



PAC Analogues of Perceptron and Winnow Via Boosting the Margin

R. SERVEDIO*

rocco@cs.harvard.edu

Division of Engineering and Applied Sciences, Harvard University, Cambridge, MA, USA

Editor: Jyrki Kivinen

Abstract. We describe a novel family of PAC model algorithms for learning linear threshold functions. The new algorithms work by boosting a simple weak learner and exhibit sample complexity bounds remarkably similar to those of known online algorithms such as Perceptron and Winnow, thus suggesting that these well-studied online algorithms in some sense correspond to instances of boosting. We show that the new algorithms can be viewed as natural PAC analogues of the online p -norm algorithms which have recently been studied by Grove, Littlestone, and Schuurmans (1997, *Proceedings of the Tenth Annual Conference on Computational Learning Theory* (pp. 171–183) and Gentile and Littlestone (1999, *Proceedings of the Twelfth Annual Conference on Computational Learning Theory* (pp. 1–11)). As special cases of the algorithm, by taking $p = 2$ and $p = \infty$ we obtain natural boosting-based PAC analogues of Perceptron and Winnow respectively. The $p = \infty$ case of our algorithm can also be viewed as a generalization (with an improved sample complexity bound) of Jackson and Craven’s PAC-model boosting-based algorithm for learning “sparse perceptrons” (Jackson & Craven, 1996, *Advances in neural information processing systems 8*, MIT Press). The analysis of the generalization error of the new algorithms relies on techniques from the theory of large margin classification.

Keywords: probably approximately correct learning, boosting, linear threshold functions

1. Introduction

One of the most fundamental problems in computational learning theory is that of learning an unknown linear threshold function from labeled examples. Many different learning algorithms for this problem have been considered over the past several decades. In particular, in recent years many researchers have studied simple online additive and multiplicative update algorithms, namely the Perceptron and Winnow algorithms and variants thereof (Auer & Warmuth, 1995; Baum, 1990; Bylander, 1998; Freund & Schapire, 1999; Gentile & Littlestone, 1999; Grove, Littlestone, & Schuurmans, 1997; Kivinen, Warmuth, & Auer, 1997; Littlestone, 1988, 1989, 1991; Servedio, 1999; Schmitt, 1998).

This paper takes a different approach. We describe a natural parameterized family of boosting-based PAC algorithms for learning linear threshold functions. The weak hypotheses used are linear functionals and the strong classifier obtained is a linear threshold function. Although our new algorithms are conceptually and algorithmically very different from Perceptron and Winnow, we establish performance bounds for the new algorithms which are

*Supported in part by an NSF Graduate Fellowship, by NSF grant CCR-95-04436 and by ONR grant N00014-96-1-0550.

remarkably similar to those of Perceptron and Winnow; we thus refer to the new algorithms as *PAC analogues* of Perceptron and Winnow. We hope that the analysis of these new algorithms will yield fresh insights into the relationship between boosting and online algorithms.

We give a unified analysis of our Perceptron and Winnow analogues which includes many other algorithms as well. Grove, Littlestone, and Schuurmans (1997) have shown that Perceptron and (a version of) Winnow can be viewed as the $p = 2$ and $p \rightarrow \infty$ cases of a general online p -norm linear threshold learning algorithm, where $p \geq 2$ is any real number. We present PAC-model boosting-based analogues of these online p -norm algorithms for any value $2 \leq p \leq \infty$. The PAC-model Perceptron and Winnow analogues mentioned above are respectively the $p = 2$ and $p = \infty$ cases of this general algorithm.

The $p = \infty$ case of our algorithm can also be viewed as a generalization of Jackson and Craven's PAC-model algorithm for learning "sparse perceptrons" (Jackson & Craven, 1996). Their algorithm boosts using weak hypotheses which are single Boolean literals; this is similar to what the $p = \infty$ case of our algorithm does. Our analysis of the $p = \infty$ case generalizes their algorithm to deal with real-valued rather than Boolean input variables, thus achieving a goal stated in Jackson and Craven (1996), and also yields a substantially stronger sample complexity bound than was established in Jackson and Craven (1996).

Section 2 of this paper contains preliminary material, including an overview of the online p -norm algorithms from Gentile and Littlestone (1999) and Grove, Littlestone, and Schuurmans (1997). In Section 3 we present a simple PAC-model p -norm algorithm and prove that it is a weak learning algorithm for all $2 \leq p < \infty$. In Section 4 we apply techniques from the theory of large margin classification to show how our weak learning algorithm can be boosted to a strong learning algorithm with small sample complexity. Finally, in Section 5 we compare our PAC algorithms with the analogous online algorithms, extend our algorithm to the case $p = \infty$, and discuss the relationship between the $p = \infty$ case of our algorithm and the Jackson–Craven algorithm for learning sparse perceptrons.

1.1. Related work

Several authors have previously studied linear threshold learning algorithms which work by combining weak predictors. Freund and Schapire have studied an algorithm which predicts using a weighted vote of the hypotheses which the Perceptron algorithm generates during its training phase (Freund & Schapire, 1999). The weight of each hypothesis in this vote is proportional to its survival time, i.e. the number of examples which elapse before it classifies an example incorrectly and causes the Perceptron algorithm to generate a new hypothesis. Freund and Schapire prove generalization error bounds on the resulting classifier which are similar to Vapnik's generalization error bounds for the "maximal margin" hyperplane (Vapnik, 1998). The Freund-Schapire algorithm differs from our approach in several ways: for one thing, their algorithm is unlike ours in that it does not use boosting to combine the weak predictors. Additionally, whereas our algorithm's final hypothesis is a single linear threshold function, their algorithm's final hypothesis is a depth-2 threshold circuit (a weighted vote over Perceptron hypotheses which are themselves linear threshold functions).

Ji and Ma have suggested that a random-search-and-test approach can be used to find weak classifier linear threshold functions for certain restricted halfspace learning problems

(Ji & Ma, 1997). They propose combining these weak classifier linear threshold functions with a simple majority vote; thus, their approach also results in a final hypothesis which is a depth 2 threshold circuit.

Our approach is closest to that of Jackson and Craven, who use boosting to combine single literals into a strong hypothesis linear threshold function. As we show in Section 5, the $p = \infty$ case of our algorithm strengthens and generalizes their results.

The close similarity in performance bounds between our boosting-based algorithms and the online p -norm algorithms suggests a relationship between boosting and online learning. Freund and Schapire (1996) and Schapire (1999) have investigated this relationship in the context of game theory.

2. Preliminaries

We start with some geometric definitions. For a point $\tilde{x} = (x_1, \dots, x_n) \in \mathfrak{R}^n$ and $p \geq 1$ we write $\|\tilde{x}\|_p$ to denote the p -norm of \tilde{x} , namely

$$\|\tilde{x}\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}.$$

The ∞ -norm of \tilde{x} is $\|\tilde{x}\|_\infty = \max_{i=1, \dots, n} |x_i|$. For $p, q \geq 1$ the q -norm is *dual* to the p -norm if $\frac{1}{p} + \frac{1}{q} = 1$; hence the 1-norm and the ∞ -norm are dual to each other and the 2-norm is dual to itself. In this paper p and q always denote dual norms. The following facts are well known (e.g. Taylor & Mann, 1972, pp. 203–204):

Hölder Inequality: $|\tilde{u} \cdot \tilde{v}| \leq \|\tilde{u}\|_p \|\tilde{v}\|_q$ for all $\tilde{u}, \tilde{v} \in \mathfrak{R}^n$ and $1 \leq p \leq \infty$.

Minkowski Inequality: $\|\tilde{u} + \tilde{v}\|_p \leq \|\tilde{u}\|_p + \|\tilde{v}\|_p$ for all $\tilde{u}, \tilde{v} \in \mathfrak{R}^n$ and $1 \leq p \leq \infty$.

Throughout this paper the *example space* X is a subset of \mathfrak{R}^n . A *linear threshold function* over X is a function f such that $f(\tilde{x}) = \text{sign}(\tilde{u} \cdot \tilde{x})$ for some $\tilde{u} \in \mathfrak{R}^n$ (recall that the function $\text{sign}(z)$ takes value 1 if $z \geq 0$ and takes value -1 if $z < 0$). We note that the standard definition of a linear threshold function allows a nonzero threshold, i.e. $f(\tilde{x}) = \text{sign}(\tilde{u} \cdot \tilde{x} - \theta)$ where θ can be any real number. However, any linear threshold function of this more general form over n variables is equivalent to a linear threshold function with threshold 0 over $n + 1$ variables, so our definition incurs no real loss of generality.

We write $\|X\|_p$ to denote $\sup_{\tilde{x} \in X} \|\tilde{x}\|_p$. We use the symbol $\delta_{\tilde{u}, X}$ to denote the quantity

$$\delta_{\tilde{u}, X} \stackrel{\text{def}}{=} \inf_{\tilde{x} \in X} (\tilde{u} \cdot \tilde{x}) (\text{sign}(\tilde{u} \cdot \tilde{x})),$$

which is a measure of the separation between examples in X and the hyperplane whose normal vector is \tilde{u} . We assume throughout the paper that $\|X\|_p < \infty$, i.e. the set X is bounded, and that $\delta_{\tilde{u}, X} > 0$, i.e. there is some nonzero lower bound on the separation between the hyperplane defined by \tilde{u} and the examples in X .

We write \log to denote logarithm base two and \ln to denote the natural logarithm.

2.1. PAC learning

For $\tilde{u} \in \mathfrak{R}^n$ let $EX(\tilde{u}, \mathcal{D})$ denote an *example oracle* which, when queried, provides a labeled example $\langle \tilde{x}, \text{sign}(\tilde{u} \cdot \tilde{x}) \rangle$ where \tilde{x} is drawn according to the distribution \mathcal{D} over X . We say that an algorithm A is a *strong learning algorithm for \tilde{u} on X* if it satisfies the following condition: there is a function $m(\epsilon, \delta, \tilde{u}, X)$ such that for any distribution \mathcal{D} over X , for all $0 < \epsilon, \delta < 1$, algorithm A makes at most $m(\epsilon, \delta, \tilde{u}, X)$ calls to $EX(\tilde{u}, \mathcal{D})$, and with probability at least $1 - \delta$ algorithm A outputs a hypothesis $h : X \rightarrow \{-1, 1\}$ such that $\Pr_{\tilde{x} \in \mathcal{D}}[h(\tilde{x}) \neq \text{sign}(\tilde{u} \cdot \tilde{x})] \leq \epsilon$. We say that such a hypothesis h is an ϵ -*accurate hypothesis for \tilde{u} under \mathcal{D}* and that the function $m(\epsilon, \delta, \tilde{u}, X)$ is the *sample complexity* of algorithm A .

As our main result we describe a strong learning algorithm and carefully analyze its sample complexity. To do this we must consider algorithms which do not satisfy the strong learning property but are still capable of generating hypotheses that have some slight advantage over random guessing (such so-called weak learning algorithms were first considered by Kearns and Valiant (1994)). Let

$$S = \langle \tilde{x}^1, \text{sign}(\tilde{u} \cdot \tilde{x}^1) \rangle, \dots, \langle \tilde{x}^m, \text{sign}(\tilde{u} \cdot \tilde{x}^m) \rangle$$

be a finite sequence of labeled examples from X and let \mathcal{D} be a distribution over S . For $0 < \gamma < 1/2$, we say that $h : X \rightarrow [-1, 1]$ is a $(1/2 - \gamma)$ -*approximator for \tilde{u} under \mathcal{D}* if

$$\frac{1}{2} \sum_{i=1}^m \mathcal{D}(\tilde{x}^i) \cdot |h(\tilde{x}^i) - \text{sign}(\tilde{u} \cdot \tilde{x}^i)| \leq \frac{1}{2} - \gamma. \quad (1)$$

We say that an algorithm A is a $(1/2 - \gamma)$ -*weak learning algorithm for \tilde{u} under \mathcal{D}* if the following condition holds: for any finite set S as described above and any distribution \mathcal{D} on S , if A is given \mathcal{D} and S as input then A outputs a hypothesis $h : X \rightarrow [-1, 1]$ which is a $(1/2 - \gamma)$ -approximator for \tilde{u} under \mathcal{D} . Thus for our purposes a weak learning algorithm is one which can always find a hypothesis that outperforms random guessing on a fixed sample.

2.2. Online learning and p -norm algorithms

In the *online* model, learning takes place over a sequence of trials. Throughout the learning process the learner maintains a hypothesis h which maps X to $\{-1, 1\}$. Each trial proceeds as follows: upon receiving an example $x \in X$ the learning algorithm outputs its prediction $\hat{y} = h(x)$ of the associated label y . The learning algorithm is then given the true label $y \in \{-1, 1\}$ and the algorithm can update its hypothesis h based on this new information before the next trial begins. The performance of an online learning algorithm on an example sequence is measured by the number of prediction mistakes which the algorithm makes.

Grove, Littlestone, and Schuurmans (1997) and Gentile and Littlestone (1999) have studied a family of online algorithms for learning linear threshold functions (see figure 1). We refer to this algorithm, which is parameterized by a real value $p \geq 2$, as the *online p -norm algorithm*. Like the well-known Perceptron algorithm, the online p -norm algorithm

Input parameters:real number $p \geq 2$ initial weight vector $\tilde{z}^0 = (z_1^0, \dots, z_n^0) \in \mathfrak{R}^n$ real number $a > 0$

1. **set** $t = 0$
2. **while** examples are available **do**
3. **get** unlabeled example \tilde{x}^t
4. **for all** $i = 1, \dots, n$ **set** $w_i^t = \text{sign}(z_i^t) |z_i^t|^{p-1}$
5. **predict** $\hat{y}_t = \text{sign}(\tilde{w}^t \cdot \tilde{x}^t)$
6. **get** label $y_t \in \{-1, +1\}$
7. **set** $\tilde{z}^{t+1} = \tilde{z}^t + a(y_t - \hat{y}_t)\tilde{x}^t$
8. **set** $t = t + 1$
9. **enddo**

Figure 1. The online p -norm algorithm.

updates its hypothesis by making an additive change to a weight vector \tilde{z} . However, as shown in steps 4 and 5 of figure 1, the p -norm algorithm does not use the \tilde{z} vector directly for prediction but rather predicts using a vector \tilde{w} which is a transformed version of the \tilde{z} vector, namely $w_i = \text{sign}(z_i) |z_i|^{p-1}$ for all $i = 1, \dots, n$. Note that when $p = 2$ we have $\tilde{z} = \tilde{w}$ and hence the online 2-norm algorithm is the Perceptron algorithm. In Grove, Littlestone and Schuurmans (1997) it is shown that as $p \rightarrow \infty$ the online p -norm algorithm approaches a version of the Winnow algorithm. More precisely, the following theorem from Grove, Littlestone, and Schuurmans (1997) gives mistake bounds for the online p -norm algorithms:

Theorem 1. Let $S = \langle \tilde{x}^1, y_1 \rangle, \dots, \langle \tilde{x}^m, y_m \rangle$ be a sequence of labeled examples where $\tilde{x} \in X$ and $y = \text{sign}(\tilde{u} \cdot \tilde{x})$ for every example $\langle \tilde{x}, y \rangle \in S$.

- (a) For any $2 \leq p < \infty$ and any $a > 0$, if the online p -norm algorithm is invoked with input parameters $(p, \tilde{z}^0 = (0, \dots, 0), a)$, then the mistake bound on the example sequence S is at most

$$\frac{(p-1) \|\tilde{u}\|_q^2 \|X\|_p^2}{\delta_{\tilde{u}, X}^2}.$$

- (b) For any $2 \leq p < \infty$, if \tilde{z}^0 satisfies $\tilde{u} \cdot \tilde{z}^0 > 0$ and $a = \frac{\delta_{\tilde{u}, X} \|\tilde{z}^0\|_p^2}{(p-1) \tilde{u} \cdot \tilde{z}^0 \|X\|_p^2}$, then the mistake bound on S is at most

$$\frac{(p-1) \|\tilde{u}\|_q^2 \|X\|_p^2}{\delta_{\tilde{u}, X}^2} \left(1 - \left(\frac{\tilde{u} \cdot \tilde{z}^0}{\|\tilde{u}\|_q \|\tilde{z}^0\|_p} \right)^2 \right).$$

- (c) Let $\tilde{z}^0 = (1, \dots, 1)$ and suppose that $u_i > 0$ for $i = 1, \dots, n$. If $p \rightarrow \infty$ and a is as described in part (b), then the mistake bound given in (b) converges to

$$\frac{2 \|\tilde{u}\|_1^2 \|X\|_\infty^2}{\delta_{\tilde{u}, X}^2} \left(\log n + \sum_{i=1}^n \frac{u_i}{\|\tilde{u}\|_1} \log \frac{u_i}{\|\tilde{u}\|_1} \right).$$

2.3. From online to PAC learning

Various generic procedures have been proposed (Angluin, 1988; Haussler, 1988; Kearns et al., 1987; Littlestone, 1989) for automatically converting on-line learning algorithms into PAC-model algorithms. In these procedures the sample complexity of the resulting PAC algorithm depends on the mistake bound of the original on-line learning algorithm. The strongest general result of this type (in terms of minimizing the sample complexity of the resulting PAC algorithm) is the conversion due to Littlestone (1989):

Theorem 2. *Let A be an online learning algorithm which changes its hypothesis only when it makes a mistake and which has a mistake bound of M for concept class C . Then there is a PAC-model learning algorithm A' for C as described above which has sample complexity*

$$O\left(\frac{1}{\epsilon} \left(\log \frac{1}{\delta} + M\right)\right).$$

By applying Theorem 2 to Theorem 1, one can obtain sample complexity bounds on a generic PAC-model conversion of the online p -norm algorithm. We now describe a completely different PAC-model algorithm which has remarkably similar sample complexity bounds.

3. A PAC-model p -norm weak learning algorithm

Our p -norm weak learning algorithm is motivated by the following simple idea: Suppose that $S = \langle \tilde{x}^1, y_1 \rangle, \dots, \langle \tilde{x}^m, y_m \rangle$ is a collection of labeled examples where $y_i = \text{sign}(\tilde{u} \cdot \tilde{x}^i)$ for each $i = 1, \dots, m$. Now imagine replacing each negative example $\langle \tilde{x}^i, -1 \rangle$ in S by the equivalent positive example $\langle -\tilde{x}^i, 1 \rangle$ to obtain a new collection S' of normalized examples. Let $\tilde{z} \in \mathfrak{R}^n$ be the average location of an example in S' , i.e. \tilde{z} is the “center of mass” of the point cloud S' . Since each example in S' is positive, each example in S' must lie on the same side of the hyperplane $\tilde{u} \cdot \tilde{x} = 0$ as the vector \tilde{u} , so clearly \tilde{z} must also lie on this side of the hyperplane. One might even hope that \tilde{z} , or some related vector, points in approximately the same direction as the vector \tilde{u} .

Our p -norm weak learning algorithm, which we call WLA, is presented in figure 2. The vector \tilde{z} is the “center of mass” of the normalized points with respect to the probability distribution \mathcal{D} which is part of the input to WLA (so running WLA repeatedly on the same data set S but with different distributions \mathcal{D} may yield different values for \tilde{z}). Like the online p -norm algorithm, the WLA algorithm transforms the vector \tilde{z} to a vector \tilde{w} using the mapping $w_i = \text{sign}(z_i) |z_i|^{p-1}$. The real-valued WLA hypothesis is a scaled version of the linear functional defined by the vector \tilde{w} . The following theorem establishes that this simple algorithm is in fact a weak learner:

Theorem 3. *WLA is a $(1/2 - \gamma)$ -weak learning algorithm for \tilde{u} under \mathcal{D} for $\gamma = \frac{\delta_{\tilde{u}, \tilde{z}}}{2\|X\|_p \|\tilde{u}\|_q}$.*

Input parameters:real number $p \geq 2$ sequence $S = \langle \tilde{x}^1, y_1 \rangle, \dots, \langle \tilde{x}^m, y_m \rangle$ of labeled examplesdistribution \mathcal{D} over S

1. **set** $\tilde{z} = \sum_{j=1}^m \mathcal{D}(\tilde{x}^j) y_j \tilde{x}^j$
2. **for all** $i = 1, \dots, n$ **set** $w_i = \text{sign}(z_i) |z_i|^{p-1}$
3. **return** hypothesis $h(\tilde{x}) \equiv \frac{\tilde{w} \cdot \tilde{x}}{\|\tilde{w}\|_q \|\tilde{x}\|_p}$

Figure 2. The p -norm weak learning algorithm WLA.

Proof: Let $S = \langle \tilde{x}^1, y_1 \rangle, \dots, \langle \tilde{x}^m, y_m \rangle$ be a sequence of labeled examples where $\tilde{x} \in X$ and $y = \text{sign}(\tilde{u} \cdot \tilde{x})$ for every pair $\langle \tilde{x}, y \rangle \in S$, and let \mathcal{D} be a distribution over S . We will show that the hypothesis h which WLA(p, S, \mathcal{D}) returns is a $(1/2 - \gamma)$ -approximator for \tilde{u} under \mathcal{D} .

To see that h maps X into $[-1, 1]$, note that for any $\tilde{x} \in X$ Hölder's inequality implies

$$|h(\tilde{x})| = \frac{|\tilde{w} \cdot \tilde{x}|}{\|\tilde{w}\|_q \|\tilde{x}\|_p} \leq \frac{\|\tilde{w}\|_q \|\tilde{x}\|_p}{\|\tilde{w}\|_q \|\tilde{x}\|_p} \leq \frac{\|\tilde{w}\|_q \|\tilde{x}\|_p}{\|\tilde{w}\|_q \|\tilde{x}\|_p} = 1.$$

Now we show that inequality (1) from Section 2.1 holds. Since $h(\tilde{x}^j) \in [-1, 1]$ and $y_j \in \{-1, 1\}$ we have that

$$|h(\tilde{x}^j) - y_j| = 1 - y_j h(\tilde{x}^j),$$

and thus

$$\begin{aligned} \frac{1}{2} \sum_{j=1}^m \mathcal{D}(\tilde{x}^j) |h(\tilde{x}^j) - y_j| &= \frac{1}{2} \sum_{j=1}^m \mathcal{D}(\tilde{x}^j) (1 - y_j h(\tilde{x}^j)) \\ &= \frac{1}{2} - \frac{1}{2\|\tilde{x}\|_p} \left(\frac{\sum_{j=1}^m \mathcal{D}(\tilde{x}^j) y_j (\tilde{w} \cdot \tilde{x}^j)}{\|\tilde{w}\|_q} \right). \end{aligned}$$

Thus it suffices to show that

$$\frac{\sum_{j=1}^m \mathcal{D}(\tilde{x}^j) y_j (\tilde{w} \cdot \tilde{x}^j)}{\|\tilde{w}\|_q} \geq \frac{\delta_{\tilde{u}, X}}{\|\tilde{u}\|_q}.$$

We first note that

$$\begin{aligned} \sum_{j=1}^m \mathcal{D}(\tilde{x}^j) y_j (\tilde{w} \cdot \tilde{x}^j) &= \tilde{w} \cdot \left(\sum_{j=1}^m \mathcal{D}(\tilde{x}^j) y_j \tilde{x}^j \right) \\ &= \tilde{w} \cdot \tilde{z} \end{aligned}$$

$$\begin{aligned}
&= \sum_{j=1}^m |z_j|^p \\
&= \|\tilde{z}\|_p^p
\end{aligned}$$

and hence the left-hand side of the desired inequality equals $\|\tilde{z}\|_p^p / \|\tilde{w}\|_q$.

We also have

$$\begin{aligned}
\|\tilde{w}\|_q &= \left(\sum_{i=1}^n (|z_i|^{p-1})^q \right)^{1/q} \\
&= \left(\sum_{i=1}^n |z_i|^p \right)^{1/q} \\
&= \|\tilde{z}\|_p^{p/q},
\end{aligned}$$

where in the second equality we used the fact that $(p-1)q = p$. Consequently the left-hand side can be further simplified to $\|\tilde{z}\|_p^p / \|\tilde{w}\|_q = \|\tilde{z}\|_p^{p-p/q} = \|\tilde{z}\|_p$, and thus our goal is to show that $\|\tilde{z}\|_p \geq \delta_{\tilde{u}, X} / \|\tilde{u}\|_q$. Since $\delta_{\tilde{u}, X} \leq \tilde{u} \cdot (y_j \tilde{x}^j)$ for $j = 1, \dots, m$, we have

$$\begin{aligned}
\delta_{\tilde{u}, X} &\leq \sum_{j=1}^m \mathcal{D}(\tilde{x}^j) \tilde{u} \cdot (y_j \tilde{x}^j) = \tilde{u} \cdot \left(\sum_{j=1}^m \mathcal{D}(\tilde{x}^j) y_j \tilde{x}^j \right) \\
&= \tilde{u} \cdot \tilde{z} \\
&\leq \|\tilde{u}\|_q \|\tilde{z}\|_p,
\end{aligned}$$

where the last line follows from the Hölder inequality, and the theorem is proved. \square

Thus, the simple WLA algorithm can serve as a weak learning algorithm for our halfspace learning problem. In the next section we use techniques from boosting and large margin classification to obtain a strong learning algorithm which has small sample complexity.

4. From weak to strong learning

4.1. Boosting to achieve high accuracy

In a series of important papers Schapire (1990) and Freund (1992, 1995) have given *boosting* algorithms which transform weak learning algorithms into strong ones. Boosting algorithms have since been the focus of intense research activity in both the applied and theoretical machine learning communities.

In this paper we use the Adaboost algorithm from Freund and Schapire (1997) which is shown in figure 3; our notation for the algorithm is similar to that of Schapire et al. (1998) and Schapire and Singer (1998). The input to Adaboost is a sequence $S = \langle x^1, y_1 \rangle, \dots, \langle x^m, y_m \rangle$ of m labeled examples, a weak learning algorithm WL, and two parameters $0 < \gamma, \mu < 1/2$.

Input parameters:sequence $S = \langle x^1, y_1 \rangle, \dots, \langle x^m, y_m \rangle$ of labeled examplesweak learning algorithm WL: $S \rightarrow [-1, 1]$ real values $0 < \gamma, \mu < 1/2$

1. **set** $T = \frac{1}{2\gamma^2} \log \frac{1}{\mu}$
2. **for all** $i = 1, \dots, m$ **set** $\mathcal{D}^1(x^i) = \frac{1}{m}$
3. **for** $t = 1, \dots, T$ **do**
4. **let** h_t be the output of WL(\mathcal{D}^t, S)
5. **set** $\epsilon_t = \frac{1}{2} \sum_{i=1}^m \mathcal{D}^t(x^i) |h_t(x^i) - y_i|$
6. **set** $\alpha_t = \frac{1}{2} \ln((1 - \epsilon_t)/\epsilon_t)$
7. **for all** $i = 1, \dots, m$ **set** $\mathcal{D}^{t+1}(x^i) = \frac{\mathcal{D}^t(x^i) \exp(-y_i \alpha_t h_t(x^i))}{Z_t}$
 where $Z_t = \sum_{i=1}^m \mathcal{D}^t(x^i) \exp(-y_i \alpha_t h_t(x^i))$ is a normalizing
 factor which ensures that \mathcal{D}^{t+1} is a distribution
9. **enddo**
10. **output** as final hypothesis $h(x) \equiv \text{sign}(f(x))$ where

$$f(x) = \frac{\sum_{t=1}^T \alpha_t h_t(x)}{\sum_{t=1}^T \alpha_t}.$$

Figure 3. The Adaboost algorithm.

Given a distribution \mathcal{D}^t over a data set S , algorithm WL outputs a hypothesis h_t which maps S to $[-1, 1]$. Adaboost works in a sequence of stages, where in stage t it generates a distribution \mathcal{D}^t and runs WL to obtain a hypothesis h_t . The final Adaboost hypothesis is a linear threshold function over the h_t s.

In Freund and Schapire (1997) prove that if the algorithm WL is a $(1/2 - \gamma)$ -weak learning algorithm, i.e. each call of WL in Adaboost generates a hypothesis h_t such that ϵ_t (as defined in line 5) is at most $1/2 - \gamma$, then the fraction of examples in S which are misclassified by the final hypothesis h is at most μ . Given this result, one straightforward way to obtain a strong learning algorithm for our halfspace learning problem is to draw a sufficiently large (as specified below) sample S from the example oracle $EX(\tilde{u}, \mathcal{D})$ and run Adaboost on S using WLA as the weak learning algorithm, γ as given in Theorem 3, and $\mu < 1/|S|$. This choice of μ ensures that Adaboost's final hypothesis makes no errors on S ; moreover, since each hypothesis generated by WLA is of the form $h_t(\tilde{x}) = \tilde{v}^t \cdot \tilde{x}$ for some $\tilde{v}^t \in \mathfrak{R}^n$, Adaboost's final hypothesis will be of the form $h(\tilde{x}) = \text{sign}(\tilde{v} \cdot \tilde{x})$ for some $\tilde{v} \in \mathfrak{R}^n$. Using the fact that the Vapnik-Chervonenkis dimension of the class of zero-threshold linear threshold functions over \mathfrak{R}^n is n , the well-known theorem of Blumer et al. (1989) implies that with probability at least $1 - \delta$ the final hypothesis h is an ϵ -accurate hypothesis for \tilde{u} under \mathcal{D} provided that $|S| \geq c(\epsilon^{-1}(n \log(\epsilon^{-1}) + \log(\delta^{-1})))$ for some constant $c > 0$.

This analysis, though attractively simple, yields a rather crude bound on sample complexity which does not depend on the particulars of the learning problem (i.e. \tilde{u} and X). In the rest of this section we use recent results on Adaboost's ability to generate a large-margin classifier and the generalization ability of large-margin classifiers to give a much tighter bound on sample complexity for this learning algorithm.

4.2. Boosting to achieve a large margin

Suppose that $h : X \rightarrow \{-1, 1\}$ is a classifier of the form $h(x) = \text{sign}(f(x))$, where f maps X into $[-1, 1]$. We say that the *margin* of h on a labeled example $\langle x, y \rangle$ is $yf(x)$; note that this quantity is nonnegative if and only if h correctly predicts the label y associated with x . The magnitude of the margin can be viewed as a measure of the confidence with which the classifier makes its prediction on x .

The following theorem, which is an extension of Theorem 5 from Schapire et al. (1998), shows that Adaboost can be used in conjunction with a real-valued weak learner to obtain large-margin hypotheses. The proof is given in Appendix A.

Theorem 4. *Suppose that Adaboost is run on an example sequence $S = \langle x^1, y_1 \rangle, \dots, \langle x^m, y_m \rangle$ using a weak learning algorithm $WL: S \rightarrow [-1, 1]$. Then for any value $\theta \geq 0$ we have*

$$\frac{|\{i \in \{1, 2, \dots, m\} : y_i f(x^i) \leq \theta\}|}{m} \leq 2^T \prod_{t=1}^T \sqrt{\epsilon_t^{1-\theta} (1 - \epsilon_t)^{1+\theta}}.$$

The results of Section 3 imply that if WLA is used as the weak learning algorithm in Adaboost, then the value ϵ_t will always be at most $1/2 - \gamma$, and the upper bound of Theorem 4 becomes $((1 - 2\gamma)^{1-\theta} (1 + 2\gamma)^{1+\theta})^{T/2}$. The following technical lemma is useful:

Lemma 5. $(1 - 4x)^{1-x} (1 + 4x)^{1+x} \leq 1 - 4x^2$ for $0 \leq x \leq 1/4$.

Proof: Using a simple convexity argument, it can be verified that $\alpha^r \leq 1 - (1 - \alpha)r$ for any $\alpha \geq 0$ and any $0 \leq r \leq 1$. This inequality implies that $(1 - 4x)^{1-x} \leq 1 - 4x + 4x^2$ and $(1 + 4x)^x \leq 1 + 4x^2$, so consequently

$$(1 - 4x)^{1-x} (1 + 4x)^{1+x} \leq (1 - 4x + 4x^2)(1 + 4x)(1 + 4x^2),$$

which is at most $1 - 4x^2$ for $0 \leq x \leq 1/4$. \square

If we set $\theta = \gamma/2$ and apply Lemma 5 with $x = \theta$, the upper bound of Theorem 4 becomes $(1 - \gamma^2)^{T/2}$ and we obtain the following:

Corollary 6. *If Adaboost is run on a sequence S of labeled examples drawn from $EX(\tilde{u}, \mathcal{D})$ using WLA as the weak learner, γ as defined in Theorem 3 and $\mu < 1/|S|^4$, then the hypothesis h which Adaboost generates will have margin at least $\gamma/2$ on every example in S .*

Proof: The bound on μ causes T to be greater than $\frac{2}{\gamma^2} \log \frac{1}{|S|}$, and consequently the upper bound of Theorem 4 is less than $1/|S|$. \square

Corollary 6 shows that a judicious choice of parameters for Adaboost enables us to obtain a final hypothesis which has a margin of at least $\gamma/2$ on every example in the training set.

In the next subsection we use Corollary 6 and the theory of large margin classification to establish a bound on the generalization error of this hypothesis in terms of the sample size m .

4.3. Large margins and generalization error

Let \mathcal{F} be a collection of real-valued functions on a set X . A finite set $\{x^1, \dots, x^k\} \subseteq X$ is said to be ξ -shattered by \mathcal{F} if there are real numbers r_1, \dots, r_k such that for all $b = (b_1, \dots, b_k) \in \{-1, 1\}^k$, there is a function $f_b \in \mathcal{F}$ such that for all $i = 1, \dots, k$,

$$f_b(x^i) \begin{cases} \geq r_i + \xi & \text{if } b_i = 1 \\ \leq r_i - \xi & \text{if } b_i = -1. \end{cases}$$

For $\xi \geq 0$, the *fat-shattering dimension of \mathcal{F} at scale ξ* , denoted $\text{fat}_{\mathcal{F}}(\xi)$, is the size of the largest set which is ξ -shattered by \mathcal{F} , if this is finite, and infinity otherwise. The fat-shattering dimension is useful for us because of the following theorem from Bartlett and Shawe-Taylor (1999):

Theorem 7. *Let \mathcal{F} be a collection of real-valued functions on X and let \mathcal{D} be a distribution over $X \times \{-1, 1\}$. Let $S = \langle \tilde{x}^1, y_1 \rangle, \dots, \langle \tilde{x}^m, y_m \rangle$ be a sequence of labeled examples drawn from \mathcal{D} . With probability at least $1 - \delta$ over the choice of S , if a classifier $h(x) \equiv \text{sign}(f(x))$ with $f \in \mathcal{F}$ has margin at least $\xi > 0$ on every example in S , then*

$$\Pr_{(x,y) \in \mathcal{D}} [h(x) \neq y] \leq \frac{2}{m} \left(d \log \frac{8em}{d} \log(32m) + \log \frac{8m}{\delta} \right),$$

where $d = \text{fat}_{\mathcal{F}}(\xi/16)$.

As noted in Section 4.1, the final hypothesis h which Adaboost outputs will be of the form $h(\tilde{x}) = \text{sign}(f(\tilde{x}))$ with $f(\tilde{x}) = \tilde{v} \cdot \tilde{x}$ for some $\tilde{v} \in \mathfrak{R}^n$. Furthermore, since each invocation of WLA generates a hypothesis of the form $h_t(\tilde{x}) = \tilde{v}^t \cdot \tilde{x}$ with $\|\tilde{v}^t\|_q \leq \frac{1}{\|\tilde{x}\|_p}$, Minkowski's inequality implies that the vector \tilde{v} must satisfy $\|\tilde{v}\|_q \leq \frac{1}{\|\tilde{x}\|_p}$. We thus consider the class of functions

$$\mathcal{F} = \left\{ \tilde{x} \mapsto \tilde{v} \cdot \tilde{x} : \|\tilde{v}\|_q \leq \frac{1}{\|\tilde{x}\|_p}, x \in X \right\}. \quad (2)$$

If we can bound $\text{fat}_{\mathcal{F}}(\xi)$, then given any sample size m , Theorem 7 immediately yields a corresponding bound on $\Pr_{x \in \mathcal{D}} [h(\tilde{x}) \neq \text{sign}(\tilde{u} \cdot \tilde{x})]$ for our halfspace learning problem. The following theorem, which is an extension of Theorem 1.6 from Bartlett and Shawe-Taylor (1999), gives the desired bound on $\text{fat}_{\mathcal{F}}(\xi)$. The proof is given in Appendix B.

Theorem 8. *Let X be a bounded region in \mathfrak{R}^n and let \mathcal{F} be the class of functions on X defined in (2) above. Then $\text{fat}_{\mathcal{F}}(\xi) \leq \frac{2 \ln 4n}{\xi^2}$.*

4.4. Putting it all together

Combining Theorem 3, Corollary 6, and Theorems 7 and 8, it follows that if our algorithm uses a sample of size $|S| = m$, then with probability at least $1 - \delta$ the hypothesis h which is generated will satisfy

$$\Pr_{\tilde{x} \in \mathcal{D}} [h(\tilde{x}) \neq \text{sign}(\tilde{u} \cdot \tilde{x})] = O\left(\frac{1}{m} \left(\frac{\|\tilde{u}\|_q^2 \|X\|_p^2}{\delta_{\tilde{u}, X}^2} \log n \log^2 m + \log \frac{m}{\delta} \right)\right).$$

Thus we have established the following (where the \tilde{O} -notation hides log factors):

Theorem 9. *The algorithm obtained by applying Adaboost to WLA using the parameter settings described in Corollary 6 is a strong learning algorithm for \tilde{u} on X with sample complexity*

$$m(\epsilon, \delta, \tilde{u}, X) = \tilde{O}\left(\frac{1}{\epsilon} \cdot \frac{\|\tilde{u}\|_q^2 \|X\|_p^2}{\delta_{\tilde{u}, X}^2}\right).$$

5. Discussion

The sample complexity of our boosting-based p -norm PAC learning algorithm is remarkably similar to that of the PAC-transformed online p -norm algorithms of Section 2.1. Both sets of bounds are essentially linear in ϵ^{-1} and quadratic in $\|\tilde{u}\|_q \|X\|_p / \delta_{\tilde{u}, X}$. Comparing the bounds in more detail, we see that our sample complexity bound contains various log factors which are not present in the bound for the online variant described in part (a) of Theorem 1. These log factors stem from the bounds given in Theorem 7 and Lemma 12; we do not know if they are inherent in the algorithm's performance or an artifact of the analysis. On the other hand, the bound of variant (a) has an extra factor of $p - 1$ which is not present in the sample complexity of our algorithm. Results of Gentile and Littlestone (1999) suggest that the most meaningful range for p is $2 \leq p \leq O(\log n)$, in which case this factor is quite small.

We note that when $p = \Omega(\log n)$ Gentile and Littlestone have given alternative expressions for the online p -norm bounds in terms of $\|X\|_\infty$ and $\|\tilde{u}\|_1$. For these values of p , using an entirely similar analysis the bounds of our algorithm can be analogously rephrased in terms of $\|X\|_\infty$ and $\|\tilde{u}\|_1$ as well.

5.1. $p = 2$ and the Perceptron algorithm

Since the $p = 2$ case of the online p -norm algorithm is precisely the Perceptron algorithm, the $p = 2$ case of our algorithm can be viewed as a natural PAC-model analogue of the online Perceptron algorithm. We note that when $p = 2$ the upper bound given in Lemma 12 can be strengthened to $\sqrt{d} \cdot \|X\|_2$ (see Lemma 1.3 of Bartlett & Shawe-Taylor, 1999 or Theorem 4.1 of Alon, Spencer, & Erdos, 1992 for a proof). This means that the fat-shattering dimension upper bound of Theorem 8 can be improved to $\frac{1}{\epsilon^2}$, which removes a log factor from the bound of Theorem 9; however this bound will still contain various log factors because of the log terms in Theorem 7.

5.2. $p = \infty$ and the Jackson-Craven algorithm

At the other extreme, we can also define a natural $p = \infty$ version of our algorithm. Consider the vectors \tilde{z} and \tilde{w} which are computed by the weak learning algorithm WLA. If we let r be the number of coordinates z_i of \tilde{z} such that $|z_i| = \|\tilde{z}\|_\infty$, then for any i we have

$$\begin{aligned} \lim_{p \rightarrow \infty} \left(\frac{w_i}{\|\tilde{w}\|_q} \right) &= \lim_{p \rightarrow \infty} \left(\frac{\text{sign}(z_i) |z_i|^{p-1}}{\left(\sum_{i=1}^n |z_i|^{(p-1)q} \right)^{1/q}} \right) \\ &= \begin{cases} \text{sign}(z_i)/r & \text{if } |z_i| = \|\tilde{z}\|_\infty \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

Hence it is natural to consider a $p = \infty$ version of WLA, which we denote WLA', in which the vector \tilde{w} is defined by taking $w_i = \text{sign}(z_i)$ if $|z_i| = \|\tilde{z}\|_\infty$ and $w_i = 0$ otherwise. All of our analysis continues to hold for the WLA' algorithm (with minor modifications as sketched below) and we obtain a $p = \infty$ strong learning algorithm:

Claim 10. Theorem 9 holds for $p = \infty$ with WLA' in place of WLA.

Proof: The proof of Theorem 3 (with WLA' in place of WLA) is unchanged up through the point where we must show that

$$\frac{\sum_{i=1}^m \mathcal{D}(\tilde{x}^i) y_i (\tilde{w} \cdot \tilde{x}^i)}{\|\tilde{w}\|_1} \geq \frac{\delta_{\tilde{u}, X}}{\|\tilde{u}\|_1}.$$

The left-hand side of this inequality can be rewritten as

$$\begin{aligned} \frac{\tilde{w} \cdot \tilde{z}}{\|\tilde{w}\|_1} &= \frac{\sum_{|z_i| = \|\tilde{z}\|_\infty} \text{sign}(z_i) z_i}{\sum_{|z_i| = \|\tilde{z}\|_\infty} 1} \\ &= \frac{\sum_{|z_i| = \|\tilde{z}\|_\infty} \|\tilde{z}\|_\infty}{\sum_{|z_i| = \|\tilde{z}\|_\infty} 1} \\ &= \|\tilde{z}\|_\infty, \end{aligned}$$

and hence it suffices to prove that $\|\tilde{z}\|_\infty \geq \delta_{\tilde{u}, X} / \|\tilde{u}\|_1$. This is established at the end of the proof of Theorem 3, so Theorem 3 holds with $p = \infty$ and WLA' substituted for WLA.

The rest of the analysis goes through unchanged except for inequalities (7) and (8) of Lemma 12. Since $\|X\|_\infty = \sup_{\tilde{x} \in X} \max_{j=1, \dots, n} |x_j|$, we have that $Y_j \leq d \|X\|_\infty^2$ for all j , and hence in place of inequalities (7) and (8) we have

$$\begin{aligned} \|\tilde{z}\|_\infty &= \max_j |z_j| \\ &\leq \max_j \sqrt{2Y_j \log 4n} \\ &\leq \sqrt{2d \log 4n} \cdot \|X\|_\infty, \end{aligned}$$

which proves Lemma 12. □

There is a close relationship between this $p = \infty$ algorithm and the work of Jackson and Craven on learning sparse perceptrons (Jackson & Craven, 1996). Note that if $r = 1$, i.e. only one coordinate of \tilde{z} has $|z_i| = \|\tilde{z}\|_\infty$, then the WLA' hypothesis is $h(\tilde{x}) = \frac{\ell}{\|\tilde{x}\|_\infty}$ where ℓ is the signed variable from

$$\{x_1, \dots, x_n, -x_1, \dots, -x_n\}$$

which is most strongly correlated under distribution \mathcal{D} with the value of $\text{sign}(\tilde{u} \cdot \tilde{x})$. This is very similar to the weak learning algorithm used by Jackson and Craven (1996), which takes the single best-correlated literal as its hypothesis (breaking ties arbitrarily).

The proof that this “best-single-literal” algorithm used in Jackson and Craven (1996) is a weak learning algorithm is due to Goldmann, Håstad, and Razborov (1992). However, the proof in Goldmann, Håstad, and Razborov (1992) assumes that the example space X is $\{0, 1\}^n$ and that the target vector \tilde{u} has all integer coefficients; thus, as noted by Jackson and Craven (1996), their algorithm for learning sparse perceptrons only applies to learning problems which are defined over discrete input domains. In contrast, our $p = \infty$ algorithm can be applied on continuous input domains—the only restrictions required by our algorithm are that the example space X and the target vector \tilde{u} satisfy $\|X\|_\infty < \infty$ and $\delta_{\tilde{u}, X} > 0$.

We also observe that Theorem 9 establishes a tighter sample complexity bound for our $p = \infty$ strong learning algorithm than was given in Jackson and Craven (1996). To see this, let $X = \{0, 1\}^n$ and suppose that the target vector $\tilde{u} \in \mathfrak{R}^n$ has all integer coefficients so the algorithm from Jackson and Craven (1996) can be applied. For this learning problem we have $\delta_{\tilde{u}, X} = \Omega(1)$ and $\|X\|_\infty = 1$; letting $s = \|\tilde{u}\|_1$, Theorem 9 implies that our $p = \infty$ strong learning algorithm has sample complexity roughly s^2/ϵ (ignoring log factors). This is essentially the same bound which can be obtained from the Balanced Winnow algorithm of Littlestone (1989) (though somewhat weaker than the bound which can be obtained from the original Winnow algorithm of Littlestone (1988)), and is an improvement over the roughly s^4/ϵ sample complexity bound given in Jackson and Craven (1996).

6. Open questions

Our results give evidence of the broad utility of boosting algorithms such as Adaboost. A natural question is how much further this utility extends: are there simple boosting-based PAC versions of other standard learning algorithms? We note in this context that Kearns and Mansour (1996) have shown that various heuristic algorithms for top-down decision tree induction can be viewed as instantiations of boosting. Another goal is to construct more powerful boosting-based PAC algorithms for linear threshold functions. All of the algorithms discussed in this paper have an inverse quadratic dependence on the separation parameter $\delta_{\tilde{u}, X}$; linear-programming based algorithms for learning linear threshold functions (see, e.g., Blum et al., 1996; Blumer et al., 1989; Cohen, 1997; Long, 1994; Maass & Turan, 1994) do not have such a dependence. Is there a natural boosting-based PAC algorithm for linear threshold functions with performance bounds similar to those of the linear-programming based algorithms?

Appendix A: Proof of Theorem 4

The proof combines ideas from Schapire et al. (1998), where it is shown that Adaboost with binary valued hypotheses generates a large margin classifier, and Schapire and Singer (1998), where an analysis is given for Adaboost's classification error with real valued hypotheses. As in Theorem 5 of Schapire et al. (1998), if $y_i f(x^i) \leq \theta$ then

$$y_i \sum_{t=1}^T \alpha_t h_t(x^i) \leq \theta \sum_{t=1}^T \alpha_t$$

which implies that

$$\exp\left(-y_i \sum_{t=1}^T \alpha_t h_t(x^i) + \theta \sum_{t=1}^T \alpha_t\right) \geq 1.$$

Following Schapire et al. (1998), we thus have that

$$\begin{aligned} & \frac{|\{i \in \{1, 2, \dots, m\} : y_i f(x^i) \leq \theta\}|}{m} \\ & \leq \sum_{i=1}^m \frac{1}{m} \cdot \left[\exp\left(-y_i \sum_{t=1}^T \alpha_t h_t(x^i) + \theta \sum_{t=1}^T \alpha_t\right) \right] \\ & = \frac{\exp\left(\theta \sum_{t=1}^T \alpha_t\right)}{m} \sum_{i=1}^m \exp\left(-y_i \sum_{t=1}^T \alpha_t h_t(x^i)\right) \\ & = \exp\left(\theta \sum_{t=1}^T \alpha_t\right) \left(\prod_{t=1}^T Z_t\right) \sum_{i=1}^m \mathcal{D}^{T+1}(x^i) \\ & = \exp\left(\theta \sum_{t=1}^T \alpha_t\right) \left(\prod_{t=1}^T Z_t\right) \end{aligned} \quad (3)$$

where the second equality follows from the definition of \mathcal{D}^{T+1} and the final equality is because \mathcal{D}^{T+1} is a distribution and hence sums to 1. Our goal is thus to bound the right side of inequality (3).

If we let

$$r_t = \sum_{i=1}^m \mathcal{D}^t(x^i) y_i h_t(x^i)$$

then using the fact that

$$|h(x^j) - y_j| = 1 - y_j h(x^j)$$

we find that $\epsilon_t = \frac{1-r_t}{2}$. Substituting into the definition of α_t , we obtain

$$\alpha_t = \frac{1}{2} \ln\left(\frac{1+r_t}{1-r_t}\right).$$

Following (Schapire & Singer, 1998) for simplicity of notation we now fix t and let $u_i = y_i h_t(x^i)$, $Z = Z_t$, $\mathcal{D} = \mathcal{D}'$, $\epsilon = \epsilon_t$, $r = r_t$, and $\alpha = \alpha_t$. As noted in Schapire and Singer (1998) a simple convexity argument shows that

$$e^{-\alpha u} \leq \frac{1+u}{2} e^{-\alpha} + \frac{1-u}{2} e^{\alpha}$$

for any $\alpha \in \Re$ and any $u \in [-1, 1]$. Since u_i always lies in the interval $[-1, 1]$, we can apply this inequality to obtain

$$\begin{aligned} Z &= \sum_{i=1}^m \mathcal{D}(x^i) e^{-\alpha u_i} \\ &\leq \sum_{i=1}^m \mathcal{D}(x^i) \left(\frac{1+u_i}{2} e^{-\alpha} + \frac{1-u_i}{2} e^{\alpha} \right). \end{aligned} \quad (4)$$

As in Section 3.5 of Schapire and Singer (1998), substituting α into inequality (4) yields

$$\begin{aligned} Z_t &\leq \sqrt{1 - r_t^2} \\ &= \sqrt{1 - (1 - 2\epsilon_t)^2} \\ &= 2\sqrt{\epsilon_t(1 - \epsilon_t)}. \end{aligned} \quad (5)$$

Substituting inequality (5) into inequality (3) and using the definition of α_t yields the desired bound of the theorem. (Theorem 4) \square

Appendix B: Proof of Theorem 8

Theorem 8 follows immediately upon combining the inequalities proved in the following two lemmas.

Lemma 11. *If the set $\{\tilde{x}^1, \dots, \tilde{x}^d\}$ is ξ -shattered by \mathcal{F} then every $b = (b_1, \dots, b_d) \in \{-1, 1\}^d$ satisfies*

$$\left\| \sum_{i=1}^d b_i \tilde{x}^i \right\|_p \geq \xi d \|X\|_p.$$

Proof: Suppose that $\{\tilde{x}^1, \dots, \tilde{x}^d\}$ is ξ -shattered by \mathcal{F} as witnessed by the real numbers r_1, \dots, r_d . Then for every $b = (b_1, \dots, b_d) \in \{-1, 1\}^d$, there is a vector $\tilde{v}_b \in \Re^n$ with $\|\tilde{v}_b\|_q \leq \frac{1}{\|X\|_p}$ such that $b_i (\tilde{v}_b \cdot \tilde{x}^i - r_i) \geq \xi$ for $i = 1, \dots, d$. Summing these d inequalities and rearranging, we obtain

$$\tilde{v}_b \cdot \left(\sum_{i=1}^d b_i \tilde{x}^i \right) \geq \xi d + \sum_{i=1}^d b_i r_i. \quad (6)$$

There are two cases to consider. Case 1 is if $\sum_{i=1}^d b_i r_i \geq 0$; if this is true, we have

$$\begin{aligned} \frac{1}{\|X\|_p} \cdot \left\| \sum_{i=1}^d b_i \tilde{x}^i \right\|_p &\geq \|\tilde{v}_b\|_q \left\| \sum_{i=1}^d b_i \tilde{x}^i \right\|_p \\ &\geq \tilde{v}_b \cdot \left(\sum_{i=1}^d b_i \tilde{x}^i \right) \\ &\geq \xi d \end{aligned}$$

(where the first inequality is by the definition of \tilde{v}_b , the second inequality is Hölder's, and the third is from inequality (6)), which yields the desired inequality $\|\sum_{i=1}^d b_i \tilde{x}^i\|_p \geq \xi d \|X\|_p$.

In the second case, $\sum_{i=1}^d b_i r_i < 0$. If this is the case then let $c = (c_1, \dots, c_d) = (-b_1, \dots, -b_d)$. We then have $\sum_{i=1}^d c_i r_i > 0$, so Case 1 implies that $\|\sum_{i=1}^d c_i \tilde{x}^i\|_p \geq \xi d \|X\|_p$, and the lemma follows since

$$\left\| \sum_{i=1}^d c_i \tilde{x}^i \right\|_p = \left\| - \sum_{i=1}^d b_i \tilde{x}^i \right\|_p = \left\| \sum_{i=1}^d b_i \tilde{x}^i \right\|_p.$$

□

Lemma 12. For any set $\{\tilde{x}^1, \dots, \tilde{x}^d\}$ with each $\|\tilde{x}^i\|_p \leq \|X\|_p$, if $p \geq 2$ then there is some $b = (b_1, \dots, b_d) \in \{-1, 1\}^d$ such that $\|\sum_{i=1}^d b_i \tilde{x}^i\|_p \leq \sqrt{2d \ln 4n} \cdot \|X\|_p$.

Proof: The proof uses the probabilistic method. We consider the random variable $\tilde{z} = \sum_{i=1}^d b_i \tilde{x}^i$ where (b_1, \dots, b_d) is uniformly distributed over $\{-1, 1\}^d$. For any coordinate $j \in \{1, \dots, n\}$ we have $z_j = \sum_{i=1}^d b_i x_j^i$ and hence $E[z_j] = 0$. Let $Y_j = |x_j^1|^2 + \dots + |x_j^d|^2$; Hoeffding's bound (Hoeffding, 1963) on sums of independent random variables states that for any $t > 0$ we have

$$\Pr[|z_j| > t] \leq 2 \exp\left(\frac{-t^2}{2Y_j}\right).$$

As a consequence, taking $t = \sqrt{2Y_j \ln 4n}$ we have that $\Pr[|z_j| \geq t] \leq 1/2n$. Using the union bound across $j = 1, 2, \dots, n$, we have that with probability at least $1/2$ every coordinate z_j of \tilde{z} satisfies $|z_j| < \sqrt{2Y_j \ln 4n}$, and hence

$$\begin{aligned} \|\tilde{z}\|_p &= \left(\sum_{j=1}^n |z_j|^p \right)^{1/p} \\ &\leq \left(\sum_{j=1}^n (\sqrt{2Y_j \ln 4n})^p \right)^{1/p} \\ &= \sqrt{2 \ln 4n} \cdot \left(\left(\sum_{j=1}^n [|x_j^1|^2 + \dots + |x_j^d|^2]^{p/2} \right)^{2/p} \right)^{1/2} \end{aligned} \quad (7)$$

Since $p \geq 2$, we have $p/2 \geq 1$ and hence Minkowski's inequality implies that

$$\begin{aligned}
 & \left(\sum_{j=1}^n [|x_j^1|^2 + \dots + |x_j^d|^2]^{p/2} \right)^{2/p} \\
 & \leq \left[\sum_{j=1}^n |x_j^1|^{2p/2} \right]^{2/p} + \dots + \left[\sum_{j=1}^n |x_j^d|^{2p/2} \right]^{2/p} \\
 & = \|\tilde{x}^1\|_p^2 + \dots + \|\tilde{x}^d\|_p^2 \\
 & \leq d\|X\|_p^2.
 \end{aligned} \tag{8}$$

The lemma follows by combining inequalities (7) and (8). \square

Acknowledgments

We thank Les Valiant and the two anonymous referees for helpful comments and suggestions.

References

- Alon, N., Spencer, J., & Erdos, P. (1992). *The probabilistic method*. New York: Wiley-Interscience.
- Angluin, D. (1988). Queries and concept learning. *Machine Learning*, 2, 319–342.
- Auer, P., & Warmuth, M. (1995). Tracking the best disjunction. In *Proceedings of the 36th Symposium on Foundations of Computer Science* (pp. 312–321).
- Bartlett, P., & Shawe-Taylor, J. (1999). Generalization performance of support vector machines and other pattern classifiers. In B. Scholkopf, C. J. C. Burges, & A. J. Smola (Eds.), *Advances in Kernel methods—Support vector learning*. Cambridge, MA: MIT Press.
- Baum, E. (1990). The perceptron algorithm is fast for nonmalicious distributions. *Neural Computation*, 2:2, 248–260.
- Blum, A., Frieze, A., Kannan, R., & Vempala, S. (1996). A polynomial time algorithm for learning noisy linear threshold functions. In *Proceedings of the 37th Symposium on Foundations of Computer Science* (pp. 330–338).
- Blumer, A., Ehrenfeucht, A., Haussler, D., & Warmuth, M. (1989). Learnability and the Vapnik-Chervonenkis dimension. *Journal of the ACM*, 36:4, 929–965.
- Bylander, T. (1998). Worst-case analysis of the perceptron and exponentiated update algorithms. *Artificial Intelligence*, 106:2, 335–352.
- Cohen, E. (1997). Learning noisy perceptrons by a perceptron in polynomial time. In *Proceedings of the 38th Symposium on Foundations of Computer Science* (pp. 514–523).
- Freund, Y. (1992). An improved boosting algorithm and its implications on learning complexity. In *Proceedings of the Fifth Annual Workshop on Computational Learning Theory* (pp. 391–398).
- Freund, Y. (1995). Boosting a weak learning algorithm by majority. *Information and Computation*, 121:2, 256–285.
- Freund, Y., & Schapire, R. (1996). Game theory, on-line prediction and boosting. In *Proceedings of the Ninth Annual Conference on Computational Learning Theory* (pp. 325–332).
- Freund, Y., & Schapire, R. (1997). A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55:1, 119–139.
- Freund, Y., & Schapire, R. (1999). Large margin classification using the perceptron algorithm. *Machine Learning*, 37:3, 277–296.
- Gentile, C., & Littlestone, N. (1999). The robustness of the p -norm algorithms. In *Proceedings of the Twelfth Annual Conference on Computational Learning Theory* (pp. 1–11).

PAC ANALOGUES OF WINNOW AND PERCEPTRON VIA BOOSTING THE MARGIN 151

- Grove, A., Littlestone, N., & Schuurmans, D. (1997). General convergence results for linear discriminant updates. In *Proceedings of the Tenth Annual Conference on Computational Learning Theory* (pp. 171–183).
- Goldmann, M., Håstad, J., & Razborov, R. (1992). Majority gates vs. general weighted threshold gates. *Computational Complexity*, 2, 277–300.
- Haussler, D. (1988). Space efficient learning algorithms. Technical Report UCSC-CRL-88-2, University of California, Santa Cruz.
- Hoeffding, W. (1963). Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58, 13–30.
- Jackson, J., & Craven, M. (1996). Learning sparse perceptrons. *Advances in neural information processing systems* 8, Cambridge, MA: MIT Press.
- Ji, C., & Ma, S. (1997). Combinations of weak classifiers. *IEEE Transactions on Neural Networks*, 8:1, 32–42.
- Kearns, M., Li, M., Pitt, L., & Valiant, L. (1987). Recent results on boolean concept learning. In *Proceedings of the Fourth International Workshop on Machine Learning* (pp. 337–352).
- Kearns, M., & Mansour, Y. (1996). On the boosting ability of top-down decision tree learning algorithms. In *Proceedings of the 28th Symposium on Theory of Computing* (pp. 459–468).
- Kearns, M., & Valiant, L. (1994). Cryptographic limitations on learning boolean formulae and finite automata. *J. ACM*, 41:1, 67–95.
- Kivinen, J., Warmuth, M., & Auer, P. (1997). The perceptron algorithm versus winnow: Linear vs. logarithmic mistake bounds when few input variables are relevant. *Artificial Intelligence*, 97:1/2, 325–343.
- Littlestone, N. (1988). Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm. *Machine Learning*, 2, 285–318.
- Littlestone, N. (1989). Mistake bounds and logarithmic linear-threshold learning algorithms. Ph.D. thesis, Technical Report UCSC-CRL-89-11, University of California, Santa Cruz.
- Littlestone, N. (1989). From online to batch learning. In *Proceedings of the Second Annual Workshop on Computational Learning Theory* (pp. 269–284).
- Littlestone, N. (1991). Redundant Noisy attributes, attribute errors, and linear-threshold learning using winnow. In *Proceedings of the Fourth Annual Conference on Computational Learning Theory* (pp. 147–156).
- Long, P. (1994). Halfspace learning, linear programming, and nonmalicious distributions. *Information Processing Letters*, 51, 245–250.
- Maass, W., & Turan, G. (1994). How fast can a threshold gate learn? In S. J. Hanson, G. Drastal, & R. Rivest (Eds.), *Computational learning theory and natural learning systems: Volume 1: Constraints and prospects*. Cambridge, MA: MIT Press.
- Schapire, R. (1990). The strength of weak learnability. *Machine Learning*, 5:2, 197–227.
- Schapire, R. (1999). Drifting games. In *Proceedings of the Twelfth Annual Conference on Computational Learning Theory* (pp. 114–124).
- Servedio, R. (1999). On PAC learning using Winnow, Perceptron, and a Perceptron-like algorithm. In *Proceedings of the Twelfth Annual Conference on Computational Learning Theory* (pp. 296–307).
- Schapire, R., Freund, Y., Bartlett, P., & Lee, W. S. (1998). Boosting the margin: A new explanation for the effectiveness of voting methods. *Annals of Statistics*, 26:5, 1651–1686.
- Schapire, R., & Singer, Y. (1998). Improved boosting algorithms using confidence-rated predictions. In *Proceedings of the Eleventh Annual Conference on Computational Learning Theory* (pp. 80–91).
- Schmitt, M. (1998). Identification criteria and lower bounds for Perceptron-like learning rules. *Neural Computation*, 10, 235–250.
- Taylor, A., & Mann, W. (1972). *Advanced calculus*. New York: John Wiley & Sons.
- Vapnik, V. (1998). *Statistical learning theory*. New York: John Wiley & Sons.

Received March 21, 2001

Revised March 21, 2001

Accepted June 22, 2001

Final manuscript September 1, 2001