

New degree bounds for polynomial threshold functions*

Ryan O’Donnell[†]

School of Computer Science

Carnegie Mellon University

Pittsburgh, PA 15213 USA

`odonnell@theory.lcs.mit.edu`

Rocco A. Servedio[‡]

Department of Computer Science

Columbia University

New York, NY 10025 USA

`rocco@cs.columbia.edu`

Abstract

A real multivariate polynomial $p(x_1, \dots, x_n)$ is said to *sign-represent* a Boolean function $f : \{0, 1\}^n \rightarrow \{-1, 1\}$ if the sign of $p(x)$ equals $f(x)$ for all inputs $x \in \{0, 1\}^n$. We give new upper and lower bounds on the degree of polynomials which sign-represent Boolean functions. Our upper bounds for Boolean formulas yield the first known subexponential time learning algorithms for formulas of *superconstant* depth. Our lower bounds for constant-depth circuits and intersections of halfspaces are the first new degree lower bounds since 1968, improving results of Minsky and Papert. The lower bounds are proved *constructively*; we give explicit dual solutions to the necessary linear programs.

*A preliminary version of these results appeared as [24].

[†]This work was done while at the Department of Mathematics, MIT, Cambridge, MA, and while supported by NSF grant 99-12342.

[‡]Corresponding author. Supported by an NSF Mathematical Sciences Postdoctoral Research Fellowship and by NSF grant CCR-98-77049. This work was done while at the Division of Engineering and Applied Sciences, Harvard University, Cambridge, MA.

1 Introduction

Let f be a Boolean function $f : \{-1, 1\}^n \rightarrow \{-1, 1\}$ and let p be a degree d multilinear polynomial in n variables with real coefficients. If the sign of $p(x)$ equals $f(x)$ for every $x \in \{-1, 1\}^n$, then we say that f is computed by a *polynomial threshold function* of degree d ; equivalently we say that p *sign-represents* f .

Polynomial threshold functions are an interesting and natural representation for Boolean functions which have many applications in complexity theory and learning theory, see, e.g., [3, 5, 6, 4, 26, 17, 16]. Positive results showing that functions have low degree polynomial threshold functions can be used to obtain efficient learning algorithms via linear programming; see, e.g., [17, 16]. Negative results showing that a function requires threshold polynomials of large degree and/or large coefficients can be used to obtain oracles separating PP from smaller classes; see, e.g., [5, 29].

In this paper we give new upper and lower bounds on polynomial threshold function degree for several interesting and natural classes of functions which have been previously considered (but not resolved) in the literature. It seems likely that both the upper and lower bound techniques we use will prove useful for broader classes of functions.

1.1 Previous work The study of polynomial threshold functions began with Minsky and Papert in their 1968 book on perceptrons [21]. Minsky and Papert gave three lower bounds on the degree of polynomial threshold functions:

- Any polynomial threshold function which computes parity on n variables must have degree at least n . This result has since been reproved many times, see, e.g., [3, 7].
- Any polynomial threshold function which computes a particular linear-size CNF formula, the “one-in-a-box” function on n variables, must have degree $\Omega(n^{1/3})$. By Boolean duality this lower bound also holds for a corresponding DNF formula.
- Any polynomial threshold function which computes the AND of two majorities each on n variables must have degree $\omega(1)$.

Despite the fact that many researchers in learning theory and complexity theory have studied polynomial threshold functions, relatively little progress has been made on improving these lower bounds since 1968. In particular, Vereshchagin [29] has a lower bound for a promise-problem extension of one-in-a-box and Beigel [5] has a lower bound for a certain linear threshold function; however, both of these show degree lower bounds for polynomial threshold functions only under the added assumption that the polynomials have small integer coefficients. More progress has been made on upper bounds; Beigel, Reingold, and Spielman

[6] proved that there is a polynomial threshold function of degree $O(\log n)$ which computes the AND of two n -bit majorities. More recently, Klivans and Servedio [17] showed that any polynomial-size DNF formula (equivalently, CNF formula) has a polynomial threshold function of degree $O(n^{1/3} \log n)$, and Klivans *et al.* [16] showed that any Boolean function of a polylogarithmic number of halfspaces with quasipolynomially-bounded weights has a polynomial threshold function of polylogarithmic degree.

We briefly note that researchers have also studied the *sparseness* of polynomial threshold functions for various types of Boolean functions, where the sparseness is simply the number of nonzero coefficients in the polynomial; see e.g. [7, 13, 18]. Sparseness bounds for polynomial threshold functions depend heavily on whether the polynomials in question are over $\{0, 1\}^n$ versus $\{-1, 1\}^n$, whereas this choice does not affect degree bounds. Sparseness bounds are in general incomparable to degree bounds, though of course a polynomial with very many nonzero coefficients over n variables cannot have too low degree. Krause and Pudlak [18] have given lower bounds on the number of nonzero coefficients which must be present in any polynomial threshold function for a particular depth-3 Boolean circuit, but their lower bounds are not strong enough to imply new lower bounds on polynomial threshold function degree.

1.2 Our results We give new upper and lower bounds on polynomial threshold functions for several interesting and natural classes of functions. Our main results are:

- We give an $\Omega(\frac{\log n}{\log \log n})$ lower bound on the degree of any polynomial threshold function which computes the AND of two n -bit majorities. Equivalently, this lower bound holds for the degree of any bivariate real polynomial $p(x, y)$ which is positive on the integer lattice points in the upper-right quadrant with coordinates bounded by n , and is negative on the lattice points in the other three quadrants with coordinates bounded in magnitude by n . This result (and our next) is the first new unconditional lower bound for polynomial threshold degree since 1968; it improves on Minsky and Papert’s lower bound of $\omega(1)$ and nearly matches the $O(\log n)$ upper bound of Beigel, Reingold and Spielman.
- We prove an “XOR lemma” for polynomial threshold function degree and use this lemma to obtain an $\Omega(n^{1/3} \log^{2d/3} n)$ lower bound on the degree of an explicit Boolean circuit of polynomial size and depth $d + 2$. This is the first improvement on Minsky and Papert’s $\Omega(n^{1/3})$ lower bound for any constant-depth circuit.
- We prove that any Boolean formula of depth d and size s is computed by a polynomial threshold function of degree $\sqrt{s}(\log s)^{O(d)}$. This gives us the first known upper bound

for Boolean formulas of superconstant depth. In particular, any Boolean formula of size $o(n^2)$ and depth $o(\frac{\log n}{\log \log n})$ has a polynomial threshold function of nontrivial (sublinear) degree. We use our upper bound to provide the first known subexponential learning algorithm for such formulas. Note that since parity on \sqrt{s} variables can be computed by a formula of size s , the best possible degree upper bound which depends only on s is \sqrt{s} .

We note that since the initial conference publication of these results [24], Ambainis *et al.* have shown that in fact any Boolean formula of size s is approximately computed by a polynomial of degree $O(\sqrt{s})$ [1], with pointwise error at most $1/3$ on each input. This immediately implies the existence of a polynomial threshold function of degree $O(\sqrt{s})$, thus improving on our degree bound.

1.3 Our techniques Perhaps surprisingly, our lower bounds are achieved constructively. The question of whether a given function has a polynomial threshold function of degree d can be formulated as the feasibility question for a certain linear program. By duality, we can show the linear program is infeasible — and hence the function has polynomial threshold degree exceeding d — by showing that the dual linear program is feasible. We construct explicit dual solutions. (Interestingly, Vereschagin’s lower bound [29] involves showing that a certain linear program *is* feasible by explicitly demonstrating the infeasibility of the dual.)

Our upper bounds build on ideas from [17, 16] and use tools from real approximation theory.

1.4 Organization Section 2 gives preliminaries on polynomial threshold functions and describes the duality technique we use for our lower bounds. In Section 3 we prove our XOR lemma for polynomial threshold functions using the duality technique, and use this lemma to obtain new lower bounds for constant depth circuits. In Section 4 we apply the lower bound technique to prove our $\Omega(\frac{\log n}{\log \log n})$ lower bound for the AND of two majorities. In Section 5 we give our upper bounds for Boolean formulas and the application to learning.

2 Preliminaries

2.1 Sign-representations of Boolean functions A *multilinear* monomial over the variables x_1, \dots, x_n is one in which each variable has degree at most one. Such a monomial is defined by the set $S \subseteq [n]$ of variables x_1, \dots, x_n that it contains; we write x_S to denote the monomial $\prod_{i \in S} x_i$. A multilinear polynomial (with coefficients in the reals) is a sum of the

form

$$p(x_1, \dots, x_n) = \sum_{S \subseteq [n]} p_S x_S$$

where each $p_S \in \mathbf{R}$. Since we will always be dealing with functions whose domain is $\{-1, 1\}^n$ or $\{0, 1\}^n$, we may consider only multilinear polynomials with no loss of generality in our results.

We remark that viewing the polynomial p as a real-valued function on $\{-1, 1\}^n$, the coefficients $\{p_S\}_{S \subseteq [n]}$ correspond to the Fourier spectrum of the function p with respect to the standard orthonormal basis formed by the collection of all 2^n monomials $\{x_S\}_{S \subseteq [n]}$. Following [26], we refer to the set $\mathcal{S} \subseteq 2^{[n]}$ of monomials on which p has nonzero coefficients as the *spectral support* of p .

We make the following standard definitions of sign-representing polynomials for a Boolean function (see [3]).

Definition 2.1 *Let $f : \{-1, 1\}^n \rightarrow \{-1, 1\}$ be a Boolean function. Let $p : \{-1, 1\}^n \rightarrow \mathbf{R}$ be a multilinear polynomial which is not identically 0.*

We say that p weakly sign-represents f if $f(x) = \text{sgn}(p(x))$ for all x such that $p(x) \neq 0$. We say that p strongly sign-represents f , or simply sign-represents f , if $f(x) = \text{sgn}(p(x))$ and $p(x) \neq 0$ for every $x \in \{-1, 1\}^n$. We write $\text{thr}(f)$ to denote the minimum degree over all polynomials strongly sign-representing f , and $\text{thr}^w(f)$ to denote the minimum degree over all polynomials weakly sign-representing f .

On occasion we will view the domain of f as $\{0, 1\}^n$ instead of $\{-1, 1\}^n$; it is easy to see that this does not change the degree of any sign-representing polynomial.

2.2 Distributions and their connection with sign-representations We will require the following notion of a *distribution* over $\{-1, 1\}^n$.

Definition 2.2 *We say that a distribution over $\{-1, 1\}^n$ is a map $w : \{-1, 1\}^n \rightarrow \mathbf{R}^{\geq 0}$ which is not identically 0. The set of points $\{x : w(x) \neq 0\}$ is called the support of the distribution w . If the support of w is all of $\{-1, 1\}^n$ then we say that w is a total distribution. If $\sum_{x \in \{-1, 1\}^n} w(x) = 1$ then we say that w is a probability distribution.*

Given a monomial x_S , $S \subseteq [n]$, we say that the correlation of x_S with f under distribution w is

$$\mathbf{E}_w[f(x)x_S] := \sum_{x \in \{-1, 1\}^n} f(x)x_S w(x).$$

Fix a Boolean function $f : \{-1, 1\}^n \rightarrow \{-1, 1\}$. As we now explain, there is an exact correspondence between sign-representing polynomials for f and distributions over $\{-1, 1\}^n$.

Let A_f be a $2^n \times 2^n$ matrix with ± 1 entries as follows. We view the rows of A_f as being indexed by inputs $x \in \{-1, 1\}^n$ and the columns of A_f as indexed by subsets of variables $S \subseteq [n]$. The entry $A_f[x, S]$ in row x and column S of the matrix is equal to $f(x)x_S$.

Note that for any $S_1, S_2 \subseteq [n]$ we have that the inner product of the S_1 and S_2 columns of A_f is (writing $S_1 \Delta S_2$ to denote the symmetric difference of S_1 and S_2)

$$\sum_{x \in \{-1, 1\}^n} f(x)x_{S_1}f(x)x_{S_2} = \sum_{x \in \{-1, 1\}^n} x_{S_1 \Delta S_2} = \begin{cases} 0 & \text{if } S_1 \neq S_2 \\ 2^n & \text{if } S_1 = S_2 \end{cases}$$

so consequently A_f is an orthogonal matrix (in fact it is a Hadamard matrix) and hence is a bijective mapping from \mathbf{R}^{2^n} to \mathbf{R}^{2^n} .

Now let p be a 2^n -dimensional column vector p whose entries are indexed by subsets $S \subseteq [n]$. Then $A_f p$ is a 2^n -dimensional column vector, which we call w , whose entries are indexed by $x \in \{-1, 1\}^n$. The x -th entry of w is

$$w(x) = \sum_{S \subseteq [n]} A_f[x, S]p_S = \sum_{S \subseteq [n]} f(x)x_S p_S = f(x) \sum_{S \subseteq [n]} p_S x_S = f(x)p(x).$$

So the vector p corresponds to a weak (strong) sign-representation of f if and only if the vector $w = A_f p$ corresponds to a distribution (total distribution) over $\{-1, 1\}^n$. We thus have

Proposition 2.3 *The mapping A_f is a bijection between weak sign-representations of f and distributions, and moreover is a bijection between strong sign-representations of f and total distributions.*

If $p(x)$ is a weak sign-representation of f and $w = A_f p$ is the corresponding distribution, we have that the S coefficient of p is proportional to the correlation of x_S with f under w . To see this, recall that this correlation is

$$\begin{aligned} \sum_{x \in \{-1, 1\}^n} f(x)x_S w(x) &= \sum_{x \in \{-1, 1\}^n} f(x)x_S \left(f(x) \sum_{T \subseteq [n]} p_T x_T \right) \\ &= \sum_x x_S \sum_T p_T x_T \\ &= \sum_T p_T \sum_x x_{S \Delta T} = 2^n p_S \end{aligned}$$

where the final equality holds because the inner sum is nonzero only if $T = S$, in which case it is 2^n . We thus have:

Proposition 2.4 *Let p be a weak sign-representation of f and $w = A_f p$ the corresponding distribution. Then the spectral support of p is \mathcal{S} if and only if*

- *f has zero correlation with x_S under w for every monomial $S \notin \mathcal{S}$, and*
- *f has non-zero correlation with x_S under w for every monomial $S \in \mathcal{S}$.*

2.3 The Theorem of the Alternative Our main tool for proving polynomial threshold degree lower bounds is the following so-called “Theorem of the Alternative.” It can be proved immediately using linear programming duality, as was essentially done by Aspnes *et al.* in [3]; see also [26] for a nice exposition of the simple proof. A completely different proof based on the distribution perspective can be given by combining the “Discriminator Lemma” of [14] with the learning-theoretic technique of boosting, see [12, 13].

Theorem 2.5 *Let $f : \{-1, 1\}^n \rightarrow \{-1, 1\}$ be a Boolean function. Let $\mathcal{S} \subseteq 2^{[n]}$ be any set of monomials. Then exactly one of the following holds:*

- *f has a strong representation with spectral support in \mathcal{S} ; or,*
- *f has a weak representation with spectral support in $2^{[n]} \setminus \mathcal{S}$.*

Given the equivalence between sign-representations and distributions from the previous subsection, there are three other ways of restating Theorem 2.5. We will need one more:

Theorem 2.6 *Let $f : \{-1, 1\}^n \rightarrow \{-1, 1\}$ be a Boolean function. Let $\mathcal{S} \subseteq 2^{[n]}$ be any set of monomials. Then exactly one of the following holds:*

- *f has a strong representation with spectral support in \mathcal{S} ; or,*
- *there is a distribution on $\{-1, 1\}^n$ under which f has zero correlation to every monomial in \mathcal{S} .*

3 An XOR lemma for PTF degree

Let f be any Boolean function $\{-1, 1\}^n \rightarrow \{-1, 1\}$ defined on variables x_1, \dots, x_n and let g be any Boolean function $\{-1, 1\}^n \rightarrow \{-1, 1\}$ defined on variables y_1, \dots, y_n . Let $f \oplus g$ denote the XOR (parity) of f and g . We will prove the following “XOR lemma:”

Theorem 3.1 *Let f and g be Boolean functions on disjoint sets of variables. Then $\text{thr}(f \oplus g) = \text{thr}(f) + \text{thr}(g)$.*

We note that Theorem 3.1 is similar in spirit (though incomparable) to a recent result of Sieling [27] which shows that $DT(f \oplus g) = DT(f) \cdot DT(g)$, where $DT(f)$ is the minimum decision tree size of f .

Proof of Theorem 3.1: The upper bound is easy; if $p_f(x)$ is a strong sign-representation of f of degree $\text{thr}(f)$ and $p_g(y)$ is a strong sign-representation of g with degree $\text{thr}(g)$ then $p_f(x)p_g(y)$ is easily seen to be a strong sign-representation of $f \oplus g$, and $\deg(p_f(x)p_g(y)) = \text{thr}(f) + \text{thr}(g)$.

For the lower bound, since f has no strong representation on the set of monomials of degree strictly less than $\text{thr}(f)$, Theorem 2.5 tells us that f has a weak representation $q_f(x)$ supported on the monomials x_S with $|S| \geq \text{thr}(f)$. Similarly, g has a weak representation $q_g(y)$ supported on the monomials y_T with $|T| \geq \text{thr}(g)$. Now $q_f(x)q_g(y)$ is a weak representation of $f \oplus g$; in particular, it is not identically zero because there is at least one x for which $q_f(x) \neq 0$ and at least one y for which $q_g(y) \neq 0$, so $q_f(x)q_g(y) \neq 0$ for these inputs. Note that $q_f(x)q_g(y)$ is supported on the set of monomials which have degree at least $\text{thr}(f)$ in x and at least $\text{thr}(g)$ in y . Applying Theorem 2.5 again we conclude that any strong representation for $f \oplus g$ must use some monomial with degree at least $\text{thr}(f)$ in x and at least $\text{thr}(g)$ in y ; this is more than sufficient to prove that $\text{thr}(f \oplus g) \geq \text{thr}(f) + \text{thr}(g)$. (Theorem 3.1) ■

For f a Boolean function let $\oplus_k f$ denote the XOR of k copies of f on disjoint sets of variables. From Theorem 3.1 we obtain:

Corollary 3.2 $\text{thr}(\oplus_k f) = k \cdot \text{thr}(f)$.

This corollary thus includes Minsky and Papert's lower bound of n for the parity function as a special case.

Corollary 3.2 also yields the following lower bound for constant depth circuits:

Theorem 3.3 *For all $d \geq 1$ there is an AND/OR/NOT circuit C of depth $d + 2$ and size $\text{poly}(n)$ which has polynomial threshold function degree $\Omega(n^{1/3}(\log n)^{2d/3})$.*

Proof: The circuit C computes the parity of $(\log n)^d$ disjoint copies of Minsky and Papert's "one-in-a-box" function, where each one-in-a-box function is defined on $n/(\log n)^d$ variables. It is well known that for any constant d , parity on $(\log n)^d$ variables can be computed by an AND/OR/NOT circuit of depth $d + 1$ and size $\text{poly}(n)$. Since the one-in-a-box function on $n/(\log n)^d$ variables is a depth-2 circuit of size $O(n/(\log n)^d)$, by substituting the appropriate one-in-a-box function for each input to the parity we see that C is a circuit of $\text{poly}(n)$ size

and depth $d + 2$ (we save one on depth by collapsing gates of the same kind on the next to bottom layer). By Minsky and Papert’s lower bound, we know that any polynomial threshold function for one-in-a-box on $n/(\log n)^d$ variables must have degree $\Omega((n/(\log n)^d)^{1/3})$. Consequently Corollary 3.2 implies that $\text{thr}(C) = \Omega(n^{1/3}(\log n)^{2d/3})$ and the theorem is proved. ■

In fact, we can actually give an alternate proof of Minsky and Papert’s lower bound for one-in-a-box by using our lower bound technique of applying the Theorem of the Alternative (Theorem 2.6) and constructing the necessary distribution explicitly. See Appendix A.

Theorem 3.3 is of interest since it gives the first $\omega(n^{1/3})$ lower bound for any function in AC^0 . We note that Theorem 3.3 also shows that the $n^{1/3} \log n$ upper bound of Klivans and Servedio for depth-2 AC^0 circuits does not hold for depth-4 AC^0 .

4 A lower bound for the AND of two majorities

Let n be odd, and let $\text{AND-MAJ}_n : \{-1, 1\}^n \times \{-1, 1\}^n \rightarrow \{-1, 1\}$ be the function which on input (x, y) , $x, y \in \{-1, 1\}^n$, outputs 1 if both $\text{MAJ}_n(x) = 1$ and $\text{MAJ}_n(y) = 1$. Here MAJ_n is the majority function on n bits, $x \mapsto \text{sgn}(\sum_{i=1}^n x_i)$. In this section we show that $\text{thr}(\text{AND-MAJ}_n) = \Omega(\frac{\log n}{\log \log n})$, improving on the $\omega(1)$ lower bound of Minsky and Papert. Note that $O(\log n)$ is an upper bound, by Beigel, Reingold, and Spielman [6].

The high-level idea of the proof is to use the Theorem of the Alternative. More precisely, we will show that there is a distribution on $\{-1, 1\}^n$ under which AND-MAJ_n has zero correlation with every “low-degree” monomial. Given this, Theorem 2.6 implies that f has no strong representation with spectral support in the set of “low-degree” monomials, so consequently the threshold degree of AND-MAJ_n must be high.

We begin by applying a simple symmetrization due to Minsky and Papert. Suppose p is a polynomial threshold function for AND-MAJ_n where n is odd. Let $\mathbf{Z}_n^{\text{odd}}$ denote the set $\{-n, -(n-2), \dots, -1, 1, \dots, n-2, n\} \subseteq \mathbf{Z}$. Let $\text{AND-sgn}_n : \mathbf{Z}_n^{\text{odd}} \times \mathbf{Z}_n^{\text{odd}} \rightarrow \{-1, 1\}$ be the function which on input (x, y) is 1 iff $x > 0$ and $y > 0$.

Claim 4.1 *There exists a polynomial threshold function for AND-MAJ_n of degree d if and only if there exists a bivariate polynomial of degree d which sign-represents AND-sgn_n .*

Proof: (if) Suppose g is a bivariate polynomial sign-representing AND-sgn_n . Let $p : \{-1, 1\}^n \times \{-1, 1\}^n \rightarrow \mathbf{R}$ be given by $p(x, y) = g(\sum x_i, \sum y_i)$. Then the multilinear reduction of p has degree d and sign-represents AND-MAJ_n .

(only if) Suppose p has degree d and sign-represents AND-MAJ_n . Let $q : \{-1, 1\}^n \times \{-1, 1\}^n \rightarrow \mathbf{R}$ be given by $(x, y) \mapsto \sum_{\pi_1, \pi_2 \in \mathbf{S}_n} p(x_{\pi_1(1)}, \dots, x_{\pi_1(n)}, y_{\pi_2(1)}, \dots, y_{\pi_2(n)})$, where \mathbf{S}_n is the symmetric group on $[n]$. Because $\text{AND-MAJ}_n(x_{\pi_1(1)}, \dots, x_{\pi_1(n)}, y_{\pi_2(1)}, \dots, y_{\pi_2(n)}) =$

AND-MAJ $_n(x, y)$ we conclude that q sign-represents AND-MAJ $_n$ and has degree d . But notice that q is symmetric in its x variables and its y variables. Hence there is a degree d bivariate polynomial $\tilde{q}(\cdot, \cdot)$ such that $\tilde{q}(\sum x_i, \sum y_i) = q(x, y)$ for all $(x, y) \in \{-1, 1\}^n \times \{-1, 1\}^n$, and thus \tilde{q} sign-represents AND-sgn $_n$. \blacksquare

It follows that if we prove a lower bound on the degree of a bivariate polynomial which sign-represents AND-sgn $_n$, we get a lower bound on $\text{thr}(\text{AND-MAJ}_n)$. Following Theorem 2.6, we shall show that there is a probability distribution over $\mathbf{Z}_n^{\text{odd}} \times \mathbf{Z}_n^{\text{odd}}$ under which every bivariate monomial of degree at most $d = \Omega(\frac{\log n}{\log \log n})$ has zero correlation with AND-sgn $_n$. To see that this is enough, suppose that \tilde{q} is a bivariate polynomial of degree d sign-representing AND-sgn $_n$ and w is a probability distribution over $\mathbf{Z}_n^{\text{odd}} \times \mathbf{Z}_n^{\text{odd}}$ with the stated property. Then on one hand $\mathbf{E}_w[\text{AND-sgn}_n(x, y)\tilde{q}(x, y)] = 0$ by linearity of expectation, since each monomial in \tilde{q} has zero correlation with AND-sgn $_n$ under w . But on the other hand, since \tilde{q} strongly sign-represents AND-sgn $_n$, $\text{AND-sgn}_n(x, y)\tilde{q}(x, y) > 0$ for all (x, y) , hence $\mathbf{E}_w[\text{AND-sgn}_n(x, y)\tilde{q}(x, y)] > 0$, contradiction.

The problem is now set up to our satisfaction. Fix an integer d . We shall try to find a support (set of points) $\mathcal{Z} \subset \mathbf{Z}^{\text{odd}} \times \mathbf{Z}^{\text{odd}}$ and a probability distribution \mathbf{w} over these points such that the function $f = \text{AND-sgn}_n$ has zero correlation under \mathbf{w} with every monomial $x^i y^j$ of total degree at most d . That is, we want $\mathbf{w} : \mathcal{Z} \rightarrow \mathbf{R}^{\geq 0}$ with $\sum_{z \in \mathcal{Z}} \mathbf{w}(z) = 1$ such that:

$$\forall 0 \leq i + j \leq d, \quad \mathbf{E}_{\mathbf{w}}[f(x, y) x^i y^j] = \sum_{(x, y) \in \mathcal{Z}} \mathbf{w}(x, y) f(x, y) x^i y^j = 0. \quad (1)$$

In addition we would like to find a solution in which $\text{size}(\mathcal{Z})$ is as small as possible, where $\text{size}(\mathcal{Z})$ denotes $\max_{(x, y) \in \mathcal{Z}} \{\max\{|x|, |y|\}\}$. Once we have such a \mathcal{Z} and \mathbf{w} , we get a lower bound of $d + 1$ for the degree of a polynomial threshold function computing AND-MAJ $_{\text{size}(\mathcal{Z})}$. In the remainder of this section we give a construction in which $\text{size}(\mathcal{Z}) = d^{O(d)}$. Since we can take $\text{size}(\mathcal{Z})$ as large as $\Theta(n)$, this means we may take $d = \Omega(\frac{\log n}{\log \log n})$, and we obtain the main result of this section:

Theorem 4.2 $\text{thr}(\text{AND-MAJ}_n) = \Omega(\frac{\log n}{\log \log n})$.

4.1 Proof of Theorem 4.2

Suppose we fix some n and wish to know if there exists a distribution \mathbf{w} supported on $\mathbf{Z}_n^{\text{odd}} \times \mathbf{Z}_n^{\text{odd}}$ satisfying (1). If we view the values $\{\mathbf{w}(x, y) : (x, y) \in \mathbf{Z}_n^{\text{odd}} \times \mathbf{Z}_n^{\text{odd}}\}$ as unknowns, this is precisely asking if a certain system of linear equations has a nonnegative solution; i.e., it is a feasibility problem for a linear program.

Thus let us say that our *constraints* are all bivariate monomials $x^i y^j$ of total degree at most d . We will refer to $x^i y^j$ as the “ (i, j) constraint monomial.” There are a total of $D = \frac{(d+1)(d+2)}{2}$ constraint monomials, and for definiteness we will consider them to be ordered as follows: $1, x, y, x^2, xy, y^2, x^3$, etc.

By basic linear programming theory, if there is a feasible solution for the $\mathbf{w}(x, y)$'s then there is one in which at most $D + 1$ of the values are nonzero. The key to our proof will be to guess an acceptable “support” set \mathcal{Z} of cardinality $D + 1$ with $\text{size}(\mathcal{Z})$ small. We will then show that the unique solution to the system of equations given by (1) and $\sum \mathbf{w}(x, y) = 1$ is indeed nonnegative.

To guess an acceptable support set, we in fact explicitly checked feasibility of the LP for small values of d . For $d = 1$ the minimum possible value of n yielding feasibility was $n = 3$; for $d = 2$ the minimum feasible n was $n = 23$. The explicit solutions were supported on points that seemed to roughly be of the form $(\pm h^i, \pm h^j)$ for a small constant h and $0 \leq i, j \leq d$.¹ By explicitly considering supports only of this form we were able to show feasibility bounds with $n \leq 10^d$ for $d = 3, 4, 5$. Further, by studying the precise solutions found by the LP solver, we were led to consider the following support set for the general case, which we will use in the remainder of the proof:

$$\mathcal{Z} = \{((-1)^\ell h^k, (-1)^k h^\ell) : 0 \leq k + \ell \leq d\} \cup \{(-1, -1)\},$$

where here h is a large quantity to be chosen later. We believe h can be taken constant, but for our proof we will eventually take $h = \Theta(d^9)$.

This support \mathcal{Z} is symmetric about the line $y = x$ and contains exactly $D + 1$ points. We will refer to $((-1)^\ell h^k, (-1)^k h^\ell)$ as the “ (k, ℓ) support point” and consider the points to be ordered in the same order as the monomials (i.e., $(1, 1), (h, -1), (-1, h), (h^2, 1), (-h, -h), (1, h^2), (h^3, -1)$, etc.), with the special point $(-1, -1)$ coming last. Note that the value of f on the (k, ℓ) support point is $(-1)^{k\ell+k+\ell}$.

Let \tilde{A} be a $D \times (D + 1)$ matrix whose columns are indexed by the support points and whose rows are indexed by the constraint monomials. Define $\tilde{A}[(i, j), (k, \ell)]$ to be the value of the (i, j) th constraint monomial at the (k, ℓ) th support point, multiplied by the value of f at the (k, ℓ) th support point. This definition shall include the case of the special $(-1, -1)$ support point, to whose column (the rightmost column of \tilde{A}) we assign the index $(0', 0')$ for reasons that will become clear soon. Let A be the $(D + 1) \times (D + 1)$ matrix given by adding a row of 1's to the bottom of \tilde{A} . For notational convenience we will also give this bottom row the index $(0', 0')$. So for $(i, j), (k, \ell) \neq (0', 0')$ we have:

¹Given the logarithmic upper bound on threshold degree for AND-MAJ _{n} proved by [6], it makes sense to see a support requiring coordinates exponential in d .

$$A[(i, j), (k, \ell)] = (-1)^{k(j+1)+\ell(i+1)+k\ell} h^{ik+j\ell}, \quad (2)$$

i.e.

$$A = \begin{bmatrix} 1 & -1 & -1 & 1 & -1 & 1 & -1 & \cdots & (-1)^d & -1 \\ 1 & -h & 1 & h^2 & h & 1 & -h^3 & \cdots & 1 & 1 \\ 1 & 1 & -h & 1 & h & h^2 & 1 & \cdots & (-h)^d & 1 \\ 1 & -h^2 & -1 & h^4 & -h^2 & 1 & -h^6 & \cdots & (-1)^d & -1 \\ \vdots & \vdots \\ 1 & (-1)^{d+1} & -h^d & 1 & (-1)^{d+1}h^d & h^{2d} & (-1)^{d+1} & \cdots & (-1)^d h^{d^2} & (-1)^{d+1} \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & \cdots & 1 & 1 \end{bmatrix}.$$

Recall that we want to find values $\mathbf{w} : \mathcal{Z} \rightarrow \mathbf{R}$ such that $\sum_{(x,y) \in \mathcal{Z}} \mathbf{w}(x,y) f(x,y) x^i y^j = 0$ for all constraints and such that $\sum_{(x,y) \in \mathcal{Z}} \mathbf{w}(x,y) = 1$. By construction these values are uniquely given by the solution to the following system of linear equations:

$$A \begin{bmatrix} w_{(0,0)} \\ w_{(1,0)} \\ w_{(0,1)} \\ w_{(2,0)} \\ \vdots \\ w_{(0,d)} \\ w_{(-1,-1)} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}. \quad (3)$$

In the remainder of the proof we show that by taking $h = \Theta(d^9)$, we can ensure that the solution to Equation (3) consists entirely of nonnegative numbers, and hence \mathbf{w} corresponds to a true probability distribution as desired. Since $h = O(d^9)$ means that $\text{size}(\mathcal{Z}) = d^{O(d)}$, and we may take h to be odd, this proves Theorem 4.2.

We shall consider solving Equation (3) via Cramer's rule. Cramer's rule tells us that Equation (3) implies:

$$w_{(u,v)} = \frac{\det A_{(u,v)}}{\det A},$$

where $A_{(u,v)}$ denotes the matrix A with the (u, v) column replaced by the right hand side of Equation (3), namely $\begin{bmatrix} 0 & 0 & 0 & \cdots & 0 & 1 \end{bmatrix}^T$. To show that each $w_{(u,v)}$ is nonnegative we will show that $\det A_{(u,v)}$ and $\det A$ have the same sign.

Let $\sigma \in \{+1, -1\}$ be the sign of the product of the diagonal entries in A . We will prove the following two lemmas and thus prove Theorem 4.2:

Lemma 4.3 $\text{sign}(\det A) = \sigma$.

Lemma 4.4 $\text{sign}(\det A_{(u,v)}) = \sigma$ for all (u, v) .

Since the proofs of these lemmas are rather technical, we give a word of intuition before entering into the details. Lemma 4.3 essentially says that the dominant contribution to $\det A$, viewed as a sum over permutations, comes from the identity permutation corresponding to the diagonal elements in A . Intuitively this is because the product of the diagonal elements of A yields the largest total exponent for h over all permutations, which is a consequence of the way we chose the support \mathcal{Z} . (As a toy example, consider the 4×4 submatrix in the top left immediately below (2). The product of the 4 diagonal elements yields an exponent of 6 while every other permissible choice of 4 elements from this submatrix yields an exponent of at most 5.) Lemma 4.4 essentially says that a similar phenomenon holds even for the matrix $A_{(u,v)}$ (which is less well-structured than A); not surprisingly the proof there is more complicated than the one for Lemma 4.3.

4.1.1 Proof of Lemma 4.3 To prove Lemma 4.3 we view $\det A$ as a polynomial in h . Let $T := \deg(\det A)$ be the degree of $\det A$. We show that the leading term of $\det A$ (corresponding to h^T) dominates all the other terms for h sufficiently large, and thus the sign of $\det A$ is the same as the sign of the leading term. More precisely, we establish the following two facts:

Claim 4.5 *The coefficient of h^T in $\det A$ is 2σ .*

Claim 4.6 *For all $u \geq 1$ the coefficient of h^{T-u} in $\det A$ is at most $2(D+2)^{4u}$ in magnitude.*

Claim 4.6 implies that the sum of the absolute values of the lower-order terms in $\det A$ is at most $\sum_{u=1}^T 2(D+2)^{4u} h^{T-u} \leq h^T \sum_{u=1}^T (2(D+2)^4/h)^u$. If we take h to be $\Theta(d^9)$ then this quantity will be strictly smaller than h^T . But by Claim 4.5 we have that the leading term of $\det A$ is $2\sigma h^T$. Thus $\text{sgn}(\det A) = \sigma$ and Lemma 4.3 holds.

We set the stage before proving Claims 4.5 and 4.6 with some notation and some observations. Let \mathbf{S} denote the permutation group on the $D+1$ indices $(0, 0), (1, 0), (0, 1), (2, 0), (1, 1), (0, 2), \dots, (0, d), (0', 0')$. Then:

$$\det A = \sum_{\pi \in \mathbf{S}} \text{sgn}(\pi) \prod_{(i,j)} A[(i, j), \pi(i, j)]. \quad (4)$$

Recall that for $(i, j), (k, \ell) \neq (0', 0')$, the entry $A[(i, j), (k, \ell)]$ is $\pm h^{ik+j\ell}$, which we will write as $\pm \exp_h((i, j) \cdot (k, \ell))$, with $\exp_h(t)$ denoting h^t and \cdot being the usual dot product. In the

case that $(i, j) = (0', 0')$ or $(k, \ell) = (0', 0')$, the entry $A[(i, j), (k, \ell)]$ is $\pm 1 = \pm h^0$. If we define $(0', 0') \cdot (a, b)$ to be 0, then we have that for any permutation $\pi \in \mathbf{S}$,

$$\prod_{(i,j)} A[(i, j), \pi(i, j)] = \pm \exp_h \left(\sum_{(i,j)} (i, j) \cdot \pi(i, j) \right).$$

Given a permutation $\pi \in \mathbf{S}$, write $t(\pi) = \sum_{(i,j)} (i, j) \cdot \pi(i, j)$, so the permutation π contributes ± 1 to the coefficient of $h^{t(\pi)}$ in $\det A$. Then the absolute value of the coefficient of h^u in $\det A$ is at most $|\{\pi \in \mathbf{S} : t(\pi) = u\}|$. We will use this fact to bound all the lower-order terms in $\det A$; for the leading term we will pay more attention to the signs.

To calculate $t(\pi)$ from π , we decompose the permutation π as a product of cycles. For each cycle $\pi_0 = ((i_1, j_1) (i_2, j_2) \cdots (i_m, j_m))$ we have by simple arithmetic:

$$\sum_{r=1}^m (i_r, j_r) \cdot (i_r, j_r) - \sum_{r=1}^m (i_r, j_r) \cdot \pi_0(i_r, j_r) = \frac{1}{2} \sum_{r=1}^m (i_r - i_{r-1})^2 + (j_r - j_{r-1})^2, \quad (5)$$

where we use the notation $(i_0, j_0) = (i_m, j_m)$. (Note that a geometric interpretation of this quantity is that it is half the sum of the squares of the lengths of the line segments which make up the cycle in the two-dimensional plane from (i_1, j_1) to (i_2, j_2) to (i_3, j_3) to \dots to (i_m, j_m) to (i_1, j_1) .) In particular, this quantity is at least 1 for every nontrivial cycle, where a trivial cycle for us is either a cycle of length 1 or the transposition exchanging $(0, 0)$ and $(0', 0')$. The quantity in Equation (5) is 0 for trivial cycles. Thus we have that the identity permutation and the transposition $((0, 0), (0', 0'))$ are the only two permutations which achieve the maximum value $t(\pi) = T$. It is easy to see that this maximum value T is $\sum_{(i,j)} i^2 + j^2$, which one easily calculates to be $T := d(d+1)^2(d+2)/6$. We further see that every other permutation ‘‘pays a penalty’’ in its t value for each nontrivial cycle it contains, and this penalty is given by the right-hand side of Equation (5). Hence to calculate $t(\pi)$ from π we simply sum up the penalties for each cycle in its cycle decomposition and subtract the total from T .

Proof of Claim 4.5: As described above, we have that there are exactly two permutations which lead to the maximum power h^T in Equation (4): the identity permutation which takes all the diagonal elements, and the $((0, 0), (0', 0'))$ transposition which takes the top-right entry of A , the bottom-left entry of A , and the diagonal elements otherwise. The product of the top-left and bottom-right entries of A is 1. The product of the top-right and bottom-left entries is -1 ; however this gets flipped to $+1$ by the sign of the permutation (it is a transposition so its sign is -1). We conclude that leading term of $\det A$ is $2\sigma h^T$ where $\sigma \in \{-1, 1\}$ is the sign of the product of the diagonal entries in A . (Claim 4.5) ■

Proof of Claim 4.6: To bound the coefficient on the lower-order term h^{T-u} in $\det A$ we simply count the number of permutations π which have $t(\pi) = T - u$. This count gives an upper bound on the magnitude of the coefficient. If $t(\pi) = T - u$ then the penalty accounting scheme from Equation (5) tells us that π has at most u nontrivial cycles. In fact we can say more: any nontrivial cycle of length m must incur a penalty of at least $\lfloor m/2 \rfloor$. (This follows from the geometric interpretation described earlier, together with the fact that any nontrivial cycle of length $m \geq 3$ can include at most one segment of length 0 between $(0, 0)$ and $(0', 0')$.) Consequently, if $t(\pi) = T - u$ then the lengths of the nontrivial cycles in π 's cycle decomposition must sum to at most $3u$ (in the worst case all its cycles may be 3-cycles each of which incurs a penalty of 1). Now observe that there are at most $(D + 2)^{4u}$ permutations on $D + 1$ elements which decompose into at most u cycles whose total length is at most $3u$. (Any such sequence of cycles can be written as a string of length $4u$ over a $D + 2$ element alphabet, where the extra symbol is used to mark the end of each cycle.) Doubling this upper bound covers the optional addition of the trivial $((0, 0), (0', 0'))$ transposition. We thus may conclude that there are at most $2(D + 2)^{4u}$ permutations $\pi \in \mathbf{S}$ which have $t(\pi) = T - u$. (Claim 4.6) ■

4.1.2 Proof of Lemma 4.4 It now remains to show that $\text{sgn}(\det A_{(u,v)}) = \sigma$ for each (u, v) . By the nature of cofactor expansion, $\det A_{(u,v)}$ is equal to a certain sign ρ , times the determinant of A with the bottom row and the (u, v) column deleted. In the case $(u, v) = (0', 0')$ we have $\rho = 1$ and we shall write $A'_{(0',0')}$ for the matrix A with its last row and column deleted. For all $(u, v) \neq (0', 0')$, let us write $A'_{(u,v)}$ for the matrix gotten by first deleting the bottom row and (u, v) column from A , and then moving the $(0', 0')$ column leftward until it is in the place where the old (u, v) used to be. Shifting the $(0', 0')$ column like this incurs a sign change equal to $-\rho$; we conclude that $\det A_{(u,v)} = -\det A'_{(u,v)}$. Hence it is sufficient for us to show that $\text{sgn}(\det A'_{(0',0')}) = \sigma$ and that $\text{sgn}(\det A'_{(u,v)}) = -\sigma$ for all $(u, v) \neq (0', 0')$.

Let us begin by dispensing with the cases $(u, v) = (0', 0')$ or $(0, 0)$. In both of these cases $A'_{(u,v)}$ is very similar to A with the last row and column deleted; when $(u, v) = (0', 0')$ this is exactly what $A'_{(u,v)}$ is, and when $(u, v) = (0, 0)$ some of the signs in the first column are changed. Hence the analysis of $\det A'_{(u,v)}$ is virtually identical to the above analysis of $\det A$, except that $(0', 0')$ is no longer present. The leading term will therefore be equal to the top-left entry of $A'_{(u,v)}$ times σh^T ; this entry is 1 when $(u, v) = (0', 0')$ and is -1 when $(u, v) = (0, 0)$, as desired. The analysis bounding the lower-order terms goes through in essentially the same way as before (again without $(0', 0')$) and we conclude that $\text{sgn}(\det A'_{(0',0')}) = \sigma$ and $\text{sgn}(\det A'_{(0,0)}) = -\sigma$ as desired.

Throughout the rest of this section we assume that $(u, v) \neq (0', 0'), (0, 0)$. Let $T_{(u,v)}$ denote the degree of $\det(A'_{(u,v)})$.

Let C denote the number of paths from (u, v) to $(1, 0)$ plus the number of paths from (u, v) to $(0, 1)$, where each path uses steps $(-1, 0)$, $(0, -1)$, and $(-1, -1)$. (Such paths are known as *Delannoy paths*, and the number of such paths between a pair of points is a *Delannoy number*, see e.g. p. 80 of [11]; hence C is a sum of two Delannoy numbers.)

We will prove the following two claims:

Claim 4.7 *The coefficient of $h^{T_{(u,v)}}$ in $\det(A'_{(u,v)})$ is $-2\sigma C$.*

Claim 4.8 *For all $s \geq 1$ the coefficient of $h^{T_{(u,v)}-s}$ in $\det(A'_{(u,v)})$ is at most $4C(D+2)^{4s}$ in magnitude.*

As in the previous subsection, these two claims show that we may take $h = \Theta(d^9)$ to obtain $\text{sgn}(\det(A'_{(u,v)})) = -\sigma$, so they suffice to prove the lemma.

Studying $\det A'_{(u,v)}$ is slightly more complex than studying $\det A$ because its rows and columns no longer have the same names; the rows of $A'_{(u,v)}$ are named $(0, 0)$, $(1, 0)$, $(0, 1)$, $(2, 0)$, \dots , (u, v) , \dots , $(0, d)$, whereas the columns are named $(0, 0)$, $(1, 0)$, $(0, 1)$, $(2, 0)$, \dots , $(0', 0')$, \dots , $(0, d)$. To deal with this, we will let \mathbf{S}' denote the permutation group on the D row indices of $A'_{(u,v)}$, and we will view (u, v) as $(0', 0')$ whenever it is the “output” of a permutation. To be precise, let ι be a mapping which maps (i, j) to (i, j) for each $(i, j) \neq (u, v)$, and maps (u, v) to $(0', 0')$. Then our determinant equation becomes:

$$\det A'_{(u,v)} = \sum_{\pi \in \mathbf{S}'} \text{sgn}(\pi) \prod_{(i,j)} A[(i, j), \iota(\pi(i, j))]. \quad (6)$$

We may write $t(\pi) = \sum_{(i,j)} (i, j) \cdot \iota(\pi(i, j))$, so we have $\prod_{(i,j)} A[(i, j), \iota(\pi(i, j))] = \pm h^{t(\pi)}$.

As before we will calculate $t(\pi)$ by considering the cycle decomposition of π and computing the penalty difference from $T = d(d+1)^2(d+2)/6$ for each cycle. Since now the “identity” permutation does not exist, the permutations maximizing $t(\pi)$ may not achieve T ; indeed, since $(u, v) \neq (0', 0')$ it is the case that maximizing permutations will not achieve $t(\pi) = T$. Let us now find the new highest value for $t(\pi)$. The cycle decomposition of π contains a unique cycle (which may be a 1-cycle) containing (u, v) , and perhaps other cycles which do not contain (u, v) . For the cycles not containing (u, v) , ι does not enter into the picture in calculating $t(\pi)$; hence Equation (5) still holds and we conclude that for any π with maximal $t(\pi)$, all its nontrivial cycles must involve (u, v) . Thus, in order to find all maximizing π 's, it is sufficient to determine which cycles containing (u, v) give the smallest penalty.

Let π^* be a cycle containing (u, v) ; say $\pi^* = ((u, v), (i_1, j_1), (i_2, j_2), \dots, (i_m, j_m))$, so according to our conventions π^* maps (u, v) to (i_1, j_1) , maps (i_r, j_r) to (i_{r+1}, j_{r+1}) for $1 \leq r \leq m-1$, and maps (i_m, j_m) to $\iota(u, v) = (0', 0')$. Write $(i_0, j_0) = (u, v)$. Then akin to Equation (5) we have:

$$\begin{aligned}
& \sum_{r=0}^m (i_r, j_r) \cdot (i_r, j_r) - \sum_{r=0}^m (i_r, j_r) \cdot \iota(\pi^*(i_r, j_r)) \\
= & \sum_{r=0}^m (i_r, j_r) \cdot (i_r, j_r) - \sum_{r=0}^m (i_r, j_r) \cdot (i_{r+1 \bmod m+1}, j_{r+1 \bmod m+1}) + i_r u + j_r v \\
= & \frac{1}{2} \left(\left(\sum_{r=1}^m (i_r - i_{r-1})^2 + (j_r - j_{r-1})^2 \right) + (u - i_r)^2 + (v - j_r)^2 \right) + i_m u + j_m v \text{ (as in Equation (5))} \\
= & \frac{1}{2} \left(\left(\sum_{r=1}^m (i_r - i_{r-1})^2 + (j_r - j_{r-1})^2 \right) + i_m^2 + j_m^2 + u^2 + v^2 \right). \tag{7}
\end{aligned}$$

The geometric interpretation of the quantity on the right-hand side of Equation (7) is that it is half the sum of the squares of the path segments on the closed path from (u, v) to (i_1, j_1) to (i_2, j_2) to \dots to (i_m, j_m) to $(0, 0)$ to (u, v) . It is immediate that in a cycle minimizing this quantity, there should be no path step which has either x or y displacement greater than 1 in magnitude (aside from the step from $(0, 0)$ to (u, v) which is forced). Consequently, the permutations π which maximize $t(\pi)$ are precisely those cycles π^* such that (1) $i_{r+1} - i_r \in \{-1, 0\}$ and $j_{r+1} - j_r \in \{-1, 0\}$ for $0 \leq r < m$, and (2) $i_m, j_m \in \{0, 1\}$. It is easy to see that each such maximizing permutation has $t(\pi) = T_{(u,v)} = T - \frac{1}{2}(u + v + u^2 + v^2)$.

Proof of Claim 4.7: Now we can compute the coefficient of $h^{T_{(u,v)}}$ in $\det A'_{(u,v)}$. Given a permutation π maximizing $t(\pi)$, let $\epsilon(\pi)$ denote the sign of π 's contribution to the determinant computation of Equation (6), i.e. $\epsilon(\pi) = \text{sgn}(\pi) \prod_{(i,j)} \text{sgn}(A[(i, j), \iota(\pi(i, j))])$. Then the leading coefficient of $\det A'_{(u,v)}$ is just the sum of $\epsilon(\pi)$ over all maximizing π .

Let $\pi = ((u, v), (i_1, j_1), (i_2, j_2), \dots, (i_m, j_m))$ be a maximizing permutation; as before we write $(i_0, j_0) = (u, v)$. By the definition of σ as the product of the signs of A 's diagonal elements, we get that:

$$\begin{aligned}
\sigma \epsilon(\pi) = \text{sgn}(\pi) & \left(\prod_{r=0}^{m-1} \text{sgn}(A[(i_r, j_r), (i_r, j_r)]) \text{sgn}(A[(i_r, j_r), (i_{r+1}, j_{r+1})]) \right) \\
& \cdot \text{sgn}(A[(i_m, j_m), (i_m, j_m)]) \text{sgn}(A[(i_m, j_m), (0', 0')]).
\end{aligned}$$

We claim that for each $0 \leq r \leq m-1$ we have $\text{sgn}(A[(i_r, j_r), (i_r, j_r)]) \text{sgn}(A[(i_r, j_r), (i_{r+1}, j_{r+1})]) =$

-1 , independent of (i_r, j_r) . For from Equation (2) we know that:

$$\begin{aligned} & \text{sgn}(A[(i_r, j_r), (i_r, j_r)])\text{sgn}(A[(i_r, j_r), (i_{r+1}, j_{r+1})]) \\ &= \exp_{-1}(i_r(j_r + 1) + j_r(i_r + 1) + i_r j_r) \exp_{-1}(i_{r+1}(j_r + 1) + j_{r+1}(i_r + 1) + i_{r+1} j_{r+1}) \\ &= \exp_{-1}(i_r j_r + i_{r+1} j_r + i_r j_{r+1} + i_{r+1} j_{r+1} + i_r + i_{r+1} + j_r + j_{r+1}) \\ &= \exp_{-1}((i_r + i_{r+1} + 1)(j_r + j_{r+1} + 1) - 1), \end{aligned}$$

which is always -1 as claimed, because $(i_r, j_r) - (i_{r+1}, j_{r+1}) \in \{(1, 0), (0, 1), (1, 1)\}$.

Thus we have:

$$\begin{aligned} \sigma \epsilon(\pi) &= \text{sgn}(\pi)(-1)^m \text{sgn}(A[(i_m, j_m), (i_m, j_m)])\text{sgn}(A[(i_m, j_m), (0', 0')]) \\ &= +\text{sgn}(A[(i_m, j_m), (i_m, j_m)])\text{sgn}(A[(i_m, j_m), (0', 0')]) \quad (*), \end{aligned}$$

because π is a cycle of length $m + 1$. If $(i_m, j_m) = (1, 1)$ then $(*) = -1$; otherwise, $(*) = +1$. Hence we conclude that $\epsilon(\pi) = \sigma$ if $(i_m, j_m) = (1, 1)$ and $\epsilon(\pi) = -\sigma$ if $(i_m, j_m) \in \{(0, 0), (1, 0), (0, 1)\}$. For each maximizing cycle π of length $m + 1$ with $(i_m, j_m) \neq (0, 0)$, there is a corresponding maximizing cycle π' of length $m + 2$ obtained by appending $(i_{m+1}, j_{m+1}) = (0, 0)$ to π . Thus we have $\epsilon(\pi) + \epsilon(\pi') = 0$ when $(i_m, j_m) = (1, 1)$ and $\epsilon(\pi) + \epsilon(\pi') = -2\sigma$ when $(i_m, j_m) = (1, 0)$ or $(0, 1)$. In conclusion, the leading term in $\det A'_{(u,v)}$ is exactly $-2\sigma C h^{T(u,v)}$, where as stated above C is the number of paths from (u, v) to $(1, 0)$ plus the number of paths from (u, v) to $(0, 1)$, where each path uses steps $(-1, 0)$, $(0, -1)$, and $(-1, -1)$. Since $(u, v) \neq (0, 0)$ we have $C \geq 1$, and the claim is proved. (Claim 4.7) ■

Proof of Claim 4.8: We must upper-bound the magnitude of the lower-order terms in $\det A'_{(u,v)}$. We do this as in the analysis of $\det A$ by upper-bounding the number of permutations π with $t(\pi) = T_{(u,v)} - s$. To each $\pi \in \mathbf{S}'$ we will associate a *maximizing* permutation π^* (i.e., one for which $t(\pi^*) = T_{(u,v)}$), and a “deviation description.” We will show that the longer the deviation description, the smaller $t(\pi)$ is compared to $t(\pi^*)$. Thus the number of permutations π with $t(\pi)$ close to $T_{(u,v)}$ will be upper-bounded by the number of optimal permutations times the number of short deviation descriptions.

Let π be an arbitrary permutation in \mathbf{S}' and write π as the product of a cycle π_0 involving (u, v) , and some other cycles π_1, \dots, π_s . The maximizing permutation π^* we associate with π will depend only on π_0 . View π_0 geometrically as a path from (u, v) to $\pi_0^{-1}(u, v)$. Call a path “optimal” if it only uses steps $(-1, 0)$, $(0, -1)$, and $(-1, -1)$, so in particular every maximizing permutation contains one nontrivial cycle containing (u, v) whose corresponding path is optimal. We will split π_0 up into its optimal and nonoptimal segments. Specifically, $a_i, b_i, c_i, d_i, \dots, a_r, b_r$ are defined as follows: π_0 proceeds optimally from (u, v) to (a_1, b_1) ,

at which point it takes a nonoptimal step. Let (c_1, d_1) be the first point it proceeds to subsequently with the property that $c_1 \leq a_1$, $d_1 \leq b_1$. Then π_0 proceeds optimally from (c_1, d_1) to (a_2, b_2) , at which point it makes a nonoptimal step. Let (c_2, d_2) be the first point it proceeds to subsequently with $c_2 \leq a_2$, $d_2 \leq b_2$. Continuing in this fashion, let (a_r, b_r) be the last point reached in the last optimal segment of π_0 ; π_0 may optionally go on and reach $\pi_0^{-1}(u, v)$. We will let the maximizing permutation π^* associated with π be any optimal path that agrees with π_0 on all steps from (u, v) to (a_1, b_1) , all steps from (c_1, d_1) to (a_2, b_2) , \dots , all steps from (c_{r-1}, d_{r-1}) to (a_r, b_r) , and then ends by proceeding optimally to $(0, 0)$.

The deviation description of π will simply be a list of all of the cycles π_1, \dots, π_s not containing (u, v) , along with a description of π_0 's deviation from π^* . This deviation consists of the path from (a_1, b_1) to (c_1, d_1) , from (a_2, b_2) to (c_2, d_2) , etc., possibly ending with some path from (a_r, b_r) to a point not in $\{0, 1\}^2$. Note that π can be recovered from π^* and the deviation description.

Now let us compute $t(\pi^*) - t(\pi)$. This difference is equal to $(T - t(\pi)) - (T - t(\pi^*))$, and Equations (5) and (7) tell us how to compute these quantities. By Equation (5), $t(\pi)$ pays an extra penalty over $t(\pi^*)$ for each of its cycles not involving (u, v) , π_1, \dots, π_s . As in the analysis of $\det A$ we know that such a cycle of length m incurs a penalty of at least $\lfloor m/2 \rfloor$. Equation (7) allows us to compare the penalties against T that each of $t(\pi^*)$ and $t(\pi)$ pays. Every time π_0 deviates from π^* it pays an extra penalty of at least 1. Indeed, just as in the analysis of extraneous cycles, a deviation path from (a_i, b_i) to (c_i, d_i) which touches m nodes must incur an extra penalty of at least $\lfloor m/2 \rfloor$. This holds also for a final deviation path which does not end up in $\{0, 1\}^2$, since it must pay for half the squared distance from the origin of its endpoint. Both π^* and π_0 pay equally for the final $\frac{1}{2}(k^2 + \ell^2)$ term.

In conclusion, if the total length of the cycles and deviation paths in π 's deviation description is m then $(T - t(\pi)) - (T - t(\pi^*))$ is at least $\lfloor m/2 \rfloor$; i.e., $t(\pi) \leq T_{(u,v)} - \lfloor m/2 \rfloor$. Hence as in the analysis of $\det A$ we can get an upper bound of $(D+2)^{4s} \cdot \#\{\text{number of maximizing } \pi_0\}$ for the number of permutations π with $t(\pi) = T_{(u,v)} - s$. But note that the leading coefficient in $\det A'_{(u,v)}$ has magnitude $2C$, and $2C$ is at least half the number of maximizing permutations π_0 . To see this, recall that C counts the number of optimal paths from (u, v) to either $(1, 0)$ or $(0, 1)$, and each maximizing permutation corresponds to an optimal path to one of $(0, 0)$, $(0, 1)$, $(1, 0)$, $(1, 1)$. The number of optimal paths to $(1, 1)$ is at most C (each such path can be extended to a path ending in $(1, 0)$ or $(0, 1)$), and hence the number of optimal paths to $(0, 0)$ is at most $2C$ (since the next to last point on any such path is either $(1, 0)$, $(0, 1)$ or $(1, 1)$). It follows that the magnitude of the sum of all lower-order terms in $\det A'_{(u,v)}$ is at most $\sum_{s=1}^{T_{(u,v)}} 4C(D+2)^{4s} h^{T_{(u,v)}-s}$, and the claim is proved. (Claim 4.8) ■

5 Upper bounds for Boolean formulas

In this section we consider Boolean formulas composed of NOT gates and unbounded fan-in AND and OR gates. The *depth* of a formula is the length of the longest path from the root to any leaf, and the *size* is the number of occurrences of variables.

We will also consider variants of polynomial threshold functions in which the polynomial is subject to a stricter requirement than just sign-representing f . Following Nisan and Szegedy [23], we write $\widetilde{\deg}(f)$ to denote the minimum degree of any polynomial which approximates f to within $1/3$ on all inputs, i.e. such a polynomial $p(x)$ must satisfy:

$$\forall x \in \{0, 1\}^n \quad |f(x) - p(x)| \leq \frac{1}{3}.$$

Clearly we have $\widetilde{\deg}(f) \geq \text{thr}(f)$ for all f . We write $|p - f|_\infty$ to denote $\max_{x \in \{0, 1\}^n} |p(x) - f(x)|$. Thus if $|p - f|_\infty < \frac{1}{3}$ we have $\deg(p) \geq \widetilde{\deg}(f) \geq \text{thr}(f)$.

We prove two similar theorems bounding the polynomial threshold degree of Boolean formulas:

Theorem 5.1 *Let f be computed by a Boolean formula of depth d and size s . Then there is a polynomial $p(x_1, \dots, x_n)$ of degree at most $2^{O(d)}(\log s)^{5d/2}\sqrt{s}$ such that $|p - f|_\infty \leq \frac{1}{s}$.*

Theorem 5.2 *Let f be computed by a Boolean formula of depth d and size s . Then there is a polynomial $p(x_1, \dots, x_n)$ of degree at most $2^{O(d)}(\log s)^{5d} s^{\frac{1}{2} - \frac{1}{2^{d+1}-2}}$ such that $\text{sgn}(p(x)) = f(x)$.*

The proof technique in both cases is to first manipulate the formula to get a more structured form, and then to apply real approximating functions (Chebyshev polynomials, the rational functions of [6]) at each gate.

Some preliminary notes: throughout this section we let 0 represent FALSE and 1 represent TRUE, and thus we view Boolean functions as mappings from $\{0, 1\}^n$ to $\{0, 1\}$. Without loss of generality we may assume that our formulas contain no NOT gates; i.e., they consist only of AND and OR gates. This is because any negations in a formula F can be pushed to the leaves using DeMorgan's laws with no increase in size or depth. Once all negations are at the leaves we can replace each negated variable $\neg x_i$ with a variable y_i to obtain a formula F' which has no negations. Given a polynomial which sign-represents or approximates F' , we can obtain a corresponding polynomial for F by replacing each y_i with $1 - x_i$, and this will not increase the degree.

5.1 Proof of Theorem 5.1 Henceforth the variables c_1, c_2, \dots refer to fixed universal constants. We will use the following lemma proved in Appendix B:

Lemma 5.3 *Let $f = \bigwedge_{i=1}^{\ell} f_i$ be a Boolean formula where $\ell \geq 2$. For $1 \leq i \leq \ell$ let p_i be a polynomial with $\deg(p_i) \leq r$ such that $|p_i - f_i|_{\infty} \leq \epsilon$, where $0 < \epsilon < \frac{1}{8\ell}$. Then there is a polynomial p with $\deg(p) \leq (4\sqrt{\ell} \log \frac{1}{\epsilon})r$ such that $|p - f|_{\infty} \leq (c_2 \ell \log \frac{1}{\epsilon})\epsilon$.*

It is easy to see that an identical result holds if $f = \bigvee_{i=1}^{\ell} f_i$, i.e. f 's top-level gate is an OR instead of an AND.

The following lemma is now easy to establish:

Lemma 5.4 *Let f be computed by a Boolean formula F of depth d and size s . Suppose that for any path from the root of F to a leaf, the product of the fanins of the gates on the path is at most t . Then there is a polynomial p with $\deg(p) \leq (c_3 \log s)^d \sqrt{t}$ such that $|p - f|_{\infty} \leq \frac{1}{s}$.*

Proof: Note first that for any Boolean formula of size s , there is a multilinear interpolating polynomial which computes the formula exactly and is of degree at most s . Consequently if $(c_3 \log s)^d \sqrt{t} \geq s$ the lemma is trivially true, so we assume that $(c_3 \log s)^d \sqrt{t} < s$.

Consider the formula F . Each leaf contains some variable x_i , so clearly there is a degree-1 polynomial which exactly computes the function at each leaf. Now apply Lemma 5.3 successively to every gate in F , going up from the leaves to the root. At each leaf we may take ϵ in Lemma 5.3 to be any positive value; we take $\epsilon = \frac{1}{s^3}$. Each time we go up through a gate of fanin ℓ the value of ϵ which we may use in Lemma 5.3 is multiplied by at most $c_2 \ell \log(s^3) = c_3 \ell \log s$. An easy induction on the depth of F shows that at the root we obtain a polynomial p such that

$$\deg(p) \leq (4 \log(s^3))^d \sqrt{t} < (c_3 \log s)^d \sqrt{t}$$

and

$$|p - f|_{\infty} \leq \frac{1}{s^3} \cdot (c_3 \log s)^d t < \frac{1}{s^3} \cdot s^2 = \frac{1}{s}$$

as desired. ■

With Lemmas 5.3 and 5.4 in hand, in order to prove Theorem 5.1 it suffices to bound the product of the fanins on any path from the root to a leaf. In an arbitrary formula this product can be quite large; it is easy to construct a formula of size s and depth d in which there is a path composed of d gates each of fanin $\frac{s}{d}$. Thus in general this product can be as large as $(\frac{s}{d})^d$; however we can remedy this situation as described below.

Lemma 5.5 *Let F be a formula of size s and depth d . There is a formula G of size s and depth $2d$ computing the same function as F such that the product of the fanins on any root-to-leaf path in G is at most $(4 \log s)^d s$.*

Proof: We prove the following slightly stronger statement: for any formula F of size s and depth d , there is a formula G of size s and depth $2d$ computing F such that the product of the fanins on any root-to-leaf path in G is at most $(2\lceil\log s\rceil)^d s$. The lemma follows since $2\log s \geq \lceil\log s\rceil$ for all s .

The proof is by induction on d . For $d = 0$ the formula must be a single variable so $s = 1$ and the claim is trivially true. Suppose without loss of generality that $F = \bigwedge_{i=1}^{\ell} F_i$ where $\ell \geq 2$, each F_i has depth at most $d - 1$, and the sum of the sizes of F_1, \dots, F_{ℓ} is s . Let $|F_i|$ denote the size of F_i . We partition the formulas F_1, \dots, F_{ℓ} into disjoint classes $C_1, \dots, C_{\lceil\log s\rceil}$ where the class C_j contains exactly those F_i such that $2^{j-1} \leq |F_i| < 2^j$. By the induction hypothesis each formula $F_i \in C_j$ has an equivalent formula G_i of size $|F_i|$ and depth at most $2d - 2$ such that the product of the fanins along any root-to-leaf path in G_i is at most $(2\lceil\log s\rceil)^{d-1} |F_i| < 2^{d+j-1} \lceil\log s\rceil^{d-1}$. Let $G = \bigwedge_{j=1}^{\lceil\log s\rceil} H_j$ where the formula H_j is defined as $H_j = \bigwedge_{i:F_i \in C_j} G_i$.

To see that this works, first observe that each C_j contains at most $s/2^{j-1}$ formulas F_i . Thus the fanin at the root of H_j is at most $s/2^{j-1}$, and hence the product of the fanins along any path in H_j is at most $2^d s \lceil\log s\rceil^{d-1}$. Thus the product of the fanins along any path in G is at most $(2\lceil\log s\rceil)^d s$ as desired and the lemma is proved. ■

Theorem 5.1 follows from combining Lemmas 5.4 and 5.5.

5.2 Proof of Theorem 5.2 Theorem 5.1 uses Chebyshev polynomials to construct polynomials which closely approximate Boolean formulas. In this section we extend this construction using rational functions to construct polynomials which only sign-represent Boolean formulas. The bound given in Theorem 5.2 is asymptotically superior to Theorem 5.1 for any constant d .

We define the degree of a rational function $p(x)/q(x)$ to be $\max\{\deg(p), \deg(q)\}$. Theorem 5.2 is a consequence of the following lemma:

Lemma 5.6 *Let f be computed by a Boolean formula of depth d and size s . Then there is a rational function r of degree at most $c_4^d (\log s)^{5d} s^{\frac{1}{2} - \frac{1}{2^{d+1}-2}}$ such that $|r - f|_{\infty} < \frac{1}{4s}$.*

The proof, which is by induction on d , is given in Appendix C. To see that Lemma 5.6 implies Theorem 5.2, let $r(x) = p(x)/q(x)$. Since $r(x) \in [-\frac{1}{s}, \frac{1}{s}]$ if $f(x) = 0$ and $r(x) \in [1 - \frac{1}{s}, 1 + \frac{1}{s}]$ if $f(x) = 1$, we have that

$$f(x) = \text{sgn}(r(x) - 1/2) = \text{sgn}((r(x) - 1/2)q(x)^2) = \text{sgn}(p(x)q(x) - q(x)^2/2)$$

for all $x \in \{0, 1\}^n$.

5.3 Discussion In earlier work Klivans and Servedio [17] showed that any Boolean formula of constant depth d and size s has a polynomial threshold function of degree $\tilde{O}(s^{1-\frac{1}{3 \cdot 2^{d-3}}})$. For even moderately large constant values of d , this bound is not far from the trivial upper bound of s . In contrast, our new bounds are considerably stronger. Theorem 5.2 gives an $o(s^{1/2})$ bound for some $d = \Omega(\log \log s)$, and Theorems 5.1 and 5.2 both give a bound of $O(s^{1/2+\epsilon})$ for any $d = o(\frac{\log s}{\log \log s})$. To our knowledge Theorems 5.1 and 5.2 are the first nontrivial upper bounds on polynomial threshold function degree for formulas of superconstant depth.

In other earlier work, Buhrman, Cleve and Wigderson [8] gave an $O(s^{1/2} \log^{d-1}(s))$ upper bound on the bounded-error quantum query complexity of certain Boolean formulas of size s and depth d . Since the bounded-error quantum query complexity upper bounds the required degree for an approximating polynomial (see Theorem 18 of [9]), their results imply an $O(s^{1/2} \log^{d-1}(s))$ upper bound on the degree of the formulas that they consider. However, their bound applies only to “Sipser functions”, namely to formulas of size s and depth d in which all of the gates at any given depth have the same fanin (the fanins can be different for gates at different depths). Our Theorem 5.1 thus generalizes their bound on the degree of approximating polynomials to a substantially broader class of formulas, since we do not make any restrictions on fanin.²

5.4 Learning Boolean formulas of superconstant depth in subexponential time

We close this section by describing some consequences of our results in computational learning theory. It is known (see [17, 16]) that if a class C of Boolean functions over $\{0, 1\}^n$ has $\text{thr}(f) \leq r$ for all $f \in C$, then C can be learned in time $n^{O(r)}$ in either of two well-studied and demanding learning models, the Probably Approximately Correct (PAC) model of learning from random examples [15, 28] and the online model of learning from adversarially generated examples [2, 20]. Thus our polynomial threshold function upper bounds from Theorems 5.1 and 5.2 immediately give a range of new subexponential time learning results for various classes of Boolean formulas. For example, we immediately obtain:

Theorem 5.7 *The class of linear-size Boolean formulas of depth $o(\frac{\log n}{\log \log n})$ can be learned in time $2^{n^{1/2+\epsilon}}$ for all $\epsilon > 0$.*

This is the first subexponential time learning algorithm for linear size formulas of superconstant depth.

We emphasize that the PAC learning results which follow from our upper bounds hold for the general PAC model of learning from random examples which are drawn from an arbitrary

²We note in passing that an easy argument shows that any Sipser function of size s has a polynomial threshold function approximator of degree at most $s^{1/2}$; the proof is based on the observation that either the product of the odd-depth fanins or the even-depth fanins in any Sipser function must be at most $s^{1/2}$.

probability distribution over $\{0, 1\}^n$. This is in contrast with many results in learning theory (such as the quasipolynomial time algorithm of Linial *et al.* [19] for learning constant-depth circuits) which require the random examples to be drawn from the uniform distribution on $\{0, 1\}^n$.

References

- [1] A. Ambainis, A. Childs, B. Reichardt, R. Spalek, and S. Zhang. Any AND-OR formula of size n can be evaluated in time $n^{1/2+o(1)}$ on a quantum computer. In *Proc. 48th IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 363–372, 2007.
- [2] D. Angluin. Queries and concept learning. *Machine Learning*, 2:319–342, 1988.
- [3] J. Aspnes, R. Beigel, M. Furst, and S. Rudich. The expressive power of voting polynomials. *Combinatorica*, 14(2):1–14, 1994.
- [4] R. Beigel. The polynomial method in circuit complexity. In *Proceedings of the Eighth Conference on Structure in Complexity Theory*, pages 82–95, 1993.
- [5] R. Beigel. Perceptrons, PP, and the Polynomial Hierarchy. *Computational Complexity*, 4:339–349, 1994.
- [6] R. Beigel, N. Reingold, and D. Spielman. PP is closed under intersection. *Journal of Computer & System Sciences*, 50(2):191–202, 1995.
- [7] J. Bruck. Harmonic analysis of polynomial threshold functions. *SIAM Journal on Discrete Mathematics*, 3(2):168–177, 1990.
- [8] H. Buhrman, R. Cleve, and A. Wigderson. Quantum vs. classical communication and computation. In *Proceedings of the Thirtieth Annual ACM Symposium on Theory of Computing*, pages 63–68. ACM Press, 1998.
- [9] H. Buhrman and R. de Wolf. Complexity measures and decision tree complexity: a survey. *Theoretical Computer Science*, 288(1):21–43, 2002.
- [10] E. Cheney. *Introduction to approximation theory*. McGraw-Hill, New York, New York, 1966.
- [11] L. Comtet. *Advanced Combinatorics: The Art of Finite and Infinite Expansions*. Reidel, Dordrecht, Netherlands, 1974.

- [12] Y. Freund. Boosting a weak learning algorithm by majority. *Information and Computation*, 121(2):256–285, 1995.
- [13] M. Goldmann. On the power of a threshold gate at the top. *Information Processing Letters*, 63(6):287–293, 1997.
- [14] 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.
- [15] M. Kearns and U. Vazirani. *An introduction to computational learning theory*. MIT Press, Cambridge, MA, 1994.
- [16] A. Klivans, R. O’Donnell, and R. Servedio. Learning intersections and thresholds of halfspaces. In *Proceedings of the 43rd Annual Symposium on Foundations of Computer Science*, pages 177–186, 2002.
- [17] 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.
- [18] M. Krause and P. Pudlak. Computing boolean functions by polynomials and threshold circuits. *Computational Complexity*, 7(4):346–370, 1998.
- [19] N. Linial, Y. Mansour, and N. Nisan. Constant depth circuits, Fourier transform and learnability. *Journal of the ACM*, 40(3):607–620, 1993.
- [20] N. Littlestone. Learning quickly when irrelevant attributes abound: a new linear-threshold algorithm. *Machine Learning*, 2:285–318, 1988.
- [21] M. Minsky and S. Papert. *Perceptrons: an introduction to computational geometry (expanded edition)*. MIT Press, Cambridge, MA, 1988.
- [22] D. J. Newman. Rational approximation to $|x|$. *Michigan Mathematical Journal*, 11:11–14, 1964.
- [23] 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.
- [24] R. O’Donnell and R. Servedio. New degree bounds for polynomial threshold functions. In *Proceedings of the 35th ACM Symposium on Theory of Computing*, pages 325–334, 2003.

- [25] R. Paturi and M. Saks. Approximating threshold circuits by rational functions. *Information and Computation*, 112(2):257–272, 1994.
- [26] M. Saks. *Slicing the hypercube*, pages 211–257. London Mathematical Society Lecture Note Series 187, 1993.
- [27] D. Sieling. Minimization of decision trees is hard to approximate. Technical Report ECCC Report TR02-054, Electronic colloquium on computational complexity, 2002.
- [28] L. Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984.
- [29] N. Vereshchagin. Lower bounds for perceptrons solving some separation problems and oracle separation of AM from PP. In *Proceedings of the Third Annual Israel Symposium on Theory of Computing and Systems*, 1995.

A A new proof of the one-in-a-box lower bound

Recall that the DNF version f of the one-in-a-box function is a read-once DNF (OR of ANDs) in which there are m ANDs (terms) each with fanin $4m^2$. Minsky and Papert [21] showed that f requires polynomial threshold degree m ; we now reprove this using our lower bound technique.

We begin by performing some of the same steps as in Section 4. Let $[4m^2]$ denote the set $\{0, 1, 2, \dots, 4m^2\}$. By symmetrization, it suffices to prove a lower bound of m for the degree of any m -variate polynomial over \mathbf{R} which sign-represents the function $g : [4m^2]^m \rightarrow \{-1, 1\}$, $g(x_1, \dots, x_m) = -1$ iff at least one of the x_i 's is 0. By the Theorem of the Alternative (as in Section 4), we can do this by constructing a distribution \mathbf{w} on $[4m^2]^m$ under which g has zero correlation with every m -variate monomial of degree at most $m - 1$.

Let $\mathbf{x}_t \in [4m^2]^m$ denote the point $((t - 1)^2, (t - 3)^2, (t - 5)^2, \dots, (t - (2m - 1))^2)$. The support of \mathbf{w} will be the following $2m + 1$ points: $\{\mathbf{x}_t : t = 0 \dots 2m\}$. (This is the same set of points Minsky and Papert considered.) The weight \mathbf{w} gives to the point \mathbf{x}_t will be $\binom{2m}{t}$. Notice that $g(\mathbf{x}_t) = 1$ iff t is odd. Therefore, to show that g has zero correlation with every monomial of degree at most $m - 1$ under \mathbf{w} , we must show that for all $0 \leq a_1 + a_2 + \dots + a_m < m$:

$$\sum_{t=0}^{2m} (-1)^t \binom{2m}{t} \prod_{i=1}^m (t - (2i - 1))^{2a_i} = 0.$$

In fact, if $h(t)$ is any univariate polynomial of degree smaller than $2m$ then:

$$\sum_{t=0}^{2m} (-1)^t \binom{2m}{t} h(t) = 0. \tag{8}$$

This follows immediately from the following well-known combinatorial identity: For all $0 \leq c < 2m$:

$$\sum_{t=0}^{2m} (-1)^t \binom{2m}{t} t^c = 0. \quad (9)$$

(To prove this identity, write $(x-1)^{2m} = \sum_{t=0}^{2m} (-1)^t \binom{2m}{t} x^t$ by the Binomial Theorem. Substitute $x=1$ to get (9) for $c=0$. Now differentiate and substitute $x=1$ to get (9) for $c=1$. Differentiate again and substitute $x=1$ to get (8) with $h(t) = t(t-1)$; by linear combination with (9) for $c=1$ we get (9) for $c=2$. Repeatedly differentiate and substitute $x=1$; this yields (8) with $h(t) = t(t-1)(t-2), t(t-1)(t-2)(t-3)$, etc., which gives (9) for $c=2, 3$, etc. by linear combination with previously derived identities. The whole process may be repeated $2m-1$ times.)

B Proof of Lemma 5.3

The following convention will be useful for this section: for P a polynomial we write “ $P(x) \in_f ([a, b], [c, d])$ ” as shorthand for

$$“\forall x \in \{0, 1\}^n : \text{if } f(x) = 0 \text{ then } P(x) \in [a, b] \text{ and if } f(x) = 1 \text{ then } P(x) \in [c, d].”$$

Proof of Lemma 5.3: By assumption we have $p_i(x) \in_{f_i} ([-\epsilon, \epsilon], [1-\epsilon, 1+\epsilon])$ for each i .

Let $P(x)$ denote $p_1(x) + \dots + p_\ell(x) + \ell\epsilon$. It is easy to verify that we have

$$P(x) \in_f ([0, \ell-1+2\ell\epsilon], [\ell, \ell+2\ell\epsilon]).$$

Let $Q(x)$ denote $P(x)/(\ell-1+2\ell\epsilon)$. We then have

$$Q(x) \in_f ([0, 1], [1 + \frac{1-2\ell\epsilon}{\ell-1+2\ell\epsilon}, 1 + \frac{1}{\ell-1+2\ell\epsilon}]).$$

Let $k = \frac{1-2\ell\epsilon}{\ell-1+2\ell\epsilon}$. We can rewrite and say $Q(x) \in_f ([0, 1], [1+k, 1+k + \frac{2\ell\epsilon}{\ell-1+2\ell\epsilon}])$. Since $\frac{2\ell\epsilon}{\ell-1+2\ell\epsilon} < \frac{2\ell\epsilon}{\ell-1} \leq 4\epsilon$ we have $Q(x) \in_f ([0, 1], [1+k, 1+k+4\epsilon])$.

Recall that the Chebyshev polynomial of the first kind $C_d(t)$ is a univariate polynomial of degree d . The following fact is proved later:

Fact B.1 *For all $d \geq 1$ we have:*

1. $C_d(t) \in [-1, 1]$ for $t \in [0, 1]$.
2. Let t_d denote $C_{\lceil \sqrt{d} \rceil}(1+1/d)$. Then $t_d > 2$.
3. For all $0 < \tau < \frac{1}{d}$ we have $C_{\lceil \sqrt{d} \rceil}(1+1/d+\tau) \in [t_d, t_d+26d\tau]$.

Let $R(x)$ denote $C_{\lceil k^{-1/2} \rceil}(Q(x))$. Since $4\epsilon < \frac{1}{2\ell} < k$, by parts 1 and 3 of Fact B.1 we have that $R(x) \in_f([-1, 1], [t_k, t_k + \frac{104\epsilon}{k}])$. Let $S(x)$ denote $(\frac{1}{t_k}R(x))^{\lceil \log \frac{1}{\epsilon} \rceil}$. Using part 2 of Fact B.1 we find that $S(x) \in_f([-\epsilon, \epsilon], [1, (1 + \frac{104\epsilon}{t_k k})^{\lceil \log \frac{1}{\epsilon} \rceil}])$. We now recall the fact that $(1 + \alpha)^r \leq 1 + 2\alpha r$ for all $\alpha, r \geq 0$ such that $\alpha r \leq 1/2$:

$$(1 + \alpha)^r = 1 + \sum_{i=1}^r \alpha^i \binom{r}{i} \leq 1 + \sum_{i=1}^r (\alpha r)^i \leq 1 + \alpha r \sum_{i=0}^{\infty} (\alpha r)^i \leq 1 + 2\alpha r.$$

Using this fact, we find that

$$\left(1 + \frac{104\epsilon}{t_k k}\right)^{\lceil \log \frac{1}{\epsilon} \rceil} \leq 1 + \frac{416 \log \frac{1}{\epsilon}}{t_k k} \epsilon.$$

Using our bounds on t_k and k , this is at most $1 + (c_2 \ell \log \frac{1}{\epsilon}) \epsilon$ as desired.

It remains only to bound $\deg(S)$. From our construction it is clear that $\deg(S) \leq r \cdot \lceil k^{-1/2} \rceil \cdot \lceil \log \frac{1}{\epsilon} \rceil$. We have that $\lceil k^{-1/2} \rceil \leq \lceil \sqrt{2\ell} \rceil \leq 2\sqrt{\ell}$ and $\lceil \log \frac{1}{\epsilon} \rceil < 2 \log \frac{1}{\epsilon}$. Thus $\deg(S) \leq 4r\sqrt{\ell} \log \frac{1}{\epsilon}$ and the lemma is proved. \blacksquare

Proof of Fact B.1: Part (1) is one of the most basic facts about Chebyshev polynomials (see [10]). Part (2) follows from the fact that $C_{\lceil \sqrt{d} \rceil}$ is convex on $[1, \infty)$ and has slope $\lceil \sqrt{d} \rceil^2 \geq d$ at 1 (see [10] or [17]).

For Part (3), since $C_{\lceil \sqrt{d} \rceil}$ is convex and increasing on $[1, \infty)$ we have that

$$t_d \leq C_{\lceil \sqrt{d} \rceil}(1 + 1/d + \tau) < t_d + \frac{\tau}{1/d} \left(C_{\lceil \sqrt{d} \rceil}(1 + 2/d) - t_d \right).$$

Thus it suffices to show that $C_{\lceil \sqrt{d} \rceil}(1 + 2/d) - t_d < 26$. To see this, we recall that $C_r(x)$ can be defined as $C_r(x) = \cosh(r \cdot \operatorname{acosh} x)$ for $|x| > 1$ (see [10]). The Taylor series expansion of $\operatorname{acosh} x$ about $x = 1$ shows that $\operatorname{acosh}(1 + \epsilon) < \sqrt{2\epsilon}$ for all $\epsilon > 0$. Thus we have that

$$\lceil \sqrt{d} \rceil \cdot \operatorname{acosh}(1 + 2/d) < \lceil \sqrt{d} \rceil \cdot \sqrt{4/d} \leq 4.$$

Hence $C_{\lceil \sqrt{d} \rceil}(1 + 2/d) \leq \cosh 4 < 28$. Since $t_d > 2$ we have $C_{\lceil \sqrt{d} \rceil}(1 + 2/d) - t_d < 26$ as desired, and Fact B.1 is proved.

C Proof of Lemma 5.6

A key tool in the proof of Lemma 5.6 is the fact that low-degree rational functions can accurately approximate the sgn function. Building on work of Newman [22] and Paturi and Saks [25], in [6] Beigel *et al.* showed the following:

Fact C.1 *Let $k \geq 1, \epsilon > 0$. There is a rational function $r_{k,\epsilon}$ of degree $O(k \log \frac{1}{\epsilon})$ such that*

- $r_{k,\epsilon}(x) \in [-1 - \epsilon, -1]$ for all $x \in [-2^k, -1]$;
- $r_{k,\epsilon}(x) \in [1, 1 + \epsilon]$ for all $x \in [1, 2^k]$.

(We note in passing that the $O(\log n)$ upper bound for polynomial threshold degree of an AND of two n -variable majorities given by Beigel *et al.* is an easy consequence of this fact.)

Proof of Lemma 5.6: The base case $d = 1$ is easy. Without loss of generality we have that f is a conjunction $f = x_1 \wedge \dots \wedge x_s$. The rational function $(r_{\log(2s), 1/4s}(2(x_1 + \dots + x_s - s + \frac{1}{2}))) + 1)/2$ is easily seen to satisfy the conditions of Lemma 5.6.

For the induction step, without loss of generality we may suppose that f is computed by a Boolean formula $F = \bigvee_{i=1}^{\ell} F_i$ where $\ell \geq 2$, each F_i has depth at most $d - 1$, and the sum of the sizes $|F_1|, \dots, |F_{\ell}|$ is s . As in Lemma 5.5, for $j = 1, \dots, \lceil \log s \rceil$ let C_j be the set of those F_i such that $2^{j-1} \leq |F_i| < 2^j$. Let $H_j = \bigvee_{i:F_i \in C_j} F_i$ (note that unlike Lemma 5.5 now the subformulas of H_j are F_i 's rather than G_i 's), so f is computed by $\bigvee_{j=1}^{\lceil \log s \rceil} H_j$. We write h_j and f_i to denote the Boolean functions computed by formulas H_j and F_i respectively.

Let $J = s^{1 - \frac{1}{2^{d-1}}}$. We will deal with the H_j 's in different ways depending on whether $2^j < J$ or $2^j \geq J$.

We first consider j such that $2^j < J$. By a minor modification of Theorem 5.1 we have that for each $F_i \in C_j$, there is a polynomial p_i such that $\deg(p_i) \leq (c_1)^{d-1} (\log s)^{5(d-1)/2} \sqrt{J}$ and $|p_i - f_i|_{\infty} \leq \frac{1}{4s}$. Let $P_j(x) = 4(\sum_{i:F_i \in C_j} p_i(x) - \frac{1}{2})$. Then we have $P_j(x) \in_{h_j}([-3, -1], [1, 4s])$, and hence

$$Q_j(x) \stackrel{\text{def}}{=} \frac{r_{\log(4s), 1/4s}(P_j(x)) + 1}{2} \in_{h_j}([-1/4s, 0], [1, 1 + 1/4s]),$$

where $\deg(Q_j) = O((c_1)^{d-1} (\log s)^{5(d-1)/2 + 2} \sqrt{J})$.

We now consider j such that $2^j \geq J$, so each $F_i \in C_j$ satisfies $|F_i| \geq J/2$. By the induction hypothesis, we have that for each $F_i \in C_j$ there is a rational function $t_i(x)$ such that

$$\begin{aligned} \deg(t_i) &\leq (c_4)^{d-1} (\log |F_i|)^{5(d-1)} |F_i|^{\frac{1}{2} - \frac{1}{2^{d-2}}} \\ &\leq (c_4)^{d-1} (\log s)^{5(d-1)} 2^{j(\frac{1}{2} - \frac{1}{2^{d-2}})} \end{aligned} \tag{10}$$

and $|t_i - f_i|_{\infty} \leq \frac{1}{4|F_i|} \leq \frac{1}{2J}$. Let $T_j(x) = 4(\sum_{i:F_i \in C_j} t_i(x) - \frac{1}{2})$. Since C_j contains at most $s/2^{j-1} \leq 2s/J$ formulas F_i , we have that

$$\sum_{i:F_i \in C_j} t_i \in_{h_j}([-s/J^2, s/J^2], [1 - s/J^2, s]).$$

Since $s/J^2 = s^{-1 + \frac{2}{2^{d-1}}} \leq s^{-1/3}$ for $d \geq 2$, we may suppose that $s/J^2 \leq \frac{1}{4}$. Consequently, we have that $T_j \in_{h_j}([-3, -1], [1, 4s])$. Since $T_j(x)$ is a sum of at most $s/2^{j-1}$ rational functions

t_i whose degrees are bounded by (10), by clearing denominators we can express T_j as a rational function of degree $O((c_4)^{d-1}(\log s)^{5(d-1)}s/2^{j(\frac{1}{2}+\frac{1}{2^{d-2}})})$. Now observe that

$$\frac{s}{2^{j(\frac{1}{2}+\frac{1}{2^{d-2}})}} \leq \frac{s}{J^{\frac{1}{2}+\frac{1}{2^{d-2}}}} = \frac{s}{s^{(1-\frac{1}{2^{d-1}})(\frac{1}{2}+\frac{1}{2^{d-2}})}} = s^{\frac{1}{2}-\frac{1}{2^{d+1-2}}} = \sqrt{J},$$

and hence $\deg(T_j) = O((c_4)^{d-1}(\log s)^{5(d-1)}\sqrt{J})$. Thus, we have

$$U_j(x) \stackrel{\text{def}}{=} \frac{r_{\log(4s), 1/4s}(S_j(x)) + 1}{2} \in_{h_j}([-1/4s, 0], [1, 1 + 1/4s]),$$

where U_j is a rational function with $\deg(U_j) = O((c_4)^{d-1}(\log s)^{5d-3}\sqrt{J})$.

Now let

$$V(x) = 4 \left(\sum_{j:2^j < J} Q_j(x) + \sum_{j:2^j \geq J} U_j - 1/2 \right).$$

Since $V(x)$ is a sum of $O(\log s)$ rational functions Q_j, T_j , by clearing denominators we have that $V(x)$ is a rational function of degree $O((c_4)^{d-1}(\log s)^{5d-2}\sqrt{J})$, and moreover $V(x) \in_f([-3, -1], [1, 4s])$. Finally, taking

$$r(x) \stackrel{\text{def}}{=} \frac{r_{\log(4s), 1/4s}(V(x)) + 1}{2},$$

we have that $|r - f|_\infty \leq \frac{1}{4s}$ and $\deg(r) \leq (c_4)^d(\log s)^{5d}\sqrt{J}$. Since $\sqrt{J} = s^{\frac{1}{2}-\frac{1}{2^{d+1-2}}}$, the lemma is proved. (Lemma 5.6) ■