

# NONUNIFORM LEARNABILITY \*

Gyora M. Benedek

ELBIT Ltd.

Haifa, Israel

Alon Itai<sup>†</sup>

Computer Science Department

Technion, Haifa, Israel

February 12, 1995

## Abstract

The learning model of Valiant is extended to allow the number of examples required for learning to depend on the particular concept to be learned, instead of requiring a uniform bound for all concepts of a concept class.

This extension, called *nonuniform learning*, enables learning many concept classes not learnable by the previous definitions. Nonuniformly learnable concept classes are characterized. Some examples (Boolean formulae, recursive and r.e. sets) are shown to be nonuniformly learnable by a polynomial (in the size of the representation of the concept and in the error parameters) number of examples, but not necessarily in polynomial time. Restricting the learning protocol such that the learner has to commit himself after a finite number of examples does not affect the concept classes which can be learned.

An extension of nonuniform learnability to nonuniform learnability with respect to specific distributions is presented.

---

\*A Preliminary version appeared in ICALP 1988

<sup>†</sup>This research was supported by the Fund for the Promotion of Research at the Technion.

# 1 Introduction

In his seminal paper [21], Valiant introduced the PAC model<sup>1</sup>, according to which the teacher selects a concept (Boolean formula) from a commonly known concept class and provides the student with examples selected at random by some distribution  $D$  unknown to the student. An example is one element of the domain and a label specifying whether it is contained in the concept. A concept class is learnable if there is an algorithm that after receiving sufficiently many examples finds, with high probability, an approximation to the concept. The success of the algorithm should be independent of the distribution used by the teacher. The learning algorithm is *polynomial* if the sample size and computation time are polynomial in the size of the representation of the concept and the error parameters.

Blumer et al. [9] generalized the notion of learnability to arbitrary domains and concept classes. However, they discuss only the number of examples, not the computation time. This number is allowed to depend on the concept class and the error parameters but is required to be independent of the concept to be learned. Necessary and sufficient conditions for learnability using the Vapnik-Chervonenkis dimension [20] are given. This dimension depends only on the structure of the concept class (and is independent of the distribution). They show that a concept class is learnable if and only if it has finite dimension.

In the next section we show intuitively learnable concept classes which are not learnable according to this definition. To broaden the definition of learnability we consider nonuniform learnability, i.e., we allow the number of examples to depend not only on the concept class and error parameters but also on the concept to be learned. We give a necessary and sufficient condition for this case: a concept class is nonuniformly learnable if and only if it is a countable union of subclasses each of finite dimension.

This definition has the weakness that, even though by subsequently increasing the number of examples a “good” approximation of the concept can be found, there is no way to determine when to stop this process. A priori, it seems that to demand that the student “know” when he has succeeded in learning strengthens the definition. We show that both definitions are equivalent.

These general principles yield some interesting results about the learnability of recursive, r.e., context free languages and general Boolean formulae. In particular we give a partial solution to an open problem first raised by Valiant [21], i.e., “is the set of Boolean formulae polynomially learnable?”. We show that Boolean formulae are learnable by polynomially many examples. However, the learning time of our algorithm is not polynomial.

The papers [21, 22, 23, 17, 18, 15] deal with nonuniform learnability of Boolean formulae only. Their definition of polynomial learnability requires both the number of examples and the running time of the learning algorithm to be polynomial in the error parameters and in the length of the Boolean formula to be learned.

Blumer et al. [10] investigate learning functions from a countable set into a finite domain. They define a class of algorithms *Occam Algorithms* as algorithms that, when receiving a sample of a function  $f$ , produce a hypothesis subject to certain constraints. The complexity of the algorithm is polynomial in the number of examples and the complexity of  $f$ . The main result of [10] is that the existence of an Occam algorithm implies polynomial learnability.

---

<sup>1</sup>The exact definition appears in the next section.

We discuss learnability of recursive functions whose range is infinite. Furthermore, by our results, some of the constraints on the Occam algorithm may be relaxed (e.g., we allow noncountable domains). None of the aforementioned papers refer to the question whether the student “knows” when to stop.

In a recent paper [7] we have defined uniform learnability for an arbitrary distribution  $D$  and have given necessary and sufficient conditions for a concept class  $C$  to be uniformly learnable with respect to  $D$ . Here we broaden this notion by considering nonuniform learnability with respect to  $D$ . We give necessary and sufficient conditions for this case also.

After the preliminary version of this paper appeared [8], Lineal et al. [16] investigated similar notions. They independently proved the “sufficient” direction of Theorem 2, while we also show that the condition is necessary.

In [6] Ben-David et al. parameterize learnability models according to sampling complexity issues. Among other problems they consider the cases where the number of examples may depend on the target concept or on the distribution or on both. They present a classification for families of learnable concept classes using results from the current paper and additional analysis.

## 2 Definitions

Following [9], let  $X$  be a set,  $R \subseteq 2^X$  a  $\sigma$ -algebra over  $X$  [12], and  $D$  a probability measure over  $R$ . A *concept class* is a set  $C \subseteq R$  of *concepts*. For  $\mathbf{x} = (x_1, \dots, x_\ell) \in X^\ell$ , the *labeled  $\ell$ -sample of  $c \in C$*  is given by  $\text{sam}_c(\mathbf{x}) = (\langle x_1, c(x_1) \rangle, \dots, \langle x_\ell, c(x_\ell) \rangle)$ , where  $c(x)$  equals 1 if  $x \in c$  and 0 otherwise. The *sample space of  $C$* , denoted  $S_C$ , is the set of all labeled  $\ell$ -samples of  $c$  over all  $c \in C$ ,  $\mathbf{x} \in X^\ell$  ( $\ell \geq 1$ ). For a concept class  $C$  on  $X$ ,  $F_C$  denotes the set of all functions  $f : S_C \rightarrow C$ .

We use the following protocol ([21] and others) between two agents, T (teacher) and L (learner): T (who wants to teach L the concept  $c$ , called the *target concept*) repeatedly picks, at random according to some distribution  $D$ , an element  $x$  from a set  $X$  and sends the pair  $\langle x, c(x) \rangle$  to L. L, after receiving sufficiently many examples, returns a concept. We may view L as a function  $f \in F_C$ , and the set L returns is  $f(\langle x_1, c(x_1) \rangle, \dots, \langle x_\ell, c(x_\ell) \rangle)$ .

Let  $Y_1, Y_2 \in R$ .  $Y_1$  and  $Y_2$  are  $\varepsilon$ -close with respect to the distribution  $D$  if  $D(Y_1 \oplus Y_2) < \varepsilon$  ( $\oplus$  denotes the symmetric difference). Otherwise,  $Y_1$  and  $Y_2$  are  $\varepsilon$ -far with respect to the distribution  $D$ .

Let  $\mathbf{x} = (x_1, \dots, x_\ell)$  be a sequence of  $\ell$  independent randomly selected elements of  $X$ . For  $f \in F_C$  and  $c \in C$  we define the following probabilities  $r_f^\dagger(D, c, \ell, \varepsilon) = \text{Prob}(f(\text{sam}_c(\mathbf{x})) \text{ is a subset of } X \text{ } \varepsilon\text{-close to } c \text{ with respect to } D)$ . This is the probability that  $f$  finds a “good” approximation for  $c$  using an  $\ell$ -sample of  $c$ . We now follow [9] and define PAC (Probably Approximately Correct) learnability:

**Uniform learnability:**[9] A function  $f \in F_C$  *uniformly learns a concept class  $C$*  if for every  $\varepsilon, \delta > 0$  there is an  $\ell = \ell(\varepsilon, \delta) > 0$  such that for every distribution  $D$  over  $R$  and every  $c \in C$ ,  $r_f^\dagger(D, c, \ell, \varepsilon) > 1 - \delta$ .  $C$  is *uniformly learnable* if there exists an  $f \in F_C$  that uniformly learns  $C$ .

The definition above is quite stringent:

**Example 1:** Let  $X$  be the open segment  $(0, 1)$  and for every  $n$  let  $C_n$  be the set of all unions of  $n$  open segments over  $X$ .  $C_n$  is uniformly learnable see [9]. Let  $C = \bigcup_{n=1}^{\infty} C_n$  then  $C$  is not uniformly learnable [9] (for further details see Section 4).

On the other hand, intuitively, since every  $c \in C$  belongs to some  $C_i$  and since  $C_i$  is uniformly learnable [9],  $c$  is also learnable by a finite number of examples. Note that here the number of examples depends on the concept. Thus, it is reasonable to expect that, the more complex the concept  $c \in C$ , the more difficult it is to learn. This leads to the definition which allows the number of examples to depend on the concept to be learned.

**Nonuniform learnability:** A function  $f \in F_C$  *nonuniformly learns* a concept class  $C \subseteq R$  if for every  $\varepsilon, \delta > 0$  and every  $c \in C$  there is an  $\ell_0 = \ell(\varepsilon, \delta, c) > 0$  such that for every  $D$  over  $R$  and  $\ell \geq \ell_0$ ,  $r_f^+(D, c, \ell, \varepsilon) > 1 - \delta$ . Note that here  $\ell_0$  depends also on  $c$ .  $C$  is *nonuniformly learnable* if there exists an  $f \in F_C$  that nonuniformly learns  $C$ . We say that  $\ell(\varepsilon, \delta, c)$  is the *sample complexity* of  $f$ .

The intuition behind this definition is that for a fixed target concept  $c \in C$ , with high probability, the hypotheses output by  $f$  get closer to  $c$  as the number of examples increases.

In Section 4 we show that the concept class of example 1, is nonuniformly learnable.

**Example 2:** Let  $X = (0, 1)$  and  $C_{\text{OPEN}}$  the set of all open sets over  $X$ . In Section 4 we show that  $C_{\text{OPEN}}$  is not nonuniformly learnable.

### 3 Previous results – Uniform learnability.

In this section we quote the result presented in [9] (see also [20]), which will be used in the next section. Let  $T$  be a subset of  $X$ , a concept class  $C \subseteq 2^X$  *shatters*  $T$  if for every subset  $T'$  of  $T$  there is a concept  $c \in C$  such that  $T \cap c = T'$ . Also,  $C$  has Vapnik-Chervonenkis dimension  $d$  (VC-dimension) ( $\dim(C) = d$ ) if there is set of  $d$  elements of  $X$  shattered by  $C$  and there is no set of  $d + 1$  elements shattered by  $C$ . If there is no such  $d$  then  $C$  has infinite VC-dimension.

**Example 3:** Let  $X = (0, 1)$  and  $C$  be the set of all open intervals over  $X$ . It is easy to see that  $C$  shatters sets of two points but there is no set of three points shattered by  $C$ , thus  $\dim(C) = 2$ . The concept class  $C$  defined in example 1 shatters any finite set thus has infinite VC-dimension.

The main result of [9] is that  $C$  is learnable if and only if  $\dim(C)$  is finite.

A slightly stronger result appeared in [11]:

**Theorem 1:** [9, 11] *If  $\dim(C) = d \geq 2$  then:*

1. *For every  $\varepsilon, \delta > 0$  any function that receives  $\max\{(4/\varepsilon) \ln(2/\delta), (8d/\varepsilon) \ln(13/\varepsilon)\}$  examples and returns a concept in  $C$  consistent with the examples uniformly learns  $C$ .*
2. *For every  $0 < \varepsilon \leq \frac{1}{8}$  and  $0 < \delta \leq \frac{1}{100}$  there exists no function that can uniformly learn  $C$  using  $\Omega((1/\varepsilon) \ln(1/\delta) + d/\varepsilon)$  examples.*

In order for part (1) of the theorem to hold, the pair  $(C, D)$ , of concept class and distribution, has to be *well behaved*. The precise definition and proofs appear in [5, 9].

## 4 A necessary and sufficient condition for nonuniform learnability

**Theorem 2:**  $C$  is nonuniformly learnable if and only if there is an infinite sequence  $C_1, C_2, \dots$  such that

1.  $C = \bigcup_{i=1}^{\infty} C_i$ ,
2.  $\dim(C_i) < \infty$  for every  $i = 1, 2, \dots$

**Proof:** If  $C$  is nonuniformly learnable then there is an  $f \in F_C$ , such that every  $c \in C$  is learnable by  $f$  with sample complexity  $\ell(\varepsilon, \delta, c)$ . Let  $\varepsilon = \delta = \frac{1}{100}$  then for  $i = 5, 6, 7, \dots$  let  $C_i$  be the set of concepts  $c$  such that  $\ell(1/100, 1/100, c) \leq i$ , i.e.,  $f$  can learn every  $c \in C_i$  from  $i$  examples with accuracy and confidence  $1/100$ . Let  $k$  be the implied constant in the lower bound of Theorem 1. Then  $i \geq k(\frac{1}{100} \ln \frac{1}{100} + \frac{\dim(C_i)}{100})$ . Therefore,  $\dim(C_i) \leq \frac{100i}{k} + \ln 100$ . Since  $C$  is nonuniformly learnable every  $c \in C$  belongs to some  $C_i$ , i.e.  $C = \bigcup_{i=1}^{\infty} C_i$ .

For the other direction, let  $C = \bigcup_{i=1}^{\infty} C_i$  and  $\dim(C_i) < \infty$ . Since the property of having finite VC-dimension is closed under union, we may assume that  $C_i \supseteq \bigcup_{j=1}^i C_j$ , and therefore  $\dim(C_i)$  is a nondecreasing series. Let  $K(\varepsilon, \delta, i) = \max\{(4/\varepsilon) \ln(2/\delta), (8\dim(C_i)/\varepsilon) \ln(13/\varepsilon)\}$ . From its definition  $K$  is nondecreasing in  $\frac{1}{\varepsilon}$ ,  $\frac{1}{\delta}$  and  $i$ .

First we describe a function  $f$  that receives  $\varepsilon$ ,  $\delta$  and  $\ell$  examples and returns a hypothesis  $h \in C$ :

1. Let  $i$  be the largest integer  $i$  such that  $\ell \geq K(\varepsilon, \delta, i)$ .
2. Return a concept  $h \in C_i$  consistent with the examples. If no such  $h$  exists then return a prespecified  $c_0 \in C$ .

Now we prove that  $f$  learns  $C$ . Let  $c$  be some concept in  $C$ . Then there exists an  $i_0$  such that for all  $i \geq i_0$ ,  $c \in C_i$ . Given  $\ell$ ,  $f$  uses 1. to compute  $i$  and returns a concept from  $C_i$  consistent with the sample. By part (1) of Theorem 1, if  $i \geq i_0$ ,  $f$  learns  $c$ .

We modify  $f$  so that it will not need to receive  $\varepsilon$  and  $\delta$  as parameters:

- 1'. Let  $i$  be the largest integer  $i$  such that  $\ell \geq K(i^{-1}, i^{-1}, i)$ .

Thus when  $\ell$  increases,  $i$  also increases and  $\varepsilon$  and  $\delta$  monotonically decrease to zero. For every  $\varepsilon, \delta > 0$  and  $c \in C$  there exists an integer  $i$  such that  $\varepsilon \geq i^{-1}$ ,  $\delta \geq i^{-1}$  and  $c \in C_i$ . It follows that,  $\ell = K(i^{-1}, i^{-1}, i)$  examples are sufficient to learn  $c$ .  $\square$

The equivalence of the models where  $f$  does not receive  $\varepsilon$  and  $\delta$  as parameters (functional model) and the case where it does was discussed also in [13].

**Corollary 1:** Any countable class of concepts is nonuniformly learnable.

**Example 4:** Consider  $X$  and  $C_n$  of example 1. We have  $\dim(C_n) = 2n$  and  $C = \bigcup_{i=1}^{\infty} C_i$ , thus  $C$  is nonuniformly learnable.

**Claim 1:** The concept class  $C_{\text{OPEN}}$  of example 2 is not nonuniformly learnable.

**Proof:** First notice that the set  $\{1/n : n \in \mathbf{N}\}$  is shattered by  $C_{\text{OPEN}}$ . The following lemma implies the claim.

**Lemma 1:** *If a concept class  $C$  shatters an infinite set then there is no sequence  $C_1, C_2, \dots$  such that*

1.  $C = \bigcup_{i=1}^{\infty} C_i$
2.  $\dim(C_i) < \infty$  for every  $i = 1, 2, \dots$

**Proof:** (Shai Ben-David) Let  $C = \bigcup_{i=1}^{\infty} C_i$  such that  $d_i = \dim(C_i) < \infty$  for every  $i = 1, 2, \dots$ . We show that every infinite set  $T \subseteq X$  has a subset  $B$  such that for every  $c \in C$   $c \cap T \neq B$  and thus  $C$  does not shatter  $T$ . For every  $n$  let  $A_n$  be a set of  $d_n + 1$  elements of  $T - \bigcup_{i=1}^{n-1} A_i$ . Since  $\dim(C_n) = d_n$ , there exists a set  $B_n \subseteq A_n$ , such that for every  $c \in C_n$   $c \cap A_n \neq B_n$ .  $B$  is defined as  $\bigcup_{i=1}^{\infty} B_i$ . From the definitions  $B_n = B \cap A_n$ . Let  $c \in C$ , we now show that  $c \cap T \neq B$ . Since  $C = \bigcup_{i=1}^{\infty} C_i$  let  $c \in C_n$  and consider the two sets  $A_n \cap (c \cap T)$  and  $A_n \cap B$ : Since  $A_n \subseteq T$ ,  $A_n \cap (c \cap T) = A_n \cap c$ . If  $c \cap T = B$  then  $A_n \cap (c \cap T) = B \cap A_n$  implying that  $A_n \cap c = B_n$  – a contradiction.  $\square$

Shelah [19] showed that the converse of Lemma 1 does not hold.

The definition of nonuniform learnability has a considerable practical disadvantage. A user does not know how many examples to give  $f$ . She might be tempted to search for better and better hypotheses (ask for more examples) even though she has already reached a sufficiently close approximation to  $c$ . In other words, an algorithm that uses such a function to learn a concept has no indication when to stop. Thus our definition of (nonuniform) learning, resembles learning in the limit of Inductive Inference [3]. To overcome this difficulty we give the following definition:

**Definition:** A concept class  $C \subseteq R$  is *nonuniformly strongly learnable* if there exists a function  $f \in F_C$  and an algorithm  $A_f(\varepsilon, \delta)$  which has access to a source of random examples and to  $f$ , such that for all  $\varepsilon, \delta > 0$  and every  $c \in C$  there is an  $\ell > 0$  such that for every distribution  $D$  over  $R$ , with probability at least  $1 - \delta$  the algorithm requires at most  $\ell$  many examples and returns a hypothesis  $\varepsilon$ -close to  $c$  (with respect to  $D$ ).

**Theorem 3:** *A concept class  $C$  is nonuniformly strongly learnable if and only if it is nonuniformly learnable.*

**Proof:** Suppose  $C$  is nonuniformly strongly learnable by  $A_f$ . Construct a function  $g \in F_C$  as follows: Let the examples be  $(\langle x_1, c(x_1) \rangle, \dots, \langle x_\ell, c(x_\ell) \rangle)$ . For  $i = 1, 2, \dots$ , evaluate  $A_f(i^{-1}, i^{-1})$  while providing the examples  $\langle x_1, c(x_1) \rangle, \dots$  to get a hypothesis  $h_i$ . Stop when  $A_f$  asks for more than  $\ell$  examples, and return the last hypothesis produced by  $A_f$ . (If  $A_f$  has already asked for more than  $\ell$  examples then return some prespecified concept  $c_0$ .) It is easy to see that  $g$  learns  $C$  nonuniformly.

For the other direction, we use the “generate and test” paradigm as in Angluin et al. [1, 2].

Let  $C$  be a nonuniformly learnable concept class, and  $f$  a learning function for  $C$ . Thus, by Theorem 2,  $C = \bigcup_{i=1}^{\infty} C_i$  where  $\dim(C_i) < \infty$  for every  $i$ . As before, we may assume that  $C_i = \bigcup_{j=1}^i C_j$ , thus  $\dim(C_i)$  is nondecreasing. Let  $L_0 = \left\lceil \frac{32}{\varepsilon} \left( \ln \frac{3}{\delta} - \ln(1 - e^{-\varepsilon/32}) \right) \right\rceil$ . The algorithm consists of iterations; the  $i$ th iteration is:

1. Let  $\mathbf{x}_i$  consist of  $2^i$  examples.
2. Let  $h_i = f(\mathbf{x}_i)$  be the hypothesis of this iteration.
3. Let  $\mathbf{y}_i$  consist of  $L_i = L_0 + i$  additional examples.
4. Let  $n_i$  equal the number of  $y_k$ 's in  $\mathbf{y}_i$  on which  $c$  and  $h_i$  disagree.
5. If  $n_i \leq \lfloor \frac{3}{4}\varepsilon L_i \rfloor$  then return  $h_i$ . Otherwise continue to the next iteration.

Since  $C$  is nonuniformly learnable, for each target concept  $c$ , and  $\varepsilon, \delta > 0$  there exists an  $\ell = \ell(\varepsilon/2, \delta/3, c)$  such that when  $f$  is given  $\ell$  random examples it produces a concept  $\varepsilon/2$ -close to  $c$  with probability at least  $1 - \delta/3$ . Let  $C_j$  be the first  $C_i$  that contains  $c$ . Let  $\ell_j = \max\left\{\frac{12}{\varepsilon} \ln \frac{6}{\delta}, \frac{16 \dim(C_j)}{\varepsilon} \ln\left(\frac{26}{\varepsilon}\right)\right\}$ . By Theorem 1,  $\ell_j$  examples are sufficient to learn any concept of  $C_j$  with accuracy  $\varepsilon/2$  and confidence  $\delta/3$ .

The algorithm may fail to learn for the following reasons:

1. It stopped with an  $h_i$  which is  $\varepsilon$ -far from  $c$ .
2. It did not halt.

We will use the following Chernoff inequalities (see [4, proposition 2.4]):

For all  $n, p, \beta$  with  $0 \leq p \leq 1$ ,  $0 \leq \beta \leq 1$

$$\sum_{k=0}^{\lfloor (1-\beta)np \rfloor} \binom{n}{k} p^k (1-p)^{n-k} \leq e^{-\beta^2 np/2}, \quad (1)$$

$$\sum_{k=\lceil (1+\beta)np \rceil}^n \binom{n}{k} p^k (1-p)^{n-k} \leq e^{-\beta^2 np/3}. \quad (2)$$

1. Let  $P_i$  be the probability that  $n_i < \frac{3}{4}\varepsilon L_i$  given that  $h_i$  is  $\varepsilon$ -far from  $c$ . From (1) we get

$$\begin{aligned} P_i &\leq \sum_{j=0}^{\lfloor \frac{3}{4}\varepsilon L_i \rfloor} \binom{L_i}{j} \varepsilon^j (1-\varepsilon)^{L_i-j} \\ &\leq e^{-\frac{1}{4^2} L_i \frac{\varepsilon}{2}} = e^{-L_i \varepsilon/32} \\ &= e^{(-L_0 - i)\varepsilon/32} = e^{-L_0 \varepsilon/32} \left(e^{-\varepsilon/32}\right)^i. \end{aligned}$$

Thus the probability of failure in this case is

$$\begin{aligned} \sum_{i=0}^{\infty} P_i &\leq e^{-L_0 \varepsilon/32} \sum_{i=0}^{\infty} \left(e^{-\varepsilon/32}\right)^i \\ &= e^{-L_0 \varepsilon/32} \frac{1}{1 - e^{-\varepsilon/32}} \\ &= e^{-\lceil \frac{32}{\varepsilon} [(\ln \frac{3}{\delta}) - \ln(1 - e^{-\varepsilon/32})] \rceil \frac{\varepsilon}{32}} \frac{1}{1 - e^{-\varepsilon/32}} \\ &\leq e^{-\frac{32}{\varepsilon} [(\ln \frac{3}{\delta}) - \ln(1 - e^{-\varepsilon/32})] \frac{\varepsilon}{32}} \frac{1}{1 - e^{-\varepsilon/32}} \\ &= \frac{\delta}{3}. \end{aligned}$$

2. Consider iteration  $r = \lceil \log_2 \ell_j \rceil$ . By Theorem 1, there is probability at most  $\delta/3$  that  $h_r$  is  $\varepsilon/2$ -far from  $c$ . Thus it remains to consider only the case that  $h_r$  is  $\varepsilon/2$ -close and the algorithm did not halt.

Let  $Q$  be the probability that  $n_r > \lfloor \frac{3}{4}\varepsilon L_r \rfloor$ . From (2) we get

$$\begin{aligned} Q &\leq \sum_{k=\lfloor \frac{3}{4}\varepsilon L_r \rfloor + 1}^{L_r} \binom{L_r}{k} \left(\frac{\varepsilon}{2}\right)^k \left(1 - \frac{\varepsilon}{2}\right)^{L_r - k} \leq \sum_{k=\lfloor \frac{3}{4}\varepsilon L_r \rfloor}^{L_r} \binom{L_r}{k} \left(\frac{\varepsilon}{2}\right)^k \left(1 - \frac{\varepsilon}{2}\right)^{L_r - k} \\ &\leq e^{-2^{-2}L_r\varepsilon/3} = e^{-L_r\varepsilon/12} \leq e^{-L_0\varepsilon/12} < \frac{\delta}{3}. \end{aligned}$$

Thus the probability of this event is less than  $\frac{2}{3}\delta$ . □

**Remark 1:** The number of examples needed is

$$\sum_{i=1}^r (2^i + L_i) < 2^{r+1} + rL_0 + \frac{1}{2}r^2 = O(2^r + rL_0) = O(2^r + r\ell_j) = O(\ell_j \log \ell_j)$$

where  $\ell_j$  was defined in the proof.

**Remark 2:** If  $f$  is computable then so is the algorithm.

## 5 Classes of languages learnable by polynomially many examples.

**Theorem 4:** Let  $X = \{0, 1\}^*$  (the set of all finite bit-strings) and  $C$  be any set of recursive languages. Then  $C$  is nonuniformly strongly learnable. Moreover, for every target concept  $c \in C$  if  $M$  is a Turing machine with  $j$  states that recognizes  $c$ , then  $c$  is learnable with accuracy  $\varepsilon$  and confidence  $\delta$  by  $O(\frac{1}{\varepsilon} \log^2 \frac{1}{\delta} + \frac{j}{\varepsilon} \log^3 \frac{j}{\varepsilon})$  many examples.

**Proof:** Let  $C_k$  be the set of recursive languages recognizable by Turing machines with  $k$  states or less. Clearly  $C = \bigcup_{i=1}^{\infty} C_i$  and since  $C_k$  is finite so is its VC-dimension. Thus  $C$  is nonuniformly learnable.

To show that the number of examples needed is polynomial we find the VC-dimension of  $C_j$ . The number of languages in  $C_j$  is bounded by the number of Turing machines with  $j$  states or less. Using the standard model for Turing machines [14] (the head can move to the left, right or remain stationary) over the alphabet  $\Sigma = \{0, 1\}$  there are at most  $(6j)^{2j}$  Turing machines with  $j$  states or less. (Notice that any Turing machine with less than  $j$  states is equivalent to one with exactly  $j$  states and a transition function that never reaches some of the states.) Since for any finite class  $C$ ,  $\dim(C) < \log_2(|C|)$ ,  $\dim(C_j) \leq 2j \log_2 6j$ .

By Remark 1,  $O(\ell_j \log \ell_j)$  examples suffice. Where

$$\ell_j = O\left(\max\left\{\frac{1}{\varepsilon} \log \frac{1}{\delta}, \frac{\dim(C_j)}{\varepsilon} \log \frac{1}{\varepsilon}\right\}\right) = O\left(\max\left\{\frac{1}{\varepsilon} \log \frac{1}{\delta}, \frac{j}{\varepsilon} (\log j) \log \frac{1}{\varepsilon}\right\}\right).$$



Case 1,  $\log \frac{1}{\delta} \geq j \log j \log \frac{1}{\varepsilon}$ :

$$\begin{aligned} O(\ell_j \log \ell_j) &= O\left(\frac{1}{\varepsilon} \log \frac{1}{\delta} \left[\log\left(\frac{1}{\varepsilon} \log \frac{1}{\delta}\right)\right]\right) = O\left(\frac{1}{\varepsilon} \log \frac{1}{\delta} \left[\log \frac{1}{\varepsilon} + \log \log \frac{1}{\delta}\right]\right) \\ &= O\left(\frac{1}{\varepsilon} \log \frac{1}{\delta} \left[\log \frac{1}{\delta} + \log \log \frac{1}{\delta}\right]\right) = O\left(\frac{1}{\varepsilon} \log^2 \frac{1}{\delta}\right). \end{aligned}$$

Case 2,  $\log \frac{1}{\delta} < j \log j \log \frac{1}{\varepsilon}$ :

$$\begin{aligned} O(\ell_j \log \ell_j) &= O\left(\frac{j}{\varepsilon} \log j \log \frac{1}{\varepsilon} \left[\log \frac{1}{\varepsilon} + \log j + \log \log j + \log \log \frac{1}{\varepsilon}\right]\right) \\ &= O\left(\frac{j}{\varepsilon} \log j \log \frac{1}{\varepsilon} \left[\log \frac{1}{\varepsilon} + \log j\right]\right). \end{aligned}$$

If  $\frac{1}{\varepsilon} \leq j$  then  $O(\ell_j \log \ell_j) = O\left(\frac{j}{\varepsilon} \log^3 j\right)$ . Otherwise,  $O(\ell_j \log \ell_j) = O\left(\frac{j}{\varepsilon} \log^3 \frac{1}{\varepsilon}\right)$ .

Therefore, in this case  $O(\ell_j \log \ell_j) = O\left(\frac{j}{\varepsilon} \log^3 \frac{j}{\varepsilon}\right)$ .  $\square$

As defined, learnability does not require computability. In the above proof we use a function  $f$  that, given  $\ell_j$  examples, returns a Turing machine with  $j$  or less states consistent with the examples. This function is related to the Kolmogorov complexity and is not recursive:

**Theorem 5:** For words  $x, y$  let  $K_2(x, y)$  denote the minimum  $k$  such that there exists a  $k$  state Turing machine that halts on both  $x$  and  $y$ , accepts  $x$  and rejects  $y$ . Then  $K_2(x, y)$  is not recursive.

**Proof:** Assume, for the sake of contradiction, that  $K_2$  is recursive. Let  $z$  be a binary word of length  $\ell$ . Define  $z^0 = [z_1, \dots, z_{\lfloor \ell/2 \rfloor}]$  and  $z^1 = [z_{\lfloor \ell/2 \rfloor + 1}, \dots, z_\ell]$ . Let  $K^*(z)$  be the minimum  $k$  such that there exists a Turing machine with  $k$  states that for  $z^0 \neq z^1$  accepts  $z^0$  and rejects  $z^1$ . The assumption that  $K_2$  is recursive implies that  $K^*$  is also recursive.

Let  $Z_n$  be the lexicographically first word  $z$  such that  $K^*(z) > n$ . ( $Z_n$  is well defined since for all  $n$  there exists a word  $z$  for which  $K^*(z) > n$ .)

For every  $n$  we can construct a Turing machine,  $T_n$ , to produce  $Z_n$ :

1. Write the binary representation of  $n$  on the working tape.
2. Consider all strings  $z$  in lexicographic order:
3. If  $z^0 \neq z^1$  compute  $K^*(z)$  and if  $K^*(z) > n$  copy  $z$  to the output tape.

Let  $k^*$  be the number of states of  $K^*$ . There exists a constant  $c'$  such that for all  $n$   $T_n$  has at most  $k^* + c' + \log n$  states.

Let  $K(z)$  be the minimum  $k$  such that there is a Turing machine with  $k$  states that produces  $z$ . From the definition of  $K$ ,

$$K(Z_n) \leq k^* + c' + \log n.$$

We now show that there exists a constant  $c$  such that

$$K^*(z) \leq K(z) + c.$$

Let  $c$  be the number of states of a Turing machine  $M^1(u, z)$  that accepts if  $u = z^0$  and rejects if  $u = z^1$ . Let  $M^z$  be the Turing machine with  $K(z)$  states that produces  $z$ , and  $M^*$  the Turing machine that on input  $u$  first applies  $M^z$  to produce  $z$  then applies  $M^1(u, z)$ .  $M^*$  accepts  $z^0$ , rejects  $z^1$  and has  $K(z) + c$  states.

However, from the definition of  $K$ ,  $n \leq K^*(Z_n)$ .

From which we get that there exist constants  $c$  and  $c'$  such that for infinitely many  $n$ ,

$$n \leq \log n + c + c' + k^*.$$

A contradiction □

The previous theorem does not imply that it is undecidable to learn, only that our proof of Theorem 4 does not yield a recursive algorithm.

**Remark 3:** Results similar to Theorem 4 may be proved for other concept classes, e.g., the class of r.e. languages.

A set  $C$  of recursive languages for which one of the following conditions holds is learnable with polynomially many examples by a computable function:

1. There exists a recursive set  $S$  of (encodings of) total Turing machines such that for every  $L \in C$  there exists a Turing machine  $T$  in  $S$  whose language is  $L$ .
2.  $C$  has bounded complexity, i.e., there is a recursive function  $g$  such that for every  $L \in C$  there exists a Turing machine  $T$  that recognizes  $L$  and for every word  $w$ ,  $T$  stops after at most  $g(|w|)$  steps, where  $|w|$  is the length of  $w$ .

**Corollary 2:** *There is a computable function that learns Boolean formulae by a polynomial number of examples.*

**Proof:** Let  $f$  be a Boolean formula with  $v$  variables such that  $f$  can be encoded by  $n$  bits. For a word  $b \in \{0, 1\}^*$  we define the value of  $f(b)$  as follows: if  $|b| \geq v$  then assign the first  $v$  bits of  $b$  to the respective variables otherwise assign 0 for the unassigned variables. We build a Turing machine  $M_f$  that on input  $b \in \{0, 1\}^*$  accepts  $b$  if and only if  $f(b) = \text{true}$ . For every  $f$  there is such an  $M_f$  with  $n + O(1)$  states. (Simply write on the input tape the description of  $f(b)$  ( $n$  states) and then activate a Turing machine that given a  $v$ -variable Boolean formula and a Boolean vector of length  $v$  evaluates the value of the formula applied to the vector.) Furthermore, the time complexity of the above  $M_f$  is bounded. Thus, by condition (2) above, Boolean formulae are learnable by a computable function. □

Note that we do not address Valiant's problem concerning the computation time. However, Pitt and Valiant [18] showed several classes of Boolean formulae for which there is no polynomial algorithm that finds a bounded function in that class consistent with a set of examples unless  $RP = P$ . In particular:  $k$ -TERM-DNF,  $k$ -CLAUSE-CNF (and their monotonic forms) and  $\mu$ -expressions. Thus even though the number of examples is polynomial, the computation time needed for learning is not necessarily polynomial.

## 6 Nonuniform learnability for a given distribution.

In order to extend the notion of learnability Benedek and Itai [7] have defined:

**Uniform Learnability for a given distribution  $D$ :**  $C$  is *learnable with respect to  $D$*  if there is a function  $f \in F_C$  such that for every  $\varepsilon, \delta > 0$  there is an  $\ell = \ell(\varepsilon, \delta) > 0$  such that for every  $c \in C$ ,  $r_f^+(D, c, \ell, \varepsilon) > 1 - \delta$ .

**Finite and countable covers:** A subset  $C_\varepsilon$  of  $2^X$  is an  $\varepsilon$ -cover of  $C$  with respect to  $D$  if for every  $c \in C$  there is a  $\tilde{c} \in C_\varepsilon$   $\varepsilon$ -close to it.  $C$  is *finitely coverable (with respect to  $D$ )* if for every  $\varepsilon > 0$  there is a finite  $\varepsilon$ -cover  $C_\varepsilon$  of  $C$ .  $C$  is *countably coverable (with respect to  $D$ )* if for every  $\varepsilon > 0$  there is a countable  $\varepsilon$ -cover  $C_\varepsilon$  of  $C$ . In the sequel we omit  $D$  when understood from the context.

The main result of [7] is:

**Theorem 6:**  $C$  is *finitely coverable (with respect to  $D$ )* if and only if  $C$  is *learnable (with respect to  $D$ )*.

Similarly to distribution-free learnability, learnability for a given distribution can also be extended to the nonuniform form.

**Nonuniform learnability for a given distribution  $D$ :**  $C$  is *nonuniformly learnable with respect to distribution  $D$*  if there is a function  $f \in F_C$  such that for every  $\varepsilon, \delta > 0$  and every  $c \in C$  there is an  $\ell = \ell(c, \varepsilon, \delta) > 0$  such that  $r_f^+(D, c, \ell, \varepsilon) > 1 - \delta$ .

**Theorem 7:**  $C$  is *nonuniformly learnable with respect to  $D$*  if and only if  $C$  is *countably coverable with respect to  $D$* .

The next lemma follows directly from the definition of countable coverable.

**Lemma 2:**  $C$  is *countably coverable with respect to  $D$*  if and only if there exists a countable class  $C^*$  such that for all  $\varepsilon > 0$ ,  $C^*$  is an  $\varepsilon$ -cover of  $C$ .

The proof of Theorem 7 is similar to that of Theorem 2.

In the same manner we define:

**Nonuniform strong learnability for a given distribution  $D$ :**  $C$  is *nonuniformly strongly learnable with respect to  $D$*  if there is a function  $f \in F_C$  such that for every  $\varepsilon, \delta > 0$  and every  $c \in C$  there is an  $\ell = \ell(c, \varepsilon, \delta) > 0$  such that  $\sum_{i=1}^{\ell-1} r_f^-(D, c, i, \varepsilon) < \delta$  and  $r_f^+(D, c, \ell, \varepsilon) > 1 - \delta$ .

The following theorem is analogous to Theorem 3:

**Theorem 8:** A concept class is *nonuniformly strongly with respect to  $D$  learnable* if and only if it is *nonuniformly learnable with respect to  $D$* .

**Example 5:** Let  $C_{\text{OPEN}}$  be as in example 2 and  $C$  be the set of all finite unions of rational intervals. Let  $D$  be the uniform distribution over  $(0, 1)$ . It can be shown that for every  $\varepsilon > 0$ ,  $C$  is a countable  $\varepsilon$ -cover of  $C_{\text{OPEN}}$  with respect to  $D$ . Thus  $C_{\text{OPEN}}$  is nonuniformly learnable with respect to  $D$ . (The above is true for every continuous distribution.)

**Remark 4:** In [7] the authors prove that for every given distribution  $D$  over  $\{0, 1\}^*$  every concept class over  $\{0, 1\}^*$  is uniformly learnable with respect to  $D$ . Thus any concept class of recursive or even r.e. concepts is learnable.

**Remark 5:** The definition of learnability may be extended to allow random algorithms and random functions. All the theorems mentioned hold also for this case.

## 7 Conclusions

The notion of PAC learnability has been extended to the case where the sample size may depend on the target concept. We have shown that a concept class can be learned under this definition if and only if it is a countable union of classes with finite VC-dimension.

The problem of characterizing concept classes whose learning time is a fixed polynomial of the size of the target concept and the error parameters still remains open.

## ACKNOWLEDGEMENTS

It is a pleasure to thank Shai Ben-David for contributing Lemma 1.

## References

- [1] D. Angluin. Queries and concept learning. *Machine Learning*, 2:319–342, 1988.
- [2] D. Angluin and P. Laird. Learning from noisy examples. *Machine Learning*, 2:343–370, 1988.
- [3] Angluin D. and Smith C.H., *Inductive Inference: Theory and Methods*, Computing Surveys, 15(3), (Sept 1983).
- [4] Angluin D. and Valiant L.G. *Fast probabilistic algorithms for Hamiltonian circuits and matchings*, JCSS 18, 155-193, (1979).
- [5] Ben-David, S. and G. M. Benedek, *Measurability constraints on PAC learnability*, TR Department of Computer Science, Technion, Haifa, 1992.
- [6] Ben-David, S., G. M. Benedek and Y. Mansour, *Parameterization scheme for classifying models of learnability*, to appear in Information and Computation
- [7] Benedek G.M. and A. Itai, *Learnability by fixed distributions*, to appear in Theoretical Computer Science. (A preliminary version appeared in COLT '88 (1988).)
- [8] Benedek G.M. and A. Itai, *Nonuniform learnability*, 15 International Colloquium on Automata, Languages and Programming, (ICALP), Tampere, Finland, 1988. Lecture notes in computer science, 317, T. Lepistö and A. Salomaa (eds.), Springer-Verlag.

- [9] Blumer A., A. Ehrenfeucht, D. Haussler and M. Warmuth, *Learnability and the Vapnik-Chervonenkis dimension*, J. ACM, 36(4), 929-965, (1989).
- [10] Blumer A., A. Ehrenfeucht, D. Haussler and M. Warmuth, *Occam's razor*, Inf. Proc. Letters 24 (1987), 377-380, North-Holland.
- [11] Ehrenfeucht A., D. Haussler, M. Kearns and L. Valiant, *A general lower bound on the number of examples needed for learning*, COLT '88.
- [12] Halmos, P. R., *Measure Theory*, Van Nostrand, (1950).
- [13] Haussler D., M. Kearns, N. Littlestone and M. Warmuth, *Equivalence of Models for Polynomial Learnability*, COLT '88 (1988).
- [14] Hopcroft J.E. and J. D. Ullman, *Introduction to automata theory, languages and computation*, Addison-Wesley (1979).
- [15] Kearns M., Ming Li, L. Pitt, L. G. Valiant, *On the learnability of Boolean formulae*, Proc. of 19<sup>th</sup> Symp. Theory of Comp., 285-295. ACM, New York, (1987).
- [16] Lineal, N., Y. Mansour and R. Rivest, *Results on learnability and the Vapnik-Chervonenkis dimension*, FOCS 1988.
- [17] Natarajan B. K., *On learning Boolean functions*, In Proc. of 19<sup>th</sup> Symp. Theory of Comp., 296-304. ACM, New York, (1987).
- [18] Pitt L. and L. G. Valiant, *Computational limitations on learning from examples*, Aiken Computation Laboratory, Harvard University, Cambridge, MA 02138, (July 1986).
- [19] Shelah S., Unpublished manuscript.
- [20] Vapnik V.N. and A.Ya. Chervonenkis, *On the uniform convergence of relative frequencies of events to their probabilities*, Th. Prob. and its Appl., 16(2), 264-80, (1971).
- [21] Valiant L.G., *A Theory of the Learnable*, Comm. ACM, 27(11), 1134-42, (1984).
- [22] Valiant L.G., *Learning disjunctions of conjunctions*, Proceedings of 9<sup>th</sup> IJCAI, vol. 1, 560-566, Los Angeles, CA., (August 1985).
- [23] Valiant L.G., *Deductive learning*, Aiken Computational Laboratory, Harvard University, (1984).