

Cryptographic Limitations on Learning Boolean Formulae and Finite Automata

Michael Kearns*

AT&T Bell Laboratories
600 Mountain Avenue, Room 2A-423
Murray Hill, New Jersey 07974

Leslie Valiant[†]

Aiken Computation Laboratory
Harvard University
Cambridge, Massachusetts 02138

Abstract

In this paper we prove the intractability of learning several classes of Boolean functions in the distribution-free model (also called the Probably Approximately Correct or PAC model) of learning from examples. These results are *representation independent*, in that they hold regardless of the syntactic form in which the learner chooses to represent its hypotheses.

Our methods reduce the problems of cracking a number of well-known public-key cryptosystems to the learning problems. We prove that a polynomial-time learning algorithm for Boolean formulae, deterministic finite automata or constant-depth threshold circuits would have dramatic consequences for cryptography and number theory: in particular, such an algorithm could be used to break the RSA cryptosystem, factor Blum integers (composite numbers equivalent to 3 modulo 4), and detect quadratic residues. The results hold even if the learning algorithm is only required to obtain a slight advantage in prediction over random guessing. The techniques used demonstrate an interesting duality between learning and cryptography.

We also apply our results to obtain strong intractability results for approximating a generalization of graph coloring.

*This research was conducted while the author was at Harvard University and supported by an A.T.& T. Bell Laboratories scholarship.

[†]Supported by grants ONR-N00014-85-K-0445, NSF-DCR-8606366 and NSF-CCR-89-02500, DAAL03-86-K-0171, DARPA AFOSR 89-0506, and by SERC.

1 Introduction

In this paper we prove *representation-independent* hardness results for the distribution-free learning of several representation classes whose efficient learnability has thus far been unresolved¹. Informally, a representation-independent hardness result states that learning is difficult regardless of the *form* in which a learning algorithm represents its hypothesis, provided this hypothesis meets the quite reasonable constraint of being evaluable in polynomial time (that is, having an equivalent polynomial-size Boolean circuit). In contrast, a *representation-based* hardness result states only that learning is difficult when the hypothesis is constrained to meet some (usually strong) structural or syntactic restrictions.

We prove representation-independent hardness results for the distribution-free learning of several simple representation classes, including polynomial-size Boolean formulae, acyclic deterministic finite automata, and constant-depth threshold circuits (which may be regarded as a form of simplified neural networks). These hardness results are based on assumptions regarding the intractability of specific number-theoretic problems of interest in cryptography, namely factoring Blum integers, inverting the RSA function, and recognizing quadratic residues. Thus, a polynomial-time learning algorithm in the distribution-free model for any of the named representation classes using *any* polynomial-time evaluable hypothesis representation would immediately yield a polynomial-time algorithm for all of these cryptographic problems, which have defied efficient solution for decades, and are widely believed to be intractable.

For practical purposes, the efficient learnability of a representation class must be considered unresolved until a polynomial-time learning algorithm is discovered or until a representation-independent hardness result is proved. This is because a representation-based result stating that the class C is not efficiently learnable by the class H (modulo some complexity-theoretic assumption such as $RP \neq NP$) still leaves open the possibility that C is efficiently learnable by a different hypothesis class H' . Indeed, this possibility has been realized for several natural target classes: for instance, it is known that for any fixed constant natural number $k \geq 2$, the problem of learning k -term disjunctive normal form (DNF) formulae in the distribution-free model is NP -hard if the learning algorithm is restricted to represent its hypothesis in k -term DNF form, but there is

¹The distribution-free model of learning that we use and will define shortly is often also referred to as the *probably approximately correct* or *PAC* model of learning.

a polynomial-time learning algorithm if we relax this restriction [32]. A similar result holds for Boolean threshold functions [32].

The only previous representation-independent hardness results for distribution-free learning follow from the elegant work of Goldreich, Goldwasser and Micali [22] on constructing random functions. Their functions have many properties stronger than those mentioned here, but for our purposes we may state their result formally as follows: let $\text{CKT}_n^{p(n)}$ denote the class of Boolean circuits over n inputs with at most $p(n)$ gates, and let $\text{CKT}^{p(n)} = \cup_{n \geq 1} \text{CKT}_n^{p(n)}$. Then it is shown by Goldreich et al. [22] that if there exists a one-way function, then for some polynomial $p(n)$, $\text{CKT}^{p(n)}$ is not learnable in polynomial time (by *any* polynomial-time evaluable representation class). Pitt and Warmuth [33] then used this result to construct other hard-to-learn representation classes. For definitions and a discussion of one-way functions we refer the reader to Yao [41], Blum and Micali [13], Levin [29], and Goldreich et al. [22].

Note that in any reasonable model of learning, we intuitively do not expect to find polynomial-time learning algorithms for classes of representations that are not polynomial-time evaluable, since a learning algorithm may not even have enough time to write down a good hypothesis. More formally, Schapire [38] has shown that any representation class that is not evaluable in polynomial time cannot be learned in polynomial time in the distribution-free model.

Thus, we may informally interpret the result of Goldreich, Goldwasser and Micali as stating that not everything with a small (polynomial-size) circuit representation is efficiently learnable (assuming there is a one-way function). However, there is a large gap in computational power between the class of polynomial-size circuits and the classes that have been the subject of intense scrutiny within the computational learning theory community of late (e.g., DNF, decision trees, Boolean formulae, classes based on finite automata, restricted classes of circuits). In this paper we prove hardness results similar to those of Goldreich, Goldwasser and Micali but for much less powerful representation classes, thus clarifying the limits of efficient learnability.

The intuition behind the approach taken to obtain these results is contained in the following analogy. Consider a computer system with two users, Alice and Bob. Alice and Bob wish to communicate via an insecure channel, and it is assumed that Eve the eavesdropper is listening to this channel. We make no assumptions about Eve's behavior other than a polynomial bound on her computing resources. In this cryptographic setting, Alice and Bob wish to communicate *privately*

in spite of Eve’s nosey presence.

A classic solution to Alice and Bob’s problem is the *one-time pad*. Here Alice and Bob would physically meet in a secure room (away from Eve) and compile a large common table of random bits. Then after separating, Bob, to send a bit b to Alice, chooses the next random bit c from the common list and sends the bit $b \oplus c$ to Alice. It is easily verified that if the bit c is uniformly distributed then the encoded bit $b \oplus c$ is also uniformly distributed, regardless of the value of the cleartext message bit b . Thus Eve, regardless of computation time, is provably unable to gain any information about the cleartext messages from listening to the channel between Alice and Bob. Alice, however, also knows the random bit c , and so may decode by computing $(b \oplus c) \oplus c = b$.

There are some obvious practical problems with the one-time pad. Foremost among these is the need for Alice and Bob to meet in person and compile the table of random bits; in a network of thousands of computers, having every pair of users meet clearly defeats the point of using computers in the first place. In response to complaints such as these and also more subtle security concerns, the field of *public-key cryptography* was created by Diffie and Hellman [18].

Public-key cryptography solves the problem of Alice and Bob via the use of *trapdoor functions*. Informally, a trapdoor function is one that can be computed in polynomial time (i.e., it is easy to compute $f(x)$ on input x) but cannot be inverted in polynomial time (i.e., it is hard to compute x on input $f(x)$) — unless one is the “creator” of the function, in which case one possesses a piece of “trapdoor” information that makes inversion possible in polynomial time. Now rather than meeting with Bob in person, Alice “creates” a trapdoor function f and *publishes* a program for computing f (which reveals no information about f^{-1}) in a directory that is available to everyone — Bob and Eve included. To send the message x to Alice, Bob simply computes $f(x)$ and sends it to Alice. Eve, seeing only $f(x)$ on the channel and not possessing the trapdoor, is unable to recover the message x in polynomial time. Alice, being the creator of f and thus having the trapdoor, can efficiently invert Bob’s ciphertext and recover x .

Our approach is based on viewing Eve as a learning algorithm. Note that since a program for f is available to Eve, she may create as many pairs of the form $(f(x), x)$ that she likes simply by choosing x and then computing $f(x)$. If we set $y = f(x)$, we see that such pairs have the form $(y, f^{-1}(y))$, and can thus be regarded as “examples” of the inverse function f^{-1} . Thus, from the learning perspective, public-key cryptography assumes the existence of functions that are not

learnable from examples, since if Eve could learn f^{-1} efficiently from examples of its input-output behavior, she could then decode messages sent from Bob to Alice! Furthermore, note that the inverse function f^{-1} is “simple” in the sense that it does have a small circuit (determined by the trapdoor, which Alice has access to and uses for decoding); thus from an information-theoretic standpoint the learning problem is “fair”, as opposed to the one-time pad, where there is no small circuit underlying the communication between Alice and Bob, just a large random bit table.

Thus we see that recent developments in the theory of cryptography provide us with simple functions that are difficult to learn. Our approach in this paper is based on refining the functions provided by cryptography in an attempt to find the *simplest* functions that are difficult to learn.

The outline of the paper is as follows: in Section 2, we provide definitions for the distribution-free model of learning, adapted from Valiant [39]. Then in Section 3, we discuss previous hardness results for learning, both of the representation-based and representation-independent type. Section 4 gives the needed definitions and background from cryptography.

In Section 5 we develop simple representation classes based on cryptographic functions and prove that learning these classes is as difficult as breaking the associated cryptosystems; in Section 6 these results are applied to prove the difficulty of learning Boolean formulae, finite automata, and threshold circuits. In Section 7 we give a generalized method for proving hardness results for learning based on *any* trapdoor function. Section 8 applies our learning results to give strong hardness results for approximating the optimal solution for various combinatorial optimization problems, including a generalization of graph coloring.

2 Definitions for Distribution-free Learning

2.1 Representing subsets of a domain

Concept classes and their representation. Let X be a set called a *domain* (also sometimes referred to as the *instance space*). We think of X as containing encodings of all objects of interest to us in our learning problem. For example, each instance in X may represent a different object in a particular room, with discrete attributes representing properties such as color, and continuous values representing properties such as height. The goal of a learning algorithm is to infer some unknown subset of X , called a *concept*, chosen from a known

concept class.

For computational purposes we always need a way of *naming* or *representing* concepts. Thus, we formally define a *representation class over X* to be a pair (σ, C) , where $C \subseteq \{0, 1\}^*$ and σ is a mapping $\sigma : C \rightarrow 2^X$ (here 2^X denotes the power set of X). In the case that the domain X has real-valued components, we sometimes assume $C \subseteq (\{0, 1\} \cup R)^*$, where R is the set of real numbers. For $c \in C$, $\sigma(c)$ is called a *concept over X* ; the image space $\sigma(C)$ is the *concept class* that is *represented* by (σ, C) . For $c \in C$, we define $pos(c) = \sigma(c)$ (the *positive examples* of c) and $neg(c) = X - \sigma(c)$ (the *negative examples* of c). The domain X and the mapping σ will usually be clear from the context, and we will simply refer to the *representation class* C . We will sometimes use the notation $c(x)$ to denote the value of the characteristic function of $\sigma(c)$ on the domain point x ; thus $x \in pos(c)$ ($x \in neg(c)$, respectively) and $c(x) = 1$ ($c(x) = 0$, respectively) are used interchangeably. We assume that domain points $x \in X$ and representations $c \in C$ are efficiently encoded using any of the standard schemes (see Garey and Johnson [20]), and denote by $|x|$ and $|c|$ the length of these encodings measured in bits (or in the case of real-valued domains, some other reasonable measure of length that may depend on the model of arithmetic computation used; see Aho, Hopcroft and Ullman [2]).

Parameterized representation classes. In this paper we study *parameterized* classes of representations. Here we have a stratified domain $X = \cup_{n \geq 1} X_n$ and representation class $C = \cup_{n \geq 1} C_n$. The parameter n can be regarded as an appropriate measure of the complexity of concepts in $\sigma(C)$ (such as the number of domain attributes), and we assume that for a representation $c \in C_n$ we have $pos(c) \subseteq X_n$ and $neg(c) = X_n - pos(c)$. For example, X_n may be the set $\{0, 1\}^n$, and C_n the class of all Boolean formulae over n variables whose length is at most n^2 . Then for $c \in C_n$, $\sigma(c)$ would contain all satisfying assignments of the formula c .

Efficient evaluation of representations. In general, we will be primarily concerned with learning algorithms that are computationally efficient. In order to prevent this demand from being vacuous, we need to insure that the *hypotheses* output by a learning algorithm can be efficiently evaluated as well. Thus if C is a representation class over X , we say that C is *polynomially evaluable* if there is a polynomial-time *evaluation algorithm* A that on input a

representation $c \in C$ and a domain point $x \in X$ outputs $c(x)$. Note that if a class C is polynomially evaluable, then each representation $c \in C$ has an equivalent polynomial-size circuit, obtained by hard-wiring the representation input of A to be c , and converting the resulting polynomial-time algorithm (now accepting the single input $x \in X$) to a polynomial-size circuit using standard techniques. All representation classes considered here are polynomially evaluable. It is worth mentioning at this point that Schapire [38] has shown that if a representation class is not polynomially evaluable, then it is not efficiently learnable in our model. Thus, perhaps not surprisingly we see that classes that are not polynomially evaluable are not only “unfair” as learning problems but also intractable.

Samples. A *labeled example* from a domain X is a pair $\langle x, b \rangle$, where $x \in X$ and $b \in \{0, 1\}$. A *labeled sample* $S = \langle x_1, b_1 \rangle, \dots, \langle x_m, b_m \rangle$ from X is a finite sequence of labeled examples from X . If C is a representation class, a *labeled example of $c \in C$* is a labeled example $\langle x, c(x) \rangle$, where $x \in X$. A *labeled sample of c* is a labeled sample S where each example of S is a labeled example of c . In the case where all labels b_i or $c(x_i)$ are 1 (0, respectively), we may omit the labels and simply write S as a list of points x_1, \dots, x_m , and we call the sample a *positive (negative, respectively) sample*.

We say that a representation h and an example $\langle x, b \rangle$ *agree* if $h(x) = b$; otherwise they *disagree*. We say that a representation h and a sample S are *consistent* if h agrees with each example in S ; otherwise they are *inconsistent*.

2.2 Distribution-free learning

Distributions on examples. On any given execution, a learning algorithm for a representation class C will be receiving examples of a single distinguished representation $c \in C$. We call this distinguished c the *target representation*. Examples of the target representation are generated probabilistically as follows: let D_c^+ be a fixed but arbitrary probability distribution over $pos(c)$, and let D_c^- be a fixed but arbitrary probability distribution over $neg(c)$. We call these distributions the *target distributions*. When learning c , learning algorithms will be given access to two oracles, POS and NEG , that behave as follows: oracle POS (NEG , respectively) returns in unit time a positive (negative, respectively) example of the target representation, drawn randomly according to the target distribution D_c^+ (D_c^- , respectively).

The distribution-free model is sometimes defined in the literature with a single target distribution over the entire domain; the learning algorithm is then given labeled examples of the target concept drawn from this distribution. We choose to explicitly separate the distributions over the positive and negative examples to facilitate the study of algorithms that learn using only positive examples or only negative examples. These models, however, are equivalent with respect to polynomial-time computation, in the sense that any class learnable in polynomial time in one model is learnable in polynomial time in the other model, as shown by Haussler et al. [24].

Given a fixed target representation $c \in C$, and given fixed target distributions D_c^+ and D_c^- , there is a natural measure of the *error* (with respect to c , D_c^+ and D_c^-) of a representation h from a representation class H . We define $e_c^+(h) = D_c^+(neg(h))$ (i.e., the weight of the set $neg(h)$ under the probability distribution D_c^+) and $e_c^-(h) = D_c^-(pos(h))$ (the weight of the set $pos(h)$ under the probability distribution D_c^-). Note that $e_c^+(h)$ (respectively, $e_c^-(h)$) is simply the probability that a random positive (respectively, negative) example of c is identified as negative (respectively, positive) by h . If both $e_c^+(h) < \epsilon$ and $e_c^-(h) < \epsilon$, then we say that h is an ϵ -good hypothesis (with respect to c , D_c^+ and D_c^-); otherwise, h is ϵ -bad. We define the *accuracy* of h to be the value $\min(1 - e_c^+(h), 1 - e_c^-(h))$.

It is worth noting that our definitions so far assume that the hypothesis h is deterministic. However, this need not be the case; for example, we can instead define $e_c^+(h)$ to be the probability that h classifies a random positive example of c as negative, where the probability is now over both the random example and the coin flips of h . All of the results presented here hold under these generalized definitions.

When the target representation c is clear from the context, we will drop the subscript c and simply write D^+ , D^- , e^+ and e^- .

In the definitions that follow, we will demand that a learning algorithm produce with high probability an ϵ -good hypothesis regardless of the target representation and target distributions. While at first this may seem like a strong criterion, note that the error of the hypothesis output is always measured with respect to the same target distributions on which the algorithm was trained. Thus, while it is true that certain examples of the target representation

may be extremely unlikely to be generated in the training process, these same examples intuitively may be “ignored” by the hypothesis of the learning algorithm, since they contribute a negligible amount of error.

Learnability. Let C and H be representation classes over X . Then C is *learnable from examples by H* if there is a (probabilistic) algorithm A with access to POS and NEG , taking inputs ϵ, δ , with the property that for any target representation $c \in C$, for any target distributions D^+ over $pos(c)$ and D^- over $neg(c)$, and for any inputs $0 < \epsilon, \delta < 1$, algorithm A halts and outputs a representation $h_A \in H$ that with probability greater than $1 - \delta$ satisfies $e^+(h_A) < \epsilon$ and $e^-(h_A) < \epsilon$.

We call C the *target class* and H the *hypothesis class*; the output $h_A \in H$ is called the *hypothesis* of A . A will be called a *learning algorithm* for C . If C and H are polynomially evaluatable, and A runs in time polynomial in $1/\epsilon, 1/\delta$ and $|c|$ then we say that C is *polynomially learnable from examples by H* ; if C is parameterized we also allow the running time of A to have polynomial dependence on the parameter n . We assume that A is given the values of n and $|c|$ as input; the latter assumption is without loss of generality [24].

We will drop the phrase “from examples” and simply say that C is *learnable by H* , and C is *polynomially learnable by H* . We say C is *polynomially learnable* to mean that C is polynomially learnable by H for some polynomially evaluatable H . We will sometimes call ϵ the *accuracy parameter* and δ the *confidence parameter*.

Thus, we ask that for any target representation and any target distributions, a learning algorithm finds an ϵ -good hypothesis with probability at least $1 - \delta$. A primary goal of research in this model is to discover which representation classes C are polynomially learnable.

Note that in the above definitions, we allow the learning algorithm to output hypotheses from some class H that is possibly different from C , as opposed to the natural choice $C = H$. While in general we assume that H is at least as powerful as C (that is, $C \subseteq H$), we will see that in some cases for computational reasons we may not wish to restrict H beyond it being polynomially evaluatable. If the algorithm produces an accurate and easily evaluated hypothesis, then our learning problem is essentially solved, and the actual form of the hypothesis is of secondary concern.

We refer to this model as the *distribution-free* model, to emphasize that we seek algorithms that work for any target distributions. It is also known in the literature as the *probably approximately correct* model. We also occasionally refer to the model as that of *strong learnability* (to mean learnability by some polynomially evaluatable representation class H), in contrast with the notion of *weak learnability* defined below.

Weak learnability. We will also consider a distribution-free model in which the hypothesis of the learning algorithm is required to perform only slightly better than random guessing.

Let C and H be representation classes over X . Then C is *weakly learnable from examples by H* if there is a polynomial p and a (probabilistic) algorithm A with access to POS and NEG , taking input δ , with the property that for any target representation $c \in C$, for any target distributions D^+ over $pos(c)$ and D^- over $neg(c)$, and for any input value $0 < \delta < 1$, algorithm A halts and outputs a representation $h_A \in H$ that with probability greater than $1 - \delta$ satisfies $e^+(h_A) < 1/2 - 1/p(|c|)$ and $e^-(h_A) < 1/2 - 1/p(|c|)$.

Thus, the accuracy of h_A must be at least $1/2 + 1/p(|c|)$. A will be called a *weak learning algorithm* for C . If C and H are polynomially evaluatable, and A runs in time polynomial in $1/\delta$ and $|c|$ we say that C is *polynomially weakly learnable by H* and C is *polynomially weakly learnable* if it is weakly learnable by H for some polynomially evaluatable H . In the case that the target class C is parameterized, we allow the polynomial p and the running time to depend on the parameter n . We will usually explicitly restrict $|c|$ to be polynomial in n , and thus may assume p depends on n alone.

We may intuitively think of weak learning as the ability to detect some slight bias separating positive and negative examples, where the advantage gained over random guessing diminishes as the complexity of the problem grows. Our main use of the weak learning model is in proving the strongest possible hardness results.

Distribution-specific learnability. The models for learnability described above demand that a learning algorithm work regardless of the distributions on the examples. We will sometimes relax this condition, and consider these models under restricted target distributions, for instance the uniform distribution. Here the definitions are the same as before, except that we ask that the performance criteria for learnability be met only under these restricted target

distributions.

2.3 Some representation classes

We now define the representation classes whose learnability we will study. In this paper the domain X_n is always $\{0, 1\}^n$ and the mapping σ simply maps each circuit to its set of satisfying assignments. The classes defined below are all parameterized; for each class we will define the subclasses C_n , and then C is defined by $C = \cup_{n \geq 1} C_n$.

Boolean Formulae: The representation class BF_n consists of all Boolean formulae over the Boolean variables x_1, \dots, x_n .

Boolean Circuits: The representation class CKT_n consists of all Boolean circuits over input variables x_1, \dots, x_n .

Threshold Circuits: A *threshold gate* over input variables x_1, \dots, x_n is defined by a value $1 \leq t \leq n$ such that the gate outputs 1 if and only if at least t of the input bits are set to 1. We let TC_n denote the class of all circuits of threshold gates over x_1, \dots, x_n . For constant d , $d\text{TC}_n$ denotes the class of all threshold circuits in TC_n with depth at most d .

Acyclic Finite Automata: The representation class ADFA_n consists of all deterministic finite automata that accept only strings of length n , that is, all deterministic finite automata M such that the language $L(M)$ accepted by M satisfies $L(M) \subseteq \{0, 1\}^n$.

We will also frequently discuss computations performed by the circuit class $\text{NC}^1 = \cup_{n \geq 1} \text{NC}_n^1$, where NC_n^1 is the class of circuits consisting of AND, OR and NOT gates of fan-in two having size polynomial in n and depth logarithmic in n .

2.4 Other definitions and notation

Chernoff bounds. We shall make extensive use of the following bounds on the area under the tails of the binomial distribution. For $0 \leq p \leq 1$ and m a positive integer, let $LE(p, m, r)$ denote the probability of at most r successes in m independent trials of a Bernoulli variable with probability of success p , and let $GE(p, m, r)$ denote the probability of at least r successes. Then for $0 \leq \alpha \leq 1$,

Fact CB1. $LE(p, m, (1 - \alpha)mp) \leq e^{-\alpha^2 mp/2}$

and

Fact CB2. $GE(p, m, (1 + \alpha)mp) \leq e^{-\alpha^2 mp/3}$

These bounds in the form they are stated are from the paper of Angluin and Valiant [9] and follow from Chernoff [17]. Although we will make frequent use of Fact CB1 and Fact CB2, we will do so in varying levels of detail, depending on the complexity of the calculation involved. However, we are primarily interested in Chernoff bounds for the following consequence of Fact CB1 and Fact CB2: given an event E of probability p , we can obtain an estimate \hat{p} of p by drawing m points from the distribution and letting \hat{p} be the frequency with which E occurs in this sample. Then for m polynomial in $1/p$ and $1/\alpha$, \hat{p} satisfies $p/2 < \hat{p} < 2p$ with probability at least $1 - \alpha$. If we also allow m to depend polynomially on $1/\beta$, we can obtain an estimate \hat{p} such that $p - \beta < \hat{p} < p + \beta$ with probability at least $1 - \alpha$ (see for instance the paper of Angluin and Laird [8]).

Notational conventions. Let $E(x)$ be an event and $\psi(x)$ a random variable that depend on a parameter x that takes on values in a set X . Then for $X' \subseteq X$, we denote by $\Pr_{x \in X'}[E(x)]$ the probability that E occurs when x is drawn uniformly at random from X' . Similarly, $\mathbf{E}_{x \in X'}[\psi(x)]$ is the expected value of ψ when x is drawn uniformly at random from X' . We also need to work with distributions other than the uniform distribution; thus if P is a distribution over X we use $\Pr_{x \in P}[E(x)]$ and $\mathbf{E}_{x \in P}[\psi(x)]$ to denote the probability of E and the expected value of ψ , respectively, when x is drawn according to the distribution P . When E or ψ depend on several parameters that are drawn from different distributions we use multiple subscripts. For example, $\Pr_{x_1 \in P_1, x_2 \in P_2, x_3 \in P_3}[E(x_1, x_2, x_3)]$ denotes the probability of event E when x_1 is drawn from distribution P_1 , x_2 from P_2 , and x_3 from P_3 (all draws being independent).

3 Previous Hardness Results for Learning

The initial paper defining the distribution-free model [39] gave the first polynomial-time learning algorithms in this model. It showed that the class of monomials is polynomially learnable, as are the classes k CNF and k DNF (with time complexity $O(n^k)$). For each of these algorithms, the

hypothesis class is the same as the target class; that is, in each case C is polynomially learnable by C .

Pitt and Valiant [32] subsequently observed that the classes k -TERM-DNF and k -CLAUSE-CNF, when viewed as functions, are properly contained within the classes k CNF and k DNF, respectively. Combined with the results above [39], this shows that for fixed k , the class k -TERM-DNF is polynomially learnable by k CNF, and the class k -CLAUSE-CNF is polynomially learnable by k DNF. More surprisingly, Pitt and Valiant prove that for any fixed $k \geq 2$, learning k -TERM-DNF by k -TERM-DNF and learning k -CLAUSE-CNF by k -CLAUSE-CNF are NP -hard problems.

These results are important in that they demonstrate the tremendous computational advantage that may be gained by a judicious change of hypothesis representation. This can be viewed as a limited but provable confirmation of the rule of thumb in artificial intelligence that *representation is important*. By moving to a more powerful hypothesis class H instead of insisting on the more “natural” choice $H = C$, we move from an NP -hard problem to a polynomial-time solution. This may be explained intuitively by the observation that while the constraint $H = C$ may be significant enough to render the learning task intractable, a richer hypothesis representation allows a greater latitude for expressing the learned formula. Later we shall see that using a larger hypothesis class inevitably requires a larger sample complexity; thus the designer of a learning algorithm may sometimes be faced with a trade-off between computation time and required sample size.

In discussing hardness results, we distinguish between two types: *representation-based* hardness results and *representation-independent* hardness results. Briefly, representation-based hardness results state that for some *fixed* representation classes C and H , learning C by H is hard in some computational sense (such as NP -hardness). Thus, the aforementioned result on the difficulty of learning k -TERM-DNF by k -TERM-DNF is representation-based. In contrast, a representation-independent hardness result says that for fixed C and *any* polynomially evaluable H , learning C by H is hard.

Representation-based hardness results are interesting for a number of reasons, two of which we have already mentioned: they can be used to give formal verification to the importance of hypothesis representation, and for practical reasons it is important to study the *least* expressive class H that can be used to learn C , since the choice of hypothesis representation can greatly affect resource complexity (such as the number of examples required) even for those classes already known

to be polynomially learnable.

However, since a representation-based hardness result dismisses the polynomial learnability of C only with respect to the *fixed* hypothesis class H , such results leave something to be desired in the quest to classify learning problems as “easy” or “hard”. For example, we may be perfectly willing to settle for an efficient algorithm learning C by H for some more expressive H if we know that learning C by C is NP -hard. Thus for practical purposes we must regard the polynomial learnability of C as not entirely resolved until we either find an efficient learning algorithm or we prove that learning C by H is hard for *any* reasonable H , that is, until we prove a representation-independent hardness result for C .

Gold [21] gave the first representation-based hardness results that apply to the distribution-free model of learning. He proves that the problem of finding the smallest deterministic finite automaton consistent with a given sample is NP -complete; the results of Haussler et al. [24] can be easily applied to Gold’s result to prove that learning deterministic finite automata of size n by deterministic finite automata of size n cannot be accomplished in polynomial time unless $RP = NP$. There are some technical issues involved in properly defining the problem of learning finite automata in the distribution-free model; see Pitt and Warmuth [33] for details. Gold’s results were improved by Li and Vazirani [30], who show that finding an automaton $9/8$ larger than the smallest consistent automaton is still NP -complete.

As we have already discussed, Pitt and Valiant [32] prove that for $k \geq 2$, learning k -TERM-DNF by k -TERM-DNF is NP -hard by giving a randomized reduction from a generalization of the graph coloring problem. Even stronger, for $k \geq 6$, they prove that even if the hypothesis DNF formulae is allowed to have $2k - 3$ terms, k -TERM-DNF cannot be learned in polynomial time unless $RP = NP$. These results hold even when the target formulae are restricted to be monotone and the hypothesis formulae is allowed to be nonmonotone. Dual results hold for the problem of learning k -CLAUSE-CNF. Pitt and Valiant also prove that μ -formulae (Boolean formulae in which each variable occurs at most once, sometimes called *read-once*) cannot be learned by μ -formulae in polynomial time, and that Boolean threshold functions cannot be learned by Boolean threshold functions in polynomial time, unless $RP = NP$.

Pitt and Warmuth [34] dramatically improved the results of Gold by proving that deterministic finite automata of size n cannot be learned in polynomial time by deterministic finite automata of

size n^α for any fixed value $\alpha \geq 1$ unless $RP = NP$. Their results leave open the possibility of an efficient learning algorithm using deterministic finite automata whose size depends on ϵ and δ , or an algorithm using some entirely different representation of the sets accepted by automata. This possibility is addressed and dismissed (modulo cryptographic assumptions) by the results in of this paper.

Hancock [23] has shown that learning decision trees of size n by decision trees of size n cannot be done in polynomial time unless $RP = NP$. Representation-based hardness results for learning various classes of neural networks can also be derived from the results of Judd [25] and Blum and Rivest [12].

The first representation-independent hardness results for the distribution-free model follow from the work of Goldreich, Goldwasser and Micali [22], whose true motivation was to find easy-to-compute functions whose output on random inputs appears random to all polynomial-time algorithms. A simplified and weakened statement of their result is that the class of polynomial-size Boolean circuits is not polynomially learnable by *any* polynomially evaluable H , provided that there exists a one-way function (see Yao [41]). Pitt and Warmuth [33] defined a general notion of reducibility for learning and gave a number of other representation classes that are not polynomially learnable under the same assumption by giving reductions from the learning problem for polynomial-size circuits. One of the main contributions of the research presented here is representation-independent hardness results for much simpler classes than those addressed by Goldreich et al. [22] or Pitt and Warmuth [33], among them the classes of Boolean formulae, acyclic deterministic finite automata and constant-depth threshold circuits.

4 Background and Definitions from Cryptography

Some basic number theory. For an introduction to number theory that is relevant to cryptography, we refer the reader to the work of Angluin [5] and Kranakis [28]. For N a natural number, Z_N will denote the ring of integers modulo N , and Z_N^* will denote the multiplicative group modulo N . Thus $Z_N = \{x : 0 \leq x \leq N-1\}$ and $Z_N^* = \{x : 1 \leq x \leq N-1 \text{ and } \gcd(x, N) = 1\}$, where $\gcd(x, N)$ denotes the greatest common divisor of x and N . The *Euler totient function* φ is defined by $\varphi(N) = |Z_N^*|$. For $x \in Z_N^*$, we say that x is a *quadratic residue* modulo N if there is an $a \in Z_N^*$ such that $x = a^2 \text{ mod } N$. We denote by QR_N the set of all

quadratic residues in Z_N^* . For a prime p and $x \in Z_p^*$, we define the *Legendre symbol* of x with respect to p by $L(x, p) = 1$ if x is a quadratic residue modulo p , and $L(x, p) = -1$ otherwise. For $N = p \cdot q$, where p and q are prime, we define the *Jacobi symbol* of $x \in Z_N^*$ with respect to N by $J(x, N) = L(x, p) \cdot L(x, q)$. Since x is a quadratic residue modulo N if and only if it is a quadratic residue modulo p and modulo q , it follows that $J(x, N) = -1$ implies that x is not a quadratic residue modulo N . However, $J(x, N) = 1$ does not necessarily imply that x is a quadratic residue mod N . For any integer N , we define the set $Z_N^*(+1) = \{x \in Z_N^* : J(x, N) = 1\}$. A *Blum integer* is an integer of the form $p \cdot q$, where p and q are primes both congruent to 3 modulo 4.

We will make use of the following facts from number theory.

Fact NT1. On inputs x and N , $\gcd(x, N)$ can be computed in polynomial time.

Fact NT2. For p a prime and $x \in Z_p^*$, $L(x, p) = x^{(p-1)/2} \bmod p$.

Fact NT3. On inputs x and N , $J(x, N)$ can be computed in polynomial time.

Fact NT4. For $N = p \cdot q$ where p and q are prime, $|Z_N^*(+1)| = |Z_N^*|/2$ and $|QR_N| = |Z_N^*|/4$.

Fact NT5. For any $x \in Z_N^*$, $x^{\varphi(N)} = 1 \bmod N$.

The RSA encryption function. Let p and q be primes of length l , and let $N = p \cdot q$. Let e be an *encrypting exponent* such that $\gcd(e, \varphi(N)) = 1$ and d a *decrypting exponent* such that $d \cdot e = 1 \bmod \varphi(N)$. The existence of such a d is guaranteed for all elements e for which $\gcd(e, \varphi(N)) = 1$. The *RSA encryption function* [37] is then defined for all $x \in Z_N$ by

$$RSA(x, N, e) = x^e \bmod N.$$

Note that decryption can be accomplished by exponentiation mod N :

$$(x^e)^d = x^{e \cdot d} \bmod N = x^{1+i \cdot \varphi(N)} \bmod N = x \bmod N$$

for some natural number i by Fact NT5 because $e \cdot d = 1 \bmod \varphi(N)$.

Thus, following the informal intuition of Section 1, we think of Alice as generating the product $N = p \cdot q$; since she also knows p and q , she can generate both e (which she publishes along with N , thus yielding an encryption program) and d (the “trapdoor”, which she keeps private).

There is currently no known polynomial-time algorithm for *inverting* the RSA encryption function — that is, the problem of computing x on inputs $RSA(x, N, e), N$ and e . Furthermore, the following result from Alexi et al. [4] indicates that determining the least significant bit of x is as hard as inverting RSA (which amounts to determining *all* the bits of x).

Theorem 1 (Alexi et al. [4]) *Let x, N and e be as above. Then with respect to probabilistic polynomial-time reducibility, the following problems are equivalent:*

- (1) *On input $RSA(x, N, e), N$ and e , output x .*
- (2) *On input $RSA(x, N, e), N$ and e , output $LSB(x)$ with probability exceeding $1/2 + 1/p(l)$, where p is any fixed polynomial, $l = \log N$ is the length of N , and $LSB(x)$ denotes the least significant bit of x . The probability is taken over x chosen uniformly from Z_N and any coin tosses of A .*

The Rabin and modified Rabin encryption functions. The *Rabin encryption function* [35] is specified by two primes p and q of length l . For $N = p \cdot q$ and $x \in Z_N^*$, we define

$$R(x, N) = x^2 \bmod N.$$

In this scheme the trapdoor is the factorization of N , which allows Alice to compute square roots modulo N , and thus to decrypt. Known results regarding the security of the Rabin function include the following:

Theorem 2 (Rabin [35]) *Let x and N be as above. Then with respect to probabilistic polynomial-time reducibility, the following problems are equivalent:*

- (1) *On input N , output a nontrivial factor of N .*
- (2) *On input N and $R(x, N)$, output a y such that $R(y, N) = R(x, N)$.*

Furthermore, this reduction still holds when N is restricted to be a Blum integer in both problems. The *modified Rabin encryption function* [4] is specified by two primes p and q of length l , both congruent to 3 modulo 4. Let $N = p \cdot q$ (thus N is a Blum integer). We define a subset M_N of Z_N^* by

$$M_N = \left\{ x : 0 \leq x \leq \frac{N}{2} \text{ and } x \in Z_N^*(+1) \right\}.$$

For $x \in M_N$, the modified Rabin encryption function is then

$$MR(x, N) = x^2 \bmod N \text{ if } x^2 \bmod N \in M_N$$

$$MR(x, N) = (N - x^2) \bmod N \text{ otherwise.}$$

This defines a 1-1 map from M_N onto M_N .

Theorem 3 (Alexi et al. [4]) *Let x and N be as above. Then with respect to probabilistic polynomial-time reducibility, the following problems are equivalent:*

- (1) *On input $MR(x, N)$ and N , output x .*
- (2) *On input $MR(x, N)$ and N , output $LSB(x)$ with probability exceeding $1/2 + 1/p(l)$, where p is any fixed polynomial and $l = \log N$ is the length of N . The probability is taken over x chosen uniformly from M_N and any coin tosses of A .*

For Blum integers, $R(x, N)$ is a 1-1 mapping of QR_N . Hence if $MR(x, N)$ is invertible then we can invert $R(x, N)$ by attempting to invert MR for both the values $R(x, N)$ and $N - R(x, N)$, and succeeding for just the right one of these. Hence Theorems 2 and 3 together imply that Problem (2) in Theorem 3 is equivalent to factoring Blum integers (with respect to probabilistic polynomial-time reducibility), a problem for which no polynomial-time algorithm is known.

The Quadratic Residue Assumption. Let $N = p \cdot q$, where p and q are primes of length l . For each $x \in Z_N^*(+1)$, define $QR(x, N) = 1$ if x is a quadratic residue mod N and $QR(x, N) = 0$ otherwise. Then the *Quadratic Residue Assumption* states that if A is any probabilistic polynomial-time algorithm that takes N and x as input, then for infinitely many N we have

$$\Pr[A(N, x) = QR(x, N)] < \frac{1}{2} + \frac{1}{p(l)}$$

where p is any fixed polynomial. The probability is taken over x chosen uniformly from the set $Z_N^*(+1)$ and any coin tosses of A . As in the Rabin scheme, knowledge of the factors of N allows Alice to compute square roots modulo N and thus to determine if an element is a quadratic residue.

5 Hard Learning Problems Based on Cryptographic Functions

In this section we construct hard learning problems based on the number-theoretic encryption functions described above. For each such function, we first define a representation class based on the function. For each possible target representation in this class, we then describe the *relevant examples* for this representation. These are the only examples with non-zero probability in the hard target distributions we define. We then proceed to prove the difficulty of even *weakly* learning the representation class under the chosen distributions, based on some standard cryptographic assumption on the security of the underlying encryption function. Finally, we show the ease of *evaluating* the representation class: more precisely, we show that each representation in the class can be computed by an NC^1 circuit (a polynomial-size, log-depth circuit of standard fan-in 2 Boolean gates). In Section 6 we apply these results to prove that weakly learning Boolean formulae, finite automata, constant-depth threshold circuits and a number of other representation classes is hard under cryptographic assumptions.

We adopt the following notation: if a_1, \dots, a_m are natural numbers, then $\text{binary}(a_1, \dots, a_m)$ is the binary representation of the sequence a_1, \dots, a_m . The relevant examples we construct will be of the form

$$\langle \text{binary}(a_1, \dots, a_m), b \rangle$$

where b is a bit indicating whether the example is positive or negative. We denote by $\text{powers}(z, N)$ the sequence of natural numbers

$$z \bmod N, z^2 \bmod N, z^4 \bmod N, \dots, z^{2^{\lceil \log N \rceil}} \bmod N$$

which are the first $\lceil \log N \rceil + 1$ successive square powers of z modulo N .

In the following subsections, we will define representation classes C_n based on the number-theoretic function families described above. Representations in C_n will be over the domain $\{0, 1\}^n$; relevant examples with length less than n will implicitly be assumed to be padded to length n . Since only the relevant examples will have non-zero probability, we assume that all non-relevant examples are negative examples of the target representation.

5.1 A learning problem based on RSA

The representation class C_n : Let l be the largest natural number satisfying $4l^2 + 8l + 2 \leq n$.

Each representation in C_n is defined by a triple (p, q, e) and this representation will be denoted $r_{(p,q,e)}$. Here p and q are primes of exactly l bits and $e \in Z_{\varphi(N)}^*$, where $N = p \cdot q$ (thus, $\gcd(e, \varphi(N)) = 1$).

Relevant examples for $r_{(p,q,e)} \in C_n$: A relevant example of $r_{(p,q,e)} \in C_n$ is of the form

$$\langle \text{binary}(\text{powers}(\text{RSA}(x, N, e), N), N, e), \text{LSB}(x) \rangle$$

where $x \in Z_N$. Note that since the length of N is at most $2l + 1$, the length of such an example in bits is at most $(2l + 1)(2l + 1) + (2l + 1) + (2l + 1) = 4l^2 + 8l + 2 \leq n$. The target distribution D^+ for $r_{(p,q,e)}$ is uniform over the relevant positive examples of $r_{(p,q,e)}$ (i.e., those for which $\text{LSB}(x) = 1$) and the target distribution D^- is uniform over the relevant negative examples (i.e., those for which $\text{LSB}(x) = 0$).

Difficulty of weakly learning $C = \cup_{n \geq 1} C_n$: Suppose that A is a polynomial-time weak learning algorithm for C . We now describe how we can use algorithm A to invert the RSA encryption function. Let N be the product of two unknown l -bit primes p and q , and let $e \in Z_{\varphi(N)}^*$. Then given only N and e , we run algorithm A . Each time A requests a positive example of $r_{(p,q,e)}$, we uniformly choose an $x \in Z_N$ such that $\text{LSB}(x) = 1$ and give the example

$$\langle \text{binary}(\text{powers}(\text{RSA}(x, N, e), N), N, e), 1 \rangle$$

to A . Note that we can generate such an example in polynomial time on input N and e . This simulation generates the target distribution D^+ . Each time that A requests a negative example of $r_{(p,q,e)}$, we uniformly choose an $x \in Z_N$ such that $\text{LSB}(x) = 0$ and give the example

$$\langle \text{binary}(\text{powers}(\text{RSA}(x, N, e), N), N, e), 0 \rangle$$

to A . Again, we can generate such an example in polynomial time, and this simulation generates the target distribution D^- . Let h_A be the hypothesis output by algorithm A following this simulation. Then given $r = \text{RSA}(x, N, e)$ for some unknown x chosen uniformly from Z_N , $h_A(\text{binary}(\text{powers}(r, N), N, e)) = \text{LSB}(x)$ with probability at least $1/2 + 1/p(l)$ for some polynomial p by the definition of weak learning because n and l are polynomially related. Thus we have a polynomial advantage for inverting the least significant bit of RSA. This allows us to invert RSA by the results of Alexi et al. [4] given as Theorem 1.

Ease of evaluating $r_{(p,q,e)} \in C_n$: For each $r_{(p,q,e)} \in C_n$, we show that $r_{(p,q,e)}$ has an equivalent NC^1 circuit. More precisely, we give a circuit that has depth $O(\log n)$ and size polynomial in n , and outputs the value of $r_{(p,q,e)}$ on inputs of the form

$$\text{binary}(\text{powers}(r, N), N, e)$$

where $N = p \cdot q$ and $r = \text{RSA}(x, N, e)$ for some $x \in Z_N$. Thus, the representation class $C = \cup_{n \geq 1} C_n$ is contained in (nonuniform) NC^1 .

Since $e \in Z_{\varphi(N)}^*$, there is a $d \in Z_{\varphi(N)}^*$ such that $e \cdot d = 1 \pmod{\varphi(N)}$ (d is just the decrypting exponent for e). Thus, $r^d \pmod N = x^{e \cdot d} \pmod N = x \pmod N$. Hence the circuit for $r_{(p,q,e)}$ simply multiplies together the appropriate powers of r (which are always explicitly provided in the input) to compute $r^d \pmod N$, and outputs the least significant bit of the resulting product. This is an NC^1 step by the iterated product circuits of Beame, Cook and Hoover [10].

5.2 A learning problem based on quadratic residues

The representation class C_n : Let l be the largest natural number satisfying $4l^2 + 6l + 2 \leq n$. Each representation in C_n is defined by a pair of l -bit primes (p, q) and this representation will be denoted $r_{(p,q)}$.

Relevant examples for $r_{(p,q)} \in C_n$: For a representation $r_{(p,q)} \in C_n$, let $N = p \cdot q$. We consider only points $x \in Z_N^*(+1)$. A relevant example of $r_{(p,q)}$ is then of the following form:

$$\langle \text{binary}(\text{powers}(x, N), N), \text{QR}(x, N) \rangle .$$

Note that the length of such an example in bits is at most $4l^2 + 6l + 2 \leq n$. The target distribution D^+ for $r_{(p,q)}$ is uniform over the relevant positive examples of $r_{(p,q)}$ (i.e., those for which $\text{QR}(x, N) = 1$) and the target distribution D^- is uniform over the relevant negative examples (i.e., those for which $\text{QR}(x, N) = 0$).

Difficulty of weakly learning $C = \cup_{n \geq 1} C_n$: Suppose that A is a polynomial-time weak learning algorithm for C . We now describe how we can use algorithm A to recognize quadratic residues. Let N be the product of two unknown l -bit primes p and q . Given only N as input, we run algorithm A . Every time A requests a positive example of $r_{(p,q)}$, we uniformly choose $y \in Z_N^*$

and give the example

$$\langle \text{binary}(\text{powers}(y^2 \bmod N, N), N), 1 \rangle$$

to A . Note that such an example can be generated in polynomial time on input N . This simulation generates the target distribution D^+ .

In order to generate the negative examples for our simulation of A , we uniformly choose $u \in Z_N^*$ until $J(u, N) = 1$. By Fact NT4, this can be done with high probability in polynomial time. The probability is $1/2$ that such a u is a non-residue modulo N . Assuming we have obtained a non-residue u , every time A requests a negative example of $r_{(p,q)}$, we uniformly choose $y \in Z_N^*$ and give to A the example

$$\langle \text{binary}(\text{powers}(uy^2 \bmod N, N), N), 0 \rangle$$

which can be generated in polynomial time. Note that if u actually is a non-residue then this simulation generates the target distribution D^- , and this run of A will with high probability produce an hypothesis h_A with accuracy at least $1/2 + 1/p(l)$ with respect to D^+ and D^- , for some polynomial p (call such a run a *good* run). On the other hand, if u is actually a residue then A has been trained improperly (that is, A has been given positive examples when it requested negative examples), and no performance guarantees can be assumed. The probability of a good run of A is at least $1/2(1 - \delta)$.

We thus simulate A as described above many times, testing each hypothesis to determine if the run was a good run. To test if a good run has occurred, we first determine if h_A has accuracy at least $1/2 + 1/2p(l)$ with respect to D^+ . This can be determined with high probability by generating D^+ as above and estimating the accuracy of h_A using Fact CB1 and Fact CB2. Assuming h_A passes this test, we now would like to test h_A against the simulated distribution D^- ; however, we do not have direct access to D^- since this requires a non-residue $\bmod N$. Thus we instead estimate the probability that h_A classifies an example as positive when this example is drawn from the uniform distribution over *all* relevant examples (both positive and negative). This can be done by simply choosing $x \in Z_N^*$ uniformly and computing $h_A(\text{binary}(\text{powers}(x, N), N))$. The probability that h_A classifies such examples as positive is near $1/2$ if and only if h_A has nearly equal accuracy on D^+ and D^- . Thus by

estimating the accuracy of h_A on D^+ , we can estimate the accuracy of h_A on D^- as well, without direct access to a simulation of D^- .

We continue to run A and test until a good run of A is obtained with high probability. Then given x chosen randomly from Z_N^* ,

$$h_A(\text{binary}(\text{powers}(x, N), N)) = QR(x, N)$$

with probability at least $1/2 + 1/p(l)$, contradicting the Quadratic Residue Assumption.

Ease of evaluating $r_{(p,q)} \in C_n$: For each $r_{(p,q)} \in C_n$, we give an NC^1 circuit for evaluating the concept represented by $r_{(p,q)}$ on an input of the form

$$\text{binary}(\text{powers}(x, N), N)$$

where $N = p \cdot q$ and $x \in Z_N^*$. This circuit has four phases.

Phase I. Compute the powers

$$x \bmod p, x^2 \bmod p, x^4 \bmod p, \dots, x^{2^{2l}} \bmod p$$

and the powers

$$x \bmod q, x^2 \bmod q, x^4 \bmod q, \dots, x^{2^{2l}} \bmod q.$$

Note that the length of N is $2l$. Since for any $a \in Z_N^*$ we have that $a \bmod p = (a \bmod N) \bmod p$, these powers can be computed from the input $\text{binary}(\text{powers}(x, N), N)$ by parallel $\bmod p$ and $\bmod q$ circuits. Each such circuit involves only a division step followed by a multiplication and a subtraction. The results of Beame et al. [10] imply that these steps can be carried out by an NC^1 circuit.

Phase II. Compute $x^{(p-1)/2} \bmod p$ and $x^{(q-1)/2} \bmod q$. These can be computed by multiplying the appropriate powers $\bmod p$ and $\bmod q$ computed in Phase I. Since the iterated product of l numbers each of length l bits can be computed in NC^1 by the results of Beame et al. [10], this is also an NC^1 step.

Phase III. Determine if $x^{(p-1)/2} = 1 \bmod p$ or $x^{(p-1)/2} = -1 \bmod p$, and if $x^{(q-1)/2} = 1 \bmod q$ or $x^{(q-1)/2} = -1 \bmod q$. That these are the only cases follows from Fact NT2; furthermore, this computation determines whether x is a residue $\bmod p$ and $\bmod q$. Given the outputs of Phase II, this is clearly an NC^1 step.

Phase IV. If the results of Phase III were $x^{(p-1)/2} = 1 \pmod p$ and $x^{(q-1)/2} = 1 \pmod q$, then output 1, otherwise output 0. This is again an NC¹ step.

5.3 A learning problem based on factoring Blum integers

The representation class C_n : Let l be the largest natural number satisfying $4l^2 + 6l + 2 \leq n$. Each representation in C_n is defined by a pair of l -bit primes (p, q) , both congruent to 3 modulo 4, and this representation will be denoted $r_{(p,q)}$. Thus the product $N = p \cdot q$ is a Blum integer.

Relevant examples for $r_{(p,q)} \in C_n$: We consider points $x \in M_N$. A relevant example of $r_{(p,q)} \in C_n$ is then of the form

$$\langle \text{binary}(\text{powers}(MR(x, N), N), N), LSB(x) \rangle .$$

The length of this example in bits is at most $4l^2 + 6l + 2 \leq n$. The target distribution D^+ for $r_{(p,q)}$ is uniform over the relevant positive examples (i.e., those for which $LSB(x) = 1$) and the target distribution D^- is uniform over the relevant negative examples (i.e., those for which $LSB(x) = 0$).

Difficulty of weakly learning $C = \cup_{n \geq 1} C_n$: Suppose that A is a polynomial-time weak learning algorithm for C . We now describe how to use A to factor Blum integers. Let N be a Blum integer. Given only N as input, we run algorithm A . Every time A requests a positive example, we choose $x \in M_N$ uniformly such that $LSB(x) = 1$, and give the example

$$\langle \text{binary}(\text{powers}(MR(x, N), N), N), 1 \rangle$$

to A . Such an example can be generated in polynomial time on input N . This simulation generates the distribution D^+ . Every time A requests a negative example, we choose $x \in M_n$ uniformly such that $LSB(x) = 0$, and give the example

$$\langle \text{binary}(\text{powers}(MR(x, N), N), N), 0 \rangle$$

to A . Again, this example can be generated in polynomial time. This simulation generates the distribution D^- . When algorithm A has halted, $h_A(\text{binary}(\text{powers}(r, N), N)) = LSB(x)$

with probability $1/2 + 1/p(l)$ for $r = MR(x, N)$ and x chosen uniformly from M_N . This implies that we can factor Blum integers by the results of Rabin [35] and Alexi et al. [4] given in Theorems 2 and 3.

Ease of evaluating $r_{(p,q)} \in C_n$: For each $r_{(p,q)} \in C_n$, we give an NC^1 circuit for evaluating the concept represented by $r_{(p,q)}$ on an input of the form

$$binary(powers(r, N), N)$$

where $N = p \cdot q$ and $r = MR(x, N)$ for some $x \in M_N$. This is accomplished by giving an NC^1 implementation of the first three steps of the root-finding algorithm of Adleman, Manders and Miller [1] as it is described by Angluin [5]. Note that if we let $a = x^2 \bmod N$, then either $r = a$ or $r = (N - a) \bmod N$ according to the definition of the modified Rabin function. The circuit has four phases.

Phase I. Determine if the input r is a quadratic residue mod N . This can be done using the given powers of r and $r_{(p,q)}$ using the NC^1 circuit described in quadratic residue-based scheme of Section 5.2. Note that since p and q are both congruent to $3 \bmod 4$, $(N - a) \bmod N$ is never a quadratic residue mod N (see Angluin [5]). If it is decided that $r = (N - a) \bmod N$, generate the intermediate output $a \bmod N$. This can clearly be done in NC^1 . Also, notice that for any z , $z^{2^i} = (N - z)^{2^i} \bmod N$ for $i \geq 1$. Hence these powers of r are identical in the two cases. Finally, recall that the NC^1 circuit for quadratic residues produced the powers of $r \bmod p$ and the powers of $r \bmod q$ as intermediate outputs, so we may assume that the powers

$$a, a^2 \bmod p, a^4 \bmod p, \dots, a^{2^{2^l}} \bmod p$$

and

$$a, a^2 \bmod q, a^4 \bmod q, \dots, a^{2^{2^l}} \bmod q$$

are also available.

Phase II. Let l_p (respectively, l_q) be the largest positive integer such that $2^{l_p} | (p - 1)$ (respectively, $2^{l_q} | (q - 1)$). Let $Q_p = (p - 1)/2^{l_p}$ (respectively, $Q_q = (q - 1)/2^{l_q}$). Using the appropriate powers of $x^2 \bmod p$ and $\bmod q$, compute $u = a^{(Q_p+1)/2} \bmod p$ and $v =$

$a^{(Q_i+1)/2} \bmod q$ with NC^1 iterated product circuits. Since p and q are both congruent to $3 \bmod 4$, u and $p - u$ are square roots of $a \bmod q$, and v and $q - v$ are square roots of $a \bmod q$ by the results of Adleman et al. [1] (see also Angluin [5]).

Phase III. Using Chinese remaindering, combine $u, p - u, v$ and $q - v$ to compute the four square roots of $a \bmod N$ (see e.g. Kranakis [28]). Given p and q , this requires only a constant number of multiplication and addition steps, and so is computed in NC^1 .

Phase IV. Find the root from Phase III that is in M_N , and output its least significant bit.

6 Learning Small Boolean Formulae, Finite Automata and Threshold Circuits is Hard

The results of Section 5 show that for some fixed polynomial $q(n)$, learning NC^1 circuits of size at most $q(n)$ is computationally as difficult as the problems of inverting RSA, recognizing quadratic residues, and factoring Blum integers. However, there is a polynomial $p(n)$ such that any NC^1 circuit of size at most $q(n)$ can be represented by a Boolean formulae of size at most $p(n)$. Thus we have proved the following:

Theorem 4 *Let $\text{BF}_n^{p(n)}$ denote the class of Boolean formulae over n variables of size at most $p(n)$, and let $\text{BF}^{p(n)} = \cup_{n \geq 1} \text{BF}_n^{p(n)}$. Then for some polynomial $p(n)$, the problems of inverting the RSA encryption function, recognizing quadratic residues and factoring Blum integers are probabilistic polynomial-time reducible to weakly learning $\text{BF}^{p(n)}$.*

In fact, we can apply the substitution arguments of Kearns et al. [26] to show that Theorem 4 holds even for the class of monotone Boolean formulae in which each variable appears at most once.

Pitt and Warmuth [33] show that if the class ADFA is polynomially weakly learnable, then the class BF is polynomially weakly learnable. Combining this with Theorem 4, we have:

Theorem 5 *Let $\text{ADFA}_n^{p(n)}$ denote the class of deterministic finite automata of size at most $p(n)$ that only accept strings of length n , and let $\text{ADFA}^{p(n)} = \cup_{n \geq 1} \text{ADFA}_n^{p(n)}$. Then for some polynomial $p(n)$, the problems of inverting the RSA encryption function, recognizing quadratic residues and factoring Blum integers are probabilistic polynomial-time reducible to weakly learning $\text{ADFA}^{p(n)}$.*

Using results of Chandra, Stockmeyer and Vishkin [16], Beame et al. [10] and Reif [36], it can be shown that the representations described in Section 5 can each be computed by a polynomial-size, constant-depth threshold circuit. Thus we have:

Theorem 6 *For some fixed constant natural number d , let $d\text{TC}_n^{p(n)}$ denote the class of threshold circuits over n variables with depth at most d and size at most $p(n)$, and let $d\text{TC}^{p(n)} = \cup_{n \geq 1} d\text{TC}_n^{p(n)}$. Then for some polynomial $p(n)$, the problems of inverting the RSA encryption function, recognizing quadratic residues and factoring Blum integers are probabilistic polynomial-time reducible to weakly learning $d\text{TC}^{p(n)}$.*

It is important to reiterate that these hardness results hold regardless of the hypothesis representation class of the learning algorithm; that is, Boolean formulae, DFA's and constant-depth threshold circuits are not weakly learnable by *any* polynomially evaluable representation class (under standard cryptographic assumptions). We note that no NP -hardness results are known for these classes even if we restrict the hypothesis class to be the same as the target class and insist on strong learnability rather than weak learnability. It is also possible to give reductions showing that many other interesting classes (e.g., CFG's and NFA's) are not weakly learnable, under the same cryptographic assumptions. In general, any representation class whose computational power subsumes that of NC^1 is not weakly learnable; however, more subtle reductions are also possible. In particular, our results resolve a problem posed by Pitt and Warmuth [33] by showing that under cryptographic assumptions, the class of all languages accepted by logspace Turing machines is not weakly learnable.

Pitt and Warmuth [33] introduce a general notion of reduction between learning problems, and a number of learning problems are shown to have equivalent computational difficulty (with respect to probabilistic polynomial-time reducibility); thus, if the learning problem for a representation class C_1 reduces to the learning problem for a representation class C_2 , then a polynomial-time learning algorithm for C_2 in the distribution-free model implies a polynomial-time learning algorithm for C_1 . Learning problems are then classified according to the complexity of their *evaluation problem*, the problem of evaluating a representation on an input example. In Pitt and Warmuth [33] the evaluation problem is treated as a uniform problem (i.e., one algorithm for evaluating all representations in the class); by treating the evaluation problem nonuniformly (e.g., a separate circuit for

each representation) we were able to show that NC^1 contains a number of presumably hard-to-learn classes of Boolean functions. By giving reductions from NC^1 to other classes of representations, we thus clarify the boundary of what is efficiently learnable.

7 A Generalized Construction Based on Any Trapdoor Function

Let us now give a brief summary of the techniques that were used in Sections 5 and 6 to obtain hardness results for learning based on cryptographic assumptions. In each construction (RSA, quadratic residue and factoring Blum integers), we began with a candidate *trapdoor function* family, informally a family of functions each of whose members f is easy to compute (that is, given x , it is easy to compute $f(x)$), hard to invert (that is, given only $f(x)$, it is difficult to compute x), but easy to invert given a secret “key” to the function [41] (the *trapdoor*). We then constructed a learning problem in which the complexity of inverting the function *given* the trapdoor key corresponds to the complexity of the representations being learned, and learning from random examples corresponds to inverting the function *without* the trapdoor key. Thus, the learning algorithm is essentially required to learn the inverse of a trapdoor function, and the small representation for this inverse is simply the secret trapdoor information.

To prove hardness results for the simplest possible representation classes, we then eased the computation of the inverse given the trapdoor key by providing the powers of the original input in each example. This additional information provably does not compromise the security of the original function. A key property of trapdoor functions exploited by our constructions is the ability to generate random examples of the target representation without the trapdoor key; this corresponds to the ability to generate encrypted messages given only the public key in a public-key cryptosystem.

By assuming that specific functions such as RSA are trapdoor functions, we were able to find modified trapdoor functions whose inverse computation given the trapdoor could be performed by very simple circuits. This allowed us to prove hardness results for specific representation classes that are of interest in computational learning theory. Such specific intractability assumptions appear necessary since the weaker and more general assumption that there exists a trapdoor family that can be computed (in the forward direction) in polynomial time does not allow us to say anything about the hard-to-learn representation class other than it having polynomial-size circuits.

However, the summary above suggests a general method for proving hardness results for learning: to show that a representation class C is not learnable, find a trapdoor function whose inverse can be computed by C given the trapdoor key. In this section we formalize these ideas and prove a theorem demonstrating that this is indeed a viable approach.

We use the following definition for a family of trapdoor functions, which can be derived from Yao [41]: let $P = \{P_n\}$ be a family of probability distributions, where for $n \geq 1$ the distribution P_n is over pairs $(k, k') \in \{0, 1\}^n \times \{0, 1\}^n$. We think of k as the n -bit *public key* and k' as the associated n -bit *private key*. Let $Q = \{Q_k\}$ be a family of probability distributions parameterized by the public key k , where if $|k| = n$ then Q_k is a distribution over $\{0, 1\}^n$. We think of Q as a distribution family over the *message space*. The function $f : \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}^n$ maps an n -bit public key k and an n -bit *cleartext message* x to the *ciphertext* $f(k, x)$. We call the triple (P, Q, f) an $\alpha(n)$ -*strong trapdoor scheme* if it has the following properties:

- (i) There is probabilistic polynomial-time algorithm G (the *key generator*) that on input 1^n outputs a pair (k, k') according to the distribution P_n . Thus, pairs of public and private keys are easily generated.
- (ii) There is a probabilistic polynomial-time algorithm M (the *message generator*) that on input k outputs x according to the distribution Q_k . Thus, messages are easily generated given the public key k .
- (iii) There is a polynomial-time algorithm E that on input k and x outputs $f(k, x)$. Thus, encryption is easy.
- (iv) Let A be any probabilistic polynomial-time algorithm. Perform the following experiment: draw a pair (k, k') according to P_n , and draw x according to Q_k . Give the inputs k and $f(k, x)$ to A . Then the probability that $A(k, f(k, x)) \neq x$ is at least $\alpha(n)$. Thus, decryption from only the public key and the ciphertext is hard.
- (v) There is a polynomial-time algorithm D that on input k, k' and $f(k, x)$ outputs x . Thus, decryption given the private key (or *trapdoor*) is easy.

As an example, consider the RSA cryptosystem [37]. Here the distribution P_n is uniform over all (k, k') where $k' = (p, q)$ for n -bit primes p and q and $k = (p \cdot q, e)$ with $e \in Z_{\varphi(p \cdot q)}^*$. The distribution

Q_k is uniform over $Z_{p \cdot q}$, and $f(k, x) = f((p \cdot q, e), x) = x^e \bmod p \cdot q$.

We now formalize the notion of the inverse of a trapdoor function being computed in a representation class. Let $C = \cup_{n \geq 1} C_n$ be a parameterized Boolean representation class. We say that a trapdoor scheme (P, Q, f) is *invertible in C given the trapdoor* if for any $n \geq 1$, for any pair of keys $(k, k') \in \{0, 1\}^n \times \{0, 1\}^n$, and for any $1 \leq i \leq n$, there is a representation $c_{(k, k')}^i \in C_n$ that on input $f(k, x)$ (for any $x \in \{0, 1\}^n$) outputs the i th bit of x .

Theorem 7 *Let p be any polynomial, and let $\alpha(n) \geq 1/p(n)$. Let (P, Q, f) be an $\alpha(n)$ -strong trapdoor scheme, and let C be a parameterized Boolean representation class. Then if (P, Q, f) is invertible in C given the trapdoor, C is not polynomially learnable.*

Proof: Let A be any polynomial-time learning algorithm for C . We use algorithm A as a subroutine in a polynomial-time algorithm A' that with high probability outputs x on input k and $f(k, x)$, thus contradicting condition (iv) in the definition of a trapdoor scheme.

Let (k, k') be n -bit public and private keys generated by the distribution P_n . Let x be an n -bit message generated according to the distribution Q_k . Then on input k and $f(k, x)$, algorithm A' behaves as follows: for $1 \leq i \leq n$, algorithm A' simulates algorithm A , choosing accuracy parameter $\epsilon = \alpha(n)/n$. For the i th run of A , each time A requests a positive example, A' generates random values x' from the distribution Q_k (this can be done in polynomial time by condition (ii) in the definition of trapdoor scheme) and computes $f(k, x')$ (this can be done in polynomial time by condition (iii) in the definition of trapdoor scheme). If the i th bit of $f(k, x')$ is 1, then A' gives x' as a positive example to A ; similarly, A' generates negative examples for the i th run of A by drawing x' such that the i th bit of $f(k, x')$ is 0. If after $O(1/\epsilon \ln n/\delta)$ draws from Q_k , A' is unable to obtain a positive (respectively, negative) example for A , then A' assumes that with high probability a random x' results in the i th bit of $f(k, x')$ being 0 (respectively, 1), and terminates this run by setting h_k^i to the hypothesis that is always 0 (respectively, 1). The probability that A' terminates the run incorrectly can be shown to be smaller than δ/n by application of Fact CB1 and Fact CB2.

Note that all of the examples given to the i th run of A are consistent with a representation in C_n , since the i th bit of $f(k, \cdot)$ is computed by the representation $c_{(k, k')}^i$. Thus with high probability A outputs an ϵ -good hypothesis h_k^i . To invert the original input $f(k, x)$, A' simply outputs the bit sequence $h_k^1(f(k, x)) \cdots h_k^n(f(k, x))$. The probability that any bit of this string differs from the

corresponding bit of x is at most $n\epsilon < \alpha(n)$, contradicting the assumption that (P, Q, f) is an $\alpha(n)$ -strong trapdoor scheme. \square

8 Application: Hardness Results for Approximation Algorithms

In this section, we digress from learning briefly and apply the results of Section 6 to prove that under cryptographic assumptions, certain combinatorial optimization problems, including a natural generalization of graph coloring, cannot be efficiently approximated even in a very weak sense. These results show that for these problems, it is difficult to find a solution that approximates the optimal solution even within a factor that grows rapidly with the input size. Such results are infrequent in complexity theory, and seem difficult to obtain for natural problems using presumably weaker assumptions such as $P \neq NP$.

We begin by stating a needed theorem of Blumer et al. known as *Occam's Razor* [14]. Their result essentially gives an upper bound on the sample size required for learning C by H , and shows that the general technique of finding an hypothesis that is both consistent with the sample drawn and significantly shorter than this sample is sufficient for distribution-free learning. Thus, if one can efficiently perform *data compression* on a random sample, then one can learn efficiently.

Theorem 8 (*Blumer et al. [14]*) *Let C and H be polynomially evaluable parameterized Boolean representation classes. Fix $\alpha \geq 1$ and $0 \leq \beta < 1$, and let A be an algorithm that on input a labeled sample S of some $c \in C_n$, consisting of m positive examples of c drawn from D^+ and m negative examples of c drawn from D^- , outputs an hypothesis $h_A \in H_n$ that is consistent with S and satisfies $|h_A| \leq n^\alpha m^\beta$, where $|h_A|$ is the length of the representation h_A in bits. Then A is a learning algorithm for C by H ; the sample size required is*

$$m = O\left(\frac{1}{\epsilon} \log \frac{1}{\delta} + \left(\frac{n^\alpha}{\epsilon} \log \frac{n^\alpha}{\epsilon}\right)^{\frac{1}{1-\beta}}\right).$$

Let $|S| = mn$ denote the number of bits in the sample S . Note that if A instead outputs h_A satisfying $|h_A| \leq n^{\alpha'} |S|^\beta$ for some fixed $\alpha' \geq 1$ and $0 \leq \beta < 1$ then $|h_A| \leq n^{\alpha'} (mn)^\beta = n^{\alpha'+\beta} m^\beta$, so A satisfies the condition of Theorem 8 for $\alpha = \alpha' + \beta$. This formulation of Occam's Razor will be of particular use to us.

Let C and H be polynomially evaluable parameterized Boolean representation classes, and define the *Consistency Problem* $Con(C, H)$ as follows:

The Consistency Problem $Con(C, H)$:

Input: A labeled sample S of some $c \in C_n$.

Output: $h \in H_n$ such that h is consistent with S and $|h|$ is minimized.

We use $opt_{Con}(S)$ to denote the size of the smallest hypothesis in H that is consistent with the sample S , and $|S|$ to denote the number of bits in S . Using the results of Section 6 and Theorem 8, we immediately obtain proofs of the following theorems.

Theorem 9 *Let BF_n denote the class of Boolean formulae over n variables, and let $BF = \cup_{n \geq 1} BF_n$. Let H be any polynomially evaluable parameterized Boolean representation class. Then the problems of inverting the RSA encryption function, recognizing quadratic residues and factoring Blum integers are probabilistic polynomial-time reducible to the problem of approximating the optimal solution of an instance S of $Con(BF, H)$ by an hypothesis h satisfying*

$$|h| \leq (opt_{Con}(S))^\alpha |S|^\beta$$

for any $\alpha \geq 1$ and $0 \leq \beta < 1$.

Theorem 10 *Let $ADFA_n$ denote the class of deterministic finite automata accepting only strings of length n , and let $ADFA = \cup_{n \geq 1} ADFA_n$. Let H be any polynomially evaluable parameterized Boolean representation class. Then inverting the RSA encryption function, recognizing quadratic residues and factoring Blum integers are probabilistic polynomial-time reducible to approximating the optimal solution of an instance S of $Con(ADFA, H)$ by an hypothesis h satisfying*

$$|h| \leq (opt_{Con}(S))^\alpha |S|^\beta$$

for any $\alpha \geq 1$ and $0 \leq \beta < 1$.

Theorem 11 *Let dTC_n denote the class of threshold circuits over n variables with depth at most d , and let $dTC = \cup_{n \geq 1} dTC_n$. Let H be any polynomially evaluable parameterized Boolean representation class. Then for some constant $d \geq 1$, the problems of inverting the RSA encryption function, recognizing quadratic residues and factoring Blum integers are probabilistic polynomial-time reducible to the problem of approximating the optimal solution of an instance S of $Con(dTC, H)$ by an hypothesis h satisfying*

$$|h| \leq (opt_{Con}(S))^\alpha |S|^\beta$$

for any $\alpha \geq 1$ and $0 \leq \beta < 1$.

These theorems demonstrate that the results of Section 6 are in some sense not dependent upon the particular models of learnability that we study, since we are able to restate the hardness of learning in terms of standard combinatorial optimization problems. Using a generalization of Theorem 8 [15], we can in fact prove Theorems 9, 10 and 11 for the *Relaxed Consistency Problem*, where the hypothesis found must agree with only a fraction $1/2 + 1/p(\text{opt}_{\text{Con}}(S), n)$ for any fixed polynomial p . The central idea of the proof is the same: since the results of Blumer et al. [15] demonstrate that for sufficient sample size, solution of the relaxed consistency problem implies weak learning, and we have shown weak learning to be as hard as the cryptographic problems for the various representation classes, the relaxed consistency problem is as hard as the cryptographic problems. Using the results of Goldreich et al. [22], it is also possible to show similar hardness results for the Boolean circuit consistency problem $\text{Con}(\text{CKT}, \text{CKT})$ using the weaker assumption that there exists a one-way function.

It is interesting to contrast Theorem 10 with similar results obtained by Pitt and Warmuth [34]. They also prove hardness results for the problem of finding small deterministic finite automata consistent with a labeled sample, but based on the weaker assumption $\mathfrak{N} \neq NP$. However (using the notation of Theorem 10), their results only hold for a more restricted range of α and β , and require the restriction that H be the class of deterministic finite automata. We refer the reader to their paper for details.

Note that Theorem 11 addresses the optimization problem $\text{Con}(d\text{TC}, \text{TC})$ as a special case. This problem is essentially that of finding a set of weights in a neural network that yields the desired input-output behavior, sometimes referred to as the *loading problem*. Theorem 11 states that even if we allow a much larger net than is actually required, finding these weights is computationally intractable, even for only a constant number of “hidden layers”. This result should be contrasted with those of Judd [25] and Blum and Rivest [12], which rely on the weaker assumption $P \neq NP$ but do not prove hardness for relaxed consistency and do not allow the hypothesis network to be substantially larger than the smallest consistent network. We also make no assumptions on the topology of the output circuit.

Theorems 9, 10 and 11 are interesting for at least two reasons. First, they suggest that it is possible to obtain stronger hardness results for combinatorial optimization approximation algorithms

by using stronger complexity-theoretic assumptions. Such results seem difficult to obtain using only the assumption $P \neq NP$. Second, these results provide us with natural examples of optimization problems for which it is hard to approximate the optimal solution even within a multiplicative factor that grows as a function of the input size. Several well-studied problems apparently have this property, but little has been proven in this direction. Perhaps the best example is graph coloring, where the best polynomial-time algorithms require approximately $n^{1-1/(k-1)}$ colors on k -colorable n -vertex graphs (see Wigderson [40] and Blum [11]) but coloring has been proven NP -hard only for $(2-\epsilon)k$ colors for any $\epsilon > 0$ (see Garey and Johnson [20]). Thus for 3-colorable graphs we only know that 5-coloring is hard, but the best algorithm requires roughly $O(n^{0.4})$ colors on n -vertex graphs! This leads us to look for approximation-preserving reductions from our provably hard optimization problems to other natural problems.

We now define a class of optimization problems that we call *formula coloring* problems. Here we have variables y_1, \dots, y_m assuming natural number values, or *colors*. We regard an assignment of colors to the y_i (called a *coloring*) as a partition P of the variable set into equivalence classes; thus two variables have the same color if and only if they are in the same equivalence class. We consider Boolean formulae that are formed using the standard basis over atomic elements of the form $(y_i = y_j)$ and $(y_i \neq y_j)$, where the predicate $(y_i = y_j)$ is satisfied if and only if y_i and y_j are assigned the same color.

A *model* for such a formula $F(y_1, \dots, y_m)$ is a coloring of the variables y_1, \dots, y_m such that F is satisfied. A *minimum model* for the F is a model using the fewest colors. For example, the formula

$$(y_1 = y_2) \vee ((y_1 \neq y_2) \wedge (y_3 \neq y_4))$$

has as a model the two-color partition $\{y_1, y_3\}, \{y_2, y_4\}$ and has as a minimum model the one-color partition $\{y_1, y_2, y_3, y_4\}$.

We will be interested in the problem of finding minimum models for certain restricted classes of formulae. For $F(y_1, \dots, y_m)$ a formula as described above, and P a model of F , we let $|P|$ denote the number of colors in P and $opt_{FC}(F)$ the number of colors in a minimum model of F .

We first show how graph coloring can be exactly represented as a formula coloring problem. If G is a graph, then for each edge (v_i, v_j) in G , we conjunct the expression $(y_i \neq y_j)$ to the formula $F(G)$. Then $opt_{FC}(F(G))$ is exactly the number of colors required to color G . Similarly,

by conjuncting expressions of the form

$$((y_1 \neq y_2) \vee (y_1 \neq y_3) \vee (y_2 \neq y_3))$$

we can also exactly represent the *3-hypergraph coloring* problem (where each hyperedge contains 3 vertices) as a formula coloring problem.

To prove our hardness results, we consider a generalization of the graph coloring problem:

The Formula Coloring Problem *FC*:

Input: A formula $F(y_1, \dots, y_m)$ which is a conjunction only of expressions of the form $(y_i \neq y_j)$ (as in the graph coloring problem) *or* of the form $((y_i \neq y_j) \vee (y_k = y_l))$.

Output: A minimum model for F .

We will show that approximating an optimal solution to this problem is as hard as approximating the consistency problem $Con(DFA, DFA)$, where DFA is the class of deterministic finite automata. Note that this problem is at least as hard to approximate as $Con(ADFA, H)$, which we have already proven an approximation hardness result in Theorem 10.

Theorem 12 *There is a polynomial-time algorithm A that on input an instance S of the problem $Con(DFA, DFA)$ outputs an instance $F(S)$ of the formula coloring problem such that S has a k -state consistent hypothesis $M \in DFA$ if and only if $F(S)$ has a model of k colors.*

Proof: Let S contain the labeled examples

$$\langle w_1, b_1 \rangle, \langle w_2, b_2 \rangle, \dots, \langle w_m, b_m \rangle$$

where each $w_i \in \{0, 1\}^n$ and $b_i \in \{0, 1\}$. Let w_i^j denote the j th bit of w_i . We create a variable z_i^j for each $1 \leq i \leq n$ and $0 \leq j \leq m$. Let M be a smallest DFA consistent with S . Then we interpret z_i^j as representing the state that M is in immediately after reading the bit w_i^j on input w_i . The formula $F(S)$ will be over the z_i^j and is constructed as follows: for each i_1, i_2 and j_1, j_2 such that $0 \leq j_1, j_2 < n$ and $w_{i_1}^{j_1+1} = w_{i_2}^{j_2+1}$ we conjunct the predicate

$$((z_{i_1}^{j_1} = z_{i_2}^{j_2}) \rightarrow (z_{i_1}^{j_1+1} = z_{i_2}^{j_2+1}))$$

to $F(S)$. Note that this predicate is equivalent to

$$((z_{i_1}^{j_1} \neq z_{i_2}^{j_2}) \vee (z_{i_1}^{j_1+1} = z_{i_2}^{j_2+1}))$$

and thus has the required form. These formulae are designed to encode the constraint that if M is in the same state in two different computations on input strings from S , and the next input symbol is the same in both strings, then the next state in each computation must be the same.

For each i_1, i_2 ($1 \leq i_1, i_2 \leq m$) such that $b_{i_1} \neq b_{i_2}$ we conjunct the predicate $(z_{i_1}^n \neq z_{i_2}^n)$. These predicates are designed to encode the constraint that the input strings in S that are accepted by M must result in different final states than those strings in S that are rejected by M .

We first prove that if M has k states, then $\text{opt}_{FC}(F(S)) \leq k$. In particular, let P be the k -color partition that assigns $z_{i_1}^{j_1}$ and $z_{i_2}^{j_2}$ the same color if and only if M is in the same state after reading $w_{i_1}^{j_1}$ on input w_{i_1} and after reading $w_{i_2}^{j_2}$ on input w_{i_2} . We show that P is a model of $F(S)$. A conjunct

$$((z_{i_1}^{j_1} = z_{i_2}^{j_2}) \rightarrow (z_{i_1}^{j_1+1} = z_{i_2}^{j_2+1}))$$

of $F(S)$ cannot be violated by P since this conjunct appears only if $w_{i_1}^{j_1+1} = w_{i_2}^{j_2+1}$; thus if state $z_{i_1}^{j_1}$ is equivalent to state $z_{i_2}^{j_2}$ then state $z_{i_1}^{j_1+1}$ must be equivalent to state $z_{i_2}^{j_2+1}$ since M is deterministic. A conjunct

$$(z_{i_1}^n \neq z_{i_2}^n)$$

of $F(S)$ cannot be violated by P since this conjunct appears only if $b_{i_1} \neq b_{i_2}$, and if state $z_{i_1}^n$ is equivalent to state $z_{i_2}^n$ then w_{i_1} and w_{i_2} are either both accepted or both rejected by M , which contradicts M being consistent with S .

For the other direction, we show that if $\text{opt}_{FC}(F(S)) \leq k$ then there is a k -state DFA M' that is consistent with S . M' is constructed as follows: the k states of M' are labeled with the k equivalence classes (colors) X_1, \dots, X_k of the variables z_i^j in a minimum model P' for $F(S)$. There is a transition from state X_p to state X_q if and only if there are i, j such that $z_i^j \in X_p$ and $z_i^{j+1} \in X_q$; this transition is labeled with the symbol w_i^{j+1} . We label X_p an *accepting* (respectively, *rejecting*) state if for some variable $z_i^n \in X_p$ we have $b_i = 1$ (respectively, $b_i = 0$).

We first argue that no state X_p of M' can be labeled both an accepting and rejecting state. For if $b_i = 1$ and $b_j = 0$ then the conjunct $(z_i^n \neq z_j^n)$ appears in $F(S)$, hence z_i^n and z_j^n must have different colors in P' .

Next we show that M is in fact deterministic. For suppose that some state X_p has transitions to X_q and X_r , and that both transitions are labeled with the same symbol. Then there exist i_1, i_2 and j_1, j_2 such that $z_{i_1}^{j_1} \in X_p$ and $z_{i_1}^{j_1+1} \in X_q$, and $z_{i_2}^{j_2} \in X_p$ and $z_{i_2}^{j_2+1} \in X_r$. Furthermore we must

have $w_{i_1}^{j_1+1} = w_{i_2}^{j_2+1}$ since both transitions have the same label. But then the conjunct

$$((z_{i_1}^{j_1} = z_{i_2}^{j_2}) \rightarrow (z_{i_1}^{j_1+1} = z_{i_2}^{j_2+1}))$$

must appear in $F(S)$, and this conjunct is violated P' , a contradiction. Thus M' is deterministic.

These arguments prove that M' is a well-defined DFA. To see that M' is consistent with S , consider the computation of M' on any w_i in S . The sequence of states visited on this computation is just $EC_{P'}(z_i^1), \dots, EC_{P'}(z_i^n)$, where $EC_{P'}(z_i^j)$ denotes the equivalence class of the variable z_i^j in the coloring P' . The final state $EC_{P'}(z_i^n)$ is by definition of M' either an accept state or a reject state according to whether $b_i = 1$ or $b_i = 0$. \square

Note that if $|S|$ is the number of bits in the sample S and $|F(S)|$ denotes the number of bits in the formula $F(S)$, then in Theorem 12 we have $|F(S)| = \Theta(|S|^2 \log |S|) = O(|S|^{2+\gamma})$ for any $\gamma > 0$ for $|S|$ sufficiently large. This means that if an algorithm colors $F(S)$ using at most $\text{opt}_{FC}(F(S))^\alpha |F(S)|^\beta$ for some $\alpha \geq 1$ and $\beta < \frac{1}{2}$, then for $|S|$ sufficiently large we can use the reduction of Theorem 12 to find a DFA consistent with S that has at most $k^\alpha |S|^{\beta'}$ for some $\beta' < 1$, contradicting Theorem 10. Thus we have:

Theorem 13 *The problems of inverting the RSA encryption function, recognizing quadratic residues and factoring Blum integers are polynomial-time reducible to approximating the optimal solution to an instance F of the formula coloring problem by a model P of F satisfying*

$$|P| \leq \text{opt}_{FC}(F)^\alpha |F|^\beta$$

for any $\alpha \geq 1$ and $0 \leq \beta < 1/2$.

Figure 7.1 summarizes hardness results for coloring a formula F using at most $f(\text{opt}_{FC}(F))g(|F|)$ colors for various functions f and g , where an entry “NP-hard” indicates that such an approximation is NP-hard, “Factoring” indicates that such an approximation is as hard as factoring Blum integers (or recognizing quadratic residues or inverting the RSA function), and “P” indicates there is a polynomial-time algorithm achieving this approximation factor. The NP-hardness results follow from Garey and Johnson [20] and Pitt and Warmuth [34].

9 Open Problems

A technical open problem is to improve the constructions given here to prove representation-independent hardness results for even simpler classes of formulae and circuits. It would also be interesting to demonstrate a partial converse to our results: for instance, if we assume there is a representation class that is hard to learn, can it be used to construct any interesting cryptographic primitives?

We encourage the reader to see the paper of Angluin and Kharitonov [7], where the methods here are extended to prove hardness results for learning with queries.

References

- [1] L. Adleman, K. Manders, G. Miller.
On taking roots in finite fields.
Proceedings of the 18th I.E.E.E. Symposium on Foundations of Computer Science, 1977, pp. 175-178.
- [2] A. Aho, J. Hopcroft, J. Ullman.
The design and analysis of computer algorithms.
Addison-Wesley, 1974.
- [3] D. Aldous.
On the Markov chain simulation method for uniform combinatorial distributions and simulated annealing.
University of California at Berkeley Statistics Department, technical report number 60, 1986.
- [4] W. Alexi, B. Chor, O. Goldreich, C.P. Schnorr.
RSA and Rabin functions: certain parts are as hard as the whole.
S.I.A.M. Journal on Computing, 17(2), 1988, pp. 194-209.
- [5] D. Angluin.
Lecture notes on the complexity of some problems in number theory.
Yale University Computer Science Department, technical report number TR-243, 1982.

- [6] D. Angluin.
Learning regular sets from queries and counterexamples.
Information and Computation, 75, 1987, pp. 87-106.
- [7] D. Angluin, M. Kharitonov.
When won't membership queries help?
Proceedings of the 23rd A.C.M. Symposium on the Theory of Computing, 1991, pp. 444-454.
- [8] D. Angluin, P. Laird.
Learning from noisy examples.
Machine Learning, 2, 1988, pp. 319-342.
- [9] D. Angluin, L.G. Valiant.
Fast probabilistic algorithms for Hamiltonian circuits and matchings.
Journal of Computer and Systems Sciences, 18, 1979, pp. 155-193.
- [10] P.W. Beame, S.A. Cook, H.J. Hoover.
Log depth circuits for division and related problems.
S.I.A.M. Journal on Computing, 15(4), 1986, pp. 994-1003.
- [11] A. Blum.
An $\tilde{O}(n^{0.4})$ -approximation algorithm for 3-coloring.
Proceedings of the 21st A.C.M. Symposium on the Theory of Computing, 1989, pp. 535-542.
- [12] A. Blum, R.L. Rivest.
Training a 3-node neural network is NP-complete.
Proceedings of the 1988 Workshop on Computational Learning Theory, Morgan Kaufmann Publishers, 1988, pp. 9-18.
- [13] M. Blum, S. Micali.
How to generate cryptographically strong sequences of pseudo-random bits.
S.I.A.M. Journal on Computing, 13(4), 1984, pp. 850-864.

- [14] A. Blumer, A. Ehrenfeucht, D. Haussler, M. Warmuth.
Occam's razor.
Information Processing Letters, 24, 1987, pp. 377-380.
- [15] A. Blumer, A. Ehrenfeucht, D. Haussler, M. Warmuth.
Learnability and the Vapnik-Chervonenkis dimension.
Journal of the A.C.M., 36(4), 1989, pp. 929-965.
- [16] A.K. Chandra, L.J. Stockmeyer, U. Vishkin.
Constant depth reducibility.
S.I.A.M. Journal on Computing, 13(2), 1984, pp. 423-432.
- [17] H. Chernoff.
A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations.
Annals of Mathematical Statistics, 23, 1952, pp. 493-509.
- [18] W. Diffie, M. Hellman.
New directions in cryptography.
I.E.E.E. Transactions on Information Theory, 22, 1976, pp. 644-654.
- [19] A. Ehrenfeucht, D. Haussler, M. Kearns. L.G. Valiant.
A general lower bound on the number of examples needed for learning.
Information and Computation, 82(3), 1989, pp. 247-261.
- [20] M. Garey, D. Johnson.
Computers and intractability: a guide to the theory of NP-completeness.
Freeman, 1979.
- [21] E.M. Gold.
Complexity of automaton identification from given data.
Information and Control, 37, 1978, pp. 302-320.
- [22] O. Goldreich, S. Goldwasser, S. Micali.
How to construct random functions.
Journal of the A.C.M., 33(4), 1986, pp. 792-807.

- [23] T. Hancock.
On the difficulty of finding small consistent decision trees.
Harvard University, unpublished manuscript, 1989.
- [24] D. Haussler, M. Kearns, N. Littlestone, M. Warmuth.
Equivalence of models for polynomial learnability.
Proceedings of the 1988 Workshop on Computational Learning Theory, Morgan Kaufmann Publishers, 1988, pp. 42-55, and University of California at Santa Cruz Information Sciences Department, technical report number UCSC-CRL-88-06, 1988.
- [25] S. Judd.
Learning in neural networks.
Proceedings of the 1988 Workshop on Computational Learning Theory, Morgan Kaufmann Publishers, 1988, pp. 2-8.
- [26] M. Kearns, M. Li, L. Pitt, L.G. Valiant.
On the learnability of Boolean formulae.
Proceedings of the 19th A.C.M. Symposium on the Theory of Computing, 1987, pp. 285-295.
- [27] M. Kearns, L. Pitt.
A polynomial-time algorithm for learning k -variable pattern languages from examples.
Proceedings of the 1989 Workshop on Computational Learning Theory, Morgan Kaufmann Publishers, 1989, pp. 57-71.
- [28] E. Kranakis.
Primality and cryptography.
John Wiley and Sons, 1986.
- [29] L. Levin.
One-way functions and pseudorandom generators.
Proceedings of the 17th A.C.M. Symposium on the Theory of Computing, 1985, pp. 363-365.
- [30] M. Li, U. Vazirani.
On the learnability of finite automata.

- Proceedings of the 1988 Workshop on Computational Learning Theory*, Morgan Kaufmann Publishers, 1988, pp. 359-370.
- [31] N. Linial, Y. Mansour, N. Nisan.
Constant depth circuits, Fourier transform and learnability.
Proceedings of the 30th I.E.E.E. Symposium on the Foundations of Computer Science, 1989, pp. 574-579.
- [32] L. Pitt, L.G. Valiant.
Computational limitations on learning from examples.
Journal of the A.C.M., 35(4), 1988, pp. 965-984.
- [33] L. Pitt, M.K. Warmuth.
Reductions among prediction problems: on the difficulty of predicting automata.
Proceedings of the 3rd I.E.E.E. Conference on Structure in Complexity Theory, 1988, pp. 60-69.
- [34] L. Pitt, M.K. Warmuth.
The minimum consistent DFA problem cannot be approximated within any polynomial.
Proceedings of the 21st A.C.M. Symposium on the Theory of Computing, 1989, pp. 421-432.
- [35] M.O. Rabin.
Digital signatures and public key functions as intractable as factoring.
M.I.T. Laboratory for Computer Science, technical report number TM-212, 1979.
- [36] J. Reif.
On threshold circuits and polynomial computations.
Proceedings of the 2nd Structure in Complexity Theory Conference, 1987, pp. 118-125.
- [37] R. Rivest, A. Shamir, L. Adleman.
A method for obtaining digital signatures and public key cryptosystems.
Communications of the A.C.M., 21(2), 1978, pp. 120-126.
- [38] R. Schapire.
On the strength of weak learnability.

Proceedings of the 30th I.E.E.E. Symposium on the Foundations of Computer Science, 1989,
pp. 28-33.

[39] L.G. Valiant.

A theory of the learnable.

Communications of the A.C.M., 27(11), 1984, pp. 1134-1142.

[40] A. Wigderson.

A new approximate graph coloring algorithm.

Proceedings of the 14th A.C.M. Symposium on the Theory of Computing, 1982, pp. 325-329.

[41] A.C. Yao.

Theory and application of trapdoor functions.

Proceedings of the 23rd I.E.E.E. Symposium on the Foundations of Computer Science, 1982,
pp. 80-91.

Difficulty of coloring F using $A \cdot B$ colors	$A = 1$	$A = F ^{1/29}$	$A = F ^{0.499\dots}$	$A = F $
$B = \text{opt}_{FC}(F)$	<i>NP</i> -hard	<i>NP</i> -hard	Factoring	<i>P</i>
$B = 1.99 \dots \text{opt}_{FC}(F)$	<i>NP</i> -hard	Factoring	Factoring	<i>P</i>
$B = (\text{opt}_{FC}(F))^\alpha$ any fixed $\alpha \geq 0$	<i>NP</i> -hard	Factoring	Factoring	<i>P</i>

Figure 1: Difficulty of approximating the formula coloring problem using at most $A \cdot B$ colors on input formula F . The constant $0.499\dots$ is intended to indicate any value strictly smaller than $\frac{1}{2}$; the constant $\frac{1}{29}$ is determined from the paper of Pitt and Warmuth.