

Mansour's Conjecture is True for Random DNF Formulas

Adam Klivans
University of Texas at Austin
klivans@cs.utexas.edu

Homin K. Lee
University of Texas at Austin
homin@cs.utexas.edu

Andrew Wan
Columbia University
atw12@cs.columbia.edu

April 2, 2010

Abstract

In 1994, Y. Mansour conjectured that for every DNF formula on n variables with t terms there exists a polynomial p with $t^{O(\log(1/\epsilon))}$ non-zero coefficients such that $\mathbf{E}_{x \in \{0,1\}^n} [(p(x) - f(x))^2] \leq \epsilon$. We make the first progress on this conjecture and show that it is true for randomly chosen DNF formulas and read-once DNF formulas.

Our result yields the first polynomial-time query algorithm for agnostically learning these subclasses of DNF formulas with respect to the uniform distribution on $\{0, 1\}^n$ (for any constant error parameter).

Applying recent work on sandwiching polynomials, our results imply that a $t^{-O(\log 1/\epsilon)}$ -biased distribution fools the above subclasses of DNF formulas. This gives pseudorandom generators for randomly chosen DNF with shorter seed length than all previous work.

1 Introduction

Let $f : \{0, 1\}^n \rightarrow \{0, 1\}$ be a DNF formula, *i.e.*, a function of the form $T_1 \vee \dots \vee T_t$ where each T_i is a conjunction of at most n literals. In this paper we are concerned with the following question: how well can a real-valued polynomial p approximate the Boolean function f ? This is an important problem in computational learning theory, as real-valued polynomials play a critical role in developing learning algorithms for DNF formulas.

Over the last twenty years, considerable work has gone into finding polynomials p with certain properties (*e.g.*, low-degree, sparse) such that

$$\mathbf{E}_{x \in \{0,1\}^n} [(p(x) - f(x))^2] \leq \epsilon.$$

In 1989, Linal *et al.* [LMN93] were the first to prove that for any t -term DNF formula f , there exists a polynomial $p : \{0, 1\}^n \rightarrow \mathbb{R}$ of degree $O(\log(t/\epsilon)^2)$ such that $\mathbf{E}_{x \in \{0,1\}^n} [(p(x) - f(x))^2] \leq \epsilon$. They showed that this type of approximation implies a quasipolynomial-time algorithm for PAC learning DNF formulas with respect to the uniform distribution. Kalai *et al.* [KKMS08] observed that this fact actually implies something stronger, namely a quasipolynomial-time agnostic learning algorithm for learning DNF formulas (with respect to the uniform distribution). Additionally, the above approximation was used in recent work due to Bazzi [Baz07] and Razborov [Raz08] to show that bounded independence fools DNF formulas.

Three years later, building on the work of Linal *et al.* Mansour [Man95] proved that for any DNF formula with t terms, there exists a polynomial p defined over $\{0, 1\}^n$ with *sparsity* $t^{O(\log \log t \log(1/\epsilon))}$ such that $\mathbf{E}_{x \in \{0,1\}^n} [(p(x) - f(x))^2] \leq \epsilon$. By sparsity we mean the number of non-zero Fourier coefficients of p .

This result implied a nearly polynomial-time *query* algorithm for PAC learning DNF formulas with respect to the uniform distribution.

Mansour conjectured [Man94] that the bound above could be improved to $t^{O(\log 1/\epsilon)}$. Such an improvement would imply a polynomial-time query algorithm for learning DNF formulas with respect to the uniform distribution (to within any constant accuracy), and learning DNF formulas in this model was a major open problem at that time.

In a celebrated work from 1994, Jeff Jackson proved that DNF formulas were learnable in polynomial time (with queries, with respect to the uniform distribution) *without* proving the Mansour conjecture. His “Harmonic Sieve” algorithm [Jac97] used boosting in combination with some weak approximation properties of polynomials. As such, for several years, Mansour’s conjecture remained open and attracted considerable interest, but its resolution did not imply any new results in learning theory.

In 2008, Gopalan *et al.* [GKK08b] proved that a positive resolution to the Mansour conjecture also implies an efficient query algorithm for *agnostically* learning DNF formulas (to within any constant error parameter). The agnostic model of learning is a challenging learning scenario that requires the learner to succeed in the presence of adversarial noise. Roughly, Gopalan *et al.* showed that if a class of Boolean functions \mathcal{C} can be ϵ -approximated by polynomials of sparsity s , then there is a query algorithm for agnostically learning \mathcal{C} in time $\text{poly}(s, 1/\epsilon)$ (since decision trees are approximated by sparse polynomials, they obtained the first query algorithm for agnostically learning decision trees with respect to the uniform distribution on $\{0, 1\}^n$). Whether DNF formulas can be agnostically learned (with queries, with respect to the uniform distribution) still remains a difficult open problem [GKK08a].

1.1 Our Results

We prove that the Mansour conjecture is true for read-once and randomly chosen DNF formulas. As far as we know, prior to this work, the Mansour conjecture was not known to be true for any interesting class of DNF formulas.

Theorem 1. *Let $f : \{0, 1\}^n \rightarrow \{0, 1\}$ be a DNF formula with t terms where each literal appears at most once. Then there exists a polynomial p with sparsity $t^{O(\log 1/\epsilon)}$ such that $\mathbf{E}[(p(x) - f(x))^2] \leq \epsilon$.*

Mansour [Man95] proves that any polynomial that approximates read-once DNF formulas to ϵ accuracy must have *degree* at least $\Omega(\log t \log(1/\epsilon) / \log \log(1/\epsilon))$. He further shows that a “low-degree” strategy of selecting all of a DNF’s Fourier coefficients of monomials up to degree d results in a polynomial p with sparsity $t^{O(\log \log t \log 1/\epsilon)}$. It is not clear, however, how to improve this to the desired $t^{O(\log 1/\epsilon)}$ bound.

Our next result shows that the Mansour conjecture is true for the class of randomly chosen DNF formulas:

Theorem 2. *Let $f : \{0, 1\}^n \rightarrow \{0, 1\}$ be a DNF formula with t terms where each term is chosen independently from the set of all terms of length $\log t$. Then with probability $1 - n^{-\Omega(\log t)}$ (over the choice of the DNF formula), there exists a polynomial p with sparsity $t^{O(\log 1/\epsilon)}$ such that $\mathbf{E}[(p(x) - f(x))^2] \leq \epsilon$.*

As mentioned earlier, by applying the result of Gopalan *et al.* [GKK08b], we obtain the first polynomial-time query algorithms for agnostically learning the above classes of DNF formulas to within any constant accuracy parameter. We consider this an important step towards agnostically learning all DNF formulas.

Corollary 3. *Let \mathcal{C} be the class of DNF formulas with t terms where each term is randomly chosen from the set of all terms of length $\log t$. Then there is a query-algorithm for agnostically learning \mathcal{C} with respect to the uniform distribution on $\{0, 1\}^n$ to accuracy ϵ in time $\text{poly}(n) \cdot t^{O(\log 1/\epsilon)}$ with probability $1 - n^{-\Omega(\log t)}$ (over the choice of the DNF formula).*

We define the notion of agnostic learning with respect to randomly chosen concept classes in Section 2. We also obtain a corresponding agnostic learning algorithm for read-once DNF formulas:

Corollary 4. *Let \mathcal{C} be the class of read-once DNF formulas with t terms. Then there is a query-algorithm for agnostically learning \mathcal{C} with respect to the uniform distribution on $\{0, 1\}^n$ to accuracy ϵ in time $\text{poly}(n) \cdot t^{O(\log 1/\epsilon)}$.*

Our sparse polynomial approximators can also be used in conjunction with recent work due to De. *et al.* to show that for any randomly chosen DNF f , a $1/t^{O(\log 1/\epsilon)}$ -biased distribution fools f :

Theorem 5. *Let f be a randomly chosen DNF formula. Then there exists a pseudorandom generator G such that*

$$\left| \Pr_{x \in \{0,1\}^s} [f(G(x)) = 1] - \Pr_{z \in \{0,1\}^n} [f(z) = 1] \right| \leq \epsilon$$

with $s = O(\log n + \log t \cdot \log(1/\epsilon))$.

Previously it was only known that these types of biased distributions fool read-once DNF formulas [DETT09].

1.2 Related Work

As mentioned earlier, Mansour, using the random restriction machinery of Håstad and Linial *et al.* [Hås86, LMN93] had shown that for any DNF formula f , there exists a p of sparsity $t^{O(\log \log t \log 1/\epsilon)}$ that approximates f .

The subclasses of DNF formulas that we show are agnostically learnable have been well-studied in the PAC model of learning. Read-once DNF formulas were shown to be PAC-learnable with respect to the uniform distribution by Kearns *et al.* [KLPV87] and random DNF formulas were recently shown to be learnable on average with respect to the uniform distribution in the following sequence of work [JS05, JLSW08, Sel08, Sel09].

Recently (and independently) De *et al.* proved that for any read-once DNF formula f , there exists an approximating polynomial p of sparsity $t^{O(\log 1/\epsilon)}$. More specifically, De *et al.* showed that for any class of functions \mathcal{C} fooled by δ -biased sets, there exist sparse, sandwiching polynomials for \mathcal{C} where the sparsity depends on δ . Since they show that $1/t^{O(\log 1/\epsilon)}$ -biased sets fool read-once DNF formulas, the existence of a sparse approximator for the read-once case is implicit in their work.

1.3 Our Approach

As stated above, our proof does not analyze the Fourier coefficients of DNF formulas, and our approach is considerably simpler than the random-restriction method taken by Mansour (we consider the lack of Fourier analysis a feature of the proof, given that all previous work on this problem has been Fourier-based). Instead, we use interpolation to construct an approximating polynomial *directly*.

Consider a DNF formula $f = T_1 \vee \dots \vee T_t$ where each T_i is on a disjoint set of exactly $\log t$ variables. The probability that each term is satisfied is $1/t$, and the expected number of satisfied terms is one. Further, since the terms are disjoint, with high probability over the choice of random input, only a few—say d —terms will be satisfied. As such, we construct a univariate polynomial p with $p(0) = 0$ and $p(i) = 1$ for $1 \leq i \leq d$. Then $p(T_1 + \dots + T_t)$ will be exactly equal to f as long as at most d terms are satisfied. A careful calculation shows that the inputs where p is incorrect will not contribute too much to $\mathbf{E}[(f - p)^2]$,

as there are few of them. Setting parameters appropriately yields a polynomial p that is both sparse and an ϵ -approximator of f .

More generally, we adopt the following strategy: given a DNF formula f (randomly chosen or read-once) either (1) with sufficiently high probability a random input does not satisfy too many terms of f or (2) f is highly biased. In the former case we can use polynomial interpolation to construct a sparse approximator and in the latter case we can simply use the constant 0 or 1 function.

The probability calculations are a bit delicate, as we must ensure that the probability of many terms being satisfied decays faster than the growth rate of our polynomial approximators. For the case of random DNF formulas, we make use of some recent work due to Jackson *et al.* on learning random monotone DNF formulas [JLSW08].

2 Preliminaries

In this paper, we will primarily be concerned with Boolean functions $f : \{0, 1\}^n \rightarrow \{0, 1\}$. Let x_1, \dots, x_n be Boolean variables. A *literal* is either a variable x_i or its negation \bar{x}_i , and a *term* is a conjunction of literals. Any Boolean function can be expressed as a disjunction of terms, and such a formula is said to be a *disjunctive normal form* (or DNF) formula. A read-once DNF formula is a DNF formula in which each variable occurs at most once.

2.1 Sparse Polynomials

Every function $f : \{0, 1\}^n \rightarrow \mathbb{R}$ can be expressed by its Fourier expansion: $f(x) = \sum_S \hat{f}(S) \chi_S(x)$ where $\chi_S(x) = \prod_{i \in S} (-1)^{x_i}$ for $S \subseteq [n]$, and $\hat{f}(S) = \mathbf{E}[f \cdot \chi_S]$. The Fourier expansion of f can be thought of as the unique polynomial representation of f over $\{+1, -1\}^n$ under the map $x_i \mapsto 1 - 2x_i$.

Mansour conjectured that polynomial-size DNF formulas could be approximated by *sparse* polynomials over $\{+1, -1\}^n$. We say a polynomial $p : \{+1, -1\}^n \rightarrow \mathbb{R}$ has sparsity s if it has at most s non-zero coefficients. We state Mansour's conjecture as originally posed in [Man94], which uses the convention of representing FALSE by $+1$ and TRUE by -1 .

Conjecture 6 ([Man94]). *Let $f : \{+1, -1\}^n \rightarrow \{+1, -1\}$ be any function computable by a t -term DNF formula. Then there exists a polynomial $p : \{+1, -1\}^n \rightarrow \mathbb{R}$ with $t^{O(\log 1/\epsilon)}$ terms such that $\mathbf{E}[(f-p)^2] \leq \epsilon$.*

We will prove the conjecture to be true for various subclasses of polynomial-size DNF formulas. In our setting, Boolean functions will output 0 for FALSE and 1 for TRUE. However, we can easily change the range by setting $f^\pm := 1 - 2 \cdot f$. Changing the range to $\{+1, -1\}$ changes the accuracy of the approximation by at most a factor of 4: $\mathbf{E}[(1 - 2f) - (1 - 2p)]^2 = 4 \mathbf{E}[(f - p)^2]$, and it increases the sparsity by at most 1.

Given a Boolean function f , we construct a sparse approximating polynomial over $\{+1, -1\}^n$ by giving an approximating polynomial $p : \{0, 1\}^n \rightarrow \mathbb{R}$ with real coefficients that has small spectral norm. The rest of the section gives us some tools to construct such polynomials and explains why doing so yields sparse approximators.

Definition 7. *The Fourier ℓ_1 -norm (also called the spectral norm) of a function $p : \{0, 1\}^n \rightarrow \mathbb{R}$ is defined to be $\|p\|_1 := \sum_S |\hat{p}(S)|$. We will also use the following minor variant, $\|p\|_1^{\neq \emptyset} := \sum_{S \neq \emptyset} |\hat{p}(S)|$.*

The following two facts about the spectral norm of functions will allow us to construct polynomials over $\{0, 1\}^n$ naturally from DNF formulas.

Fact 8. Let $p : \{0, 1\}^m \rightarrow \mathbb{R}$ be a polynomial with coefficients $p_S \in \mathbb{R}$ for $S \subseteq [m]$, and $q_1, \dots, q_m : \{0, 1\}^n \rightarrow \{0, 1\}$ be arbitrary Boolean functions. Then $p(q_1, \dots, q_m) = \sum_S p_S \prod_{i \in S} q_i$ is a polynomial over $\{0, 1\}^n$ with spectral norm at most

$$\sum_{S \subseteq [m]} |p_S| \prod_{i \in S} \|q_i\|_1.$$

Proof. The fact follows by observing that for any $p, q : \{0, 1\}^n \rightarrow \mathbb{R}$, we have $\|p + q\|_1 \leq \|p\|_1 + \|q\|_1$ and $\|pq\|_1 \leq \|p\|_1 \|q\|_1$. ■

Fact 9. Let $T : \{0, 1\}^n \rightarrow \{0, 1\}$ be an AND of a subset of its literals. Then $\|T\|_1 = 1$.

Finally, using the fact below, we show why approximating polynomials with small spectral norm give sparse approximating polynomials.

Fact 10 ([KM93]). Given any function $f : \{0, 1\}^n \rightarrow \mathbb{R}$ and $\epsilon > 0$, let $\mathcal{S} = \{S \subseteq [n] : |\hat{f}(S)| \geq \epsilon / \|f\|_1\}$, and $g(x) = \sum_{S \in \mathcal{S}} \hat{f}(S) \chi_S(x)$. Then $\mathbf{E}[(f - g)^2] \leq \epsilon$, and $|\mathcal{S}| \leq \|f\|_1^2 / \epsilon$.

Now, given functions $f, p : \{0, 1\}^n \rightarrow \mathbb{R}$ such that $E[(f - p)^2] \leq \epsilon$, we may construct a 4ϵ -approximator for f with sparsity $\|p\|_1^2 / \epsilon$ by defining $p'(x) = \sum_{S \in \mathcal{S}} \hat{p}(S) \chi_S(x)$ as in Fact 10. Clearly p' has sparsity $\|p\|_1^2 / \epsilon$, and

$$E[(f - p')^2] = E[(f - p + p - p')^2] \leq E[2((f - p)^2 + (p - p')^2)] \leq 4\epsilon,$$

where the first inequality follows from the inequality $(a + b)^2 \leq 2(a^2 + b^2)$ for any reals a and b .

2.2 Agnostic learning

We first describe the traditional framework for agnostically learning concept classes with respect to the uniform distribution and then give a slightly modified definition for an “average-case” version of agnostic learning where the unknown concept (in this case a DNF formula) is randomly chosen.

Definition 11 (Standard agnostic model). Let \mathcal{D} be the uniform distribution on $\{+1, -1\}^n$, and let $f : \{+1, -1\}^n \rightarrow \{+1, -1\}$ be an arbitrary function. Define

$$\text{opt} = \min_{c \in \mathcal{C}} \Pr_{x \sim \mathcal{D}} [c(x) \neq f(x)].$$

That is, opt is the error of the best fitting concept in \mathcal{C} with respect to \mathcal{D} . We say that an algorithm A agnostically learns \mathcal{C} with respect to \mathcal{D} if the following holds for any f : if A is given black-box access to f then with high probability A outputs a hypothesis h such that $\Pr_{x \sim \mathcal{D}} [h(x) \neq f(x)] \leq \text{opt} + \epsilon$.

The intuition behind the above definition is that a learner—given access to a concept $c \in \mathcal{C}$ where an η fraction of c 's inputs have been adversarially corrupted—should still be able to output a hypothesis with accuracy $\eta + \epsilon$ (achieving error better than η may not be possible, as the adversary could embed a completely random function on an η fraction of c 's inputs). Here η plays the role of opt .

This motivates the following definition for agnostically learning a randomly chosen concept from some class \mathcal{C} :

Definition 12 (Agnostically learning random concepts). *Let \mathcal{C} be a concept class and choose c randomly from \mathcal{C} . We say that an algorithm A agnostically learns random concepts from \mathcal{C} if with probability at least $1 - \delta$ over the choice of c the following holds: if the learner is given black-box access to c' and $\Pr_{x \in \{+1, -1\}^n} [c(x) \neq c'(x)] \leq \eta$, then A outputs a hypothesis h such that $\Pr_{x \in \{+1, -1\}^n} [h(x) \neq c'(x)] \leq \eta + \epsilon$.*

We are unaware of any prior work defining an agnostic framework for learning randomly chosen concepts.

The main result we use to connect the approximation of DNF formulas by sparse polynomials with agnostic learning is due to Gopalan *et al.* [GKK08b]:

Theorem 13 ([GKK08b]). *Let \mathcal{C} be a concept class such that for every $c \in \mathcal{C}$ there exists a polynomial p such that $\|p\|_1 \leq s$ and $\mathbf{E}_{x \in \{+1, -1\}^n} [|p(x) - c(x)|^2] \leq \epsilon^2/2$. Then there exists an algorithm B such that the following holds: given black-box access to any Boolean function $f : \{+1, -1\}^n \rightarrow \{+1, -1\}$, B runs in time $\text{poly}(n, s, 1/\epsilon)$ and outputs a hypothesis $h : \{+1, -1\}^n \rightarrow \{+1, -1\}$ with*

$$\Pr_{x \in \{+1, -1\}^n} [h(x) \neq f(x)] \leq \text{opt} + \epsilon.$$

3 Approximating DNFs using univariate polynomial interpolation

Let $f = T_1 \vee T_2 \vee \dots \vee T_t$ be any DNF formula. We say $T_i(x) = 1$ if x satisfies the term T_i , and 0 otherwise. Let $y_f : \{0, 1\}^n \rightarrow \{0, \dots, t\}$ be the function that outputs the number of terms of f satisfied by x , i.e., $y_f(x) = T_1(x) + T_2(x) + \dots + T_t(x)$.

Our constructions will use the following univariate polynomial P_d to interpolate the values of f on inputs $\{x : y_f(x) \leq d\}$.

Fact 14. *Let*

$$P_d(y) := (-1)^{d+1} \frac{(y-1)(y-2) \cdots (y-d)}{d!} + 1. \quad (1)$$

Then, (1) the polynomial P_d is a degree- d polynomial in y ; (2) $P_d(0) = 0$, $P_d(y) = 1$ for $y \in [d]$, and for $y \in [t] \setminus [d]$, $P_d(y) = -\binom{y-1}{d} + 1 \leq 0$ if d is even and $P_d(y) = \binom{y-1}{d} + 1 > 1$ if d is odd; and (3) the sum of the magnitudes of P_d 's coefficients is d .

Proof. Properties (1) and (2) can be easily verified by inspection. Expanding the falling factorial, we get that $(y-1)(y-2) \cdots (y-d) = \sum_{j=0}^d (-1)^{d-j} \begin{bmatrix} d+1 \\ j+1 \end{bmatrix} y^j$, where $\begin{bmatrix} a \\ b \end{bmatrix}$ denotes a Stirling number of the first kind. The Stirling numbers of the first kind count the number of permutations of a elements with b disjoint cycles. Therefore, $\sum_{j=0}^d \begin{bmatrix} d+1 \\ j+1 \end{bmatrix} = (d+1)!$ [GKP94]. The constant coefficient of P_d is 0 by Property (2), thus the sum of the absolute values of the other coefficients is $((d+1)! - d!)/d! = d$. ■

For any t -term DNF formula f , we can construct a polynomial $p_{f,d} : \{0, 1\}^n \rightarrow \mathbb{R}$ defined as $p_{f,d} := P_d \circ y_f$. A simple calculation, given below, shows that the ℓ_1 -norm of $p_{f,d}$ is polynomial in t and exponential in d .

Lemma 15. *Let f be a t -term DNF formula, then $\|p_{f,d}\|_1 \leq t^{O(d)}$.*

Proof. By Fact 14, P_d is a degree- d univariate polynomial with d non-zero coefficients of magnitude at most d . We can view the polynomial $p_{f,d}$ as the polynomial $P'_d(T_1, \dots, T_t) := P_d(T_1 + \dots + T_t)$ over variables $T_i \in \{0, 1\}$. Expanding out P'_d gives us at most dt^d monomials with coefficients of magnitude at most d . Now each monomial of P'_d is a product of T_i 's, so applying Facts 9 and 8 we have that $\|p_{f,d}\|_1 \leq t^{O(d)}$. ■

The next section will show that the polynomial $p_{f,d}$ (for $d = \Theta(\log 1/\epsilon)$) is in fact a good approximation for random DNF formulas. As a warm-up, we will show the simple case of read-once DNF formulas.

3.1 A Simple Case: Read-Once DNF Formulas

For read-once DNF formulas, the probability that a term is satisfied is independent of whether or not any of the other terms are satisfied, and thus f is unlikely to have many terms satisfied simultaneously.

Lemma 16. *Let $f = T_1 \vee \dots \vee T_t$ be a read-once DNF formula of size t such that $\Pr[f] < 1 - \epsilon$. Then the probability over the uniform distribution on $\{0, 1\}^n$ that some set of $j > e \ln 1/\epsilon$ terms is satisfied is at most $\left(\frac{e \ln 1/\epsilon}{j}\right)^j$.*

Proof. For any assignment x to the variables of f , let $y_f(x)$ be the number terms satisfied in f . By linearity of expectation, we have that $\mathbf{E}_x[y_f(x)] = \sum_{i=1}^t \Pr[T_i = 1]$. Note that $\Pr[\neg f] = \prod_{i=1}^t (1 - \Pr[T_i])$, which is maximized when each $\Pr[T_i] = \mathbf{E}[y_f]/t$, hence $\Pr[\neg f] \leq (1 - \mathbf{E}[y_f]/t)^t \leq e^{-\mathbf{E}[y_f]}$. Thus we may assume that $\mathbf{E}[y_f] \leq \ln 1/\epsilon$, otherwise $\Pr[f] \geq 1 - \epsilon$.

Assuming $\mathbf{E}[y_f] \leq \ln 1/\epsilon$, we now bound the probability that some set of $j > e \ln 1/\epsilon$ terms of f is satisfied. Since all the terms are disjoint, this probability is $\sum_{S \subseteq [t], |S|=j} \prod_{i \in S} \Pr[T_i]$, and the arithmetic-geometric mean inequality gives that this is maximized when every $\Pr[T_i] = \mathbf{E}[y_f]/t$. Then the probability of satisfying some set of j terms is at most:

$$\binom{t}{j} \left(\frac{\ln 1/\epsilon}{t}\right)^j \leq \left(\frac{et}{j}\right)^j \left(\frac{\ln 1/\epsilon}{t}\right)^j = \left(\frac{e \ln 1/\epsilon}{j}\right)^j,$$

which concludes the proof of the lemma. ■

The following lemma shows that we can set d to be fairly small, $\Theta(\log 1/\epsilon)$, and the polynomial $p_{f,d}$ will be a good approximation for any DNF formula f , as long as f is unlikely to have many terms satisfied simultaneously.

Lemma 17. *Let f be any t -term DNF formula, and let $d = 4e^3 \ln 1/\epsilon$. If*

$$\Pr[y_f(x) = j] \leq \left(\frac{e \ln 1/\epsilon}{j}\right)^j,$$

for every $d \leq j \leq t$, then the polynomial $p_{f,d}$ satisfies $\mathbf{E}[(f - p_{f,d})^2] \leq \epsilon$.

Proof. We condition on the values of $y_f(x)$, controlling the magnitude of $p_{f,d}$ by the unlikelihood of y_f being large. By Fact 14, $p_{f,d}(x)$ will output 0 if x does not satisfy f , $p_{f,d}(x)$ will output 1 if $y_f(x) \in [d]$, and $|p_{f,d}(x)| < \binom{y_f}{d}$ for $y_f(x) \in [t] \setminus [d]$. Hence:

$$\begin{aligned} \|f - p_{f,d}\|^2 &< \sum_{j=d+1}^t \binom{j}{d}^2 \left(\frac{e \ln 1/\epsilon}{j}\right)^j \\ &< \sum_{j=d+1}^t 2^{2j} \left(\frac{e \ln 1/\epsilon}{4e^3 \ln 1/\epsilon}\right)^j \\ &< \epsilon \sum_{j=d+1}^t \frac{1}{e^j} < \epsilon. \end{aligned}$$

■

Combining Lemmas 15, 16, and 17 gives us Mansour’s conjecture for read-once DNF formulas.

Theorem 18. *Let f be any read-once DNF formula with t terms. Then there is a polynomial $p_{f,d}$ with $\|p_{f,d}\|_1 \leq t^{O(\log 1/\epsilon)}$ and $\mathbf{E}[(f - p_{f,d})^2] \leq \epsilon$ for all $\epsilon > 0$.*

4 Mansour’s Conjecture for Random DNF Formulas

In this section, we establish various properties of random DNF formulas and use these properties to show that for almost all f , Mansour’s conjecture holds. Roughly speaking, we will show that a random DNF formula behaves like a read-once DNF formula, in that any “large” set of terms is unlikely to be satisfied by a random assignment. This notion is formalized in Lemma 21. For such DNF formulas, we may use the construction from Section 3 to obtain a good approximating polynomial for f with small spectral norm (Theorem 23).

Throughout the rest of this section we assume that $n^a \leq t(n) \leq n^b$ for any constants $a, b > 0$. For brevity we write t for $t(n)$. Let \mathcal{D}_n^t be the probability distribution over t -term DNF formulas induced by the following process: each term is independently and uniformly chosen at random from all $t \binom{n}{\log t}$ possible terms of size exactly $\log t$ over $\{x_1, \dots, x_n\}$.

If t grows very slowly relative to n , say $t = O(n^{1/4})$, then with high probability a random f drawn from \mathcal{D}_n^t will be a read-once DNF formula, in which case the results of Section 3.1 hold. If the terms are not of size $\Theta(\log n)$, then the DNF will be biased, and thus be easy to learn. We refer the reader to [JS05] for a full discussion of the model.

To prove Lemma 21, we require two Lemmas, which are inspired by the results of [JS05] and [JLSW08]. Lemma 19 shows that with high probability the terms of a random DNF formula are close to being disjoint, and thus cover close to $j \log t$ variables.

Lemma 19. *With probability at least $1 - t^j e^{j \log t} (j \log t)^{\log t} / n^{\log t}$ over the random draw of f from \mathcal{D}_n^t , at least $j \log t - (\log t)/4$ variables occur in every set of j distinct terms of f . The failure probability is at most $1/n^{\Omega(\log t)}$ for any $j < c \log n$, for some constant c .*

Proof. Let $k := \log t$. Fix a set of j terms, and let $v \leq jk$ be the number of distinct variables (negated or not) that occur in these terms. We will bound the probability that $v > w := jk - k/4$. Consider any particular fixed set of w variables. The probability that none of the j terms include any variable outside of the w variables is precisely $\left(\frac{\binom{w}{k}}{\binom{n}{k}}\right)^j$. Thus, the probability that $v \leq w$ is by the union bound:

$$\binom{n}{w} \left(\frac{\binom{w}{k}}{\binom{n}{k}}\right)^j < \left(\frac{en}{w}\right)^w \left(\frac{w}{n}\right)^{jk} = \frac{e^{jk-k/4} (jk - k/4)^{k/4}}{n^{k/4}} < \frac{e^{jk} (jk)^{k/4}}{n^{k/4}}.$$

Taking a union bound over all (at most t^j) sets, we have that with the correct probability every set of j terms contains at least w distinct variables. ■

We will use the method of bounded differences (a.k.a., McDiarmid’s inequality) to prove Lemma 21.

Proposition 20 (McDiarmid’s inequality). *Let X_1, \dots, X_m be independent random variables taking values in a set \mathcal{X} , and let $f : \mathcal{X}^m \rightarrow \mathbb{R}$ be such that for all $i \in [m]$, $|f(a) - f(a')| \leq d_i$, whenever $a, a' \in \mathcal{X}^m$ differ in just the i th coordinate. Then for all $\tau > 0$,*

$$\Pr[f > \mathbf{E}f + \tau] \leq \exp\left(-\frac{2\tau^2}{\sum_i d_i^2}\right) \text{ and } \Pr[f < \mathbf{E}f - \tau] \leq \exp\left(-\frac{2\tau^2}{\sum_i d_i^2}\right).$$

The following lemma shows that with high probability over the choice of random DNF formula, the probability that exactly j terms are satisfied is close to that for the “tribes” function: $\binom{t}{j}t^{-j}(1-1/t)^{t-j}$.

Lemma 21. *There exists a constant c such that for any $j < c \log n$, with probability at least $1 - 1/n^{\Omega(\log t)}$ over the random draw of f from \mathcal{D}_n^t , the probability over the uniform distribution on $\{0, 1\}^n$ that an input satisfies exactly j distinct terms of f is at most $2\binom{t}{j}t^{-j}(1-1/t)^{t-j}$.*

Proof. Let $f = T_1 \vee \dots \vee T_t$, and let $\beta := t^{-j}(1-1/t)^{t-j}$. Fix any $J \subset [t]$ of size j , and let U_J be the probability over $x \in \{0, 1\}^n$ that the terms T_i for $i \in J$ are satisfied and no other terms are satisfied. We will show that $U_J < 2\beta$ with high probability; a union bound over all possible sets J of size j in $[t]$ gives that $U_J \leq 2\beta$ for every J with high probability. Finally, a union bound over all $\binom{t}{j}$ possible sets of j terms (where the probability is taken over x) proves the lemma.

Without loss of generality, we may assume that $J = [j]$. For any fixed x , we have:

$$\Pr_{f \in \mathcal{D}_n^t} [x \text{ satisfies exactly the terms in } J] = \beta,$$

and thus by linearity of expectation, we have $\mathbf{E}_{f \in \mathcal{D}_n^t} [U_J] = \beta$. Now we show that with high probability that the deviation of U_J from its expected value is low.

Applying Lemma 19, we may assume that the terms T_1, \dots, T_j contain at least $j \log t - (\log t)/4$ many variables, and that $J \cup T_i$ for all $i = j+1, \dots, t$ includes at least $(j+1) \log t - (\log t)/4$ many unique variables, while increasing the failure probability by only $1/n^{\Omega(\log t)}$. Note that conditioning on this event can change the value of U_J by at most $1/n^{\Omega(\log t)} < \frac{1}{2}\beta$, so under this conditioning we have $\mathbf{E}[P_j] \geq \frac{1}{2}\beta$. Conditioning on this event, fix the terms T_1, \dots, T_j . Then the terms T_{j+1}, \dots, T_t are chosen uniformly and independently from the set of all terms T of length $\log t$ such that the union of the variables in J and T includes at least $(j+1) \log t - (\log t)/4$ unique variables. Call this set \mathcal{X} .

We now use McDiarmid’s inequality where the random variables are the terms T_{j+1}, \dots, T_t randomly selected from \mathcal{X} , letting $g(T_{j+1}, \dots, T_i) = U_J$ and $g(T_{j+1}, \dots, T_{i-1}, T'_i, T_{i+1}, \dots, T_t) = U'_J$ for all $i = j+1, \dots, t$. We claim that:

$$|U_J - U'_J| \leq d_i := \frac{t^{1/4}}{t^{j+1}}.$$

This is because U'_J can only be larger than U_J by assignments which satisfy T_1, \dots, T_j and T_i . Similarly, U'_J can only be smaller than U_J by assignments which satisfy T_1, \dots, T_j and T'_i . Since T_i and T'_i come from \mathcal{X} , we know that at least $(j+1)t - (\log t)/4$ variables must be satisfied.

Thus we may apply McDiarmid’s inequality with $\tau = \frac{3}{2}\beta$, which gives that $\Pr_f[U_J > 2\beta]$ is at most

$$\exp\left(\frac{-2\frac{9}{4}\beta^2}{t^{3/2}/t^{2j+2}}\right) \leq \exp\left(\frac{-9\sqrt{t}(1-1/t)^{2(t-j)}}{2}\right).$$

Combining the failure probabilities over all the $\binom{t}{j}$ possible sets, we get that with probability at least

$$\binom{t}{j} \left(\frac{1}{n^{\Omega(\log t)}} + e^{-9\sqrt{t}(1-1/t)^{2(t-j)}/2} \right) = \frac{1}{n^{\Omega(\log t)}},$$

over the random draw of f from \mathcal{D}_n^t , U_J for all $J \subseteq [t]$ of size j is at most 2β . Thus, the probability that a random input satisfies exactly some j distinct terms of f is at most $2\binom{t}{j}\beta$. \blacksquare

Using these properties of random DNF formulas we can now show a lemma analogous to Lemma 17 for random DNF formulas.

Lemma 22. *Let f be any DNF formula with $t = n^{O(1)}$ terms, and let $\epsilon > 0$ which satisfies $1/\epsilon = o(\log \log n)$. Then set $d = 4e^3 \ln 1/\epsilon$ and $\ell = c \log n$, for any constant c . Lemma 21. If*

$$\Pr[y_f(x) = j] \leq \left(\frac{e \ln 1/\epsilon}{j} \right)^j,$$

for every $d \leq j \leq \ell$, then the polynomial $p_{f,d}$ satisfies $\mathbf{E}[(f - p_{f,d})^2] \leq \epsilon$.

Proof. We condition on the values of $y_f(x)$, controlling the magnitude of $p_{f,d}$ by the unlikelihood of y_f being large. By Fact 14, $p_{f,d}(x)$ will output 0 if x does not satisfy f , $p_{f,d}(x)$ will output 1 if $y_f(x) \in [d]$, and $|p_{f,d}(x)| < \binom{y_f}{d}$ for $y_f(x) \in [t] \setminus [d]$. Hence:

$$\begin{aligned} \|f - p_{f,d}\|^2 &< \sum_{j=d+1}^{\ell-1} \binom{j}{d}^2 \left(\frac{e \ln 1/\epsilon}{j} \right)^j + \binom{t}{d}^2 \cdot \Pr[y_f \geq \ell] \\ &< \sum_{j=d+1}^{\ell-1} 2^{2j} \left(\frac{e \ln 1/\epsilon}{4e^3 \ln 1/\epsilon} \right)^j + n^{-\Omega(\log \log n)} \\ &< \epsilon \sum_{j=d+1}^{\ell-1} \frac{1}{e^j} + n^{-\Omega(\log \log n)} < \epsilon. \end{aligned}$$

■

We can now show that Mansour's conjecture [Man94] is true with high probability over the choice of f from \mathcal{D}_n^t .

Theorem 23. *Let $f : \{0, 1\}^n \rightarrow \{0, 1\}$ be a $t = n^{O(1)}$ -term DNF formula where each term is chosen independently from the set of all terms of length $\log t$. Then with probability at least $1 - n^{-\Omega(\log t)}$ over the choice of f , there exists a polynomial p with $\|p\|_1 \leq t^{O(\log 1/\epsilon)}$ such that $\mathbf{E}[(p(x) - f(x))^2] \leq \epsilon$.*

Proof. Recall that if $t = O(n^{1/4})$, f is a read-once DNF formula with high probability and thus Theorem 18 holds.

Let $d := 4e^3 \ln(1/\epsilon)$ and $p_{f,d}$ be as defined in Section 3. Lemma 15 tells us that $\|p_{f,d}\|_1 \leq t^{O(\log 1/\epsilon)}$. We show that with probability at least $1 - n^{-\omega(1)}$ over the random draw of f from \mathcal{D}_n^t , $p_{f,d}$ will be a good approximator for f . This follows by Lemma 21; with probability at least $1 - (c \log(n) - d - 1)/n^{\Omega(\log t)} = 1 - n^{-\Omega(\log t)}$, we have $\Pr[y = j]$ for all $d < j \leq c \log(n)$. Thus for such f Lemma 17 tells us that $\mathbf{E}[(f - p_{f,d})^2] \leq \epsilon$. ■

5 Pseudorandomness

De et al. [DETT09] recently improved long-standing pseudorandom generators against DNF formulas.

Definition 24. *A probability distribution X over $\{0, 1\}^n$ ϵ -fools a real function $f : \{0, 1\}^n \rightarrow \mathbb{R}$ if*

$$|\mathbf{E}[f(X)] - \mathbf{E}[f(U_n)]| \leq \epsilon.$$

If \mathcal{C} is a class of functions, then we say that X ϵ -fools \mathcal{C} if X ϵ -fools every function $f \in \mathcal{C}$.

We say a probability distribution X over $\{0, 1\}^n$ is ϵ -biased if it ϵ -fools the character function χ_S for every $S \subseteq [n]$.

De *et al.* observed that the result of Bazzi [Baz07] implied a pseudorandom generator that ϵ -fools t -term DNF formulas over n variables with seed length $O(\log n \cdot \log^2(t/\beta))$, which already improves the long-standing upper bound of $O(\log^4(tn/\epsilon))$ of Luby *et al.* [LVW93]. They go on to show a pseudorandom generator with seed length $O(\log n + \log^2(t/\epsilon) \log \log(t/\epsilon))$.

They prove that a sufficient condition for a function f to be ϵ -fooled by an ϵ -biased distribution is that the function be “sandwiched” between two bounded real-valued functions whose Fourier transform has small ℓ_1 norm:

Lemma 25 (Sandwich Bound [DETT09]). *Suppose $f, f_\ell, f_u : \{0, 1\}^n \rightarrow \mathbb{R}$ are three functions such that for every $x \in \{0, 1\}^n$, $f_\ell(x) \leq f(x) \leq f_u(x)$, $\mathbf{E}[f_u(U_n)] - \mathbf{E}[f(U_n)] \leq \epsilon$, and $\mathbf{E}[f(U_n)] - \mathbf{E}[f_\ell(U_n)] \leq \epsilon$. Let $L = \max(\|f_\ell\|_1^{\neq 0}, \|f_u\|_1^{\neq 0})$. Then any β -biased probability distribution $(\epsilon + \beta L)$ -fools f .*

Naor and Naor [NN93] prove that an ϵ -biased distribution over n bits can be sampled using a seed of $O(\log(n/\epsilon))$ bits. Using our construction from Section 4, we show that random DNF formulas are ϵ -fooled by a pseudorandom generator with seed length $O(\log n + \log(t) \log(1/\epsilon))$:

Theorem 26. *Let $f = T_1 \vee \dots \vee T_t$ be a random DNF formula chosen from \mathcal{D}_n^t . For $1 \leq d \leq t$, with probability $1 - 1/n^{\Omega(\log t)}$ over the choice of f , β -biased distributions $O(2^{-\Omega(d)} + \beta t^d)$ -fool f . In particular, we can ϵ -fool most $f \in \mathcal{D}_n^t$ by a $t^{-O(\log(1/\epsilon))}$ -biased distribution.*

Proof. Let d^+ be the first odd integer greater than d , and let d^- be the first even integer greater than d . Let $f_u = p_{f,d^+}$ and $f_\ell = p_{f,d^-}$ (where $p_{f,d}$ is defined as in Section 3). By Lemma 15, the ℓ_1 -norms of f_u and f_ℓ are $t^{O(d)}$. By Fact 14, we know that $P_{d^+}(y) = \binom{y-1}{d} + 1 > 1$ and $P_{d^-}(y) = -\binom{y-1}{d} + 1 \leq 0$ for $y \in [t] \setminus [d]$, hence:

$$\mathbf{E}[f_u(U_n)] - \mathbf{E}[f(U_n)] = \sum_{j=d+1}^t \left(\binom{j-1}{d} + 1 - 1 \right) \Pr[y_f = j],$$

which with probability $1 - 1/n^{\Omega(\log t)}$ over the choice of f is at most $2^{-\Omega(d)}$ by the analysis in Lemma 17. The same analysis applies for f_ℓ , thus applying Lemma 25 gives us the theorem. ■

De *et al.* match our bound for random DNF formulas for the special case of read-once DNF formulas. We remark that our construction from Section 3.1 can be used to recover the bound for read-once DNF formulas as well.

6 Discussion

On the relationship between Mansour’s Conjecture and the Entropy-Influence Conjecture. As a final note, we would like to make a remark on the relationship between Mansour’s conjecture and the entropy-influence conjecture. The *spectral entropy* of a function is defined to be $E(f) := \sum_S -\hat{f}(S)^2 \log(\hat{f}(S)^2)$ and the *total influence* to be $I(f) := \sum_S |S| \hat{f}(S)^2$. The *entropy-influence* conjecture is that $E(f) =$

$O(I(f))$ [FK96].¹ Boppana showed that the total influence of t -term DNF formulas is $O(\log t)$ [Bop97]. From this it follows that Mansour’s conjecture is implied by the entropy-influence conjecture.

It can be shown that for $n^{O(1)}$ -size DNF formulas Mansour’s conjecture implies an upper bound on the spectral entropy of $O(\log n)$. Thus, for the class of DNF formulas we consider in Section 4 (which have total influence $\Omega(\log n)$), our results imply that the entropy-influence conjecture is true.

Acknowledgments. Thanks to Sasha Sherstov for important contributions at an early stage of this work, and Omid Etesami for pointing out an error in an earlier version of this paper.

References

- [Baz07] Louay Bazzi. Polylogarithmic independence can fool DNF formulas. In *Proc. 48th IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 63–73, 2007.
- [Bop97] Ravi B. Boppana. The average sensitivity of bounded-depth circuits. *Information Processing Letters*, 63(5):257–261, 1997.
- [DETT09] Anindya De, Omid Etesami, Luca Trevisan, and Madhur Tulsiani. Improved pseudorandom generators for depth 2 circuits. Technical Report 141, Electronic Colloquium on Computational Complexity (ECCC), 2009.
- [FK96] Ehud Friedgut and Gil Kalai. Every monotone graph property has a sharp threshold. *Proceedings of the American Mathematical Society*, 124(10), 1996.
- [GKK08a] Parikshit Gopalan, Adam Kalai, and Adam R. Klivans. A query algorithm for agnostically learning DNF? In *21st Annual Conference on Learning Theory - COLT 2008, Helsinki, Finland, July 9-12, 2008*, pages 515–516. Omnipress, 2008.
- [GKK08b] Parikshit Gopalan, Adam Tauman Kalai, and Adam R. Klivans. Agnostically learning decision trees. In *Proceedings of the 40th Annual ACM Symposium on Theory of Computing, Victoria, British Columbia, Canada, May 17-20, 2008*, pages 527–536. ACM, 2008.
- [GKP94] Ronald L. Graham, Donald E. Knuth, and Oren Patashnik. *Concrete Mathematics: A Foundation for Computer Science*. Addison-Wesley, 1994.
- [Hås86] Johan Håstad. *Computational Limitations for Small Depth Circuits*. MIT Press, 1986.
- [Jac97] Jeffrey C. Jackson. An efficient membership-query algorithm for learning DNF with respect to the uniform distribution. *Journal of Computer and System Sciences*, 55(3):414–440, 1997.
- [JLSW08] Jeffrey C. Jackson, Homin K. Lee, Rocco A. Servedio, and Andrew Wan. Learning random monotone DNF. In *11th International Workshop on Approximation Algorithms for Combinatorial Optimization Problems and 12th International Workshop on Randomization and Computation (RANDOM-APPROX)*, pages 483–497. Springer-Verlag, 2008.

¹<http://terrytao.wordpress.com/2007/08/16/gil-kalai-the-entropyinfluence-conjecture/>

- [JS05] Jeffrey C. Jackson and Rocco A. Servedio. On learning random DNF formulas under the uniform distribution. In *8th International Workshop on Approximation Algorithms for Combinatorial Optimization Problems and 9th International Workshop on Randomization and Computation (RANDOM-APPROX)*, volume 3624 of *Lecture Notes in Computer Science*, pages 342–353. Springer-Verlag, 2005.
- [KKMS08] A. Kalai, A. Klivans, Y. Mansour, and R. Servedio. Agnostically learning halfspaces. *SIAM Journal on Computing*, 37(6):1777–1805, 2008.
- [KLPV87] M. Kearns, M. Li, L. Pitt, and L. Valiant. On the learnability of Boolean formulae. In *Proceedings of the nineteenth annual ACM symposium on Theory of computing*, pages 285–295. ACM, 1987.
- [KM93] Eyal Kushilevitz and Yishay Mansour. Learning decision trees using the Fourier spectrum. *SIAM Journal on Computing*, 22(6):1331–1348, December 1993. Prelim. ver. in *Proc. of STOC'91*.
- [LMN93] N. Linial, Y. Mansour, and N. Nisan. Constant depth circuits, Fourier transform and learnability. *Journal of the ACM*, 40(3):607–620, 1993.
- [LVW93] Michael Luby, Boban Velickovic, and Avi Wigderson. Deterministic approximate counting of depth-2 circuits. In *ISTCS 1993*, pages 18–24, 1993.
- [Man94] Y. Mansour. *Learning Boolean functions via the Fourier transform*, pages 391–424. Kluwer Academic Publishers, 1994.
- [Man95] Yishay 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. Prelim. ver. in *Proc. of COLT'92*.
- [NN93] Joseph Naor and Moni Naor. Small-bias probability spaces: Efficient constructions and applications. *SIAM Journal on Computing*, 22(4):838–856, 1993.
- [Raz08] A. Razborov. A simple proof of Bazzi's theorem. Technical Report 81, Electronic Colloquium on Computational Complexity (ECCC), 2008.
- [Sel08] Linda Sellie. Learning random monotone DNF under the uniform distribution. In *Proc. of the 21th Annual Conference on Computational Learning Theory (COLT)*, pages 181–192, 2008.
- [Sel09] Linda Sellie. Exact learning of random DNF over the uniform distribution. In *Proc. 41st Annual ACM Symposium on Theory of Computing (STOC)*, pages 45–54, 2009.