



**University of Calgary**

**PRISM: University of Calgary's Digital Repository**

---

Science

Science Research & Publications

---

1995-04-01

# ON THE FOURIER SPECTRUM OF MONOTONE FUNCTIONS

Bshouty, N.; Tamon, T.

---

<http://hdl.handle.net/1880/45749>

unknown

---

*Downloaded from PRISM: <https://prism.ucalgary.ca>*

# On the Fourier Spectrum of Monotone Functions

Nader H. Bshouty    Christino Tamon

Department of Computer Science  
University of Calgary  
2500 University Drive NW  
Calgary, AB Canada T2N 1N4  
e-mail: {bshouty, tamon}@cpsc.ucalgary.ca

## Abstract

In this paper, monotone Boolean functions are studied using harmonic analysis on the cube. The main result is that *any* monotone Boolean function has most of its power spectrum on its Fourier coefficients of “degree” at most  $O(\sqrt{n})$  under *any* product distribution. This is similar to a result of Linial, Mansour, and Nisan [LMN93] which showed that  $AC^0$  functions have almost all of its power spectrum on the coefficients of degree at most  $(\log n)^{O(1)}$  under the uniform distribution. As a consequence of the main result the following two corollaries are obtained:

- For any  $\epsilon > 0$ , monotone Boolean functions are PAC learnable with error  $\epsilon$  under product distributions in time  $2^{\tilde{O}(\frac{1}{\epsilon}\sqrt{n})}$ .
- Any monotone Boolean function can be approximated within error  $\epsilon$  under product distributions by a non-monotone Boolean circuit of size  $2^{\tilde{O}(\frac{1}{\epsilon}\sqrt{n})}$  and depth  $\tilde{O}(\frac{1}{\epsilon}\sqrt{n})$ .

The learning algorithm runs in time subexponential as long as the required error is  $\Omega(1/(\sqrt{n} \log n))$ . It is shown that this is tight in the sense that for any subexponential time algorithm there is a monotone Boolean function for which this algorithm cannot approximate with error better than  $\tilde{O}(1/\sqrt{n})$ .

The main result is also applied to other problems in learning and complexity theory. In learning theory, several polynomial-time algorithms for learning some classes of monotone Boolean functions, such as Boolean functions with  $O(\log^2 n / \log \log n)$  relevant variables, are presented. In complexity theory, some questions regarding monotone NP-complete problems are addressed.

# 1 Introduction

In recent years, harmonic analysis on the cube or the discrete Fourier transform of Boolean functions has emerged as one of the most versatile tools in theoretical computer science. It has found various applications in areas such as circuit complexity [KKL88, BS92], computational learning theory [LMN93, FJS91, KM91, AM91, B92, M92, BFJ+, J94], cryptography [CKM93], and others.

Although harmonic analysis on the cube was originally introduced by Kahn, Kalai, and Linial [KKL88] for the purpose of studying the “influence” of variables on Boolean functions, the first paper which introduced the Fourier transform to learning theory was the beautiful paper of Linial, Mansour, and Nisan [LMN93]. In the latter paper, they proved that  $AC^0$  functions have almost all of its power spectrum on the Fourier coefficients of Hamming weight  $(\log n)^{O(1)}$ . This result led to the PAC-learnability of  $AC^0$  functions in time  $n^{\text{poly}(\log n)}$  under the uniform distribution. The impact of the Fourier transform on learning theory cannot be underestimated judging from the sequence of papers that followed the paper [LMN93]. This technique alone has made some outstanding results possible in recent years culminating in Jackson’s result [J94] on the PAC learnability of DNF formulas under the uniform distribution with membership queries.

In this paper, we study monotone Boolean functions using harmonic analysis on the cube. Our main result is that *any* monotone Boolean function has most of its power spectrum on its Fourier coefficients of *degree* at most  $O(\sqrt{n})$  under *any* product distribution. Based on our main result we derive two important implications in learning theory and circuit complexity.

- Given any  $\epsilon > 0$ , the class of monotone Boolean functions is PAC learnable with error  $\epsilon$  under product distributions in time  $2^{\tilde{O}(\frac{1}{\epsilon}\sqrt{n})}$ .
- Given any  $\epsilon > 0$  and any monotone Boolean function  $f$ , there exists a non-monotone Boolean circuit of size  $2^{\tilde{O}(\frac{1}{\epsilon}\sqrt{n})}$  and depth  $\tilde{O}(\frac{1}{\epsilon}\sqrt{n})$  that approximates  $f$  with error  $\epsilon$  under product distributions.

The time complexity of our learning algorithm is subexponential as long as  $\epsilon = \Omega(1/(\sqrt{n} \log n))$ . We will show that this is the best possible error rate for any subexponential time learning algorithm. To the best of our knowledge, the above learning result is the first subexponential PAC learning algorithm for monotone Boolean functions (even for monotone Boolean functions which require an exponential circuit size) and the second result is the the first approximation result for monotone Boolean functions using non-monotone Boolean circuit of subexponential size and sublinear depth.

We also introduce and study a new measure of complexity for probability distributions called the *convex dimension* of a distribution. We prove that if a concept class is approximable under a collection of distributions then it is also approximable under any distribution that belongs to the convex combination of that collection. The convex dimension  $cdim(D)$  of a distribution  $D$  is the minimal number of product distributions such that  $D$  is in their convex combination. We note that our previous results hold also for any distribution  $D$  modulo a complexity factor of  $cdim(D)$  and in particular, our algorithm is subexponential for any distribution  $D$  with a subexponential convex dimension.

Further applications of the main result include some improvements on two efficient, i.e., polynomial time, learning algorithms for monotone Boolean functions. Kearns and Valiant [KV89] proved that any monotone Boolean function is efficiently weakly PAC learnable with error  $\frac{1}{2} - \frac{1}{2n}$  under the uniform distribution. We improve their result by showing that there is an efficient weak PAC learning algorithm under the uniform distribution with error  $\frac{1}{2} - \Omega(\frac{\log^2 n}{n})$ . Our result is based on a direct application of a result due to Kahn, Kalai, and Linial [KKL88]. Furthermore, under any

product distribution, we show that there is an efficient weak PAC learning algorithm with error  $\frac{1}{2} - \frac{c}{n}$ , for any constant  $c$ .

Our second improvement is on Sakai and Maruoka's [SM94] PAC learnability result of monotone  $O(\log n)$ -term DNF under the uniform distribution. We improve their result in two ways. First, we prove that their result extends to any constant-bounded product distribution, and second, we give an extension to a more general concept class. In particular, let  $A(k)$  be the class of Boolean functions of the form  $f(T_1, \dots, T_k)$  where  $f$  is a monotone Boolean function on  $k$  inputs and each  $T_i$  is a monotone conjunction or disjunction. Our other results state that there are polynomial time PAC learning algorithms for  $A(\log n)$ , for  $A(O(\log^2 n / (\log \log n)^3))$ , and for monotone Boolean functions that depend on  $O(\log^2 n / \log \log n)$  relevant variables (the last two results require the error  $\epsilon$  to be constant).

Finally we apply our results to monotone graph properties. We show that there is a Boolean circuit of size  $2^{\tilde{O}(\frac{n}{\epsilon} \sqrt{r(n)})}$  and depth  $\tilde{O}(\frac{n}{\epsilon} \sqrt{r(n)})$  that approximates within error  $\epsilon$  any monotone graph property with a threshold function of  $r(n)$ . For example, there is a Boolean circuit of size  $n^{O(\frac{1}{\epsilon} \sqrt{\log n})}$  that approximates to within error  $\epsilon$  the Hamiltonian property on random graphs  $G(n, p)$ , for any  $p$ . We also discuss the connection of this work with a related result [BFF87].

All of the algorithms considered in this papers fall in the category of a statistical query algorithm, and hence, by the result of Kearns [K93], are noise-tolerant.

The paper is organized as follows: Section 2 is devoted to a description of the learning models considered in this paper. Section 3 describes notations and the basic theory of the Fourier transform for Boolean functions. It also contains some basic facts that will be required in proving some of the later results. In Section 4, we prove our main spectral characterization of Boolean functions that is based on influence and the average sensitivity. The next section, Sections 5, is devoted to learning monotone Boolean functions. This section also contains some lower bound results that shows near optimality of the learning results. Section 6 describes the application of the main result to the circuit approximation of monotone Boolean functions whereas Section 7 considers the approximation of monotone graph properties. Finally, Section 8 discusses applications of the main result in deriving some efficient learning algorithms.

## 2 The Learning Model

The learning model considered in this paper is the *Probably Approximately Correct* (PAC) learning model introduced by Valiant [V84] and its weak variant introduced by Kearns and Valiant [KV89]. Let  $C_n$  be a class of Boolean functions over  $n$  variables, let  $D$  be a probability distribution over  $\{0, 1\}^n$ , and let  $f \in C_n$  be a target function. The learning algorithm has access to an example oracle  $EX(D, f)$  which generates random labeled examples  $(a, f(a))$ , where  $a \in \{0, 1\}^n$  is drawn according the distribution  $D$ . Given any positive  $\epsilon$  and  $\delta$ , after observing some random examples, the learning algorithm must output a hypothesis  $h$  that satisfies

$$\Pr[D(h\Delta f) \leq \epsilon] \geq 1 - \delta,$$

where  $h\Delta f = \{x : h(x) \neq f(x)\}$ . The above probability is taken with respect to the random examples seen by the algorithm and some internal randomization in the algorithm. The running time of the learner will depend on  $n, 1/\epsilon, 1/\delta$ , and the size  $s(f)$  of  $f$  (under some predetermined representation). If there is such a learning algorithm that succeeds for all  $f \in C_n$  then  $C_n$  is PAC learnable under distribution  $D$ . We sometimes refer to such a learning algorithm as an  $(\epsilon, \delta)$  PAC learning algorithm. We say that  $C_n$  is weakly PAC learnable under distribution  $D$  if there is a fixed polynomial  $p$  and a learning algorithm that succeeds for an error  $\epsilon = \frac{1}{2} - \frac{1}{p(n, s(f))}$ .

### 3 Preliminaries

In this section we review some notation and some standard facts about the discrete Fourier transform of Boolean functions.

When we write  $\log$  we mean  $\log_2$ . We will use the shorthand  $[n]$  for the set  $\{1, 2, \dots, n\}$  and the Iversonian  $I[\text{statement}]$  to mean 1 if the statement is true and 0 otherwise. For  $a \in \{0, 1\}^n$ , let  $a_i$  denote the  $i$ -th bit of  $a$ . The vector  $e_i \in \{0, 1\}^n$  denotes the vector with all zeros except for the  $i$ -th bit which is 1. The Hamming weight of  $a$ , i.e. the number of ones in  $a$ , is denoted by  $|a|$ .

Let  $f : \{0, 1\}^n \rightarrow \{-1, +1\}$  be a Boolean function. Let  $D$  be a product distribution over  $\{0, 1\}^n$  with  $\Pr[x_i = 1] = \mu_i$ . Thus for  $a \in \{0, 1\}^n$  we have the distribution of  $a$  is  $D(a) = \prod_{a_i=1} \mu_i \prod_{a_i=0} (1 - \mu_i)$ . The distribution  $D$  is called *constant bounded* if there exists a constant  $c \in (0, 1)$  independent of  $n$  such that for all  $i$  we have  $\mu_i \in [c, 1 - c]$ . The standard deviation of  $x_i$  is defined as

$$\sigma_i = \sqrt{\mu_i(1 - \mu_i)}.$$

The *influence* of variable  $x_i$  on  $f$  (see [KKL88, HM91]) over a product distribution  $D$  is defined as the probability that  $f(x)$  differs from  $f(x \oplus e_i)$  when  $x$  is chosen according to  $D$ . Here  $x \oplus e_i$  means  $x$  with its  $i$ -th bit flipped. We will use the notation  $I_{D,i}(f)$  to denote the above probability. Often we will use the restriction notation of functions,  $f_0 = f|_{x_i=0}$  and  $f_1 = f|_{x_i=1}$ . With this notation we have that for any Boolean function  $f$

$$I_{D,i}(f) = E_D[I[f(x) \neq f(x \oplus e_i)]] = \frac{1}{2}E_D[|f_1 - f_0|] = \frac{1}{4}E_D[(f_1 - f_0)^2].$$

If  $f$  is a monotone Boolean function, i.e.,  $f_0 \leq f_1$  always, then this simplifies to

$$I_{D,i}(f) = \frac{1}{2}E_D[f_1 - f_0].$$

To facilitate stating some of our results we introduce the following notion of *influence norm*  $I_D(f)$  of  $f$  with respect to a product distribution  $D$ :

$$I_D(f) = \sqrt{\sum_{i=1}^n (2\sigma_i I_{D,i}(f))^2}.$$

The *sensitivity* of  $f$  at a point  $a \in \{0, 1\}^n$ , denoted by  $s_a(f)$ , is defined as the number of neighbors of  $a$  (in the standard  $n$ -cube ordering) whose  $f$ -values differ from  $f(a)$ . More formally,

$$s_a(f) = |\{i \in [n] : f(a \oplus e_i) \neq f(a)\}| = \sum_{i=1}^n I[f(a) \neq f(a \oplus e_i)].$$

The *average sensitivity* of  $f$  with respect to a product distribution  $D$ , denoted  $s_D(f)$ , is defined as

$$s_D(f) = E_D[s_x(f)].$$

It is well-known that the average sensitivity is equivalent to the sum of the influences, as seen from the following simple derivation:

$$s_D(f) = E_D[s_x(f)] = \sum_{i=1}^n E_D[I[f(x) \neq f(x \oplus e_i)]] = \sum_{i=1}^n I_{D,i}(f).$$

The Fourier transform of Boolean functions over product distribution is defined as follows (see [FJS91]). First we define the inner product over the  $2^n$ -dimensional vector space of all real-valued functions over  $\{0, 1\}^n$  as follows:

$$(f, g)_D = \sum_x D(x) f(x) g(x) = E_D[fg].$$

Now let  $z_i(x) = (\mu_i - x_i)/\sigma_i$ . Note that  $z_i$  has mean zero and variance one (i.e. it is standard normal). Next we define the basis function

$$\phi_a(x) = \prod_{a_i=1} z_i(x).$$

These basis functions satisfy the following properties.

1. *decomposable* :  $\phi_{ab}(xy) = \phi_a(x)\phi_b(y)$ , where  $xy$  is the concatenation of strings  $x$  and  $y$  (possibly of different lengths).

2. *orthonormal* :

$$(\phi_a, \phi_b)_D = \begin{cases} 1 & \text{if } a \neq b \\ 0 & \text{otherwise} \end{cases}$$

Given the orthonormality of these  $\phi_a$ 's we get the Fourier representation of any Boolean function  $f$  as

$$f = \sum_a \tilde{f}(a) \phi_a,$$

where  $\tilde{f}(a) = (f, \phi_a)_D = E_D[f\phi_a]$ . Also because of orthonormality we have Parseval's equation:

$$1 = E_D[f^2] = \sum_a \tilde{f}^2(a).$$

Finally note  $\tilde{f}(0_n) = E_D[f]$ , where  $0_n$  is the  $n$ -bit all-zero vector.

For the case of  $D$  being the uniform distribution, the following notations are used:  $\chi_a$  in place of  $\phi_a$  and  $\hat{f}(a)$  in place of  $\tilde{f}(a)$ . Note that in this case,  $\mu_i = \sigma_i = 1/2$ , for all  $i \in [n]$ .

In most cases we will appeal to the following version of Chernoff-Hoeffding bounds (see [HR89, M95]).

**Theorem 3.1** (Chernoff-Hoeffding bounds)

Let  $x_1, \dots, x_m$  be independent identically distributed random variables with  $E[x_i] = p$ ,  $|x_i| \leq B$ , and let  $s_m = x_1 + \dots + x_m$ . Then

$$m \geq \frac{2B^2}{\epsilon^2} \ln \frac{2}{\delta} \quad \text{implies} \quad \Pr \left[ \left| \frac{s_m}{m} - p \right| > \epsilon \right] \leq \delta.$$

## 4 Spectral Lemmas

We are now ready to a lemma which relates the influence and the Fourier transform of Boolean functions. The next lemma is a folklore result whose proof we include for completeness.

**Lemma 1** For any Boolean function  $f$ , for any product distribution  $D$  and for any  $i \in [n]$ ,

$$4\sigma_i^2 I_{D,i}(f) = \sum_{a:a_i=1} \tilde{f}^2(a).$$

*Proof* Without loss of generality, let  $i = 1$ . First recall that  $I_{D,1}(f) = \frac{1}{4}E_D[(f_0 - f_1)^2]$ . Applying Parseval to the right hand side gives

$$I_{D,1}(f) = \frac{1}{4} \sum_{b \in \{0,1\}^{n-1}} (\widetilde{f_0 - f_1})^2(b) = \frac{1}{4} \sum_{b \in \{0,1\}^{n-1}} (\tilde{f}_0(b) - \tilde{f}_1(b))^2.$$

We now find some relation for  $\tilde{f}_0$  and  $\tilde{f}_1$ . Recall that

$$f(x) = \sum_a \tilde{f}(a) \phi_a(x) = \sum_{0b} \tilde{f}(0b) \phi_b(y) + \sum_{1b} \tilde{f}(1b) \phi_b(y) \frac{\mu_1 - x_1}{\sigma_1}.$$

The last step uses the decomposable property of  $\phi_a$ . From this we will get

$$f_0 \equiv f|_{x_1=0} = \sum_b \left( \tilde{f}(0b) + \frac{\mu_1}{\sigma_1} \tilde{f}(1b) \right) \phi_b(y),$$

and

$$f_1 \equiv f|_{x_1=1} = \sum_b \left( \tilde{f}(0b) - \frac{(1-\mu_1)}{\sigma_1} \tilde{f}(1b) \right) \phi_b(y).$$

This implies

$$\tilde{f}_0(b) = \tilde{f}(0b) + \frac{\mu_1}{\sigma_1} \tilde{f}(1b), \quad \text{and} \quad \tilde{f}_1(b) = \tilde{f}(0b) - \frac{(1-\mu_1)}{\sigma_1} \tilde{f}(1b).$$

So continuing with  $I_{D,1}(f)$ .

$$\begin{aligned} I_{D,1}(f) &= \frac{1}{4} \sum_b (\tilde{f}_0(b) - \tilde{f}_1(b))^2 \\ &= \frac{1}{4} \sum_b \left( \tilde{f}(0b) + \frac{\mu_1}{\sigma_1} \tilde{f}(1b) - \tilde{f}(0b) + \frac{(1-\mu_1)}{\sigma_1} \tilde{f}(1b) \right)^2 \\ &= \frac{1}{4} \sum_b \left( \frac{\mu_1 + (1-\mu_1)}{\sigma_1} \right)^2 \tilde{f}^2(1b) = \sum_b \left( \frac{\tilde{f}(1b)}{2\sigma_1} \right)^2. \end{aligned}$$

□

The following definition extends the definition of Hamming weight of a Boolean vector to the case of product distributions.

**Definition 1** (Hamming weight)

Let  $D$  be a product distribution. We define the weight of  $a \in \{0,1\}^n$  under  $D$  to be

$$\|a\|_D = \log \prod_{a_i=1} \frac{1}{\sigma_i}.$$

Note that  $\|a\|_D \geq |a|$  for any product distribution, and equality is attained precisely when  $D$  is the uniform distribution. When the context is clear we will drop the subscript  $D$  from  $\|a\|_D$  and just write  $\|a\|$ .

**Theorem 4.1** (Main Theorem)

For any product distribution  $D$ , for any Boolean function  $f$ , for all positive integer  $k$ ,

$$\sum_{\|a\| \geq k} \tilde{f}^2(a) \leq \frac{2}{k} I_D(f) \sqrt{\sum_{i=1}^n \left( \sigma_i \log \frac{1}{\sigma_i} \right)^2} \leq 1.062 \frac{\sqrt{n}}{k} I_D(f).$$

*Proof* Note that from Lemma 1

$$\begin{aligned}
\sum_{i=1}^n 4\sigma_i^2 I_{D,i}(f) \log \sigma_i^{-1} &= \sum_{i=1}^n \sum_{a: a_i=1} \log \sigma_i^{-1} \tilde{f}^2(a) \\
&= \sum_{a \in \{0,1\}^n} \tilde{f}^2(a) \sum_{i: a_i=1} \log \sigma_i^{-1} \\
&= \sum_{a \in \{0,1\}^n} \|a\| \tilde{f}^2(a).
\end{aligned}$$

Recall that the Cauchy-Schwartz inequality is

$$\left( \sum_{i=1}^n a_i b_i \right)^2 \leq \left( \sum_{i=1}^n a_i^2 \right) \left( \sum_{i=1}^n b_i^2 \right).$$

Now we let  $a_i = 2\sigma_i I_{D,i}(f)$  and  $b_i = 2\sigma_i \log \sigma_i^{-1}$  to get

$$\begin{aligned}
I_D(f)^2 &= \sum_{i=1}^n 4\sigma_i^2 I_{D,i}(f)^2 \\
&\geq \frac{(\sum_{i=1}^n 4\sigma_i^2 I_{D,i}(f) \log \sigma_i^{-1})^2}{4 \sum_{i=1}^n (\sigma_i \log \sigma_i^{-1})^2}, \quad \text{by Cauchy-Schwartz} \\
&= \frac{1}{4 \sum_{i=1}^n (\sigma_i \log \sigma_i^{-1})^2} \left( \sum_a \|a\| \tilde{f}^2(a) \right)^2 \\
&\geq \frac{1}{4 \sum_{i=1}^n (\sigma_i \log \sigma_i^{-1})^2} \left( k \sum_{\|a\| \geq k} \tilde{f}^2(a) \right)^2
\end{aligned}$$

for any positive integer  $k$ , which proves the first inequality. The second inequality can be seen using simple calculus since  $(x \log x^{-1})^2 \leq e^{-2} \log^2 e \leq 0.2817$ , for all  $x \in [0, 1/2]$  (simply note that  $2\sqrt{0.2817} = 1.062$ ).  $\square$

An alternative way of viewing and deriving the above result is via the notion of average sensitivity. On occasion we will drop the subscript  $D$  from  $s_D(f)$  when the product distribution  $D$  is clear from context.

**Theorem 4.2** (Alternative Main Theorem)

For any Boolean function  $f$  and any product distribution  $D$

$$\sum_a \{ \tilde{f}(a)^2 : \|a\|_D \geq (6/5)s_D(f)/\epsilon \} \leq \epsilon.$$

*Proof* We start with the identity

$$\sum_a \|a\| \tilde{f}(a)^2 = \sum_i 4\sigma_i (\sigma_i \log \frac{1}{\sigma_i}) I_{D,i}(f),$$

and use the fact  $\sigma \leq 1/2$  and  $x \log \frac{1}{x} \leq 3/5$  (not the best), for  $x \in [0, 1]$ , to get

$$\sum_a \|a\| \tilde{f}(a)^2 \leq \frac{6}{5} s(f)$$



(since  $s(f) = \sum_i I_{D,i}(f)$ ). But the left-hand side is bounded from below by  $k \sum_{a: \|a\| \geq k} \tilde{f}(a)^2$  for any  $k$ . So in particular choose  $k = \frac{6}{5}s(f)/\epsilon$ .  $\square$

The following lemma states a key link between the Fourier spectrum and the influences in monotone Boolean functions.

**Lemma 2** *For any monotone Boolean function  $f$ , for any product distribution  $D$ , and for any  $i \in [n]$*

$$\tilde{f}(e_i) = -2\sigma_i I_{D,i}(f).$$

*Proof* Let  $D_i$  be the induced distribution over all the variables except  $x_i$ . We have the following derivation.

$$\begin{aligned} \tilde{f}(e_i) &= E_D[f\phi_{e_i}] \\ &= E_{D_i} E_{x_i}[fz_i] \\ &= E_{D_i} \left[ (1 - \mu_i) \frac{\mu_i}{\sigma_i} f_0 + \mu_i \frac{\mu_i - 1}{\sigma_i} f_1 \right] \\ &= E_{D_i} \left[ \frac{\mu_i(1 - \mu_i)}{\sigma_i} (f_0 - f_1) \right] \\ &= \sigma_i E_{D_i}[f_0 - f_1]. \end{aligned}$$

Now recall that for monotone Boolean functions  $I_{D,i}(f) = \frac{1}{2} E_{D_i}[f_1 - f_0]$ .  $\square$

## 5 Learning Monotone Boolean Functions

The main results of this section are a subexponential time PAC learning algorithm for any monotone Boolean function and some nearly matching lower bounds on the error rate and some other parameters. But first we review the connections between Fourier transform and PAC learning that were first given in [LMN93, BFJ+] (see also Mansour's excellent survey [M94]).

**Fact 1** (PAC Learning and Fourier Spectrum)

*Let  $D$  be a product distribution, let  $f \in \{-1, +1\}$  be a Boolean function, and let  $A \subset \{0, 1\}^n$  be a set of assignments. The real-valued function  $g(x) = \sum_{a \in A} \tilde{f}(a) \phi_a(x)$  satisfies*

$$\Pr_D[f \neq \text{sgn}(g)] \leq E_D[(f - g)^2] = \sum_{a \notin A} \tilde{f}^2(a),$$

where  $\text{sgn}(g)(x) = (-1)^{I[g(x) < 0]}$  returns the sign of  $g(x)$ .

*Given  $g$ , the randomized Boolean function  $h$  defined as (see [BFJ+])*

$$h(x) = \begin{cases} -1 & \text{with probability } p(x) \\ +1 & \text{with probability } 1 - p(x) \end{cases}$$

where  $p(x) = \frac{(1-g(x))^2}{2(1+g^2(x))}$ , satisfies a slightly better error bound

$$\Pr_D[f \neq h] \leq \frac{1}{2} E_D[(f - g)^2] = \frac{1}{2} \sum_{a \notin A} \tilde{f}^2(a).$$

Using this fact, Boolean functions can be learned by collecting the Fourier coefficients  $\tilde{f}(a)$  for all  $a \in A$ . This can be done by finding an approximation  $h_a$  of  $E_D[f\phi_a] = \tilde{f}(a)$  that satisfies

$$|h_a - \tilde{f}(a)| \leq \sqrt{\epsilon/(2|A|)}$$

with confidence  $1 - \delta/|A|$ , for every  $a \in A$ . Then we define the hypothesis  $h = \sum_{a \in A} h_a \phi_a$ . This hypothesis will be an approximation of  $f$  that satisfies

$$\Pr_D[f(x) \neq h(x)] \leq \sum_{a \notin A} \tilde{f}^2(a) + \frac{\epsilon}{2}$$

with probability at least  $1 - \delta$ . Now if  $\sum_{a \notin A} \tilde{f}^2(a) \leq \epsilon/2$  then the hypothesis  $h$  will satisfy

$$\Pr_D[f(x) \neq h(x)] \leq \epsilon.$$

To approximate the Fourier coefficients we use sampling to find  $E_D[f\phi_a]$  for all  $a \in A$ . By the Chernoff-Hoeffding bounds (Theorem 3.1), if  $|f\phi_a| = |\phi_a| \leq B$  then we will need a sample of size at least

$$\frac{4B^2|A|}{\epsilon} \ln \frac{|A|}{\delta}. \quad (1)$$

When  $A$  is the set of all assignments of Hamming weight less than or equal to  $k$ , the above algorithm is called the *k-lowdegree Fourier algorithm*. In this case the number of coefficients that the algorithm must estimate is

$$|A| = \sum_{i=0}^k \binom{n}{i} \leq \left(\frac{n\epsilon}{k}\right)^k = 2^{k \log(n\epsilon/k)}. \quad (2)$$

So  $\exp(k \log(n\epsilon/k))$  will dictate the running time of the *k-lowdegree Fourier algorithm*.

## 5.1 A Subexponential Time Algorithm

The following theorems will show that monotone Boolean functions are PAC learnable under product distributions in subexponential time. For ease of analysis we will assume that the learner knows the product distribution, i.e., the parameters  $\mu_i$  is exactly known, for all  $i \in [n]$ . Later in a separate subsection we will discuss the case when the learner does not know the distribution and show that this only incurs a  $\log n$  blow up in the exponent of the time complexity. Since in most cases, we deal with  $k = \omega(\log n)$ , the time complexity is still  $2^{\tilde{O}(k)}$ .

**Theorem 5.1** *For any  $\epsilon, \delta > 0$ , any monotone Boolean function is PAC learnable under any product distribution with error  $\epsilon$  and confidence  $1 - \delta$  in time*

$$2^{O(\frac{1}{\epsilon} \sqrt{n} \log(\epsilon \sqrt{n}))}.$$

*Proof* We will use the *k-lowdegree Fourier algorithm* with

$$k = 1.062 I_D(f) \frac{\sqrt{n}}{(\epsilon/2)}$$

and with the hypothesis set to  $h = \sum_{||a|| \leq k} h_a \phi_a$ . By Theorem 4.1,  $h$  is an  $(\epsilon/2)$ -approximation of  $f$ . Since  $||a|| \geq |a|$ , we have that

$$\{a : ||a|| \leq k\} \subseteq \{a : |a| \leq k\}$$

and hence the  $k$ -lowdegree Fourier algorithm only need to collect (estimate) the Fourier coefficients of Hamming weight at most  $k$ . From the definitions of  $\|a\|$  and  $\phi_a$ , if  $\|a\| \leq k$  we get that

$$|\phi_a| = \left| \prod_{a_i=1} \frac{\mu_i - x_i}{\sigma_i} \right| = 2^{\|a\|} \left| \prod_{a_i=1} (\mu_i - x_i) \right| \leq 2^k,$$

since  $|\mu_i - x_i| \leq 1$ . Now by (1),(2) and the above, the algorithm outputs a hypothesis that is an approximation of  $f$  to within error  $\epsilon$  with sample size and time complexity of

$$2^{O\left(\frac{\sqrt{n}I_D(f)}{\epsilon} \ln \frac{\sqrt{n}\epsilon}{I_D(f)}\right)} \log \frac{1}{\delta}. \quad (3)$$

By Lemma 2, we note that  $I_D(f) \leq 1$ , for any monotone function  $f$ , because

$$I_D(f)^2 = \sum_i (2\sigma_i I_{D,i}(f))^2 = \sum_i \tilde{f}(e_i)^2 \leq 1$$

by Parseval's. Thus we have  $I_D(f) \log \frac{1}{I_D(f)} \leq 1$ , and therefore

$$I_D(f) \log \frac{\epsilon\sqrt{n}}{I_D(f)} = I_D(f) \log \frac{1}{I_D(f)} + I_D(f) \log(\epsilon\sqrt{n}) = O(\log(\epsilon\sqrt{n})).$$

This analysis proves the time complexity stated in the theorem.  $\square$

*Remark.* Note that using the above algorithm with subexponential time, the best achievable error rate is  $\epsilon = \frac{1}{\sqrt{n}}$ . In the next subsection we show that this is the best possible error rate up to a  $O(\log n)$  factor.

We can alternatively derive the above theorem using the average sensitivity instead of the influence norm, i.e., apply Theorem 4.2 instead of Theorem 4.1. Note that in this case,  $s_D(f) \leq \sqrt{n}$ , for any monotone Boolean function  $f$ .

## 5.2 Learning when the Product Distribution is Unknown

In this subsection we address the issue of learning when the parameters of the product distribution, i.e., the  $\mu_i$ 's, are unknown. First we argue that we may ignore all  $\mu_i$ 's that are less than  $n^{-2}$  or greater than  $1 - n^{-2}$  since this will add only an additive factor of  $n^{-1}$  to the final error. Formally, let  $A \subseteq \{0, 1\}^n$  be the set of all *good* assignments  $x$ , i.e., ones that satisfy  $x_i = 0$ , for all  $\mu_i < n^{-2}$ , and  $x_i = 1$ , for all  $\mu_i > 1 - n^{-2}$ . Note that the probability of an example  $x$  being good is at least  $1 - n^{-1}$ . Suppose that  $h$  is a hypothesis that approximates the target  $f$  quite well on examples from  $A$ . Then

$$\Pr_D[h(x) \neq f(x)] \leq \Pr[h(x) \neq f(x) | x \text{ good}] + \Pr[x \text{ not good}] \leq \epsilon + n^{-1}.$$

So we may assume that  $n^{-2} < \mu_i < 1 - n^{-2}$ , for all  $i \in [n]$ .

We will estimate each  $\mu_i$  up to an error of  $M^{-(\log n + 4)}$  (using Chernoff-Hoeffding bounds this takes only  $\text{poly}(M^{\log n})$  time), where, setting  $k = O(\sqrt{n}/\epsilon)$ ,

$$M = 2^k = 2^{O(\sqrt{n}/\epsilon)}.$$

We will consider three quantities

$$h_A = \sum_{a \in S} E_D[f\phi'_a]\phi'_a, \quad h_B = \sum_{a \in S} E_{D'}[f\phi'_a]\phi'_a, \quad h_C = \sum_{a \in R} E_D[f\phi_a]\phi_a,$$

where  $D'$  is the estimated distribution (using the estimated  $\mu_i$ 's) and  $\phi'_a$ 's are the basis functions according to  $D'$ , and  $R$  and  $S$  are the sets of assignments for which the algorithm needs to estimate the  $\phi_a$ 's and  $\phi'_a$ 's, respectively. Notice that from the proof of Theorem 5.1 and since  $I_D(f) \leq 1$ , for all monotone Boolean function  $f$ , we have

$$|\phi'_a| \leq M.$$

We wanted to learn  $h_C$  but because only an approximation of  $\mu_i$  can be found we will try to learn  $h_B$ . Now since the example oracle gives the examples according to the distribution  $D$  and not  $D'$  we will instead learn  $h_A$ . Since the learning parameters (the number of coefficients) are done for  $D'$  we have

$$E_{D'}[(h_B - f)^2] \leq \epsilon.$$

Suppose  $\mu_i + \tau_i$  is the estimation for  $\mu_i$ , where  $|\tau_i| < M^{-(\log n + 4)}$ . Notice that

$$\begin{aligned} \frac{D'(x)}{D(x)} &= \prod_{x_i=1} \left(1 + \frac{\tau_i}{\mu_i}\right) \prod_{x_i=0} \left(1 - \frac{\tau_i}{1 - \mu_i}\right) \\ &= \left(1 + O\left(\frac{n^2}{M^{\log n + 4}}\right)\right)^n \\ &= 1 + O\left(\frac{1}{M^{\log n + 3}}\right). \end{aligned}$$

Now

$$\begin{aligned} E_D[(h_B - f)^2] &= E_{D'} \left[ \frac{D}{D'} (h_B - f)^2 \right] \\ &\leq E_{D'}[(h_B - f)^2] \max \frac{D}{D'} \\ &\leq \epsilon \left(1 + O\left(\frac{1}{M^{\log n + 3}}\right)\right) \\ &\leq \epsilon + O\left(\frac{\epsilon}{M^{\log n + 3}}\right). \end{aligned}$$

Therefore  $h_B$  is also a good approximation of  $f$  with respect to the distribution  $D$ . Now we show that  $h_A$  is good enough. We have

$$\begin{aligned} |h_A - h_B| &= \left| \sum_{a \in S} E_D[f \phi'_a] \phi'_a - \sum_{a \in S} E_{D'}[f \phi'_a] \phi'_a \right| \\ &= \left| \sum_{a \in S} E_D \left[ f \phi'_a \left(1 - \frac{D'}{D}\right) \right] \phi'_a \right| \\ &\leq |S| M^2 \left| 1 - \min \frac{D'}{D} \right| \\ &\leq O\left(\frac{1}{M}\right), \end{aligned}$$

because

$$|S| \leq \sum_{i=0}^k \binom{n}{i} \leq n^k = 2^{k \log n} = M^{\log n}.$$

Therefore

$$\begin{aligned}
E_D[(h_A - f)^2] &\leq 2E_D[(h_A - h_B)^2 + (h_B - f)^2] \\
&= 2E_D[(h_A - h_B)^2] + 2E_D[(h_B - f)^2] \\
&\leq 2\epsilon + O\left(\frac{1}{M}\right).
\end{aligned}$$

This completes the analysis for learning when the product distribution is unknown.

### 5.3 Lower bounds

In this subsection we give several lower bounds showing that our algorithm is nearly optimal in terms of running time complexity and error rate. The following theorem shows that the error rate achieved in our algorithm is the best possible for a subexponential time algorithm.

**Theorem 5.2** *Any learning algorithm which PAC-learns any monotone Boolean function under the uniform distribution and which runs in subexponential time (even with time  $2^{cn}$ , for any  $c < 1$ ) will output an approximation with an error of at least  $\Omega\left(\frac{1}{\sqrt{n} \log n}\right)$ .*

*Proof* There are at least  $m(n) = 2^{\binom{n}{2}} \geq 2^{d2^n \sqrt{n}}$ , monotone Boolean functions over  $n$  variables, for some constant  $d < 1$  (see [W87]). Suppose  $A$  is the  $\epsilon$ -approximation algorithm for any monotone Boolean function. If  $A$  outputs a hypothesis  $h$  then  $h$  can  $\epsilon$ -approximate at most

$$k(n) = \sum_{i \leq \epsilon 2^n} \binom{2^n}{i} \leq 2^{2^n \epsilon \log(e/\epsilon)}$$

Boolean functions. Assuming  $A$  runs in time  $2^{cn}$ , for some constant  $c < 1$ , then  $A$  can output at most  $2^{2^{cn}}$  possible hypotheses. Therefore we must have  $2^{2^{cn}} k(n) \geq m(n)$  which implies

$$2^{cn} + \epsilon 2^n \log \frac{e}{\epsilon} \geq \frac{d2^n}{\sqrt{n}}.$$

This implies  $\epsilon = \Omega\left(\frac{1}{\sqrt{n} \log n}\right)$ .  $\square$

The next corollary gives a lower bound for the error rate of any learning algorithm that runs in time bounded by  $2^{\sqrt{n}}$ .

**Corollary 1** *Any learning algorithm for monotone Boolean functions under the uniform distribution with a running time bounded by  $2^{\sqrt{n}}$  cannot achieve an error smaller than*

$$\Omega\left(\frac{1}{n^{1/4} \log n}\right).$$

*Proof* Let  $\mathcal{A}$  be an algorithm that runs in time  $2^{\sqrt{n}}$  and achieves error  $\epsilon(n)$ , for any  $n$ . We will construct an algorithm  $\mathcal{B}$  that learns any monotone Boolean function over  $n$  variables in time  $2^{cn}$ , for some  $c < 1$ , and achieves an error of  $2\epsilon((cn)^2)$ . Now since by Theorem 5.2

$$2\epsilon((cn)^2) = \Omega\left(\frac{1}{\sqrt{cn} \log cn}\right)$$

we get the result of the corollary.

The algorithm  $\mathcal{B}$  with input  $f(x_1, \dots, x_n)$  will *pretend* that the input is over  $(cn)^2$  variables and runs the algorithm  $\mathcal{A}$  to learn the function  $f$ . The error rate achievable by this algorithm is  $\epsilon = \epsilon((cn)^2)$ . This algorithm outputs a hypothesis  $h$  that satisfies

$$E_{x_1, \dots, x_{(cn)^2}} [f \neq h] < \epsilon.$$

This does not imply that  $h$  is a good approximation to  $f$  under the original domain  $\{x_1, \dots, x_n\}$ . We now proceed by randomly and uniformly choosing values  $x_{n+1}^0, \dots, x_{(cn)^2}^0$  and return the hypothesis

$$h(x_1, \dots, x_n, x_{n+1}^0, \dots, x_{(cn)^2}^0).$$

Since the expectation is over the uniform distribution we have

$$E_{x_1, \dots, x_{(cn)^2}} [f \neq h] = E_{x_{n+1}, \dots, x_{(cn)^2}} E_{x_1, \dots, x_n} [f \neq h] < \epsilon,$$

and therefore with probability at least  $1/2$  (by Markov's inequality) a random  $x_{n+1}^0, \dots, x_{(cn)^2}^0$  gives a  $2\epsilon = 2\epsilon((cn)^2)$  approximation to  $f$ .  $\square$

We now investigate the best error rate of the low-degree algorithm. The first theorem shows that there exists a monotone Boolean function such that to approximate this function within error  $\epsilon$  using its low-degree Fourier coefficients, for  $\epsilon = n^{-1/2}$ , we need to collect all coefficients of weight less or equal to  $cn$ , for some constant  $c < 1$ . The second theorem shows that to approximate the majority function with the same error we need to collect all of its Fourier coefficients of weight up to  $O(\sqrt{n})$ .

**Theorem 5.3** *For any constant  $c < 1$  there is a monotone Boolean function  $f$  which satisfies*

$$\sum_{|a| \geq cn} \hat{f}^2(a) \geq \Omega\left(\frac{1}{\sqrt{n} \log n}\right).$$

*Proof* Assume for contradiction that there is some constant  $c < 1$  such that for any monotone function  $f$

$$\sum_{|a| \geq cn} \hat{f}^2(a) \leq O\left(\frac{1}{\sqrt{n} \log n}\right).$$

But this implies that the low-degree algorithm which searches all coefficients of degree at most  $cn$  will approximate  $f$  within an error of  $O(1/(\sqrt{n} \log n))$ . This contradicts Theorem 5.2 modulo constant factors.  $\square$

**Theorem 5.4** *The majority function  $f$  satisfies  $\sum_{|a| \geq \Omega(\sqrt{n})} \hat{f}^2(a) \geq \Omega(1/\sqrt{n})$ .*

*Proof* Let  $f(x) = \text{MAJ}(x) = I[\sum_{i=1}^n x_i \geq n/2]$ . Since  $f$  is a symmetric function, the influence of all variables are equal. From the first equality in the proof of Theorem 4.1 we have

$$\sum_a |a| \hat{f}^2(a) = \sum_{i=1}^n I_i(f) = nI_1(f).$$

To get a bound on  $I_1(MAJ)$ , note that  $I_1(MAJ) \geq 2^{-n} \binom{n}{n/2} \geq \frac{c}{\sqrt{n}}$ , for some constant  $c$ . Thus

$$\begin{aligned} c\sqrt{n} &\leq \sum_a |a| \hat{f}(a)^2 = \sum_{|a| \geq \frac{c}{2}\sqrt{n}} |a| \hat{f}(a)^2 + \sum_{|a| < \frac{c}{2}\sqrt{n}} |a| \hat{f}(a)^2 \\ &\leq n \sum_{|a| \geq \frac{c}{2}\sqrt{n}} \hat{f}(a)^2 + \frac{c}{2}\sqrt{n}. \end{aligned}$$

This implies that

$$\sum_{|a| \geq \frac{c}{2}\sqrt{n}} \hat{f}(a)^2 \geq \frac{c}{2\sqrt{n}}.$$

□

We now use the Vapnik-Chernovenkis dimension to find lower bounds on the sample size. Kearns [K] had also observed that the VC dimension can be used to prove negative learning results for monotone Boolean functions. Recall that if  $C$  is a class of Boolean functions then  $C$  shatters  $A \subseteq \{0, 1\}^n$  if for every Boolean function  $g : A \rightarrow \{0, 1\}$  there exists a Boolean function  $f \in C$  such that  $f|_A = g$ . The *Vapnik-Chernovenkis* dimension of  $C$ ,  $VCdim(C)$ , is the cardinality of the largest subset  $A$  which is shattered by  $C$ . Ehrenfeucht *et al.* [EHKV88] proved a sample complexity lower bound of

$$\Omega\left(\frac{1}{\epsilon} \ln \frac{1}{\delta} + \frac{1}{\epsilon} VCdim(C)\right)$$

for PAC learning any class  $C$  with error  $\epsilon$  and confidence  $\delta$ . It is easy to see that the VC-dimension of monotone functions is at least  $\binom{n}{n/2} \sim 2^n/\sqrt{n}$ . Hence we get the following easy corollary.

**Corollary 2** *Any PAC learning algorithm for all monotone Boolean functions under an arbitrary distribution with error  $\epsilon$  and confidence  $\delta$  (for sufficiently small  $\epsilon$  and  $\delta$ ) requires at least  $\Omega\left(\frac{2^n}{\epsilon\sqrt{n}} + \frac{1}{\epsilon} \ln \frac{1}{\delta}\right)$  examples.*

## 6 Circuit Approximations of Monotone Boolean Functions

We study the circuit complexity of approximating monotone Boolean functions using Boolean circuits (non-monotone). We prove that any monotone Boolean function can be approximated by a non-monotone Boolean circuit of subexponential size and sublinear depth. This result is a consequence of Theorem 4.1.

**Theorem 6.1** *For any monotone Boolean function  $f$  on  $n$  variables and for any constant  $\epsilon > 0$ , there is a Boolean circuit of size*

$$2^{O(\frac{1}{\epsilon}\sqrt{n} \log n)}$$

*and depth*

$$O\left(\frac{1}{\epsilon}\sqrt{n} \log n\right)$$

*which approximates  $f$  to within  $\epsilon$  under the uniform distribution.*

*Proof* Note that the low-degree algorithm outputs a hypothesis

$$h(x) = \sum_{|a| \leq O(\frac{1}{\epsilon}\sqrt{n})} c_a \chi_a(x),$$

where  $c_a \sim E[f\chi_a]$ . Note that each  $c_a$  can be at most  $2^{O(\frac{1}{\epsilon}\sqrt{n}\log n)}$  bits. So essentially we need to add  $m$  number each being  $m$  bits, where  $m(n) = 2^{O(\frac{1}{\epsilon}\sqrt{n}\log n)}$ . This problem is known to be in  $NC^1$  (Boolean functions computable by a bounded fan-in, logarithmic depth, and polynomial size Boolean circuit).  $\square$

In the next section we will investigate the circuit approximation of monotone graph properties. For this we will need the following theorem.

**Theorem 6.2** *For any product distribution  $D$  with  $\mu_i = \mu$ , for all  $i$ , any monotone Boolean function  $f$  on  $n$  variables and any constant  $\epsilon > 0$ , there is a Boolean circuit of size*

$$2^{O(\frac{2}{\epsilon}\sqrt{n}\log(\frac{\epsilon\sqrt{n}}{\sigma}))}$$

and depth

$$O\left(\frac{1}{\epsilon}\sigma\sqrt{n}\log(\epsilon\sqrt{n}/\sigma)\right)$$

which approximates  $f$  to within  $\epsilon$ .

*Proof* (Sketch) We have  $\|a\| = |a|\log\frac{1}{\sigma}$ , where  $\sigma = \sqrt{\mu(1-\mu)}$ . By Theorem 4.1,

$$\sum_{|a|\geq k} \tilde{f}(a)^2 \leq \frac{2\sigma}{k} I_D(f)\sqrt{n}.$$

To get an error  $\epsilon$  we need  $k \geq \frac{2\sigma}{\epsilon} I_D(f)\sqrt{n}$ . Recall that for monotone Boolean functions,  $I_D(f) \leq 1$ . Notice that we can truncate the  $\mu_i$ 's and  $\sigma_i$ 's while only incurring a polynomial blow-up in error. When these are truncated,  $\phi_a(x)$  can be computed in polynomial time and therefore there is a polynomial size circuit that computes them. Now we proceed as in Theorem 6.1.  $\square$

## 6.1 Circuit Approximation of Monotone Graph Properties

We consider some monotone graph properties on the random graph  $G(n, p)$ , where  $n$  is the number of vertices of  $G$  and  $p$  is the edge existence probability. Some well-known graph properties are monotone: the clique function  $CLIQUE_n^k$  which is one if and only if the graph has a clique of size at least  $k$ , the hamiltonicity function  $HAM_n$  which is one if and only if the graph has a Hamiltonian cycle, the planarity function  $PLANAR_n$  which is one if and only if the graph is planar etc. We investigate the problem of approximating these monotone graph-theoretic functions.

We adopt the *probabilistic* model of the random graph on  $n$  vertices (see [S94] for other models). The random graph  $G = G(n, p)$  is a probability distribution on the edges of the complete graph  $K_n$  on  $n$  vertices, where each edge exists independently with probability  $p \in [0, 1]$ . A Boolean function  $f$  on the edge set  $E(G)$  is called a *monotone graph property* if  $f$  is a monotone (or antimonotone) Boolean function over  $E(G)$ . Any monotone graph property exhibits a *threshold phenomena* (see [B85, S94]). Let  $f$  be a monotone graph property on  $G(n, p)$ . A function  $r(n)$  is called a *threshold function* for  $f$  if it satisfies

1. if  $\lim_{n \rightarrow \infty} \frac{p(n)}{r(n)} = 0$  then  $\lim_{n \rightarrow \infty} \Pr[f(G(n, p)) = 1] = 0$ .
2. if  $\lim_{n \rightarrow \infty} \frac{r(n)}{p(n)} = 0$  then  $\lim_{n \rightarrow \infty} \Pr[f(G(n, p)) = 1] = 1$ .



**Definition 2** ( $\epsilon$ - $G(n, p)$  circuit)

For a graph property  $A$ , for fixed  $\epsilon > 0$  and  $p \in [0, 1]$ , and random inputs drawn from  $G(n, p)$ , we call a Boolean circuit an  $\epsilon$ - $G(n, p)$  circuit if it outputs a correct answer to the property  $A$  with probability at least  $1 - \epsilon$ . Here the probability is with respect to the distribution  $G(n, p)$  on the input graphs.

**Theorem 6.3** Let  $f$  be a monotone graph property with a threshold function  $r(n)$  (the number of inputs to  $f$  is  $m = \binom{n}{2}$ ). Then for any fixed  $\epsilon > 0$ ,  $p \in [0, 1]$ , there is an  $\epsilon$ - $G(n, p)$  circuit which approximates  $f$  with respect to the distribution  $G(n, p)$  of size

$$2^{O\left(\frac{\sqrt{nr(n)}}{\epsilon} \log \frac{\epsilon\sqrt{n}}{\sqrt{r(n)}}\right)}.$$

*Proof* Fix  $\epsilon, p > 0$ . By (1),(2) there exists a constant  $c > 1$  such that for  $p \geq cr(n)$  we have  $\Pr[f(G(n, p)) = 1] \geq 1 - \epsilon$  and for  $p \leq r(n)/c$  we have  $\Pr[f(G(n, p)) = 0] \geq 1 - \epsilon$ . Therefore for  $p \geq cr(n)$  we build the constant circuit 1 and for  $p \leq r(n)/c$  we build the constant circuit 0. For  $p$  in the range  $(r(n)/c, cr(n))$  we build the circuit in Theorem 6.2 with  $\mu = p$  and error  $\epsilon/c$ . This circuit has the required size and depth.  $\square$

Next we consider two monotone graph properties that are NP-complete:  $CLIQUE_k^n$  and  $HAM_n$ . The clique function has a threshold of  $r(n) = n^{-2/(k-1)}$  and the Hamiltonicity function has a threshold of  $r(n) = \ln n/n$ . Our approach failed to give an interesting bound for the clique function. For the Hamiltonicity function we have the following.

**Corollary 3** For any  $\epsilon$  and a fixed  $p$ , there is a Boolean circuit that approximates  $HAM_n$  to within error  $\epsilon$  and has size

$$2^{O(\frac{1}{\epsilon}\sqrt{n} \log^{1.5} n)}$$

and depth

$$O\left(\frac{1}{\epsilon}\sqrt{n} \log^{1.5} n\right).$$

*Proof* Apply Theorem 6.3 with  $r(n) = \ln n/n$ .  $\square$

*Remark.* The above results can be contrasted with several works on solving the search versions of NP-hard problems on random graphs. For example, Bollobás, Fenner, and Frieze [BFF87] described a expected polynomial time algorithm for finding Hamiltonian cycles in the random graph  $G(n, m)$  with  $n$  vertices and

$$m = n \log n/2 + n \log \log n/2 + c_n n$$

edges, where  $c_n$  is some sequence of integers. Their model of the random graph  $G(n, m)$  is a uniform distribution on undirected graphs with  $n$  vertices and  $m$  edges. They proved that the success probability of their algorithm equals to that of the existence of such Hamiltonian cycles. Their result differs from ours in several ways. First, our result only concerns the approximation of the decision problem, and second, we don't need to impose any restriction on the number of edges in our random graph (also our random graph model is slightly different than theirs, i.e.,  $G(n, p)$  vs.  $G(n, m)$ ), but third, their algorithm runs in expected polynomial time whereas our algorithm runs in subexponential time.

## 7 Approximating over a Convex Mixture of Distributions

The result in this section states that approximability over a collection of distributions  $\{D_i\}$  implies the approximability over any distribution in the convex space of  $\{D_i\}$ . First, we need to introduce a notion of *convex dimension* of a probability distribution.

**Definition 3** (Convex Dimension of a Distribution)

Let  $I$  be some index set. Let  $\mathcal{D} = \{D_i\}_{i \in I}$  be the set of all product distributions over  $\{0, 1\}^n$ . Consider a distribution  $D$  of the form

$$D = \sum_{i \in I} \lambda_i D_i,$$

where  $\lambda_i \in [0, 1]$ , for each  $i \in I$ , and  $\sum_{i \in I} \lambda_i = 1$ . We call  $D$  the convex linear combination of distributions  $\{D_i : i \in I\}$ . The convex dimension  $\text{cdim}(D)$  of  $D$  is the least  $m$  such that  $D$  can be represented as a convex linear combination of  $m$  product distributions.

**Lemma 3** Let  $\tau_1, \dots, \tau_m \in \{-1, 1\}$  and let  $d_1, \dots, d_m$  be positive real numbers. If  $\sum_{i=1}^m d_i \tau_i < 0$  then

$$\sum_{\tau_i = -1} d_i \geq \frac{1}{2} \sum_{i=1}^m d_i.$$

*Proof* If  $\sum_{i=1}^m d_i \tau_i < 0$  then  $\sum_{i: \tau_i = 1} d_i < \sum_{i: \tau_i = -1} d_i$ . Thus  $\sum_{i=1}^m d_i < 2 \sum_{i: \tau_i = -1} d_i$ .  $\square$

**Theorem 7.1** Let  $f$  be a Boolean function that can be approximated over each  $D_i$  from the set  $\{D_i\}_{i=1}^m$  of distributions. Then  $f$  can be approximated over any distribution  $D$  that is a convex linear combination of the  $D_i$ 's, say  $D = \sum_{i=1}^m \lambda_i D_i$ , provided that the  $\lambda_i$ 's are known and each distribution  $D_i$  is known and is polynomial time computable.

*Proof* Suppose  $\sum_{i=1}^m \lambda_i = 1$  and  $D = \sum_{i=1}^m \lambda_i D_i$ . Let  $h_i$  be a hypothesis over distribution  $D_i$  satisfying  $E_{D_i}[h_i \neq f] \leq \epsilon/2$ . Define the hypothesis  $H$  over  $D$  to be

$$H(x) = \text{sgn} \left( \sum_{i=1}^m \lambda_i D_i(x) h_i(x) \right).$$

Then

$$E_D[H \neq f] = E_D[\lambda_1 D_1 h_1 f + \dots + \lambda_m D_m h_m f < 0].$$

Using Lemma 3, if  $\lambda_1 D_1 h_1 f + \dots + \lambda_m D_m h_m f < 0$  for some  $x$  then

$$\sum_{i: h_i(x) f(x) = -1} \lambda_i D_i \geq \frac{\lambda_1 D_1 + \dots + \lambda_m D_m}{2}.$$

Therefore we have the following.

$$\begin{aligned} E_D[H \neq f] &= \sum_{x_0: H(x_0) \neq f(x_0)} \sum_i \lambda_i D_i(x_0) \\ &\leq 2 \sum_{x_0: H(x_0) \neq f(x_0)} \sum_{i: h_i(x_0) \neq f(x_0)} \lambda_i D_i(x_0), \text{ Lemma 3} \\ &\leq 2 \sum_{x_0} \sum_{i: h_i(x_0) \neq f(x_0)} \lambda_i D_i(x_0) \\ &= 2 \sum_i \lambda_i E_{D_i}[h_i \neq f] \leq \epsilon. \end{aligned}$$

$\square$

## 8 Learning in Polynomial Time

In this section we describe some polynomial time learning results on some subclasses of monotone Boolean functions.

### 8.1 Weak Learning

Kearns and Valiant [KV89] proved that monotone Boolean functions are weakly learnable under the uniform distribution with error  $1/2 - 1/(2n)$ . In the following we improve their result slightly. In particular, there is a weak learner with error  $1/2 - \Omega(\log^2 n/n)$  under the uniform distribution and there are weak learners with error  $1/2 - c/n$ , for any constant  $c$ , under any product distribution.

For learning under the uniform distribution we will use the following result of Kahn, Kalai and Linial [KKL88] on the lower bound of the sum of the squares of the influences of variables.

**Lemma 4** [KKL88] *Let  $f \in \{0, 1\}$  be a Boolean function with  $p = \Pr[f(x) = 1] \leq 1/2$ . Then*

$$\sum_{i=1}^n I_i(f)^2 \geq \frac{p^2 (\log n)^2}{5n}.$$

**Theorem 8.1** *There is a polynomial time weak PAC learning algorithm with error  $\epsilon = \frac{1}{2} - \Omega\left(\frac{\log^2 n}{n}\right)$  for any monotone Boolean function under the uniform distribution.*

*Proof* We will assume that without loss of generality that  $p = \Pr[f(x) = 1] \leq 1/2$ , since we can take  $\neg f(\neg x_1, \dots, \neg x_n)$ . This transformation does not effect the influences.

If  $p < 1/4$  we already have a weak learning using the all-zero hypothesis, otherwise since  $I_i(f)^2 = \hat{f}^2(e_i)$  and using Lemma 4,  $\sum_{i=1}^n \hat{f}^2(e_i) \geq p^2 \log^2 n / (5n) \geq \log^2 n / 80n$ . Combining this with Fact 1 with  $A = \{e_i : i \in [n]\}$ , we get a weak learner with the claimed accuracy.  $\square$

**Theorem 8.2** *For any constant  $k$  there is a polynomial time weak PAC learning algorithm with error  $\epsilon = \frac{1}{2} - \frac{k}{n}$  for monotone Boolean functions under any product distribution.*

*Proof* By Fact 1, if  $\sum_{||a|| \geq k} \hat{f}^2(a) \leq 1/2$  then we get a  $1/4$ -approximator by the standard low-degree algorithm (using again the fact that  $||a|| \geq |a|$ ). Otherwise by Theorem 4.1, we have  $1.062 \frac{\sqrt{n}}{k} I_D(f) \geq \frac{1}{2}$  and hence  $\sum_i \hat{f}^2(e_i) \geq \frac{k^2}{4(1.062)^{2n}}$ .  $\square$

### 8.2 Strong Learning

Sakai and Maruoka [SM94] proved that monotone  $O(\log n)$ -term DNF is PAC learnable under the uniform distribution. We improve their result in two ways. First we extend the class to a larger subclass of the monotone Boolean functions and second we extend the distribution to constant-bounded product distributions.

A variable  $x_i$  is *relevant* for  $f$  if there are  $a, b \in \{0, 1\}^n$  with  $a = b \oplus e_i$  and  $f(a) \neq f(b)$ . Note that  $x_i$  is relevant if and only if  $I_i(f) > 0$ .

**Definition 4** (The concept class  $A(k)$ )

*Let  $A(k)$  be the class of Boolean functions of the form  $f(T_1, \dots, T_k)$ , where  $f$  is an arbitrary monotone Boolean function on  $k$  inputs and each  $T_i$  is a monotone conjunction or a disjunction over  $n$  variables.*

**Theorem 8.3** *The class  $A(\log n)$  is PAC learnable under constant bounded product distributions.*

*Proof* From Lemma 1, the influence of a variable is

$$I_{D,i}(f) = \frac{1}{4\sigma_i^2} \sum_{a:a_i=1} \tilde{f}(a)^2.$$

Since the product distribution is constant-bounded,  $\sigma_i$  is a constant and therefore if the influence is small then we may assume that  $x_i$  is not relevant since this incur only a small additional error to the hypothesis. Now the divide-conquer learning algorithm presented in [B95] can be used in the same way.  $\square$

**Theorem 8.4** *For any constant  $\epsilon$ , the class of monotone functions which depend on  $O\left(\frac{\log^2 n}{\log^2 \log n}\right)$  variables is PAC learnable with error  $\epsilon$  under constant bounded product distributions.*

*Proof* By the assumption, there are at most  $m(n) = \frac{\log^2 n}{\log^2 \log n}$  variables having nonzero influence. Note that we may ignore all the variables with very small influence since we are dealing with constant bounded product distributions. So we can apply the low-degree algorithm which will run in time  $2\sqrt{m(n)\log m(n)} = n^{O(1)}$ .  $\square$

**Theorem 8.5** *For any constant  $\epsilon$ , the class  $A\left(\frac{\log^2 n}{\log^3 \log n}\right)$  is PAC learnable with error  $\epsilon$  under constant bounded product distributions.*

*Proof* First, we claim that any term with size  $\Omega(\log \log n)$  may be ignored without incurring an error of more than  $O(1)$ . Now observe that with this simplification, there are at most  $\log^2 n / (\log \log n)^2$  variables. This problem reduces to Theorem 8.4.  $\square$

*Remark.* The learning algorithms discussed so far fit into the *statistical query* learning model introduced by Kearns [K93]. Hence by Kearns' results, these algorithms are robust against *classification noise* in the example oracle.

## 9 Acknowledgments

This paper would not have existed without Jeff Jackson telling us about all the *neat* things related to Fourier transform. His enthusiastic lectures during his visit at the Department of Computer Science, University of Calgary has more impact on this work than anything else. We wish to thank him for his generosity and helpful contributions during our work in this paper. We also thank Dan Boneh, for suggesting the notion of influence norm, to Yishay Mansour, for suggesting the average sensitivity viewpoint, and to Uriel Feige, for pointing out to us some references on Hamiltonian cycles in random graphs.

## References

- [AM91] William Aiello and Milena Mihail. Learning the Fourier Spectrum of Probabilistic Lists and Trees. In *Proceedings of the Second Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 291–299, 1991.
- [AS92] Noga Alon and Joel Spencer. *The Probabilistic Method*. Wiley-Interscience, 1992.

- [B85] Béla Bollobás. *Random Graphs*. Academic Press Inc., 1985.
- [B92] Mihir Bellare. A Technique for Upper Bounding the Spectral Norm with Applications to Learning. In *Proceedings of the Fifth Annual ACM Workshop on Computational Learning Theory*, pages 62-70, 1992.
- [B95] Nader H. Bshouty. Simple Learning Algorithms Using Divide and Conquer. In *Proceedings of the Eighth Annual ACM Conference on Computational Learning Theory*, pages 62-70, 1995.
- [BFF87] B. Bollobás, T. Fenner, and A. Frieze. An Algorithm for Finding Hamiltonian Paths and Cycles in Random Graphs. In *Combinatorica*:7, 1987.
- [BFJ+] Avrim Blum, Merrick Furst, Jeffrey Jackson, Michael Kearns, Yishay Mansour, and Steven Rudich. Weakly Learning DNF and Characterizing Statistical Query using Fourier Analysis. In *Proceedings of the Twenty Sixth Annual ACM Symposium on Theory of Computing*, pages 253–262, 1994.
- [BS92] Jehoshua Bruck and Roman Smolensky. Polynomial Threshold Functions,  $AC^0$  Functions, and Spectral Norms. In *SIAM Journal on Computing*, **21**:1, pages33–42, 1992.
- [CKM93] Don Coppersmith, Hugo Krawczyk, and Yishay Mansour. The Shrinking Generator. In *Advances in Cryptology CRYPTO*, 1993.
- [EHKV88] Andrzej Ehrenfeucht, David Haussler, Michael Kearns, and Leslie Valiant. A General Lower Bound on the Number of Examples Needed for Learning. In *Proceedings of the 1988 Workshop on Computational Learning Theory*, pages 139–154, 1988.
- [FJS91] Merrick Furst, Jeffrey Jackson, and Sean Smith. Improved Learning of  $AC^0$  Functions. In *Proceedings of the Fourth Annual Workshop on Computational Learning Theory*, pages 317–325, 1991.
- [HM91] Thomas Hancock and Yishay Mansour. Learning Monotone  $k$ - $\mu$  DNF Formulas on Product Distributions. In *Proceedings of the Fourth Annual Workshop on Computational Learning Theory*, pages 179–183, 1991.
- [HR89] Torben Hagerup and Christine Rub. A Guided Tour to Chernoff Bounds. In *Information Processing Letters*, **33**, pages 305–308, 1989.
- [J94] Jeffrey Jackson. An Efficient Membership-Query Algorithm for Learning DNF with Respect to the Uniform Distribution. In *Proceedings of the 35th Annual Symposium on Foundations of Computer Science*, pages 42–53, 1994.
- [K] Michael Kearns. *The Computational Complexity of Machine Learning*. MIT Press, 1990.
- [K93] Michael Kearns. Efficient Noise Tolerant Learning from Statistical Queries. In *Proceedings of the Twenty Fifth Annual ACM Symposium on the Theory of Computing*, pages 392–401, 1993.
- [KKL88] Jeff Kahn, Gil Kalai, and Nathan Linial. The Influence of Variables on Boolean Functions. In *Proceedings of the 29th Annual Symposium on Foundations of Computer Science*, pages 68–80, 1988.

- [KM91] Eyal Kushilevitz and Yishay Mansour. Learning Decision Trees using the Fourier Spectrum. In *Proceedings of the 23rd Annual ACM Symposium on Theory of Computing*, pages 455–464, 1991. Also in *SIAM Journal on Computing*, **22**:6, pages 1331–1348, 1993.
- [KV89] Michael Kearns and Leslie Valiant. Cryptographic Limitations on Learning Boolean Formulae and Finite Automata. In *Proceedings of the 21st Annual ACM Symposium on Theory of Computing*, pages 433–444, 1989.
- [LMN93] Nathan Linial, Yishay Mansour, and Noam Nisan. Constant Depth Circuits, Fourier Transform and Learnability. *Journal of ACM*, **40**:3, pages 607–620, 1993. An earlier version appeared in *Proceedings of the 30th Annual IEEE Symposium on the Foundations of Computer Science*, 1989.
- [M92] Yishay Mansour. An  $O(n^{\log \log n})$  Learning Algorithm for DNF. In *Proceedings of the Fifth Annual ACM Workshop on Computational Learning Theory*, pages 53–61, 1992.
- [M94] Yishay Mansour. Learning Boolean Functions via the Fourier Transform. Tutorial Notes for the *Workshop on Computational Learning Theory*, 1994.
- [M95] Yishay Mansour. Randomized Interpolation and Approximation of Sparse Polynomials. In *SIAM Journal on Computing*, **24**:2, pages 357–368, 1995.
- [S94] Joel Spencer. *Ten Lectures on the Probabilistic Method*. *SIAM CBMS-NSF Regional Conference Series in Applied Mathematics*, second edition, 1994.
- [SM94] Yoshifumi Sakai and Akira Maruoka. Learning  $O(\log n)$ -term Monotone DNF. In *Proceedings of the Seventh Annual ACM Conference on Computational Learning Theory*, pages 165–172, 1994.
- [V84] Leslie Valiant. A Theory of the Learnable. *Communications of the ACM*, 27(11):1134–1142, 1984.
- [W87] Ingo Wegener. *The Complexity of Boolean Functions*. *Wiley-Teubner*, 1987.