

Learning Intersections and Thresholds of Halfspaces

Adam R. Klivans*
Department of Mathematics
MIT

Cambridge, MA 02139
klivans@math.mit.edu

Ryan O'Donnell†
Department of Mathematics
MIT

Cambridge, MA 02139
odonnell@theory.lcs.mit.edu

Rocco A. Servedio‡
Division of Engineering and Applied Sciences
Harvard University
Cambridge, MA 02138
rocco@deas.harvard.edu

Abstract

We give the first polynomial time algorithm to learn any function of a constant number of halfspaces under the uniform distribution to within any constant error parameter. We also give the first quasipolynomial time algorithm for learning any function of a polylog number of polynomial-weight halfspaces under any distribution. As special cases of these results we obtain algorithms for learning intersections and thresholds of halfspaces. Our uniform distribution learning algorithms involve a novel non-geometric approach to learning halfspaces; we use Fourier techniques together with a careful analysis of the noise sensitivity of functions of halfspaces. Our algorithms for learning under any distribution use techniques from real approximation theory to construct low degree polynomial threshold functions.

1 Introduction

Let h be a hyperplane in \mathbb{R}^n , i.e. $h = \{x : \sum_{i=1}^n w_i x_i = \theta\}$. Such a hyperplane naturally induces a Boolean function $f(x) = \text{sgn}(\sum_{i=1}^n w_i x_i - \theta)$ which is called a *linear threshold function* or simply a *halfspace*. Learning an unknown halfspace from labeled data is one of the oldest problems in machine learning, dating back to the 1950s [7, 35]. This problem has been intensively studied over the years, and as described in Section 2.3 efficient algorithms are now known for several different learning models.

While the problem of learning a single halfspace is fairly well understood, learning more complicated functions which depend on several halfspaces seems to be quite difficult; in particular learning an intersection of several unknown halfspaces stands as a central open problem in computational learning theory. Intersections of halfspaces are attractive for many reasons: any convex body can be expressed as an intersection of halfspaces, and several well-studied classes of Boolean functions such as DNF formulas can be naturally viewed as special cases of intersections of halfspaces over the Boolean cube. Finally, we hope that learning an intersection (AND) of halfspaces will be a first step towards learning richer and more expressive functions of halfspaces.

1.1 Previous Work

Given the apparent difficulty of learning intersections of halfspaces from random examples, several researchers have considered learning algorithms which are

*Supported in part by NSF grant CCR-97-01304.

†Supported by NSF grant CCR-99-12342.

‡Supported by an NSF Mathematical Sciences Postdoctoral Research Fellowship and by NSF grant CCR-98-77049.

allowed to make *membership queries* for the value of the unknown function at points of the algorithm's choice. Building on work of Blum et. al. [8] and Baum [2], Kwek and Pitt [25] have given a membership query algorithm for learning the intersection of k halfspaces in \mathbb{R}^n (with respect to any probability distribution) in time polynomial in n and k .

Progress has been much more limited for learning intersections of halfspaces from random examples only; all such results to date require that the examples be drawn from some restricted class of probability distributions. Baum [3] gave a polynomial time algorithm for learning an intersection of two origin-centered halfspaces under any "symmetric" distribution (which satisfies $\mathcal{D}(x) = \mathcal{D}(-x)$ for all $x \in \mathbb{R}^n$). His algorithm is essentially a reduction to the problem of learning a single halfspace. Building on work of Blum and Kannan [10], Vempala [41] gave a polynomial time algorithm which can learn an intersection of $\log n / \log \log n$ halfspaces under "near-uniform" distributions on the Euclidean ball in \mathbb{R}^n . Vempala's algorithm uses random sampling to identify the subspace spanned by the normal vectors of the unknown halfspaces.

1.2 Our Results

We give new results for learning intersections of halfspaces, thresholds of halfspaces, and arbitrary functions of halfspaces over the Boolean cube. All of our algorithms learn from random examples only, and we obtain results for learning both from uniformly distributed examples and from examples drawn from an arbitrary probability distribution.

1.2.1 Uniform Distribution Learning

Our main learning result for the uniform distribution is a polynomial time algorithm for learning any function of any constant number of halfspaces to within any constant error parameter. More precisely, we prove:

Theorem 1 *Let $g : \{+1, -1\}^k \rightarrow \{+1, -1\}$ be any Boolean function on k bits and let h_1, \dots, h_k be arbitrary halfspaces on $\{+1, -1\}^n$. The class of all functions $\{g(h_1(x), \dots, h_k(x))\}$ can be learned under the uniform distribution to accuracy ϵ in time $n^{O(k^2/\epsilon^2)}$, assuming $\epsilon < 1/k^2$.*

For $k = O(1)$ and $\epsilon = \Omega(1)$ this time bound is polynomial in n . We note that prior to our work no polynomial time algorithm was known which could learn even an intersection of two arbitrary halfspaces under the uniform distribution on $\{+1, -1\}^n$.

We can substantially improve the dependence on k for the special case of learning a read-once intersection or majority of halfspaces:

Theorem 2 *Let h_1, \dots, h_k be arbitrary halfspaces on $\{+1, -1\}^n$ which depend on disjoint sets of variables. The class $\{h_1(x) \wedge h_2(x) \wedge \dots \wedge h_k(x)\}$ of read-once intersections of k halfspaces can be learned under the uniform distribution to accuracy ϵ in time $n^{O((\log(k)/\epsilon)^2)}$, assuming $\epsilon < 1/\log k$.*

Theorem 3 *Let h_1, \dots, h_k be arbitrary halfspaces on $\{+1, -1\}^n$ which depend on disjoint sets of variables. The class $\{\text{sgn}(h_1(x) + h_2(x) + \dots + h_k(x))\}$ of read-once majorities of k halfspaces can be learned under the uniform distribution to accuracy ϵ in time $n^{\tilde{O}((\log(k)/\epsilon)^4)}$, assuming $\epsilon < 1/\log k$.*

1.2.2 Learning under Arbitrary Distributions

Our algorithms for learning under arbitrary distributions use different techniques. In this scenario, our time bounds for learning depend chiefly on two parameters: the number of halfspaces k and the weight w of each halfspace (i.e. the magnitude of its integer coefficients). Our main learning result for arbitrary distributions is:

Theorem 4 *Let $g : \{+1, -1\}^k \rightarrow \{+1, -1\}$ be any Boolean function on k bits and let h_1, \dots, h_k be weight- w halfspaces on $\{+1, -1\}^n$. The class of all functions $\{g(h_1(x), \dots, h_k(x))\}$ can be learned to accuracy ϵ under any distribution in time $n^{O(k^2 \log w)/\epsilon}$.*

Thus we can learn any function of $\text{polylog}(n)$ many halfspaces of $\text{poly}(n)$ weight in quasipolynomial time under any probability distribution.

When we restrict attention to intersections of halfspaces, we can get better time bounds.

Theorem 5 *Let h_1, \dots, h_k be weight- w halfspaces on $\{+1, -1\}^n$. The class $\{h_1(x) \wedge h_2(x) \wedge \dots \wedge h_k(x)\}$ of intersections of k weight- w halfspaces can be learned to accuracy ϵ under any distribution in time $n^{O(k \log k \log w)/\epsilon}$.*

This theorem does well for the intersection of a fairly small number of halfspaces of reasonable weight. For the intersection of a large number of halfspaces of very small weight, we obtain a different bound:

Theorem 6 *Let h_1, \dots, h_k be weight- w halfspaces on $\{+1, -1\}^n$. The class $\{h_1(x) \wedge h_2(x) \wedge \dots \wedge h_k(x)\}$ of intersections of k weight- w halfspaces can be learned to accuracy ϵ under any distribution in time $n^{O(\sqrt{w} \log k)/\epsilon}$.*

Theorems 5 and 6 exhibit a tradeoff between the number of halfspaces k and the weight w of the halfspaces; as discussed in Section 4 this tradeoff is essentially optimal given our techniques.

We can generalize the bound of Theorem 5 to thresholds of halfspaces:

Theorem 7 *Let h_1, \dots, h_ℓ be weight- w halfspaces on $\{+1, -1\}^n$ and let g be a weight- k halfspace on $\{+1, -1\}^\ell$. The class $\{g(h_1(x), \dots, h_\ell(x))\}$ of weight- k thresholds of weight- w halfspaces can be learned to accuracy ϵ under any distribution in time $n^{O(k \log k \log w)} / \epsilon$.*

Finally, we extend recent results of Klivans and Servedio [24] on learning DNF formulas (ORs of ANDs) to thresholds of ANDs:

Theorem 8 *Let C_1, \dots, C_ℓ be arbitrary Boolean conjunctions over $\{+1, -1\}^n$ and let $g : \{+1, -1\}^\ell \rightarrow \{+1, -1\}$ be a weight- w halfspace. The class $\{g(C_1, \dots, C_\ell)\}$ of weight- w thresholds of ANDs can be learned to accuracy ϵ under any distribution in time $n^{O(n^{1/3} \log w)} / \epsilon$.*

We thus achieve the same running time bound from [24] for learning DNF formulas while learning a strictly more expressive class of functions.

1.3 Our Approach

The techniques we use for learning under the uniform distribution are quite different from those we use for learning under an arbitrary distribution. In the uniform distribution case, we show tight concentration bounds for the Fourier spectra of functions of halfspaces; this lets us learn using a Fourier based sampling algorithm from Linial *et al.* [26]. In the arbitrary distribution case, we show how to represent functions of halfspaces as low-degree polynomial threshold functions, which lets us learn using linear programming.

1.3.1 Uniform Distribution Learning: Fourier Analysis and Noise Sensitivity

The centerpiece of our uniform-distribution learning algorithms is a new Fourier concentration bound for functions of halfspaces. Recall that a Fourier concentration bound for a class of functions \mathcal{C} is a statement of the following form: For every function $f \in \mathcal{C}$ on n inputs, all but an ϵ fraction of the Fourier spectrum of f is concentrated on degree up to $\alpha(\epsilon, n)$. Given such a bound, the “low degree algorithm” of Linial *et al.* [26] provides a uniform-distribution learning algorithm for \mathcal{C} running in time $n^{O(\alpha(\epsilon, n))}$. The main result in [26] is a Fourier

concentration bound for the class of functions expressible by \mathcal{C}^0 circuits, with $\alpha(\epsilon, n) = \text{polylog}(n/\epsilon)$. Our new concentration bound is $\alpha(\epsilon, n) = O(k\sqrt{\epsilon})$ for the class of arbitrary functions of k halfspaces. We also give tighter bounds for more restricted classes of functions of halfspaces.

Our technique for proving these concentration bounds is to study the *noise sensitivity* of halfspaces. The noise sensitivity of a function f is simply the probability that $f(x)$ differs from $f(y)$ where x is a randomly chosen point and y is a slight perturbation of x . The noise sensitivity of Boolean functions was studied extensively by Benjamini *et al.* [6]; they showed that functions with low noise sensitivity have good Fourier concentration bounds. By analyzing the noise sensitivity of halfspaces rather than studying their Fourier properties directly, we are able to get Fourier concentration bounds for various classes of functions of halfspaces.

1.3.2 Learning under Arbitrary Distributions: Polynomial Threshold Functions

Our results for learning under an arbitrary distribution begin with the fact that a polynomial threshold function of degree d can be learned in time $n^{O(d)}$. Our contribution is showing that various classes of functions of halfspaces can be expressed as polynomial threshold functions of low degree. This technique has previously been used for learning by Klivans and Servedio [24], who showed that any polynomial-size DNF formula can be represented as a polynomial threshold function of degree $\tilde{O}(n^{1/3})$. They thus obtained a learning algorithm for DNF which works under any distribution and runs in time $2^{\tilde{O}(n^{1/3})}$.

We show that any function of k halfspaces where the sum of the (integer) coefficients of each halfspace is bounded by w can be represented as a polynomial threshold function of degree $O(k^2 \log w)$. We prove this using rational function approximation tools which were first used in a complexity theory context by Beigel *et al.* [4]. Roughly, we use rational functions which approximate the function sgn to closely approximate the ± 1 output values of our halfspaces. Having done this we obtain a single polynomial threshold function computing an arbitrary function g of halfspaces by composing these approximations with an interpolating polynomial for g .

In certain circumstances, we can trade off the dependences on k and w by using extremal polynomials in place of rational functions. By using Chebychev polynomials (previously used by Klivans and Servedio [24] and Nisan and Szegedy [32] in a somewhat similar context), we obtain a polynomial threshold function of degree $O(w^{1/2} \log k)$ computing the intersection of k half-

spaces where the sum of the coefficients of each linear threshold function is at most w .

2 Preliminaries

2.1 Definitions and learning model

Throughout the paper unless otherwise indicated we represent TRUE and FALSE by -1 and $+1$ respectively. A *linear threshold function* or *halfspace* on $\{+1, -1\}^n$ is a function $f : \{+1, -1\}^n \rightarrow \{+1, -1\}$, $f(x) = \text{sgn}(\sum_{i=1}^n w_i x_i - \theta)$ where $w_1, \dots, w_n, \theta \in \mathbb{R}$ and $\text{sgn}(z) = +1$ iff $z \geq 0$. We say that w_1, \dots, w_n, θ represent f ; since the domain is the discrete Boolean cube, every halfspace f has some representation in which $w_1, \dots, w_n, \theta \in \mathbb{Z}$. The *weight* of f is the smallest integer w for which there exist a representation $w_1, \dots, w_n, \theta \in \mathbb{Z}$ with $|\theta| + \sum_{i=1}^n |w_i| = w$. The *majority* function on inputs x_1, \dots, x_n is $\text{MAJ}(x) = \text{sgn}(\sum_{i=1}^n x_i)$.

We use Valiant's well-studied *Probably Approximately Correct* (PAC) model of learning Boolean functions from random examples [40]. In this model a *concept class* C is a collection $\cup_{n \geq 1} C_n$ of Boolean functions where each $c \in C_n$ is a function on n bits. In the PAC model a learning algorithm has access to an *example oracle* $EX(c, \mathcal{D})$ which, when queried, provides a labeled example $\langle x, c(x) \rangle$ where x is drawn from distribution \mathcal{D} over $\{+1, -1\}^n$ and $c \in C_n$ is the unknown target concept which the algorithm is trying to learn. Given Boolean functions c, h on $\{+1, -1\}^n$ we say that h is an ϵ -approximator for c under \mathcal{D} if $\Pr_{x \in \mathcal{D}}[h(x) = c(x)] \geq 1 - \epsilon$; the goal of a PAC learning algorithm is to generate an ϵ -approximator for the unknown target concept c . More precisely, an algorithm A is a *PAC learning algorithm for concept class* C if the following condition holds: for all $n \geq 1$, all $c \in C_n$, any distribution \mathcal{D} on $\{+1, -1\}^n$, and any $0 < \epsilon < \frac{1}{2}, 0 < \delta < 1$, if A is given ϵ and δ as input and has access to $EX(c, \mathcal{D})$, then with probability at least $1 - \delta$ algorithm A outputs an ϵ -approximator for c under \mathcal{D} . We say that A *PAC learns* C in time t if A runs for at most t time steps and outputs a hypothesis h which can be evaluated on any point $x \in \{+1, -1\}^n$ in time t .

If the above condition holds only for the uniform distribution \mathcal{U} on $\{+1, -1\}^n$, we say that A is a *uniform distribution PAC learning algorithm for* C . Uniform distribution PAC learnability of Boolean functions has been studied by many authors; see, e.g., [5, 9, 12, 14, 13, 18, 17, 19, 20, 21, 22, 23, 26, 29, 36, 37, 38, 42, 43].

All of our learning algorithms, both for the uniform distribution and for arbitrary distributions, have running time bounds with a $\log(1/\delta)$ dependence on δ , and hence we typically omit mention of this dependence.

2.2 Fourier Analysis and Uniform Distribution Learning

We briefly review some facts about Fourier analysis on the Boolean cube and its relation to uniform distribution learning. For a detailed treatment with proofs see [28].

Let the space $\{+1, -1\}^n$ be endowed with the uniform probability measure, and define an inner product on functions $f, g : \{+1, -1\}^n \rightarrow \mathbb{R}$ by $\langle f, g \rangle = \mathbb{E}[fg]$. For $S \subseteq [n]$ the parity function $\chi_S : \{+1, -1\}^n \rightarrow \{+1, -1\}$ is defined by $\chi_S(x) = \prod_{i \in S} x_i$. The set of all 2^n parity functions $\{\chi_S\}_{S \subseteq [n]}$ forms an orthonormal basis for the vector space of real-valued functions on $\{+1, -1\}^n$, and hence every real-valued function $f : \{+1, -1\}^n \rightarrow \mathbb{R}$ can be uniquely expressed as a linear combination $f(x) = \sum_{S \subseteq [n]} \hat{f}(S) \chi_S(x)$. The coefficients $\hat{f}(S)$ are called the *Fourier coefficients* of f ; collectively, they are its *Fourier spectrum*. Orthonormality implies that $\hat{f}(S) = \langle f, \chi_S \rangle$ and thus $\hat{f}(S)$ measures the correlation between f and χ_S . Orthonormality also implies *Parseval's identity*, which states that for every real-valued function f we have $E[f^2] = \sum_{S \subseteq [n]} \hat{f}(S)^2$. For Boolean functions we thus have $\sum_{S \subseteq [n]} \hat{f}(S)^2 = 1$.

Let $\alpha(\epsilon, n)$ be a function $\alpha : (0, \frac{1}{2}) \times \mathbb{N} \rightarrow \mathbb{N}$. We say that a concept class C has a *Fourier concentration bound* of $\alpha(\epsilon, n)$ if for all $n \geq 1$, all $0 < \epsilon < \frac{1}{2}$, and all $f \in C_n$ we have

$$\sum_{|S| \geq \alpha(\epsilon, n)} \hat{f}(S)^2 \leq \epsilon,$$

i.e. at most an ϵ fraction of the Fourier spectrum weight resides in coefficients of degree more than $\alpha(\epsilon, n)$.

An important connection between Fourier analysis and learning theory is the following result, implicit in [26] (see [28] for a nice exposition and proof):

Fact 9 *Let C be a concept class with a Fourier concentration bound of $\alpha(\epsilon, n)$. Then there is a uniform distribution PAC learning algorithm for C which runs in time $n^{O(\alpha(\epsilon, n))}$.*

2.3 Polynomial Threshold Functions and Learning under Arbitrary Distributions

A *polynomial threshold function* is a function $f : \{+1, -1\}^n \rightarrow \{+1, -1\}$, $f(x) = \text{sgn}(p(x))$ where $p(x_1, \dots, x_n)$ is a real-valued multivariate polynomial. Equivalently, we say that the polynomial p *sign-represents* f . The *degree* of a polynomial threshold function is the degree of the polynomial p . Note that a linear threshold function is a polynomial threshold function of degree one.

The following fact is well known (see [11]):

Fact 10 *The concept class of linear threshold functions over $\{+1, -1\}^n$ can be PAC learned under any distribution in time $\text{poly}(n)/\epsilon$.*

The algorithm of Fact 10 is based on polynomial time linear programming. We will use the following easy extension of Fact 10 (see e.g. [24]):

Fact 11 *Let C be a concept class over $\{+1, -1\}^n$ such that each $f \in C_n$ can be expressed as a polynomial threshold function of degree at most $d(n)$. Then C can be PAC learned under any distribution in time $n^{O(d(n))}/\epsilon$.*

Remark 12 *By results of Maass and Turan [27] the class of linear threshold functions is also known to be learnable in polynomial time in the model of exact learning from equivalence queries. An analogue of Fact 11 holds in this model as well, and in fact all of our distribution-independent PAC learning results hold in this model as well.*

3 Learning Functions of Halfspaces under Uniform Distributions

3.1 From Noise Sensitivity to Fourier Concentration

Let $f : \{+1, -1\}^n \rightarrow \{+1, -1\}$ be a Boolean function. We define the *noise sensitivity* of f at ϵ to be:

$$\text{NS}_\epsilon(f) = \Pr_{x, \text{noise}} [f(x) \neq f(N_\epsilon(x))]$$

(cf. [6]). Here x is uniformly chosen from $\{+1, -1\}^n$, and N_ϵ is the *noise operator* which flips each bit of its input independently with probability ϵ . The following formula from [33] relates the noise sensitivity of f to its Fourier spectrum; related formulas are given in [13, 6].

Proposition 13 *For any $f : \{+1, -1\}^n \rightarrow \{+1, -1\}$:*

$$\text{NS}_\epsilon(f) = \frac{1}{2} - \frac{1}{2} \sum_{S \subseteq [n]} (1 - 2\epsilon)^{|S|} \hat{f}(S)^2.$$

It follows that a strong upper bound on the noise sensitivity of f implies a good Fourier concentration bound for f (the simple proof is omitted):

Proposition 14 *For any $f : \{+1, -1\}^n \rightarrow \{+1, -1\}$, $0 < \gamma < 1/2$:*

$$\sum_{|S| \geq 1/\gamma} \hat{f}(S)^2 < 2.32 \text{NS}_\gamma(f).$$

Corollary 15 *Let $f : \{+1, -1\}^n \rightarrow \{+1, -1\}$ be any Boolean function, and let $\alpha : [0, 1/2] \rightarrow [0, 1]$ be an increasing function such that $\text{NS}_\epsilon(f) \leq \alpha(\epsilon)$. Then:*

$$\sum_{|S| \geq m} \hat{f}(S)^2 \leq \epsilon \quad \text{for} \quad m = \frac{1}{\alpha^{-1}(\epsilon/2.32)}.$$

In Section 3.2 we first give a detailed analysis of the noise sensitivity of a single linear threshold function. We then extend this analysis to arbitrary functions of several halfspaces and thus obtain Fourier concentration bounds and learning results using Corollary 15 and Fact 9. In Section 3.3 we give noise sensitivity bounds (and thus obtain Fourier concentration bounds and learning results) for read-once intersections and majorities of halfspaces.

3.2 Noise Sensitivity of Halfspaces and Functions of Halfspaces

3.2.1 Noise Sensitivity of a Halfspace

Benjamini *et al.* [6] were the first to analyze the noise sensitivity of linear threshold functions. They proved that every halfspace $h : \{+1, -1\}^n \rightarrow \{+1, -1\}$ has $\text{NS}_\epsilon(h) \leq C\epsilon^{1/4}$ where C is an absolute constant (note that this bound does not depend on n). Y. Peres subsequently improved this bound to $O(\sqrt{\epsilon})$ (unpublished, [34]). Since the majority function has noise sensitivity $\Omega(\sqrt{\epsilon})$ [6, 33], this is the best bound possible that depends only on ϵ .

We will use a refinement of Peres's bound that takes into account how "balanced" the halfspace is between outputs $+1$ and -1 . To motivate our bound, let h be a halfspace with $p = \Pr[h(x) = +1] < \frac{1}{2}$. If x is an input to h and $y = N_\epsilon(x)$, then both x and y are uniformly distributed, so by a union bound $\Pr[h(x) = +1 \text{ or } h(y) = +1] \leq 2p$ and hence $\text{NS}_\epsilon(h) \leq 2p$. Thus if $2p \ll \sqrt{\epsilon}$ it should be possible to get a much stronger bound upper bound on noise sensitivity which depends on p .

Our bound on the noise sensitivity of an arbitrary linear threshold function is as follows:

Theorem 16 *Let $h : \{+1, -1\}^n \rightarrow \{+1, -1\}$ be a halfspace, i.e. $h(x) = \text{sgn}(\sum_{i=1}^n w_i x_i - \theta)$ for some $w \in \mathbb{R}^n, \theta \in \mathbb{R}$. Let $p = \min\{\Pr[h = +1], \Pr[h = -1]\}$. Then:*

$$\text{NS}_\epsilon(h) \leq \min\{2p, 20.5 p \sqrt{\epsilon \ln(1/p)}\}.$$

Theorem 16 is proved in Appendix A. Since $p \sqrt{\ln(1/p)}$ is maximized when $p = 1/2$, and $20.5 \cdot \frac{1}{2} \sqrt{\ln 2} < 8.54$, we get:

Corollary 17 *For any halfspace $h : \{+1, -1\}^n \rightarrow \{+1, -1\}$, we have $\text{NS}_\epsilon(h) \leq 8.54 \sqrt{\epsilon}$.*

Since $(8.54 \cdot 2.32)^2 < 400$, Corollary 15 yields:

Theorem 18 Let $f : \{+1, -1\}^n \rightarrow \{+1, -1\}$ be any linear threshold function. Then for all $0 < \epsilon < 1/2$,

$$\sum_{|S| \geq (20/\epsilon)^2} \hat{f}(S)^2 \leq \epsilon.$$

This bound is essentially tight as shown by the following theorem, whose proof we omit:

Theorem 19 For any $0 < \epsilon < 1/2$, the following inequality holds for all sufficiently large odd n :

$$\sum_{|S| > \alpha} \widehat{\text{MAJ}}_n(S)^2 > \epsilon,$$

whenever $\alpha < \frac{8}{\pi^3 \epsilon^2} - 2$.

3.2.2 Bounds on Functions of k Halfspaces

Corollary 17 immediately implies an upper bound on the noise sensitivity of any function of k halfspaces:

Theorem 20 Let $f : \{+1, -1\}^n \rightarrow \{+1, -1\}$ be any function of k halfspaces. Then $\text{NS}_\epsilon(f) \leq O(k\sqrt{\epsilon})$.

Proof: Let h_1, \dots, h_k be arbitrary Boolean halfspaces, and let $g : \{+1, -1\}^k \rightarrow \{+1, -1\}$ be an arbitrary Boolean function. Write $f = g(h_1, \dots, h_k)$. By Corollary 17, $\Pr[h_i(x) \neq h_i(N_\epsilon(x))]$ is $O(\sqrt{\epsilon})$ for each $i = 1 \dots k$. By a union bound, the probability that at least one h_i 's value changes under noise is at most $O(k\sqrt{\epsilon})$. But f 's value changes only if at least one of the h_i 's values changes. Hence $\Pr[f(x) \neq f(N_\epsilon(x))] \leq O(k\sqrt{\epsilon})$. ■

By applying Corollary 15, we get a Fourier concentration bound of $O(k^2/\epsilon^2)$ for the class of functions of k halfspaces (assuming $\epsilon < 1/k^2$). We now get our main uniform distribution learning result, Theorem 1, using Fact 9.

3.3 Extensions: Learning Read-Once Functions of Halfspaces

Theorem 20 bounds $\text{NS}_\epsilon(f)$ by $O(k\sqrt{\epsilon})$ for f an arbitrary function of k halfspaces. Can stronger bounds be obtained if f is restricted to be a simple function (such as intersection) of k halfspaces? While we do not know the answer to this question, we can give substantially improved bounds if f is a read-once intersection or majority of halfspaces, i.e. each variable x_i occurs as input to at most one halfspace.

For read-once intersection of halfspaces we have the following noise sensitivity bound (whose proof we omit):

Theorem 21 Let $f : \{+1, -1\}^n \rightarrow \{+1, -1\}$, $f(x) = h_1(x) \wedge h_2(x) \wedge \dots \wedge h_k(x)$ where h_1, \dots, h_k are halfspaces on disjoint sets of variables. Then $\text{NS}_\epsilon(f) = O(\sqrt{\epsilon \log k})$.

Note that $k \leq n$ for any read-once intersection of halfspaces. By using Corollary 15 and Fact 9, we get Theorem 2.

In earlier work Hancock *et al.* [15] gave a polynomial time algorithm for learning a read-once intersection of majorities (i.e. each weight w_i of each halfspace is 1) under the uniform distribution. We emphasize that we are learning a much richer class of functions since our threshold functions can have arbitrary (even exponentially large) weights.

For the more expressive class of read-once majority of halfspaces we have the following noise sensitivity bound, whose proof we omit:

Theorem 22 Let $f(x) = \text{sgn}(h_1(x) + \dots + h_k(x) - \theta)$ where h_1, \dots, h_k are halfspaces on disjoint sets of variables. Then $\text{NS}_\epsilon(f) = \tilde{O}((\epsilon \log k)^{1/4})$.

Once again, we get Theorem 3 by applying Corollary 17 and Fact 9.

4 Learning Intersections of Halfspaces under Arbitrary Distributions

For our distribution independent results we use a set of completely different techniques. For our first set of results we use rational function approximations to the sgn function. For our second set of results, we use the extremal properties of the Chebyshev polynomials.

Throughout this section for convenience we let $+1$ denote TRUE and -1 denote FALSE.

4.1 Bounding PTF degree via rational function approximation

We use the rational function approximation to the sgn function introduced by Beigel *et al.* [4], building on work of Newman [31]. The motivation of Beigel *et al.* was the study of the complexity class PP. Later, Siu *et al.* [39] used the techniques for proving lower bounds against small-weight threshold circuits. The main technical theorem we need is:

Theorem 23 [4] For every $\ell, t \geq 1$ there exists a rational function P_t^ℓ in one variable with the following properties:

- For any $x \in [1, 2^t]$ the value $P_t^\ell(x) \in [1, 1 + 1/\ell]$.
- For any $x \in [-2^t, -1]$ the value $P_t^\ell(x) \in [-1 - 1/\ell, -1]$.

- $P_t^\ell(x)$ has degree $O(t \log \ell)$.

Here the degree of a rational function is the sum of the degrees of its numerator and denominator. Using this tool, we can get a polynomial threshold function for the intersection of k halfspaces:

Theorem 24 *Let h_1, \dots, h_k be linear threshold functions $\{+1, -1\}^n \rightarrow \{+1, -1\}$ each of weight at most w . The function $h_1(x) \wedge \dots \wedge h_k(x)$ can be expressed as a polynomial threshold function of degree $O(k \log k \log w)$.*

Proof: For $i = 1, \dots, k$ let $H_i(x) = \sum_{j=1}^n w_j^i x_j - \theta_i$ be a minimum weight representation of $h_i(x)$, so w_j^i, θ_j are integers and $h_i(x) = \text{sgn}(H_i(x))$. Without loss of generality we may suppose that $|H_i(x)| \geq 1$ for all $x \in \{+1, -1\}^n$, since this can be achieved by doubling each w_j^i and replacing θ_i by $2\theta_i - 1$.

Let $\ell = 2k$ and $t = \log w$ and consider the sum of rational functions

$$I = P_{\log w}^{2k}(H_1(x)) + \dots + P_{\log w}^{2k}(H_k(x)).$$

Now on input x , if $h_i(x) = 1$ (TRUE) for every i , then $I \geq k$. On the other hand, if for some i we have $h_i(x) = -1$ (FALSE), then $I \leq k - 1$. Thus $\text{sgn}(I - k)$ computes $h_1(x) \wedge \dots \wedge h_k(x)$ on all inputs. Now $I - k$ is the sum of k rational functions each of degree $O(\log k \log w)$. Hence if we bring these to a common denominator, we can write $I - k$ as a single rational function $Q(x)/R(x)$ of degree $O(k \log k \log w)$, and $Q(x)/R(x)$ sign-represents $h_1(x) \wedge \dots \wedge h_k(x)$. Hence so does $Q(x)R(x)$, which is a polynomial threshold function of degree $O(k \log k \log w)$. ■

Using this theorem and Fact 11, we get Theorem 5. We note that for $k = w^2$ the $O(k \log k \log w)$ bound of Theorem 24 is nearly optimal; Minsky and Papert have shown that the “one-in-a-box” function requires polynomial threshold function degree at least k , and the one-in-a-box function can be expressed as an AND of k halfspaces of weight k^2 (in fact the halfspaces are ORs of fanin k^2).

By a virtually identical proof we can generalize Theorem 24 as follows:

Theorem 25 *Let h_1, \dots, h_ℓ be linear threshold functions $\{+1, -1\}^n \rightarrow \{+1, -1\}$ each of weight at most w . Let $f : \{+1, -1\}^\ell \rightarrow \{+1, -1\}$ be a weight k threshold function. Then $f(h_1, \dots, h_\ell)$ can be expressed as a polynomial threshold function of degree $O(k \log k \log w)$.*

Again, combining this with Fact 11 yields Theorem 7. We note that the degree bound given by Theorem

25 is nearly optimal for $k = w$. Hajnal *et al.* [16] show how to compute the parity function on k variables as a weight- $O(k)$ threshold of weight- $O(k)$ halfspaces, but it is well known that any polynomial threshold function computing parity on k variables requires degree at least k (see [1, 30]).

These rational function techniques can also be used to give learning algorithms for arbitrary functions of halfspaces; we get Theorem 4 from the following, whose proof we omit:

Theorem 26 *Let h_1, \dots, h_k be linear threshold functions $\{+1, -1\}^n \rightarrow \{+1, -1\}$ each of weight at most w . Let $f : \{+1, -1\}^k \rightarrow \{+1, -1\}$ be any function. Then $f(h_1, \dots, h_k)$ can be expressed as a polynomial threshold function of degree $O(k^2 \log w)$.*

4.2 Bounding PTF Degree via Chebychev Polynomials

4.2.1 Learning Intersections of Small-Weight Halfspaces

The degree bounds obtained in the previous section are most interesting for cases when the number of underlying halfspaces is small (say a constant) and each halfspace is of weight $n^{O(1)}$. We can also give an alternate construction of polynomial threshold functions based on Chebychev polynomials and obtain improved degree bounds for cases involving a polynomial number of small weight (e.g. weight $O(1)$) halfspaces. The proof is omitted.

Theorem 27 *Let h_1, \dots, h_k be linear threshold functions $\{+1, -1\}^n \rightarrow \{+1, -1\}$ each of weight at most w . The function $h_1(x) \wedge \dots \wedge h_k(x)$ can be expressed as a polynomial threshold function of degree $O(\sqrt{w} \log k)$.*

4.2.2 Learning a Threshold of ANDs

In the last section we used Chebychev polynomials to obtain polynomial threshold function degree bounds for ANDs of thresholds. We can also use them to obtain bounds for thresholds of ANDs. The following theorem (whose proof we omit) generalizes the bound of Klivans and Servedio (their Theorem 1) for OR of ANDs [24]:

Theorem 28 *Let C_1, \dots, C_ℓ be Boolean conjunctions, each over at most k literals, and let g be a weight- w threshold over $\{+1, -1\}^\ell$. The function $g(C_1, \dots, C_\ell)$ can be expressed as a polynomial threshold function of degree $O(\sqrt{k} \log w)$.*

Using Theorem 28 in place of Klivans and Servedio’s Theorem 1, the arguments of [24] can be used to prove the following strengthened version of Klivans and Servedio’s main structural result for DNF formulas:

Theorem 29 Let $f : \{+1, -1\}^n \rightarrow \{+1, -1\}$ be a weight- w threshold of ANDs. Then f can be expressed as a polynomial threshold function of degree $O(n^{1/3} \log w)$.

Applying Fact 11 we obtain Theorem 8 which is a strict generalization of the main learning result of [24].

References

- [1] J. Aspnes, R. Beigel, M. Furst, and S. Rudich. The expressive power of voting polynomials. *Combinatorica*, 14(2):1–14, 1994.
- [2] E. Baum. Neural net algorithms that learn in polynomial time from examples and queries. *IEEE Trans. on Neural Networks*, 2:5–19, 1991.
- [3] E. Baum. A polynomial time algorithm that learns two hidden unit nets. *Neural Computation*, 2:510–522, 1991.
- [4] R. Beigel, N. Reingold, and D. Spielman. Pp is closed under intersection. *Journal of Computer and System Sciences*, 50(2):191–202, 1995.
- [5] M. Bellare. A technique for upper bounding the spectral norm with applications to learning. In *Proceedings of the Fifth Annual Workshop on Computational Learning Theory*, pages 62–70, 1992.
- [6] I. Benjamini, G. Kalai, and O. Schramm. Noise sensitivity of boolean functions and applications to percolation. *Inst. Hautes Études Sci. Publ. Math.*, 90:5–43, 1999.
- [7] H. Block. The perceptron: a model for brain functioning. *Reviews of Modern Physics*, 34:123–135, 1962.
- [8] A. Blum, P. Chalasani, S. A. Goldman, and D. K. Slonim. Learning with unreliable boundary queries. *Journal of Computer and System Sciences*, 56(2):209–222, 1998.
- [9] A. Blum, M. Furst, J. Jackson, M. Kearns, Y. Mansour, and S. Rudich. Weakly learning dnf and characterizing statistical query learning using fourier analysis. In *Proceedings of the Twenty-Sixth Annual Symposium on Theory of Computing*, pages 253–262, 1994.
- [10] A. Blum and R. Kannan. Learning an intersection of a constant number of halfspaces under a uniform distribution. *Journal of Computer and System Sciences*, 54(2):371–380, 1997.
- [11] A. Blumer, A. Ehrenfeucht, D. Haussler, and M. Warmuth. Learnability and the vapnik-chervonenkis dimension. *Journal of the ACM*, 36(4):929–965, 1989.
- [12] D. Boneh and R. Lipton. Amplification of weak learning over the uniform distribution. In *Proceedings of the Sixth Annual Workshop on Computational Learning Theory*, pages 347–351, 1993.
- [13] N. Bshouty, J. Jackson, and C. Tamon. More efficient pac learning of dnf with membership queries under the uniform distribution. In *Proceedings of the Twelfth Annual Conference on Computational Learning Theory*, pages 286–295, 1999.
- [14] N. Bshouty and C. Tamon. On the fourier spectrum of monotone functions. *Journal of the ACM*, 43(4):747–770, 1996.
- [15] M. Golea, M. Marchand, and T. Hancock. On learning μ -perceptron networks on the uniform distribution. *Neural Networks*, 9:67–82, 1994.
- [16] A. Hajnal, W. Maass, P. Pudlak, M. Szegedy, and G. Turan. Threshold circuits of bounded depth. *Journal of Computer and System Sciences*, 46:129–154, 1993.
- [17] T. Hancock. *The complexity of learning formulas and decision trees that have restricted reads*. PhD thesis, Harvard University, 1992.
- [18] T. Hancock and Y. Mansour. Learning monotone k - μ dnf formulas on product distributions. In *Proceedings of the Fourth Annual Conference on Computational Learning Theory*, pages 179–193, 1991.
- [19] J. Jackson. An efficient membership-query algorithm for learning dnf with respect to the uniform distribution. *Journal of Computer and System Sciences*, 55:414–440, 1997.
- [20] J. Jackson, A. Klivans, and R. Servedio. Learnability beyond ac^0 . In *Proceedings of the 34th ACM Symposium on Theory of Computing*, 2002.
- [21] R. Khardon. On using the fourier transform to learn disjoint dnf. *Information Processing Letters*, 49:219–222, 1994.
- [22] M. Kharitonov. Cryptographic hardness of distribution-specific learning. In *Proceedings of the Twenty-Fifth Annual Symposium on Theory of Computing*, pages 372–381, 1993.
- [23] A. Klivans and R. Servedio. Boosting and hard-core sets. In *Proceedings of the Fortieth Annual Symposium on Foundations of Computer Science*, pages 624–633, 1999.
- [24] A. Klivans and R. Servedio. Learning dnf in time $2^{\tilde{O}(n^{1/3})}$. In *Proceedings of the Thirty-Third Annual Symposium on Theory of Computing*, pages 258–265, 2001.
- [25] S. Kwek and L. Pitt. Pac learning intersections of halfspaces with membership queries. *Algorithmica*, 22(1/2):53–75, 1998.
- [26] N. Linial, Y. Mansour, and N. Nisan. Constant depth circuits, fourier transform and learnability. *Journal of the ACM*, 40(3):607–620, 1993.
- [27] W. Maass and G. Turan. *How fast can a threshold gate learn?*, pages 381–414. MIT Press, 1994.
- [28] Y. Mansour. *Learning Boolean functions via the Fourier transform*, pages 391–424. 1994.
- [29] Y. Mansour. An $o(n^{\log \log n})$ learning algorithm for dnf under the uniform distribution. *Journal of Computer and System Sciences*, 50:543–550, 1995.
- [30] M. Minsky and S. Papert. *Perceptrons: an introduction to computational geometry*. MIT Press, Cambridge, MA, 1968.
- [31] D. J. Newman. Rational approximation to $|x|$. *Michigan Mathematical Journal*, 11:11–14, 1964.
- [32] N. Nisan and M. Szegedy. On the degree of boolean functions as real polynomials. In *Proceedings of the Twenty-Fourth Annual Symposium on Theory of Computing*, pages 462–467, 1992.
- [33] R. O’Donnell. Hardness amplification within np. In *Proceedings of the 34th ACM Symposium on Theory of Computing*, 2002.

- [34] Y. Peres. personal communication, 2001.
- [35] F. Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological Review*, 65:386–407, 1958.
- [36] Y. Sakai and A. Maruoka. Learning monotone log-term dnf formulas under the uniform distribution. *Theory of Computing Systems*, 33:17–33, 2000.
- [37] R. Servedio. On pac learning using winnow, perceptron, and a perceptron-like algorithm. In *Proceedings of the Twelfth Annual Conference on Computational Learning Theory*, pages 296–307, 1999.
- [38] R. Servedio. On learning monotone dnf under product distributions. In *Proceedings of the Fourteenth Annual Conference on Computational Learning Theory*, pages 473–489, 2001.
- [39] K.-Y. Siu, V. Roychowdhury, and T. Kailath. Rational approximation techniques for analysis of neural networks. *IEEE Transactions on Information Theory*, 40(2):445–474, 1994.
- [40] L. Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984.
- [41] S. Vempala. A random sampling based algorithm for learning the intersection of halfspaces. In *Proceedings of the 38th Annual Symposium on Foundations of Computer Science*, pages 508–513, 1997.
- [42] K. Verbeurgt. Learning dnf under the uniform distribution in quasi-polynomial time. In *Proceedings of the Third Annual Workshop on Computational Learning Theory*, pages 314–326, 1990.
- [43] K. Verbeurgt. Learning sub-classes of monotone dnf on the uniform distribution. In *Proceedings of the Ninth Conference on Algorithmic Learning Theory*, pages 385–399, 1998.

A Proof of Theorem 16

As we argued before, $NS_\epsilon(h) \leq 2p$ always. Hence it suffices to show $NS_\epsilon(h) \leq 20.5p\sqrt{\epsilon \ln(1/p)}$. Without loss of generality we assume that $\sum_{i=1}^n w_i x_i \neq \theta$ for all $x \in \{+1, -1\}^n$ (if not we can slightly perturb the weights without changing the boolean function h .)

When the noise operator N_ϵ changes the sign of an input bit, we call this a *flip*. When all the flips taken together cause the value of h to change, we call this a *flop*. Let $P(k)$ be the probability of a flop, conditioned on exactly k flips occurring. We will upper bound $P(k)$ and then take an appropriate binomial average over k in the end.

So let us suppose that there are exactly $k > 0$ flips. Write $m = \lfloor n/k \rfloor$. Let $x \in \{+1, -1\}^n$ be chosen uniformly at random, and let π be a uniformly random permutation on $[n]$. Define $X_1 = \sum_{i=1}^k x_i w_{\pi(i)}$, $X_2 = \sum_{i=k+1}^{2k} x_i w_{\pi(i)}$, \dots , $X_m = \sum_{i=(m-1)k+1}^{mk} x_i w_{\pi(i)}$, and finally $Z = \sum_{i=mk+1}^n x_i w_{\pi(i)}$. Write $S = \sum_{j=1}^m X_j + Z$. Because of the random permutation π , we can “assume that the weights of X_1 were flipped.”

In other words, $(S, S - 2X_1)$ has exactly the same joint distribution as $(w \cdot y, w \cdot y')$, where $y \in \{+1, -1\}^n$ is uniform and y' is y with exactly k randomly selected bits flipped.

Put $S' = S - X_1 = \sum_{j=2}^m X_j + Z$, so the “sum before flipping” is $S' + X_1$ and the “sum after flipping” is $S' - X_1$. Hence a flop occurs iff $|S' - \theta| < |X_1|$. (Note that $|S' - \theta| = |X_1|$ is impossible, by our first assumption.)

Suppose that we condition on there being *no* flop; i.e., we condition on the event $|S' - \theta| > |X_1|$. Then since $\text{sgn}(X_1)$ is independent from both $|X_1|$ and S' , we have that $\Pr[\text{sgn}(X_1) = \text{sgn}(S' - \theta)] = \Pr[\text{sgn}(X_1) \neq \text{sgn}(S' - \theta)] = 1/2$. But, since we are conditioning on the event $|S' - \theta| > |X_1|$, we have that $\text{sgn}(S' - \theta)$ and $\text{sgn}(S - \theta)$ are always the same. Therefore we may conclude that under this conditioning, $\Pr[\text{sgn}(S - \theta) \neq \text{sgn}(X_1)] = 1/2$; i.e.,

$$\Pr[\text{sgn}(S - \theta) \neq \text{sgn}(X_1) \text{ \& no flop}] = \frac{1}{2} \Pr[\text{no flop}].$$

Yet the event $[\text{sgn}(S - \theta) \neq \text{sgn}(X_1) \text{ \& no flop}]$ is exactly the same event as $[\text{sgn}(S - \theta) \neq \text{sgn}(X_1)]$, since if there *is* a flop, then $\text{sgn}(S - \theta)$ must be the same as $\text{sgn}(X_1)$. Hence we have $\frac{1}{2} \Pr[\text{no flop}] = \Pr[\text{sgn}(S - \theta) \neq \text{sgn}(X_1)]$; i.e.,

$$\frac{1}{2}(1 - P(k)) = \Pr[\text{sgn}(S - \theta) \neq \text{sgn}(X_1)].$$

Now note that we could have derived this statement with X_2 in place of X_1 , or indeed any of X_2, \dots, X_m in place of X_1 , simply because once we apply the random permutation π , we could have picked any of these blocks to “be the flips.” So in fact,

$$\forall i = 1 \dots m, \quad \frac{1}{2}(1 - P(k)) = \Pr[\text{sgn}(S - \theta) \neq \text{sgn}(X_i)].$$

Write τ for the random variable $\text{sgn}(S - \theta)$, and σ_i for the random variable $\text{sgn}(X_i)$. Converting probabilities to expectations of indicator variables, we have:

$$\forall i = 1 \dots m, \quad \frac{1}{2}(1 - P(k)) = \mathbf{E}[\mathbf{1}_{\tau \neq \sigma_i}].$$

Summing these equations over all i yields:

$$\frac{m}{2}(1 - P(k)) = \mathbf{E}\left[\sum_{i=1}^m \mathbf{1}_{\tau \neq \sigma_i}\right] \quad (1)$$

Recall that this expectation is over the choices of x and π , which force the values of τ and σ_i .

The above arguments are due to Peres, and indeed from this point it is fairly easy to obtain an $O(\sqrt{\epsilon})$ upper bound. Some more work is required to obtain our desired bound which depends on p .

Suppose without loss of generality that $p = \Pr[h = -1]$, so $p = \Pr[\tau = -1]$. For $t = 1 \dots m$, define:

$$p_t = \Pr[\tau = -1 \mid \text{exactly } t \text{ of the } \sigma_i \text{'s are } +1]$$

Since the event $S < \theta$ is negatively correlated with the events $X_i > 0$, we conclude that:

$$1 \geq p_0 \geq p_1 \geq p_2 \geq \dots \geq p_m \geq 0 \quad (2)$$

But also note that $\sigma_1, \dots, \sigma_m$ are all independent and uniformly distributed in $\{+1, -1\}$. Therefore:

$$p = \sum_{t=0}^m \frac{\binom{m}{t}}{2^m} p_t \quad (3)$$

Continuing from (1), we have:

$$\begin{aligned} & (m/2)(1 - P(k)) \\ &= \mathbf{E}[\# \text{ of } \sigma_i \text{'s differing from } \tau] \\ &= \sum_{t=0}^m \frac{\binom{m}{t}}{2^m} \mathbf{E} \left[\sum_{i=1}^m \mathbf{1}_{\tau \neq \sigma_i} \mid \text{exactly } t \text{ } \sigma_i \text{'s are } +1 \right] \\ &= \sum_{t=0}^m \frac{\binom{m}{t}}{2^m} \mathbf{E}[t \mathbf{1}_{\tau=-1} + (m-t) \mathbf{1}_{\tau=+1} \\ & \quad \mid \text{exactly } t \text{ } \sigma_i \text{'s are } +1] \\ &= \sum_{t=0}^m \frac{\binom{m}{t}}{2^m} (t \Pr[\tau = -1 \mid \text{exactly } t \text{ } \sigma_i \text{'s are } +1] \\ & \quad + (m-t) \Pr[\tau = -1 \mid \text{exactly } t \text{ } \sigma_i \text{'s are } +1]) \\ &= \sum_{t=0}^m \frac{\binom{m}{t}}{2^m} (tp_t + (m-t)(1-p_t)). \end{aligned}$$

So $\frac{m}{2}(1 - P(k)) = \mathbf{E}_T[Tp_T + (m-T)(1-p_T)]$, where $T \sim \text{Binomial}(m, 1/2)$. Some arithmetic gives:

$$\begin{aligned} P(k) &= \frac{2}{m} \mathbf{E}_T[T] - 1 + 2\mathbf{E}_T[p_T] - \frac{4}{m} \mathbf{E}_T[Tp_T] \\ &= 2 \left(p - \frac{2}{m} \sum_{t=0}^m \frac{\binom{m}{t}}{2^m} tp_t \right) \quad (4) \end{aligned}$$

We will obtain an upper bound for $P(k)$ by maximizing (4) subject to (2) and (3). This is a linear programming problem. Hence the maximum occurs at a vertex, which in this case means the maximum occurs when, for an appropriate $1 \leq b \leq m/2$, we have $p_0 = p_1 = \dots = p_{b-1} = 1$, $p_{b+1} = p_{b+2} = \dots = p_m = 0$, and p_b is such that (3) is tight. (We have $b \leq m/2$ since $p \leq 1/2$.)

It can be shown (proof omitted) that when the p_t 's take on these values, we have the following upper bound:

$$(4) \leq \frac{(14.48) p \sqrt{\ln(1/p)}}{\sqrt{m}}. \quad (5)$$

Given this upper bound, it remains only to take the appropriate binomial average over k to obtain our final bound on $\text{NS}_\epsilon(h)$. Since $m = \lfloor n/k \rfloor$, we have

$$\begin{aligned} \Pr[\text{flop}] &= \mathbf{E}_k[P(k)] \\ &\leq (14.48) p \sqrt{\ln(1/p)} \mathbf{E}[\sqrt{1/\lfloor n/k \rfloor}] \\ &\leq (14.48) p \sqrt{\ln(1/p)} \sqrt{\mathbf{E}[1/\lfloor n/k \rfloor]} \\ &\leq (14.48) p \sqrt{\ln(1/p)} \sqrt{(2/n) \mathbf{E}[k]} \\ &= (14.48) \sqrt{2} p \sqrt{\epsilon \ln(1/p)} \\ &< 20.5 p \sqrt{\epsilon \ln(1/p)} \end{aligned}$$

and Theorem 16 is proved. ■