that  $\sum_{\mathbf{x}_i \in X} v_i = 0$ , it follows that

$$\tilde{D}(v, \xi_h) = 2D(v, h^+).$$

We conclude that

$$\epsilon(P, \xi_h) = \frac{1}{2} - D(v, h^+),$$

where  $v_i = y_i P_i$  and  $\sum_{i=1}^m v_i = 0$ . Therefore, in order to derive an *upper* bound on  $\inf_{\xi} \epsilon(P, \xi)$  it suffices to obtain a *lower* bound on the geometric discrepancy  $D_s(v) = \sup_h D(v, h^+)$ , as claimed in (6).

We now show that it suffices to consider the case  $\sum_i v_i = 0$  (equivalently,  $\sum_{i \in \mathcal{I}^+} P_i = \sum_{i \in \mathcal{I}^-} P_i$ ). The proof of the following Lemma is given in the Appendix.

**Lemma 4.1.** *Let  $m$  arbitrary, but distinct, points  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$  in  $\mathbb{R}^d$  be given. Let  $\Xi$  be a class of binary classifiers, containing the constant classifiers  $\xi(\mathbf{x}) = +1$  and  $\xi(\mathbf{x}) = -1$ . Assume that for any symmetric probability distribution  $P$ , a classifier  $\xi \in \Xi$  can be found such that*

$$\epsilon(P, \xi) \leq \frac{1}{2} - \gamma, \tag{17}$$

with  $\gamma$  depending only on the points  $\mathbf{x}_1, \dots, \mathbf{x}_m$ . Then for any probability distribution (not necessarily symmetric) on  $X$ , there exists a classifier  $\xi \in \Xi$  for which

$$\epsilon(P, \xi) \leq \frac{1}{2} - \frac{\gamma}{2}.$$

Note that Lemma 4.1 is more general than we need, as it applies to any class of classifiers  $\Xi$  that contains the constant classifiers, rather than the class of linear classifiers, which is our main concern.

#### 4.1. Main theorem

The main result of this paper is the demonstration that even though effective linear weak learners do not exist in general (see (15)), one can find natural conditions under which they do exist. We begin with the more general theorem, specializing to some specific cases in two corollaries.

For a given set of points  $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\} \in (\mathbb{R}^d \times \{-1, +1\})^m$  denote by  $X^+$  the subset of positively labeled points, and similarly for  $X^-$ . Let

$$\eta \triangleq \min_{i,j} \{\|\mathbf{x}_i - \mathbf{x}_j\|_2, \mathbf{x}_i \in X^+, \mathbf{x}_j \in X^-\},$$

$$L \triangleq \max_{i,j} \{\|\mathbf{x}_i - \mathbf{x}_j\|_2, \mathbf{x}_i, \mathbf{x}_j \in X^+ \cup X^-\}.$$

Here  $\eta$  measures the smallest Euclidean distance between oppositely labeled points, while  $L$  is the diameter of the set of all points. We refer to  $\eta$  as the *gap*. The motivation for introducing  $\eta$  is the expectation that small error may be achieved for configurations in which the oppositely labeled points are well separated. This gap should be contrasted with the parameter  $\gamma$  in Definition 2, which characterizes a classifier, rather than a set of points (the points themselves enter through (3)).

In addition to the characterization of a data set in terms of the gap  $\eta$ , we introduce a further geometric parameter related to the homogeneity of the data. Let  $X^+/X^-$ , the sets of positively/negatively labeled points, be partitioned into  $K_+/K_-$  disjoint (not necessarily convex) subsets. The intuition is that if either  $K_+$  or  $K_-$  is small, and the convex hulls of  $\{X_k^+\}$  or  $\{X_k^-\}$  do not contain oppositely labeled points, a relatively low error can be achieved by some hyperplane. For example, in the case of figure 2,  $K_+ = K_- = 1$  and the convex hull of  $X^+$  contains no negatively labeled points, leading to low error. We discuss this issue further in Remark 2 below.

We further introduce an auxiliary function that plays a major part in the proof of Theorem 4.1 below.

$$\rho^{-1} F_n(\rho, \eta) \triangleq -\binom{2n}{n} - 2 \sum_{\ell=1}^n (-1)^\ell (1 + \ell^2 (\eta/4n\rho)^2)^{1/2} \binom{2n}{n-\ell}. \quad (18)$$

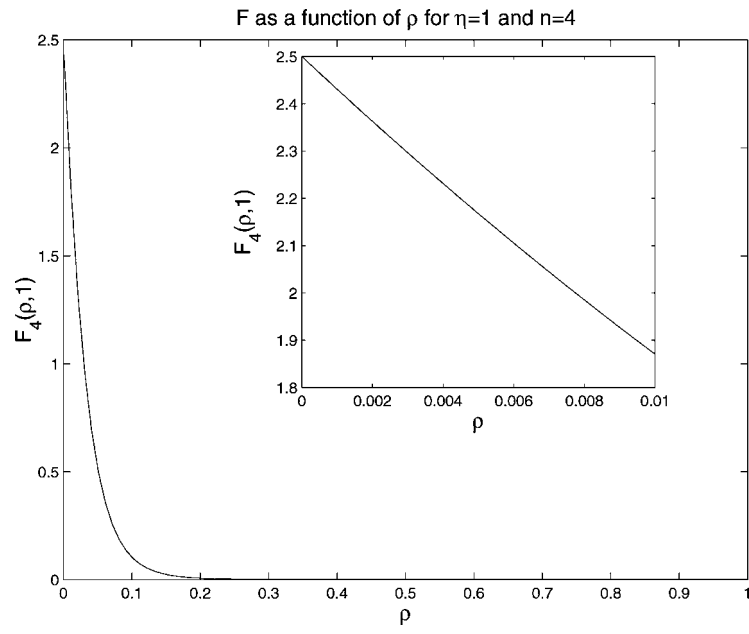
This function is in fact given by  $-J_\Phi(\mathbf{x}, \mathbf{x}')$ , where  $\|\mathbf{x} - \mathbf{x}'\|_2 = \rho$ , for a specific choice of the signed measure  $\Phi$  (see (25) in the appendix). We argue in Lemma C.2 that  $F_n(\rho, \eta)$  is a positive and monotonically decreasing function of  $\rho$ . Further properties of  $F_n(\rho, \eta)$  are also provided in Lemma C.2. In figure 5 we plot  $F_n(\rho, \eta)$  as a function of  $\eta$  for  $\rho = 2$  and  $n = 4$ , and as a function of  $\rho$  for  $\eta = 1$  and  $n = 4$ .

Before introducing the main theorem, we formally define a partition of the sets  $X^+$  and  $X^-$ .

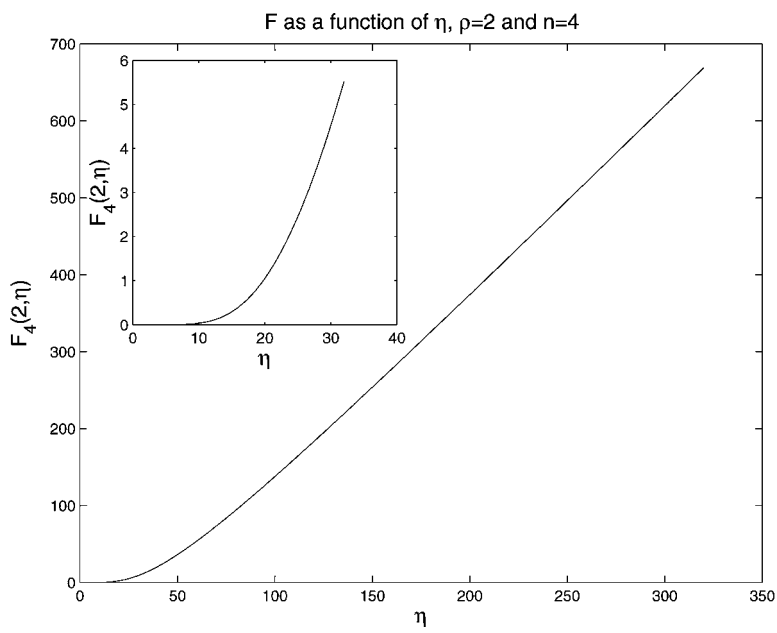
$$\begin{aligned} X^+ &= X_1^+ \cup \dots \cup X_{K_+}^+ & (X_i^+ \cap X_j^+ &= \emptyset, i \neq j) \\ X^- &= X_1^- \cup \dots \cup X_{K_-}^- & (X_i^- \cap X_j^- &= \emptyset, i \neq j) \end{aligned} \quad (19)$$

Furthermore, denote by  $\mathcal{I}_k^+$  the indices corresponding to points in  $X_k^+$ , and similarly for  $\mathcal{I}_k^-$ . We are now ready to state our main result.

**Theorem 4.1.** *Let  $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $y_i \in \{\pm 1\}$  be a fixed sample such that  $|X^\pm| = m_\pm$ , and let non-negative weights  $P_i$  be assigned to each point in  $S$  so*



(a)



(b)

Figure 5.  $F_4(\rho, \eta)$  as a function of  $\rho$  (for fixed  $\eta$ ) and as a function of  $\eta$  (for fixed  $\rho$ ).

that  $\sum_{i=1}^m P_i = 1$ . Let the points  $X = X^+ \cup X^-$  be partitioned as in (19). Then there exists a linear classifier  $\xi$  such that

$$\epsilon(P, \xi) \leq \frac{1}{2} - \gamma$$

where

$$\gamma \geq \sqrt{-(C_d/4L)I(v)} \quad (v_i = y_i P_i), \quad (20)$$

and where  $C_d = \Gamma(d/2)/[\pi(d-1)\Gamma((d-1)/2)]$  depends only on the dimension  $d$ , and

$$\begin{aligned} -I(v) \geq & 0.004\eta \left[ \log_2(492\sqrt{m_+m_-}) \right]^{-3/2} \sum_{i=1}^m v_i^2 \\ & + \sum_{k=1}^{K_+} F_n(\rho_k^+, \eta) \sum_{\substack{i \neq j \\ i, j \in \mathcal{I}_k^+}} |v_i||v_j| + \sum_{k=1}^{K_-} F_n(\rho_k^-, \eta) \sum_{\substack{i \neq j \\ i, j \in \mathcal{I}_k^-}} |v_i||v_j|. \end{aligned} \quad (21)$$

Here  $\rho_k^\pm$  is the diameter of  $X_k^\pm$ , and  $n = \lceil (1/2) \log_2(492\sqrt{m_+m_-}) \rceil$ .

A note is in order concerning the dependence of (20) on the dimension  $d$ . Since the r.h.s. of (21) does not depend on the dimension, it follows that  $\gamma$  depends on  $d$  only through  $C_d$ , which behaves like  $1/\sqrt{d}$  for large  $d$  (see the discussion following Lemma 3.1). It is possible to remove the logarithmic factor in  $m$  from (21), but at the price of a dramatic decrease in the constant  $C_d$ . For further discussion of this issue, refer to Alexander (1994).

*Remark 1.* As can be expected, the lower bound improves for larger values of the gap parameter  $\eta$ . A careful inspection of the bound shows that it depends only on  $\eta/L$ , rather than  $\eta$  or  $L$  individually. This fact arises from the scale invariance of the problem.

*Remark 2.* The statement of the theorem requires the partition of the sets  $X^\pm$  into  $K_\pm$  subsets. Clearly there are many such partitions, for values of  $K_\pm$  between 1 (corresponding to a single cluster) and  $m_\pm$  (corresponding to a single cluster per data point). A careful inspection of the second and third terms in (21) implies that there is a trade-off to be considered. First, observe that if one of the subsets of  $X^\pm$  is a singleton, the corresponding second or third term in (21) is absent. Next, assume that the number of regions  $K_\pm$ , is large, and that each subset  $X_k^\pm$  is composed of a relatively small set of nearby points. In this case, the number of summands is large, and the functions  $F_n(\rho_k^\pm, \eta)$  are also large since  $\rho_k^\pm$  are small and the function  $F_n(\rho, \eta)$  is monotonically decreasing with  $\rho$ . On the other hand, the terms  $\sum_{i \neq j} |v_i||v_j|$  are small, since they contain a small number of terms. In the extreme case where the size of the partition is equal to the number of points, the second and third terms in (21) are entirely absent. Consider now the other extreme situation, where  $K_\pm = 1$ , i.e., no partition of the points is constructed. In this case, the number of summands in the final two terms in (21) is reduced, and  $F_n(\rho_k^\pm, \eta) \approx F_n(L, \eta)$  is also smaller. On the

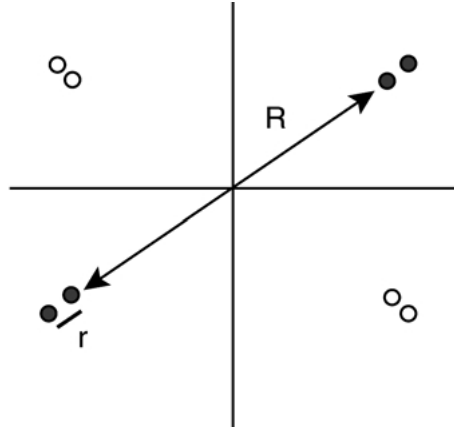


Figure 6. 8 points in a XOR position,  $R$  is the maximal distance between equally labeled points, and  $r$  is the distance between equally labeled points within a localized region.

other hand, the terms  $\sum_{i \neq j} |v_i||v_j|$  are large, since they contain a large number of terms. Thus, the optimal partition of the sets  $X^\pm$  depends in each case on the specific geometry of the points, as well as on their assigned measures  $v_i$ . For example, we show in example 2 below, a clear case where it is advantageous to choose  $K_\pm = 2$  rather than  $K_\pm = 1$ .

*Example 2.* We give an example where a better bound may be obtained using a larger partition of the set  $X$ . Consider the XOR configuration of 8 equally weighted points in figure 6. Since the first term in the bound (21) is not affected by the partition we compare only the final two terms, for the case  $K_\pm = 1$  and  $K_\pm = 2$ , where in the latter case the partition is into the separate four quadrants. Simple algebra leads to the following two bounds:

$$\begin{aligned}
 -I(v) &\geq \frac{3}{8} F_n(R, \eta) \quad (K_\pm = 1) \\
 -I(v) &\geq \frac{1}{16} F_n(r, \eta) \quad (K_\pm = 2).
 \end{aligned}$$

From the first property in Lemma C.2,  $F_n(\rho, \eta)$  is positive and strictly decreasing function of  $\rho$ . By picking  $r$  small enough and  $R$  large enough, the bound (21) can be made larger for the case  $K_\pm = 2$ .

In order to obtain a better understanding of the structure of the bounds, we consider the special case where all the weights are equal in magnitude, namely  $|v_i| = 1/m$  for all  $i$ . We then have the following corollary.

**Corollary 4.1.** *Let the conditions of Theorem 4.1 hold, with the additional constraint that  $|v_i| = 1/m$  for all  $i = 1, 2, \dots, m$ . Denote the number of positively/negatively labeled*

points by  $m_{\pm} = \alpha_{\pm}m$ ,  $2 < m_+ < m - 2$ . Then there exists a linear classifier  $\xi$  such that

$$\epsilon(P, \xi) \leq \frac{1}{2} - \gamma$$

where

$$\gamma \geq \sqrt{\frac{C_d}{4L}} \left\{ 0.063\sqrt{\eta}(\log_2(492\sqrt{\alpha_+\alpha_-}m))^{-3/4}m^{-1/2} + \frac{1}{2}\sqrt{F_n(L, \eta)} \left( \frac{\alpha_+^2}{K_+} + \frac{\alpha_-^2}{K_-} \right)^{1/2} \right\},$$

where  $C_d$  and  $n$  are defined in Theorem 4.1. For small values of  $\eta/L$  we find that  $F_n(L, \eta) \geq \Omega(m^{-2\log(2Le/\eta)}(\log m)^{-3/2})$ , implying that

$$\gamma = \Omega(\sqrt{\eta}(\log m)^{-3/4}m^{-1/2}) + \Omega\left((\log m)^{-3/4}m^{-\log(\frac{2Le}{\eta})} \left[ \frac{\alpha_+}{\sqrt{K_+}} + \frac{\alpha_-}{K_-} \right]\right).$$

**Proof:** Consider the third term in (21). Since  $F_n(\rho, \eta)$  is monotonically decreasing with  $\rho$ , we may replace  $F_n(\rho_k^+, \eta)$  by  $F_n(L, \eta)$ . We also assume for simplicity that  $m_{+,k} \geq 2$  and similarly for  $m_{-,k}$ . Setting  $|v_i|$  to  $1/m$  we see that

$$\begin{aligned} \sum_{k=1}^{K_+} F_n(\rho_k^+, \eta) \sum_{\substack{i \neq j \\ i, j \in \mathcal{I}_k^+}} |v_i||v_j| &\geq \frac{1}{m^2} F_n(L, \eta) \sum_{k=1}^{K_+} m_{+,k}(m_{+,k} - 1) \\ &\geq \frac{1}{2m^2} F_n(L, \eta) \sum_{k=1}^{K_+} m_{+,k}^2 \\ &\geq \frac{1}{2K_+} F_n(L, \eta) \alpha_+^2, \end{aligned}$$

where we have used the Cauchy-Schwartz inequality in the final step. Using the same arguments for the final term in (21) we find that in this case

$$-I(v) \geq 0.004\eta[\log_2(492\sqrt{m_+m_-})]^{-3/2}m^{-1} + F_n(L, \eta) \left[ \frac{\alpha_+^2}{2K_+} + \frac{\alpha_-^2}{2K_-} \right].$$

Next, we use the inequality  $\sqrt{a+b} \geq (\sqrt{a} + \sqrt{b})/\sqrt{2}$ , which follows from Jensen's inequality and the concavity of the square root function, obtaining

$$\gamma \geq \sqrt{\frac{C_d}{8L}} \left\{ 0.063\sqrt{\eta} \log_2(492\sqrt{m_+m_-})^{-3/4}m^{-1/2} + \frac{\sqrt{F_n(L, \eta)}}{2} \left[ \frac{\alpha_+^2}{2K_+} + \frac{\alpha_-^2}{2K_-} \right]^{1/2} \right\}.$$

To complete the proof we need a lower bound on  $F_n(L, \eta)$ . This follows by using the results of Lemma C.2 to show that  $F_n(L, \eta) \geq \Omega((\eta/n)^{2n} c_n (2n!))$ , where  $c_n$  is the  $(n + 1)$ -th coefficient in the expansion of  $\sqrt{1 + x}$ , and noting that  $c_n \geq (n + 3)^{-2}$ . Finally, from the proof of Theorem 4.1 we recall that  $n \geq (1/2) \log_2(492\sqrt{m_+ m_-})$ .  $\square$

The main significance of Corollary 4.1 is the explicit dependence on the number of regions  $K_{\pm}$ . As can be expected, the bound improves when the number of regions is small. Note though that the more precise bound given in (21) requires a more delicate treatment of the geometry of regions, as discussed in Remark 2.

We comment on a hypothesis class which is widely used in many of the applications of boosting, namely the class of *stumps*. This class consists of linear classifiers for which the decision boundary is parallel to one of the axes. We show that under the conditions of Theorem 4.1 stumps cannot achieve an error lower than  $1/2$ , while this is clearly possible using general linear classifiers. Let  $\Xi$  be the class of axis parallel linear classifiers (stumps). Then we can show that there exists a set of points  $S$ , obeying the conditions of Theorem 4.1, for which  $\epsilon(P, \xi) = 1/2$  for any  $\xi \in \Xi$ . This can be seen by simply considering the set of points in figure 7, composed of a XOR configuration of equally weighted points. It is clear by symmetry that no axis-parallel line can achieve an error lower than  $1/2$ , for any choice of the distance between positive and negative points.

Finally, we comment on the necessity of a positive value for  $\eta$ . One can easily show that in the special case where the hyperplanes are homogeneous, i.e., constrained to pass through the origin, the positivity of  $\eta$  is essential. To see this, consider a spherically symmetric two-dimensional configuration of points, such that all points in the first and third quadrants are positively labeled, while points in the second and fourth quadrants are negatively labeled. Moreover, assume that  $\eta$ , the distance between the closest points from the two classes is vanishingly small. Then it is easy to see by symmetry that no homogeneous hyperplane can achieve an error smaller than  $1/2$ .

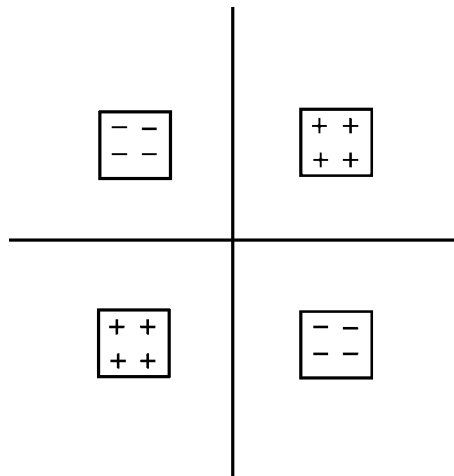


Figure 7. A configuration of points that cannot be classified by stumps with an error smaller than  $1/2$ .

4.2. Application to boosting

We now apply Theorem 4.1 to boosting algorithms, and in particular to the AdaBoost algorithm (Freund & Schapire, 1996a). The basic claims are derived by combining the generalization error bound (1) and the bound (21). First, note that the variable  $\theta$  in (1) is a free parameter, which may be tuned at will. We wish to consider situations where the bound converges to zero as  $m$  increases, and thus need to guarantee that both terms in (1) vanish for large  $m$ . Note that this behavior is required when the Bayes error is zero. Consider first the margin error  $\mathbf{P}_S[Yf(X) \leq \theta]$ . In view of (4) and the comments following it, it suffices that  $\theta < \gamma$ , where each weak learner achieves an error which is smaller than  $1/2 - \gamma$ . Retaining only the first term of the lower bound (21), we observe that it behaves like  $\Omega(\eta(\log m)^{-3/2} \sum_{i=1}^m v_i^2)$ . The worst possible situation occurs when all the weights are equal, namely  $|v_i| = 1/m$  for all  $i$ , in which case the lower bound on  $\gamma$  is (see (20))  $\Omega([\eta(\log m)^{-3/2} m^{-1}]^{1/2})$ . In this case, upon substitution in (1), we obtain that the second (complexity penalty) term does not decay to zero, due to the fact that the gap parameter  $\gamma$  decays to zero too quickly. However, it is well known that in many successful applications of boosting, only a small number of weights, say  $s$ , retain non-negligible values, while most weights shrink (exponentially fast) to zero. In this situation, keeping in mind the normalization condition  $\sum_{i=1}^m |v_i| = 1$ , the term  $\sum_{i=1}^m v_i^2$  is in fact of order  $1/s$ , instead of order  $1/m$ . As long as  $s = o(m)$ , the existence of an effective weak learner is established. In this situation we may in fact obtain much better rates of convergence, as attested by Theorem 4.2 below.

**Theorem 4.2.** *Let  $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $y_i \in \{\pm 1\}$  be a sample of  $m$  points drawn independently at random from a distribution  $D$  over  $\mathbb{R}^d \times \{-1, +1\}$ , and such that the minimal distance between differently labeled points in  $S$  is  $\eta$ . Assume that the AdaBoost algorithm has been run for  $T$  steps, and that there exists an integer  $T_0 \leq T$ , such that the probabilities  $P_i$  assigned to the points in  $S$  obey  $\sum_{i=1}^m P_i^2 \geq 1/s$  for  $t > T_0$ . Then for sufficiently large  $m$ , with probability at least  $1 - \delta$ ,*

$$\mathbf{P}_D[Yf(X) \leq 0] \leq c \exp\left\{-\frac{2\eta(1-\mu)(T-T_0)}{s(\log m)^{3/2}}\right\} + O\left(\frac{1}{\sqrt{m}} \left(\frac{d_{\mathcal{H}} s}{\eta \mu^2} \left(\log \frac{m}{d_{\mathcal{H}}}\right)^{7/2} + \log \frac{1}{\delta}\right)^{1/2}\right),$$

where the constant  $c$  depends on  $T_0$ , and  $\mu \in (0, 1)$  is a free parameter.

**Proof:** First, observe that the generalization bound (1) applies to all  $f \in \text{co}(\mathcal{H})$ , and in particular to the  $f$  obtained via the boosting algorithm. Second, under the conditions of the theorem, we know that there exists a constant  $a$  and a positive number  $\gamma_0 = [a\eta s^{-1}(\log m)^{-3/2}]^{1/2}$  such that each weak learner achieves an error smaller than  $1/2 - \gamma_0$  for  $T > T_0$ . Set  $\theta = \mu\gamma_0$ ,

$0 < \mu < 1$ . Then from (4), we find that the empirical error of the classifier  $f$  obeys

$$\mathbf{P}_S[Yf(X) \leq \theta] \leq \prod_{t=1}^{T_0} \sqrt{4\epsilon_t^{1-\theta} (1 - \epsilon_t)^{1+\theta}} \times \exp\left\{\frac{T - T_0}{2} [(1 - \mu\gamma_0) \ln(1 - 2\gamma_0) + (1 + \mu\gamma_0) \ln(1 + 2\gamma_0)]\right\},$$

where we have split the product (4) into  $T \leq T_0$  and  $T > T_0$ . For large values of  $m$ , i.e., small values of  $\gamma_0$ , we expand the logarithm, substitute for  $\gamma_0$  as a function of  $m$ , and obtain the desired claim for the first term in (1). The second term is obtained by straightforward substitution of  $\theta = \mu\gamma_0$  in the second term of (1).  $\square$

*Remark 3.* The bound of Theorem 4.2 is interesting as it relates the number of boosting iterations to the sample size and to the properties of the data (through  $\eta$ ). In particular, as long as the gap parameter  $\eta$  is strictly positive (independently of  $m$ ), and as long as

$$T = \omega((\log m)^{3/2} s / \eta)$$

we are guaranteed that the classification error will converge to zero. For example, Let  $T = \Omega((\log m)^{5/2} s / \eta)$ . Then we find that the bound behaves like

$$\mathbf{P}_D[Yf(X) \leq 0] \leq O\left(\frac{1}{m^{2\eta(1-\mu)/s}}\right) + O\left(\left(\frac{d_{\mathcal{H}} s}{\eta \mu^2}\right)^{1/2} \left(\frac{(\log(m/d_{\mathcal{H}}))^{7/2}}{m}\right)^{1/2}\right),$$

which decays to zero as a function of  $m$  at a rate depending on  $s$  and on the gap parameter  $\eta$ . It is important to observe that an appropriate value of  $T$  is determined *dynamically* here, and not preset in advance, since the condition that  $\sum_i P_i^2 \geq 1/s$  is only known after the algorithm is run.

It should be clear that Theorem 4.2 and Remark 3 may be used in order to select an appropriate value for the number of boosting iterations needed to guarantee that the total error bound approaches zero. As far as we aware there has been no systematic method to-date to address this problem. A particularly simple situation occurs when the lower bound on  $\gamma$  is independent of  $m$ , as in the example given at the end of Section 2. An argument very similar to the one above shows that the choice  $T = \Omega(1/\eta)$ , independently of  $m$ , suffices to guarantee that both the empirical error and the generalization error converge to zero.

We comment that Kearns and Mansour (1996) have shown that several widely used algorithms for the construction of top-down decision trees, are in fact boosting algorithms. Considering decision trees formed by oblique splits at each node, we conclude that our results establish bounds on the generalization error for decision trees as well.

Finally, we briefly address the computational problem of finding an effective linear weak learner. Recall from Lemma 3.1 that  $-I(v) = \int v(h^+)^2 d\mu(h)$ , where  $\mu$  is the motion-invariant measure and the integration is over all hyperplanes  $h$  cutting the convex hull of  $S$ . Since  $v$  is bounded and  $\mu$  possesses a finite second order moment (being a bounded measure

over a compact space), it follows from standard arguments (e.g., Motwani & Raghavan, 1995) that by randomly selecting  $h$  according to  $\mu$  an effective linear weak learner is found with arbitrarily high probability.

## 5. Conclusions

We have considered the existence of a weak learner in boosting. The existence of such a learner is a crucial requirement for the success of boosting algorithms, but there has been no general geometric proof to-date that such a learner exists. In this work we have shown that an effective linear weak learner exists, and have provided bounds on its performance. It turns out that the only condition needed for the existence of a weak linear learner is that the positively/negatively labeled points be separated by some nonzero gap. Combining our results with the standard generalization bounds for boosting, we are able to establish new bounds on the performance of boosting algorithms. Since the existence of a weak learner has been established for linear classifiers, the results clearly hold for more complex classifiers for which linear classifiers are a sub-class. Two notable examples are neural networks and decision trees with oblique splits. We also argued that our bounds can be directly used in order to determine the number of boosting iteration needed in order to guarantee convergence of the expected error to zero. Whether the bounds are sufficiently tight in order for this procedure to be practically useful remains to be seen.

It is important to stress that the only condition needed in order to guarantee the existence of a weak linear classifier is that  $\eta/L$  be strictly positive, where  $\eta$  is the minimal distance between oppositely labeled points and  $L$  is the size of the support of the data. In practical pattern recognition applications it may be hard to guarantee a sizable gap as the sample size increases. A challenging open problem here would be to investigate under what conditions a weak learner may be shown to exist without this assumption. We have argued that, at least for homogeneous hyperplanes, this requirement is necessary. In ongoing work we have been able to show that the requirement is also necessary for general hyperplanes, although the argument is more subtle. One possible approach to eliminating the requirement of a finite gap involves discarding oppositely labeled points from the original data set which are ‘too close’, thus effectively increasing the gap at the expense of an added error. If the number of discarded points is not large, we expect the bound to improve, as indeed verified by some recent unpublished results.

It is interesting to comment that the main result of this paper required some rather advanced tools from the field of combinatorial geometry, specifically the sub-field of geometric discrepancy. It would seem extremely difficult to establish the existence of effective weak learners using elementary techniques. An immediate question that arises here relates to the establishment of weak learnability for other types of classifiers, which are *not* based on hyperplanes. This issue is currently under investigation.

We have only briefly touched upon the issue of algorithmic design. As mentioned in Section 4.2 the construction of the motion-invariant measure of Section 3 immediately suggests a stochastic algorithm for this purpose. It would also be interesting to see whether deterministic algorithms exist, for which the property of weak learnability can be established.

Finally, we mention some recent work (Koltchinskii, Panchenko, & Lozano, 2001), which suggests non-trivial improvements to the generalization error bounds in (1). We are currently investigating how these improved bounds affect our results on weak learners.

### Appendix A: Proof of Lemma 4.1

**Proof:** Since

$$\epsilon(P, \xi) = \frac{1}{2} - \frac{1}{2} \sum_{i=1}^m v_i \xi(x_i),$$

it suffices to prove that if for every symmetric  $\nu$  there exists a  $\xi \in \Xi$  such that

$$\sum_{i=1}^m v_i \xi(x_i) > 2\epsilon \quad \left( \sum_i v_i = 0 \right),$$

then for any (not necessarily symmetric) distribution  $\nu$  there exists a classifier  $\xi \in \Xi$  such that

$$\sum_{i=1}^m v_i \xi(x_i) > \epsilon.$$

Let  $\nu$  be an arbitrary measure on  $X$ . If  $|\sum_{i=1}^m v_i| > \epsilon$  then one can take  $\xi$  to be the constant classifier  $\xi(\mathbf{x}) = \text{sgn}(\sum_{i=1}^m v_i)$  and the result follows. We therefore can assume, without loss of generality, that

$$\left| \sum_{i=1}^m v_i \right| \leq \epsilon.$$

Let us construct the following signed measure on  $X$ . Denote by  $\mathcal{I}^+$  the set of points for which  $v_i > 0$  and by  $\mathcal{I}^-$  the set of points for which  $v_i < 0$ , and let

$$v'_i = \begin{cases} \frac{1}{2} \frac{v_i}{\sum_{i \in \mathcal{I}^+} v_i} & \text{if } i \in \mathcal{I}^+, \\ -\frac{1}{2} \frac{v_i}{\sum_{i \in \mathcal{I}^-} v_i} & \text{if } i \in \mathcal{I}^-. \end{cases}$$

Note that  $v'_i \geq 0$  for  $i \in \mathcal{I}^\pm$ . Moreover,  $\sum_{i=1}^m v'_i = 0$ , so by assumption there exists a classifier  $\xi \in \Xi$  for which  $\sum_{i=1}^m v'_i \xi(x_i) > 2\epsilon$ . Clearly

$$\begin{aligned} \sum_{i=1}^m v_i \xi(x_i) &= \sum_{i=1}^m v'_i \xi(x_i) + \sum_{i=1}^m (v_i - v'_i) \xi(x_i) \\ &\geq \sum_{i=1}^m v'_i \xi(x_i) - \left| \sum_{i=1}^m (v_i - v'_i) \xi(x_i) \right|. \end{aligned} \quad (22)$$

By the construction of  $v'$  from  $v$  we can now bound the second term of (22):

$$\begin{aligned}
\left| \sum_{i=1}^m (v_i - v'_i) \xi(x_i) \right| &\leq \sum_{i=1}^m |v_i - v'_i| \\
&= \sum_{i \in \mathcal{I}^+} |v_i - v'_i| + \sum_{i \in \mathcal{I}^-} |v_i - v'_i| \\
&\stackrel{(a)}{=} \left| \sum_{i \in \mathcal{I}^+} (v_i - v'_i) \right| + \left| \sum_{i \in \mathcal{I}^-} (v_i - v'_i) \right| \\
&= \left| \sum_{i \in \mathcal{I}^+} v_i - \frac{1}{2} \right| + \left| \sum_{i \in \mathcal{I}^-} v_i + \frac{1}{2} \right| \\
&\leq \frac{\epsilon}{2} + \frac{\epsilon}{2}. \tag{23}
\end{aligned}$$

where (a) used the observation that for all  $i \in \mathcal{I}^+$ ,  $v_i - v'_i$  are of the same sign (and similarly for  $\mathcal{I}^-$ ). In the final line we used the following simple observation. If  $a$  and  $b$  are two real numbers such that  $|a| + |b| = 1$ ,  $|a + b| \leq \epsilon$  and  $a > 0$ , then  $|a - 1/2| \leq \epsilon/2$  &  $|b + 1/2| \leq \epsilon/2$ . The claim follows upon setting  $a = \sum_{i \in \mathcal{I}^+} v_i$  and  $b = \sum_{i \in \mathcal{I}^-} v_i$ . Substituting (23) in (22) yields  $\sum_{i=1}^m v_i \xi(\mathbf{x}_i) \geq 2\epsilon - \epsilon = \epsilon$ , which establishes the claim.  $\square$

### Appendix B: Proof of Theorem 4.1

Let  $\mathcal{I}^+(\mathcal{I}^-)$  denote the indices of the positively (negatively) labeled points. Due to Lemma 4.1 we assume without loss of generality that  $P$  is symmetric, namely  $\sum_{i \in \mathcal{I}^+} P_i = \sum_{i \in \mathcal{I}^-} P_i$ , implying that  $\sum_i v_i = 0$ , where  $v_i = y_i P_i$ . We start from the basic identity (9). Recall that the measure  $\Phi$  is defined over the set  $R = \{r_1, r_2, \dots, r_n\}$  and the convolution measure over  $X \times R$  is given by  $(v \star \Phi)(\mathbf{x}_i, r_k) = v(\mathbf{x}_i) \Phi(r_k)$ . Let  $X^+ = X_1^+ \cup \dots \cup X_{K_+}^+$ , and similarly for  $X^-$ , be partitions of  $X^+$  and  $X^-$ , respectively. Denote the indices of points in  $X_k^\pm$  by  $\mathcal{I}_k^\pm$ . Keeping in mind the negativity of  $J_\Phi(\mathbf{x}_i, \mathbf{x}_j)$  (Lemma 3.3) and  $I(\Phi)$  (Lemma 3.1), we have the expression,

$$\begin{aligned}
-I(v) &\geq |I(\Phi)| \sum_i v_i^2 - \sum_{i \in \mathcal{I}^+} \sum_{j \in \mathcal{I}^-} |J_\Phi(\mathbf{x}_i, \mathbf{x}_j)| |v_i| |v_j| \\
&\quad + \sum_{k=1}^{K_+} \sum_{\substack{i \neq j \\ i, j \in \mathcal{I}_k^+}} |J_\Phi(\mathbf{x}_i, \mathbf{x}_j)| v_i v_j + \sum_{k=1}^{K_-} \sum_{\substack{i \neq j \\ i, j \in \mathcal{I}_k^-}} |J_\Phi(\mathbf{x}_i, \mathbf{x}_j)| v_i v_j, \tag{24}
\end{aligned}$$

where we have used the fact that for points  $\mathbf{x}_i$  and  $\mathbf{x}_j$  that belong to the same class the weights  $v_i$  and  $v_j$  are of equal sign, while  $v_i v_j < 0$  if  $i \in \mathcal{I}^+$  and  $v_j \in \mathcal{I}^-$ . Note that we have replaced the equality sign in (9) by an inequality, as we have removed the term corresponding to the contribution from the regions of equal sign. The remainder of the proof proceeds by first defining an appropriate measure  $\Phi$ , and then bounding each of the terms in (21).

We begin with a useful lemma from Alexander (1994).

**Lemma B.1** (Lemma 6, Alexander, 1994). *Let  $\nu$  be an atomic measure concentrated on a set of  $m$  points in  $\mathbb{R}^d$ . Let  $\delta$  be the minimal distance between any two points. Then*

$$-I(\nu) \geq c_d \delta \sum_{i=1}^m \nu_i^2,$$

where  $c_2 = 0.02$  for  $d = 2$ . Clearly this is also a lower bound for  $d \geq 2$ .

Defining the measure  $\Phi$

Similarly to (Alexander, 1994) define the signed one-dimensional measure  $\tilde{\Phi}_{n-1}$ , concentrated on  $\{0, 1, \dots, n\}$  as follows:

$$\tilde{\Phi}_{n-1}(k) = (-1)^k \binom{n}{k}.$$

It is easily seen that  $\tilde{\Phi}$  obeys the following following two conditions:

$$\begin{aligned} \|\tilde{\Phi}_{n-1}\|_1 &= 2^n, \\ \sum_i i^s \tilde{\Phi}_{n-1}(i) &= 0 \quad (0 \leq s \leq n-1), \end{aligned}$$

namely,  $\tilde{\Phi}_{n-1}$  possesses  $n-1$  vanishing moments. The last property is easy to check by looking at the derivatives of the expression

$$(1-x)^n = \sum_{k=0}^n (-1)^k \binom{n}{k} x^k,$$

and setting  $x = 1$ .

We then define a normalized measure  $\Phi$ , supported on the set  $\{0, \eta/4n, 2\eta/4n, \dots, \eta/4\}$  by

$$\Phi(k\eta/4n) \triangleq 2^{-n} \tilde{\Phi}(k) = (-1)^k \binom{n}{k} 2^{-n} \quad (k = 1, 2, \dots, n). \tag{25}$$

Clearly  $\|\Phi\|_1 = 1$  and the  $n-1$  first moments of  $\Phi$  vanish.

Before moving to the actual proof, we recall an additional result from Alexander (1991), which will be used in the sequel.

**Lemma B.2** (Alexander, 1991, Corollary 2). *Let  $\Phi$  be a measure in  $\Psi(\mathbb{R})$  for which the first  $n$  moments vanish. Then  $I^{2k}(\Phi) = 0$  for  $1 \leq k \leq n$ , where  $I^{2k}(\Phi)$  is given in (12).*

Lower bound on the first term in (21)

From Lemma B.1 we conclude that

$$\begin{aligned} |I(\Phi)| &\geq 0.02(\eta/4n) \sum_{k=1}^n |\phi_k|^2 \\ &= (0.005\eta/n)2^{-2n} \sum_{k=0}^n \binom{n}{k}^2 \\ &\stackrel{(a)}{=} (0.005\eta/n)2^{-2n} \binom{2n}{n} \\ &\geq 0.001\eta n^{-3/2}. \end{aligned}$$

In (a) we used a standard identity for binomial coefficients (see (0.157.1) in Gradshteyn & Ryzhik, 1994), and the Stirling bound

$$\sqrt{2\pi n} \left(\frac{n}{e}\right)^n \leq n! \leq \sqrt{2\pi n} \left(\frac{n}{e}\right)^{n+1/2}$$

was used in the final step, in the form  $\binom{2n}{n} \geq (1/2)(\pi n)^{-1/2}2^{2n}$ .

Lower bound on the second term in (21)

Using a Taylor expansion for the square root function one may show (see Lemma 3.2) that

$$J_{\Phi}(\mathbf{x}, \mathbf{x}') = \rho \sum_{k=1}^{\infty} c_k I^{2k}(\Phi) \rho^{-2k}$$

where  $\rho = \|\mathbf{x} - \mathbf{x}'\|_2$  and  $I^{2k}(\Phi)$  is defined in (12). Let  $h = \eta/4$  be the diameter of the support of  $\Phi$ . First, recall that  $|c_k| < 1$  and  $I^k(\Phi) = 0$  for  $k = 0, 1, \dots, n-1$  (see Alexander, 1994). Moreover, from Lemma B.2 we have that  $I^{2k}(\Phi) = 0$  for  $1 \leq k \leq n-1$ , and  $|I^{2k}(\Phi)| \leq \|\Phi\|_1^2 h^{2k}$  (see (13)) in general. We find that

$$\begin{aligned} |J_{\Phi}(\mathbf{x}, \mathbf{x}')| &\leq \rho \sum_{k=n}^{\infty} h^{2k} \rho^{-2k} \\ &= \rho [1 - h^2 \rho^{-2}]^{-1} \left(\frac{h}{\rho}\right)^{2n} \end{aligned}$$

where we have used the fact that  $\|\Phi\|_1 = 1$ . Since  $\rho \geq \eta$  (recall that  $\mathbf{x}_i \in X^+$  and  $\mathbf{x}_j \in X^-$ ) and  $h = \eta/4$ , it follows that  $h/\rho \leq 1/4$  and thus  $[1 - h^2 \rho^{-2}]^{-1} \leq 16/15$ . We then find that

$$\begin{aligned} |J_{\Phi}(\mathbf{x}, \mathbf{x}')| &\leq 1.07 \left(\frac{h}{\rho}\right)^{2n-1} h \\ &\leq 1.07 \cdot 4^{-(2n-1)} (\eta/4) \end{aligned}$$

Continuing we find that

$$\begin{aligned} -\sum_{i \in \mathcal{I}^+} \sum_{j \in \mathcal{I}^-} |J_\Phi(\mathbf{x}_i, \mathbf{x}_j)| |v_i| |v_j| &\geq -1.07 \cdot 2^{-4n} \eta \sum_{i \in \mathcal{I}^+} \sum_{j \in \mathcal{I}^-} |v_i| |v_j| \\ &\geq -1.07 \cdot 2^{-4n} \eta \sqrt{m_+ m_-} \left( \sum_{i \in \mathcal{I}^+} v_i^2 \right)^{1/2} \left( \sum_{j \in \mathcal{I}^-} v_j^2 \right)^{1/2}, \end{aligned}$$

where  $m_\pm = |\mathcal{I}^\pm|$ . In the second step we used the Cauchy-Schwartz inequality to show that  $\sum_i \sum_j |v_i| |v_j| \leq \sqrt{(\sum_i \sum_j v_i^2 v_j^2)(\sum_i \sum_j 1)}$ .

*Lower bound on the third and fourth terms in (21)*

We now move on to deal with the final terms in (21). For this purpose we need a lower bound on  $|J_\Phi(\mathbf{x}_i, \mathbf{x}_j)|$ . From (10) we have the explicit expression

$$J_\Phi(\mathbf{x}_i, \mathbf{x}_j) = \sum_{k=0}^n \sum_{l=0}^n \|\mathbf{q}_{ik} - \mathbf{q}_{jl}\|_2 \phi_k \phi_l,$$

where  $\phi_k = \Phi(r_k)$  and  $\mathbf{q}_{ik} = (\mathbf{x}_i^T, r_k)^T \in \mathbb{R}^{d+1}$ .

Let  $\mathbf{x}$  and  $\mathbf{x}'$  be points in  $\mathbb{R}^d$  such that  $\|\mathbf{x} - \mathbf{x}'\|_2 = \rho$ . Let  $\zeta = \eta/4n\rho$ , then the following relations hold for the measure  $\Phi$  defined in (25).

$$\begin{aligned} \rho^{-1} J_\Phi(\mathbf{x}, \mathbf{x}') &= \rho^{-1} \sum_{k=0}^n \sum_{\ell=0}^n (-1)^{k+\ell} (\rho^2 + (\eta/4n)^2 (k-\ell)^2)^{1/2} \binom{n}{k} \binom{n}{\ell} \\ &= \sum_{k=0}^n \sum_{\ell=0}^n (-1)^{k+\ell} (1 + \zeta^2 (k-\ell)^2)^{1/2} \binom{n}{k} \binom{n}{\ell} \\ &\stackrel{(a)}{=} \sum_{k=0}^n \binom{n}{k}^2 - 2(1 + \zeta^2)^{1/2} \sum_{k=1}^n \binom{n}{k} \binom{n}{k-1} \\ &\quad + 2(1 + 4\zeta^2)^{1/2} \sum_{k=2}^n \binom{n}{k} \binom{n}{k-2} + \dots \\ &\stackrel{(b)}{=} \sum_{k=0}^n \binom{n}{k}^2 + 2 \sum_{\ell=1}^n (-1)^\ell (1 + \ell^2 \zeta^2)^{1/2} \sum_{k=\ell}^n \binom{n}{k} \binom{n}{k-\ell} \\ &\stackrel{(c)}{=} \binom{2n}{n} + 2 \sum_{\ell=1}^n (-1)^\ell (1 + \ell^2 \zeta^2)^{1/2} \binom{2n}{n-\ell}. \end{aligned} \tag{26}$$

In step (a) we separated the sum into terms with  $k = \ell$ ,  $|k - \ell| = 1$ ,  $|k - \ell| = 2$ , etc., while in step (b) a re-arrangement of the sum was performed. Step (c) follows from the

identity (0.156.1) in Gradshteyn and Ryzhik (1994), given by

$$\sum_{k=0}^{n-\ell} \binom{n}{k} \binom{n}{\ell+k} = \binom{2n}{n-\ell}.$$

For the specific measure  $\Phi$  introduced in (25), define

$$F_n(\rho, \eta) = -J_\Phi(\mathbf{x}, \mathbf{x}') > 0 \quad (\|\mathbf{x} - \mathbf{x}'\|_2 = \rho). \tag{27}$$

From Lemma 3.3 we conclude that  $F_n(\rho, \eta)$  is a positive and strictly decreasing function of  $\rho$ . Thus, the third term in (24) is bounded from below by

$$\sum_{k=1}^{K_+} F_n(\rho_k^+, \eta) \sum_{\substack{i \neq j \\ i, j \in \mathcal{I}_k^+}} |v_i| |v_j|$$

where

$$\rho_k^+ = \max_{i, j \in \mathcal{I}_k^+} \|\mathbf{x}_i - \mathbf{x}_j\|_2.$$

A similar bound holds for the fourth term in (24).

*Concluding the proof*

Combining all the results derived, we have the following lower bound,

$$\begin{aligned} -I(v) \geq & \eta \left[ 0.001n^{-3/2} \sum_i v_i^2 - 1.07 \times 2^{-4n} \sqrt{m_+ m_-} \|v\|_+ \|v\|_- \right] \\ & + \sum_{k=1}^{K_+} F_n(\rho_k^+, \eta) \sum_{\substack{i \neq j \\ i, j \in \mathcal{I}_k^+}} |v_i| |v_j| + \sum_{k=1}^{K_-} F_n(\rho_k^-, \eta) \sum_{\substack{i \neq j \\ i, j \in \mathcal{I}_k^-}} |v_i| |v_j|, \end{aligned} \tag{28}$$

where  $\|v\|_\pm = (\sum_{i \in \mathcal{I}^\pm} v_i^2)^{1/2}$ . The parameter  $n$  has been free up to now. In order to guarantee that the second term is smaller than the first, we need to select  $n$  large enough so that the second term is smaller than (say) half the first term, namely

$$0.0005n^{-3/2} \sum_i v_i^2 > 1.07 \cdot 2^{-4n} \sqrt{m_+ m_-} \|v\|_+ \|v\|_-.$$

Simple mathematical manipulation shows that this is equivalent to the requirement

$$n \geq \frac{3}{8} \log_2 \frac{n}{a}$$

$$a = \left[ \frac{1}{2140} \frac{\|v\|^2}{\sqrt{m_+ m_-} \|v\|_+ \|v\|_-} \right]^{2/3}$$

We use the following simple lemma.

**Lemma B.3** (Vidyasagar, 1996, Lemma 4.4). *Suppose  $\alpha, \beta > 0$ ,  $\alpha\beta > 2$ ,  $n \geq 1$ . Then*

$$n \geq 2\alpha \log \alpha\beta \Rightarrow n > \alpha \log \beta n.$$

We thus conclude that a sufficient condition on  $n$  is that

$$n \geq \frac{1}{2} \log_2 \frac{492\sqrt{m_+ m_-} \|v\|_+ \|v\|_-}{\|v\|^2}$$

Since  $\|v\|_+ \|v\|_- \leq \|v\|^2$ , it suffices that

$$n \geq \frac{1}{2} \log_2 (492\sqrt{m_+ m_-})$$

Thus, the first two terms in (28) can be lower bounded by  $0.004 \log_2 (492\sqrt{m_+ m_-}) \|v\|^2$ . Substituting in (28) we obtain the desired lower bound on  $-I(v)$ . The bound on the error  $\epsilon$  is obtained by making use of (6) and (8).  $\square$

### Appendix C: Main properties of $F_n(\rho, \eta)$

We establish some of the main properties of the function  $F_n(\rho, \eta)$ . We begin with a combinatorial lemma.

**Lemma C.1.**

$$2 \sum_{\ell=1}^n (-1)^\ell \ell^{2k} \binom{2n}{n-\ell} = \begin{cases} 0 & \text{if } 1 \leq k < n, \\ (2n)! & \text{if } k = n. \end{cases}$$

**Proof:** The proof is by induction over  $k$ . First, recall from Gradshteyn and Ryzhik (1994) (0.154.3) that  $\sum_{\ell=0}^{2n} (-1)^\ell \ell^2 \binom{2n}{\ell} = 0$ . We therefore have that:

$$0 = \sum_{\ell=0}^{2n} (-1)^\ell \ell^2 \binom{2n}{\ell}$$

$$= (-1)^n n^2 \binom{2n}{n} + (-1)^n \sum_{\ell=1}^n ((n-\ell)^2 + (n+\ell)^2) (-1)^\ell \binom{2n}{n-\ell}$$

$$\begin{aligned}
&= (-1)^n n^2 \binom{2n}{n} + 2n^2 (-1)^n \sum_{\ell=1}^n (-1)^\ell \binom{2n}{n-\ell} \\
&\quad + 2(-1)^n \sum_{\ell=1}^n (-1)^\ell \ell^2 \binom{2n}{n-\ell}
\end{aligned}$$

Next, note that the first two terms vanish. To see this observe that

$$\binom{2n}{n} + 2 \sum_{\ell=1}^n (-1)^\ell \binom{2n}{n-\ell} = 0. \tag{29}$$

This claim follows by noting that (29) is simply the binomial expansion of  $(-1)^n (1-1)^{2n}$ , where the sum has been re-arranged. Thus  $\sum_{\ell=1}^n (-1)^\ell \ell^2 \binom{2n}{n-\ell} = 0$ .

We proceed to the induction step. Assume the result is true for  $k' = 1, \dots, k-1$ , then from Gradshteyn and Ryzhik (1994) (0.154.3) for  $k < n$

$$\sum_{\ell=0}^{2n} (-1)^\ell \ell^{2k} \binom{2n}{\ell} = 0.$$

Rearranging and collecting terms yields:

$$\begin{aligned}
0 &= (-1)^n n^{2k} \binom{2n}{n} + (-1)^n \sum_{\ell=1}^n ((n-\ell)^{2k} + (n+\ell)^{2k}) (-1)^\ell \binom{2n}{n-\ell} \\
&= (-1)^n n^{2k} \binom{2n}{n} + 2n^{2k} (-1)^n \sum_{\ell=1}^n (-1)^\ell \binom{2n}{n-\ell} \\
&\quad + 2(-1)^n \sum_{m=1}^{k-1} \binom{2k}{2m} n^{2k-2m} \sum_{\ell=1}^n (-1)^\ell \ell^{2m} \binom{2n}{n-\ell} \\
&\quad + 2(-1)^n \sum_{\ell=1}^n (-1)^\ell \ell^{2k} \binom{2n}{n-\ell} \tag{30}
\end{aligned}$$

Note that all the odd powers vanish, since the sum of binomial coefficients for  $(n-\ell)^{2k} + (n+\ell)^{2k}$  consists of even powers only. The sum of the two first terms vanish as in (29), and the third term vanishes too by the induction step. The last term is the desired element and is therefore 0.

For the case  $k = n$  we proceed as follows. From (0.154.4) in Gradshteyn and Ryzhik (1994)

$$\sum_{\ell=0}^{2n} (-1)^\ell \ell^{2n} \binom{2n}{\ell} = (2n)!.$$

Repeating the argumentation above (see (30)) results in

$$2 \sum_{\ell=1}^n (-1)^\ell \ell^{2n} \binom{2n}{n-\ell} = (2n)!, \tag{31}$$

which establishes the claim. □

We are now ready to state and prove the main properties of  $F_n(\rho, \eta)$ .

**Lemma C.2.** *The function  $F_n(\rho, \eta)$  defined in (18) satisfies the following properties:*

1. *For fixed  $\eta$ ,  $F_n(\rho, \eta)$  is a positive and monotonically decreasing function of  $\rho$ .*
2. *When  $\eta/\rho$  is large (relative to  $n$ ),  $\rho^{-1}F_n(\rho, \eta) = \Theta(a + b\eta/\rho)$ , where  $a$  and  $b$  are positive constants.*
3. *For small values of  $\eta/\rho$ ,*

$$F_n(\rho, \eta) = \Theta\left(\left[\left(\frac{1}{4n}\right)^{2n} \frac{1 \cdot 1 \cdot 3 \cdots (2n-3)}{2 \cdot 4 \cdot 6 \cdots (2n)} (2n)!\right] \rho \left(\frac{\eta}{\rho}\right)^{2n}\right).$$

**Proof:** We begin by recalling the definition of  $F_n(\rho, \eta)$ .

$$\rho^{-1}F_n(\rho, \eta) = -\binom{2n}{n} - 2 \sum_{\ell=1}^n (-1)^\ell (1 + \ell^2(\eta/4n\rho)^2)^{1/2} \binom{2n}{n-\ell}. \tag{32}$$

The first property follows from Lemma 3.3 (taken from Alexander, 1991).

For large  $\eta/\rho$  we have that  $(1 + \ell^2(\eta/4n\rho)^2)^{1/2} = \Theta(\ell\eta/4n\rho)$ . Thus, (32) yields:

$$\rho^{-1}F_n(\rho, \eta) = \Theta\left[-\binom{2n}{n} - (\eta/2n\rho) \sum_{\ell=1}^n (-1)^\ell \ell \binom{2n}{n-\ell}\right]. \tag{33}$$

This proves the second property for  $b = -(1/2n) \sum_{\ell=1}^n (-1)^\ell \ell \binom{2n}{n-\ell}$  and  $a = -\binom{2n}{n}$ .

For small  $\eta/\rho$  one can use the following Taylor series (e.g., Gradshteyn & Ryzhik, 1994)

$$\sqrt{1+x} = 1 + \sum_{k=1}^{\infty} (-1)^{k+1} c_k x^k,$$

where  $c_1 = 1$  and

$$c_k = \frac{1 \cdot 1 \cdot 3 \cdots (2k-3)}{2 \cdot 4 \cdot 6 \cdots (2k)} \quad (k > 1).$$

Substituting this in (32) results in the following expression:

$$\rho^{-1}F_n(\rho, \eta) = -\binom{2n}{n} - 2 \sum_{\ell=1}^n (-1)^\ell \left(1 + \sum_{k=1}^{\infty} (-1)^{k+1} c_k (\ell\eta/4n\rho)^{2k}\right) \binom{2n}{n-\ell}. \tag{34}$$

Changing the order of summation yields the following equation:

$$\begin{aligned} \rho^{-1} F_n(\rho, \eta) = & -\binom{2n}{n} - 2 \sum_{\ell=1}^n (-1)^\ell \binom{2n}{n-\ell} \\ & + 2 \sum_{k=1}^{\infty} c_k \sum_{\ell=1}^n (-1)^\ell (\ell\eta/4n\rho)^{2k} \binom{2n}{n-\ell}. \end{aligned} \quad (35)$$

The sum of the first two terms vanishes, as in (29). We proceed to study the elements in the second sum. As a result of Lemma C.1 we know that the first  $n - 1$  elements in the sum over  $k$  in the third term on the r.h.s. of (35) vanish. From Lemma C.1 we have that

$$2 \sum_{\ell=1}^n (-1)^\ell \ell^{2n} \binom{2n}{n-\ell} = (2n)!. \quad (36)$$

By taking the first non-vanishing element in (35) as the approximation we have that:

$$\rho^{-1} F_n(\rho, \eta) = \Theta(c_n(\eta/4n\rho)^{2n}(2n)!) \quad (37)$$

□

### Acknowledgments

We are grateful to Ron Aharoni for directing us to the field of geometric discrepancy, and to Jiří Matoušek, Allen Rogers and Shai Ben-David for very helpful discussions and comments. Special thanks to the insightful anonymous reviewers who have done an excellent job in providing us with comments and suggestions which have greatly improved the manuscript.

### Note

1. To see this, pick a random classifier  $h$ . If the error is smaller or equal to  $1/2$ , we are done; if not, the error of the classifier  $-h$  is smaller than  $1/2$ .

### References

- Alexander, R. (1975). Generalized sums of distances. *Pacific J. Math.*, 56, 297–304.
- Alexander, R. (1990). Geometric methods in the study of irregularities of distribution. *Combinatorica*, 10:2, 115–136.
- Alexander, R. (1991). Principles of a new method in the study of irregularities of distribution. *Invent. Math.*, 103, 279–296.
- Alexander, R. (1994). The effect of dimension on certain geometric problems of irregularities of distribution. *Pac. J. Math.*, 165:1, 1–15.
- Anthony, M., & Bartlett, P.L. (1999). *Neural Network Learning; Theoretical Foundations*. Cambridge: Cambridge University Press.

- Bartlett, P., & Ben-David, S. (1999). On the hardness of learning with neural networks. In *Proceedings of the Fourth European Conference on Computational Learning Theory'99*.
- Breiman, L. (1996). Bagging predictors. *Machine Learning*, 26:2, 123–140.
- Breiman, L. (1998). *Prediction games and arcing algorithms*. Technical Report 504, Berkeley.
- Cristianini, N., & Shawe-Taylor, J. (2000). *An introduction to support vector machines: And other kernel-based learning methods*. Cambridge, England: Cambridge University Press.
- Freund, Y. (1995). Boosting a weak learning algorithm by majority. *Information and Computation*, 121, 256–285.
- Freund, Y., & Schapire, R. E. (1996). Experiments with a new boosting algorithm. In *Proceeding of the Thirteenth International Conference on Machine Learning* (pp. 148–156).
- Freund, Y., & Schapire, R. E. (1996). Game theory, on-line prediction and boosting. In *Proceeding of the Thirteenth International Conference on Machine Learning* (pp. 148–156).
- Friedman, J. Hastie, T., & Tibshirani, R. (2000). Additive logistic regression: A statistical view of boosting. *The Annals of Statistics*, 38:2, 337–374.
- Gradshteyn, I. S., & Ryzhik, I. M. (1994). *Tables of Integrals, Series and Products*. 5th edn. New York: Academic Press.
- Johnson, D. S., & Preparata, F. P. (1978). The densest hemisphere problem. *Theoretical Computer Science*, 6, 93–107.
- Jordan, M. I., & Jacobs, R. A. (1994). Hierarchical mixtures of experts and the EM algorithm. *Neural Computation*, 6:2, 181–214.
- Kearns, M., & Mansour, Y. (1996). On the boosting ability of top-down decision tree learning algorithms. In *Proc. 28th ACM Symposium on the Theory of Computing* (pp. 459–468). New York: ACM Press.
- Koltchinskii, V., Panchenko, D., & Lozano, F. (2001). Some new bounds on the generalization error of combined classifiers. In T. Dietterich (Eds.), *Advances in neural information processing systems 14*, Boston, MIT Press.
- Mason, L., Baxter, J., Bartlett, P., & Frean, M. (2000). Functional gradient techniques for combining hypotheses. In B. Schölkopf, A. Smola, P. Bartlett, & D. Schuurmans (Eds.), *Advances in Large Margin Classifiers*. Boston: MIT Press.
- Matoušek, J. (1995). Tight upper bound on the discrepancy of half-spaces. *Discrete & Computational Geometry*, 13, 593–601.
- Matoušek, J. (1999). *Geometric Discrepancy: An Illustrated Guide*. New York: Springer Verlag.
- Meir, R. El-Yaniv, R., & Ben-David, S. (2000). Localized boosting. In N. Cesa-Bianchi and S. Goldman (Eds.), *Proc. Thirteenth Annual Conference on Computational Learning Theory* (pp. 190–199). San Francisco, CA: Morgan Kaufman.
- Mannor S., & Meir, R. (2001). Weak learners and improved rates of convergence in boosting. In T. Dietterich (Ed.), *Advances in Neural Information Processing Systems 14*, Boston. MIT Press.
- Motwani R., & Raghavan, P. (1995). *Randomized Algorithms*. Cambridge, England: Cambridge University Press.
- Santaló, L. A. (1976). *Integral geometry and geometric probability*. Reading, MA: Addison-Wesley.
- Schoenberg, I. J. (1937). On certain metric spaces arising from euclidean spaces by change of metric and their embedding into hilbert space. *Ann. Math.*, 38, 787–793.
- Schapire, R. (1990). The strength of weak learnability. *Machine Learning*, 5:2, 197–227.
- Schapire, R. E., Freund, Y., Bartlett, P., & Lee, W. S. (1998). Boosting the margin: A new explanation for the effectiveness of voting methods. *The Annals of Statistics*, 26:5, 1651–1686.
- Schapire, R. E., & Singer, Y. (1999). Improved boosting algorithms using confidence-rated predictions. *Machine Learning*, 37:3, 297–336.
- Vapnik, V. N. (1982). *Estimation of dependences based on empirical data*. New York: Springer Verlag.
- Vapnik, V. N. (1998). *Statistical learning theory*. New York: Wiley Interscience.
- Vapnik, V., & Chervonenkis, A. (1971). On the uniform convergence of relative frequencies of events to probabilities. *Theory of Probability and its Applications*, 16, 264–280.
- Vidyasagar, M. (1996). *A Theory of Learning and Generalization*. New York: Springer Verlag.

Received August 9, 2000

Revised January 11, 2001

Accepted January 12, 2001

Final manuscript February 9, 2001