

TCS-TR-A-05-08

TCS Technical Report

From Learning in the Limit to Stochastic Finite Learning

by

THOMAS ZEUGMANN

Division of Computer Science

Report Series A

August 27, 2005



Hokkaido University
Graduate School of
Information Science and Technology

Email: thomas@ist.hokudai.ac.jp

Phone: +81-011-706-7684

Fax: +81-011-706-7684

From Learning in the Limit to Stochastic Finite Learning

THOMAS ZEUGMANN

Division of Computer Science

Hokkaido University

N-14, W-9, Sapporo 060-0814, Japan

thomas@ist.hokudai.ac.jp

August 27, 2005

Abstract

Inductive inference can be considered as one of the fundamental paradigms of algorithmic learning theory. We survey results recently obtained and show their impact to potential applications.

Since the main focus is put on the efficiency of learning, we also deal with postulates of naturalness and their impact to the efficiency of limit learners. In particular, we look at the learnability of the class of all pattern languages and ask whether or not one can design a learner within the paradigm of learning in the limit that is nevertheless efficient.

For achieving this goal, we deal with iterative learning and its interplay with the hypothesis spaces allowed. This interplay has also a severe impact to postulates of naturalness satisfiable by any learner.

Furthermore, since a limit learner is only supposed to converge, one never knows at any particular learning stage whether or not the learner did already succeed. The resulting uncertainty may be prohibitive in many applications. We survey results to resolve this problem by outlining a new learning model, called *stochastic finite learning*. Though pattern languages can neither be finitely inferred from positive data nor PAC-learned, our approach can be extended to a stochastic finite learner that *exactly* infers all pattern languages from positive data with *high confidence*.

Finally, we apply the techniques developed to the problem of learning conjunctive concepts.

1. Introduction

Inductive inference can be considered as one of the fundamental paradigms of algorithmic learning theory. In particular, inductive inference of recursive functions and of recursively enumerable languages has been studied intensively within the last

four decades (cf., e.g., [3, 4, 30, 16]). The basic model considered within this framework is learning in the limit which can be informally described as follows. The learner receives more and more data about the target and maps these data to hypotheses. Of special interest is the investigation of scenarios in which the sequence of hypotheses *stabilizes* to an *accurate* and *finite* description (e.g. a grammar, a program) of the target. Clearly, then some form of learning must have taken place. Here by data we mean either any infinite sequence of pairs argument-value (in case of learning recursive functions) such that all arguments appear eventually or any infinite sequence of all members of the target language (in case of language learning from positive data). Alternatively, one can also study language learning from both positive and negative data.

Most of the work done in the field has been aimed at the following goals: showing what general collections of function or language classes are learnable, characterizing those collections of classes that can be learned, studying the impact of several postulates on the behavior of learners to their learning power, and dealing with the influence of various parameters to the efficiency of learning. However, defining an appropriate measure for the complexity of learning in the limit has turned out to be quite difficult (cf. Pitt [31]). Moreover, whenever learning in the limit is done, in general one never knows whether or not the learner has already converged. This is caused by the fact that it is in general undecidable whether or not convergence already occurred. But even if it is decidable, it is practically infeasible to do so. Thus, there is always an uncertainty which may be prohibitive in many applications of learning.

Therefore, different learning models have been proposed. In particular, Valiant's [46] model of *probably approximately correct* (abbr. PAC) learning has been very influential. As a matter of fact, this model puts strong emphasis on the efficiency of learning and avoids the problem of convergence at all. In the PAC model, the learner receives a finite labeled sample of the target concept and outputs, with high probability, a hypothesis that is approximately correct. The sample is drawn with respect to an unknown probability distribution and the error of as well as the confidence in the hypothesis are measured with respect to this distribution, too. Thus, if a class is PAC learnable, one obtains nice performance guarantees. Unfortunately, many interesting concept classes are not PAC learnable.

Consequently, one has to look for other models of learning or one is back to learning in the limit. So, let us assume that learning in the limit is our method of choice. What we would like to present in this survey is a rather general way to transform learning in the limit into *stochastic finite learning*. It should also be noted that our ideas may be beneficial even in case that the considered concept class is PAC learnable.

Furthermore, we aim to outline how a thorough study of limit learnability of concept classes may nicely contribute to support our new approach. We exemplify the research undertaken by looking at the class of all pattern languages introduced by Angluin [1]. As Salomaa [37] has put it "Patterns are everywhere" and thus we believe that our research is worth the effort undertaken.

There are several problems that have to be addressed when dealing with the learnability of pattern languages. First, the nice thing about patterns is that they are very intuitive. Therefore, it seems desirable to design learners outputting patterns as their hypotheses. Unfortunately, membership is known to be \mathcal{NP} -complete for the pattern languages (cf. [1]). Thus, many of the usual approaches used in machine learning will directly lead to infeasible learning algorithms. As a consequence, we shall ask what kind of appropriate hypothesis spaces can be used at all to learn the pattern languages, and what are the appropriate learning strategies.

In particular, we shall deal with the problem of redundancy in the hypothesis space chosen, with consistency, conservativeness, and iterative learning. Here consistency means that the intermediate hypotheses output by the learner do correctly reflect the data seen so far. Conservativeness addresses the problem to avoid overgeneralization, i.e., preventing the learner from guessing a proper superset of the target language. These requirements are naturally arising desiderata, but this does not mean that they can be fulfilled. With *iterative* learning, the learning machine, in making a conjecture, has access to only its previous conjecture and the latest data item coming in. Iterative learning is also a natural requirement whenever learning in the limit is concerned, since no practical learner can process at every learning stage all examples provided so far, it may even not be able to store them.

Then, we address the question how efficient the overall learning process can be performed, and how we can get rid of the uncertainty of not knowing whether or not the learner has already converged.

Finally, we show our ideas to be beneficial for a class known to be PAC learnable by looking at the class of all concepts describable by a monomial. A concept is describable by a monomial m if its elements are precisely the assignments satisfying m .

2. Preliminaries

Unspecified notation follows Rogers [35]. By $\mathbb{N} = \{0, 1, 2, \dots\}$ we denote the set of all natural numbers. We set $\mathbb{N}^+ = \mathbb{N} \setminus \{0\}$. The cardinality of a set S is denoted by $|S|$. Let \emptyset , \in , \subset , \subseteq , \supset , and \supseteq denote the empty set, element of, proper subset, subset, proper superset, and superset, respectively.

Let $\varphi_0, \varphi_1, \varphi_2, \dots$ denote any fixed *acceptable programming system* for all (and only) the partial recursive functions over \mathbb{N} (cf. Rogers [35]). Then φ_k is the partial recursive function computed by *program* k .

In the following subsection we define the main learning models considered within this paper.

2.1. Learning in the Limit

Gold's [12] model of learning in the limit allows one to formalize a rather general class of learning problems, i.e., learning from examples. For defining this model we

assume any recursively enumerable set \mathcal{X} and refer to it as the *learning domain*. By $\wp(\mathcal{X})$ we denote the power set of \mathcal{X} . Let $\mathcal{C} \subseteq \wp(\mathcal{X})$, and let $c \in \mathcal{C}$ be non-empty; then we refer to \mathcal{C} and c as a *concept class* and a *concept*, respectively. Let c be a concept, and let $t = (x_j)_{j \in \mathbb{N}}$ be any infinite sequence of elements $x_j \in c$ such that $\text{range}(t) := \{x_j \mid j \in \mathbb{N}\} = c$. Then t is said to be a *positive presentation* or, synonymously, a *text* for c . By $\text{text}(c)$ we denote the set of all positive presentations for c . Moreover, let t be a positive presentation, and let $y \in \mathbb{N}$. Then, we set $t_y = x_0, \dots, x_y$, i.e., t_y is the initial segment of t of length $y+1$, and $t_y^+ := \{x_j \mid j \leq y\}$. We refer to t_y^+ as the *content* of t_y .

Furthermore, let $\sigma = x_0, \dots, x_{n-1}$ be any finite sequence. Then we use $|\sigma|$ to denote the *length* n of σ , and let σ^+ denote the content of σ . Additionally, let t be a text and let τ be a finite sequence; then we use $\sigma \diamond t$ and $\sigma \diamond \tau$ to denote the sequence obtained by *concatenating* σ onto the front of t and τ , respectively.

Alternatively, one can also consider *complete presentations* or, synonymously, *informants*. An informant for a concept c is an infinite sequence of all elements of the underlying learning domain that are classified with respect to their containment in c . More formally, we define informants as follows. Let c be a concept; then any sequence $i = (x_j, b_j)_{j \in \mathbb{N}}$ of labeled examples, where $b_j \in \{0, 1\}$ such that $\{x_j \mid j \in \mathbb{N}\} = \mathcal{X}$ and $i^+ = \{x_j \mid (x_j, b_j) = (x_j, 1), j \in \mathbb{N}\} = c$ and $i^- = \{x_j \mid (x_j, b_j) = (x_j, 0), j \in \mathbb{N}\} = \mathcal{X} \setminus c$ is called an informant for c . For the sake of presentation, the following definitions are only given for the text case, the generalization to the informant case should be obvious. We sometimes use the term *data sequence* to refer to both text and informant, respectively.

An *inductive inference machine* (abbr. IIM) is an algorithm that takes as input larger and larger initial segments of a text and outputs, after each input, a hypothesis from a prespecified *hypothesis space* $\mathcal{H} = (h_j)_{j \in \mathbb{N}}$. The indices j are regarded as suitable finite encodings of the concepts described by the hypotheses. A hypothesis h is said to describe a concept c iff $c = h$.

A sequence $(j_n)_{n \in \mathbb{N}}$ of natural numbers is said to converge to number j if $j_n = j$ for all but finitely many $n \in \mathbb{N}$.

DEFINITION 1. Let \mathcal{C} be any concept class, and let $\mathcal{H} = (h_j)_{j \in \mathbb{N}}$ be a hypothesis space for it. \mathcal{C} is called *learnable in the limit from text with respect to \mathcal{H}* iff there is an IIM M such that for every $c \in \mathcal{C}$ and every text t for c ,

- (1) for all $n \in \mathbb{N}^+$, $M(t_n)$ is defined,
- (2) there is a j such that $c = h_j$ and the sequence $(M(t_n))_{n \in \mathbb{N}}$ converges to j .

The set of all concepts classes that are learnable in the limit with respect to \mathcal{H} is denoted by $\text{LimTxt}_{\mathcal{H}}$. By LimTxt we denote the collection of all concepts classes \mathcal{C} for which there is a hypothesis space \mathcal{H} such that \mathcal{C} is learnable in the limit from text¹ with respect to \mathcal{H} .

¹If learning from informant is considered we use $\text{LimInf}_{\mathcal{H}}$ and LimInf in an analogous way.

Note that instead of *LimTxt* sometimes *TxtEx* is used. In our notation, *Lim* stands for “limit.” Suppose, an IIM learns some concept c . That means, after having seen only finitely many data of c the IIM reached its (unknown) point of convergence and it computed a *correct* and *finite* description of the target concept. Hence, some form of learning must have taken place.

Note that Definition 1 does not contain any requirement concerning efficiency. Before we are going to deal with efficiency, we want to point to another crucial parameter of our learning model, i.e., the hypothesis space \mathcal{H} . Since our goal is *algorithmic* learning, we can consider the special case that $\mathcal{X} = \mathbb{N}$ and let \mathcal{C} be any subset of the collection of all recursively enumerable sets over \mathbb{N} . Let $W_k = \text{domain}(\varphi_k)$, where φ_k is the partial recursive function computed by program k in the fixed acceptable programming system. In this case, $(W_k)_{k \in \mathbb{N}}$ is the most general hypothesis space.

Within this setting many learning problems can be described. Moreover, this setting has been used to study the general capabilities of different learning models which can be obtained by suitable modifications of Definition 1. There are numerous papers performing studies along this line of research (cf., e.g., [16, 30] and the references therein). On the one hand, the results obtained considerably broaden our general understanding of algorithmic learning. On the other hand, one has also to ask what kind of consequences one may derive from these results for practical learning problems. This is a non-trivial question, since the setting of learning recursively enumerable languages is very rich. Thus, it is conceivable that several of the phenomena observed hold in this setting due to the fact that too many sets are recursively enumerable and that there are no counterparts within the world of efficient computability.

As a first step to address this question we mainly consider the scenario that *indexable concept* classes with uniformly decidable membership have to be learned (cf. Angluin [2]). A class of non-empty concepts \mathcal{C} is said to be an ***indexable class*** with uniformly decidable membership provided there are an effective enumeration c_0, c_1, c_2, \dots of all and only the concepts in \mathcal{C} and a recursive function f such that for all $j \in \mathbb{N}$ and all elements $x \in \mathcal{X}$ we have

$$f(j, x) = \begin{cases} 1, & \text{if } x \in c_j, \\ 0, & \text{otherwise.} \end{cases}$$

In the following we refer to indexable classes with uniformly decidable membership as to indexable classes for short. Furthermore, we call any enumeration $(c_j)_{j \in \mathbb{N}}$ of \mathcal{C} with uniformly decidable membership problem an *indexed family*.

Since the paper of Angluin [2], learning of indexable concept classes has attracted much attention (cf., e.g., Zeugmann and Lange [52]). Let us shortly provide some well known indexable classes. Let Σ be any finite alphabet of symbols, and let \mathcal{X} be the free monoid over Σ , i.e., $\mathcal{X} = \Sigma^*$. We set $\Sigma^+ = \Sigma^* \setminus \{\lambda\}$, where λ denotes the empty string. As usual, we refer to subsets $L \subseteq \mathcal{X}$ as languages. Then the set of all regular languages, context-free languages, and context-sensitive languages are indexable classes.

Next, let $X_n = \{0, 1\}^n$ be the set of all n -bit Boolean vectors. We consider $\mathcal{X} = \bigcup_{n \geq 1} X_n$ as learning domain. Then, the set of all concepts expressible as a monomial, a k -CNF, a k -DNF, and a k -decision list form indexable classes.

When learning indexable classes \mathcal{C} , it is generally assumed that the hypothesis space \mathcal{H} has to be an indexed family, too. We distinguish *class preserving learning* and *class comprising learning* defined by $\mathcal{C} = \text{range}(\mathcal{H})$ and $\mathcal{C} \subseteq \text{range}(\mathcal{H})$, respectively. When dealing with class preserving learning, one has the freedom to choose as hypothesis space a possibly *different enumeration* of the target class \mathcal{C} . In contrast, when class comprising learning is concerned, the hypothesis space may enumerate, additionally, languages not belonging to \mathcal{C} . Note that, in general, one has to allow class comprising hypothesis spaces to obtain the maximum possible learning power (cf. Lange and Zeugmann [20, 22]). Finally, we call an hypothesis space *redundant* if it is larger than necessary, i.e., there is at least one hypothesis in \mathcal{H} not describing any concept from the target class or one concept possesses at least two different descriptions in \mathcal{H} . Thus, *non-redundant* hypothesis spaces are as small as possible.

Formally, a hypothesis space $\mathcal{H} = (h_j)_{j \in \mathbb{N}}$ is *non-redundant* for some target concept class \mathcal{C} iff $\text{range}(\mathcal{H}) = \mathcal{C}$ and $h_i \neq h_j$ for all $i, j \in \mathbb{N}$ with $i \neq j$. Otherwise, \mathcal{H} is a *redundant* hypothesis space for \mathcal{C} .

Next, let us come back to the issue of efficiency. Looking at Definition 1 we see that an IIM M has always access to the whole history of the learning process, i.e., in order to compute its actual guess M is fed all examples seen so far. In contrast to that, next we define *iterative IIMs*. An iterative IIM is only allowed to use its last guess and the next element in the positive presentation of the target concept for computing its actual guess. Conceptionally, an iterative IIM M defines a sequence $(M_n)_{n \in \mathbb{N}}$ of machines each of which takes as its input the output of its predecessor.

DEFINITION 2 (Wiehagen [47]). Let \mathcal{C} be a concept class, let c be a concept, and let $\mathcal{H} = (h_j)_{j \in \mathbb{N}}$ be a hypothesis space. An IIM M *ItLimTxt $_{\mathcal{H}}$ -learns* c iff for every $t = (x_j)_{j \in \mathbb{N}} \in \text{text}(c)$ the following conditions are satisfied:

- (1) for all $n \in \mathbb{N}$, $M_n(t)$ is defined, where $M_0(t) := M(x_0)$ and for all $n \geq 0$:

$$M_{n+1}(t) := M(M_n(t), x_{n+1}),$$
- (2) the sequence $(M_n(t))_{n \in \mathbb{N}}$ converges to a number j such that $c = h_j$.

Finally, M *ItLimTxt $_{\mathcal{H}}$ -learns* \mathcal{C} iff, for each $c \in \mathcal{C}$, M *ItLimTxt $_{\mathcal{H}}$ -learns* c . By *ItLimTxt* we denote the collection of all concept classes for which there are an IIM M and a hypothesis space \mathcal{H} such that M *ItLimTxt $_{\mathcal{H}}$ -learns* \mathcal{C} .

In the latter definition $M_n(t)$ denotes the $(n+1)$ st hypothesis output by M when successively fed the text t . So, it is justified to make the following convention. Let $\sigma = x_0, \dots, x_n$ be any finite sequence of elements from \mathcal{X} . Moreover, let \mathcal{C} be any concept class over \mathcal{X} , and let M be any IIM that iteratively learns \mathcal{C} . Then we denote by $M_y(\sigma)$ the $(y+1)$ st hypothesis output by M when successively fed σ

provided $y \leq n$, and there exists a concept $c \in \mathcal{C}$ with $\sigma^+ \subseteq c$. Furthermore, we let $M_*(\sigma)$ denote $M_{|\sigma|-1}(\sigma)$.

Moreover, when learning a concept class from text, a major problem one has to deal with is avoiding or detecting *overgeneralization*. An overgeneralization occurs if the learner is guessing a superconcept of the target concept. Clearly, such an overgeneralized guess cannot be detected by using the incoming positive data only. Hence, one may be tempted to disallow overgeneralized guesses at all. Learners behaving thus are called conservative. Intuitively speaking a conservative IIM maintains its actual hypothesis at least as long as it has not seen data contradicting it. More formally, an IIM M is said to be *conservative* iff for all concepts c in the target class \mathcal{C} , all texts t for c and all $y, z \in \mathbb{N}$ the condition “if $M(t_y) \neq M(t_{y+z})$ then $t_{y+z}^+ \not\subseteq h_{M(t_y)}$ ” is fulfilled.

Another property of learners quite often found in the literature is consistency (cf., e.g., Wiehagen and Zeugmann [48, 49] and the references therein). A learner is called consistent if all its intermediate hypotheses correctly reflect the data seen so far. Formally, an IIM M is said to be *consistent* iff $t_x^+ \subseteq h_{M(t_x)}$ for all $x \in \mathbb{N}$ and every text t for every concept c in the target class \mathcal{C} .

Whenever one talks about the efficiency of learning besides the storage needed by the learner one has also to consider the time complexity of the learner. When talking about the time complexity of learning, it does not suffice to consider the time needed to compute the actual guess. What really counts in applications is the overall time needed until successful learning. Therefore, following Daley and Smith [10] we define the total learning time as follows.

Let \mathcal{C} be any concept class, and let M be any IIM that learns \mathcal{C} in the limit. Then, for every $c \in \mathcal{C}$ and every text t for c , let

$$\begin{aligned} \text{Conv}(M, t) &:= \text{the least number } m \in \mathbb{N}^+ \\ &\text{such that for all } n \geq m, M(t_n) = M(t_m) \end{aligned}$$

denote the *stage of convergence* of M on t (cf. [12]). Moreover, by $T_M(t_n)$ we denote the time to compute $M(t_n)$. We measure this time as a function of the length of the input and call it the *update time*. Finally, the total learning time taken by the IIM M on successive input t is defined as

$$TT(M, t) := \sum_{n=1}^{\text{Conv}(M, t)} T_M(t_n).$$

Clearly, if M does not learn the target concept from text t then the total learning time is *infinite*.

Two more remarks are in order here. First, it has been argued elsewhere that within the learning in the limit paradigm a learning algorithm is invoked only when the current hypothesis has some problem with the latest observed data. However, such

a viewpoint implicitly assumes that membership in the target concept is decidable in time polynomial in the length of the actual input. This may be not the case. Therefore, directly testing consistency would immediately lead to a non-polynomial update time provided membership is not known to be in \mathcal{P} .

Second, Pitt [31] addresses the question with respect to what parameter one should measure the total learning time. In the definition given above this parameter is the length of all examples seen so far. Clearly, now one could try to play with this parameter by waiting for a large enough input before declaring success. However, when dealing with the learnability of non-trivial concept classes, in the worst-case the total learning time will be anyhow unbounded. This effect is caused by the requirement to learn from all input sequences. Clearly, there are input sequences that start with many repetitions of data not containing enough information for successful learning. Thus, it does not make much sense to deal with the worst-case. Instead, we shall study the *expected* total learning time. In such a setting one cannot simply wait for long enough inputs. We shall then restrict ourselves to probability distributions that generate data sequences from which the target can be learned. Therefore, using the definition of total learning time given above seems to be reasonable.

Next, we define important concept classes which we are going to consider throughout this survey.

2.2. The Pattern Languages

Following Angluin [1] we define patterns and pattern languages as follows. Let $\mathcal{A} = \{0, 1, \dots\}$ be any finite alphabet containing at least two elements. Let $X = \{x_i \mid i \in \mathbb{N}\}$ be an infinite set of variables such that $\mathcal{A} \cap X = \emptyset$. *Patterns* are non-empty strings over $\mathcal{A} \cup X$, e.g., 01 , $0x_0111$, $1x_0x_00x_1x_2x_0$ are patterns. The length of a string $s \in \mathcal{A}^*$ and of a pattern π is denoted by $|s|$ and $|\pi|$, respectively. A pattern π is in *canonical form* provided that if k is the number of different variables in π then the variables occurring in π are precisely x_0, \dots, x_{k-1} . Moreover, for every j with $0 \leq j < k - 1$, the leftmost occurrence of x_j in π is left to the leftmost occurrence of x_{j+1} . The examples given above are patterns in canonical form. In the sequel we assume, without loss of generality, that all patterns are in canonical form. By *Pat* we denote the set of all patterns in canonical form.

If k is the number of different variables in π then we refer to π as to a *k-variable pattern*. By Pat_k we denote the set of all *k-variable patterns*. Furthermore, let $\pi \in Pat_k$, and let $u_0, \dots, u_{k-1} \in \mathcal{A}^+$; then we denote by $\pi[x_0/u_0, \dots, x_{k-1}/u_{k-1}]$ the string $w \in \mathcal{A}^+$ obtained by substituting u_j for each occurrence of x_j , $j = 0, \dots, k - 1$, in the pattern π . For example, let $\pi = 0x_01x_1x_0$. Then $\pi[x_0/10, x_1/01] = 01010110$. The tuple (u_0, \dots, u_{k-1}) is called a *substitution*. Furthermore, if $|u_0| = \dots = |u_{k-1}| = 1$, then we refer to (u_0, \dots, u_{k-1}) as to a *shortest substitution*. Let $\pi \in Pat_k$; we define the *language generated by pattern π* by

$$L(\pi) = \{\pi[x_0/u_0, \dots, x_{k-1}/u_{k-1}] \mid u_0, \dots, u_{k-1} \in \mathcal{A}^+\}.$$

By PAT_k we denote the set of all k -variable pattern languages. Finally, $PAT = \bigcup_{k \in \mathbb{N}} PAT_k$ denotes the set of all pattern languages over \mathcal{A} .

Furthermore, we let Q range over finite sets of patterns and define $L(Q) = \bigcup_{\pi \in Q} L(\pi)$, i.e., the union of all pattern languages generated by patterns from Q . Moreover, we use $Pat(k)$ and $PAT(k)$ to denote the family of all unions of at most k canonical patterns and the family of all unions of at most k pattern languages, respectively. That is, $Pat(k) = \{Q \mid Q \subseteq Pat, |Q| \leq k\}$ and $PAT(k) = \{L(Q) \mid Q \in Pat(k)\}$. Finally, let $L \subseteq \mathcal{A}^+$ be a language, and let $k \in \mathbb{N}^+$; we define $Club(L, k) = \{Q \mid |Q| \leq k, L \subseteq L(Q), (\forall Q')[Q' \subset Q \Rightarrow L \not\subseteq L(Q')]\}$. *Club* stands for consistent least upper bounds.

The pattern languages and variations thereof have been intensively investigated (cf., e.g., Salomaa [37, 38], and Shinohara and Arikawa [43] for an overview). Nix [29] as well as Shinohara and Arikawa [43] outlined interesting applications of pattern inference algorithms. For example, pattern language learning algorithms have been successfully applied for solving problems in molecular biology (cf., e.g., Shimozono *et al.* [39], Shinohara and Arikawa [43]).

As it turned out, pattern languages and finite unions of pattern languages are subclasses of Smullyan's [45] elementary formal systems (abbr. EFS). Arikawa *et al.* [5] have shown that EFS can also be treated as a logic programming language over strings. Recently, the techniques for learning finite unions of pattern languages have been extended to show the learnability of various subclasses of EFS (cf. Shinohara [42]). The investigations of the learnability of subclasses of EFSs are interesting because they yield corresponding results about the learnability of subclasses of logic programs. Hence, these results are also of relevance for Inductive Logic Programming (ILP) [28, 23, 8, 24]. Miyano *et al.* [26] intensively studied the polynomial-time learnability of EFSs.

Therefore, we may consider the learnability of pattern languages and of unions thereof as a nice test bed for seeing what kind of results one may obtain by considering the corresponding learning problems within the setting of learning in the limit.

3. Results Concerning Patterns

Within this section we ask whether or not the pattern languages and finite unions thereof can be learned efficiently. The principal learnability of the pattern languages from text with respect to the hypothesis space Pat has been established by Angluin [1]. However, her algorithm is based on computing descriptive patterns for the data seen so far. Here a pattern π is said to be descriptive (for the set S of strings contained in the input provided so far) if π can generate all strings contained in S and no other pattern with this property generates a proper subset of the language generated by π . Since no efficient algorithm is known for computing descriptive patterns, and finding a descriptive pattern of *maximum* length is \mathcal{NP} -hard, its update time is practically *intractable*.

There are also serious difficulties when trying to learn the pattern languages within the PAC model introduced by Valiant [46]. In the original model, the sample complexity depends exclusively on the VC dimension of the target concept class and the error and confidence parameters ε and δ , respectively. Recently, Mitchell *et al.* [25] have shown that even the class of all one-variable pattern languages has infinite VC dimension. Consequently, even this special subclass of PAT is not uniformly PAC learnable. Moreover, Schapire [40] has shown that pattern languages are not PAC learnable in the generalized model provided $\mathcal{P}/poly \neq \mathcal{NP}/poly$ with respect to every hypothesis space for PAT that is uniformly polynomially evaluable. Though this result highlights the difficulty of PAC learning PAT it has no clear application to the setting considered in this paper, since we aim to learn PAT with respect to the hypothesis space Pat . Since the membership problem for this hypothesis space is \mathcal{NP} -complete, it is *not* polynomially evaluable (cf. [1]).

In contrast, Kearns and Pitt [18] have established a PAC learning algorithm for the class of all k -variable pattern languages. Positive examples are generated with respect to arbitrary product distributions while negative examples are allowed to be generated with respect to any distribution. In their algorithm the length of substitution strings is required to be polynomially related to the length of the target pattern. Finally, they use as hypothesis space all unions of polynomially many patterns that have k or fewer variables². The overall learning time of their PAC learning algorithm is polynomial in the length of the target pattern, the bound for the maximum length of substitution strings, $1/\varepsilon$, $1/\delta$, and $|\mathcal{A}|$. The constant in the running time achieved depends *doubly exponential* on k , and thus, their algorithm becomes rapidly impractical when k increases.

Finally, Lange and Wiehagen [19] have proposed an inconsistent but iterative and conservative algorithm that learns PAT with respect to Pat . We shall study this algorithm below in some more detail.

But before doing it, we aim to figure out under which circumstances iterative learning of PAT is possible at all. A first answer is given by the following theorems from Case *et al.* [9]. Note that Pat is a non-redundant hypothesis space for PAT .

THEOREM 1 (CASE *et al.* [9]). *Let \mathcal{C} be any concept class, and let $\mathcal{H} = (h_j)_{j \in \mathbb{N}}$ be any non-redundant hypothesis space for \mathcal{C} . Then, every IIM M that $ItLimText_{\mathcal{H}}$ -learns \mathcal{C} is conservative.*

Proof. Recall that we use $M_*(\sigma)$ to denote $M_{|\sigma|-1}(\sigma)$ for any finite sequence $\sigma = x_0, \dots, x_n$ of elements from \mathcal{X} (cf. Definition 2).

Suppose the converse, i.e., there are a concept $c \in \mathcal{C}$, a text $t = (x_j)_{j \in \mathbb{N}} \in text(c)$, and a $y \in \mathbb{N}$ such that, for $j = M_*(t_y)$ and $k = M_*(t_{y+1}) = M(j, x_{y+1})$, both $j \neq k$ and $t_{y+1}^+ \subseteq h_j$ are satisfied. The latter implies $x_{y+1} \in h_j$, and thus we may

²More precisely, the number of allowed unions is at most $poly(|\pi|, s, 1/\varepsilon, 1/\delta, |\mathcal{A}|)$, where π is the target pattern, s the bound on the length on substitution strings, ε and δ are the usual error and confidence parameter, respectively, and \mathcal{A} is the alphabet of constants over which the patterns are defined.

consider the following text $\tilde{t} \in \text{text}(h_j)$. Let $\hat{t} = (\hat{x}_j)_{j \in \mathbb{N}}$ be any text for h_j and let $\tilde{t} = \hat{x}_0, x_{y+1}, \hat{x}_1, x_{y+1}, \hat{x}_2, \dots$. Since M has to learn h_j from \tilde{t} there must be a $z \in \mathbb{N}$ such that $M_*(\tilde{t}_{z+r}) = j$ for all $r \geq 0$. But $M_*(\tilde{t}_{2z+1}) = M(j, x_{y+1}) = k$, a contradiction. \blacksquare

Next, we point to another peculiarity of *PAT*, i.e., it meets the superset condition defined as follows. Let \mathcal{C} be any indexable class. \mathcal{C} meets the *superset condition* if, for all $c, c' \in \mathcal{C}$, there is some $\hat{c} \in \mathcal{C}$ being a superset of both c and c' .

THEOREM 2. (CASE *et al.* [9]). *Let \mathcal{C} be any indexable class meeting the superset condition, and let $\mathcal{H} = (h_j)_{j \in \mathbb{N}}$ be any non-redundant hypothesis space for \mathcal{C} . Then, every consistent IIM M that *ItLimText \mathcal{H}* -learns \mathcal{C} may be used to decide the inclusion problem for \mathcal{H} .*

Proof. Let \mathcal{X} be the underlying learning domain, and let $(w_j)_{j \in \mathbb{N}}$ be an effective enumeration of all elements in \mathcal{X} . Then, for every $i \in \mathbb{N}$, $t^i = (x_j^i)_{j \in \mathbb{N}}$ is the following computable text for h_i . Let z be the least index such that $w_z \in h_i$. Recall that, by definition, $h_i \neq \emptyset$, since \mathcal{H} is an indexed family, and thus w_z must exist. Then, for all $j \in \mathbb{N}$, we set $x_j^i = w_j$, if $w_j \in h_i$, and $x_j^i = w_z$, otherwise.

We claim that the following algorithm *Inc* decides, for all $i, k \in \mathbb{N}$, whether or not $h_i \subseteq h_k$.

Algorithm *Inc*: “On input $i, k \in \mathbb{N}$ do the following:

Determine the least $y \in \mathbb{N}$ with $i = M_*(t_y^i)$. Test whether or not $t_y^{i,+} \subseteq h_k$.

In case it is, output ‘Yes,’ and stop. Otherwise, output ‘No,’ and stop.”

Clearly, since \mathcal{H} is an indexed family and t^i is a computable text, *Inc* is an algorithm. Moreover, M learns h_i on every text for it, and \mathcal{H} is a non-redundant hypothesis space. Hence, M has to converge on text t^i to i , and therefore *Inc* has to terminate.

It remains to verify the correctness of *Inc*. Let $i, k \in \mathbb{N}$.

Clearly, if *Inc* outputs ‘No,’ a string $s \in h_i \setminus h_k$ has been found, and $h_i \not\subseteq h_k$ follows.

Next, consider the case that *Inc* outputs ‘Yes.’ Suppose to the contrary that $h_i \not\subseteq h_k$. Then, there is some $s \in h_i \setminus h_k$. Now, consider M when fed the text $t = t_y^i \diamond t^k$. Since $t_y^{i,+} \subseteq h_k$, t is a text for h_k . Since M learns h_k , there is some $r \in \mathbb{N}$ such that $k = M_*(t_y^i \diamond t_r^k)$. By assumption, there are some $\hat{c} \in \mathcal{C}$ with $h_i \cup h_k \subseteq \hat{c}$, and some text \hat{t} for \hat{c} having the initial segment $t_y^i \diamond s \diamond t_r^k$. By Theorem 1, M is conservative. Since $s \in h_i$ and $i = M_*(t_y^i)$, we obtain $M_*(\hat{t}_{y+1}) = M(i, s) = i$. Consequently, $M_*(t_y^i \diamond s \diamond t_r^k) = M_*(t_y^i \diamond t_r^k)$. Finally, since $s \in \hat{t}_{y+r+2}^+$, $k = M_*(t_y^i \diamond t_r^k)$, and $s \notin h_k$, M fails to consistently learn \hat{c} from text \hat{t} , a contradiction. This proves the theorem. \blacksquare

Taking into account that the inclusion problem for Pat is undecidable (cf. Jiang *et al.* [17]) and that PAT meets the superset condition, since $L(x_0) = \mathcal{A}^+$, by Theorem 2, we immediately arrive at the following corollary.

COROLLARY 3 (CASE *et al.* [9]). *If an IIM M $ItLimTxt_{Pat}$ -learns PAT then M is inconsistent.*

As a matter of fact, the latter corollary generalizes to *all* non-redundant hypothesis spaces for PAT . All the ingredients to prove this can be found in Zeugmann *et al.* [53]. Consequently, if one wishes to learn the pattern languages or unions of pattern languages iteratively, then either redundant hypothesis spaces or inconsistent learners cannot be avoided.

As for unions, the first result goes back to Shinohara [41] who proved the class of all unions of at most two pattern languages to be in $LimTxt_{Pat(2)}$. Wright [50] extended this result to $PAT(k) \in LimTxt_{Pat(k)}$ for all $k \geq 1$. Moreover, Theorem 4.2 in Shinohara and Arimura's [44] together with a lemma from Blum and Blum [6] shows that $\bigcup_{k \in \mathbb{N}} PAT(k)$ is not $LimTxt_{\mathcal{H}}$ -learnable for every hypothesis space \mathcal{H} .

The iterative learnability of $PAT(k)$ has been established by Case *et al.* [9]. Our learner is also consistent. Thus, the hypothesis space used had to be designed to be *redundant*. We only sketch the proof here.

THEOREM 4.

- (1) $Club(L, k)$ is finite for all $L \subseteq \mathcal{A}^+$ and all $k \in \mathbb{N}^+$,
- (2) If $L \in PAT(k)$, then $Club(L, k)$ is non-empty and contains a set Q , such that $L(Q) = L$.

Proof. Part (2) is obvious. Part (1) is easy for finite L . For infinite L , it follows from the lemma below.

LEMMA 1. *Let $k \in \mathbb{N}^+$, let $L \subseteq \mathcal{A}^+$ be any language, and suppose $t = (s_j)_{j \in \mathbb{N}} \in text(L)$. Then,*

- (1) $Club(t_0^+, k)$ can be obtained effectively from s_0 , and $Club(t_{n+1}^+, k)$ is effectively obtainable from $Club(t_n^+, k)$ and s_{n+1}
(* note the iterative nature *).
- (2) The sequence $Club(t_0^+, k), Club(t_1^+, k), \dots$ converges to $Club(L, k)$.

Putting it all together, one directly gets the following theorem.

THEOREM 5. *For all $k \geq 1$, $PAT(k) \in ItLimTxt$.*

Proof. Let $can(\cdot)$, be some computable bijection from finite classes of finite sets of patterns onto \mathbb{N} . Let pad be a 1–1 padding function such that, for all $x, y \in \mathbb{N}$, $W_{pad(x,y)} = W_x$. For a finite class \mathcal{S} of sets of patterns, let $g(\mathcal{S})$ denote a grammar obtained, effectively from \mathcal{S} , for $\bigcap_{Q \in \mathcal{S}} L(Q)$.

Let $L \in PAT(k)$, and let $t = (s_j)_{j \in \mathbb{N}} \in text(L)$. The desired IIM M is defined as follows. We set

$$M_0(t) = M(s_0) = \text{pad}(g(\text{Club}(t_0^+, k)), \text{can}(\text{Club}(t_0^+, k))), \text{ and for all } n > 0,$$

$$\begin{aligned} M_{n+1}(t) &= M(M_n(t), s_{n+1}) \\ &= \text{pad}(g(\text{Club}(t_{n+1}^+, k)), \text{can}(\text{Club}(t_{n+1}^+, k))). \end{aligned}$$

Using Lemma 1 it is easy to verify that $M_{n+1}(t) = M(M_n(t), s_{n+1})$ can be obtained effectively from $M_n(t)$ and s_{n+1} . Therefore, M *ItLimTxt*-identifies $PAT(k)$. \blacksquare

So far, the general theory provided substantial insight into the iterative learnability of the pattern languages. But still, we do not know anything about the number of examples needed until successful learning and the total amount of time to process them. Therefore, we address this problem in the following subsection.

3.1. Stochastic Finite Learning

As we have already mentioned, it does not make much sense to study the worst-case behavior of learning algorithms with respect to their total learning time. The reason for this phenomenon should be clear, since an arbitrary text may provide the information needed for learning very late. Therefore, in the following we always assume a class \mathcal{D} of admissible probability distributions over the relevant learning domain. Ideally, this class should be parameterized. Then, the data fed to learner are generated randomly with respect to one of the probability distributions from the class \mathcal{D} of underlying probability distributions. Furthermore, we introduce a random variable $CONV$ for the stage of convergence. Note that $CONV$ can be also interpreted as the total number of examples read by the IIM M until convergence. We therefore also refer to $CONV$ as to the *sample complexity*. The first major step to be performed consists now in determining the expectation $E[CONV]$. Clearly, $E[CONV]$ should be finite for all concepts $c \in \mathcal{C}$ and all distributions $D \in \mathcal{D}$. Second, one has to deal with tail bounds for $E[CONV]$. The easiest way to perform this step is to use Markov's inequality, i.e., we always know that

$$\Pr(CONV \geq t \cdot E[CONV]) \leq \frac{1}{t} \text{ for all } t \in \mathbb{N}^+.$$

However, quite often one can obtain much better tail bounds. If the underlying learner is known to be *conservative* and *rearrangement-independent* we always get *exponentially* shrinking tail bounds. A learner is said to be rearrangement-independent if its output depends exclusively on the range and length of its input (cf. [21] and the references therein). These tail bounds are established by the following theorem.

THEOREM 6 (ROSSMANITH AND ZEUGMANN [36]). *Let $CONV$ be the sample complexity of a conservative and rearrangement-independent learning algorithm. Then*

$$\Pr(CONV \geq 2t \cdot E[CONV]) \leq 2^{-t} \text{ for all } t \in \mathbb{N}.$$

Proof. First, recall the definition of *median*. If X is a random variable then μX is a median of X iff

$$\Pr(X \geq \mu X) \geq 1/2 \text{ and } \Pr(X \leq \mu X) \geq 1/2.$$

A nonempty set of medians exists for each random variable and consists either of a single real number or of a closed real interval. We will denote the smallest median of X by μX , since this choice gives the best upper bounds.

Now, we can establish the following claim.

Claim 1. *Let X be the sample complexity of a conservative and rearrangement-independent learning algorithm. Then $\Pr(X \geq t \cdot \mu X) \leq 2^{-t}$ for all $t \in \mathbb{N}$.*

We divide the text s_0, s_1, \dots into blocks of length μX . The probability that the algorithm converges after reading any of the blocks is then at least $1/2$. Since the algorithm is rearrangement-independent the order of the blocks does not matter and since the algorithm is conservative it does not change its hypothesis after computing once the right hypothesis. This proves Claim 1.

Claim 2. $\mu X \leq 2E[X]$ for every positive random variable X .

Claim 2 is a direct consequence of the Markov inequality and the definition of median.

Putting Claims 1 and 2 together directly yields the theorem. ■

Theorem 6 puts the importance of rearrangement-independent and conservative learners into the right perspective. As long as the learnability of indexed families is concerned, these results have a wide range of potential applications, since every conservative learner can be transformed into a learner that is both conservative and rearrangement-independent provided the hypothesis space is appropriately chosen (cf. Lange and Zeugmann [21]).

Furthermore, since the distribution of $CONV$ decreases geometrically for all conservative and rearrangement-independent learning algorithms, all higher moments of $CONV$ exist in this case, too. Thus, instead of applying Theorem 6 directly, one can hope for further improvements by applying even sharper tail bounds using, for example, Chebyshev's inequality.

Additionally, the learner takes a confidence parameter δ as input. But in contrast to learning in the limit, the learner itself decides how many examples it wants to read. Then it computes a hypothesis, outputs it and stops. The hypothesis output is correct for the target with probability at least $1 - \delta$.

The explanation given so far explains how it works, but not why it does. Intuitively, the stochastic finite learner simulates the limit learner until an upper bound for twice the expected total number of examples needed until convergence has been met. Assuming this to be true, by Markov's inequality the limit learner has now

converged with probability $1/2$. All what is left, is to decrease the probability of failure. This is done by using the tail bounds for $CONV$. Applying Theorem 6, one easily sees that increasing the sample complexity by a factor of $O(\log \frac{1}{\delta})$ results in a probability of $1 - \delta$ for having reached the stage of convergence. If Theorem 6 is not applicable, one can still use Markov's inequality but then the sample complexity needed will increase by a factor of $1/\delta$.

It remains to explain how the stochastic finite learner can calculate the upper bound for $E[CONV]$. This is precisely the point where we need the parameterization of the class \mathcal{D} of underlying probability distributions. Since in general, it is not known which distribution from \mathcal{D} has been chosen, one has to assume a bit of *prior knowledge* or *domain knowledge* provided by suitable upper and/or lower bounds for the parameters involved. A more serious difficulty is to incorporate the unknown target concept into this estimate. This step depends on the concrete learning problem on hand, and requires some extra effort. We shall exemplify it below.

Now we are ready to formally define stochastic finite learning.

DEFINITION 3 ([33, 34, 36]). Let \mathcal{D} be a set of probability distributions on the learning domain, \mathcal{C} a concept class, \mathcal{H} a hypothesis space for \mathcal{C} , and $\delta \in (0, 1)$. $(\mathcal{C}, \mathcal{D})$ is said to be *stochastically finitely learnable with δ -confidence* with respect to \mathcal{H} iff there is an IIM M that for every $c \in \mathcal{C}$ and every $D \in \mathcal{D}$ performs as follows. Given any random data sequence θ for c generated according to D , M stops after having seen a finite number of examples and outputs a single hypothesis $h \in \mathcal{H}$. With probability at least $1 - \delta$ (with respect to distribution D) h has to be correct, that is $c = h$.

If stochastic finite learning can be achieved with δ -confidence for every $\delta > 0$ then we say that $(\mathcal{C}, \mathcal{D})$ can be learned stochastically finite *with high confidence*.

Note that there are subtle differences between our model and PAC learning. By its definition, stochastic finite learning is not completely distribution independent. A bit of *additional knowledge* concerning the underlying probability distributions is required. Thus, from that perspective, stochastic finite learning is weaker than the PAC-model. On the other hand, we do *not* measure the quality of the hypothesis with respect to the underlying probability distribution. Instead, we require the hypothesis computed to be exactly correct with high probability. Note that exact identification with high confidence has been considered within the PAC paradigm, too (cf., e.g., Goldman *et al.* [13]). Conversely, we also can easily relax the requirement to learn *probably exactly correct* but whenever possible we shall not do it.

Furthermore, in the uniform PAC model as introduced in Valiant [46] the sample complexity depends exclusively on the VC dimension of the target concept class and the error and confidence parameters ε and δ , respectively. This model has been generalized by allowing the sample size to depend on the concept complexity, too (cf., e.g., Blumer *et al.* [7] and Haussler *et al.* [15]). Provided no upper bound for the concept complexity of the target concept is given, such PAC learners decide themselves

how many examples they wish to read (cf. [15]). This feature is also adopted to our setting of stochastic finite learning. However, all variants of PAC learning we are aware of require that all hypotheses from the relevant hypothesis space are uniformly polynomially evaluable. Though this requirement may be necessary in some cases to achieve (efficient) stochastic finite learning, it is not necessary in general as we shall see below.

Next, let us exemplify our model by looking at the concept class of all pattern languages. The results presented below have been obtained by Zeugmann [51] and Rossmanith and Zeugmann [36]. Our stochastic finite learner uses Lange and Wiehagen's [19] pattern language learner as a main ingredient. We consider here learning from positive data only.

Recall that every string of a particular pattern language is generated by at least one substitution. Therefore, it is convenient to consider probability distributions over the set of all possible substitutions. That is, if $\pi \in Pat_k$, then it suffices to consider any probability distribution D over $\underbrace{\mathcal{A}^+ \times \dots \times \mathcal{A}^+}_{k\text{-times}}$. For $(u_0, \dots, u_{k-1}) \in \mathcal{A}^+ \times \dots \times \mathcal{A}^+$ we denote by $D(u_0, \dots, u_{k-1})$ the probability that variable x_0 is substituted by u_0 , variable x_1 is substituted by u_1 , ..., and variable x_{k-1} is substituted by u_{k-1} .

In particular, we mainly consider a special class of distributions, i.e., *product distributions*. Let $k \in \mathbb{N}^+$, then the class of all product distributions for Pat_k is defined as follows. For each variable x_j , $0 \leq j \leq k-1$, we assume an arbitrary probability distribution D_j over \mathcal{A}^+ on substitution strings. Then we call $D = D_0 \times \dots \times D_{k-1}$ product distribution over $\mathcal{A}^+ \times \dots \times \mathcal{A}^+$, i.e., $D(u_0, \dots, u_{k-1}) = \prod_{j=0}^{k-1} D_j(u_j)$. Moreover, we call a product distribution *regular* if $D_0 = \dots = D_{k-1}$. Throughout this paper, we restrict ourselves to deal with *regular* distributions. We therefore use d to denote the distribution over \mathcal{A}^+ on substitution strings, i.e., $D(u_0, \dots, u_{k-1}) = \prod_{j=0}^{k-1} d(u_j)$. We call a regular distribution *admissible* if $d(a) > 0$ for at least two different elements $a \in \mathcal{A}$. As a special case of an admissible distribution we consider the *uniform* distribution over \mathcal{A}^+ , i.e., $d(u) = 1/(2 \cdot |\mathcal{A}|)^\ell$ for all strings $u \in \mathcal{A}^+$ with $|u| = \ell$. Note that here only strings of equal length have the same probability and not each elementary event.

We will express all estimates with the help of the following parameters: $E[\Lambda]$, α and β , where Λ is a random variable for the length of the examples drawn. α and β are defined below. To get concrete bounds for a concrete implementation one has to obtain c from the algorithm and has to compute $E[\Lambda]$, α , and β from the admissible probability distribution D . Let u_0, \dots, u_{k-1} be independent random variables with distribution d for substitution strings. Whenever the index i of u_i does not matter, we simply write u or u' .

The two parameters α and β are now defined via d . First, α is simply the probability that u has length 1, i.e.,

$$\alpha = \Pr(|u| = 1) = \sum_{a \in \mathcal{A}} d(a).$$

Second, β is the conditional probability that two random strings that get substituted into π are identical under the condition that both have length 1, i.e.,

$$\beta = \Pr(u = u' \mid |u| = |u'| = 1) = \sum_{a \in \mathcal{A}} d(a)^2 / \left(\sum_{a \in \mathcal{A}} d(a) \right)^2.$$

Note that we have *omitted* the assumption of a text to exhaust the target language. Instead, we only demand the data sequence fed to the learner to contain “enough” information to recognize the target pattern. The meaning of “enough” is mainly expressed by the parameter α .

The model of computation as well as the representation of patterns we assume is the same as in Angluin [1]. In particular, we assume a random access machine that performs a reasonable menu of operations each in unit time on registers of length $O(\log n)$ bits, where n is the input length.

Lange and Wiehagen’s [19] algorithm (abbr. LWA) works as follows. Let h_n be the hypothesis computed after reading s_0, \dots, s_n , i.e., $h_n = M(s_0, \dots, s_n)$. Then $h_0 = s_0$ and for all $n \geq 1$:

$$h_n = \begin{cases} h_{n-1}, & \text{if } |h_{n-1}| < |s_n| \\ s_n, & \text{if } |h_{n-1}| > |s_n| \\ h_{n-1} \cup s_n, & \text{if } |h_{n-1}| = |s_n| \end{cases}$$

The algorithm computes the new hypothesis only from the latest example and the old hypothesis. If the latest example is longer than the old hypothesis, the example is ignored, i.e., the hypothesis does not change. If the latest example is shorter than the old hypothesis, the old hypothesis is ignored and the new example becomes the new hypothesis. If, however, $|h_{n-1}| = |s_n|$ the new hypothesis is the *union* of h_{n-1} and s_n . In order to explain the union, we need the following notation. Let $\pi \in Pat$, $1 \leq i \leq |\pi|$; we use $\pi(i)$ to denote the i -th symbol in π . Now, the union $\varrho = \pi \cup s$ of a canonical pattern π and a string s of the same length is defined as

$$\varrho(i) = \begin{cases} \pi(i), & \text{if } \pi(i) = s(i) \\ x_j, & \text{if } \pi(i) \neq s(i) \ \& \ \exists k < i : [\varrho(k) = x_j, s(k) = s(i), \\ & \quad \pi(k) = \pi(i)] \\ x_m, & \text{otherwise, where } m = \#var(\varrho(1) \dots \varrho(i-1)) \end{cases}$$

where $\varrho(0) = \lambda$ for notational convenience. Note that the resulting pattern is again canonical.

If the target pattern does not contain any variable then the LWA converges after having read the first example. Hence, this case is trivial and we therefore assume in the following always $k \geq 1$, i.e., the target pattern has to contain at least one variable.

Figure 1 displays the union operation for $\pi = 01x_0x_121x_0x_201x_0x_1$ and $s = 120021010212$. Since the letters in the first column are different and there is no previous column, $\varrho(1) = x_0$. The letters in the second column are different, and the second column is not equal to the first column, so $\varrho(2) = x_1$. Next, $\pi(3) = x_0$ and $\pi(4) = x_1$, and thus ϱ must also contain different variables at positions 3 and 4. Consequently, these variables get renamed, i.e., $\varrho(3) = x_2$ and $\varrho(4) = x_3$. The letters in the 5th and 6th column are identical, hence $\varrho(5) = 2$ and $\varrho(6) = 1$ (cf. the first case in the definition of the union operation). In the 7th column, we have x_0 and 0 and this column is equal to the third column. Therefore, the second case in the definition of the union operation applies and $\varrho(7) = x_2$. Now, $\varrho(8) = x_4$ and $\varrho(9) = 0$ are obvious. The 10th column is identical to the second one, thus $\varrho(10) = x_1$. Next, we have x_0 and 1 while both the third and 7th column contain x_0 and 0. Therefore, a new variable has to be introduced and $\varrho(11) = x_5$ (cf. the third case in the definition of the union operation). Analogously, the x_1 in the 12th column has to be distinguished from the x_1 in 4th column resulting in $\varrho(12) = x_6$.

π	0	1	x_0	x_1	2	1	x_0	x_2	0	1	x_0	x_1
s	1	2	0	0	2	1	0	1	0	2	1	2
$\varrho = \pi \cup s$	x_0	x_1	x_2	x_3	2	1	x_2	x_4	0	x_1	x_5	x_6

Figure 1

Our next theorem analyzes the complexity of the union operation.

THEOREM 7 (Rossmannith and Zeugmann [36]). *The union operation can be computed in linear time.*

Furthermore, the following bound for the stage of convergence for every target pattern from Pat_k can be shown.

THEOREM 8 (Rossmannith and Zeugmann [36]).
 $E[CONV] = O\left(\frac{1}{\alpha^k} \cdot \log_{1/\beta}(k)\right)$ for all $k \geq 2$.

Hence, by Theorem 7, the expected total learning time can be estimated by
 $E[TT] = O\left(\frac{1}{\alpha^k} E[\Lambda] \log_{1/\beta}(k)\right)$ for all $k \geq 2$.

For a better understanding of the bound obtained we evaluate it for the uniform distribution and compare it to the minimum number of examples needed for learning a pattern language via the LWA.

THEOREM 9 (Rossmannith and Zeugmann [36]). *To learn a pattern $\pi \in Pat_k$, $k \geq 2$, from texts randomly generated with respect to the uniform distribution, the LWA has the expected total learning time $E[TT] = O(2^k |\pi| \log_{|\mathcal{A}|}(k))$.*

THEOREM 10 (Zeugmann [51]). *To learn a pattern $\pi \in Pat_k$ the LWA needs exactly $\lfloor \log_{|\mathcal{A}|}(|\mathcal{A}| + k - 1) \rfloor + 1$ examples in the best case.*

The main difference between the two bounds just given is the factor 2^k which precisely reflects the time the LWA has to wait until it has seen the first shortest string from the target pattern language. Moreover, in the best-case the LWA is processing shortest examples only. Thus, we introduce M_C to denote the number of minimum length examples read until convergence. Then, one can show that

$$\mathbb{E}[M_C] \leq \frac{2 \ln(k) + 3}{\ln(1/\beta)} + 2.$$

Note that Theorem 8 is shown by using the bound for $\mathbb{E}[M_C]$ just given. More precisely, we have $\mathbb{E}[CONV] = (1/\alpha^k) \mathbb{E}[M_C]$. Now, we are ready to transform the LWA into a stochastic finite learner.

THEOREM 11 (Rossmannith and Zeugmann [36]). *Let $\alpha_*, \beta_* \in (0, 1)$. Assume \mathcal{D} to be a class of admissible probability distributions over \mathcal{A}^+ such that $\alpha \geq \alpha_*$, $\beta \leq \beta_*$ and $\mathbb{E}[\Lambda]$ is finite for all distributions $D \in \mathcal{D}$. Then (PAT, \mathcal{D}) is stochastically finitely learnable with high confidence from text.*

Proof. Let $D \in \mathcal{D}$, and let $\delta \in (0, 1)$ be arbitrarily fixed. Furthermore, let $t = s_0, s_1, s_2, \dots$ be any randomly generated text with respect to D for the target pattern language. The wanted learner M uses the LWA as a subroutine. Additionally, it has a counter for memorizing the number of examples already seen. Now, we exploit the fact that the LWA produces a sequence $(\tau_n)_{n \in \mathbb{N}}$ of hypotheses such that $|\tau_n| \geq |\tau_{n+1}|$ for all $n \in \mathbb{N}$.

The learner runs the LWA until for the first time C many examples have been processed, where

$$C = \left(\frac{1}{\alpha_*}\right)^{|\tau|} \cdot \left(\frac{2 \ln(|\tau|) + 3}{\ln(1/\beta_*)} + 2\right) \quad (\text{A})$$

and τ is the actual output made by the LWA.

Finally, in order to achieve the desired confidence, the learner sets $\gamma = \lceil \log \frac{1}{\delta} \rceil$ and runs the LWA for a total of $2 \cdot \gamma \cdot C$ examples. This is the reason we need the counter for the number of examples processed. Now, it outputs the last hypothesis τ produced by the LWA, and stops thereafter.

Clearly, the learner described above is finite. Let L be the target language and let $\pi \in Pat_k$ be the unique pattern such that $L = L(\pi)$. It remains to argue that $L(\pi) = L(\tau)$ with probability at least $1 - \delta$.

First, the bound in (A) is an upper bound for the expected number of examples needed for convergence by the LWA that has been established in Theorem 8 (via the reformulation using $\mathbb{E}[M_C]$ given above). On the one hand, this follows from our assumptions about the allowed α and β as well as from the fact that $|\tau| \geq |\pi|$ for every hypothesis output. On the other hand, the learner does not know k , but the estimate $\#var(\pi) \leq |\pi|$ is sufficient. Note that we have to use in (A) the bound

for $E[M_C]$ given above, since the target pattern may contain zero or one different variables.

Therefore, after having processed C many examples the LWA has already converged on average. The desired confidence is then an immediate consequence of Theorem 6. \blacksquare

The latter theorem allows a nice corollary which we state next. Making the same assumption as done by Kearns and Pitt [18], i.e., assuming the additional prior knowledge that the target pattern belongs to Pat_k , the complexity of the stochastic finite learner given above can be considerably improved. The resulting learning time is linear in the expected string length, and the constant depending on k grows only exponentially in k in contrast to the doubly exponentially growing constant in Kearns and Pitt's [18] algorithm. Moreover, in contrast to their learner, our algorithm learns from positive data only, and outputs a hypothesis that is correct for the target language with high probability.

Again, for the sake of presentation we shall assume $k \geq 2$. Moreover, if the prior knowledge $k = 1$ is available, then there is also a much better stochastic finite learner for PAT_1 (cf. [34]).

COROLLARY 12. *Let $\alpha_*, \beta_* \in (0, 1)$. Assume \mathcal{D} to be a class of admissible probability distributions over \mathcal{A}^+ such that $\alpha \geq \alpha_*$, $\beta \leq \beta_*$ and $E[\Lambda]$ is finite for all distributions $D \in \mathcal{D}$. Furthermore, let $k \geq 2$ be arbitrarily fixed. Then there exists a learner M such that*

- (1) M learns (PAT_k, \mathcal{D}) stochastically finitely with high confidence from text, and
- (2) The running time of M is $O((1/\alpha_*^k)E[\Lambda] \log_{1/\beta_*}(k) \log_2(1/\delta))$.
 (* Note that $1/\alpha_*^k$ and $\log_{1/\beta_*}(k)$ now are constants. *)

This finishes our exposition concerning the pattern languages and unions thereof.

In the following section we show our ideas to be beneficial for a class known to be PAC learnable by looking at the class of all concepts describable by a monomial.

4. Learning Conjunctive Concepts

In this section we exemplify the general approach outlined above by using the class of all concepts describable by a monomial. For all details omitted the reader is referred to Reischuk and Zeugmann [33].

For defining the classes of concepts we shall deal with in this section, let $\mathcal{L}_n = \{x_1, \bar{x}_1, x_2, \bar{x}_2, \dots, x_n, \bar{x}_n\}$ be a set of literals. x_i is a *positive* literal and \bar{x}_i a *negative* one. A conjunction of literals defines a *monomial*. For a monomial m let $\#(m)$ denote its length, that is the number of literals in it.

A monomial m describes a subset $L(m)$ of \mathcal{X}_n , in other words a concept, in the obvious way: the concept contains exactly those binary vectors for which the

monomial evaluates to 1, that is $L(m) := \{b \in \mathcal{X}_n \mid m(b) = 1\}$. The collection of objects we are going to learn is the set \mathcal{C}_n of all concepts that are describable by monomials over \mathcal{X}_n . There are two trivial concepts, the empty subset and \mathcal{X}_n itself. \mathcal{X}_n , which will also be called “TRUE”, can be represented by the empty monomial. The concept “FALSE” has several descriptions. To avoid ambiguity, we always represent “FALSE” by the monomial $x_1\bar{x}_1 \dots x_n\bar{x}_n$. Furthermore, we often identify the set of all monomials over \mathcal{L}_n and the concept class \mathcal{C}_n . Note that $|\mathcal{C}_n| = 3^n + 1$.

For the concept class \mathcal{C}_n we choose as hypothesis space the set of all monomials over \mathcal{L}_n and denote this hypothesis space by \mathcal{H}_n . We shall distinguish learning from positive data only and learning from both positive and negative data. Note that when considering learning from positive data only, one cannot decide whether or not the learner has already converged. When learning from positive and negative data is considered, the stage of convergence is *decidable*, but one would have to read the data sequence until all Boolean vectors did appear. Thus, for any interesting n , decidability is practically infeasible.

The learner used is essentially Haussler’s [14] Wholist algorithm. We present it here as an iterative learner. The iterative learner is defined in stages, where Stage ℓ conceptually describes M_ℓ .

Algorithm \mathcal{IML} : “Let $c \in \mathcal{C}_n$, let $i = (b_0, m(b_0)), (b_1, m(b_1)), \dots$ be any informant for c . Go to Stage 0.

Stage 0. \mathcal{IML} receives as input $(b_0, m(b_0))$.

Initialize $h_{ini} = x_1\bar{x}_1 \dots x_n\bar{x}_n$.

If $m(b_0) = 0$ then h_{ini} remains unchanged; else

for $j := 1$ to n do

if $b_0^j = 1$ then delete \bar{x}_j in h_{ini} else delete x_j in h_{ini} .

Denote the result by h_0 , output h_0 and go to Stage 1.

Stage ℓ , $\ell \geq 1$. \mathcal{IML} receives as input $h_{\ell-1}$ and the $(\ell + 1)$ th element $(b_\ell, m(b_\ell))$ of i .

If $m(b_\ell) = 0$ then set $h_\ell = h_{\ell-1}$; else

for $j := 1$ to n do

if $b_\ell^j = 1$ then delete \bar{x}_j in $h_{\ell-1}$ else delete x_j in $h_{\ell-1}$.

Denote the result by h_ℓ , output h_ℓ and go to Stage $\ell + 1$.

By convention, if all literals have been removed, then $h_\ell = \emptyset$, and $h_\ell(b) = 1$ for all $b \in \mathcal{X}_n$.

end.”

The following example illustrates Algorithm \mathcal{IML} . Let $n = 7$, and let $m = x_1\bar{x}_2x_4x_7$ be the target monomial. Suppose the input sequence to start as follows: $\langle 1001111, 1, 0110110, 0, 1011101, 1, 1011001, 1, \dots \rangle$. In Stage 0 the loop is executed and all literals that do not evaluate to 1 are removed from h_{ini} . Thus, the hypothesis computed is $h_0 = x_1\bar{x}_2\bar{x}_3x_4x_5x_6x_7$. Next, h_0 and $0110110, 0$ are read and since the label is 0, the hypothesis output is again $x_1\bar{x}_2\bar{x}_3x_4x_5x_6x_7$, i.e., $h_1 = h_0$. In Stage 2, \mathcal{IML} receives h_1 and $1011101, 1$ and executes the loop resulting in $h_2 = x_1\bar{x}_2x_4x_5x_7$. Now, \mathcal{IML} reads h_2 and $1011001, 1$ and executes the loop in Stage 3. The result is $h_3 = x_1\bar{x}_2x_4x_7$ which equals the target monomial. Consequently, Algorithm \mathcal{IML} has reached the stage of convergence.

Now, we can directly state the following theorem.

THEOREM 13. *Algorithm \mathcal{IML} $ItLimInf_{\mathcal{H}_n}$ -learns \mathcal{C}_n .*

Moreover, Algorithm \mathcal{IML} can be easily adapted to learn from positive data only. We have just to omit the tests whether or not $m(b_n) = 0$. We call the resulting algorithm \mathcal{IMLP} . Now, the following corollary is obvious.

COROLLARY 14. *Algorithm \mathcal{IMLP} $ItLimTxt_{\mathcal{H}_n}$ -learns \mathcal{C}_n .*

If the target monomial is the concept “FALSE”, then Algorithms \mathcal{IML} and \mathcal{IMLP} immediately converge. Thus, we call “FALSE” the minimal concept. If the target concept contains precisely n literals, then one positive example suffices. This positive example is unique. Thus, for these two cases everything is clear and therefore, we call these concepts *trivial*. Also, the set of probability distributions D on the set of positive examples for these trivial c are trivial, too.

To study the general case, let us call the literals appearing in a non-minimal monomial m *relevant*. All the other literals in \mathcal{L}_n will be called *irrelevant* for m . There are $2n - \#(m)$ irrelevant literals. Recall that $\#(m)$ denotes the number of literals in monomial m .

We call bit i *relevant* for monomial m if x_i or \bar{x}_i is relevant for m and use $\mathbf{k} := \mathbf{k}(m) = n - \#(m)$ to denote the number of irrelevant bits.

4.1. Learning Monomials from Positive Data

First, we consider learning from positive data. To avoid the trivial cases, we let $c = L(m) \in \mathcal{C}_n$ be a concept with monomial $m = \bigwedge_{j=1}^{\#(m)} \ell_{i_j}$ such that $k = k(m) = n - \#(m) > 0$. There are 2^k positive examples for c . For the sake of presentation, we assume these examples to be *binomially distributed*. That is, in a random positive example all entries corresponding to irrelevant bits are selected independently of each other. With some probability p this will be a 1, and with probability q a 0, where $q = 1 - p$. We shall consider only nontrivial distributions where $0 < p < 1$. Note that otherwise the data sequence does not contain all positive examples. We aim to compute the expected number of examples taken by \mathcal{IMLP} until convergence. Again

we use $CONV$ to denote a random variable counting the number of examples read till convergence.

The first example received forces \mathcal{IMLP} to delete precisely n of the $2n$ literals in h_{ini} . Thus, this example always plays a *special* role. Note that the resulting hypothesis h_0 depends on b_0 , but the number k of literals that remain to be deleted from h_0 until convergence is *independent* of b_0 . Using tail bound techniques, we can show the following theorem.

THEOREM 15 (Reischuk and Zeugmann [33]). *Let $c = L(m)$ be a non-minimal concept in \mathcal{C}_n , and let the positive examples for c be binomially distributed with parameter p . Define $\psi := \min\{\frac{1}{1-p}, \frac{1}{p}\}$ and $\tau := \max\{\frac{p}{1-p}, \frac{1-p}{p}\}$. Then the expected number of positive examples needed by algorithm \mathcal{IMLP} until convergence can be bounded by*

$$E[CONV] \leq \lceil \log_{\psi} k(m) \rceil + \tau + 2 .$$

Proof. Let $k = k(m)$ and let $q = 1-p$. The first positive example contains ν times a 1 and $k-\nu$ times a 0 with probability $\binom{k}{\nu} p^{\nu} q^{k-\nu}$ at the positions not corresponding to a literal in the target monomial m . Now, assuming any such vector, we easily see that h_0 contains ν positive irrelevant literals and $k-\nu$ negative irrelevant literals. Therefore, in order to achieve convergence, the Algorithm \mathcal{IMLP} now needs positive examples that contain at least one 0 for each positive irrelevant literal and at least one 1 for each negative irrelevant literal. Thus, the probability that at least one irrelevant literal survives μ subsequent positive examples is bounded by $\nu p^{\mu} + (k-\nu)q^{\mu}$. Therefore,

$$\Pr(CONV - 1 > \mu) \leq \sum_{\nu=0}^k \binom{k}{\nu} p^{\nu} q^{k-\nu} \cdot (\nu p^{\mu} + (k-\nu)q^{\mu}) .$$

Next, we derive a closed formula for the sum given above.

$$\text{Claim 1. } \sum_{\nu=0}^k \binom{k}{\nu} p^{\nu} q^{k-\nu} \cdot \nu = kp \quad \text{and} \quad \sum_{\nu=0}^k \binom{k}{\nu} p^{\nu} q^{k-\nu} \cdot (k-\nu) = kq$$

The first equality can be shown as follows.

$$\begin{aligned} \sum_{\nu=0}^k \binom{k}{\nu} p^{\nu} q^{k-\nu} \cdot \nu &= \sum_{\nu=1}^k \binom{k}{\nu} p^{\nu} q^{k-\nu} \cdot \nu \\ &= \sum_{\nu=0}^{k-1} \binom{k}{\nu+1} p^{\nu+1} q^{k-1-\nu} \cdot (\nu+1) \\ &= \sum_{\nu=0}^{k-1} k \cdot \binom{k-1}{\nu} p^{\nu+1} q^{k-(\nu+1)} \\ &= kp \cdot \sum_{\nu=0}^{k-1} \binom{k-1}{\nu} p^{\nu} q^{(k-1)-\nu} \\ &= kp \cdot (p+q)^{k-1} = kp . \end{aligned}$$

The other equality can be proved analogously, which yields Claim 1.

Now, proceeding as above, we obtain

$$\begin{aligned}
\mathbb{E}[\text{CONV} - 1] &= \sum_{\mu=0}^{\infty} \Pr(\text{CONV} - 1 > \mu) \\
&\leq \lambda + \sum_{\mu=\lambda}^{\infty} \sum_{\nu=0}^k \binom{k}{\nu} p^{\nu} q^{k-\nu} \cdot (\nu p^{\mu} + (k - \nu) q^{\mu}) \\
&= \lambda + \sum_{\mu=\lambda}^{\infty} \sum_{\nu=0}^k \binom{k}{\nu} p^{\nu} q^{k-\nu} \cdot \nu p^{\mu} + \sum_{\mu=\lambda}^{\infty} \sum_{\nu=0}^k \binom{k}{\nu} p^{\nu} q^{k-\nu} \cdot (k - \nu) q^{\mu} \\
&= \lambda + \sum_{\mu=\lambda}^{\infty} p^{\mu} \cdot \underbrace{\sum_{\nu=0}^k \binom{k}{\nu} p^{\nu} q^{k-\nu} \nu}_{=kp \text{ by Claim 1}} + \sum_{\mu=\lambda}^{\infty} q^{\mu} \cdot \underbrace{\sum_{\nu=0}^k \binom{k}{\nu} p^{\nu} q^{k-\nu} (k - \nu)}_{=kq \text{ by Claim 1}} \\
&= \lambda + kp \cdot \sum_{\mu=\lambda}^{\infty} p^{\mu} + kq \cdot \sum_{\mu=\lambda}^{\infty} q^{\mu} = \lambda + k \cdot \left(\frac{p}{q} \cdot p^{\lambda} + \frac{q}{p} \cdot q^{\lambda} \right) \\
&\leq \lambda + k \cdot \left(\frac{p}{q} \cdot \psi^{-\lambda} + \frac{q}{p} \cdot \psi^{-\lambda} \right) \\
&\leq \lambda + k \psi^{-\lambda} \cdot (1 + \tau) .
\end{aligned}$$

Finally, choosing $\lambda = \lceil \log_{\psi} k \rceil$ gives the statement of the theorem as an easy calculation shows. \blacksquare

Now, taking into account that τ does not depend on the dimension n of the learning domain $\{0, 1\}^n$, we can easily determine the expected total learning time.

COROLLARY 16 (Reischuk and Zeugmann [33]). *For every binomially distributed text with parameter $0 < p < 1$ the average total learning time of Algorithm \mathcal{IMLP} for concepts in \mathcal{C}_n is at most $O(n \log n)$.*

Also, Algorithm \mathcal{IMLP} possesses the two favorable properties needed to apply Theorem 6, i.e., it is rearrangement-independent and conservative. Thus, we can conclude

$$\Pr(\text{CONV} > 2t \cdot \mathbb{E}[\text{CONV}]) \leq 2^{-t} \text{ for all } t \in \mathbb{N} .$$

Next, we turn our attention to the design of a stochastic finite learner. We study the case that the positive examples are binomially distributed with parameter p . But we do not require precise knowledge about the underlying distribution. Instead, we reasonably assume that *prior knowledge* is provided by parameters p_{low} and p_{up} such that $p_{low} \leq p \leq p_{up}$ for the true parameter p . Binomial distributions fulfilling this requirement are called **(p_{low}, p_{up}) -admissible distributions**. Let $\mathcal{D}_n[p_{low}, p_{up}]$ denote the set of such distributions on \mathcal{X}_n .

If bounds p_{low} and p_{up} are available, the Algorithm \mathcal{IMLP} can be transformed into a stochastic finite learner inferring all concepts from \mathcal{C}_n with high confidence.

THEOREM 17 (Reischuk and Zeugmann [33]). *Let $0 < p_{low} \leq p_{up} < 1$ and $\psi := \min\{\frac{1}{1-p_{low}}, \frac{1}{p_{up}}\}$. Then $(\mathcal{C}_n, \mathcal{D}_n[p_{low}, p_{up}])$ is stochastically finitely learnable with high confidence from text. To achieve δ -confidence no more than $O(\log_2 1/\delta \cdot \log_\psi n)$ many examples are necessary.*

Proof. The stochastic finite learner is based on Algorithm \mathcal{IMLP} and a counter for the number of examples already processed. We set

$$\tau_{\max} = \left\lceil \max \left\{ \frac{p_{low}}{1-p_{low}}, \frac{1-p_{low}}{p_{low}}, \frac{p_{up}}{1-p_{up}}, \frac{1-p_{up}}{p_{up}} \right\} \right\rceil.$$

If Algorithm \mathcal{IMLP} is run for $\vartheta := \lceil \log_\psi n \rceil + \tau_{\max} + 2$ many examples, Theorem 15 implies that ϑ is at least as large as the expected convergence stage $E[CONV]$.

In order to achieve the desired confidence, the learner sets $\gamma := \lceil \log \frac{1}{\delta} \rceil$ and runs Algorithm \mathcal{IMLP} for a total of $2 \gamma \cdot \vartheta$ examples. This is the reason why we need a counter for the number of examples processed. The algorithm outputs the last hypothesis $h_{2 \gamma \cdot \vartheta}$ produced by Algorithm \mathcal{IMLP} and stops thereafter. The reliability follows from the tail bounds established in Theorem 6. \blacksquare

4.2. Learning Monomials from Informant

Finally, we ask how the results obtained so far translate to the case of learning from informant. Since Algorithm \mathcal{IML} does not learn anything from negative examples, one may expect that it behaves much poorer in this setting. First, we investigate the uniform distribution over \mathcal{X}_n . Again, we have the trivial cases that the target concept is “FALSE” or m is a monomial without irrelevant bits. In the first case, no example is needed at all, while in the latter one, there is only one positive example having probability 2^{-n} . Thus the expected number of examples needed until successful learning is $2^n = 2^{\#(m)}$.

THEOREM 18 (Reischuk and Zeugmann [33]). *Let $c = L(m) \in \mathcal{C}_n$ be a nontrivial concept. If a data sequence for c is generated from the uniform distribution on the learning domain by independent draws, the expected number of examples needed by Algorithm \mathcal{IML} until convergence is bounded by*

$$E[CONV] \leq 2^{\#(m)} (\lceil \log_2 k(m) \rceil + 3).$$

Proof. Let $CONV+$ be a random variable for the number of positive examples needed until convergence. Every positive example is preceded by a possibly empty block of negative examples. Thus, we can partition the initial segment of any randomly drawn informant read until convergence into $CONV+$ many blocks B_j containing a certain number of negative examples followed by precisely one positive example. Let Λ_j be a random variable for the length of block B_j . Then $CONV = \Lambda_1 + \Lambda_2 + \dots + \Lambda_{CONV+}$, where the Λ_j are independently identically distributed. In order to

compute the distribution of Λ_j , it suffices to calculate the probabilities to draw a negative and a positive example, respectively. Since the overall number of positive examples for c is 2^k with $k = k(m)$, the probability to generate a positive example is 2^{k-n} . Hence, the probability to draw a negative example is $1 - 2^{k-n}$. Consequently,

$$\Pr(\Lambda_j = \mu + 1) = (1 - 2^{k-n})^\mu \cdot 2^{k-n}.$$

Therefore,

$$\begin{aligned} \mathbb{E}[\text{CONV}] &= \mathbb{E}[\Lambda_1 + \Lambda_2 + \dots + \Lambda_{\text{CONV}+}] \\ &= \sum_{\zeta=0}^{\infty} \mathbb{E}[\Lambda_1 + \Lambda_2 + \dots + \Lambda_\zeta \mid \text{CONV}+ = \zeta] \cdot \Pr(\text{CONV}+ = \zeta) \\ &= \sum_{\zeta=0}^{\infty} \zeta \cdot \mathbb{E}[\Lambda_1] \cdot \Pr(\text{CONV}+ = \zeta) \\ &= \mathbb{E}[\text{CONV}+] \cdot \mathbb{E}[\Lambda_1] \end{aligned}$$

By Theorem 15, we have $\mathbb{E}[\text{CONV}+] \leq \lceil \log_2 k \rceil + 3$, and thus it remains to estimate $\mathbb{E}[\Lambda_1]$. A simple calculation shows

LEMMA 2. *For every $0 < a < 1$, it holds:*

$$\sum_{\mu=0}^{\infty} (\mu + 1) \cdot a^\mu = (1 - a)^{-2}.$$

Using this estimation we can conclude

$$\begin{aligned} \mathbb{E}[\Lambda_1] &= \sum_{\mu=0}^{\infty} (\mu + 1) \cdot \Pr(\Lambda_1 = \mu + 1) \\ &= 2^{k-n} \sum_{\mu=0}^{\infty} (\mu + 1) \cdot (1 - 2^{k-n})^\mu = 2^{n-k}, \end{aligned}$$

and thus the theorem follows. ▀

Hence, as long as the length of m is constant, and therefore $k(m) = n - O(1)$, we still achieve an expected total learning time of order $n \log n$. But if $\#(m)$ grows linearly the expected total learning becomes exponential. On the other hand, if there are many relevant literals then even h_0 may be considered as a not too bad *approximation* for c . Consequently, it is natural at this point to introduce an error parameter $\varepsilon \in (0, 1)$ as in the PAC model, and to ask whether one can achieve an expected sample complexity for computing an ε -approximation that is bounded by a function depending on $\log n$ and $1/\varepsilon$.

To answer this question, let us formally define $\text{error}_m(h_j) = D(L(h_j) \Delta L(m))$ to be the error made by hypothesis h_j with respect to monomial m . Here $L(h_j) \Delta L(m)$ stands for the symmetric difference of $L(h_j)$ and $L(m)$ and D for the underlying

probability distribution with respect to which the examples are drawn. Note that by construction of Algorithm \mathcal{IML} we can conclude $error_m(h_j) = D(L(m) \setminus L(h_j))$.

We call h_j an ε -approximation for m if $error_m(h_j) \leq \varepsilon$. Furthermore, we redefine the stage of convergence. Let m be any monomial, and let $d = (d_j)_{j \in \mathbb{N}}$ be an informant for $L(m)$, then

$$\begin{aligned} \mathbf{CONV}_\varepsilon(\mathbf{d}) &:= \text{the least number } j \text{ such that} \\ &\quad error_m(\mathcal{IML}(d_i)) \leq \varepsilon \text{ for all } i \geq j. \end{aligned}$$

Note that once the Algorithm \mathcal{IML} has reached an ε -approximate hypothesis all further hypotheses will also be at least that close to the target monomial.

The following theorem gives an affirmative answer to the question posed above.

THEOREM 19 (Reischuk and Zeugmann [33]). *Let $c = L(m) \in \mathcal{C}_n$ be a nontrivial concept. Assuming that examples are drawn at random independently from the uniform distribution, the expected number of examples needed by Algorithm \mathcal{IML} until converging to an ε -approximation for c can be bounded by*

$$\mathbb{E}[\mathbf{CONV}_\varepsilon] \leq \frac{1}{\varepsilon} \cdot (\lceil \log_2 k(m) \rceil + 3).$$

Proof. It holds $error_m(h_{ini}) = 2^{k(m)-n}$, since h_{ini} misclassifies exactly the positive examples. Therefore, if $error_m(h_0) \leq \varepsilon$, we are already done. Now suppose $error_m(h_0) > \varepsilon$. Consequently, $1/\varepsilon > 2^{n-k(m)}$, and thus the bound stated in the theorem is larger than $2^{n-k(m)}(\lceil \log_2 k(m) \rceil + 3)$, which, by Theorem 18 is the expected number of examples needed until convergence to a correct hypothesis. ■

Thus, additional knowledge concerning the underlying probability distribution pays off again. Applying Theorem 6 and modifying the stochastic finite learner presented above *mutatis mutandis*, we get a learner identifying ε -approximations for all concepts in \mathcal{C}_n stochastically with high confidence using $O(\frac{1}{\varepsilon} \cdot \log \frac{1}{\delta} \cdot \log n)$ many examples. Comparing this bound with the sample complexity given in the PAC model, we see that it is reduced exponentially, i.e., instead of a factor n now we have the factor $\log n$.

Finally, we can generalize the last results to the case that the data sequences are binomially distributed for some parameter $p \in (0, 1)$. This means that any particular vector containing ν times a 1 and $n - \nu$ a 0 has probability $p^\nu(1-p)^{n-\nu}$ since a 1 is drawn with probability p and a 0 with probability $1-p$. First, Theorem 18 generalizes as follows.

THEOREM 20 (Reischuk and Zeugmann [33]). *Let $c = L(m) \in \mathcal{C}_n$ be a nontrivial concept. Let m contain precisely π positive literals and ν negative literals. If the labeled examples for c are independently binomially distributed with parameter p and*

$\psi := \min\{\frac{1}{1-p}, \frac{1}{p}\}$ and $\tau := \max\{\frac{p}{1-p}, \frac{1-p}{p}\}$, then the expected number of examples needed by Algorithm \mathcal{IML} until convergence can be bounded by

$$\mathbb{E}[\text{CONV}] \leq \frac{1}{p^\pi(1-p)^\nu} (\lceil \log_\psi k(m) \rceil + \tau + 2) .$$

Proof. Assuming the same notation as in the proof of Theorem 18, it is easy to see that we only have to recompute $\mathbb{E}[\Lambda_1]$, and thus $\Pr(\Lambda_1 = \mu + 1)$, too. Since m contains precisely π positive literals and ν negative literals, the probability to draw a positive example is clearly $p^\pi(1-p)^\nu$, and thus the probability to randomly draw a negative example is $1 - p^\pi(1-p)^\nu$. Consequently,

$$\Pr(\Lambda_1 = \mu + 1) = (1 - p^\pi(1-p)^\nu)^\mu \cdot p^\pi(1-p)^\nu ,$$

and Lemma 2 gives $\mathbb{E}[\Lambda_1] = \frac{1}{p^\pi(1-p)^\nu}$. ■

Theorem 19 directly translates into the setting of binomially distributed inputs, too.

THEOREM 21 (Reischuk and Zeugmann [33]). *Let $c = L(m) \in \mathcal{C}_n$ be a nontrivial concept. Assume that the examples are drawn with respect to a binomial distribution with parameter p , and let $\psi := \min\{\frac{1}{1-p}, \frac{1}{p}\}$ and $\tau := \max\{\frac{p}{1-p}, \frac{1-p}{p}\}$. Then the expected number of examples needed by Algorithm \mathcal{IML} until converging to an ε -approximation for c can be bounded by*

$$\mathbb{E}[\text{CONV}] \leq \frac{1}{\varepsilon} \cdot (\lceil \log_\psi k(m) \rceil + \tau + 2) .$$

Finally, one can also get stochastic finite approximations with high confidence from informant with an exponentially smaller sample complexity.

THEOREM 22 (Reischuk and Zeugmann [33]). *Let $0 < p_{low} \leq p_{up} < 1$ and $\psi := \min\{\frac{1}{1-p_{low}}, \frac{1}{p_{up}}\}$. For $(\mathcal{C}_n, \mathcal{D}_n[p_{low}, p_{up}])$ ε -approximations are stochastically finitely learnable with δ -confidence from informant for any $\varepsilon, \delta \in (0, 1)$.*

For this purpose, $O(\frac{1}{\varepsilon} \cdot \log_2 1/\delta \cdot \log_\psi n)$ many examples suffice.

5. Conclusions

The present paper surveyed results recently obtained concerning the iterative learnability of the class of all pattern languages and finite unions thereof. In particular, it could be shown that there are strong dependencies between iterative learning, the class of admissible hypothesis spaces and additional requirements to the learner such as consistency, conservativeness and the decidability of the inclusion problem for the hypothesis space chosen. Looking at these results, we have seen that the LWA is in some sense optimal.

Moreover, by analyzing the average-case behavior of Lange and Wiehagen's pattern language learning algorithm with respect to its total learning time and by establishing exponentially shrinking tail bounds for a rather rich class of limit learners, we have been able to transform the LWA into a stochastic finite learner. The price paid is the incorporation of a bit prior knowledge concerning the class of underlying probability distributions. When applied to the class of all k -variable pattern languages, where k is *a priori* known, the resulting total learning time is linear in the expected string length.

Thus, the present paper provides evidence that analyzing the average-case behavior of limit learners with respect to their total learning time may be considered as a promising path towards a new theory of efficient algorithmic learning. Recently obtained results along the same path as outlined in Erlebach *et al.*[11] as well as in Reischuk and Zeugmann [32, 34] provide further support for the fruitfulness of this approach.

In particular, in Reischuk and Zeugmann [32, 34] we have shown that one-variable pattern languages are learnable for basically all meaningful distributions within an optimal linear total learning time on the average. Furthermore, this learner can also be modified to maintain the *incremental* behavior of Lange and Wiehagen's [19] algorithm. Instead of memorizing the pair (PRE, SUF), it can also store just the two or three examples from which the prefix PRE and the suffix SUF of the target pattern has been computed. While it is no longer *iterative*, it is still a *bounded example memory* learner. A bounded example memory learner is essentially an iterative learner that is additionally allowed to memorize an *a priori* bounded number of examples (cf. [9] for a formal definition).

While the one-variable pattern language learner from [34] is highly practical, our stochastic finite learner for the class of all pattern languages is still not good enough for practical purposes. But our results surveyed point to possible directions for potential improvements. However, much more effort seems necessary to design a stochastic finite learner for $PAT(k)$.

Additionally, we have applied our techniques to design a stochastic finite learner for the class of all concepts describable by a monomial which is based on Haussler's [14] Wholist algorithm. Here we have assumed the examples to be binomially distributed. The sample size of our stochastic finite learner is mainly bounded by $\log(1/\delta) \log n$, where δ is again the confidence parameter and n is the dimension of the underlying Boolean learning domain. Thus, the bound obtained is exponentially better than the bound provided within the PAC model.

Our approach also differs from U-learnability introduced by Muggleton [27]. First of all, our learner is fed with positive examples only, while in Muggleton's [27] model examples labeled with respect to their containment in the target language are provided. Next, we do not make any assumption concerning the distribution of the target patterns. Furthermore, we do *not* measure the expected total learning time with respect to a given class of distributions over the targets and a given class of distributions

for the sampling process, but exclusively in dependence on the length of the target. Finally, we require exact learning and not approximately correct learning.

References

- [1] D. Angluin, Finding patterns common to a set of strings, *Journal of Computer and System Sciences* 21 (1) (1980) 46–62.
- [2] D. Angluin, Inductive inference of formal languages from positive data, *Information and Control* 45 (2) (1980) 117–135.
- [3] D. Angluin and C.H. Smith, Inductive inference: Theory and methods. *Computing Surveys* 15 (3) (1983) 237–269.
- [4] D. Angluin and C.H. Smith, Formal inductive inference, in: St.C. Shapiro (Ed.), *Encyclopedia of Artificial Intelligence*, Vol. 1, Wiley-Interscience Publication, New York, 1987, pp. 409–418.
- [5] S. Arikawa, T. Shinohara and A. Yamamoto, Learning elementary formal systems, *Theoretical Computer Science* 95 (1) (1992) 97–113.
- [6] L. Blum and M. Blum, Toward a mathematical theory of inductive inference, *Information and Control* 28 (2) (1975) 125–155.
- [7] A. Blumer, A. Ehrenfeucht, D. Haussler and M. Warmuth, Learnability and the Vapnik-Chervonenkis dimension, *Journal of the ACM* 36 (1989) 929–965.
- [8] I. Bratko and S. Muggleton, Applications of inductive logic programming, *Communications of the ACM* 38 (11) (1995) 65–70.
- [9] J. Case, S. Jain, S. Lange and T. Zeugmann, Incremental concept learning for bounded data mining, *Information and Computation* 152 (1) (1999) 74–110.
- [10] R.P. Daley and C.H. Smith, On the complexity of inductive inference, *Information and Control* 69 (1-3) (1986) 12–40.
- [11] T. Erlebach, P. Rossmanith, H. Stadtherr, A. Steger and T. Zeugmann, Learning one-variable pattern languages very efficiently on average, in parallel, and by asking queries, *Theoretical Computer Science* 261 (1/2) (2001) 119–156.
- [12] E.M. Gold, Language identification in the limit, *Information and Control* 10 (5) (1967) 447–474.
- [13] S.A. Goldman, M.J. Kearns and R.E. Schapire, Exact identification of circuits using fixed points of amplification functions, *SIAM Journal of Computing* 22 (4) (1993) 705–726.

- [14] D. Haussler, Bias, version spaces and Valiant's learning framework, in: Proceedings of the Fourth International Workshop on Machine Learning, Morgan Kaufmann, San Mateo, CA, 1987, pp. 324–336.
- [15] D. Haussler, M. Kearns, N. Littlestone and M.K. Warmuth, Equivalence of models for polynomial learnability, *Information and Computation* 95 (2) (1991) 129–161.
- [16] S. Jain, D. Osherson, J.S. Royer and A. Sharma, "Systems That Learn: An Introduction to Learning Theory," MIT-Press, Boston, Massachusetts, 1999.
- [17] T. Jiang, A. Salomaa, K. Salomaa and S. Yu, Inclusion is undecidable for pattern languages, in: A. Lingas, R. Karlsson, and S. Carlsson (Eds.), Automata, Languages and Programming, 20nd International Colloquium, ICALP 93, Lund, Sweden, July 5-9, 1993, Proceedings, Lecture Notes in Computer Science, Vol. 700, Springer-Verlag, Berlin, 1993, pp. 301–312.
- [18] M. Kearns and L. Pitt, A polynomial-time algorithm for learning k -variable pattern languages from examples, in: R. Rivest, D. Haussler, and M.K. Warmuth (Eds.), Proceedings of the Second Annual Workshop on Computational Learning Theory, Santa Cruz, California, Morgan Kaufmann, San Mateo, CA, 1989, pp. 57–71.
- [19] S. Lange and R. Wiehagen, Polynomial-time inference of arbitrary pattern languages, *New Generation Computing* 8 (4) (1991) 361–370.
- [20] S. Lange and T. Zeugmann, Language learning in dependence on the space of hypotheses, in: L. Pitt (Ed.), Proceedings of the Sixth Annual ACM Conference on Computational Learning Theory, July 26th-28th, 1993, Santa Cruz, California, ACM Press, New York, 1993, pp. 127–136.
- [21] S. Lange and T. Zeugmann, Set-driven and rearrangement-independent learning of recursive languages, *Mathematical Systems Theory* 29 (6) (1996) 599–634.
- [22] S. Lange and T. Zeugmann, Incremental learning from positive data, *Journal of Computer and System Sciences* 53 (1) (1996) 88–103.
- [23] N. Lavrač and S. Džeroski, "Inductive Logic Programming: Techniques and Applications," Ellis Horwood, 1994.
- [24] T. Mitchell. "Machine Learning," McGraw-Hill, New York, 1997.
- [25] A. Mitchell, A. Sharma, T. Scheffer and F. Stephan, The VC-dimension of subclasses of pattern languages, in: O. Watanabe and T. Yokomori (Eds.), Algorithmic Learning Theory, 10th International Conference, ALT '99, Tokyo, Japan, December 1999, Proceedings, Lecture Notes in Artificial Intelligence, Vol. 1720, Springer, Berlin, 1999, pp. 93 - 105.

- [26] S. Miyano, A. Shinohara and T. Shinohara, Polynomial-time learning of elementary formal systems, *New Generation Computing* 18 (3) (2000) 217–242.
- [27] S. Muggleton, Bayesian Inductive Logic Programming, in: M. Warmuth (Ed.), Proceedings of the Seventh Annual ACM Conference on Computational Learning Theory, July 12th-15th, 1994, New Brunswick, New Jersey, ACM Press, New York, 1994, pp. 3–11.
- [28] S. Muggleton and L. De Raedt, Inductive logic programming: Theory and methods, *Journal of Logic Programming* 19/20 (1994) 669–679.
- [29] R.P. Nix, Editing by examples, Yale University, Dept. Computer Science, Technical Report 280, 1983.
- [30] D.N. Osherson, M. Stob and S. Weinstein, “Systems that Learn, An Introduction to Learning Theory for Cognitive and Computer Scientists,” MIT-Press, Cambridge, Massachusetts, 1986.
- [31] L. Pitt, Inductive inference, DFAs and computational complexity, in: K.P. Jantke (Ed.), Analogical and Inductive Inference, International Workshop AII ’89, Reinhardtsbrunn Castle, GDR, October 1989, Proceedings, Lecture Notes in Artificial Intelligence, Vol. 397, Springer-Verlag, Berlin, 1989, pp. 18–44.
- [32] R. Reischuk and T. Zeugmann, Learning one-variable pattern languages in linear average time, in: Proceedings of the Eleventh Annual Conference on Computational Learning Theory, July 24th - 26th 1998, Madison, Wisconsin, ACM Press, New York, 1998, pp. 198–208.
- [33] R. Reischuk and T. Zeugmann, A complete and tight average-case analysis of learning monomials, in: C. Meinel and S. Tison (Eds.), STACS 99, 16th Annual Symposium on Theoretical Aspects of Computer Science, Trier, Germany, March 1999, Proceedings, Lecture Notes in Computer Science, Vol. 1563, Springer 1999, pp. 414–423.
- [34] R. Reischuk and T. Zeugmann, An average-case optimal one-variable pattern language learner, *Journal of Computer and System Sciences* 60 (2) (2000) 302–335.
- [35] H. Rogers, Jr., “Theory of Recursive Functions and Effective Computability,” McGraw–Hill, New York, 1967.
- [36] P. Rossmanith and T. Zeugmann, Stochastic finite learning of the pattern languages, *Machine Learning* 44 (1-2) (2001) 67–91.
- [37] A. Salomaa, Patterns (The Formal Language Theory Column), EATCS Bulletin 54 (1994) 46–62.

- [38] A. Salomaa, Return to patterns (The Formal Language Theory Column), EATCS Bulletin 55 (1994) 144–157.
- [39] S. Shimozono, A. Shinohara, T. Shinohara, S. Miyano, S. Kuhara and S. Arikawa, Knowledge acquisition from amino acid sequences by machine learning system BONSAI, *Trans. Information Processing Society of Japan* 35 (1994) 2009–2018.
- [40] R.E. Schapire, Pattern languages are not learnable, in: M.A. Fulk and J. Case (Eds.), Proceedings of the Third Annual Workshop on Computational Learning Theory, Morgan Kaufmann, San Mateo, CA, 1990, pp. 122–129.
- [41] T. Shinohara, Inferring unions of two pattern languages, *Bulletin of Informatics and Cybernetics* 20 (1983) 83–88.
- [42] T. Shinohara, Inductive inference of monotonic formal systems from positive data, *New Generation Computing* 8 (4) (1991) 371–384.
- [43] T. Shinohara and S. Arikawa, Pattern inference, in: K.P. Jantke and S. Lange (Eds.), Algorithmic Learning for Knowledge-Based Systems, Lecture Notes in Artificial Intelligence, Vol. 961, Springer-Verlag, Berlin, 1995, pp. 259–291.
- [44] T. Shinohara and H. Arimura, Inductive inference of unbounded unions of pattern languages from positive data, in: S. Arikawa and A.K. Sharma (Eds.), Algorithmic Learning Theory, 7th International Workshop, ALT '96, Sydney, Australia, October 1996, Proceedings, Lecture Notes in Artificial Intelligence, Vol. 1160, Springer-Verlag, Berlin, 1996, pp. 256–271.
- [45] R. Smullyan, “Theory of Formal Systems,” *Annals of Mathematical Studies*, No. 47, Princeton, NJ, 1961.
- [46] L.G. Valiant, A theory of the learnable, *Communications of the ACM* 27 (11) (1984) 1134–1142.
- [47] R. Wiehagen, Limes-Erkennung rekursiver Funktionen durch spezielle Strategien, *Journal of Information Processing and Cybernetics (EIK)* 12 (1/2) (1976) 93–99.
- [48] R. Wiehagen and T. Zeugmann, Ignoring data may be the only way to learn efficiently, *Journal of Experimental and Theoretical Artificial Intelligence* 6 (1) (1994) 131–144.
- [49] R. Wiehagen and T. Zeugmann, Learning and Consistency, in: K.P. Jantke and S. Lange (Eds.), Algorithmic Learning for Knowledge-Based Systems, Lecture Notes in Artificial Intelligence, Vol. 961, Springer-Verlag, Berlin, 1995, pp. 1–24.
- [50] K. Wright, Identification of unions of languages drawn from an identifiable class, in: R. Rivest, D. Haussler, and M.K. Warmuth (Eds.), Proceedings of the Second Annual Workshop on Computational Learning Theory, Santa Cruz, California, Morgan Kaufmann, San Mateo, CA, 1989, pp. 328–333.

- [51] T. Zeugmann, Lange and Wiehagen's pattern language learning algorithm: An average-case analysis with respect to its total learning time, *Annals of Mathematics and Artificial Intelligence* 23 (1-2) (1998) 117–145.
- [52] T. Zeugmann and S. Lange, A guided tour across the boundaries of learning recursive languages, in: K.P. Jantke and S. Lange (Eds.), *Algorithmic Learning for Knowledge-Based Systems*, Lecture Notes in Artificial Intelligence, Vol. 961, Springer-Verlag, Berlin, 1995, pp. 190–258.
- [53] T. Zeugmann, S. Lange and S. Kapur, Characterizations of monotonic and dual monotonic language learning, *Information and Computation* 120 (2) (1995) 155–173.