

# Learning Halfspaces with Malicious Noise

**Adam R. Klivans**

*Computer Science Department, University of Texas at Austin*

KLIVANS@CS.UTEXAS.EDU

**Philip M. Long**

*Google*

PLONG@GOOGLE.COM

**Rocco A. Servedio**

*Computer Science Department, Columbia University*

ROCCO@CS.COLUMBIA.EDU

**Editor:** Manfred K. Warmuth

## Abstract

We give new algorithms for learning halfspaces in the challenging *malicious noise* model, where an adversary may corrupt both the labels and the underlying distribution of examples. Our algorithms can tolerate malicious noise rates exponentially larger than previous work in terms of the dependence on the dimension  $n$ , and succeed for the fairly broad class of all isotropic log-concave distributions.

We give poly( $n, 1/\epsilon$ )-time algorithms for solving the following problems to accuracy  $\epsilon$ :

- Learning origin-centered halfspaces in  $\mathbf{R}^n$  with respect to the uniform distribution on the unit ball with malicious noise rate  $\eta = \Omega(\epsilon^2 / \log(n/\epsilon))$ . (The best previous result was  $\Omega(\epsilon / (n \log(n/\epsilon))^{1/4})$ .)
- Learning origin-centered halfspaces with respect to any isotropic log-concave distribution on  $\mathbf{R}^n$  with malicious noise rate  $\eta = \Omega(\epsilon^3 / \log^2(n/\epsilon))$ . This is the first efficient algorithm for learning under isotropic log-concave distributions in the presence of malicious noise.

We also give a poly( $n, 1/\epsilon$ )-time algorithm for learning origin-centered halfspaces under any isotropic log-concave distribution on  $\mathbf{R}^n$  in the presence of *adversarial label noise* at rate  $\eta = \Omega(\epsilon^3 / \log(1/\epsilon))$ . In the adversarial label noise setting (or agnostic model), labels can be noisy, but not example points themselves. Previous results could handle  $\eta = \Omega(\epsilon)$  but had running time exponential in an unspecified function of  $1/\epsilon$ .

Our analysis crucially exploits both concentration and anti-concentration properties of isotropic log-concave distributions. Our algorithms combine an iterative outlier removal procedure using Principal Component Analysis together with “smooth” boosting.

**Keywords:** PAC learning, noise tolerance, malicious noise, agnostic learning, label noise, half-space learning, linear classifiers.

## 1. Introduction

A *halfspace* is a Boolean-valued function of the form  $f = \text{sign}(\sum_{i=1}^n w_i x_i - \theta)$ . Learning halfspaces in the presence of noisy data is a fundamental problem in machine learning. In addition

to its practical relevance, the problem has connections to many well-studied topics such as kernel methods (Shawe-Taylor and Cristianini, 2000), cryptographic hardness of learning (Klivans and Sherstov, 2006), hardness of approximation (Feldman et al., 2006; Guruswami and Raghavendra, 2006), learning Boolean circuits (Blum et al., 1997), and additive/multiplicative update learning algorithms (Littlestone, 1991; Freund and Schapire, 1999).

Learning an unknown halfspace from correctly labeled (non-noisy) examples is one of the best-understood problems in learning theory, with work dating back to the famous Perceptron algorithm of the 1950s (Rosenblatt, 1958) and a range of efficient algorithms known for different settings (Novikoff, 1962; Littlestone, 1987; Blumer et al., 1989; Maass and Turan, 1994). Much less is known, however, about the more difficult problem of learning halfspaces in the presence of noise.

Important progress was made by Blum *et al.* (Blum et al., 1997) who gave a polynomial-time algorithm for learning a halfspace under *classification noise*. In this model each label is flipped independently with some fixed probability; the noise does not affect the actual example points themselves, which are generated according to an arbitrary probability distribution over  $\mathbf{R}^n$ .

In the current paper we consider a much more challenging *malicious noise* model. In this model, introduced by Valiant (1985) (see also (Kearns and Li, 1993)), there is an unknown target function  $f$  and distribution  $\mathcal{D}$  over examples. Each time the learner receives an example, independently with probability  $1 - \eta$  it is drawn from  $\mathcal{D}$  and labeled correctly according to  $f$ , but with probability  $\eta$  it is an arbitrary pair  $(x, y)$  which may be generated by an omniscient adversary. The parameter  $\eta$  is known as the “noise rate.”

Malicious noise is a notoriously difficult model with few positive results. It was already shown by Kearns and Li (1993) that for essentially all concept classes, it is information-theoretically impossible to learn to accuracy  $1 - \epsilon$  if the noise rate  $\eta$  is greater than  $\epsilon/(1 + \epsilon)$ . Indeed, known algorithms for learning halfspaces (Servedio, 2003; Kalai et al., 2008) or even simpler target functions (Mansour and Parnas, 1998) with malicious noise typically make strong assumptions about the underlying distribution  $\mathcal{D}$ , and can learn to accuracy  $1 - \epsilon$  only for noise rates  $\eta$  much smaller than  $\epsilon$ . We describe the most closely related work that we know of in Section 1.2.

In this paper we consider learning under the uniform distribution on the unit ball in  $\mathbf{R}^n$ , and more generally under any isotropic log-concave distribution. The latter is a fairly broad class of distributions that includes spherical Gaussians and uniform distributions over a wide range of convex sets. Our algorithms can learn from malicious noise rates that are quite high, as we now describe.

## 1.1 Main Results

Our first result is an algorithm for learning halfspaces in the malicious noise model with respect to the uniform distribution on the  $n$ -dimensional unit ball:

**Theorem 1** *There is a  $\text{poly}(n, 1/\epsilon)$ -time algorithm that learns origin-centered halfspaces to accuracy  $1 - \epsilon$  with respect to the uniform distribution on the unit ball in  $n$  dimensions in the presence of malicious noise at rate  $\eta = \Omega(\epsilon^2 / \log(n/\epsilon))$ .*

The condition on  $\eta$  is expressed using  $\Omega$  and not  $O$  because we are showing that a weak upper bound on the noise rate suffices to achieve accuracy  $1 - \epsilon$ .

Via a more sophisticated algorithm, we can learn in the presence of malicious noise under any isotropic log-concave distribution:

**Theorem 2** *There is a  $\text{poly}(n, 1/\epsilon)$ -time algorithm that learns origin-centered halfspaces to accuracy  $1 - \epsilon$  with respect to any isotropic log-concave distribution over  $\mathbf{R}^n$  and can tolerate malicious noise at rate  $\eta = \Omega(\epsilon^3 / \log^2(n/\epsilon))$ .*

We are not aware of any previous polynomial-time algorithms for learning under isotropic log-concave distributions in the presence of malicious noise.

Finally, we also consider a related noise model known as *adversarial label noise*. In this model there is a fixed probability distribution  $P$  over  $\mathbf{R}^n \times \{-1, 1\}$  (i.e., over labeled examples) for which a  $1 - \eta$  fraction of draws are labeled according to an unknown halfspace. The marginal distribution over  $\mathbf{R}^n$  is assumed to be isotropic log-concave; so the idea is that an “adversary” chooses an  $\eta$  fraction of examples to mislabel, but unlike the malicious noise model she cannot change the (isotropic log-concave) distribution of the actual example points in  $\mathbf{R}^n$ . Learning with adversarial label noise is clearly harder than with independent misclassification noise – the ability to choose which labels to corrupt allows the adversary to coordinate their effects to an extent.

For the adversarial label noise model we prove:

**Theorem 3** *There is a  $\text{poly}(n, 1/\epsilon)$ -time algorithm that learns origin-centered halfspaces to accuracy  $1 - \epsilon$  with respect to any isotropic log-concave distribution over  $\mathbf{R}^n$  and can tolerate adversarial label noise at rate  $\eta = \Omega(\epsilon^3 / \log(1/\epsilon))$ .*

## 1.2 Previous Work

**Malicious noise.** General-purpose tools developed by Kearns and Li (1993) (see also (Kearns et al., 1994)) directly imply that halfspaces can be learned for any distribution over the domain in randomized  $\text{poly}(n, 1/\epsilon)$  time with malicious noise at a rate  $\Omega(\epsilon/n)$ ; the algorithm repeatedly picks a random subsample of the training data, hoping to miss all the noisy examples. Kannan (see Arora et al. (1993)) devised a deterministic algorithm with a  $\Omega(\epsilon/n)$  bound that repeatedly exploits Helly’s Theorem to find a group of  $n + 1$  examples that includes a noisy example, then removes the group. Kalai et al. (2008) showed that the  $\text{poly}(n, 1/\epsilon)$ -time averaging algorithm (Servedio, 2001) tolerates noise at a rate  $\Omega(\epsilon/\sqrt{n})$  when the distribution is uniform. They also described an improvement to  $\tilde{\Omega}(\epsilon/n^{1/4})$  based on the observation that uniform examples will tend to be well-separated, so that pairs of examples that are too close to one another can be removed.

**Adversarial label noise.** Kalai, et al. showed that if the distribution over the instances is uniform over the unit ball, the averaging algorithm tolerates adversarial label noise at a rate  $\Omega(\epsilon/\sqrt{\log(1/\epsilon)})$  in  $\text{poly}(n, 1/\epsilon)$  time. (In that paper, learning in the presence of adversarial label noise was called “agnostic learning”.) They also described an algorithm that fits low-degree polynomials that tolerates noise at a rate within an additive  $\epsilon$  of the accuracy, but in  $\text{poly}(n^{1/\epsilon^4})$  time; for log-concave distributions, their algorithm took  $\text{poly}(n^{d(1/\epsilon)})$  time, for an unspecified function  $d$ . The latter algorithm does not require that the distribution is isotropic, as ours does.

**Robust PCA.** Independently of this work, Xu et al. (2009) designed and analyzed an algorithm that performs principal component analysis when some of the examples are corrupted arbitrarily, as in the malicious noise model studied here. Also, the thesis of Brubaker (2009) presents a “Robust PCA” algorithm which is a PCA variant aimed at ameliorating the effects of noisy examples.

### 1.3 Techniques

**Outlier Removal.** Consider first the simplest problem of learning an origin-centered halfspace with respect to the uniform distribution on the  $n$ -dimensional ball. A natural idea is to use a simple “averaging” algorithm that takes the vector average of the positive examples it receives and uses this as the normal vector of its hypothesis halfspace. Servedio (2001) analyzed this algorithm for the random classification noise model, and Kalai et al. (2008) extended the analysis to the adversarial label noise model.

Intuitively the “averaging” algorithm can only tolerate low malicious noise rates because the adversary can generate noisy examples which “pull” the average vector far from its true location. Our main insight is that the adversary does this most effectively when the noisy examples are coordinated to pull in roughly the same direction. We use a form of outlier detection based on Principal Component Analysis to detect such coordination. This is done by computing the direction  $\mathbf{w}$  of maximal variance of the data set; if the variance in direction  $\mathbf{w}$  is suspiciously large, we remove from the sample all points  $\mathbf{x}$  for which  $(\mathbf{w} \cdot \mathbf{x})^2$  is large. Our analysis shows that this causes many noisy examples, and only a few non-noisy examples, to be removed.

We repeat this process until the variance in every direction is not too large. (This cannot take too many stages since many noisy examples are removed in each stage.) While some noisy examples may remain, we show that their scattered effects cannot hurt the algorithm much.

Thus, in a nutshell, our overall algorithm for the uniform distribution is to first do outlier removal<sup>1</sup> by an iterated PCA-type procedure, and then simply run the averaging algorithm on the remaining “cleaned-up” data set.

**Extending to Log-Concave Distributions via Smooth Boosting.** We are able to show that the iterative outlier removal procedure described above is useful for isotropic log-concave distributions as well as the uniform distribution: if examples are removed in a given stage, then many of the removed examples are noisy and only a few are non-noisy (the analysis here uses concentration bounds for isotropic log-concave distributions). However, even if there were no noise in the data, the average of the positive examples under an isotropic log-concave distribution need not give a high-accuracy hypothesis. Thus the averaging algorithm alone will not suffice after outlier removal.

To get around this, we show that after outlier removal the average of the positive examples gives a (real-valued) *weak* hypothesis that has some nontrivial predictive accuracy. (Interestingly, the proof of this relies heavily on *anti*-concentration properties of isotropic log-concave distributions!) A natural approach is then to use a boosting algorithm to convert this weak learner into a strong learner. This is not entirely straightforward because boosting “skews” the distribution of examples; this has the undesirable effects of both increasing the effective malicious noise rate, and causing the distribution to no longer be isotropic log-concave. However, by using a “smooth” boosting algorithm (Servedio, 2003) that skews the distribution as little as possible, we are able to control these undesirable effects and make the analysis go through. (The extra factor of  $\epsilon$  in the bound of Theorem 2 compared with Theorem 1 comes from the fact that the boosting algorithm constructs “ $1/\epsilon$ -skewed” distributions.)

We note that our approach of using smooth boosting is reminiscent of earlier work (Servedio, 2002, 2003), but the current algorithm goes well beyond that. Servedio (2002) did not consider a

---

1. We note briefly that the sophisticated outlier removal techniques of (Blum et al., 1997; Dunagan and Vempala, 2004) do not seem to be useful in our setting; those works deal with a strong notion of outliers, which is such that no point on the unit ball can be an outlier if a significant fraction of points are uniformly distributed on the unit ball.

noisy scenario, and Servedio (2003) only considered the averaging algorithm without any outlier removal as the weak learner (and thus could only handle quite low rates of malicious noise in our isotropic log-concave setting).

**Tolerating adversarial label noise.** Finally, our results for learning under isotropic log-concave distributions with adversarial label noise are obtained using a similar approach. The algorithm here is in fact simpler than the malicious noise algorithm: since the adversarial label noise model does not allow the adversary to alter the distribution of the examples in  $\mathbf{R}^n$ , we can dispense with the outlier removal and simply use smooth boosting with the averaging algorithm as the weak learner. (This is why we get a slightly better quantitative bound in Theorem 3 than Theorem 2).

**Organization.** For completeness we review the precise definitions of isotropic log-concave distributions and the various learning models in Section 2. We present the simpler and more easily understood uniform distribution analysis in Section 3. We extend the algorithm and analysis to isotropic log-concave distributions in Section 4. Learning with adversarial label noise is treated in Section 5. We conclude in Section 6.

## 2. Definitions and Preliminaries

### 2.1 Learning with Malicious Noise

Given a probability distribution  $\mathcal{D}$  over  $\mathbf{R}^n$ , and a target function  $f : \mathbf{R}^n \rightarrow \{-1, 1\}$ , we define the oracle  $EX_\eta(f, \mathcal{D})$  as follows:

- with probability  $1 - \eta$  the oracle draws  $\mathbf{x}$  according to  $\mathcal{D}$ , and outputs  $(\mathbf{x}, f(\mathbf{x}))$ , and
- with probability  $\eta$  the oracle outputs an arbitrary  $(\mathbf{x}, y)$  pair. This “noisy” example can be thought of as being generated adversarially and can depend on the state of the learning algorithm and previous draws from the oracle.

Given a data set drawn from  $EX_\eta(f, \mathcal{D})$ , we often refer to the examples  $(\mathbf{x}, f(\mathbf{x}))$  (that came from  $\mathcal{D}$ ) as “clean” examples and the remaining examples  $(\mathbf{x}, y)$  as “dirty” examples.

For a set  $\mathcal{S}$  of probability distributions and a set  $F$  of possible target functions, we say that a learning algorithm  $A$  learns  $F$  to accuracy  $1 - \epsilon$  with respect to  $\mathcal{S}$  in the presence of malicious noise at a rate  $\eta$  if the following holds: for any  $f \in F$ , and  $\mathcal{D} \in \mathcal{S}$ , given access to  $EX_\eta(f, \mathcal{D})$ , with probability at least  $1/2$ , the output hypothesis  $h$  generated by  $A$  satisfies  $\Pr_{\mathbf{x} \sim \mathcal{D}}[h(\mathbf{x}) \neq f(\mathbf{x})] \leq \epsilon$ . (The probability of success may be amplified arbitrarily close to 1 using standard techniques (Haussler et al., 1991).)

Since scaling  $\mathbf{x}$  by a positive constant does not affect its classification by a linear classifier, drawing examples uniformly from the unit ball is equivalent to drawing them uniformly from the surface  $\mathbb{S}^{n-1}$  of the unit sphere. When this is the distribution, we may also assume w.l.o.g. that even noisy examples  $(\mathbf{x}, y)$  have  $\mathbf{x} \in \mathbb{S}^{n-1}$  – this is simply because a learning algorithm can trivially identify and ignore any noisy example  $(\mathbf{x}, y)$  that has  $\|\mathbf{x}\| \neq 1$ .

### 2.2 Log-concave distributions

A probability distribution over  $\mathbf{R}^n$  is said to be *log-concave* if its density function is  $\exp(-\psi(\mathbf{x}))$  for a convex function  $\psi$ .

A probability distribution over  $\mathbf{R}^n$  is *isotropic* if the mean of the distribution is 0 and the covariance matrix is the identity, i.e.,  $\mathbf{E}[x_i x_j] = 1$  for  $i = j$  and 0 otherwise.

Isotropic log-concave (henceforth abbreviated i.l.c.) distributions are a fairly broad class of distributions. It is well known that any distribution induced by taking a uniform distribution over an arbitrary convex set and applying a suitable linear transformation to make it isotropic is then isotropic and log-concave. For an excellent treatment on basic properties of log-concave distributions, see Lovász and Vempala (2007).

We will use the following facts:

**Lemma 4 ((Lovász and Vempala, 2007))** *Let  $\mathcal{D}$  be an isotropic log-concave distribution over  $\mathbf{R}^n$  and  $\mathbf{a} \in \mathbb{S}^{n-1}$  any direction. Then for  $\mathbf{x}$  drawn according to  $\mathcal{D}$ , the distribution of  $\mathbf{a} \cdot \mathbf{x}$  is an isotropic log-concave distribution over  $\mathbf{R}$ .*

**Lemma 5 ((Lovász and Vempala, 2007))** *Any isotropic log-concave distribution  $\mathcal{D}$  over  $\mathbf{R}^n$  has light tails,*

$$\Pr_{\mathbf{x} \sim \mathcal{D}} [|\mathbf{x}| > \beta \sqrt{n}] \leq e^{-\beta+1}.$$

If  $n = 1$ , the density of  $\mathcal{D}$  is bounded:

$$\Pr_{x \sim \mathcal{D}} [x \in [a, b]] \leq |b - a|.$$

### 3. The uniform distribution and malicious noise

In this section we prove Theorem 1. As described above, our algorithm first does outlier removal using PCA and then applies the “averaging algorithm.”

We may assume throughout that the noise rate  $\eta$  is smaller than some absolute constant, and that the dimension  $n$  is larger than some absolute constant.

#### 3.1 The Algorithm: Removing Outliers and Averaging

Consider the following Algorithm  $A_{\text{mu}}$ :

**Algorithm  $A_{\text{mu}}$ :**

1. Draw a sample  $S$  of  $m = \text{poly}(n/\epsilon)$  many examples from the malicious oracle.
2. Identify the direction  $\mathbf{w} \in \mathbb{S}^{n-1}$  that maximizes

$$\sigma_{\mathbf{w}}^2 \stackrel{\text{def}}{=} \sum_{(\mathbf{x}, y) \in S} (\mathbf{w} \cdot \mathbf{x})^2.$$

If  $\sigma_{\mathbf{w}}^2 < \frac{10m \log m}{n}$  then go to Step 4 otherwise go to Step 3.

3. Remove from  $S$  every example that has  $(\mathbf{w} \cdot \mathbf{x})^2 \geq \frac{10 \log m}{n}$ . Go to Step 2.
4. For the examples  $S$  that remain let  $\mathbf{v} = \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} y \mathbf{x}$  and output the linear classifier  $h_{\mathbf{v}}$  defined by  $h_{\mathbf{v}}(\mathbf{x}) = \text{sgn}(\mathbf{v} \cdot \mathbf{x})$ .

We first observe that Step 2 can be carried out in polynomial time:

**Lemma 6** *There is a polynomial-time algorithm that, given a finite collection  $S$  of points in  $\mathbf{R}^n$ , outputs  $\mathbf{w} \in \mathbb{S}^{n-1}$  that maximizes  $\sum_{\mathbf{x} \in S} (\mathbf{w} \cdot \mathbf{x})^2$ .*

*Proof.* By applying Lagrange multipliers, we can see that the optimal  $\mathbf{w}$  is an eigenvector of  $A = \sum_{\mathbf{x} \in S} \mathbf{x}\mathbf{x}^T$ . Further, if  $\lambda$  is the eigenvalue of  $\mathbf{w}$ , then  $\sum_{\mathbf{x} \in S} (\mathbf{w} \cdot \mathbf{x})^2 = \mathbf{w}^T A \mathbf{w} = \mathbf{w}^T (\lambda \mathbf{w}) = \lambda$ . The eigenvector  $\mathbf{w}$  with the largest eigenvalue can be found in polynomial time (see, e.g., (Jolliffe, 2002)). ■

Before embarking on the analysis we establish a terminological convention. Much of our analysis deals with high-probability statements over the draw of the  $m$ -element sample  $S$ ; it is straightforward but quite cumbersome to explicitly keep track of all of the failure probabilities. Thus we write “with high probability” (or “w.h.p.”) in various places below as a shorthand for “with probability at least  $1 - 1/\text{poly}(n/\epsilon)$ .” The interested reader can easily verify that an appropriate  $\text{poly}(n/\epsilon)$  choice of  $m$  makes all the failure probabilities small enough so that the entire algorithm succeeds with probability at least  $1/2$  as required.

### 3.2 Properties of the clean examples

In this subsection we establish properties of the clean examples that were sampled in Step 1 of  $A_{\text{mu}}$ . The first says that no direction has much more variance than the expected variance of  $1/n$ :

**Lemma 7** *W.h.p. over a random draw of  $\ell$  clean examples  $S_{\text{clean}}$ , we have*

$$\max_{\mathbf{a} \in \mathbb{S}^{n-1}} \left\{ \frac{1}{\ell} \sum_{(\mathbf{x}, y) \in S_{\text{clean}}} (\mathbf{a} \cdot \mathbf{x})^2 \right\} \leq \frac{1}{n} + \sqrt{\frac{O(n + \log \ell)}{\ell}}.$$

*Proof.* The proof uses standard tools from VC theory and is in Appendix A. ■

The next lemma says that in fact no direction has too many clean examples lying far out in that direction:

**Lemma 8** *For any  $\beta > 0$  and  $\kappa > 1$ , if  $S_{\text{clean}}$  is a random set of  $\ell \geq \frac{O(1) \cdot n^2 \beta^2 e^{\beta^2 n/2}}{(1+\kappa) \ln(1+\kappa)}$  clean examples then w.h.p. we have*

$$\max_{\mathbf{a} \in \mathbb{S}^{n-1}} \frac{1}{\ell} |\{\mathbf{x} \in S_{\text{clean}} : (\mathbf{a} \cdot \mathbf{x})^2 > \beta^2\}| \leq (1 + \kappa) e^{-\beta^2 n/2}.$$

*Proof.* In Appendix B. ■

### 3.3 What is removed

In this section, we provide bounds on the number of clean and dirty examples removed in Step 3.

The first bound is a Corollary of Lemma 8.

**Corollary 9** *W.h.p. over the random draw of the  $m$ -element sample  $S$ , the number of clean examples removed during any one execution of Step 3 in  $A_{\text{mu}}$  is at most  $6n \log m$ .*

*Proof.* Since the noise rate  $\eta$  is sufficiently small, w.h.p. the number  $\ell$  of clean examples is at least (say)  $m/2$ . We would like to apply Lemma 8 with  $\kappa = 5\ell^4 n \log \ell$  and  $\beta = \sqrt{\frac{10 \log m}{n}}$ , and indeed we may do this because we have

$$\frac{O(1) \cdot n^2 \beta^2 e^{\beta^2 n/2}}{(1 + \kappa) \ln(1 + \kappa)} \leq \frac{O(1) \cdot n (\log m) m^5}{(1 + \kappa) \ln(1 + \kappa)} \leq O\left(\frac{m}{\log m}\right) \leq \frac{m}{2} \leq \ell$$

for  $n$  sufficiently large. Since clean points are only removed if they have  $(\mathbf{a} \cdot \mathbf{x})^2 > \beta^2$ , Lemma 8 gives us that the number of clean points removed is at most

$$m(1 + \kappa)e^{-\beta^2 n/2} \leq 6m^5 n \log(\ell)/m^5 \leq 6n \log m.$$

■

The counterpart to Corollary 9 is the following lemma. It tells us that if examples are removed in Step 3, then there must be many *dirty* examples removed. It exploits the fact that Lemma 7 bounds the variance in *all* directions  $\mathbf{a}$ , so that it can be reused to reason about what happens in different executions of step 3.

**Lemma 10** *W.h.p. over the random draw of  $S$ , whenever  $A_{\text{mu}}$  executes step 3, it removes at least  $\frac{4m \log m}{n}$  noisy examples from  $S_{\text{dirty}}$ , the set of dirty examples in  $S$ .*

*Proof.* As stated earlier we may assume that  $\eta \leq 1/4$ . This implies that w.h.p. the fraction  $\hat{\eta}$  of noisy examples in the initial set  $S$  is at most  $1/2$ . Finally, Lemma 7 implies that  $m = \tilde{\Omega}(n^3)$  suffices for it to be the case that w.h.p., for all  $\mathbf{a} \in \mathbb{S}^{n-1}$ , for the original multiset  $S_{\text{clean}}$  of clean examples drawn in step 1, we have

$$\sum_{(\mathbf{x}, y) \in S_{\text{clean}}} (\mathbf{a} \cdot \mathbf{x})^2 \leq \frac{2m}{n}. \quad (1)$$

We shall say that a random sample  $S$  that satisfies all these requirements is “reasonable”. We will show that for any reasonable dataset, the number of noisy examples removed during the execution of step 3 of  $A_{\text{mu}}$  is at least  $\frac{4m \log m}{n}$ .

If we remove examples using direction  $\mathbf{w}$  then it means  $\sum_{(\mathbf{x}, y) \in S} (\mathbf{w} \cdot \mathbf{x})^2 \geq \frac{10m \log m}{n}$ . Since  $S$  is reasonable, by (1) the contribution to the sum from the clean examples that survived to the current stage is at most  $2m/n$  so we must have

$$\sum_{(\mathbf{x}, y) \in S_{\text{dirty}}} (\mathbf{w} \cdot \mathbf{x})^2 \geq 10m \log(m)/n - 2m/n > 9m \log(m)/n.$$

Let us decompose  $S_{\text{dirty}}$  into  $N \cup F$  where  $N$  (“near”) consists of those points  $x$  s.t.  $(\mathbf{w} \cdot \mathbf{x})^2 \leq 10 \log(m)/n$  and  $F$  (“far”) is the remaining points for which  $(\mathbf{w} \cdot \mathbf{x})^2 > 10 \log(m)/n$ . Since  $|N| \leq |S_{\text{dirty}}| \leq \hat{\eta}m$ , (any dirty examples removed in earlier rounds will only reduce the size of  $S_{\text{dirty}}$ ) we have

$$\sum_{(\mathbf{x}, y) \in N} (\mathbf{w} \cdot \mathbf{x})^2 \leq (\hat{\eta}m)10 \log(m)/n$$

and so

$$|F| \geq \sum_{(\mathbf{x}, y) \in F} (\mathbf{w} \cdot \mathbf{x})^2 \geq 9m \log(m)/n - (\hat{\eta}m)10 \log(m)/n \geq 4m \log(m)/n$$

(the last line used the fact that  $\hat{\eta} < 1/2$ ). Since the points in  $F$  are removed in Step 3, the lemma is proved. ■



### 3.4 Exploiting limited variance in any direction

In this section, we show that if all directional variances are small, then the algorithm’s final hypothesis will have high accuracy.

We first recall a simple lemma which shows that a sample of “clean” examples results in a high-accuracy hypothesis for the averaging algorithm:

**Lemma 11 ((Servedio, 2001))** *Suppose  $\mathbf{x}_1, \dots, \mathbf{x}_m$  are chosen uniformly at random from  $\mathbb{S}^{n-1}$ , and a target weight vector  $\mathbf{u} \in \mathbb{S}^{n-1}$  produces labels  $y_1 = \text{sign}(\mathbf{u} \cdot \mathbf{x}_1), \dots, y_m = \text{sign}(\mathbf{u} \cdot \mathbf{x}_m)$ . Let  $\mathbf{v} = \frac{1}{m} \sum_{t=1}^m y_t \mathbf{x}_t$ . Then w.h.p.  $\mathbf{u} \cdot \mathbf{v} = \Omega(\frac{1}{\sqrt{n}})$ , while  $\|\mathbf{v} - (\mathbf{u} \cdot \mathbf{v})\mathbf{u}\| = O(\sqrt{\log(n)/m})$ .*

Now we can state Lemma 12.

**Lemma 12** *Let  $S = S_{\text{clean}} \cup S_{\text{dirty}}$  be the sample of  $m$  examples drawn from the noisy oracle  $\text{EX}_\eta(f, \mathcal{U})$ . Let*

- $S'_{\text{clean}}$  be those clean examples that were never removed during step 3 of  $A_{\text{mu}}$ ,
- $S'_{\text{dirty}}$  be those dirty examples that were never removed during step 3 of  $A_{\text{mu}}$ ,
- $\eta' = \frac{|S'_{\text{dirty}}|}{|S'_{\text{clean}} \cup S'_{\text{dirty}}|}$ , i.e., the fraction of dirty examples among the examples that survive step 3, and
- $\alpha = \frac{|S_{\text{clean}} - S'_{\text{clean}}|}{|S'_{\text{clean}} \cup S'_{\text{dirty}}|}$ , the ratio of the number of clean points that were erroneously removed to the size of the final surviving data set.

Let  $S' \stackrel{\text{def}}{=} S'_{\text{clean}} \cup S'_{\text{dirty}}$ . Suppose that  $|S'| \geq m/2$  (i.e., fewer than half the total points were removed) and that, for every direction  $\mathbf{w} \in \mathbb{S}^{n-1}$  we have

$$\sigma_{\mathbf{w}}^2 \stackrel{\text{def}}{=} \sum_{(\mathbf{x}, y) \in S'} (\mathbf{w} \cdot \mathbf{x})^2 \leq \frac{10m \log m}{n}.$$

Then w.h.p. over the draw of  $S$ , the halfspace with normal vector  $\mathbf{v} \stackrel{\text{def}}{=} \frac{1}{|S'|} \sum_{(\mathbf{x}, y) \in S'} y \mathbf{x}$  has error rate

$$O\left(\sqrt{\eta' \log m} + \alpha \sqrt{n} + \sqrt{\frac{n \log n}{m}}\right).$$

*Proof.* The claimed bound is trivial unless  $\eta' \leq o(1)/\log m$  and  $\alpha \leq o(1)/\sqrt{n}$ , so we shall freely use these bounds in what follows.

Let  $\mathbf{u}$  be the unit length normal vector for the target halfspace. Let  $\mathbf{v}_{\text{clean}}$  be the average of all the clean examples,  $\mathbf{v}'_{\text{dirty}}$  be the average of the dirty (noisy) examples that were not deleted (i.e., the examples in  $S'_{\text{dirty}}$ ), and  $\mathbf{v}_{\text{del}}$  be the average of the clean examples that were deleted. Then

$$\begin{aligned} \mathbf{v} &= \frac{1}{|S'_{\text{clean}} \cup S'_{\text{dirty}}|} \sum_{(\mathbf{x}, y) \in S'_{\text{clean}} \cup S'_{\text{dirty}}} y \mathbf{x} \\ &= \frac{1}{|S'_{\text{clean}} \cup S'_{\text{dirty}}|} \left( \left( \sum_{(\mathbf{x}, y) \in S_{\text{clean}}} y \mathbf{x} \right) + \left( \sum_{(\mathbf{x}, y) \in S'_{\text{dirty}}} y \mathbf{x} \right) - \left( \sum_{(\mathbf{x}, y) \in S_{\text{clean}} - S'_{\text{clean}}} y \mathbf{x} \right) \right) \\ \mathbf{v} &= (1 - \eta' + \alpha) \mathbf{v}_{\text{clean}} + \eta' \mathbf{v}'_{\text{dirty}} - \alpha \mathbf{v}_{\text{del}}. \end{aligned} \tag{2}$$

Let us begin by exploiting the bound on the variance in every direction to bound the length of  $\mathbf{v}'_{\text{dirty}}$ . For any  $\mathbf{w} \in \mathbb{S}^{n-1}$  we know that

$$\sum_{(\mathbf{x}, y) \in S'} (\mathbf{w} \cdot \mathbf{x})^2 \leq \frac{10m \log m}{n}, \quad \text{and hence} \quad \sum_{(\mathbf{x}, y) \in S'_{\text{dirty}}} (\mathbf{w} \cdot \mathbf{x})^2 \leq \frac{10m \log m}{n}$$

since  $S'_{\text{dirty}} \subseteq S'$ . Since  $|S'_{\text{dirty}}| \leq \eta' m$ , the fact that  $\|\mathbf{r}\|_1 \leq \sqrt{k} \|\mathbf{r}\|_2$  for any vector  $\mathbf{r} \in \mathbf{R}^k$  gives

$$\sum_{(\mathbf{x}, y) \in S'_{\text{dirty}}} |\mathbf{w} \cdot \mathbf{x}| \leq \sqrt{\frac{10m |S'_{\text{dirty}}| \log m}{n}}.$$

Taking  $\mathbf{w}$  to be the unit vector in the direction of  $\mathbf{v}'_{\text{dirty}}$ , we have  $\|\mathbf{v}'_{\text{dirty}}\| =$

$$\mathbf{w} \cdot \mathbf{v}'_{\text{dirty}} = \mathbf{w} \cdot \frac{1}{|S'_{\text{dirty}}|} \sum_{(\mathbf{x}, y) \in S'_{\text{dirty}}} y \mathbf{x} \leq \frac{1}{|S'_{\text{dirty}}|} \sum_{(\mathbf{x}, y) \in S'_{\text{dirty}}} |\mathbf{w} \cdot \mathbf{x}| \leq \sqrt{\frac{10m \log m}{|S'_{\text{dirty}}| n}}. \quad (3)$$

Because the domain distribution is uniform, the error of  $h_{\mathbf{v}}$  is proportional to the angle between  $\mathbf{v}$  and  $\mathbf{u}$ , in particular,

$$\Pr[h_{\mathbf{v}} \neq f] = \frac{1}{\pi} \arctan \left( \frac{\|\mathbf{v} - (\mathbf{v} \cdot \mathbf{u})\mathbf{u}\|}{\mathbf{u} \cdot \mathbf{v}} \right) \leq (1/\pi) \frac{\|\mathbf{v} - (\mathbf{v} \cdot \mathbf{u})\mathbf{u}\|}{\mathbf{u} \cdot \mathbf{v}}. \quad (4)$$

We have that  $\|\mathbf{v} - (\mathbf{v} \cdot \mathbf{u})\mathbf{u}\|$  equals

$$\begin{aligned} & \|(1 - \eta' + \alpha)(\mathbf{v}_{\text{clean}} - (\mathbf{v}_{\text{clean}} \cdot \mathbf{u})\mathbf{u}) + \eta'(\mathbf{v}'_{\text{dirty}} - (\mathbf{v}'_{\text{dirty}} \cdot \mathbf{u})\mathbf{u}) - \alpha(\mathbf{v}_{\text{del}} - (\mathbf{v}_{\text{del}} \cdot \mathbf{u})\mathbf{u})\| \\ & \leq 2\|\mathbf{v}_{\text{clean}} - (\mathbf{v}_{\text{clean}} \cdot \mathbf{u})\mathbf{u}\| + \eta'\|\mathbf{v}'_{\text{dirty}}\| + \alpha\|\mathbf{v}_{\text{del}}\| \end{aligned}$$

where we have used the triangle inequality and the fact that  $\alpha, \eta'$  are “small.” Lemma 11 lets us bound the first term in the sum by  $O(\sqrt{\log(n)/m})$ , and the fact that  $\mathbf{v}_{\text{del}}$  is an average of vectors of length 1 lets us bound the third by  $\alpha$ . For the second term, Equation (3) gives us

$$\eta'\|\mathbf{v}'_{\text{dirty}}\| \leq \sqrt{\frac{10m(\eta')^2 \log m}{|S'_{\text{dirty}}| n}} = \sqrt{\frac{10m\eta' \log m}{|S'| n}} \leq \sqrt{\frac{20\eta' \log m}{n}},$$

where for the last equality we used  $|S'| \geq m/2$ . We thus get

$$\|\mathbf{v} - (\mathbf{v} \cdot \mathbf{u})\mathbf{u}\| \leq O\left(\sqrt{\log(n)/m}\right) + \sqrt{20\eta' \log(m)/n} + \alpha. \quad (5)$$

Now we consider the denominator of (4). We have

$$\mathbf{u} \cdot \mathbf{v} = (1 - \eta' + \alpha)(\mathbf{u} \cdot \mathbf{v}_{\text{clean}}) + \eta'\mathbf{u} \cdot \mathbf{v}'_{\text{dirty}} - \alpha\mathbf{u} \cdot \mathbf{v}_{\text{del}}.$$

Similar to the above analysis, we again use Lemma 11 (but now the lower bound  $\mathbf{u} \cdot \mathbf{v} \geq \Omega(1/\sqrt{n})$ ), Equation (3), and the fact that  $\|\mathbf{v}_{\text{del}}\| \leq 1$ . Since  $\alpha$  and  $\eta'$  are “small,” we get that there is an absolute constant  $c$  such that  $\mathbf{u} \cdot \mathbf{v} \geq c/\sqrt{n} - \sqrt{20\eta' \log(m)/n} - \alpha$ . Combining this with (5) and (4), we get

$$\Pr[h_{\mathbf{v}} \neq f] \leq \frac{O\left(\sqrt{\frac{\log n}{m}}\right) + \sqrt{\frac{20\eta' \log m}{n}} + \alpha}{\pi \left(\frac{c}{\sqrt{n}} - \sqrt{\frac{20\eta' \log m}{n}} - \alpha\right)} = O\left(\sqrt{\frac{n \log n}{m}} + \sqrt{\eta' \log m} + \alpha\sqrt{n}\right).$$

■

### 3.5 Proof of Theorem 1

By Corollary 9, w.h.p. each outlier removal stage removes at most  $6n \log m$  clean points.

Since, by Lemma 10, each outlier removal stage removes at least  $\frac{4m \log m}{n}$  noisy examples, there must be at most  $O(n/(\log m))$  such stages. Consequently the total number of clean examples removed across all stages is  $O(n^2)$ . Since w.h.p. the initial number of clean examples is at least  $3m/4$ , this means that the final data set (on which the averaging algorithm is run) contains at least  $3m/4 - O(n^2)$  clean examples, and hence at least  $3m/4 - O(n^2)$  examples in total. The condition  $m \gg n^2$  means that the number of surviving examples will be at least  $m/2$ . Consequently the value of  $\alpha$  from Lemma 12 after the final outlier removal stage (the ratio of the total number of clean examples deleted, to the total number of surviving examples) is at most  $\frac{O(n^2)}{m}$ .

The standard Hoeffding bound implies that w.h.p. the actual fraction of noisy examples in the original sample  $S$  is at most  $\eta + \sqrt{O(\log m)/m}$ . It is easy to see that w.h.p. the fraction of dirty examples does not increase (since each stage of outlier removal removes more dirty points than clean points, for a suitably large  $\text{poly}(n/\epsilon)$  value of  $m$ ), and thus the fraction  $\eta'$  of dirty examples among the remaining examples after the final outlier removal stage is at most  $\eta + \sqrt{O(\log m)/m}$ .

Applying Lemma 12, for a suitably large value  $m = \text{poly}(n/\epsilon)$ , we obtain  $\Pr[h_{\forall} \neq f] \leq O(\sqrt{\eta \log m})$ . Rearranging this bound, we can learn to accuracy  $\epsilon$  even for  $\eta = \Omega(\epsilon^2/\log(n/\epsilon))$ . This completes the proof of the theorem.  $\blacksquare$

## 4. Isotropic log-concave distributions and malicious noise

Our algorithm  $A_{\text{mlc}}$  that works for arbitrary isotropic log-concave distributions uses smooth boosting.

### 4.1 Smooth Boosting

A boosting algorithm uses a subroutine, called a *weak learner*, that is only guaranteed to output hypotheses with a non-negligible advantage over random guessing.<sup>2</sup> The boosting algorithm that we consider uses a *confidence-rated* weak learner (Schapire and Singer, 1999), which predicts  $\{-1, 1\}$  labels using continuous values in  $[-1, 1]$ . Formally, the *advantage* of a hypothesis  $h'$  with respect to a distribution  $\mathcal{D}'$  is defined to be  $\mathbf{E}_{x \sim \mathcal{D}'}[h'(x)f(x)]$ , where  $f$  is the target function.

For the purposes of this paper, a boosting algorithm makes use of the weak learner, an example oracle (possibly corrupted with noise), a desired accuracy  $\epsilon$ , and a bound  $\gamma$  on the advantage of the hypothesis output by the weak learner.

A boosting algorithm that is trying to learn an unknown target function  $f$  with respect to some distribution  $\mathcal{D}$  repeatedly simulates a (possibly noisy) example oracle for  $f$  with respect to some other distribution  $\mathcal{D}'$  and calls a subroutine  $A_{\text{weak}}$  with respect to this oracle, receiving a *weak hypothesis*, which maps  $\mathbf{R}^n$  to the continuous interval  $[-1, 1]$ .

After repeating this for some number of stages, the boosting algorithm combines the weak hypotheses generated during its various calls to the weak learner into a final aggregate hypothesis which it outputs.

Let  $\mathcal{D}, \mathcal{D}'$  be two distributions over  $\mathbf{R}^n$ . We say that  $\mathcal{D}'$  is  $(1/\epsilon)$ -smooth with respect to  $\mathcal{D}$  if  $\mathcal{D}'(E) \leq (1/\epsilon)\mathcal{D}(E)$  for all events  $E$ .

---

2. For simplicity of presentation we ignore the confidence parameter of the weak learner in our discussion; this can be handled in an entirely standard way.

The following lemma from (Servedio, 2003) (similar results can be readily found elsewhere, see, e.g., (Gavinsky, 2003)) identifies the properties that we need from a boosting algorithm for our analysis.

**Lemma 13 ((Servedio, 2003))** *There is a boosting algorithm  $B$  and a polynomial  $p$  such that, for any  $\epsilon, \gamma > 0$ , the following properties hold. When learning a target function  $f$  using  $\text{EX}_\eta(f, \mathcal{D})$ , we have: (a) If each call to  $A_{\text{weak}}$  takes time  $t$ , then  $B$  takes time  $p(t, 1/\gamma, 1/\epsilon)$ . (b) The weak learner is always called with an oracle  $\text{EX}_{\eta'}(f, \mathcal{D}')$  where  $\mathcal{D}'$  is  $(1/\epsilon)$ -smooth with respect to  $\mathcal{D}$  and  $\eta' \leq \eta/\epsilon$ . (c) Suppose that for each distribution  $\text{EX}_{\eta'}(f, \mathcal{D}')$  passed to  $A_{\text{weak}}$  by  $B$ , the output of  $A_{\text{weak}}$  has advantage  $\gamma$ . Then the final output  $h$  of  $B$  satisfies  $\Pr_{x \in \mathcal{D}}[h(x) \neq f(x)] \leq \epsilon$ .*

## 4.2 The Algorithm

Our algorithm for learning under isotropic log-concave distributions with malicious noise, Algorithm  $A_{\text{mlc}}$ , applies the smooth booster from Lemma 13 with the following weak learner, which we call Algorithm  $A_{\text{mlcw}}$ . (The value  $c_0$  is an absolute constant that will emerge from our analysis.)

### Algorithm $A_{\text{mlcw}}$ :

1. Draw  $m = \text{poly}(n/\epsilon)$  examples from the oracle  $\text{EX}_{\eta'}(f, \mathcal{D}')$ .
2. Remove all those examples  $(\mathbf{x}, y)$  for which  $\|\mathbf{x}\| > \sqrt{3n \log m}$ .
3. Repeatedly
  - find a direction (unit vector)  $\mathbf{w}$  that maximizes  $\sum_{(\mathbf{x}, y) \in S} (\mathbf{w} \cdot \mathbf{x})^2$  (see Lemma 6)
  - if  $\sum_{(\mathbf{x}, y) \in S} (\mathbf{w} \cdot \mathbf{x})^2 \leq c_0^2 m \log^2(n/\epsilon)$  then move on to Step 4, and otherwise
  - remove from  $S$  all examples  $(\mathbf{x}, y)$  for which  $|\mathbf{w} \cdot \mathbf{x}| > c_0 \log(n/\epsilon)$ , and iterate again.
4. Let  $\mathbf{v} = \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} y \mathbf{x}$ , and return  $h$  defined by  $h(\mathbf{x}) = \frac{\mathbf{v} \cdot \mathbf{x}}{3n \log m}$ , if  $|\mathbf{v} \cdot \mathbf{x}| \leq 3n \log m$ , and  $h(\mathbf{x}) = \text{sgn}(\mathbf{v} \cdot \mathbf{x})$  otherwise.

## 4.3 The key claim: the weak learner is effective

Our main task is to analyze the weak learner. Given the following Lemma, Theorem 2 will be an immediate consequence of Lemma 13.

**Lemma 14** *Suppose Algorithm  $A_{\text{mlcw}}$  is run using  $\text{EX}_{\eta'}(f, \mathcal{D}')$  where  $f$  is an origin-centered half-space,  $\mathcal{D}'$  is  $(1/\epsilon)$ -smooth w.r.t. an isotropic log-concave distribution  $\mathcal{D}$ ,  $\eta' \leq \eta/\epsilon$ , and  $\eta \leq \Omega(\epsilon^3 / \log^2(n/\epsilon))$ . Then w.h.p. the hypothesis  $h$  returned by  $A_{\text{mlcw}}$  has advantage  $\Omega\left(\frac{\epsilon^2}{n \log(n/\epsilon)}\right)$ .*

Before proving Lemma 14, we need to prove some uniformity results on non-noisy examples drawn from an isotropic, log-concave distribution. This will enable us to use outlier removal and averaging to find a weak learner.

## 4.4 Lemmas in support of Lemma 14

In this section, let us consider a single call to the weak learner with an oracle  $\text{EX}_{\eta'}(f, \mathcal{D}')$  where  $\mathcal{D}'$  is  $(1/\epsilon)$ -smooth with respect to an isotropic log-concave distribution  $\mathcal{D}$  and  $\eta' \leq \eta/\epsilon$ . Our analysis will follow the same basic steps as Section 3.

A preliminary observation is that w.h.p. all clean examples drawn in Step 1 of Algorithm  $A_{\text{mlcw}}$  have  $\|\mathbf{x}\| \leq \sqrt{3n \log m}$ ; indeed, for any given draw of  $\mathbf{x}$  from  $\mathcal{D}'$ , the probability that  $\|\mathbf{x}\| > \sqrt{3n \log m}$  is at most  $\frac{\epsilon}{em^3}$  by Lemma 5 together with the fact that  $\mathcal{D}'$  is  $1/\epsilon$ -smooth with respect to an i.l.c. distribution. Therefore, w.h.p., only noisy examples are removed in Step 2 of the algorithm, and we shall assume that the distributions  $\mathcal{D}$  and  $\mathcal{D}'$  are in fact supported entirely on  $\{\mathbf{x} : \|\mathbf{x}\| \leq \sqrt{3n \log m}\}$ . This assumption affects us in two ways: first, it costs us an additional  $\frac{\epsilon}{em^2}$  in the failure probability analysis below (which is not a problem and is in fact swallowed up by our ‘‘w.h.p.’’ notation). Second, it means that the overall  $1 - \epsilon$  accuracy bound we establish for the entire learning algorithm may be slightly worse than the true value. This is because our final hypothesis may always be wrong on the examples  $\mathbf{x}$  that have  $\|\mathbf{x}\| > \sqrt{3n \log m}$  and are ignored in our analysis; however such examples have probability mass at most  $\frac{\epsilon}{m^3}$  under the isotropic log-concave distribution  $\mathcal{D}$  (again by Lemma 5), and thus the additional accuracy cost is at most  $\frac{\epsilon}{m^3}$ . Since  $\epsilon \gg \frac{\epsilon}{m^3}$ , this does not affect the overall correctness of our analysis. Note that a consequence of this assumption is that we can just take  $h(\mathbf{x}) = \frac{\mathbf{v} \cdot \mathbf{x}}{3n \log m}$ .

The remarks about high-probability statements and failure probabilities from Section 3.1 apply here as well, and as in Section 3 we write ‘‘w.h.p.’’ as shorthand for ‘‘with probability  $1 - 1/\text{poly}(n/\epsilon)$ .’’

We first show that the variance of  $\mathcal{D}'$  in every direction is not too large:

**Lemma 15** *For any  $\mathbf{a} \in \mathbb{S}^{n-1}$  we have  $E_{\mathbf{x} \sim \mathcal{D}'}[(\mathbf{a} \cdot \mathbf{x})^2] = O(\log^2(1/\epsilon))$ .*

*Proof.* For  $\mathbf{x}$  chosen according to  $\mathcal{D}$ , the distribution of  $\mathbf{a} \cdot \mathbf{x}$  is a unit variance log-concave distribution by Lemma 4. Thus, for any positive integer  $k$ ,

$$\begin{aligned} E_{\mathbf{x} \sim \mathcal{D}'}[(\mathbf{a} \cdot \mathbf{x})^2] &\leq k^2 + \sum_{i=k}^{\infty} (i+1)^2 \Pr_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \in (i, i+1]] \\ &\leq k^2 + \sum_{i=k}^{\infty} (i+1)^2 (1/\epsilon) \Pr_{\mathbf{x} \sim \mathcal{D}}[|\mathbf{a} \cdot \mathbf{x}| \in (i, i+1]] \\ &\leq k^2 + (1/\epsilon) \sum_{i=k}^{\infty} (i+1)^2 \Pr_{\mathbf{x} \sim \mathcal{D}}[|\mathbf{a} \cdot \mathbf{x}| > i] \\ &\leq k^2 + (1/\epsilon) \sum_{i=k}^{\infty} (i+1)^2 e^{-i+1} \leq k^2 + (1/\epsilon) \cdot \Theta(k^2 e^{-k}) \end{aligned}$$

where the first inequality in the last line uses Lemmas 4 and 5.

Setting  $k = \ln(1/\epsilon)$  completes the proof. ■

The following anticoncentration bound will be useful for proving that clean examples drawn from  $\mathcal{D}'$  tend to be classified correctly with a large margin.

**Lemma 16** *Let  $\mathbf{u} \in \mathbb{S}^{n-1}$ . Then*

$$E_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{u} \cdot \mathbf{x}|] \geq \epsilon/8.$$

*Proof.* Clearly

$$E_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{u} \cdot \mathbf{x}|] \geq (\epsilon/4) \Pr_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{u} \cdot \mathbf{x}| > \epsilon/4].$$

But by Lemma 5,

$$\Pr_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{u} \cdot \mathbf{x}| \leq \epsilon/4] \leq \frac{1}{\epsilon} \Pr_{\mathbf{x} \sim \mathcal{D}}[|\mathbf{u} \cdot \mathbf{x}| \leq \epsilon/4] \leq \frac{\epsilon/2}{\epsilon} = 1/2. \quad \blacksquare$$

The next two lemmas are isotropic log-concave analogues of the uniform distribution Lemmas 7 and 8 respectively. The first one says that w.h.p. no direction  $\mathbf{a}$  has much more variance than the expected variance in any direction:

**Lemma 17** *W.h.p. over a random draw of  $\ell$  clean examples  $S_{\text{clean}}$  from  $\mathcal{D}'$ , we have*

$$\max_{\mathbf{a} \in \mathbb{S}^{n-1}} \left\{ \frac{1}{\ell} \sum_{(\mathbf{x}, y) \in S_{\text{clean}}} (\mathbf{a} \cdot \mathbf{x})^2 \right\} \leq O(1) \left( \log^2 \frac{1}{\epsilon} + \frac{n^{3/2} \log^2 \ell}{\sqrt{\ell}} \right).$$

*Proof.* By Lemma 15, for any  $\mathbf{a} \in \mathbb{S}^{n-1}$  we have

$$\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[(\mathbf{a} \cdot \mathbf{x})^2] = \Theta(\log^2(1/\epsilon)).$$

Since as remarked earlier we may assume  $\mathcal{D}'$  is supported on  $\{\mathbf{x} : \|\mathbf{x}\| \leq \sqrt{3n \log m}\}$ , we may apply Lemmas 25 and 27 (see Appendix A) with functions  $f_{\mathbf{a}}$  defined by  $f_{\mathbf{a}} = \frac{(\mathbf{a} \cdot \mathbf{x})^2}{3n \log m}$ . This completes the proof.  $\blacksquare$

The second lemma says that for a sufficiently large clean data set, w.h.p. no direction has too many examples lying too far out in that direction:

**Lemma 18** *For any  $\beta > 0$  and  $\kappa > 1$ , if  $S_{\text{clean}}$  is a set of  $\ell \geq \frac{O(1)\epsilon e^\beta (n \ln(e^{-\beta}/\epsilon) + \log m)}{(1+\kappa) \ln(1+\kappa)}$  clean examples drawn from  $\mathcal{D}'$ , then w.h.p. we have*

$$\max_{\mathbf{a} \in \mathbb{S}^{n-1}} \frac{1}{\ell} |\{\mathbf{x} \in S_{\text{clean}} : |\mathbf{a} \cdot \mathbf{x}| > \beta\}| \leq (1+\kappa) \left(\frac{1}{\epsilon}\right) e^{-\beta+1}.$$

*Proof.* Lemma 5 implies that for the original isotropic log-concave distribution  $\mathcal{D}$ , we have

$$\Pr_{\mathbf{x} \sim \mathcal{D}}[|\mathbf{a} \cdot \mathbf{x}| > \beta] \leq e^{-\beta+1}.$$

Since  $\mathcal{D}'$  is  $(1/\epsilon)$ -smooth with respect to  $\mathcal{D}$ , this implies that

$$\Pr_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| > \beta] \leq \frac{e^{-\beta+1}}{\epsilon}. \quad (6)$$

In the proof of Lemma 8, we observed that the VC-dimension of

$$\{\{\mathbf{x} : |\mathbf{a} \cdot \mathbf{x}| > \beta\} : \mathbf{a} \in \mathbf{R}^n, \beta \in \mathbf{R}\}$$

is  $O(n)$ , so applying Lemma 28 with (6) completes the proof of this lemma.  $\blacksquare$

The following is an isotropic log-concave analogue of Corollary 9, establishing that not too many clean examples are removed in the outlier removal step:

**Corollary 19** *W.h.p. over the random draw of the  $m$ -element sample  $S$  from  $EX_{\eta'}(f, \mathcal{D}')$ , the number of clean examples removed during any one execution of the outlier removal step (final substep of Step 2) in Algorithm  $A_{\text{mlcw}}$  is at most  $6m\epsilon^3/n^4$ .*

*Proof.* Since the true noise rate  $\eta$  is assumed sufficiently small, the value  $\eta' \leq \eta/\epsilon$  is at most  $\epsilon/4$ , and thus w.h.p. the number  $\ell$  of clean examples in  $S$  is at least (say)  $m/2$ . We would like to apply Lemma 18 with  $\kappa = (n/\epsilon)^{c_0-4}$  and  $\beta = c_0 \log(n/\epsilon)$ , and we may do this since we have

$$\frac{O(1)\epsilon e^\beta (n \ln(\epsilon e^\beta) + \log m)}{(1 + \kappa) \ln(1 + \kappa)} \leq \frac{O(1)\epsilon(n/\epsilon)^{c_0} n \log m}{(n/\epsilon)^{c_0-4} \log m} \leq O(1)n^5/\epsilon^3 \ll \frac{m}{2} \leq \ell$$

for a suitable fixed  $\text{poly}(n/\epsilon)$  choice of  $m$ . Since clean points are only removed if they have  $|\mathbf{a} \cdot \mathbf{x}| \geq \beta$ , Lemma 18 gives us that the number of clean points removed is at most

$$m(1 + \kappa) \cdot \frac{1}{\epsilon} e^{-\beta+1} \leq m \frac{(6/\epsilon)(n/\epsilon)^{c_0-4}}{(n/\epsilon)^{c_0}} \leq 6m\epsilon^3/n^4. \quad \blacksquare$$

The following lemma is an analogue of Lemma 10; it lower bounds the number of dirty examples that are removed in the outlier removal step.

**Lemma 20** *W.h.p. over the random draw of  $S$ , any time Algorithm  $A_{\text{mlcw}}$  executes the outlier removal step it removes at least  $\frac{m}{O(n)}$  noisy examples.*

*Proof.* Since our ultimate goal is only to prove that the algorithm succeeds for some  $\eta$  which is  $o(\epsilon)$ , we may assume without loss of generality that the original noise rate  $\eta$  is less than  $\epsilon/4$ . This means that  $\eta' < 1/4$ , and consequently a Chernoff bound gives that w.h.p. the fraction  $\hat{\eta}'$  of noisy examples in  $S$  at the beginning of the weak learner's training is at most  $1/2$ . And Lemma 17 implies that for a sufficiently large polynomial choice of  $m$ , we have that w.h.p. for all  $\mathbf{a} \in \mathbb{S}^{n-1}$ , the following holds for all the clean examples in the data before any examples were removed:

$$\sum_{(\mathbf{x}, y) \in S_{\text{clean}}} (\mathbf{a} \cdot \mathbf{x})^2 \leq cm \log^2(1/\epsilon) \quad (7)$$

where  $c$  is an absolute constant. We say that a random sample that meets all these requirements is “reasonable.” We now set the constant  $c_0$  that is used in the specification of  $A_{\text{mlcw}}$  to be  $\sqrt{2(c+1)}$ . We will now show that, for any reasonable sample  $S$ , the number of noisy examples removed during the first execution of the outlier removal step of  $A_{\text{mlcw}}$  is at least  $\frac{m}{O(n)}$ .

If we remove examples using direction  $\mathbf{w}$  then it means  $\sum_{x \in S} (\mathbf{w} \cdot \mathbf{x})^2 \geq c_0^2 m \log^2(n/\epsilon)$ . Since  $S$  is reasonable, by (7) the contribution to the sum from the clean examples that have survived until this point is at most  $cm \log^2(1/\epsilon)$  so we must have

$$\sum_{(\mathbf{x}, y) \in S_{\text{dirty}}} (\mathbf{w} \cdot \mathbf{x})^2 \geq (c_0^2 - c)m \log^2(n/\epsilon).$$

Let  $S_{\text{dirty}} = N \cup F$  where  $N$  is the examples  $(\mathbf{x}, y)$  for which  $\mathbf{x}$  satisfies  $(\mathbf{w} \cdot \mathbf{x})^2 \leq c_0^2 \log^2(n/\epsilon)$  and  $F$  is the other points. We have

$$\sum_{(\mathbf{x}, y) \in N} (\mathbf{w} \cdot \mathbf{x})^2 \leq c_0^2 \hat{\eta}' m \log^2(n/\epsilon).$$

and so, since  $\|\mathbf{x}\| \leq \sqrt{3n \log m}$  implies that  $(\mathbf{w} \cdot \mathbf{x})^2 \leq 3n \log m$  for all unit length  $\mathbf{w}$ , we have

$$\begin{aligned}
 |F| &\geq \sum_{(\mathbf{x}, y) \in F} \frac{(\mathbf{w} \cdot \mathbf{x})^2}{3n \log m} = \sum_{(\mathbf{x}, y) \in S_{\text{dirty}}} \frac{(\mathbf{w} \cdot \mathbf{x})^2}{3n \log m} - \sum_{(\mathbf{x}, y) \in N} \frac{(\mathbf{w} \cdot \mathbf{x})^2}{3n \log m} \\
 &\geq \frac{(c_0^2 - c)m \log^2(n/\epsilon) - c_0^2 \eta' m \log^2(n/\epsilon)}{3n \log m} \\
 &\geq \frac{m \log^2(n/\epsilon)}{3n \log m} \\
 &\geq \frac{m}{O(n)}
 \end{aligned}$$

where the next-to-last inequality uses  $\eta' \leq 1/2$  and  $c_0 = \sqrt{2(c+1)}$ , and the final one uses  $m = O(\text{poly}(n/\epsilon))$ . The points in  $F$  are precisely the ones that are removed, and thus the lemma is proved.  $\blacksquare$

#### 4.5 Proof of Lemma 14

We first note that Lemma 20 implies that w.h.p. the weak learner must terminate after at most  $O(n)$  iterations of outlier removal.

Let  $\mathbf{u}$  be the unit length normal vector of the separating halfspace for the target function  $f$ . Recall that we have assumed without loss of generality that  $\|\mathbf{x}\| \leq \sqrt{3n \log m}$  for all  $\mathbf{x}$  in the training set, so that  $\|\mathbf{v}\| \leq \sqrt{3n \log m}$ , and thus the advantage of  $h$  with respect to  $\mathcal{D}'$  can be expressed as

$$\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[h(\mathbf{x})f(\mathbf{x})] = \frac{\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[(\mathbf{v} \cdot \mathbf{x})f(\mathbf{x})]}{3n \log m} \quad (8)$$

and so we shall work on lower bounding  $\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[(\mathbf{v} \cdot \mathbf{x})f(\mathbf{x})]$ .

As in the proof of Lemma 12, let

- $S_{\text{clean}}$  be all of the clean examples in the initial sample  $S$ , and  $S'_{\text{clean}}$  be those that are not removed in any stage of outlier removal;
- $S_{\text{dirty}}$  be all of the dirty examples in the initial sample  $S$ , and  $S'_{\text{dirty}}$  be those that are not removed in any stage of outlier removal;
- $\eta' = \frac{|S'_{\text{dirty}}|}{|S'_{\text{clean}} \cup S'_{\text{dirty}}|}$ , i.e., the noise rate among the examples that survive until the end of training of the weak learner, and
- $\alpha = \frac{|S_{\text{clean}} - S'_{\text{clean}}|}{|S'_{\text{clean}} \cup S'_{\text{dirty}}|}$ , the ratio of the number of clean points that were erroneously removed to the size of the final surviving data set.

As before we write  $S'$  for  $S'_{\text{clean}} \cup S'_{\text{dirty}}$ . Also as before, let  $\mathbf{v}_{\text{clean}}$  be the average of *all* the clean examples,  $\mathbf{v}'_{\text{dirty}}$  be the average of the dirty (noisy) examples that were not deleted, and  $\mathbf{v}_{\text{del}}$  be the average of the clean examples that were deleted. Then arguing exactly as before, we have

$$\mathbf{v} = (1 - \eta' + \alpha)\mathbf{v}_{\text{clean}} + \eta'\mathbf{v}'_{\text{dirty}} - \alpha\mathbf{v}_{\text{del}}. \quad (9)$$



The expectation of  $\mathbf{v}_{\text{clean}}$  will play a special role in the analysis:

$$\mathbf{v}_{\text{clean}}^* \stackrel{\text{def}}{=} \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[f(\mathbf{x})\mathbf{x}].$$

Once again, we will demonstrate the limited effect of  $\mathbf{v}'_{\text{dirty}}$  by bounding its length. This time, the outlier removal enforces the fact that, for any  $\mathbf{w} \in \mathbb{S}^{n-1}$ , we have

$$\sum_{(\mathbf{x}, y) \in S} (\mathbf{w} \cdot \mathbf{x})^2 \leq c_0^2 m \log^2(n/\epsilon).$$

Applying this for the unit vector  $\mathbf{w}$  in the direction of  $\mathbf{v}'_{\text{dirty}}$  as was done in Lemma 12, this implies

$$\|\mathbf{v}'_{\text{dirty}}\| \leq c_0 \log(n/\epsilon) \sqrt{\frac{m}{|S'_{\text{dirty}}|}}.$$

Next, let us apply this to bound an expression that captures the average harm done by  $\mathbf{v}'_{\text{dirty}}$ .

$$\begin{aligned} |\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[f(\mathbf{x})(\mathbf{v}'_{\text{dirty}} \cdot \mathbf{x})]| &= |\mathbf{v}'_{\text{dirty}} \cdot \mathbf{v}_{\text{clean}}^*| \\ &\leq c_0 \log(n/\epsilon) \sqrt{\frac{m}{|S'_{\text{dirty}}|}} \|\mathbf{v}_{\text{clean}}^*\|. \end{aligned} \quad (10)$$

To show that  $\mathbf{v}_{\text{clean}}$  plays a relatively large role, it is helpful to lower bound the length of  $\mathbf{v}_{\text{clean}}^*$ . We do this by lower bounding the length of its projection onto the unit normal vector  $\mathbf{u}$  of the target as follows:

$$\mathbf{v}_{\text{clean}}^* \cdot \mathbf{u} = \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[(f(\mathbf{x})\mathbf{x}) \cdot \mathbf{u}] = \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[\text{sgn}(\mathbf{u} \cdot \mathbf{x})(\mathbf{x} \cdot \mathbf{u})] = \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{x} \cdot \mathbf{u}|] \geq \epsilon/8,$$

by Lemma 16. Since  $\mathbf{u}$  is unit length, this implies

$$\|\mathbf{v}_{\text{clean}}^*\| \geq \epsilon/8. \quad (11)$$

Armed with this bound, we can now lower bound the benefit imparted by  $\mathbf{v}_{\text{clean}}$ :

$$\begin{aligned} \mathbf{E}_{\mathbf{z} \sim \mathcal{D}'}[f(\mathbf{z})(\mathbf{v}_{\text{clean}} \cdot \mathbf{z})] &= \frac{1}{S_{\text{clean}}} \sum_{(\mathbf{x}, y) \in S_{\text{clean}}} \mathbf{E}_{\mathbf{z} \sim \mathcal{D}'}[yf(\mathbf{z})(\mathbf{x} \cdot \mathbf{z})] \\ &= \frac{1}{S_{\text{clean}}} \sum_{(\mathbf{x}, y) \in S_{\text{clean}}} (y\mathbf{x}) \cdot \mathbf{v}_{\text{clean}}^*. \end{aligned}$$

Since  $\mathbf{E}[(y\mathbf{x}) \cdot \mathbf{v}_{\text{clean}}^*] = \|\mathbf{v}_{\text{clean}}^*\|^2$ , and  $(y\mathbf{x}) \cdot \mathbf{v}_{\text{clean}}^* \in [-3n \log m, 3n \log m]$ , a Hoeffding bound implies that w.h.p.

$$\mathbf{E}_{\mathbf{z} \sim \mathcal{D}'}[f(\mathbf{z})(\mathbf{v}_{\text{clean}} \cdot \mathbf{z})] \geq \|\mathbf{v}_{\text{clean}}^*\|^2 - O(n \log^{3/2} m) / \sqrt{|S_{\text{clean}}|}.$$

Since the noise rate  $\eta'$  is at most  $\eta/\epsilon$  and  $\eta$  certainly less than  $\epsilon/4$  as discussed above, another Hoeffding bound gives that w.h.p.  $|S_{\text{clean}}|$  is at least  $m/2$ ; thus for a suitably large polynomial choice of  $m$ , using (11) we have

$$\mathbf{E}_{\mathbf{z} \sim \mathcal{D}'}[f(\mathbf{z})(\mathbf{v}_{\text{clean}} \cdot \mathbf{z})] \geq \|\mathbf{v}_{\text{clean}}^*\|^2 - O(n \log^{3/2} m) / \sqrt{m/2} \geq \frac{\|\mathbf{v}_{\text{clean}}^*\|^2}{2}. \quad (12)$$

Now we are ready to put our bounds together and lower bound the advantage of  $\mathbf{v}$ . We have

$$\begin{aligned} \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[f(\mathbf{x})(\mathbf{v} \cdot \mathbf{x})] &= (1 - \eta' + \alpha)\mathbf{E}[f(\mathbf{x})(\mathbf{v}_{\text{clean}} \cdot \mathbf{x})] \\ &\quad + \eta'\mathbf{E}[f(\mathbf{x})(\mathbf{v}'_{\text{dirty}} \cdot \mathbf{x})] - \alpha\mathbf{E}[f(\mathbf{x})(\mathbf{v}_{\text{del}} \cdot \mathbf{x})]. \end{aligned}$$

We bound each of the three contributions in turn. First, using  $1 - \eta' \geq 1/2$  and (12), we have  $(1 - \eta' + \alpha)\mathbf{E}[f(\mathbf{x})(\mathbf{v}_{\text{clean}} \cdot \mathbf{x})] \geq \frac{\|\mathbf{v}_{\text{clean}}^*\|^2}{4}$ .

Next, by (10), we have

$$|\eta'\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[f(\mathbf{x})(\mathbf{v}'_{\text{dirty}} \cdot \mathbf{x})]| \leq c_0 \log(n/\epsilon) \sqrt{2\eta'} \|\mathbf{v}_{\text{clean}}^*\|.$$

Since we may assume that  $\eta \leq c'\epsilon^3/\log^2(n/\epsilon)$  for as small a fixed constant  $c'$  as we like (recall the overall bound of Theorem 2), we get

$$c_0 \log(n/\epsilon) \sqrt{2\eta'} \|\mathbf{v}_{\text{clean}}^*\| \leq (\epsilon/64) \|\mathbf{v}_{\text{clean}}^*\|$$

(for a suitably small constant choice of  $c'$ ), and this is less than  $\frac{\|\mathbf{v}_{\text{clean}}^*\|^2}{8}$  since  $\|\mathbf{v}_{\text{clean}}^*\| \geq \epsilon/8$ .

Finally Corollary 19, together with the fact that there are at most  $O(n)$  iterations of outlier removal and the final surviving data set is of size at least  $m/4$ , gives us that  $\alpha \leq \frac{O(n)(6m\epsilon^3/n^4)}{m/4}$ , which (recalling that both  $\mathbf{v}_{\text{del}}$  and all  $\mathbf{x}$  in the support of  $\mathcal{D}'$  have norm at most  $\sqrt{3n \log m}$ ) means that  $|\alpha\mathbf{E}[f(\mathbf{x})(\mathbf{v}_{\text{del}} \cdot \mathbf{x})]| = o(\epsilon^2)$ .

Combining all these bounds, we get

$$\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[f(\mathbf{x})(\mathbf{v} \cdot \mathbf{x})] \geq \frac{\|\mathbf{v}_{\text{clean}}^*\|^2}{4} - \frac{\|\mathbf{v}_{\text{clean}}^*\|^2}{8} - o(\epsilon^2) \geq \frac{\epsilon^2}{1024}$$

by (11). Together with (8), the proof of Lemma 14 is completed.

## 5. Learning under isotropic log-concave distributions with adversarial label noise

### 5.1 The Model

We now define the model of learning with adversarial label noise under isotropic log-concave distributions. In this model the learning algorithm has access to an oracle that provides independent random examples drawn according to a fixed distribution  $P$  on  $\mathbf{R}^n \times \{-1, 1\}$ , where

- the marginal distribution over  $\mathbf{R}^n$  is isotropic log-concave, and
- there is a halfspace  $f$  such that  $\Pr_{(\mathbf{x}, y) \sim P}[f(\mathbf{x}) \neq y] = \eta$ .

The parameter  $\eta$  is the *noise rate*. As usual, the goal of the learner is to output a hypothesis  $h$  such that  $\Pr_{(\mathbf{x}, y) \sim \mathcal{D}}[h(\mathbf{x}) \neq y] \leq \epsilon$ ; if an algorithm achieves this goal, we say it learns to accuracy  $1 - \epsilon$  in the presence of adversarial label noise at rate  $\eta$ .

## 5.2 The Algorithm

Like the algorithm  $A_{\text{mlc}}$  considered in the last section, the algorithm  $A_{\text{alc}}$  studied in this section applies the smooth boosting algorithm of Lemma 13 to a weak learner that performs averaging. The weak learner  $A_{\text{alcw}}$  behaves as follows:

**Algorithm  $A_{\text{alcw}}$ :**

1. Draw a set  $S$  of  $m$  examples according to  $P'$  (the oracle for a modified distribution provided by the boosting algorithm).
2. Remove all examples  $(\mathbf{x}, y)$  such that  $\|\mathbf{x}\| > \sqrt{3n \log m}$  from  $S$ .
3. Let  $\mathbf{v} = \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} y \mathbf{x}$ . Return the confidence-rated classifier  $h$  defined by  $h(\mathbf{x}) = \frac{\mathbf{v} \cdot \mathbf{x}}{3n \log m}$  if  $|\mathbf{v} \cdot \mathbf{x}| \leq 3n \log m$ , and  $h(\mathbf{x}) = \text{sgn}(\mathbf{v} \cdot \mathbf{x})$  otherwise.

## 5.3 Claim about the weak learner

As in the previous section, the heart of our analysis will be to analyze the weak learner. We omit discussing the application of the smooth boosting algorithm here, as it is nearly identical to Section 4.

**Lemma 21** *Suppose Algorithm  $A_{\text{alcw}}$  is run using  $P'$  as the source of labeled examples, where  $P'$  is a distribution that is  $(1/\epsilon)$ -smooth with respect to a joint distribution  $P$  on  $\mathbf{R}^n \times \{-1, 1\}$  whose marginal  $\mathcal{D}'$  on  $\mathbf{R}^n$  is isotropic and log-concave. Further, assume there exists a linear threshold function  $f$  such that  $\Pr_{(\mathbf{x}, y) \sim P'}[f(\mathbf{x}) \neq y] \leq \eta/\epsilon$  and  $\eta \leq \Omega(\frac{\epsilon^3}{\log(1/\epsilon)})$ . Then with high probability,  $A_{\text{alcw}}$  outputs a hypothesis with advantage  $\Omega(\frac{\epsilon^2}{n \log(n/\epsilon)})$ .*

## 5.4 Lemmas in support of Lemma 21

During this section, let us focus our attention on a single call to the weak learner. Let  $P'$  be a distribution as in Lemma 21 and let  $\mathcal{D}'$  be the marginal on  $\mathbf{R}^n$ . We observe that since  $P'$  is  $(1/\epsilon)$ -smooth with respect to  $P$ , the marginal  $\mathcal{D}'$  of  $P'$  is  $(1/\epsilon)$ -smooth with respect to the marginal  $\mathcal{D}$  of  $P$ .

As in Section 4, we may assume that the support of  $\mathcal{D}'$  lies entirely on  $\mathbf{x}$  such that  $\|\mathbf{x}\| \leq \sqrt{3n \log m}$  (this negligibly affects the final bounds obtained in our analyses).

The following technical lemma will be used to limit the extent to which the distribution  $P'$  can concentrate a lot of noise in one direction.

**Lemma 22** *Let  $E$  be any event with positive probability under  $\mathcal{D}'$ , and let  $\kappa = \mathcal{D}'(E)$ . For any unit length  $\mathbf{a} \in \mathbf{R}^n$ ,  $\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid E] = O(\log \frac{1}{\kappa \epsilon})$ .*

*Proof.* Let  $\beta$  be such that  $\Pr_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| > \beta] = \kappa$ . By Lemmas 4 and 5, together with the fact that  $\mathcal{D}'$  is  $(1/\epsilon)$  smooth with respect to  $\mathcal{D}$ , we have

$$\kappa \leq \frac{1}{\epsilon} e^{-\beta+1}$$

which implies  $\beta \leq 1 + \ln(\frac{1}{\epsilon \kappa})$ .

Let  $F$  be the event that  $|\mathbf{a} \cdot \mathbf{x}| > \beta$ . We will show that  $\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid E] \leq \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid F]$ , and then bound  $\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid F]$ . If  $\Pr[(E - F) \cup (F - E)] = 0$ , then, obviously,  $\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid E] = \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid F]$ . Suppose  $\Pr[(E - F) \cup (F - E)] > 0$ . Then

$$\begin{aligned} & \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid E] \\ &= \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid E \cap F] \Pr[E \cap F] + \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid E - F] \Pr[E - F] \\ &= \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid E \cap F] \Pr[E \cap F] + \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid E - F] \Pr[F - E] \\ &\quad \text{(because } \Pr[E] = \Pr[F]) \\ &< \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid E \cap F] \Pr[E \cap F] + \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid F - E] \Pr[F - E], \end{aligned}$$

because for every  $\mathbf{x} \in E - F$  and every  $\mathbf{x}' \in F - E$ ,

$$|\mathbf{a} \cdot \mathbf{x}| \leq \beta < |\mathbf{a} \cdot \mathbf{x}'|.$$

But

$$\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid E \cap F] \Pr[E \cap F] + \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid F - E] \Pr[F - E] = \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid F],$$

so

$$\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid E] < \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid F]. \quad (13)$$

Now, setting  $b = \lfloor \beta \rfloor$ , we have

$$\begin{aligned} \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \mid F] &\leq \frac{1}{\mathcal{D}'(F)} \sum_{i=b}^{\infty} (i+1) \Pr_{\mathbf{x} \sim \mathcal{D}'}[|\mathbf{a} \cdot \mathbf{x}| \in (i, i+1]] \\ &\leq \frac{1}{\mathcal{D}'(F)} \sum_{i=b}^{\infty} (i+1) e^{-i+1} \\ &= \frac{1}{\mathcal{D}'(F)} \left( O\left(\frac{e^{-b} b}{\epsilon}\right) \right) \\ &= O(b), \end{aligned}$$

since  $\mathcal{D}'(F) = \Theta(e^{-b}/\epsilon)$ . Combining with (13) completes the proof.  $\blacksquare$

## 5.5 Proof of Lemma 21

Fix some halfspace  $f$  such that  $\Pr_{(\mathbf{x}, y) \sim P}[f(\mathbf{x}) \neq y] = \eta$ , and let  $\mathbf{u}$  be the unit normal vector of its separating hyperplane.

Let  $P'$  be the joint distribution given to  $A_{\text{alcw}}$  and let  $\mathcal{D}'$  be its marginal on  $\mathbf{R}^n$ . As noted in the previous subsection,  $\mathcal{D}'$  is  $(1/\epsilon)$ -smooth with respect to the original marginal distribution  $\mathcal{D}$  of  $P$ .

First, we bound the advantage of the hypothesis  $h$  with respect to  $P'$  in terms of the tendency of  $h$  to agree with the best linear function  $f$ :

$$\mathbf{E}_{(\mathbf{x}, y) \sim P'}[h(\mathbf{x})y] \geq \mathbf{E}_{(\mathbf{x}, y) \sim P'}[h(\mathbf{x})f(\mathbf{x})] - \eta = \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[h(\mathbf{x})f(\mathbf{x})] - \eta. \quad (14)$$

Furthermore, as we have assumed without loss of generality that  $\|\mathbf{x}\| \leq \sqrt{3n \log m}$  for all examples in the training set, and therefore that  $\|\mathbf{v}\| \leq \sqrt{3n \log m}$ , we have

$$\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[h(\mathbf{x})f(\mathbf{x})] = \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'} \left[ \frac{f(\mathbf{x})(\mathbf{x} \cdot \mathbf{v})}{3n \log m} \right] \quad (15)$$

so we will work on bounding  $\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[f(\mathbf{x})(\mathbf{x} \cdot \mathbf{v})]$ .

Let  $P'_{\text{clean}}$  be obtained by conditioning a random draw  $(\mathbf{x}, y)$  from  $P'$  on the event that  $f(\mathbf{x}) = y$ . Define  $P'_{\text{dirty}}$  analogously, and let  $\mathcal{D}'_{\text{clean}}$  and  $\mathcal{D}'_{\text{dirty}}$  be the corresponding marginals on  $\mathbf{R}^n$ . Let

$$\begin{aligned} \mathbf{v}_{\text{dirty}}^* &= \mathbf{E}_{(\mathbf{x}, y) \sim P'_{\text{dirty}}}[y\mathbf{x}] \\ \mathbf{v}_{\text{correct}}^* &= \mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[f(\mathbf{x})\mathbf{x}]. \end{aligned}$$

Note that the linearity of expectation implies that

$$\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[f(\mathbf{x})(\mathbf{x} \cdot \mathbf{v})] = (\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[f(\mathbf{x})(\mathbf{x})]) \cdot \mathbf{v} = \mathbf{v}_{\text{correct}}^* \cdot \mathbf{v} = \frac{1}{m} \sum_{(\mathbf{x}, y) \in S} \mathbf{v}_{\text{correct}}^* \cdot (y\mathbf{x}). \quad (16)$$

Equation (16) expresses  $\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[f(\mathbf{x})(\mathbf{x} \cdot \mathbf{v})]$ , which is closely related to the advantage of  $h$  through (15) and (14), as a sum of independent random variables, one for each example. We will bound  $\mathbf{E}_{\mathbf{x} \sim \mathcal{D}'}[f(\mathbf{x})(\mathbf{x} \cdot \mathbf{v})]$  by bounding the expected effect of a random example on its value, and applying a Hoeffding bound.

Let  $\eta' = \Pr_{(\mathbf{x}, y) \sim P'}[f(\mathbf{x}) \neq y]$ . Since  $P'$  is  $1/\epsilon$ -smooth with respect to  $P$ , we have  $\eta' \leq \eta/\epsilon$ . We can rearrange the effect of a random example as follows

$$\begin{aligned} \mathbf{E}_{(\mathbf{x}, y) \sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (y\mathbf{x})] &= (1 - \eta')\mathbf{E}_{(\mathbf{x}, y) \sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (f(\mathbf{x})\mathbf{x})|y = f(\mathbf{x})] \\ &\quad + \eta'\mathbf{E}_{(\mathbf{x}, y) \sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (-f(\mathbf{x})\mathbf{x})|y \neq f(\mathbf{x})] \\ &= (1 - \eta')\mathbf{E}_{(\mathbf{x}, y) \sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (f(\mathbf{x})\mathbf{x})|y = f(\mathbf{x})] \\ &\quad + \eta'\mathbf{E}_{(\mathbf{x}, y) \sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (f(\mathbf{x})\mathbf{x})|y \neq f(\mathbf{x})] \\ &\quad - \eta'\mathbf{E}_{(\mathbf{x}, y) \sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (f(\mathbf{x})\mathbf{x})|y \neq f(\mathbf{x})] \\ &\quad + \eta'\mathbf{E}_{(\mathbf{x}, y) \sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (-f(\mathbf{x})\mathbf{x})|y \neq f(\mathbf{x})]. \end{aligned} \quad (17)$$

Since

$$\begin{aligned} &\mathbf{E}_{(\mathbf{x}, y) \sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (f(\mathbf{x})\mathbf{x})] \\ &= \eta'\mathbf{E}_{(\mathbf{x}, y) \sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (f(\mathbf{x})\mathbf{x})|y \neq f(\mathbf{x})] + (1 - \eta')\mathbf{E}_{(\mathbf{x}, y) \sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (f(\mathbf{x})\mathbf{x})|y = f(\mathbf{x})], \end{aligned}$$

by replacing the first two terms of (17) with  $\mathbf{E}_{(\mathbf{x}, y) \sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (f(\mathbf{x})\mathbf{x})]$ , we get

$$\begin{aligned} \mathbf{E}_{(\mathbf{x}, y) \sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (y\mathbf{x})] &= \mathbf{E}_{(\mathbf{x}, y) \sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (f(\mathbf{x})\mathbf{x})] \\ &\quad - \eta'\mathbf{E}_{(\mathbf{x}, y) \sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (f(\mathbf{x})\mathbf{x})|y \neq f(\mathbf{x})] \\ &\quad + \eta'\mathbf{E}_{(\mathbf{x}, y) \sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (-f(\mathbf{x})\mathbf{x})|y \neq f(\mathbf{x})] \\ &= \mathbf{E}_{(\mathbf{x}, y) \sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (f(\mathbf{x})\mathbf{x})] \\ &\quad - 2\eta'\mathbf{E}_{(\mathbf{x}, y) \sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (f(\mathbf{x})\mathbf{x})|y \neq f(\mathbf{x})] \end{aligned}$$

Twice applying the linearity of expectation, we get

$$\begin{aligned}
 \mathbf{E}_{(\mathbf{x},y)\sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (y\mathbf{x})] &= \|\mathbf{v}_{\text{correct}}^*\|^2 - 2\eta' \mathbf{E}_{(\mathbf{x},y)\sim P'}[\mathbf{v}_{\text{correct}}^* \cdot (f(\mathbf{x})\mathbf{x}) | y \neq f(\mathbf{x})] \\
 &= \|\mathbf{v}_{\text{correct}}^*\|^2 - 2\eta' \mathbf{v}_{\text{correct}}^* \cdot \mathbf{v}_{\text{dirty}}^* \\
 &\geq \|\mathbf{v}_{\text{correct}}^*\|^2 - 2\eta' \|\mathbf{v}_{\text{correct}}^*\| \cdot \|\mathbf{v}_{\text{dirty}}^*\| \\
 &\geq \frac{1}{2} \|\mathbf{v}_{\text{correct}}^*\|^2 - 4(\eta')^2 \|\mathbf{v}_{\text{dirty}}^*\|^2,
 \end{aligned}$$

The last line follows from the fact that  $q^2 - qr \geq (q^2 - r^2)/2$  for all real  $q, r$ .

So now our goals are a lower bound on  $\|\mathbf{v}_{\text{correct}}^*\|$  and an upper bound on  $\|\mathbf{v}_{\text{dirty}}^*\|$ .

We can lower bound  $\|\mathbf{v}_{\text{correct}}^*\|$  essentially the same way we did before, by lower bounding its projection onto the “target” normal vector  $\mathbf{u}$ :

$$\mathbf{v}_{\text{correct}}^* \cdot \mathbf{u} = \mathbf{E}_{(\mathbf{x},y)\sim P'}[(f(\mathbf{x})\mathbf{x}) \cdot \mathbf{u}] = \mathbf{E}_{(\mathbf{x},y)\sim P'}[\text{sgn}(\mathbf{u} \cdot \mathbf{x})(\mathbf{x} \cdot \mathbf{u})] = \mathbf{E}_{(\mathbf{x},y)\sim P'}[|\mathbf{x} \cdot \mathbf{u}|] \geq \epsilon/16, \tag{18}$$

by Lemma 16.

We upper bound  $\|\mathbf{v}_{\text{dirty}}^*\|$  as follows:

$$\begin{aligned}
 \|\mathbf{v}_{\text{dirty}}^*\|^2 &= \mathbf{v}_{\text{dirty}}^* \cdot E_{\mathbf{x}\sim \mathcal{D}'_{\text{dirty}}}[-f(\mathbf{x})\mathbf{x}] \\
 &= \|\mathbf{v}_{\text{dirty}}^*\| \cdot E_{\mathbf{x}\sim \mathcal{D}'_{\text{dirty}}} \left[ \left( \frac{\mathbf{v}_{\text{dirty}}^*}{\|\mathbf{v}_{\text{dirty}}^*\|} \right) \cdot (-f(\mathbf{x})\mathbf{x}) \right] \\
 &\leq \|\mathbf{v}_{\text{dirty}}^*\| \cdot E_{\mathbf{x}\sim \mathcal{D}'_{\text{dirty}}} \left[ \left| \left( \frac{\mathbf{v}_{\text{dirty}}^*}{\|\mathbf{v}_{\text{dirty}}^*\|} \right) \cdot \mathbf{x} \right| \right] \\
 &\leq \|\mathbf{v}_{\text{dirty}}^*\| O(\log(1/(\eta'\epsilon)))
 \end{aligned}$$

by Lemma 22. Thus  $\|\mathbf{v}_{\text{dirty}}^*\| \leq O(\log(1/(\eta'\epsilon)))$ .

Combining this with (18) and (16) we have that if

$$\eta' \sqrt{\log(1/(\eta'\epsilon))} \leq c\epsilon^2$$

for a suitably small constant  $c$ , then  $\mathbf{E}_{\mathbf{x}\sim \mathcal{D}'}[f(\mathbf{x})(\mathbf{x} \cdot \mathbf{v})]$  is a sum of  $m$  i.i.d. random variables, each with mean at least  $\Omega(\epsilon^2)$ , and coming from an interval of length  $O(n \log m)$ . Applying the standard Hoeffding bound, polynomially many examples suffice for  $\mathbf{E}_{\mathbf{x}\sim \mathcal{D}'}[f(\mathbf{x})(\mathbf{x} \cdot \mathbf{v})] \geq \Omega(\epsilon^2)$ . Combining with (15) and (14) completes the proof.

## 6. Conclusion

Our algorithms use boosting together with a confidence-rated weak learner that perform a simple averaging of labeled examples. As shown in earlier work (Servedio, 2002, 2003) there are close connections between such an approach and the Perceptron algorithm. It seems likely that the Perceptron could be used as an alternative to boosting and averaging in our algorithms; it would be interesting to see if a Perceptron-based approach has any theoretical or empirical advantages over the algorithms we give in this paper.

More generally, there are relatively few algorithms for learning interesting classes of functions in the presence of malicious noise. We hope that our results will help lead to the development of more efficient algorithms for this challenging noise model.

As a challenge for future work, we pose the following question: do there exist computationally efficient algorithms for learning halfspaces under *arbitrary* distributions in the presence of malicious noise? As of now no better results are known for this problem than the generic conversions of (Kearns and Li, 1993), which can be applied to any concept class. We feel that even a small improvement in the malicious noise rate that can be handled for halfspaces would be a very interesting result.

## Acknowledgement

We are grateful to the anonymous reviewers for their comments.

## References

- S. Arora, L. Babai, J. Stern, and Z. Sweedyk. The hardness of approximate optima in lattices, codes, and systems of linear equations. *Proceedings of the 34th Annual Symposium on Foundations of Computer Science*, pages 724–733, 1993.
- A. Blum, A. Frieze, R. Kannan, and S. Vempala. A polynomial time algorithm for learning noisy linear threshold functions. *Algorithmica*, 22(1/2):35–52, 1997.
- A. Blumer, A. Ehrenfeucht, D. Haussler, and M. Warmuth. Learnability and the Vapnik-Chervonenkis dimension. *Journal of the ACM*, 36(4):929–965, 1989.
- S. C. Brubaker. *Extensions of Principle Components Analysis*. PhD thesis, Georgia Institute of Technology, 2009.
- N. H. Bshouty, Y. Li, and P. M. Long. Using the doubling dimension to analyze the generalization of learning algorithms. *Journal of Computer & System Sciences*, 75(6):323–335, 2009.
- J. Dunagan and S. Vempala. Optimal outlier removal in high-dimensional spaces. *J. Computer & System Sciences*, 68(2):335–373, 2004.
- V. Feldman, P. Gopalan, S. Khot, and A. Ponnuswami. New results for learning noisy parities and halfspaces. In *Proc. 47th IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 563–576, 2006.
- Yoav Freund and Robert E. Schapire. Large margin classification using the perceptron algorithm. *Machine Learning*, 37(3):277–296, 1999.
- Dmitry Gavinsky. Optimally-smooth adaptive boosting and application to agnostic learning. *Journal of Machine Learning Research*, 4:101–117, 2003.
- V. Guruswami and P. Raghavendra. Hardness of learning halfspaces with noise. In *Proc. 47th IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 543–552. IEEE Computer Society, 2006.
- D. Haussler, M. Kearns, N. Littlestone, and M. Warmuth. Equivalence of models for polynomial learnability. *Information and Computation*, 95(2):129–161, 1991.

- I.T. Jolliffe. *Principal Component Analysis*. Springer Series in Statistics, 2002.
- A. Kalai, A. Klivans, Y. Mansour, and R. Servedio. Agnostically learning halfspaces. *SIAM Journal on Computing*, 37(6):1777–1805, 2008.
- M. Kearns and M. Li. Learning in the presence of malicious errors. *SIAM Journal on Computing*, 22(4):807–837, 1993.
- M. Kearns, R. Schapire, and L. Sellie. Toward Efficient Agnostic Learning. *Machine Learning*, 17(2/3):115–141, 1994.
- A. Klivans and A. Sherstov. Cryptographic hardness for learning intersections of halfspaces. In *Proc. 47th IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 553–562, 2006.
- N. Littlestone. Learning quickly when irrelevant attributes abound: a new linear-threshold algorithm. *Machine Learning*, 2:285–318, 1987.
- N. Littlestone. Redundant noisy attributes, attribute errors, and linear-threshold learning using Winnow. In *Proceedings of the Fourth Annual Workshop on Computational Learning Theory*, pages 147–156, 1991.
- L. Lovász and S. Vempala. The geometry of logconcave functions and sampling algorithms. *Random Structures and Algorithms*, 30(3):307–358, 2007.
- W. Maass and G. Turan. How fast can a threshold gate learn? In *Computational Learning Theory and Natural Learning Systems: Volume I: Constraints and Prospects*, pages 381–414. MIT Press, 1994.
- Y. Mansour and M. Parnas. Learning conjunctions with noise under product distributions. *Information Processing Letters*, 68(4):189–196, 1998.
- A. Novikoff. On convergence proofs on perceptrons. In *Proceedings of the Symposium on Mathematical Theory of Automata*, volume XII, pages 615–622, 1962.
- D. Pollard. *Convergence of Stochastic Processes*. Springer Verlag, 1984.
- F. Rosenblatt. The Perceptron: a probabilistic model for information storage and organization in the brain. *Psychological Review*, 65:386–407, 1958.
- R. Schapire and Y. Singer. Improved boosting algorithms using confidence-rated predictions. *Machine Learning*, 37:297–336, 1999.
- R. Servedio. *Efficient Algorithms in Computational Learning Theory*. PhD thesis, Harvard University, 2001.
- R. Servedio. PAC analogues of Perceptron and Winnow via boosting the margin. *Machine Learning*, 47(2/3):133–151, 2002.
- R. Servedio. Smooth boosting and learning with malicious noise. *Journal of Machine Learning Research*, 4:633–648, 2003.



- J. Shawe-Taylor and N. Cristianini. *An introduction to support vector machines*. Cambridge University Press, 2000.
- M. Talagrand. Sharper bounds for Gaussian and empirical processes. *Annals of Probability*, 22: 28–76, 1994.
- L. Valiant. Learning disjunctions of conjunctions. In *Proceedings of the Ninth International Joint Conference on Artificial Intelligence*, pages 560–566, 1985.
- H. Xu, C. Caramanis, and S. Mannor. Principal component analysis with contaminated data: The high dimensional case. *Journal of Machine Learning Research*, 2009. To appear.

## Appendix A. Proof of Lemma 7

Let us start with a couple of definitions and a couple of bounds from the literature.

**Definition 23 (VC-dimension)** A set  $F$  of  $\{-1, 1\}$ -valued functions defined on a common domain  $X$  shatters  $x_1, \dots, x_d$  if every sequence  $y_1, \dots, y_d \in \{-1, 1\}$  of function values has a function  $f$  such that  $f(x_1) = y_1, \dots, f(x_d) = y_d$ . The VC-dimension of  $F$  is the size of the largest set shattered by  $F$ .

**Definition 24 (pseudo-dimension)** For a set  $F$  of real-valued functions defined on a common domain  $X$ , the pseudo-dimension of  $F$  is the VC-dimension of  $\{\text{sign}(f(\cdot) - \theta) : f \in F, \theta \in \mathbf{R}\}$ .

**Lemma 25 ((Pollard, 1984; Talagrand, 1994))** Let  $F$  be a set of real-valued functions defined on a common domain  $X$  taking values in  $[0, 1]$ , and let  $d$  be the pseudo-dimension of  $F$ . Let  $\mathcal{D}$  be a probability distribution over  $X$ . Then if  $x_1, \dots, x_m$  are obtained by drawing  $m$  times independently according to  $\mathcal{D}$ , for any  $\delta > 0$ ,

$$\Pr \left[ \exists f \in F, \frac{1}{m} \sum_{s=1}^m f(x_s) > E_{\mathcal{D}}[f] + c \sqrt{\frac{d + \log(1/\delta)}{m}} \right] \leq \delta,$$

where  $c > 0$  is an absolute constant.

**Lemma 26 (see Blumer et al. (1989))** The VC-dimension of unions of two halfspaces is  $O(n)$ .

Now, let us bound the pseudo-dimension of the class of functions that we need.

**Lemma 27** Let  $F_n$  consist of the functions  $f$  from  $\mathbf{R}^n$  to  $\mathbf{R}$  which can be defined by  $f(\mathbf{x}) = (\mathbf{a} \cdot \mathbf{x})^2$  for some  $\mathbf{a} \in \mathbf{R}^n$ . The pseudo-dimension of  $F_n$  is at most  $O(n)$ .

*Proof.* According to the definition, the pseudo dimension of  $F_n$  is the VC-dimension of the set  $G_n$  of  $\{-1, 1\}$ -valued functions  $g_{\mathbf{a}, \theta}$  defined by  $g_{\mathbf{a}, \theta}(\mathbf{x}) = \text{sign}((\mathbf{a} \cdot \mathbf{x})^2 - \theta)$ . Each  $g_{\mathbf{a}, \theta}$  is equivalent to an OR of two halfspaces:

$$\mathbf{a} \cdot \mathbf{x} \geq \sqrt{\theta} \quad \text{OR} \quad (-\mathbf{a}) \cdot \mathbf{x} \geq \sqrt{\theta}$$

Thus the VC-dimension of  $G_n$  is at most the VC-dimension of the class of all ORs of two halfspaces. Applying Lemma 26 completes the proof. ■

Applying Lemmas 25 and 27, we obtain Lemma 7.

## Appendix B. Proof of Lemma 8

We will use the following, which strengthens bounds like Lemma 25 when the expectations being estimated are small. It differs from most bounds of this type by providing an especially strong bound on the probability that the estimates are *much* larger than the true expectations.

**Lemma 28 ((Bshouty et al., 2009))** *Suppose  $F$  is a set of  $\{0, 1\}$ -valued functions with a common domain  $X$ . Let  $d$  be the VC-dimension of  $F$ . Let  $\mathcal{D}$  be a probability distribution over  $X$ . Choose  $\alpha > 0$  and  $K \geq 4$ . Then if*

$$m \geq \frac{c \left( d \log \frac{1}{\alpha} + \log \frac{1}{\delta} \right)}{\alpha K \log K},$$

where  $c$  is an absolute constant, then

$$\Pr_{\mathbf{u} \sim \mathcal{D}^m} [\exists f \in F, \mathbf{E}_{\mathcal{D}}(f) \leq \alpha \text{ but } \hat{\mathbf{E}}_{\mathbf{u}}(f) > K\alpha] \leq \delta,$$

where  $\hat{\mathbf{E}}_{\mathbf{u}}(f) = \frac{1}{m} \sum_{i=1}^m f(u_i)$ .

To prove Lemma 8, we first use the fact that, for any fixed  $\mathbf{a} \in \mathbb{S}^{n-1}$  and  $\beta > 0$ , it is known (see (Kalai et al., 2008)) that

$$\Pr_{x \in \mathbb{S}^{n-1}} [|\mathbf{a} \cdot \mathbf{x}| > \beta] \leq e^{-\beta^2 n/2}.$$

Further, as in the proof of Lemma 7, we have that

$$|\mathbf{a} \cdot \mathbf{x}| > \beta \quad \text{if and only if} \quad \mathbf{a} \cdot \mathbf{x} > \beta \text{ OR } (-\mathbf{a}) \cdot \mathbf{x} > \beta,$$

so that the set of events whose probabilities we need to estimate is contained in the set of unions of pairs of halfspaces. Applying Lemma 26 and Lemma 28 completes the proof.