# On Learning Embedded Midbit Functions

Rocco A. Servedio[1]

Division of Engineering and Applied Sciences
Harvard University
Cambridge, MA 02138 USA
rocco@deas.harvard.edu
http://www.cs.harvard.edu/ rocco

**Abstract.** A midbit function on $\ell$ binary inputs $x_1, \ldots, x_\ell$ outputs the middle bit in the binary representation of $x_1 + \cdots + x_\ell$. We consider the problem of PAC learning *embedded* midbit functions, where the set $S \subset \{x_1, \ldots, x_n\}$ of relevant variables on which the midbit depends is unknown to the learner.

To motivate this problem, we first show that a polynomial time learning algorithm for the class of embedded midbit functions would immediately yield a fairly efficient (quasipolynomial time) PAC learning algorithm for the entire complexity class *ACC*. We then give two different subexponential learning algorithms, each of which learns embedded midbit functions under any probability distribution in $2^{\sqrt{n}\log n}$ time. Finally, we give a polynomial time algorithm for learning embedded midbit functions under the uniform distribution.

## 1 Introduction

A central goal of computational learning theory is to understand the computational complexity of learning various classes of Boolean functions. While much research has been devoted to learning syntactic classes such as decision trees, DNF formulas, and constant depth circuits, researchers have also considered various "semantically defined" classes as well. A natural and important class of this sort is the class of *embedded symmetric functions* which was studied by Blum *et al.* [5]. (Recall that a Boolean function is symmetric if its value depends only on the number of input bits which are set to 1.) An embedded symmetric function is a Boolean function which depends only on some subset of its input variables and is a symmetric function on this subset, i.e., it is a symmetric function whose domain is "embedded" in a larger domain containing irrelevant variables.

In this paper we give a detailed PAC learning analysis of an interesting and natural family of embedded symmetric functions, namely the *embedded midbit functions*. An embedded midbit function is defined by a subset $i_1, \ldots, i_s$ of variables from $\{1, \ldots, n\}$. The value of this embedded midbit function on an input $x \in \{0,1\}^n$ is the value of the middle bit in the binary representation of $x_{i_1} + x_{i_2} + \cdots + x_{i_s}$. As described below, we show that the class of embedded midbit functions has many interesting properties from a PAC learning perspective.

## 1.1   Our Results

We first give a hardness result (Theorem 2) for learning embedded midbit functions in the standard PAC model of learning from random examples drawn from an arbitrary probability distribution. Using Green *et al.*'s characterization of the complexity class $ACC$ [9], we show that if there is a PAC learning algorithm for the class of embedded midbit functions which runs in polynomial time (or even quasipolynomial time), then the class $ACC$ of constant-depth, polynomial-size circuits of unbounded fanin AND/OR/MOD$_m$ gates can also be PAC learned in quasipolynomial time. This would be a major breakthrough since, as described in Section 3, the fastest PAC learning algorithms to date for even very restricted subclasses of $ACC$ require much more than quasipolynomial time. Our hardness result strengthens an earlier hardness result of Blum *et al.* for embedded symmetric functions, and establishes an interesting connection between learning the "semantic" class of embedded midbit functions and learning rich syntactic classes.

While Theorem 2 implies that it may be difficult to learn embedded midbit functions efficiently under an arbitrary distribution, this does not mean that PAC learning algorithms for embedded midbit functions must require exponential time. In Section 4 we give two different subexponential time PAC learning algorithms, each of which can learn embedded midbit functions over $n$ variables in time $n^{O(\sqrt{n})}$.

Finally, by means of a careful analysis of the correlation of single variables and pairs of variables with embedded midbit functions, we show in Section 5 that embedded midbit functions can be learned in polynomial time under the uniform distribution. Embedded midbit functions thus give a simple and natural concept class which seems to exhibit a large gap between the complexity of learning in the uniform distribution PAC model and the general (arbitrary distribution) PAC model.

## 2   Preliminaries

Throughout this paper $S$ denotes a subset of the variables $\{x_1, \ldots, x_n\}$ and $s$ denotes $|S|$. All logarithms are base 2.

**Definition 1.** *For $S \neq \emptyset$ the embedded midbit function $M_S : \{0,1\}^n \to \{0,1\}$ is defined as $M_S(x) =$ the value of the $\lfloor \log(s)/2 \rfloor$-th bit in the binary representation of $\sum_S x_i$, where we consider the least significant bit to be the 0-th bit. (We take $M_\emptyset(x)$ to be identically 0.) The class $C_{mid}$ of embedded midbit functions is $C_{mid} = \{M_S\}_{S \subseteq \{x_1, \ldots, x_n\}}$.*

We write $C_{sym}$ to denote the class of all embedded symmetric functions on $\{0,1\}^n$ as described in Section 1; note that $C_{mid} \subset C_{sym}$.

**Definition 2.** *Given an embedded midbit function $M_S(x)$, let $f_s : \{0, 1, \ldots, s\} \to \{0,1\}$ be the unique function such that $M_S(x) = f_s(\sum_S x_i)$ for all $x \in \{0,1\}^n$. We say that $f_s$ is the* basis function *of $M_S(x)$ and we refer to the $(s+1)$-bit string $f_s(0)f_s(1) \ldots f_s(s)$ as the* pattern *of $f_s$.*

If $f_s$ is the basis function for $M_S$ then the pattern for $f_s$ is a concatenation of strings of the form $0^{k(s)} 1^{k(s)}$, where $k(s) = 2^{\lfloor \log(s)/2 \rfloor}$ and the concatenation is truncated to be of length precisely $s+1$. It is easy to see that $\sqrt{s}/2 < k(s) \leq \sqrt{s}$.

A function $f$ is *quasipolynomial* if $f(n) = 2^{(\log n)^{O(1)}}$. We write $[a \bmod b]$ to denote the unique real number $r \in [0, b)$ such that $a = kb + r$ for some integer $k$.

## 2.1 The learning model

We work in the standard Probably Approximately Correct (PAC) learning model [17] and the uniform distribution variant of the PAC model. Let $C$ be a class of Boolean functions over $\{0, 1\}^n$. In the PAC model, a learning algorithm has access to a random example oracle $EX(c, \mathcal{D})$ which when invoked in one time step provides a labeled example $\langle x, c(x) \rangle \in \{0, 1\}^n \times \{0, 1\}$ where $x$ is drawn from the distribution $\mathcal{D}$ over $\{0, 1\}^n$. An algorithm $A$ is a PAC learning algorithm for class $C$ if the following holds: for all $c \in C$ and all distributions $\mathcal{D}$ over $\{0, 1\}^n$, if $A$ is given as input $\epsilon, \delta > 0$ and $A$ is given access to $EX(c, \mathcal{D})$, then with probability at least $1 - \delta$ the output of $A$ is a hypothesis $h : \{0, 1\}^n \to \{0, 1\}$ such that $\Pr_{x \in \mathcal{D}}[c(x) \neq h(x)] \leq \epsilon$. (Strictly speaking, the output of $A$ is some particular representation of $h$ such as a Boolean circuit.) Algorithm $A$ is said to run in time $t$ if (i) the worst case running time of $A$ (over all choices of $c \in C$ and all distributions $\mathcal{D}$) is at most $t$, and (ii) for every output $h$ of $A$ and all $x \in \{0, 1\}^n$, $h(x)$ can be evaluated in time $t$.

If $A$ satisfies the above definition only for some fixed distribution $\mathcal{D}$ (such as the uniform distribution on $\{0, 1\}^n$), then we say that $A$ is a PAC learning algorithm for $C$ under distribution $\mathcal{D}$.

## 3 Hardness of Learning Embedded Midbit Functions

In this section we show that learning embedded midbit functions is almost as difficult as learning a rich syntactic class which contains decision trees, DNF formulas, and constant depth circuits.

## 3.1 Background: Hardness of Learning $C_{sym}$

We first describe a result of Blum *et al.* which gives some evidence that the broader class $C_{sym}$ of embedded symmetric functions may be hard to PAC learn in polynomial time. Let $C_{log}$ denote the class of Boolean functions on $n$ bits which have at most $\log n$ relevant variables. Note that like $C_{sym}$, the class $C_{log}$ has the property that learning is no more difficult than finding relevant variables – in either case, once the set of relevant variables has been identified, learning is simply a matter of observing and filling in at most $n$ "table entries" which define the function (these entries are the bits of the pattern for a function from $C_{sym}$, and are the values of the function on all $2^{\log n}$ inputs for a function from $C_{log}$).

Building on this intuition, Blum *et al.* gave a polynomial time prediction-preserving reduction from $C_{log}$ to $C_{sym}$, thus showing that if $C_{sym}$ can be PAC learned in polynomial time then $C_{log}$ can also be PAC learned in polynomial time. Since no polynomial time learning algorithm is yet known for $C_{log}$, this gives some evidence that $C_{sym}$ may not be learnable in polynomial time.

### 3.2   Hardness of Learning $C_{mid}$

The class $ACC$ was introduced by Barrington [2] and since been studied by many researchers, e.g. [1, 3, 4, 9, 12, 18, 19]. $ACC$ consists of languages recognized by a family of constant-depth polynomial-size circuits with NOT gates and unbounded fanin AND, OR and $MOD_m$ gates, where $m$ is fixed for each circuit family. In the context of learning theory $ACC$ is quite an expressive class, containing as it does polynomial size decision trees, polynomial size DNF formulas, and the well-studied class $AC^0$ of constant-depth polynomial-size AND/OR/NOT circuits.

Building on work of Beigel and Tarui [4], Green *et al.* [9] have given the following characterization of $ACC$ :

**Theorem 1.** *For each $L \in ACC$ there is a depth-2 circuit which recognizes $L \cap \{0, 1\}^n$ and has the following structure: the top-level gate computes a midbit function of its inputs, and the bottom level consists of $2^{(\log n)^{O(1)}}$ AND gates each of fanin $(\log n)^{O(1)}$.*

Using this characterization we obtain the following hardness result for learning $C_{mid}$ :

**Theorem 2.** *If $C_{mid}$ can be PAC learned in polynomial (or even quasipolynomial) time, then $ACC$ can be PAC learned in quasipolynomial time.*

*Proof.* Let $f : \{0, 1\}^n \to \{0, 1\}$ be the target ACC function. Let $q(n) = 2^{(\log n)^{O(1)}}$ be an upper bound on the number of AND gates on the bottom level of the Green *et al.* representation for $f$, and let $\ell(n) = (\log n)^{O(1)}$ be an upper bound on the fanin of each bottom level AND gate. Given an instance $x \in \{0, 1\}^n$ we generate a new instance $x' \in \{0, 1\}^m$ where $m = 2^{(\log n)^{O(1)}}$ by listing $q(n)$ copies of each $AND$ of at most $\ell(n)$ variables from $x_1, \ldots, x_n$. Theorem 1 implies that there is an embedded midbit function $f'$ on $m$ bits such that $f(x) = f'(x')$ for all $x \in \{0, 1\}^n$. By assumption we can PAC learn this function $f'$ in $2^{(\log m)^{O(1)}} = 2^{(\log n)^{O(1)}}$ time, so the theorem is proved.

We note that while our reduction only establishes quasipolynomial time learnability for $ACC$ from learnability of $C_{mid}$, whereas the Blum reduction would establish polynomial time learnability of $C_{log}$, the class $ACC$ is likely to be much harder to learn than $C_{log}$. While $C_{log}$ can be PAC learned in $n^{\log n}$ time by doing an exhaustive search for the set of $\log n$ relevant variables, no learning algorithm for $ACC$ is known which runs in subexponential time. In fact, no such algorithm

is known even for the subclass of polynomial-size, depth 3 AND/OR/NOT circuits; to date the most expressive subclass of $ACC$ which is known to be PAC learnable in subexponential time is the class of polynomial-size AND/OR/NOT circuits of depth 2, which has recently been shown by Klivans and Servedio [11] to be PAC learnable in time $2^{\tilde{O}(n^{1/3})}$.

# 4   Learning embedded midbit functions in $n^{O(\sqrt{n})}$ time

The results of Section 3 suggest that the class of embedded midbit functions may not be PAC learnable in quasipolynomial time. However, we will show that it is possible to learn this class substantially faster than a naive exponential time algorithm. In this section we describe two different algorithms each of which PAC learns $C_{mid}$ in time $n^{O(\sqrt{n})}$.

## 4.1   An algorithm based on learning linear threshold functions

Our first approach is a variant of an algorithm given by Blum et al. in section 5.2 of [5].

**Definition 3.** *Let $f : \{0,1\}^n \to \{0,1\}$ be a Boolean function and $p(x_1, \ldots, x_n)$ a real-valued polynomial. We say that $p(x)$ sign-represents $f(x)$ if for all $x \in \{0,1\}^n$, $p(x) \geq 0$ iff $f(x) = 1$.*

*Claim.* Let $M_S$ be an embedded midbit function. Then there is a polynomial $p_S(x_1, \ldots, x_n)$ of degree $O(\sqrt{n})$ which sign-represents $M_S(x)$.

*Proof.* Let $f_s$ be the basis function for $M_S$. Since $k(s) = \Omega(\sqrt{s})$, the number of "flip" positions in the pattern of $f_s$ where $f_s(i) \neq f_s(i+1)$ is $O(\sqrt{s})$. Since the pattern for $f_s$ has $O(\sqrt{s})$ flips, there is some polynomial $P(X)$ of degree $O(\sqrt{s})$ which is nonnegative on precisely those $i \in \{0, 1, \ldots, s\}$ which have $f_s(i) = 1$. This implies that $p_S(x_1, \ldots, x_n) = P(\sum_S x_i)$ sign-represents $M_S(x)$. Since the degree of $p_S$ is $O(\sqrt{s})$ and $s \leq n$ the claim is proved.

Consider the expanded feature space consisting of all monotone conjunctions of at most $O(\sqrt{n})$ variables. (Note that this feature space contains $\sum_{i=1}^{O(\sqrt{n})} n^i = n^{O(\sqrt{n})}$ features.) Claim 4.1 implies that $M_S(x)$ is equivalent to some linear threshold function over this space. Thus we can use known polynomial time PAC learning algorithms for linear threshold functions [6] over this expanded feature space to learn embedded midbit functions in $n^{O(\sqrt{n})}$ time.

We note that one can show that the sign-representing polynomial $p_S(x_1, \ldots, x_n)$ described in Claim 4.1 can be taken without loss of generality to have integer coefficients of total magnitude $n^{O(\sqrt{n})}$. This implies that simple algorithms such as Winnow or Perceptron can be used to learn in $n^{O(\sqrt{n})}$ time (instead of the more sophisticated algorithm of [6] which is based on polynomial time linear programming). We also note that in [13] Minsky and Papert used a symmetrization technique to give a lower bound on the degree of any polynomial which sign-represents the parity function. The same technique can be used to show that the $O(\sqrt{n})$ degree bound of Claim 4.1 is optimal for embedded midbit functions.

## 4.2   An algorithm based on learning parities

We have seen that any embedded midbit function is equivalent to some linear threshold function over the feature space of all $O(\sqrt{n})$-size monotone conjunctions. We now show that any embedded midbit function is equivalent to some parity over this feature space as well.

**Lemma 1.** *Let $r, \ell \geq 0$. Then $\binom{r}{2^\ell}$ is even if and only if*

$$[r \bmod 2^{\ell+1}] \in \{0, 1, \ldots, 2^\ell - 1\}.$$

*Proof.* By induction on $\ell$. The base case $\ell = 0$ is trivial; we suppose that the claim holds for $\ell = 0, \ldots, i-1$ for some $i \geq 1$. For the induction step we use the fact (Exercise 5.61 of [8]) that

$$\binom{r}{m} \equiv \binom{\lfloor r/p \rfloor}{\lfloor m/p \rfloor} \binom{[r \bmod p]}{[m \bmod p]} \pmod{p}$$

for all primes $p$ and all $r, m \geq 0$. Taking $p = 2$ and $m = 2^i$, since $i \geq 1$ we have

$$\binom{r}{2^i} \equiv \binom{\lfloor r/2 \rfloor}{2^{i-1}} \binom{[r \bmod 2]}{0} \equiv \binom{\lfloor r/2 \rfloor}{2^{i-1}} \pmod{2}$$

By the induction hypothesis this is 0 if and only if $[\lfloor r/2 \rfloor \bmod 2^i] \in \{0, 1, \ldots, 2^{i-1} - 1\}$, which holds if and only if $[r \bmod 2^{i+1}] \in \{0, 1, \ldots, 2^i - 1\}$.

*Claim.* Let $M_S$ be an embedded midbit function. Then $M_S(x)$ is equivalent to some parity of monotone conjunctions each of which contains at most $O(\sqrt{n})$ variables.

*Proof.* Let $\oplus$ denote the parity function. We have

$$M_S(x) = 0 \iff \lfloor \log(s)/2 \rfloor\text{-th bit of } \sum_S x_i \text{ is } 0$$

$$\iff \left[ \sum_S x_i \bmod 2^{\lfloor \log(s)/2 \rfloor + 1} \right] \in \{0, 1, \ldots, 2^{\lfloor \log(s)/2 \rfloor} - 1\}$$

$$\iff \binom{\sum_S x_i}{2^{\lfloor \log(s)/2 \rfloor}} = \binom{\sum_S x_i}{k(s)} \text{ is even}$$

$$\iff \bigoplus_{A \subseteq S, |A| = k(s)} \left( \bigwedge_{i \in A} x_i \right) = 0.$$

The third step is by Lemma 1 and the last step is because for any $x$ exactly $\binom{\sum_S x_i}{k(s)}$ of the conjunctions $\{\bigwedge_{i \in A} x_i\}_{A \subseteq S, |A| = k(s)}$ take value 1. Since $k(s) = O(\sqrt{n})$ the claim is proved.

As in the discussion following Claim 4.1, Claim 4.2 implies that we can use known PAC learning algorithms for parity [7, 10] over an expanded feature space to learn embedded midbit functions in $n^{O(\sqrt{n})}$ time.

## 5   A polynomial time algorithm for learning embedded midbits under uniform

In [5] Blum *et al.* posed as an open problem the question of whether embedded symmetric concepts can be learned under the uniform distribution in polynomial time. In this section, we show that embedded midbit functions can be PAC learned under the uniform distribution in polynomial time. This is in strong contrast to the results of Section 3 which indicate that embedded midbit functions probably cannot be PAC learned (in even quasipolynomial time) under arbitrary probability distributions.

Throughout this section we let $t(s)$ denote $\lfloor \frac{s}{k(s)} \rfloor$.

### 5.1   First approach: testing single variables

To learn $M_S$ it is sufficient to identify the set $S \subseteq \{x_1, \ldots, x_n\}$ of relevant variables. A natural first approach is to test the correlation of each individual variable with $M_S(x)$; clearly variables not in $S$ will have zero correlation, and one might hope that variables in $S$ will have nonzero correlation. However this hope is incorrect as shown by Lemma 3 below.

For $1 \leq i \leq n$ define $p_i = \Pr[M_S(x) = 1 | x_i = 1] - \Pr[M_S(x) = 1]$. The following fact is easily verified:

**Fact 3** *If $i \notin S$ then $p_i = 0$.*

**Lemma 2.** *If $i \in S$ then*

$$p_i = \frac{1}{2^s} \sum_{\ell=1}^{t(s)} (-1)^{\ell-1} \binom{s-1}{\ell k(s) - 1}. \tag{1}$$

*Proof.* Since the distribution on examples is uniform over $\{0,1\}^n$, the probability that exactly $\ell$ of the $s$ relevant variables are 1 is exactly $\binom{s}{\ell}/2^s$. Hence we have

$$p_i = \frac{1}{2^{s-1}} \sum_{\ell : f_s(\ell)=1} \binom{s-1}{\ell-1} - \frac{1}{2^s} \sum_{\ell : f_s(\ell)=1} \binom{s}{\ell}.$$

Using the identity $\binom{s}{\ell} = \binom{s-1}{\ell-1} + \binom{s-1}{\ell}$ we find that

$$p_i = \frac{1}{2^s} \sum_{f_s(\ell)=1} \left( \binom{s-1}{\ell-1} - \binom{s-1}{\ell} \right).$$

Cancelling terms where possible we obtain (1).

**Lemma 3.** *There are embedded midbit functions $M_S(x)$ with $S$ a proper subset of $\{x_1, \ldots, x_n\}$ such that $p_i = 0$ for all $1 \leq i \leq n$.*

*Proof.* By Fact 3 for $i \notin S$ we have $p_i = 0$. Suppose that $t(s)$ is even and $t(s)k(s) - 1 = s - 1 - (k(s) - 1)$. Then the expression for $p_i$ given in (1) is exactly 0 since the positive and negative binomial coefficients $\pm \binom{s-1}{\ell k(s)-1}$ and $\mp \binom{s-1}{(t(s)-\ell+1)k(s)-1}$ cancel each other out (e.g. take $s = 27, k(s) = 4, t(s) = 6$).

Thus the correlation of individual variables with $M_S(x)$ need not provide information about membership in $S$. However, we will show that by testing correlations of *pairs* of variables with $M_S(x)$ we can efficiently determine whether or not a given variable belongs to $S$.

## 5.2   Second approach: testing pairs of variables

For $1 \leq i, j \leq n, i \neq j$ let $p_{i,j} = \Pr[M_S(x) = 1 | x_i = x_j = 1] - \Pr[M_S(x) = 1 | x_j = 1]$. Similar to Fact 3 we have:

**Fact 4** *If $i \notin S$ then $p_{i,j} = 0$.*

**Lemma 4.** *If $i \in S$ and $j \in S$ then*

$$p_{i,j} = \frac{1}{2^{s-1}} \sum_{\ell=1}^{t(s)} (-1)^{\ell-1} \binom{s-2}{\ell k(s) - 2}. \tag{2}$$

*Proof.* We have

$$p_{i,j} = \frac{1}{2^{s-2}} \sum_{\ell : f_s(\ell)=1} \binom{s-2}{\ell - 2} - \frac{1}{2^{s-1}} \sum_{\ell : f_s(\ell)=1} \binom{s-1}{\ell - 1}.$$

Rearranging the sum as in Lemma 2 proves the lemma.

Our algorithm is based on the fact (Theorem 5 below) that quantities (1) and (2) cannot both be extremely close to 0.

**Theorem 5.** *Let $k$ be even and $\sqrt{s}/2 < k \leq \sqrt{s}$. Let*

$$A = \frac{1}{2^s} \sum_{\ell} (-1)^{\ell-1} \binom{s-1}{\ell k - 1} \quad and \quad B = \frac{1}{2^{s-1}} \sum_{\ell} (-1)^{\ell-1} \binom{s-2}{\ell k - 2}.$$

*Then* $\max\{|A|, |B|\} \geq \frac{1}{1000s}$.

The proof of Theorem 5 is somewhat involved and is deferred to Section 5.3.

With Theorem 5 in hand we can prove our main positive learning result for $C_{mid}$.

**Theorem 6.** *The class of embedded midbit functions is learnable under the uniform distribution in polynomial time.*

**Input:** variable $x_i \in \{x_1, \ldots, x_n\}$
**Output:** either "$x_i \in S$" or "$x_i \notin S$" correct with probability $1 - \frac{\delta}{n}$

1.  **let** $T$ be a sample of $m = \text{poly}(n, \log \frac{1}{\delta})$ labeled examples $\langle x, M_S(x) \rangle$
2.  **let** $\hat{p}_i$ be an empirical estimate of $p_i$ obtained from $T$
3.  **for all** $j \in \{1, \ldots, n\} - \{i\}$
4.      **let** $\hat{p}_{i,j}$ be an empirical estimate of $p_{i,j}$ obtained from $T$
5.  **if** $|\hat{p}_i| > \frac{1}{2000n}$ or $|\hat{p}_{i,j}| > \frac{1}{2000n}$ for some $j \in \{1, \ldots, n\} - \{i\}$
6.      **then output** "$i \in S$"
7.      **else output** "$i \notin S$"

**Fig. 1.** An algorithm to determine whether $x_i$ is relevant for $M_S(x)$.

*Proof.* Since there are fewer than $n^3$ midbit functions $M_S(x)$ which have $s \leq 3$ we can test each of these for consistency with a polynomial size random sample in polynomial time, and thus we can learn in polynomial time if $s \leq 3$. We henceforth assume that $s \geq 4$ and thus that $k(s) \geq 2$ is even.

We show that the algorithm in Figure 1 correctly determines whether or not $x_i \in S$ with probability $1 - \frac{\delta}{n}$. By running this algorithm $n$ times on variables $x_1, \ldots, x_n$ we can identify the set $S$ and thus learn $M_S$ correctly with probability $1 - \delta$.

**Case 1:** $x_i \notin S$. In this case by Facts 3 and 4 we have $p_i = p_{i,j} = 0$. Hence for a suitably chosen value of $m = \text{poly}(n, \log(\frac{1}{\delta}))$ each of the $n$ empirical estimates $\hat{p}_i, \hat{p}_{i,j}$ will satisfy $|\hat{p}_i| < \frac{1}{2000n}$ and $|\hat{p}_{i,j}| < \frac{1}{2000n}$ with probability $1 - \frac{\delta}{n^2}$. Thus in this case the algorithm outputs "$x_i \notin S$" with probability at least $1 - \frac{\delta}{n}$.

**Case 2:** $x_i \in S$. Since $s \geq 4$ there is some $x_j \neq x_i$ such that $x_j \in S$. Lemmas 2 and 4 and Theorem 5 imply that the true value of at least one of $|p_i|, |p_{i,j}|$ will be at least $\frac{1}{1000s} \geq \frac{1}{1000n}$. As before, for a suitably chosen value of $m = \text{poly}(n, \log(\frac{1}{\delta}))$, each of the $n$ empirical estimates $\hat{p}_i, \hat{p}_{i,j}$ will differ from its true value by less than $\frac{1}{2000n}$ with probability $1 - \frac{\delta}{n}$. Thus in this case the algorithm outputs "$x_i \in S$" with probability at least $1 - \frac{\delta}{n}$.

## 5.3   Proof of Theorem 5

The following lemma gives a useful expression for sums in the form of (1) and (2).

**Lemma 5.** *Let $r, j, k > 0$ with $k$ even. Then*

$$\sum_{\ell} (-1)^{\ell-1} \binom{r}{\ell k - j} =$$

$$\frac{-2}{k} \left( \sum_{\ell=1,3,5,\ldots,k-1} \left( 2 \cos \frac{\ell \pi}{2k} \right)^r \cos \left( \frac{(r+2j)\ell\pi}{2k} \right) \right). \tag{3}$$

*Proof.* We reexpress the left side as

$$\sum_\ell \binom{r}{\ell(2k)+(k-j)} - \sum_\ell \binom{r}{\ell(2k)-j}. \tag{4}$$

The following well known identity (see e.g. [15, 16]) is due to Ramus [14]:

$$\sum_\ell \binom{r}{\ell k - j} = \frac{1}{k} \sum_{\ell=1}^{k} \left(2\cos\frac{\ell\pi}{k}\right)^r \cos\left(\frac{(r+2j)\ell\pi}{k}\right).$$

Applying this identity to (4) we obtain

$$\frac{1}{2k}\left(\sum_{\ell=1}^{2k}\left(2\cos\frac{\ell\pi}{2k}\right)^r \cos\left(\frac{(r-2k+2j)\ell\pi}{2k}\right) - \sum_{\ell=1}^{2k}\left(2\cos\frac{\ell\pi}{2k}\right)^r \cos\left(\frac{(r+2j)\ell\pi}{2k}\right)\right)$$

Since even terms cancel out in the two sums above, we obtain

$$\frac{-1}{k}\left(\sum_{\ell=1,3,\dots,2k-1}\left(2\cos\frac{\ell\pi}{2k}\right)^r \cos\left(\frac{(r+2j)\ell\pi}{2k}\right)\right). \tag{5}$$

Consider the term of this sum obtained when $\ell = 2k - h$ for some odd value $h$:

$$\left(2\cos\frac{(2k-h)\pi}{2k}\right)^r \cos\left(\frac{(r+2j)(2k-h)\pi}{2k}\right)$$

$$= (-1)^{r+(r+2j)}\left(2\cos\frac{-h\pi}{2k}\right)^r \cos\left(\frac{(r+2j)(-h)\pi}{2k}\right)$$

$$= \left(2\cos\frac{h\pi}{2k}\right)^r \cos\left(\frac{(r+2j)h\pi}{2k}\right)$$

This equals the term obtained when $\ell = h$. Since $k = 2m$ is even we have that (5) equals the right side of (3).

The following two technical lemmas will help us analyze the right hand side of equation (3). No attempt has been made to optimize constants in the bounds.

**Lemma 6.** *Let $r, k$ be such that $k \geq 4$ is even and $k^2 - 2 \leq r < 4k^2 - 1$. Then*

*(i) for $\ell = 1, 3, \dots, k - 3$ we have $0 < \left(\cos\frac{(\ell+2)\pi}{2k}\right)^r < \left(\cos\frac{\ell\pi}{2k}\right)^r / 16$,*

*(ii) $\left(\cos\frac{\pi}{2k}\right)^r \geq \frac{1}{200}$.*

*Proof.* By considering the Taylor series of $\cos x$ one can show that $1 - \frac{x^2}{2} \leq \cos x \leq 1 - \frac{x^2}{3}$ for all $x \in [0, \frac{\pi}{2}]$.

Part (i): since $0 < \frac{\ell\pi}{2k} < \frac{(\ell+2)\pi}{2k} < \frac{\pi}{2}$, we have

$$\cos\frac{(\ell+2)\pi}{2k} = \cos\frac{\ell\pi}{2k}\cos\frac{\pi}{k} - \sin\frac{\ell\pi}{2k}\sin\frac{\pi}{k}$$

$$< \left(1 - \frac{\pi^2}{3k^2}\right)\cos\frac{\ell\pi}{2k}$$

and hence

$$
\begin{aligned}
\left(\cos \frac{(\ell+2)\pi}{2k}\right)^r &\leq \left(1 - \frac{\pi^2}{3k^2}\right)^r \left(\cos \frac{\ell\pi}{2k}\right)^r \\
&\leq \left(1 - \frac{\pi^2}{3k^2}\right)^{k^2-2} \left(\cos \frac{\ell\pi}{2k}\right)^r \\
&\leq \frac{e^{-\pi^2/3}}{(1 - \pi^2/(3k^2))^2} \cdot \left(\cos \frac{\ell\pi}{2k}\right)^r \\
&\leq \frac{1}{16} \cdot \left(\cos \frac{\ell\pi}{2k}\right)^r.
\end{aligned}
$$

Here the third inequality uses $\left(1 - \frac{1}{x}\right)^x \leq e^{-1}$ and the fourth inequality uses $k \geq 4$.

Part (ii): we have

$$
\begin{aligned}
\left(\cos \frac{\pi}{2k}\right)^r &> \left(\cos \frac{\pi}{2k}\right)^{4k^2} \\
&> \left(1 - \frac{\pi^2}{8k^2}\right)^{4k^2}.
\end{aligned}
$$

This is an increasing function of $k$ so for $k \geq 4$ the value is at least $\left(1 - \frac{\pi^2}{128}\right)^{64} \geq \frac{1}{200}$.

**Lemma 7.** *For all real $x$ and all odd $\ell \geq 3$, we have $|\cos(\ell x)| \leq \ell |\cos x|$.*

*Proof.* Fix $\ell \geq 3$. Let $y = \frac{\pi}{2} - x$ so $\ell|\cos x| = \ell|\sin y|$ and

$$
\begin{aligned}
|\cos(\ell x)| &= \left|\cos \frac{\ell\pi}{2} \cos(\ell y) - \sin \frac{\ell\pi}{2} \sin(\ell y)\right| \\
&= |\sin(\ell y)|
\end{aligned}
$$

(note that we have used the fact that $\ell$ is odd). Thus we must show that $|\sin(\ell y)| \leq \ell|\sin y|$. This is clearly true if $|\sin y| \geq \frac{1}{\ell}$; otherwise we may suppose that $0 \leq y < \sin^{-1}\frac{1}{\ell}$ (the other cases are entirely similar) so $0 \leq \ell y \leq \frac{\pi}{2}$. Now $\sin(\ell y) \leq \ell \sin y$ follows from the concavity of $\sin y$ on $[0, \frac{\pi}{2}]$ and the fact that the derivative of $\sin y$ is 1 at $y = 0$.

Using these tools we can now prove Theorem 5.

**Theorem 5** *Let $k$ be even and $\sqrt{s}/2 < k \leq \sqrt{s}$. Let*

$$
A = \frac{1}{2^s} \sum_\ell (-1)^{\ell-1} \binom{s-1}{\ell k - 1} \qquad and \qquad B = \frac{1}{2^{s-1}} \sum_\ell (-1)^{\ell-1} \binom{s-2}{\ell k - 2}.
$$

*Then $\max\{|A|, |B|\} \geq \frac{1}{1000s}$.*

*Proof.* By Lemma 5 we have

$$A = \frac{-1}{k}\left(\sum_{\ell=1,3,\ldots,k-1}\left(\cos\frac{\ell\pi}{2k}\right)^{s-1}\cos\left(\frac{(s+1)\ell\pi}{2k}\right)\right)\qquad(6)$$

and

$$B = \frac{-1}{k}\left(\sum_{\ell=1,3,\ldots,k-1}\left(\cos\frac{\ell\pi}{2k}\right)^{s-2}\cos\left(\frac{(s+2)\ell\pi}{2k}\right)\right).\qquad(7)$$

First the easy case: if $k=2$ then $4 \le s \le 15$ and $A = \frac{-1}{2}(\cos\frac{\pi}{4})^{s-1}\cos\left(\frac{(s+1)\pi}{4}\right)$, $B = \frac{-1}{2}(\cos\frac{\pi}{4})^{s-2}\cos\left(\frac{(s+2)\pi}{4}\right)$. Since either $\left|\cos\left(\frac{(s+1)\pi}{4}\right)\right|$ or $\left|\cos\left(\frac{(s+2)\pi}{4}\right)\right|$ must be $\frac{\sqrt{2}}{2}$ we have $\max\{|A|,|B|\} \ge \frac{1}{2^{\frac{s}{2}+1}}$ which is easily seen to be at least $\frac{1}{1000s}$ for $4 \le s \le 15$.

Now suppose $k \ge 4$. For $\ell = 3,\ldots,k-1$ we have

$$\left|\left(\cos\frac{\ell\pi}{2k}\right)^{s-1}\cos\left(\frac{(s+1)\ell\pi}{2k}\right)\right| \le \frac{\left(\cos\frac{\pi}{2k}\right)^{s-1}}{4^{\ell-1}}\cdot\left|\cos\left(\frac{(s+1)\ell\pi}{2k}\right)\right|$$

$$\le \left|\frac{\ell}{4^{\ell-1}}\cdot\left(\cos\frac{\pi}{2k}\right)^{s-1}\cos\left(\frac{(s+1)\pi}{2k}\right)\right|$$

where the first inequality is by repeated application of part (i) of Lemma 6 and the second is by Lemma 7. We thus have

$$\sum_{\ell=3,5,\ldots,k-1}\left|\left(\cos\frac{\ell\pi}{2k}\right)^{s-1}\cos\left(\frac{(s+1)\ell\pi}{2k}\right)\right|$$

$$\le \sum_{\ell=3,5,\ldots,k-1}\left|\frac{\ell}{4^{\ell-1}}\cdot\left(\cos\frac{\pi}{2k}\right)^{s-1}\cos\left(\frac{(s+1)\pi}{2k}\right)\right|$$

$$< \left|\left(\cos\frac{\pi}{2k}\right)^{s-1}\cos\left(\frac{(s+1)\pi}{2k}\right)\right|\cdot\sum_{\ell=3}^{\infty}\frac{\ell}{4^{\ell-1}}$$

$$= \frac{5}{18}\cdot\left|\left(\cos\frac{\pi}{2k}\right)^{s-1}\cos\left(\frac{(s+1)\pi}{2k}\right)\right|.$$

Thus the $\ell = 1$ term in the sum (6) dominates the sum and we have

$$|A| \ge \frac{13}{18}\cdot\frac{1}{k}\left|\left(\cos\frac{\pi}{2k}\right)^{s-1}\cos\left(\frac{(s+1)\pi}{2k}\right)\right|$$

$$\ge \frac{13}{3600k}\cdot\left|\cos\left(\frac{(s+1)\pi}{2k}\right)\right|$$

by part (ii) of Lemma 6. An identical analysis for $B$ shows that

$$|B| \geq \frac{13}{3600k} \cdot \left| \cos\left( \frac{(s+2)\pi}{2k} \right) \right|$$

as well.

We now observe that

$$\max\left\{ \left| \cos \frac{(s+1)\pi}{2k} \right|, \left| \cos \frac{(s+2)\pi}{2k} \right| \right\} \geq \cos\left( \frac{\pi}{2} - \frac{\pi}{4k} \right) = \sin \frac{\pi}{4k}.$$

Using Taylor series this is easily seen to be at least $\frac{\pi}{8k}$. Hence we have

$$\max\{|A|, |B|\} \geq \frac{13}{3600k} \cdot \frac{\pi}{8k} > \frac{1}{1000k^2} \geq \frac{1}{1000s}$$

and the theorem is proved.

## 6    Acknowledgement

## References

1. E. Allender and U. Hertrampf. Depth reduction for circuits of unbounded fan-in. *Information and Computation*, 112(2):217–238, 1994.
2. D. Barrington. Bounded-width polynomial-size branching programs recognize exactly those languages in $NC^1$. *Journal of Computer and System Sciences*, 38(1):150–164, 1989.
3. D. Barrington and D. Therien. Finite monoids and the fine structure of $NC^1$. *J. ACM*, 35(4):941–952, 1988.
4. R. Beigel and J. Tarui. On *ACC*. *Computational Complexity*, 4:350–366, 1994.
5. A. Blum, P. Chalasani, and J. Jackson. On learning embedded symmetric concepts. In *Proceedings of the Sixth Annual Conference on Computational Learning Theory*, pages 337–346, 1993.
6. A. Blumer, A. Ehrenfeucht, D. Haussler, and M. Warmuth. Learnability and the Vapnik-Chervonenkis dimension. *Journal of the ACM*, 36(4):929–965, 1989.
7. P. Fischer and H.U. Simon. On learning ring-sum expansions. *SIAM Journal on Computing*, 21(1):181–192, 1992.
8. R. L. Graham, D. E. Knuth, and O. Patashnik. *Concrete Mathematics*. Addison-Wesley, Reading, MA, 1994.
9. F. Green, J. Kobler, K. Regan, T. Schwentick, and J. Toran. The power of the middle bit of a #P function. *Journal of Computer and System Sciences*, 50(3):456–467, 1998.
10. D. Helmbold, R. Sloan, and M. Warmuth. Learning integer lattices. *SIAM Journal on Computing*, 21(2):240–266., 1992.
11. A. Klivans and R. Servedio. Learning DNF in time $2^{\tilde{O}(n^{1/3})}$. In *Proceedings of the Thirty-Third Annual Symposium on Theory of Computing*, pages 258–265, 2001.

12. P. McKenzie and D. Therien. Automata theory meets circuit complexity. In *Proceedings of the International Colloquium on Automata, Languages and Programming*, pages 589–602, 1989.

13. M. Minsky and S. Papert. *Perceptrons: an introduction to computational geometry.* MIT Press, Cambridge, MA, 1968.

14. C. Ramus. Solution générale d'un problème d'analyse combinatoire. *J. Reine Agnew. Math.*, 11:353–355, 1834.

15. J. Riordan. *An Introduction to Combinatorial Analysis.* Wiley, New York, 1958.

16. J. Riordan. *Combinatorial Identities.* Wiley, New York, 1968.

17. L. Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984.

18. A. Yao. Separating the polynomial time hierarchy by oracles. In *Proceedings of the Twenty-Sixth Annual Symposium on Foundations of Computer Science*, pages 1–10, 1985.

19. A. Yao. On *ACC* and threshold circuits. In *Proceedings of the Thirty-First Annual Symposium on Foundations of Computer Science*, pages 619–627, 1990.