

A Nearly Optimal Variant of the Perceptron Algorithm for the Uniform Distribution on the Unit Sphere

Marco Schmalhofer

SCHMALHOFER@EM.UNI-FRANKFURT.DE

Institute of Computer Science, Goethe University Frankfurt

Editors: Jacob Abernethy and Shivani Agarwal

Abstract

We show a simple perceptron-like algorithm to learn origin-centered halfspaces in \mathbb{R}^n with accuracy $1 - \epsilon$ and confidence $1 - \delta$ in time

$$\mathcal{O}\left(\frac{n^2}{\epsilon} \left(\log \frac{1}{\epsilon} + \log \frac{1}{\delta}\right)\right)$$

using

$$\mathcal{O}\left(\frac{n}{\epsilon} \left(\log \frac{1}{\epsilon} + \log \frac{1}{\delta}\right)\right)$$

labeled examples drawn uniformly from the unit n -sphere. This improves upon algorithms given in [Baum \(1990\)](#), [Long \(1994\)](#) and [Servedio \(1999\)](#). The time and sample complexity of our algorithm match the lower bounds given in [Long \(1995\)](#) up to logarithmic factors.

Keywords: Halfspace learning, perceptron, uniform distribution, n -sphere

1. Introduction

Learning halfspaces from labeled examples is one of the central challenges in machine learning. In [Blumer et al. \(1989\)](#) it is shown that n -dimensional halfspaces can be learned to accuracy $1 - \epsilon$ with confidence $1 - \delta$ in the classical PAC model, and hence for arbitrary distributions, using $\mathcal{O}((n/\epsilon) \log(1/\epsilon) + (1/\epsilon) \log(1/\delta))$ examples. Therefore it suffices to find a halfspace consistent with the given examples, which can be accomplished in time polynomial in n , $1/\epsilon$ and $1/\delta$ (e.g. by linear programming). In [Ehrenfeucht et al. \(1989\)](#) a lower bound of $\Omega((n/\epsilon) + (1/\epsilon) \log(1/\delta))$ on the number of examples is derived, which also holds if the examples are drawn uniformly from the unit sphere, see [Long \(1995\)](#). In this case the bound is even tight ([Long \(2003\)](#)). In [Balcan and Long \(2013\)](#) polynomial time algorithms are constructed which achieve that bound even for any log-concave distribution.

The classical perceptron algorithm by Rosenblatt ([Rosenblatt \(1958\)](#)) determines a consistent halfspace given sufficiently many correctly classified examples (see e.g. [Novikoff \(1962\)](#)). Furthermore, in [Baum \(1990\)](#) a variant of the perceptron algorithm was provided, which learns halfspaces in time $\tilde{\mathcal{O}}(n^2/\epsilon^3)$ using $\tilde{\mathcal{O}}(n/\epsilon^3)$ examples. This was improved by [Servedio \(1999\)](#). The algorithm proposed there achieves time complexity $\tilde{\mathcal{O}}(n^2/\epsilon^2)$ with a sample size of $\tilde{\mathcal{O}}(n/\epsilon^2)$. The perceptron algorithm was also shown to be able to solve linear programs in polynomial time, see [Dunagan and Vempala \(2004\)](#). Table 1 summarizes related work considering halfspace learning. In the right column uniform means uniform on the unit sphere.

In this paper we show that the classical perceptron algorithm can be supplemented with an adaptive learning rate such that it (ϵ, δ) -learns halfspaces in time $\mathcal{O}\left(\frac{n^2}{\epsilon} \left(\log \frac{1}{\epsilon} + \log \frac{1}{\delta}\right)\right)$ using $\mathcal{O}\left(\frac{n}{\epsilon} \left(\log \frac{1}{\epsilon} + \log \frac{1}{\delta}\right)\right)$ examples drawn uniformly from the unit sphere. The extremely simple algorithm has nice properties, its error is monotonically decreasing, its hypothesis always has norm one even without rebalancing and it is conservative, i.e. it updates its hypotheses only for counterexamples.

Table 1: Related work on (ϵ, δ) -learning of halfspaces.

article	sample complexity	time complexity	distribution
Blumer et al. (1989)	$\mathcal{O}\left(\frac{1}{\epsilon} \left(n \log \frac{1}{\epsilon} + \log \frac{1}{\delta}\right)\right)$	poly	arbitrary
Ehrenfeucht et al. (1989)	$\Omega\left(\frac{1}{\epsilon} \left(n + \log \frac{1}{\delta}\right)\right)$	–	arbitrary
Haussler et al. (1994)	$\mathcal{O}\left(\frac{n}{\epsilon} \log \frac{1}{\delta}\right)$	poly	arbitrary
Baum (1990)	$\tilde{\mathcal{O}}(n/\epsilon^3)$	$\tilde{\mathcal{O}}(n^2/\epsilon^3)$	uniform
Long (1994) (for $\delta = \epsilon$)	$\tilde{\mathcal{O}}(n/\epsilon)$	$\tilde{\mathcal{O}}(n^2/\epsilon + n^{3.38})$	uniform
Long (1995)	$\Omega\left(\frac{1}{\epsilon} \left(n + \log \frac{1}{\delta}\right)\right)$	–	uniform
Servedio (1999)	$\tilde{\mathcal{O}}(n/\epsilon^2)$	$\tilde{\mathcal{O}}(n^2/\epsilon^2)$	uniform
Long (2003)	$\mathcal{O}\left(\frac{1}{\epsilon} \left(n + \log \frac{1}{\delta}\right)\right)$	–	uniform
Balcan and Long (2013)	$\mathcal{O}\left(\frac{1}{\epsilon} \left(n + \log \frac{1}{\delta}\right)\right)$	poly	log-concave
our paper	$\mathcal{O}\left(\frac{n}{\epsilon} \left(\log \frac{1}{\epsilon} + \log \frac{1}{\delta}\right)\right)$	$\mathcal{O}\left(\frac{n^2}{\epsilon} \left(\log \frac{1}{\epsilon} + \log \frac{1}{\delta}\right)\right)$	uniform

We present basic definitions and notations in Section 2 and motivate the algorithm in Section 3, where first properties and experimental results are also presented. Finally, in Section 4 we derive Theorem 2 as our main result. Its proof crucially depends on Lemma 6, which provides the conditional expectation $\mathbb{E}[dd^* \mid \beta]$, where d and d^* are the distances of a randomly drawn counterexample to the current hyperplane and the target hyperplane, respectively, assuming β is the angle between them. In Section 5 we summarize our results and suggest some open problems.

2. Preliminaries

We study the classical problem of learning *homogeneous halfspaces*

$$f_{\mathbf{w}} : \mathbb{R}^n \rightarrow \{-1, 1\}, f_{\mathbf{w}}(\mathbf{x}) = \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle)$$

represented by a *weight vector* $\mathbf{w} \in \mathbb{R}^n$. We denote the unknown target halfspace by f^* and its normalized weight vector by \mathbf{w}^* . The learner is given *labeled examples* of the form

$$(\mathbf{x}, f^*(\mathbf{x})) \in \mathbb{R}^n \times \{-1, 1\},$$

where each example \mathbf{x} is drawn independently according to the *uniform distribution* on the *unit n -sphere*

$$S^{n-1} := \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| = 1\}.$$

An example \mathbf{x} is called *positive* if $f^*(\mathbf{x}) = 1$ and *negative* otherwise. After receiving examples the learner outputs a *hypothesis* $\mathbf{w} \in \mathbb{R}^n$. The error of \mathbf{w} is measured by the probability of misclassifying a randomly drawn example, i.e.

$$\text{err}(\mathbf{w}) := \mathbb{P}[f^*(\mathbf{x}) \neq f_{\mathbf{w}}(\mathbf{x})],$$

where \mathbf{x} is drawn uniformly from S^{n-1} . By rotational symmetry it is easy to see that the error of a hypothesis \mathbf{w} is determined by the angle $\angle(\mathbf{w}^*, \mathbf{w})$ between \mathbf{w}^* and \mathbf{w} , i.e.

$$\text{err}(\mathbf{w}) = \frac{\angle(\mathbf{w}^*, \mathbf{w})}{\pi} = \frac{1}{\pi} \arccos \frac{\langle \mathbf{w}^*, \mathbf{w} \rangle}{\|\mathbf{w}\|}. \quad (1)$$

3. Adaptive Perceptron – The Algorithm

In this section we present our algorithm, which is in fact the classical perceptron learning rule supplemented with a variable learning rate $\eta > 0$. First, let us consider the single perceptron update

$$\mathbf{w}' = \mathbf{w} + \eta b \mathbf{x} \quad (2)$$

of a hypothesis $\mathbf{w} \neq \mathbf{0}$ through a counterexample $(\mathbf{x}, b) \in S^{n-1} \times \{-1, 1\}$. Then according to Equation (1) the error of the updated hypothesis \mathbf{w}' is

$$\text{err}(\mathbf{w}') = \frac{1}{\pi} \arccos \frac{\langle \mathbf{w}^*, \mathbf{w} + \eta b \mathbf{x} \rangle}{\|\mathbf{w} + \eta b \mathbf{x}\|}. \quad (3)$$

Now we minimize $\text{err}(\mathbf{w}')$ as a function of η . Since \arccos is monotonically decreasing we find a global minimum of $\text{err}(\mathbf{w}')$ by maximizing its argument $g(\eta) := \langle \mathbf{w}^*, \mathbf{w} + \eta b \mathbf{x} \rangle / \|\mathbf{w} + \eta b \mathbf{x}\|$. We determine a zero of g' by applying the quotient rule and forgetting its denominator:

$$\begin{aligned} 0 &\stackrel{!}{=} \frac{d}{d\eta} \left(\frac{\langle \mathbf{w}^*, \mathbf{w} + \eta b \mathbf{x} \rangle}{\|\mathbf{w} + \eta b \mathbf{x}\|} \right) \\ \Leftrightarrow 0 &= \frac{d \langle \mathbf{w}^*, \mathbf{w} + \eta b \mathbf{x} \rangle}{d\eta} \|\mathbf{w} + \eta b \mathbf{x}\| - \frac{d \|\mathbf{w} + \eta b \mathbf{x}\|}{d\eta} \langle \mathbf{w}^*, \mathbf{w} + \eta b \mathbf{x} \rangle \\ \Leftrightarrow 0 &= b \langle \mathbf{w}^*, \mathbf{x} \rangle \|\mathbf{w} + \eta b \mathbf{x}\| - \frac{b \langle \mathbf{w}, \mathbf{x} \rangle + \eta}{\|\mathbf{w} + \eta b \mathbf{x}\|} \langle \mathbf{w}^*, \mathbf{w} + \eta b \mathbf{x} \rangle \end{aligned} \quad (4)$$

By setting $d^* := b \langle \mathbf{w}^*, \mathbf{x} \rangle$ and $d := -b \langle \mathbf{w}, \mathbf{x} \rangle / \|\mathbf{w}\|$ for the distances from \mathbf{x} to the target hyperplane and the actual hyperplane we obtain from Equation (4)

$$\begin{aligned} \left(\|\mathbf{w}\|^2 - \eta \|\mathbf{w}\| d \right) d^* &= (\eta - \|\mathbf{w}\| d) \langle \mathbf{w}^*, \mathbf{w} \rangle \\ \Leftrightarrow \eta &= \|\mathbf{w}\| \frac{d^* + d \langle \mathbf{w}^*, \frac{\mathbf{w}}{\|\mathbf{w}\|} \rangle}{\langle \mathbf{w}^*, \frac{\mathbf{w}}{\|\mathbf{w}\|} \rangle + d d^*}. \end{aligned}$$

Now with $\langle \mathbf{w}^*, \mathbf{w} / \|\mathbf{w}\| \rangle = \cos \beta$, where β is the angle between \mathbf{w} and \mathbf{w}^* , we get the locally optimal learning rate as a function of $\|\mathbf{w}\|$, d^* , d and β , namely

$$\eta_{opt} = \|\mathbf{w}\| \frac{d^* + d \cos \beta}{\cos \beta + d d^*}. \quad (5)$$

Of course this learning rate is useless in practice, since it depends on d^* and β , which are unknown to the algorithm. Nevertheless it motivates a useful choice of η : Assume \mathbf{w} has a small error. Then

β is small and thus $\cos \beta$ close to one. The expected distances $\mathbb{E}[d^*]$ and $\mathbb{E}[d]$ should therefore be also small. Since by symmetry these expected distances are equal, so

$$\eta := \|\mathbf{w}\| 2d \tag{6}$$

is hopefully a good and certainly an easily computable choice. This directly provides our algorithm:

Algorithm 1: ADAPTIVEPERCEPTRON

Input: Number of examples s to be drawn uniformly from S^{n-1}

Output: Hypothesis \mathbf{w}

Get first labeled example (\mathbf{x}, b) where $b = f^*(\mathbf{x})$;

$\mathbf{w} \leftarrow b\mathbf{x}$;

for $i = 2 \dots s$ **do**

 Get labeled example (\mathbf{x}, b) ;

if $f_{\mathbf{w}}(\mathbf{x}) \neq b$ **then**

$\mathbf{w}' \leftarrow \mathbf{w} - 2 \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{x}$;

$\mathbf{w} \leftarrow \mathbf{w}'$;

end

end

return \mathbf{w} ;

Suprisingly, the error of ADAPTIVEPERCEPTRON is monotonically decreasing. Moreover, all hypotheses are unit vectors. These two properties turn out to be crucial for the analysis.

Proposition 1 (first properties) *For any hypothesis \mathbf{w} determined by ADAPTIVEPERCEPTRON,*

(a) $\|\mathbf{w}\| = 1$, as well as

(b) $\text{err}(\mathbf{w}') \leq \text{err}(\mathbf{w}) \leq 1/2$.

Proof After the first example \mathbf{x} we have $\|\mathbf{w}\| = \|b\mathbf{x}\| = 1$. Also $\text{err}(\mathbf{w}) = \frac{1}{\pi} \arccos \langle \mathbf{w}^*, \mathbf{w} \rangle \leq 1/2$, since $\langle \mathbf{w}^*, b\mathbf{x} \rangle \geq 0$. Moreover, after updating \mathbf{w} with a counterexample \mathbf{x} we have

$$\|\mathbf{w}'\|^2 = \|\mathbf{w} - 2 \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{x}\|^2 = \|\mathbf{w}\|^2 - 4 \langle \mathbf{w}, \mathbf{x} \rangle \langle \mathbf{w}, \mathbf{x} \rangle + 4 \langle \mathbf{w}, \mathbf{x} \rangle^2 \|\mathbf{x}\|^2 = \|\mathbf{w}\|^2 = 1.$$

For the error of \mathbf{w}' we have

$$\begin{aligned} \text{err}(\mathbf{w}') &= \frac{1}{\pi} \arccos \langle \mathbf{w}^*, \mathbf{w}' \rangle \\ &= \frac{1}{\pi} \arccos(\langle \mathbf{w}^*, \mathbf{w} \rangle - 2 \langle \mathbf{w}, \mathbf{x} \rangle \langle \mathbf{w}^*, \mathbf{x} \rangle) \\ &\leq \frac{1}{\pi} \arccos \langle \mathbf{w}^*, \mathbf{w} \rangle = \text{err}(\mathbf{w}), \end{aligned}$$

since \arccos is monotonically decreasing and $\langle \mathbf{w}, \mathbf{x} \rangle \langle \mathbf{w}^*, \mathbf{x} \rangle \leq 0$ for the counterexample \mathbf{x} . ■

We conducted experiments for dimension $n = 2^{10}$ and up to $s = 10^{10}$ examples, comparing the learning curves of the perceptron algorithm of Baum (1990), the average algorithm of Serfaty (1999) and ADAPTIVEPERCEPTRON. It turns out that for up to $s \approx n$ examples the three

algorithms do not differ significantly with ADAPTIVEPERCEPTRON even lagging behind. However after 10^5 examples asymptotics seems to take over, ADAPTIVEPERCEPTRON clearly pulls ahead and continues to stay in front from there on (see Figure 1).

We also investigated the “hypothetical” OPTADAPTIVEPERCEPTRON, which in each step uses the locally optimal learning rate η_{opt} (see Equation (5)). Observe that ADAPTIVEPERCEPTRON and OPTADAPTIVEPERCEPTRON are almost indistinguishable.

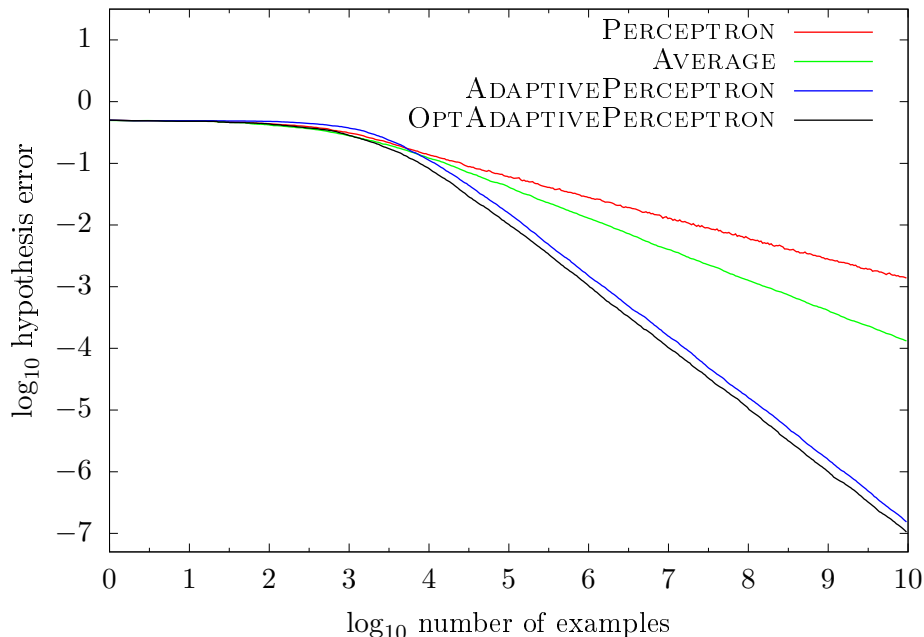


Figure 1: log-log plot of the learning curve for different algorithms at $n = 1024$ dimensions.

4. Adaptive Perceptron – Analysis

In this section we prove our main result:

Theorem 2 *After $s = \Theta\left(\frac{n}{\epsilon}(\log \frac{1}{\epsilon} + \log \frac{1}{\delta})\right)$ examples ADAPTIVEPERCEPTRON outputs a hypothesis which with probability at least $1 - \delta$ has error at most ϵ for each $\epsilon, \delta \in (0, 1]$, $n \geq 2$.*

Note that our simulation is consistent with the results of Theorem 2, i.e. the error of ADAPTIVEPERCEPTRON behaves asymptotically as s^{-1} , whereas the error of PERCEPTRON and AVERAGE is roughly $s^{-1/3}$ and $s^{-1/2}$, respectively (see Figure 1 again).

Proof [Theorem 2] Let w_k be the k -th hypothesis determined by ADAPTIVEPERCEPTRON for $k \geq 1$. We write $\beta_k = \angle(w^*, w_k)$ for the angle between w^* and w_k . The theorem follows if the expected cosine of β_k is “exponentially close” to one, i.e. we later show the following lemma.

Lemma 3 (expected cosine of hypothesis angle) *For each $k \geq 1$ we have*

$$\mathbb{E}[\cos \beta_k] \geq 1 - e^{-\frac{2}{3n}(k-1)},$$

where the expectation is taken over the random sequence of examples.

4.1. Proof of Theorem 2 with Lemma 3

We show the theorem with the help of two inequalities.

Lemma 4 For each $0 < z \leq \frac{\pi}{2}$ we have

$$\frac{\sin z}{z} - \cos z \geq \frac{1}{3}(1 - \cos z), \quad (7)$$

and for each $0 \leq y \leq 1$ we have

$$\arccos y \leq 2\sqrt{1-y}. \quad (8)$$

Proof Let $0 < z \leq \pi/2$. We show $\frac{\sin z}{z} - \frac{2}{3}\cos z - \frac{1}{3} \geq 0$ using Taylor approximations of sine and cosine. With Lagrange remainder we have

$$\sin z = z - \frac{1}{6}z^3 + \frac{\cos \xi_1}{120}z^5 \quad \text{and} \quad \cos z = 1 - \frac{1}{2}z^2 + \frac{\cos \xi_2}{24}z^4$$

with constants $0 \leq \xi_1, \xi_2 \leq z$. Since $\cos \xi_1 \geq 0$ and $\cos \xi_2 \leq 1$, we may conclude

$$\sin z \geq z - \frac{1}{6}z^3 \quad \text{and} \quad \cos z \leq 1 - \left(\frac{1}{2} - \frac{1}{24}z^2\right)z^2 \leq 1 - \frac{1}{3}z^2, \quad (*)$$

where the last inequality follows from the fact that $0 < z < 2$. Hence we have

$$\frac{\sin z}{z} - \frac{2}{3}\cos z - \frac{1}{3} \geq 1 - \frac{1}{6}z^2 - \frac{2}{3} + \frac{2}{9}z^2 - \frac{1}{3} = \frac{1}{18}z^2 \geq 0$$

and (7) is proven. To prove (8) one can solve (*) for z . This yields

$$z \leq \sqrt{3 - 3\cos z} \leq 2\sqrt{1 - \cos z}.$$

Substituting $z := \arccos y$ shows (8). ■

Since \arccos is concave in $[0, 1]$, we can apply Jensen's inequality in addition to applying Lemma 3 and thus bound the expected error of the k -th hypothesis by

$$\begin{aligned} \mathbb{E}[\text{err}(\mathbf{w}_k)] &= \frac{1}{\pi} \mathbb{E}[\arccos(\cos \beta_k)] \\ &\leq \frac{1}{\pi} \arccos\left(\mathbb{E}[\cos \beta_k]\right) && \text{(Jensen's inequality)} \\ &\leq \frac{2}{\pi} \sqrt{1 - \mathbb{E}[\cos \beta_k]} && \text{(Inequality (8))} \\ &\leq \frac{2}{\pi} e^{-\frac{1}{3n}(k-1)} && \text{(Lemma 3)} \\ &\leq e^{-\frac{k}{3n}}. && (n \geq 2) \end{aligned}$$

Hence after

$$k \geq k_0 := 3n \ln \frac{2}{\epsilon \delta} \quad (9)$$

counterexamples have occurred, we have $\mathbb{E}[\text{err}(\mathbf{w}_k)] \leq \frac{\epsilon\delta}{2}$. Now suppose the output hypothesis \mathbf{w}_k has error greater than ϵ . Then due to monotonicity (Proposition 1, (b)) the error has always been greater than ϵ . So if we draw

$$s := \frac{2k_0}{\epsilon} = \Theta\left(\frac{n}{\epsilon} \left(\log \frac{1}{\epsilon} + \log \frac{1}{\delta}\right)\right) \quad (10)$$

examples, we expect to have at least $2k_0$ counterexamples in this case. We apply the Chernoff bound in the following version to bound the probability of getting less than k_0 counterexamples in this case.

Lemma 5 (Lower tail Chernoff bound) *If Y_1, \dots, Y_s are $\{0, 1\}$ -valued random variables with $\mathbb{P}[Y_i = 1 \mid Y_1, \dots, Y_{i-1}] \geq p$, then for all $0 < c < 1$,*

$$\mathbb{P}\left[\sum_i Y_i < (1 - c)sp\right] \leq e^{-c^2 sp/2}.$$

Proof The lemma follows from the standard Chernoff bound. ■

To bound the probability of getting less than k_0 counterexamples in s trials let the 0-1 variable Y_i indicate whether the i -th example is a counterexample. Set $p := \epsilon$, $c := 1/2$ and we obtain

$$\mathbb{P}\left[\sum_i Y_i < k_0\right] \leq (\epsilon\delta/2)^{3n/4} \leq \delta/2.$$

If the error of the output hypothesis \mathbf{w}_k is greater than ϵ , less than k_0 counterexamples were encountered or the random variable $\text{err}(\mathbf{w}_k)$ exceeds its expected value by at least a factor of $2/\delta$. Hence, by the union bound and Markov's inequality, the probability that the output hypothesis \mathbf{w}_k has error greater than ϵ is at most $\delta/2 + \delta/2 = \delta$, which proves Theorem 2. ■

4.2. Proof of Lemma 3 with Lemma 6

Proof [Lemma 3] Since \mathbf{w}^* and \mathbf{w}_k have norm one (Proposition 1, (a)), the cosine of β_k is given as

$$\cos \beta_k = \langle \mathbf{w}^*, \mathbf{w}_k \rangle. \quad (11)$$

Recalling the update rule of ADAPTIVEPERCEPTRON for $k \geq 2$, we have

$$\begin{aligned} \langle \mathbf{w}^*, \mathbf{w}_k \rangle &= \langle \mathbf{w}^*, \mathbf{w}_{k-1} - 2 \langle \mathbf{w}_{k-1}, \mathbf{x}_k \rangle \mathbf{x}_k \rangle \\ &= \langle \mathbf{w}^*, \mathbf{w}_{k-1} \rangle - 2 \langle \mathbf{w}_{k-1}, \mathbf{x}_k \rangle \langle \mathbf{w}^*, \mathbf{x}_k \rangle, \end{aligned} \quad (12)$$

where \mathbf{x}_k is the k -th counterexample. We combine Equations (11) and (12) and set $d_k^* := \pm \langle \mathbf{w}^*, \mathbf{x}_k \rangle$, $d_k := \mp \langle \mathbf{w}_{k-1}, \mathbf{x}_k \rangle$ for the distances from \mathbf{x}_k to the target hyperplane and the current hyperplane of the algorithm. This yields

$$\cos \beta_k = \cos \beta_{k-1} + 2d_k d_k^*. \quad (13)$$

By symmetry the probability distribution of $d_k d_k^*$ only depends on the hypothesis angle β_{k-1} . Thus we can form total expectation to obtain

$$\mathbb{E}[\cos \beta_k] = \mathbb{E}[\cos \beta_{k-1}] + 2 \mathbb{E}\left[\mathbb{E}[d_k d_k^* \mid \beta_{k-1}]\right]. \quad (14)$$

The following key lemma provides the conditional expectation.

Lemma 6 (expected product of distances) *Let \mathbf{w} be a hypothesis with $\angle(\mathbf{w}^*, \mathbf{w}) = \beta > 0$. Assume \mathbf{x} is a randomly drawn counterexample. Let d^* and d be the distances from \mathbf{x} to the target hyperplane and the hyperplane represented by \mathbf{w} . Then we have*

$$\mathbb{E}[dd^* \mid \beta] = \frac{1}{n} \left(\frac{\sin \beta}{\beta} - \cos \beta \right).$$

We show Lemma 6 later. In combination with Lemma 4 we see

$$\mathbb{E}[d_k d_k^* \mid \beta_{k-1}] \stackrel{\text{Lemma 6}}{=} \frac{1}{n} \left(\frac{\sin \beta_{k-1}}{\beta_{k-1}} - \cos \beta_{k-1} \right) \stackrel{\text{Inequality (7)}}{\geq} \frac{1}{3n} (1 - \cos \beta_{k-1}).$$

Thus, we may bound $\mathbb{E}[\cos \beta_k]$ in Equation (14) as follows

$$\begin{aligned} \mathbb{E}[\cos \beta_k] &\geq \mathbb{E}[\cos \beta_{k-1}] + 2 \mathbb{E} \left[\frac{1}{3n} (1 - \cos \beta_{k-1}) \right] \\ &\geq \frac{2}{3n} + \left(1 - \frac{2}{3n} \right) \mathbb{E}[\cos \beta_{k-1}]. \end{aligned}$$

Now expand this inequality recursively and notice that $\cos \beta_1 = \langle \mathbf{w}^*, \mathbf{w}_1 \rangle = \langle \mathbf{w}^*, b_1 \mathbf{x}_1 \rangle \geq 0$ for the first example (\mathbf{x}_1, b_1) . Hence we have for all $k \geq 1$,

$$\mathbb{E}[\cos \beta_k] \geq \frac{2}{3n} \sum_{i=0}^{k-2} \left(1 - \frac{2}{3n} \right)^i = 1 - (1 - 2/3n)^{k-1} \geq 1 - e^{-\frac{2}{3n}(k-1)}.$$

■

4.3. Proof of Lemma 6

Proof [Lemma 6] Let \mathbf{w} , β , \mathbf{x} and d , d^* be given as stated in the lemma. Without loss of generality (rotational symmetry) let $\mathbf{w}^* = (1, 0, \dots, 0)$ and $\mathbf{w} = (\cos \beta, -\sin \beta, 0, \dots, 0)$. Note that $\angle(\mathbf{w}^*, \mathbf{w}) = \beta$, $\|\mathbf{w}^*\| = \|\mathbf{w}\| = 1$ and $dd^* = -\langle \mathbf{w}, \mathbf{x} \rangle \langle \mathbf{w}^*, \mathbf{x} \rangle = x_1(x_2 \sin \beta - x_1 \cos \beta)$. Also by symmetry we may assume that \mathbf{x} is a positive counterexample. Now consider \mathbf{x} to be the angular part of a standard normal vector $\mathbf{u} = r\mathbf{x}$, where r is its length. Note that r and \mathbf{x} are independent and thus it holds

$$\mathbb{E}_{\mathbf{x}}[-\langle \mathbf{w}, \mathbf{x} \rangle \langle \mathbf{w}^*, \mathbf{x} \rangle] \mathbb{E}_r[r^2] = \mathbb{E}_{\mathbf{u}}[-\langle \mathbf{w}, \mathbf{u} \rangle \langle \mathbf{w}^*, \mathbf{u} \rangle].$$

Since r^2 has a chi-squared distribution with n degrees of freedom, its expected value is $\mathbb{E}_r[r^2] = n$. Hence it remains to show that $\mathbb{E}_{\mathbf{u}}[-\langle \mathbf{w}, \mathbf{u} \rangle \langle \mathbf{w}^*, \mathbf{u} \rangle] = \sin(\beta)/\beta - \cos \beta$. This can be done by

calculating a simple Gaussian integral:

$$\begin{aligned}
& \mathbb{E}_{\mathbf{u}} \left[- \langle \mathbf{w}, \mathbf{u} \rangle \langle \mathbf{w}^*, \mathbf{u} \rangle \right] \\
&= \frac{2\pi}{\beta} \int_{\substack{\mathbf{u} \in \mathbb{R}^n, u_1 \geq 0, \\ u_1 \cos \beta - u_2 \sin \beta < 0}} \frac{1}{\sqrt{(2\pi)^n}} e^{-\frac{1}{2}(u_1^2 + \dots + u_n^2)} u_1 (u_2 \sin \beta - u_1 \cos \beta) d(u_1 \dots u_n) \\
&= \frac{1}{\beta} \int_{\substack{u_1 \geq 0, \\ u_1 \cos \beta - u_2 \sin \beta < 0}} e^{-\frac{1}{2}(u_1^2 + u_2^2)} u_1 (u_2 \sin \beta - u_1 \cos \beta) d(u_1, u_2) \\
&= \frac{1}{\beta} \int_{r=0}^{\infty} \int_{\varphi=0}^{\beta} e^{-r^2/2} r \sin \varphi (r \cos \varphi \sin \beta - r \sin \varphi \cos \beta) r d\varphi dr \\
&= \frac{1}{\beta} \int_{r=0}^{\infty} r^3 e^{-r^2/2} dr \int_{\varphi=0}^{\beta} \sin \varphi \cos \varphi \sin \beta - \sin^2 \varphi \cos \beta d\varphi.
\end{aligned}$$

Now substituting the two integrals

$$\int_{r=0}^{\infty} r^3 e^{-r^2/2} dr = \left[-r^2 e^{-r^2/2} \right]_0^{\infty} + \int_0^{\infty} 2r e^{-r^2/2} dr = \left[-2e^{-r^2/2} \right]_0^{\infty} = 2$$

and

$$\begin{aligned}
& \int_{\varphi=0}^{\beta} \sin \varphi \cos \varphi \sin \beta - \sin^2 \varphi \cos \beta d\varphi \\
&= \left[-\frac{1}{2} \cos^2 \varphi \right]_0^{\beta} \sin \beta - \left[\frac{1}{2} (\varphi - \sin \varphi \cos \varphi) \right]_0^{\beta} \cos \beta = \frac{1}{2} (\sin \beta - \beta \cos \beta)
\end{aligned}$$

shows the claim. ■

5. Conclusions and Open Problems

The classical perceptron algorithm – with adaptive learning rate – turns out to be nearly optimal for learning homogeneous halfspaces against the uniform distribution on the unit sphere. The algorithm is fast, extremely simple, strictly error-decreasing and even conservative, i.e. it performs updates only on counterexamples.

Experiments suggest that OPTADAPTIVEPERCEPTRON performs only slightly better than ADAPTIVEPERCEPTRON. It would be interesting to investigate if there exist (conservative) learning algorithms which perform better than ADAPTIVEPERCEPTRON.

Moreover, it would be interesting to search for possible generalizations. For which distributions does Theorem 2 still hold? Is it possible to find a version of the algorithm which fits a given class of distributions?

Acknowledgments

I would like to thank Georg Schnitger and Martin Hoefer for a number of helpful discussions. Also thanks to all anonymous referees for their useful comments, especially for observing a simplification of the proof of Lemma 6 and for bringing [Dunagan and Vempala \(2004\)](#) to our attention.

References

- Maria-Florina Balcan and Philip M. Long. Active and passive learning of linear separators under log-concave distributions. In Shai Shalev-Shwartz and Ingo Steinwart, editors, *COLT 2013 - The 26th Annual Conference on Learning Theory, June 12-14, 2013, Princeton University, NJ, USA*, volume 30 of *JMLR Workshop and Conference Proceedings*, pages 288–316. JMLR.org, 2013. URL <http://proceedings.mlr.press/v30/Balcan13.html>.
- Eric B. Baum. The perceptron algorithm is fast for nonmalicious distributions. *Neural Computation*, 2(2):248–260, 1990. doi: 10.1162/neco.1990.2.2.248. URL <https://doi.org/10.1162/neco.1990.2.2.248>.
- Anselm Blumer, Andrzej Ehrenfeucht, David Haussler, and Manfred K. Warmuth. Learnability and the vapnik-chervonenkis dimension. *J. ACM*, 36(4):929–965, 1989. doi: 10.1145/76359.76371. URL <https://doi.org/10.1145/76359.76371>.
- John Dunagan and Santosh S. Vempala. A simple polynomial-time rescaling algorithm for solving linear programs. In László Babai, editor, *Proceedings of the 36th Annual ACM Symposium on Theory of Computing, Chicago, IL, USA, June 13-16, 2004*, pages 315–320. ACM, 2004. ISBN 1-58113-852-0. doi: 10.1145/1007352.1007404. URL <https://doi.org/10.1145/1007352.1007404>.
- Andrzej Ehrenfeucht, David Haussler, Michael J. Kearns, and Leslie G. Valiant. A general lower bound on the number of examples needed for learning. *Inf. Comput.*, 82(3):247–261, 1989. doi: 10.1016/0890-5401(89)90002-3. URL [https://doi.org/10.1016/0890-5401\(89\)90002-3](https://doi.org/10.1016/0890-5401(89)90002-3).
- David Haussler, Nick Littlestone, and Manfred K. Warmuth. Predicting $\{0,1\}$ -functions on randomly drawn points. *Inf. Comput.*, 115(2):248–292, 1994. doi: 10.1006/inco.1994.1097. URL <https://doi.org/10.1006/inco.1994.1097>.
- Philip M. Long. Halfspace learning, linear programming, and nonmalicious distributions. *Inf. Process. Lett.*, 51(5):245–250, 1994. doi: 10.1016/0020-0190(94)90003-5. URL [https://doi.org/10.1016/0020-0190\(94\)90003-5](https://doi.org/10.1016/0020-0190(94)90003-5).
- Philip M. Long. On the sample complexity of PAC learning half-spaces against the uniform distribution. *IEEE Trans. Neural Networks*, 6(6):1556–1559, 1995. doi: 10.1109/72.471352. URL <https://doi.org/10.1109/72.471352>.
- Philip M. Long. An upper bound on the sample complexity of pac-learning halfspaces with respect to the uniform distribution. *Inf. Process. Lett.*, 87(5):229–234, 2003. doi: 10.1016/S0020-0190(03)00311-9. URL [https://doi.org/10.1016/S0020-0190\(03\)00311-9](https://doi.org/10.1016/S0020-0190(03)00311-9).
- A. B. Novikoff. On convergence proofs on perceptrons. In *Proceedings of the Symposium on the Mathematical Theory of Automata*, volume 12, pages 615–622, New York, NY, USA, 1962. Polytechnic Institute of Brooklyn.
- F. Rosenblatt. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, pages 65–386, 1958.

Rocco A. Servedio. On PAC learning using winnow, perceptron, and a perceptron-like algorithm. In Shai Ben-David and Philip M. Long, editors, *Proceedings of the Twelfth Annual Conference on Computational Learning Theory, COLT 1999, Santa Cruz, CA, USA, July 7-9, 1999*, pages 296–307. ACM, 1999. ISBN 1-58113-167-4. doi: 10.1145/307400.307474. URL <https://doi.org/10.1145/307400.307474>.