

Learning random monotone DNF

Jeffrey C. Jackson^{1*}, Homin K. Lee², Rocco A. Servedio^{2**}, and Andrew Wan²

¹ Duquesne University, Pittsburgh, PA 15282
jacksonj@duq.edu

² Columbia University, New York, NY 10027
homin@cs.columbia.edu, rocco@cs.columbia.edu, atw12@cs.columbia.edu

Abstract. We give an algorithm that with high probability properly learns random monotone DNF with $t(n)$ terms of length $\approx \log t(n)$ under the uniform distribution on the Boolean cube $\{0, 1\}^n$. For any function $t(n) \leq \text{poly}(n)$ the algorithm runs in time $\text{poly}(n, 1/\epsilon)$ and with high probability outputs an ϵ -accurate monotone DNF hypothesis. This is the first algorithm that can learn monotone DNF of arbitrary polynomial size in a reasonable average-case model of learning from random examples only.

1 Introduction

Motivation and background. Any Boolean function $f : \{0, 1\}^n \rightarrow \{0, 1\}$ can be expressed as a disjunction of conjunctions of Boolean literals, i.e. as an OR of ANDs. Such a logical formula is said to be a *disjunctive normal formula*, or DNF. Learning polynomial-size DNF formulas (disjunctions of $\text{poly}(n)$ many conjunctions) from random examples is an outstanding open question in computational learning theory, dating back more than 20 years to Valiant’s introduction of the PAC (Probably Approximately Correct) learning model [Val84].

The most intensively studied variant of the DNF learning problem is PAC learning DNF under the uniform distribution. In this problem the learner must generate a high-accuracy hypothesis with high probability when given uniform random examples labeled according to the unknown target DNF. Despite much effort, no polynomial-time algorithms are known for this problem.

A tantalizing question that has been posed as a goal by many authors (see e.g. [Jac97, JT97, BBL98, Blu03b, Ser04]) is to learn *monotone* DNF, which only contain unnegated Boolean variables, under the uniform distribution. Besides being a natural restriction of the uniform distribution DNF learning problem, this problem is interesting because several impediments to learning general DNF under uniform – known lower bounds for Statistical Query based algorithms [BFJ⁺94], the apparent hardness of learning the subclass of $\log(n)$ -juntas [Blu03a] – do not apply in the monotone case. This paper solves a natural average-case version of this problem.

* Supported in part by NSF award CCF-0209064

** Supported in part by NSF award CCF-0347282, by NSF award CCF-0523664, and by a Sloan Foundation Fellowship.

Previous work. Many partial results have been obtained on learning monotone DNF under the uniform distribution. Verbeurgt [Ver90] gave an $n^{O(\log n)}$ -time uniform distribution algorithm for learning any $\text{poly}(n)$ -term DNF, monotone or not. Several authors [KMSP94,SM00,BT96] have given results on learning monotone t -term DNF for larger and larger values of t ; most recently, [Ser04] gave a uniform distribution algorithm that learns any $2^{O(\sqrt{\log n})}$ -term monotone DNF to any constant accuracy $\epsilon = \Theta(1)$ in $\text{poly}(n)$ time. O’Donnell and Servedio [OS06] have recently shown that $\text{poly}(n)$ -leaf *decision trees* that compute monotone functions (a subclass of $\text{poly}(n)$ -term monotone DNF) can be learned to any constant accuracy under uniform in $\text{poly}(n)$ time. Various other problems related to learning different types of monotone functions under uniform have also been studied, see e.g. [KLV94,BBL98,Ver98,HM91,AM02].

Aizenstein and Pitt [AP95] first proposed a model of random DNF formulas and gave an exact learning algorithm that learns random DNFs generated in this way. As noted in [AP95] and [JS06], this model admits a trivial learning algorithm in the uniform distribution PAC setting. Jackson and Servedio [JS05] gave a uniform distribution algorithm that learns log-depth decision trees on average in a natural random model. Previous work on average-case uniform PAC DNF learning, also by Jackson and Servedio, is described below.

Our results. The main result of this paper is a polynomial-time algorithm that can learn random $\text{poly}(n)$ -term monotone DNF with high probability. (We give a full description of the exact probability distribution defining our random DNFs in Section 4; briefly, the reader should think of our random t -term monotone DNFs as being obtained by independently drawing t monotone conjunctions uniformly from the set of all conjunctions of length $\log_2 t$ over variables x_1, \dots, x_n . Although many other distributions could be considered, this seems a natural starting point. Some justification for the choice of term length is given in Sections 4 and 6.)

Theorem 1. [Informally] *Let $t(n) \leq \text{poly}(n)$, and let $c > 0$ be any fixed constant. Then random monotone $t(n)$ -term DNFs are PAC learnable (with failure probability $\delta = n^{-c}$) to accuracy ϵ in $\text{poly}(n, 1/\epsilon)$ time under the uniform distribution. The algorithm outputs a monotone DNF as its hypothesis.*

In independent and concurrent work, Sellie [Sel08] has given an alternate proof of this theorem using different techniques.

Our technique. Jackson and Servedio [JS06] showed that for any $\gamma > 0$, a result similar to Theorem 1 holds for random t -term monotone DNF with $t \leq n^{2-\gamma}$. The main open problem stated in [JS06] was to prove Theorem 1. Our work solves this problem by using the previous algorithm to handle $t \leq n^{3/2}$, developing new Fourier lemmas for monotone DNF, and using these lemmas together with more general versions of techniques from [JS06] to handle $t \geq n^{3/2}$.

The crux of our strategy is to establish a connection between the term structure of certain monotone DNFs and their low-order Fourier coefficients. There is an extensive body of research on Fourier properties of monotone Boolean functions [BT96,MO03,BBL98], polynomial-size DNF [Jac97,Man95], and related

classes. These results typically establish that *every* function in the class has a Fourier spectrum with certain properties; unfortunately, the Fourier properties that have been obtainable to date for general statements of this sort have not been sufficient to yield polynomial-time learning algorithms.

We take a different approach by carefully defining a set of conditions, and showing that *if a monotone DNF f satisfies these conditions then the structure of the terms of f will be reflected in the low-order Fourier coefficients of f* . In [JS06], the degree two Fourier coefficients were shown to reveal the structure of the terms for certain (including random) monotone DNFs having at most $n^{2-\gamma}$ terms. In this work we develop new lemmas about the Fourier coefficients of more general monotone DNF, and use these new lemmas to establish a connection between term structure and constant degree Fourier coefficients for monotone DNFs with any polynomial number of terms. Roughly speaking, this connection holds for monotone DNF that satisfy the following conditions:

- each term has a reasonably large fraction of assignments which satisfy it and no other term;
- for each small tuple of distinct terms, only a small fraction of assignments simultaneously satisfy all terms in the tuple; and
- for each small tuple of variables, only a few terms contains the entire tuple.

The “small” tuples referred to above should be thought of as tuples of constant size. The constant degree coefficients capture the structure of the terms in the following sense: tuples of variables that all co-occur in some term will have a large magnitude Fourier coefficient, and tuples of variables that do not all co-occur in some term will have a small magnitude Fourier coefficient (even if subsets of the tuple do co-occur in some terms). We show this in Section 2.

Next we show a reconstruction procedure for obtaining the monotone DNF from tuple-wise co-occurrence information. Given a hypergraph with a vertex for each variable, the procedure turns each co-occurrence into a hyperedge, and then searches for all hypercliques of size corresponding to the term length. The hypercliques that are found correspond to the terms of the monotone DNF hypothesis that the algorithm constructs. This procedure is described in Section 3; we show that it succeeds in constructing a high-accuracy hypothesis if the monotone DNF f satisfies a few additional conditions. This generalizes a reconstruction procedure from [JS06] that was based on finding cliques in a graph (in the $n^{2-\gamma}$ -term DNF setting, the algorithm deals only with co-occurrences of pairs of variables so it is sufficient to consider only ordinary graphs rather than hypergraphs).

The ingredients described so far thus give us an efficient algorithm to learn any monotone DNF that satisfies all of the required conditions. Finally, we show that random monotone DNF satisfy all the required conditions with high probability. We do this in Section 4 via a fairly delicate probabilistic argument. Section 5 combines the above ingredients to prove Theorem 1. We close the paper by showing that our technique lets us easily recapture the result of [HM91] that read- k monotone DNF are uniform-distribution learnable in polynomial time.

Preliminaries. We write $[n]$ to denote the set $\{1, \dots, n\}$ and use capital letters for subsets of $[n]$. We will use calligraphic letters such as \mathcal{C} to denote sets of

sets and script letters such as \mathcal{X} to denote sets of sets of sets. We write \log to denote \log_2 and \ln to denote the natural log. We write U_n to denote the uniform distribution over the Boolean cube $\{0, 1\}^n$.

A Boolean function $f : \{0, 1\}^n \rightarrow \{0, 1\}$ is *monotone* if changing the value of an input bit from 0 to 1 never causes the value of f to change from 1 to 0. We denote the input variables to f as x_1, \dots, x_n . A *t-term monotone DNF* is a t -way OR of ANDs of Boolean variables (no negations). Recall that every monotone Boolean function has a unique representation as a reduced monotone DNF. We say that a term T of such a monotone DNF is *uniquely satisfied* by input x if x satisfies T and no other term of f .

Our learning model is an “average-case” variant of the well-studied uniform distribution PAC learning model. Let D_C be a probability distribution over some fixed class C of Boolean functions over $\{0, 1\}^n$, and let f (drawn from D_C) be an unknown target function. A learning algorithm A for D_C takes as input an accuracy parameter $0 < \epsilon < 1$ and a confidence parameter $0 < \delta < 1$. During its execution, algorithm A has access to a *random example oracle* $EX(f, U_n)$, which, when queried generates a random labeled example $(x, f(x))$, where x is drawn from U_n . The learning algorithm outputs a hypothesis h , which is a Boolean function over $\{0, 1\}^n$. The error of this hypothesis is defined to be $\Pr_{U_n}[h(x) \neq f(x)]$. We say that A *learns* D_C *under* U_n if for every $0 < \epsilon, \delta < 1$, with probability at least $1 - \delta$ (over both the random examples used for learning and the random draw of f from D_C) algorithm A outputs a hypothesis h which has error at most ϵ .

2 Fourier coefficients and monotone DNF term structure

Throughout this section let $f(x_1, \dots, x_n)$ be a monotone DNF and let $S \subseteq \{1, \dots, n\}$ be a fixed subset of variables. We write s to denote $|S|$ throughout this section. The Fourier coefficient, written $\hat{f}(S)$, measures the correlation between f and the parity of the variables in S .

The main result of this section is Lemma 3, which shows that under suitable conditions on f , the value $|\hat{f}(S)|$ is “large” if and only if f has a term containing all the variables of S . To prove this, we observe that the inputs which uniquely satisfy such a term will make a certain contribution to $\hat{f}(S)$. (In Section 2.1 we explain this in more detail and show how to view $\hat{f}(S)$ as a sum of contributions from inputs to f .) It remains then to show that the contribution from other inputs is small. The main technical novelty comes in Sections 2.2 and 2.3, where we show that all other inputs which make a contribution to $\hat{f}(S)$ must satisfy the terms of f in a special way, and use this property to show that under suitable conditions on f , the fraction of such inputs must be small.

2.1 Rewriting $\hat{f}(S)$.

We observe that $\hat{f}(S)$ can be expressed in terms of 2^s conditional probabilities, each of which is the probability that f is satisfied conditioned on a particular

setting of the variables in S . That is:

$$\begin{aligned}\hat{f}(S) &\stackrel{\text{def}}{=} \mathbf{E}_{x \in U^n} \left[(-1)^{\sum_{i \in S} x_i} \cdot f(x) \right] = \frac{1}{2^n} \sum_{x \in \{0,1\}^n} (-1)^{\sum_{i \in S} x_i} \cdot f(x) \\ &= \frac{1}{2^n} \sum_{U \subseteq S} (-1)^{|U|} \sum_{x \in Z_S(U)} f(x) = \frac{1}{2^s} \sum_{U \subseteq S} (-1)^{|U|} \Pr_x[f(x) = 1 \mid x \in Z_S(U)],\end{aligned}$$

where $Z_S(U)$ denotes the set of those $x \in \{0,1\}^n$ such that $x_i = 1$ for all $i \in U$ and $x_i = 0$ for all $i \in S \setminus U$. If f has some term T containing all the variables in S , then $\Pr_x[f(x) = 1 \mid x \in Z_S(S)]$ is at least as large as $\Pr_x[T \text{ is uniquely satisfied in } f \mid x \in Z_S(S)]$. On the other hand, if f has no such term, then $\Pr_x[f(x) = 1 \mid x \in Z_S(S)]$ does not receive this contribution. We will show that this contribution is the chief determinant of the magnitude of $\hat{f}(S)$.

It is helpful to rewrite $\hat{f}(S)$ as a sum of contributions from each input $x \in \{0,1\}^n$. To this end, we decompose f according to the variables of S . Given a subset $U \subseteq S$, we will write g_U to denote the disjunction of terms in f that contain every variable indexed by $U \subseteq S$ and no variable indexed by $S \setminus U$, but with the variables indexed by U removed from each term. (So for example if $f = x_1x_2x_4x_6 \vee x_1x_2x_5 \vee x_1x_2x_3 \vee x_3x_5 \vee x_1x_5x_6$ and $S = \{1, 2, 3\}$ and $U = \{1, 2\}$, then $g_U = x_4x_6 \vee x_5$.) Thus we can split f into disjoint sets of terms: $f = \bigvee_{U \subseteq S} (t_U \wedge g_U)$, where t_U is the term consisting of exactly the variables indexed by U .

Suppose we are given $U \subseteq S$ and an x that belongs to $Z_S(U)$. We have that $f(x) = 1$ if and only if $g_{U'}(x)$ is true for some $U' \subseteq U$. (Note that $t_{U'}(x)$ is true for every $U' \subseteq U$ since x belongs to $Z_S(U)$.) Thus we can rewrite the Fourier coefficients $\hat{f}(S)$ as follows: (Below we write $I(P)$ to denote the indicator function that takes value 1 if predicate P is true and value 0 if P is false.)

$$\begin{aligned}\hat{f}(S) &= \frac{1}{2^n} \sum_{U \subseteq S} (-1)^{|U|} \sum_{x \in Z_S(U)} f(x) = \sum_{U \subseteq S} (-1)^{|U|} \frac{1}{2^n} \sum_{x \in Z_S(U)} I \left(\bigvee_{U' \subseteq U} g_{U'}(x) \right) \\ &= \sum_{x \in \{0,1\}^n} \frac{1}{2^s} \frac{1}{2^n} \sum_{U \subseteq S} (-1)^{|U|} I \left(\bigvee_{U' \subseteq U} g_{U'}(x) \right).\end{aligned}$$

We can rewrite this as $\hat{f}(S) = \sum_{x \in \{0,1\}^n} \text{Cons}_S(x)$, where

$$\text{Cons}_S(x) \stackrel{\text{def}}{=} \frac{1}{2^s} \frac{1}{2^n} \sum_{U \subseteq S} (-1)^{|U|} I \left(\bigvee_{U' \subseteq U} g_{U'}(x) \right). \quad (1)$$

The value $\text{Cons}_S(x)$ may be viewed as the ‘‘contribution’’ that x makes to $\hat{f}(S)$. Recall that when f has a term T which contains all the variables in S , those $x \in Z_S(S)$ which uniquely satisfy T will contribute to $\hat{f}(S)$. We will show that under suitable conditions on f , the other x 's make little or no contribution.

2.2 Bounding the contribution to $\hat{f}(S)$ from various inputs.

The variable \mathcal{C} will denote a subset of $\mathcal{P}(S)$, the power set of S ; i.e. \mathcal{C} denotes a collection of subsets of S . We may view \mathcal{C} as defining a set of g_U 's (those g_U 's for which U belongs to \mathcal{C}).

We may partition the set of inputs $\{0,1\}^n$ into $2^{|\mathcal{P}(S)|} = 2^{2^s}$ parts according to what subset of the 2^s functions $\{g_U\}_{U \subseteq S}$ each $x \in \{0,1\}^n$ satisfies. For \mathcal{C} a subset of $\mathcal{P}(S)$ we denote the corresponding piece of the partition by $P_{\mathcal{C}}$; so $P_{\mathcal{C}}$ consists of precisely those $x \in \{0,1\}^n$ that satisfy $(\bigwedge_{U \in \mathcal{C}} g_U) \wedge (\bigwedge_{U \notin \mathcal{C}} \bar{g}_U)$. Note that for any given fixed \mathcal{C} , each x in $P_{\mathcal{C}}$ has exactly the same contribution $\text{Con}_S(x)$ to the Fourier coefficient $\hat{f}(S)$ as every other x' in $P_{\mathcal{C}}$; this is simply because x and x' will satisfy exactly the same set of g_U 's in (1). More generally, we have the following (proved in the full version):

Lemma 1. *Let \mathcal{C} be any subset of $\mathcal{P}(S)$. Suppose that there exist $U_1, U_2 \in \mathcal{C}$ such that $U_1 \subsetneq U_2$. Then for any y, z where $y \in P_{\mathcal{C}}$ and $z \in P_{\mathcal{C} \setminus U_2}$, we have that: $\text{Con}_S(y) = \text{Con}_S(z)$.*

Given a collection \mathcal{C} of subsets of S , let $\text{Con}_S(\mathcal{C})$ denote $\sum_{x \in P_{\mathcal{C}}} \text{Con}_S(x)$, and we refer to this quantity as the contribution that \mathcal{C} makes to the Fourier coefficient $\hat{f}(S)$. It is clear that we have $\hat{f}(S) = \sum_{\mathcal{C} \subseteq \mathcal{P}(S)} \text{Con}_S(\mathcal{C})$.

The following lemma, proved in the full version establishes a broad class of \mathcal{C} 's for which $\text{Con}_S(\mathcal{C})$ is zero:

Lemma 2. *Let \mathcal{C} be any collection of subsets of S . If $\bigcup_{U \in \mathcal{C}} U \neq S$ then $\text{Con}_S(x) = 0$ for each $x \in P_{\mathcal{C}}$ and hence $\text{Con}_S(\mathcal{C}) = 0$.*

It remains to analyze those \mathcal{C} 's for which $\bigcup_{U \in \mathcal{C}} U = S$; for such a \mathcal{C} we say that \mathcal{C} covers S .

Recall from the previous discussion that $\text{Con}_S(\mathcal{C}) = |P_{\mathcal{C}}| \cdot \text{Con}_S(x)$ where x is any element of $P_{\mathcal{C}}$. Since $|\text{Con}_S(x)| \leq \frac{1}{2^n}$ for all $x \in \{0,1\}^n$, for any collection \mathcal{C} , we have that

$$|\text{Con}_S(\mathcal{C})| \leq \Pr_{x \in U_n} [x \in P_{\mathcal{C}}] = \Pr_{x \in U_n} [(\bigwedge_{U \in \mathcal{C}} g_U) \wedge (\bigwedge_{U \notin \mathcal{C}} \bar{g}_U)] \leq \Pr_{x \in U_n} [(\bigwedge_{U \in \mathcal{C}} g_U)].$$

We are interested in bounding this probability for $\mathcal{C} \neq \{S\}$ (we will deal with the special case $\mathcal{C} = \{S\}$ separately later). Recall that each g_U is a disjunction of terms; the expression $\bigwedge_{U \in \mathcal{C}} g_U$ is satisfied by precisely those x that satisfy at least one term from each g_U as U ranges over all elements of \mathcal{C} . For $j \geq 1$ let us define a quantity B_j as follows

$$B_j \stackrel{\text{def}}{=} \max_{i_1, \dots, i_j} \Pr_{x \in U_n} [x \text{ simultaneously satisfies terms } T_{i_1}, \dots, T_{i_j} \text{ in } \bigvee_{U \subseteq S} (g_U)]$$

where the max is taken over all j -tuples of distinct terms in $\bigvee_{U \subseteq S} (g_U)$. Then it is not hard to see that by a union bound, we have

$$|\text{Con}_S(\mathcal{C})| \leq B_{|\mathcal{C}|} \prod_{U \in \mathcal{C}} (\#g_U), \quad (2)$$

where $\#g_U$ denotes the number of terms in the monotone DNF g_U .

The idea of why (2) is a useful bound is as follows. Intuitively, one would expect that the value of B_j decreases as j (the number of terms that must be satisfied) increases. One would also expect the value of $\#g_U$ to decrease as the size of U increases (if U contains more variables then fewer terms in f will contain all of those variables). This means that there is a trade-off which helps us bound (2): if $|\mathcal{C}|$ is large then $B_{|\mathcal{C}|}$ is small, but if $|\mathcal{C}|$ is small then (since we know that $\bigcup_{U \in \mathcal{C}} U = S$) some U is large and so $\prod_{U \in \mathcal{C}} \#g_U$ will be smaller.

2.3 Bounding $\hat{f}(S)$ based on whether S co-occurs in a term of f .

We are now ready to state formally the conditions on \hat{f} that allow us to detect a co-occurrence of variables in the value of the corresponding Fourier coefficient.

Lemma 3. *Let $f : \{0, 1\}^n \rightarrow \{-1, 1\}$ be a monotone DNF. Fix a set $S \subset [n]$ of size $|S| = s$ and let*

$$\mathcal{Y} = \{\mathcal{C} \subseteq \mathcal{P}(S) : \mathcal{C} \text{ covers } S \text{ and } S \notin \mathcal{C}\}.$$

Suppose that we define $\alpha, \beta_1, \dots, \beta_{2^s}$ and $\Phi : \mathcal{Y} \rightarrow \mathbb{R}$ so that:

- C1** Each term in f is uniquely satisfied with probability at least α ;
- C2** For $1 \leq j \leq 2^s$, each j -tuple of terms in f is simultaneously satisfied with probability at most β_j ; and
- C3** For every $\mathcal{C}_{\mathcal{Y}} \in \mathcal{Y}$ we have $\prod_{U \in \mathcal{C}_{\mathcal{Y}}} (\#g_U) \leq \Phi(\mathcal{C}_{\mathcal{Y}})$.

Then

1. If the variables in S do not simultaneously co-occur in any term of f , then

$$|\hat{f}(S)| \leq \Upsilon \quad \text{where} \quad \Upsilon := \sum_{\mathcal{C}_{\mathcal{Y}} \in \mathcal{Y}} (2^s \beta_{|\mathcal{C}_{\mathcal{Y}}|} \Phi(\mathcal{C}_{\mathcal{Y}}));$$

2. If the variables in S do simultaneously co-occur in some term of f , then $|\hat{f}(S)| \geq \frac{\alpha}{2^s} - 2 \cdot \Upsilon$.

Using Lemma 3, if f satisfies conditions **C1** through **C3** with values of β_j and $\Phi(\cdot)$ so that there is a ‘‘gap’’ between $\alpha/2^s$ and 3Υ , then we can determine whether all the variables in S simultaneously co-occur in a term by estimating the magnitude of $\hat{f}(S)$.

Proof. Let \mathcal{C}^* denote the ‘special’ element of $\mathcal{P}(S)$ that consists solely of the subset S , i.e. $\mathcal{C}^* = \{S\}$, and let $\mathcal{X} = \{\mathcal{C} \subseteq \mathcal{P}(S) : \mathcal{C} \text{ covers } S \text{ and } S \in \mathcal{C} \text{ and } \mathcal{C} \neq \mathcal{C}^*\}$. Using Lemma 2, we have

$$\hat{f}(S) = \text{Con}_S(\mathcal{C}^*) + \sum_{\mathcal{C}_{\mathcal{Y}} \in \mathcal{Y}} \text{Con}_S(\mathcal{C}_{\mathcal{Y}}) + \sum_{\mathcal{C}_{\mathcal{X}} \in \mathcal{X}} \text{Con}_S(\mathcal{C}_{\mathcal{X}}). \quad (3)$$

We first prove point 1. Suppose that the variables of S do not simultaneously co-occur in any term of f . Then g_S is the empty disjunction and $\#g_S = 0$,

so $\text{Con}_S(\mathcal{C}) = 0$ for any \mathcal{C} containing S . Thus in this case we have $\hat{f}(S) = \sum_{\mathcal{C}_{\mathcal{Y}} \in \mathcal{Y}} \text{Con}_S(\mathcal{C}_{\mathcal{Y}})$; using (2) and condition **C3**, it follows that $|\hat{f}(S)|$ is at most $\sum_{\mathcal{C}_{\mathcal{Y}} \in \mathcal{Y}} B_{|\mathcal{C}_{\mathcal{Y}}|} \Phi(\mathcal{C}_{\mathcal{Y}})$. It is not hard to see that $B_{|\mathcal{C}_{\mathcal{Y}}|} \leq 2^s \beta_{|\mathcal{C}_{\mathcal{Y}}|}$ (we give a proof in the full version). So in this case we have

$$|\hat{f}(S)| \leq \sum_{\mathcal{C}_{\mathcal{Y}} \in \mathcal{Y}} |\text{Con}_S(\mathcal{C}_{\mathcal{Y}})| \leq \sum_{\mathcal{C}_{\mathcal{Y}} \in \mathcal{Y}} B_{|\mathcal{C}_{\mathcal{Y}}|} \Phi(\mathcal{C}_{\mathcal{Y}}) \leq \sum_{\mathcal{C}_{\mathcal{Y}} \in \mathcal{Y}} (2^s \beta_{|\mathcal{C}_{\mathcal{Y}}|} \Phi(\mathcal{C}_{\mathcal{Y}})) = \Upsilon.$$

Now we turn to point 2. Suppose that the variables of S do co-occur in some term of f . Let x be any element of $P_{\mathcal{C}^*}$, so x satisfies g_U if and only if $U = S$. It is easy to see from (1) that for such an x we have $\text{Con}_S(x) = (-1)^{|S|}/(2^n 2^s)$. We thus have that

$$\text{Con}_S(\mathcal{C}^*) = \frac{(-1)^{|S|}}{2^s} \cdot \Pr[x \in P_{\mathcal{C}^*}] = \frac{(-1)^{|S|}}{2^s} \Pr[g_S \wedge \left(\bigwedge_{U \subsetneq S} \bar{g}_U \right)]. \quad (4)$$

Since S co-occurs in some term of f , we have that g_S contains at least one term T . By condition **C1**, the corresponding term $(T \wedge (\bigwedge_{i \in S} x_i))$ of f is uniquely satisfied with probability at least α . Since each assignment that uniquely satisfies $(T \wedge (\bigwedge_{i \in S} x_i))$ (among all the terms of f) must satisfy $g_S \wedge (\bigwedge_{U \subsetneq S} \bar{g}_U)$, we have that the magnitude of (4) is at least $\alpha/2^s$.

We now show that $|\sum_{\mathcal{C}_{\mathcal{X}} \in \mathcal{X}} \text{Con}_S(\mathcal{C}_{\mathcal{X}})| \leq \Upsilon$, which completes the proof, since we already have that $|\sum_{\mathcal{C}_{\mathcal{Y}} \in \mathcal{Y}} \text{Con}_S(\mathcal{C}_{\mathcal{Y}})| \leq \sum_{\mathcal{C}_{\mathcal{Y}} \in \mathcal{Y}} |\text{Con}_S(\mathcal{C}_{\mathcal{Y}})| \leq \Upsilon$. First note that if the set $\mathcal{C}_{\mathcal{X}} \setminus \{S\}$ does not cover S , then by Lemmas 1 and 2 we have that $\text{Con}_S(x) = 0$ for each $x \in P_{\mathcal{C}_{\mathcal{X}}}$ and thus $\text{Con}_S(\mathcal{C}_{\mathcal{X}}) = 0$. So we may restrict our attention to those $\mathcal{C}_{\mathcal{X}}$ such that $\mathcal{C}_{\mathcal{X}} \setminus \{S\}$ covers S . Now since such a $\mathcal{C}_{\mathcal{X}} \setminus \{S\}$ is simply some $\mathcal{C}_{\mathcal{Y}} \in \mathcal{Y}$, and each $\mathcal{C}_{\mathcal{Y}} \in \mathcal{Y}$ is obtained as $\mathcal{C}_{\mathcal{X}} \setminus \{S\}$ for at most one $\mathcal{C}_{\mathcal{X}} \in \mathcal{X}$, we have

$$\left| \sum_{\mathcal{C}_{\mathcal{X}} \in \mathcal{X}} \text{Con}_S(\mathcal{C}_{\mathcal{X}}) \right| \leq \sum_{\mathcal{C}_{\mathcal{Y}} \in \mathcal{Y}} |\text{Con}_S(\mathcal{C}_{\mathcal{Y}})| \leq \Upsilon.$$

3 Hypothesis formation

In this section, we show that if a target monotone DNF f satisfies the conditions of Lemma 3 and two other simple conditions stated below (see Theorem 2), then it is possible to learn f from uniform random examples.

Theorem 2. *Let f be a t -term monotone DNF. Fix $s \in [n]$. Suppose that*

- *For all sets $S \subset [n], |S| = s$, conditions **C1** through **C3** of Lemma 3 hold for certain values α, β_j , and $\Phi(\cdot)$ satisfying $\Delta > 0$, where $\Delta := \alpha/2^s - 3 \cdot \Upsilon$. (Recall that $\Upsilon := \sum_{\mathcal{C}_{\mathcal{Y}} \in \mathcal{Y}} (2^s \beta_{|\mathcal{C}_{\mathcal{Y}}|} \Phi(\mathcal{C}_{\mathcal{Y}}))$, where $\mathcal{Y} = \{\mathcal{C} \subseteq \mathcal{P}(S) : \mathcal{C} \text{ covers } S \text{ and } S \notin \mathcal{C}\}$.)*

C4 *Every set S of s co-occurring variables in f appears in at most γ terms (here $\gamma \geq 2$); and*

C5 Every term of f contains at most κ variables (note that $s \leq \kappa \leq n$).

Then algorithm \mathcal{A} (described formally in the full version) PAC learns f to accuracy ϵ with confidence $1 - \delta$ given access to $EX(f, U_n)$, and runs in time $\text{poly}(n^{s+\gamma}, t, 1/\Delta, \gamma^\kappa, 1/\epsilon, \log(1/\delta))$.

Proof. Lemma 3 implies that for each set $S \subset [n], |S| = s$,

- if the variables in S all co-occur in some term of f , then $|\hat{f}(S)|$ is at least $\Delta/2$ larger than $\Upsilon + \Delta/2$;
- if the variables in S do not all co-occur in some term of f , then $|\hat{f}(S)|$ is at least $\Delta/2$ smaller than $\Upsilon + \Delta/2$.

A straightforward application of Hoeffding bounds (to estimate the Fourier coefficients using a random sample of uniformly distributed examples) shows that Step 1 of Algorithm \mathcal{A} can be executed in $\text{poly}(n^s, 1/\Delta, \log(1/\delta))$ time, and that with probability $1 - \delta/3$ the S 's that are marked as “good” will be precisely the s -tuples of variables that co-occur in some term of f .

Conceptually, the algorithm next constructs the hypergraph G_f that has one vertex per variable in f and that includes an s -vertex hyperedge if and only if the corresponding s variables co-occur in some term of f . Clearly there is a k -hyperclique in G_f for each term of k variables in f . So if we could find all of the k -hypercliques in G_f (where again k ranges between s and κ), then we could create a set HC_f of monotone conjunctions of variables such that f could be represented as an OR of t of these conjunctions. Treating each of the conjunctions in HC_f as a variable in the standard elimination algorithm for learning disjunctions (see e.g. Chapter 1 of [KV94]) would then enable us to properly PAC learn f to accuracy ϵ with probability at least $1 - \delta/3$ in time polynomial in $n, t, |HC_f|, 1/\epsilon$, and $\log(1/\delta)$. Thus, \mathcal{A} will use a subalgorithm \mathcal{A}' to find all the k -hypercliques in G_f and will then apply the elimination algorithm over the corresponding conjunctions to learn the final approximator h .

We now explain the subalgorithm \mathcal{A}' for locating the set HC_f of k -hypercliques. For each set S of s co-occurring variables, let $N_S \subseteq ([n] \setminus S)$ be defined as follows: a variable x_i is in N_S if and only if x_i is present in some term that contains all of the variables in S . Since by assumption there are at most γ terms containing such variables and each term contains at most κ variables, this means that $|N_S| < \kappa\gamma$. The subalgorithm will use this bound as follows. For each set S of s co-occurring variables, \mathcal{A}' will determine the set N_S using a procedure \mathcal{A}'' described shortly. Then, for each $s \leq k \leq \kappa$ and each $(k - s)$ -element subset N' of N_S , \mathcal{A}' will test whether or not $N' \cup S$ is a k -hyperclique in G_f . The set of all k -hypercliques found in this way is HC_f . For each S , the number of sets tested in this process is at most

$$\sum_{i=0}^{\kappa} \binom{|N_S|}{i} \leq \sum_{i=0}^{\kappa} \binom{\kappa\gamma}{i} \leq \left(\frac{e\kappa\gamma}{\kappa}\right)^\kappa = (e\gamma)^\kappa.$$

Thus, $|HC_f| = O(n^s(e\gamma)^\kappa)$, and this is an upper bound on the time required to execute Step 2 of subalgorithm \mathcal{A}' .

Finally, we need to define the procedure \mathcal{A}'' for finding N_S for a given set S of s co-occurring variables. Fix such an S and let N_γ be a set of at most γ variables in $([n] \setminus S)$ having the following properties:

- P1** In the projection $f_{N_\gamma \leftarrow 0}$ of f in which all of the variables of N_γ are fixed to 0, the variables in S do not co-occur in any term; and
- P2** For every set $N'_\gamma \subset N_\gamma$ such that $|N'_\gamma| = |N_\gamma| - 1$, the variables in S do co-occur in at least one term of $f_{N'_\gamma \leftarrow 0}$.

We will use the following claim (proved in the full version):

Claim. N_S is the union of all sets N_γ of cardinality at most γ that satisfy **P1** and **P2**.

There are only $O(n^\gamma)$ possible candidate sets N_γ to consider, so our problem now reduces to the following: given a set N of at most γ variables, determine whether the variables in S co-occur in $f_{N \leftarrow 0}$.

Recall that since f satisfies the three conditions **C1**, **C2** and **C3**, Lemma 3 implies that $|\hat{f}(S)|$ is either at most \mathcal{Y} (if the variables in S do not co-occur in any term of f) or at least $\frac{\alpha}{2^s} - 2 \cdot \mathcal{Y}$ (if the variables in S do co-occur in some term).

We now claim that the function $f_{N \leftarrow 0}$ has this property as well: i.e., $|\widehat{f_{N \leftarrow 0}}(S)|$ is either at most the same value \mathcal{Y} (if the variables in S do not co-occur in any term of $f_{N \leftarrow 0}$) or at least the same value $\frac{\alpha}{2^s} - 2 \cdot \mathcal{Y}$ (if the variables in S do co-occur in some term of $f_{N \leftarrow 0}$). To see this, observe that the function $f_{N \leftarrow 0}$ is just f with some terms removed. Since each term in f is uniquely satisfied with probability at least α (this is condition **C1**), the same must be true of $f_{N \leftarrow 0}$ since removing terms from f can only increase the probability of being uniquely satisfied for the remaining terms. Since each j -tuple of terms in f is simultaneously satisfied with probability at most β_j (this is condition **C2**), the same must be true for j -tuples of terms in $f_{N \leftarrow 0}$. Finally, for condition **C3**, the value of $\#g_U$ can only decrease in passing from f to $f_{N \leftarrow 0}$. Thus, the upper bound of \mathcal{Y} that follows from applying Lemma 3 to f is also a legitimate upper bound when the lemma is applied to $|\widehat{f_{N \leftarrow 0}}(S)|$, and similarly the lower bound of $\frac{\alpha}{2^s} - 2 \cdot \mathcal{Y}$ is also a legitimate lower bound when the lemma is applied to $f_{N \leftarrow 0}$. Therefore, for every $|N| \leq \gamma$, a sufficiently accurate (within $\Delta/2$) estimate of $|\widehat{f_{N \leftarrow 0}}(S)|$ (as obtained in Step 1 of subalgorithm \mathcal{A}'') can be used to determine whether or not the variables in S co-occur in any term of $f_{N \leftarrow 0}$.

To obtain the required estimate for $|\widehat{f_{N \leftarrow 0}}(S)|$, observe that for a given set N , we can simulate a uniform example oracle for $f_{N \leftarrow 0}$ by filtering the examples from the uniform oracle for f so that only examples setting the variables in N to 0 are accepted. Since $|N| \leq \gamma$, the filter accepts with probability at least $1/2^\gamma$. A Hoeffding bound argument then shows that the Fourier coefficients $|\widehat{f_{N \leftarrow 0}}(S)|$ can be estimated (with probability of failure no more than a small fraction of δ) from an example oracle for f in time polynomial in n , 2^γ , $1/\Delta$, and $\log(1/\delta)$.

Algorithm \mathcal{A}'' , then, estimates Fourier coefficients of restricted versions of f , using a sample size sufficient to ensure that all of these coefficients are sufficiently

accurate over all calls to \mathcal{A}' with probability at least $1 - \delta/3$. These estimated coefficients are then used by \mathcal{A}' to locate the set N_S as just described. The overall algorithm \mathcal{A} therefore succeeds with probability at least $1 - \delta$, and it is not hard to see that it runs in the time bound claimed.

Required parameters. In the above description of Algorithm \mathcal{A} , we assumed that it is given the values of $s, \alpha, \mathcal{Y}, \gamma$, and κ . In fact it is not necessary to assume this; a standard argument gives a variant of the algorithm which succeeds without being given the values of these parameters.

The idea is simply to have the algorithm “guess” the values of each of these parameters, either exactly or to an adequate accuracy. The parameters s, γ and κ take positive integer values bounded by $\text{poly}(n)$. The other parameters α, \mathcal{Y} take values between 0 and 1; a standard argument shows that if approximate values α' and \mathcal{Y}' (that differ from the true values by at most $1/\text{poly}(n)$) are used instead of the true values, the algorithm will still succeed. Thus there are at most $\text{poly}(n)$ total possible settings for $(s, \gamma, \kappa, \alpha, \mathcal{Y})$ that need to be tried. We can run Algorithm \mathcal{A} for each of these candidate parameter settings, and test the resulting hypothesis; when we find the “right” parameter setting, we will obtain a high-accuracy hypothesis (and when this occurs, it is easy to recognize that it has occurred, simply by testing each hypothesis on a new sample of random labeled examples). This parameter guessing incurs an additional polynomial factor overhead. Thus Theorem 2 holds true for the extended version of Algorithm \mathcal{A} that takes only ϵ, δ as input parameters.

4 Random Monotone DNF

The random monotone DNF model. Let $\mathcal{M}_n^{t,k}$ be the probability distribution over monotone t -term DNF induced by the following process: each term is independently and uniformly chosen at random from all $\binom{n}{k}$ monotone ANDs of size exactly k over x_1, \dots, x_n .

Given a value of t , throughout this section we consider the $\mathcal{M}_n^{t,k}$ distribution where $k = \lfloor \log t \rfloor$ (we will relax this and consider a broader range of values for k in Section 6). To motivate this choice, consider a random draw of f from $\mathcal{M}_n^{t,k}$. If k is too large relative to t then a random $f \in \mathcal{M}_n^{t,k}$ will likely have $\Pr_{x \in U_n}[f(x) = 1] \approx 0$, and if k is too small relative to t then a random $f \in \mathcal{M}_n^{t,k}$ will likely have $\Pr_{x \in U_n}[f(x) = 1] \approx 1$; such functions are trivial to learn to high accuracy using either the constant-0 or constant-1 hypothesis. A straightforward analysis (see e.g. [JS06]) shows that for $k = \lfloor \log t \rfloor$ we have that $\mathbf{E}_{f \in \mathcal{M}_n^{t,k}}[\Pr_{x \in U_n}[f(x) = 1]]$ is bounded away from both 0 and 1, and thus we feel that this is an appealing and natural choice.

Probabilistic analysis. In this section we will establish various useful probabilistic lemmas regarding random monotone DNF of polynomially bounded size.

Assumptions: Throughout the rest of Section 4 we assume that $t(n)$ is any function such that $n^{3/2} \leq t(n) \leq \text{poly}(n)$. To handle the case when $t(n) \leq n^{3/2}$, we will use the results from [JS06]. Let $a(n)$ be such that $t(n) = n^{a(n)}$. For

brevity we write t for $t(n)$ and a for $a(n)$ below, but the reader should keep in mind that a actually denotes a function $\frac{3}{2} \leq a = a(n) \leq O(1)$. Because of space limitations all proofs are given in the full version.

The first lemma provides a bound of the sort needed by condition **C3** of Lemma 3:

Lemma 4. *Let $|S| = s = \lfloor a \rfloor + 2$. Fix any $\mathcal{C}_{\mathcal{Y}} \in \mathcal{Y}$. Let $\delta_{\text{terms}} = n^{-\Omega(\log n)}$. With probability at least $1 - \delta_{\text{terms}}$ over the random draw of f from $\mathcal{M}_n^{t,k}$, we have that for some absolute constant c and all sufficiently large n ,*

$$\prod_{U \in \mathcal{C}_{\mathcal{Y}}} (\#g_U) \leq c \cdot \frac{t^{|\mathcal{C}_{\mathcal{Y}}|-1} k^{2^s}}{\sqrt{n}}. \quad (5)$$

The following lemma shows that for f drawn from $\mathcal{M}_n^{t,k}$, with high probability each term is “uniquely satisfied” by a noticeable fraction of assignments as required by condition **C1**. (Note that since $k = O(\log n)$ and $t > n^{3/2}$, we have $\delta_{\text{usat}} = n^{-\Omega(\log \log n)}$ in the following.)

Lemma 5. *Let $\delta_{\text{usat}} := \exp(\frac{-tk}{3n}) + t^2(\frac{k}{n})^{\log \log t}$. For n sufficiently large, with probability at least $1 - \delta_{\text{usat}}$ over the random draw of $f = T_1 \vee \dots \vee T_t$ from $\mathcal{M}_n^{t,k}$, f is such that for all $i = 1, \dots, t$ we have $\Pr_x[T_i \text{ is satisfied by } x \text{ but no other } T_j \text{ is satisfied by } x] \geq \frac{\Theta(1)}{2^k}$.*

We now upper bound the probability that any j distinct terms of a random DNF $f \in \mathcal{M}_n^{t,k}$ will be satisfied simultaneously (condition **C2**). (In the following lemma, note that for $j = \Theta(1)$, since $t = n^{\Theta(1)}$ and $k = \Theta(\log n)$ we have that the quantity δ_{simult} is $n^{-\Theta(\log \log n)}$.)

Lemma 6. *Let $1 \leq j \leq 2^s$, and let $\delta_{\text{simult}} := \frac{t^j e^{jk - \log k} (jk - \log k)^{\log k}}{n^{\log k}}$. With probability at least $1 - \delta_{\text{simult}}$ over the random draw of $f = T_1 \vee \dots \vee T_t$ from $\mathcal{M}_n^{t,k}$, for all $1 \leq \iota_1 < \dots < \iota_j \leq t$ we have $\Pr[T_{\iota_1} \wedge \dots \wedge T_{\iota_j}] \leq \beta_j$, where $\beta_j := \frac{k}{2^{jk}}$.*

Finally, the following lemma shows that for all sufficiently large n , with high probability over the choice of f , every set S of s variables appears in at most γ terms, where γ is independent of n (see condition **C4**).

Lemma 7. *Fix any constant $c > 0$. Let $s = \lfloor a \rfloor + 2$ and let $\gamma = a + c + 1$. Let $\delta_\gamma = n^{-c}$. Then for n sufficiently large, with probability at least $1 - \delta_\gamma$ over the random draw of f from $\mathcal{M}_n^{t,k}$, we have that every s -tuple of variables appears in at most γ terms of f .*

5 Proof of Theorem 1

Theorem 1 [Formally] *Let $t(n)$ be any function such that $t(n) \leq \text{poly}(n)$, let $a(n) = O(1)$ be such that $t(n) = n^{a(n)}$, and let $c > 0$ be any fixed constant. Then for any $n^{-c} < \delta < 1$ and $0 < \epsilon < 1$, $\mathcal{M}_n^{t(n), \lfloor \log t(n) \rfloor}$ is PAC learnable under U_n in $\text{poly}(n^{2a(n)+c+3}, (a(n)+c+1)^{\log t(n)}, t(n), 1/\epsilon, \log 1/\delta)$ time.*

Proof. The result is proved for $t(n) \leq n^{3/2}$ already in [JS06], so we henceforth assume that $t(n) \geq n^{3/2}$. We use Theorem 2 and show that for $s = \lfloor a(n) \rfloor + 2$, random monotone $t(n)$ -term DNFs, with probability at least $1 - \delta$, satisfy conditions **C1–C5** with values $\alpha, \beta_j, \Phi(\cdot), \Delta, \gamma$, and κ such that $\Delta > 0$ and the quantities $n^{s+\gamma}, 1/\Delta$, and γ^κ are polynomial in n . This will show that the extended version of Algorithm \mathcal{A} defined in Section 3 PAC learns random monotone $t(n)$ -term DNFs in time $\text{poly}(n, 1/\epsilon)$. Let $t = t(n)$ and $k = \lfloor \log t \rfloor$, and let f be drawn randomly from $\mathcal{M}_n^{t,k}$. By Lemmas 4–7, with probability at least $1 - \delta_{\text{usat}} - \delta_\gamma - 2^{2^s} \delta_{\text{terms}} - \delta_{\text{simult}}$, f will satisfy **C1–C5** with the following values:

$$\begin{aligned} \mathbf{C1} \quad & \alpha > \frac{\Theta(1)}{2^k}; \quad \mathbf{C2} \quad \beta_j \leq \frac{k}{2^{jk}} \text{ for } 1 \leq j \leq 2^s; \\ \mathbf{C3} \quad & \Phi(\mathcal{C}_{\mathcal{Y}}) \leq O(1) \frac{t^{|\mathcal{C}_{\mathcal{Y}}|-1} k^{2^s}}{\sqrt{n}} \text{ for all } \mathcal{C}_{\mathcal{Y}} \in \mathcal{Y}; \quad \mathbf{C4} \quad \gamma \leq a(n) + c + 1; \\ \mathbf{C5} \quad & \kappa = k = \lfloor \log t \rfloor, \end{aligned}$$

which gives us that $n^{s+\gamma} = n^{2a+c+3}$ and $\gamma^\kappa = (a+c+1)^{\lfloor \log t \rfloor}$. Finally, we show that $\Delta = \Omega(1/t)$ so $1/\Delta$ is polynomial in n :

$$\begin{aligned} \Delta &= \alpha/2^s - 3 \cdot \mathcal{R} = \frac{\Theta(1)}{t2^s} - 3 \sum_{\mathcal{C}_{\mathcal{Y}} \in \mathcal{Y}} 2^s \beta_{|\mathcal{C}_{\mathcal{Y}}|} \Phi(\mathcal{C}_{\mathcal{Y}}) \\ &\geq \frac{\Theta(1)}{t2^s} - \Theta(1) \sum_{\mathcal{C}_{\mathcal{Y}} \in \mathcal{Y}} 2^s \frac{k}{t^{|\mathcal{C}_{\mathcal{Y}}|}} \cdot \frac{t^{|\mathcal{C}_{\mathcal{Y}}|-1} k^{2^s}}{\sqrt{n}} \\ &= \frac{\Theta(1)}{t2^s} - \frac{\Theta(1)k^{2^s+1}}{t\sqrt{n}} = \Omega(1/t). \end{aligned}$$

6 Discussion

Robustness of parameter settings. Throughout Sections 4 and 5 we have assumed for simplicity that the term length k in our random t -term monotone DNF is exactly $\lfloor \log t \rfloor$. In fact, the results extend to a broader range of k 's; one can straightforwardly verify that by very minor modifications of the given proofs, Theorem 1 holds for $\mathcal{M}_n^{t,k}$ for any $(\log t) - O(1) \leq k \leq O(\log t)$.

Relation to previous results. Our results are powerful enough to subsume some known “worst-case” results on learning restricted classes of monotone DNF formulas. Hancock and Mansour [HM91] have shown that read- k monotone DNF (in which each Boolean variable x_i occurs in at most k terms) are learnable under the uniform distribution in $\text{poly}(n)$ time for constant k . Their result extends an earlier result of Kearns *et al.* [KLV94] showing that read-once DNF (which can be assumed monotone without loss of generality) are polynomial-time learnable under the uniform distribution. It is not hard to see that (a very restricted special case of) our algorithm can be used to learn read- k monotone DNF in polynomial time; we give some details in the full version.

References

- [AM02] K. Amano and A. Maruoka. On learning monotone boolean functions under the uniform distribution. In *Proc. 13th ALT*, pages 57–68, 2002.
- [AP95] H. Aizenstein and L. Pitt. On the learnability of disjunctive normal form formulas. *Machine Learning*, 19:183–208, 1995.
- [BBL98] A. Blum, C. Burch, and J. Langford. On learning monotone boolean functions. In *Proc. 39th FOCS*, pages 408–415, 1998.
- [BFJ⁺94] A. Blum, M. Furst, J. Jackson, M. Kearns, Y. Mansour, and S. Rudich. Weakly learning DNF and characterizing statistical query learning using Fourier analysis. In *Proc. 26th STOC*, pages 253–262, 1994.
- [Blu03a] A. Blum. Learning a function of r relevant variables (open problem). In *Proc. 16th COLT*, pages 731–733, 2003.
- [Blu03b] A. Blum. Machine learning: a tour through some favorite results, directions, and open problems. FOCS 2003 tutorial slides, 2003.
- [BT96] N. Bshouty and C. Tamon. On the Fourier spectrum of monotone functions. *Journal of the ACM*, 43(4):747–770, 1996.
- [HM91] T. Hancock and Y. Mansour. Learning monotone k - μ DNF formulas on product distributions. In *Proc. 4th COLT*, pages 179–193, 1991.
- [Jac97] J. Jackson. An efficient membership-query algorithm for learning DNF with respect to the uniform distribution. *JCSS*, 55:414–440, 1997.
- [JS05] J. Jackson and R. Servedio. Learning random log-depth decision trees under the uniform distribution. *SICOMP*, 34(5):1107–1128, 2005.
- [JS06] J. Jackson and R. Servedio. On learning random DNF formulas under the uniform distribution. *Theory of Computing*, 2(8):147–172, 2006.
- [JT97] J. Jackson and C. Tamon. Fourier analysis in machine learning. ICML/COLT 1997 tutorial slides, 1997.
- [KLV94] M. Kearns, M. Li, and L. Valiant. Learning Boolean formulas. *Journal of the ACM*, 41(6):1298–1328, 1994.
- [KMSP94] L. Kučera, A. Marchetti-Spaccamela, and M. Protassi. On learning monotone DNF formulae under uniform distributions. *Information and Computation*, 110:84–95, 1994.
- [KV94] M. Kearns and U. Vazirani. *An introduction to computational learning theory*. MIT Press, Cambridge, MA, 1994.
- [Man95] Y. Mansour. An $O(n^{\log \log n})$ learning algorithm for DNF under the uniform distribution. *JCSS*, 50:543–550, 1995.
- [MO03] E. Mossel and R. O’Donnell. On the noise sensitivity of monotone functions. *Random Structures and Algorithms*, 23(3):333–350, 2003.
- [OS06] R. O’Donnell and R. Servedio. Learning monotone decision trees in polynomial time. In *Proc. 21st CCC*, pages 213–225, 2006.
- [Sel08] L. Sellie. Learning Random Monotone DNF Under the Uniform Distribution. In *Proc. 21st COLT*, to appear, 2008.
- [Ser04] R. Servedio. On learning monotone DNF under product distributions. *Information and Computation*, 193(1):57–74, 2004.
- [SM00] Y. Sakai and A. Maruoka. Learning monotone log-term DNF formulas under the uniform distribution. *Theory of Computing Systems*, 33:17–33, 2000.
- [Val84] L. Valiant. A theory of the learnable. *CACM*, 27(11):1134–1142, 1984.
- [Ver90] K. Verbeugt. Learning DNF under the uniform distribution in quasipolynomial time. In *Proc. 3rd COLT*, pages 314–326, 1990.
- [Ver98] K. Verbeugt. Learning sub-classes of monotone DNF on the uniform distribution. In *Proc. 9th ALT*, pages 385–399, 1998.