

Learning Languages from Bounded Resources: The Case of the DFA and the Balls of Strings^{*}

Colin de la Higuera, Jean-Christophe Janodet, and Frédéric Tantini

Universities of Lyon, 18 r. Pr. Lauras, F-42000 St-Etienne
{cdlh,janodet,frederic.tantini}@univ-st-etienne.fr

Abstract. Comparison of standard language learning paradigms (identification in the limit, query learning, PAC learning) has always been a complex question. Moreover, when to the question of converging to a target one adds computational constraints, the picture becomes even less clear: how much do queries or negative examples help? Can we find good algorithms that change their minds very little or that make very few errors? In order to approach these problems we concentrate here on two classes of languages, the topological balls of strings (for the edit distance) and the deterministic finite automata (DFA), and (re-)visit the different learning paradigms to sustain our claims.

Keywords: Polynomial learnability, deterministic finite automata, balls of strings, edit distance.

1 Introduction

The study of the properties of the learning algorithms, particularly those in grammatical inference, can be either empirical (based on experiments from datasets), or theoretical. In the latter, the goal is to study the capacity of the algorithm to retrieve, exactly or approximately, a target language. Often, the goal is also to measure the resources (time, amount of data) necessary to achieve this task. Different paradigms have been proposed to take into account notions of convergence from bounded resources, but none has really imposed itself, and few comparison exists between these definitions.

In this paper, we visit standard criteria for *polynomial* identification and compare them by considering two fundamentally different classes of languages: the regular languages represented by deterministic finite automata, and the balls of strings *w.r.t.* the edit distance [1]. These balls of strings are formed by choosing one specific string, called the centre, and all its neighbours up to a given length for the edit distance, called the radius.

From a practical standpoint, the balls of strings appear in a variety of settings: in approximate string matching tasks, the goal is to find all close matches to some

* This work was supported in part by the IST Programme of the European Community, under the PASCAL Network of Excellence, IST-2006-216886. This publication only reflects the authors' views.

target string [2,3]; in noisy settings, garbled versions of an unidentified string are given and the task is to recover the original string [4]; when using dictionaries, the task can be described as that of finding the intersection between two languages, the dictionary itself and a ball around the target string [5]; in the field of bioinformatics, extracting valid models from large datasets of DNA or proteins can involve looking for substrings at distance less than some given bound, and the set of these approximate substrings can also be represented by balls [6].

When aiming to prove that a class of languages is learnable, there are typically three different settings. The first one, identification in the limit [7], mimics the cognitive process of a child that would acquire his native language by picking up the sentences that are broadcasted in his environment. More formally, information keeps on arriving about a target language and the Learner keeps making new hypotheses. We say that convergence takes place if there is a moment when the process is stationary and the hypothesis is correct.

The second one, query learning [8], looks like a game of riddles where the Learner (pupil) can ask questions (queries) to an Oracle (teacher) about the target language. The game ends when the Learner guesses the target. Of course, the learning results strongly depends on the sort of queries that the Learner is allowed to ask. Both previous paradigms are probably at least as interesting for the negative results they induce as for the positive ones. Indeed, concerning query learning, if a Learner cannot identify an unknown concept by choosing and testing examples, how could he hope to succeed to learn from examples that are imposed by an application?

The last paradigm, PAC learning (for Probably Approximately Correct) [9] is intended to be a more pragmatic setting. It formalizes a situation where one tries to build automatically a predictive model from the data. In this setting, one assumes that there is a (unknown) distribution \mathcal{D} over the strings of the target language, which is used to sample learning and testing examples. Two parameters are fixed: ϵ is related to the error of the model (*i.e.*, the probability of a string to be misclassified) and δ is related to the confidence one has in the sampling. Ideally, a good PAC-Learner returns, with high confidence ($> 1 - \delta$), hypotheses that have small error rates ($< \epsilon$).

The three settings are usually difficult to compare, in particular when complexity issues are discussed. Some exceptions are the work by Angluin comparing PAC-learning and using equivalence queries [10], the work by Pitt relating equivalence queries and implicit prediction errors [11], comparisons between learning with characteristic samples, simple PAC [12,13] and MAT in [14]. Other analysis of polynomial aspects of learning grammars, automata and languages can be found in [11,15,16,17]. If the customary comparative approach is to introduce a learning paradigm and survey a variety of classes of languages for this paradigm, we choose here to fix the classes of languages and to proceed to a horizontal analysis of their learnability by visiting the paradigms systematically.

Concerning the DFA, we complete a long list of known results. Concerning the balls of strings, our results are generally negative: identification in the limit from examples and counter-examples is impossible in most cases, even from

membership and equivalence queries. PAC-learning is also impossible in polynomial time, unless $\mathcal{RP} = \mathcal{NP}$. Yet, the errors are usually due to the counterexamples. Hence, we show that it is sometimes (and surprisingly) easier to learn from positive examples only than from positive and negative examples.

Section 2 is devoted to preliminary definitions. In Sections 3, 4 and 5, we focus on the so-called good balls and on the DFA, and we present the results concerning PAC-learning, query learning and polynomial identification in the limit, respectively. We conclude in Section 6.

2 Definitions

An *alphabet* Σ is a finite nonempty set of symbols called *letters*. In the sequel, we suppose that $|\Sigma| \geq 2$. A *string* $w = a_1 \cdots a_n$ is any finite sequence of letters. We write λ for the empty string and $|w|$ for the length of w . Let Σ^* denote the set of all strings over Σ . We say that u is a *subsequence* of v , denoted $u \preceq v$, if $_{def} u = a_1 \cdots a_n$ and there exist $u_0, \dots, u_n \in \Sigma^*$ s.t. $v = u_0 a_1 u_1 \cdots a_n u_n$. We introduce the set $lcs(u, v)$ of all longest common subsequences of u and v . We also introduce the *hierarchical order*: $u \preceq v$ if $_{def} |u| < |v|$ or ($|u| = |v|$ and $u \leq_{lex} v$), where \leq_{lex} denotes the standard lexicographic order. A *language* is any subset $L \subseteq \Sigma^*$. Let \mathbb{N} denote the set of non negative integers. For all $k \in \mathbb{N}$, let $\Sigma^{\leq k}$ (respectively $\Sigma^{> k}$) be the set of all strings of length at most k (respectively of length more than k). We define $A \oplus B = (A \setminus B) \cup (B \setminus A)$.

Grammatical inference aims at learning the languages of a fixed class \mathcal{L} represented by the grammars of a class \mathcal{G} . \mathcal{L} and \mathcal{G} are related by a naming function $\mathbb{L} : \mathcal{G} \rightarrow \mathcal{L}$ that is total ($\forall G \in \mathcal{G}, \mathbb{L}(G) \in \mathcal{L}$) and surjective ($\forall L \in \mathcal{L}, \exists G \in \mathcal{G}$ s.t. $\mathbb{L}(G) = L$). For any string $w \in \Sigma^*$ and language $L \in \mathcal{L}$, we shall write $L \models w$ if $_{def} w \in L$. Concerning the grammars, they may be understood as any piece of information allowing some parser to recognize the strings. For any string $w \in \Sigma^*$ and grammar $G \in \mathcal{G}$, we shall write $G \vdash w$ if the parser recognizes w . Basically, the parser must be sound and complete *w.r.t.* the semantics: $G \vdash w \iff \mathbb{L}(G) \models w$. In the following, we will mainly consider learning paradigms subject to complexity constraints. In their definitions, $\|G\|$ will denote the size of the grammar G (e.g., the number of states in the case of DFA). Moreover, given a set X of strings, we will write $|X|$ for the cardinality of X and $\|X\|$ for the sum of the lengths of the strings in X .

The *edit distance* $d(w, w')$ is the minimum number of *primitive edit operations* needed to transform w into w' [1]. The operation is either (1) a *deletion*: $w = uav$ and $w' = uv$, or (2) an *insertion*: $w = uv$ and $w' = uav$, or (3) a *substitution*: $w = uav$ and $w' = ubv$, where $u, v \in \Sigma^*$, $a, b \in \Sigma$ and $a \neq b$. E.g., $d(abaa, aab) = 2$ since $\underline{a}baa \rightarrow a\underline{a}a \rightarrow aab$ and the rewriting of $abaa$ into aab cannot be achieved with less than two steps. $d(w, w')$ can be computed in $\mathcal{O}(|w| \cdot |w'|)$ time by dynamic programming [18].

The edit distance is a metric, so we can introduce the *balls* over Σ . The *ball of centre* $o \in \Sigma^*$ and *radius* $r \in \mathbb{N}$, denoted $B_r(o)$, is the set of all strings whose distance is at most r from o : $B_r(o) = \{w \in \Sigma^* : d(o, w) \leq r\}$. E.g., if $\Sigma = \{a, b\}$,

then $B_1(ba) = \{a,b,aa,ba,bb,aba,baa,bab,baa\}$ and $B_r(\lambda) = \Sigma^{\leq r}$ for all $r \in \mathbb{N}$. We will write $\mathbf{BALL}(\Sigma)$ for the family of all the balls.

To the purpose of grammatical inference, we are going to represent any ball $B_r(o)$ by the pair (o, r) that will play the role of a grammar. Indeed, its size is $|o| + \log r$ (which corresponds to the number of bits necessary to encode the grammar¹). Moreover, the parser able to decide whether $w \in B_r(o)$ or not is simple: (1) it computes $d(o, w)$ and (2) it checks if this distance is $\leq r$, that can be achieved in time $\mathcal{O}(|o| \cdot |w| + \log r)$. Finally, as $|\Sigma| \geq 2$, we can show that (o, r) is a unique thus *canonical* grammar of $B_r(o)$ [20]. In consequence, we shall also denote by $\mathbf{BALL}(\Sigma)$ the class of grammars associated to the balls.

A ball $B_r(o)$ is called *good* if_{def} $r \leq |o|$. The advantage of using good balls is that there is a polynomial relation between the size of the centre and the size of the longest strings in the ball. We will write $\mathbf{GB}(\Sigma)$ for the class of all the good balls (and that of the corresponding grammars).

A *deterministic finite automaton* (DFA) is a 5-tuple $A = \langle \Sigma, Q, q_0, F, \delta \rangle$ s.t. Q is a set of states, $q_0 \in Q$ is an initial state, $F \subseteq Q$ is a set of final states and $\delta : Q \times \Sigma \rightarrow Q$ is a transition function. Every DFA can be *completed* with one sink state s.t. δ is a total function. As usual, δ is extended to Σ^* . The *language recognized by A* is $\mathbb{L}(A) = \{w \in \Sigma^* : \delta(q_0, w) \in F\}$. The size of A is $|Q|$. We will write $\mathbf{DFA}(\Sigma)$ for the class of all DFA over the alphabet Σ .

The inference of DFA has been intensively studied for forty years at least [7,11,21]. On the other hand, the learnability of the balls is a recent issue [20] motivated by the problem of identifying languages from noisy data. However, our approach raises a preliminary question: if any ball whose grammar would have size n could be recognized by a DFA with $p(n)$ states (for some polynomial $p()$), then one could deduce learnability results on the former from (known) learnability results on the latter. Yet it is generally believed (albeit still an open question) that the transformation of a ball into a DFA is not polynomial [2].

3 PAC-Learnability

The PAC paradigm [9] has been widely used in machine learning. It aims at building, with high confidence, good approximations of an unknown concept.

Definition 1 (ϵ -good hypothesis). *Let G be the target grammar and H be a hypothesis grammar. Let \mathcal{D} be a distribution over Σ^* and $\epsilon > 0$. We say that H is an ϵ -good hypothesis w.r.t. G if_{def} $Pr_{\mathcal{D}}(x \in \mathbb{L}(G) \oplus \mathbb{L}(H)) < \epsilon$.*

The PAC-learnability of grammars from strings of unbounded size has always been tricky [22,16,23]. Indeed, with the standard definition, a PAC-Learner can ask an Oracle to return a sample randomly drawn according to the distribution \mathcal{D} . However, in the case of strings, there is always the risk (albeit small) to sample a string too long to account for in polynomial time. In order to avoid this

¹ Notice that $|o| + r$ is not a correct measure of the size as implicitly it would mean encoding the radius in unary, something unreasonable [19].

problem, we will sample from a distribution restricted to strings shorter than a specific value given by the following lemma:

Lemma 1. *Let \mathcal{D} be a distribution over Σ^* . Then $\forall \epsilon, \delta > 0$, with probability at least $1 - \delta$, if one draws a sample X of at least $\frac{1}{\epsilon} \ln \frac{1}{\delta}$ strings following \mathcal{D} , then the probability of any new string x to be longer than all the strings of X is less than ϵ . Formally, let $\mu_X = \max\{|y| : y \in X\}$, then $Pr_{x \sim \mathcal{D}}(|x| > \mu_X) < \epsilon$.*

Proof. Let ℓ be the smallest integer s.t. $Pr_{\mathcal{D}}(\Sigma^{>\ell}) < \epsilon$. A sufficient condition for $Pr_{\mathcal{D}}(|x| > \mu_X) < \epsilon$ is that we take a sample X large enough to be nearly sure (with probability $> 1 - \delta$) to have one string $\geq \ell$. Basically, the probability of drawing n strings in X of length $< \ell$ is $\leq (1 - \epsilon)^n$. So the probability of getting at least one string of length $\geq \ell$ is $> 1 - (1 - \epsilon)^n$. In order to build X , we thus need $1 - (1 - \epsilon)^n > 1 - \delta$, that is to say, $(1 - \epsilon)^n < \delta$. As $(1 - \epsilon)^n \leq e^{-n\epsilon}$, it is sufficient to take $n \geq \frac{1}{\epsilon} \ln \frac{1}{\delta}$ to reach a convenient value for μ_X . \square

An algorithm is now asked to learn a grammar given a *confidence* parameter δ and an *error* parameter ϵ . The algorithm must also be given an upper bound n on the size of the target grammar and an upper bound m on the length of the examples it is going to get (perhaps computed using Lemma 1). The algorithm can query an Oracle for an example randomly drawn according to the distribution \mathcal{D} . The query of an example or a counter-example will be denoted $\text{EX}()$. When the Oracle is only queried for a positive example, we will write $\text{POS-EX}()$. And when the Oracle is only queried for string of length $\leq m$, we will write $\text{EX}(m)$ and $\text{POS-EX}(m)$ respectively. Formally, the Oracle will then return a string drawn from \mathcal{D} , or $\mathcal{D}(\mathbb{L}(G))$, or $\mathcal{D}(\Sigma^{\leq m})$, or $\mathcal{D}(\mathbb{L}(G) \cap \Sigma^{\leq m})$, respectively, where $\mathcal{D}(L)$ is the restriction of \mathcal{D} to the strings of L : $Pr_{\mathcal{D}(L)}(x) = Pr_{\mathcal{D}}(x)/Pr_{\mathcal{D}}(L)$ if $x \in L$, 0 otherwise. $Pr_{\mathcal{D}(L)}(x)$ is not defined if $L = \emptyset$.

Definition 2 (Polynomial PAC-learnability). *Let \mathcal{G} be a class of grammars. \mathcal{G} is PAC-learnable if_{def} there exists an algorithm \mathfrak{A} s.t. $\forall \epsilon, \delta > 0$, for any distribution \mathcal{D} over Σ^* , $\forall n \in \mathbb{N}$, $\forall G \in \mathcal{G}$ of size $\leq n$, for any upper bound $m \in \mathbb{N}$ on the size of the examples, if \mathfrak{A} has access to $\text{EX}()$, ϵ , δ , n and m , then with probability $> 1 - \delta$, \mathfrak{A} returns an ϵ -good hypothesis w.r.t. G . If \mathfrak{A} runs in time polynomial in $\frac{1}{\epsilon}$, $\frac{1}{\delta}$, n and m , we say that \mathcal{G} is polynomially PAC-learnable.*

Typical techniques proving non PAC-learnability depend on complexity assumptions [24]. Let us recall that \mathcal{RP} (Randomised Polynomial Time) is the complexity class of decision problems for which a probabilistic Turing machine exists which (1) runs in time polynomial in the input size, (2) on a negative instance, always returns NO and (3) on a positive instance, returns YES with probability $> \frac{1}{2}$ (otherwise, it returns NO). The algorithm is randomised: it is allowed to flip a random coin while it is running. The algorithm does not make any error on negative instances, and it is important to remark that on positive instances, since the error is $< \frac{1}{2}$, by repeating the run of the algorithm as many times as necessary, the actual error can be brought to be as small as one wants. We will use the strong belief and assumption that $\mathcal{RP} \neq \mathcal{NP}$ [19].

We are going to show that the good balls are not polynomially PAC-learnable. The proof follows the classical lines for such results: we first prove that the associated consistency problem is \mathcal{NP} -hard, through reductions from a well known \mathcal{NP} -complete problem (*Longest Common Subsequence*). Then it follows that if a polynomial PAC-learning algorithm for balls existed, this algorithm would provide us with a proof that this \mathcal{NP} -complete problem would also be in \mathcal{RP} .

Lemma 2. *The following problems are \mathcal{NP} -complete:*

1. **Longest Common Subsequence (LCS):** *Given n strings x_1, \dots, x_n and an integer k , does there exist a string w which is a subsequence of each x_i and is of length k ?*
2. **Longest Common Subsequence of Strings of a Given Length (LCSSGL):** *Given n strings x_1, \dots, x_n all of length $2k$, does there exist a string w which is a subsequence of each x_i and is of length k ?*
3. **Consistent ball (CB):** *Given two sets X_+ and X_- of strings, does there exist a good ball containing X_+ and which does not intersect X_- ?*

Proof. (1) See [25]. (2) See [26, page 42], Problem LCS0. (3) We use a reduction of Problem LCSSGL. We take the strings of length $2k$, and put these with string λ into the set X_+ . We build X_- by taking each string of length $2k$ and inserting every possible symbol once only (hence constructing at most $n(2k+1)|\Sigma|$ strings of size $2k+1$). It follows that a ball that contains X_+ but no element of X_- has necessarily a centre of length k and a radius of k (since we focus on good balls only). The centre is then a subsequence of all the strings of length $2k$ that were given. Conversely, if a ball is built using a subsequence of length k as centre, this ball is of radius k , contains also λ , and because of the radius, contains no element of X_- . Finally the problem is in \mathcal{NP} , since given a centre o , it is easy to check if $\max_{x \in X_+} d(o, x) < \min_{x \in X_-} d(o, x)$. \square

Theorem 1. *Unless $\mathcal{RP} = \mathcal{NP}$, $\mathcal{GB}(\Sigma)$ is not polynomially PAC-learnable.*

Proof. Suppose that $\mathcal{GB}(\Sigma)$ is polynomially PAC-learnable with \mathfrak{A} and take an instance $\langle X_+, X_- \rangle$ of Problem CB. We write $h = |X_+| + |X_-|$ and define over Σ^* the distribution $Pr(x) = \frac{1}{h}$ if $x \in X_+ \cup X_-$, 0 if not. Let $\epsilon = \frac{1}{h+1}$, $\delta < \frac{1}{2}$, $m = n = \max\{|w| : w \in X_+\}$. Let $B_r(o)$ be the ball returned by $\mathfrak{A}(\epsilon, \delta, n, m)$ and test if $(X_+ \subseteq B_r(o) \text{ and } X_- \cap B_r(o) = \emptyset)$. If there is no consistent ball, then $B_r(o)$ is inconsistent with the data, so the test is false. If there is a consistent ball, then $B_r(o)$ is ϵ -good, with $\epsilon < \frac{1}{h}$. So, with probability at least $1 - \delta > \frac{1}{2}$, there is no error at all and the test is true. This procedure runs in polynomial time in $\frac{1}{\epsilon}$, $\frac{1}{\delta}$, n and m . So if the good balls were PAC-learnable, there would be a randomized algorithm for the \mathcal{NP} -complete CB Problem (by Lemma 2). \square

Concerning the PAC-learnability of the DFA, a lot of studies have been done [11,16,23,24] The associated consistency problem is hard [27] and an efficient learning algorithm could be used to invert the RSA encryption function [16]:

Theorem 2 ([16]). *$\mathcal{DFA}(\Sigma)$ is not polynomially PAC-learnable.*

In certain cases, it may even be possible to PAC-learn from positive examples only. In this setting, during the learning phase, the examples are sampled following POS-EX() whereas during the testing phase, the sampling is done following EX(), but in both cases the distribution is identical. Again, we can sample using POS-EX(m), where m is obtained by using Lemma 1 and little additional cost. For any class \mathcal{L} of languages, we get:

Lemma 3. *If \mathcal{L} contains 2 languages L_1 and L_2 s.t. $L_1 \cap L_2 \neq \emptyset$, $L_1 \not\subseteq L_2$ and $L_2 \not\subseteq L_1$, then \mathcal{L} is not polynomially PAC-learnable from positive examples only.*

Proof. Let $w_1 \in L_1 - L_2$, $w_2 \in L_2 - L_1$ and $w_3 \in L_1 \cap L_2$. Consider the distribution \mathcal{D}_1 s.t. $\Pr_{\mathcal{D}_1}(w_1) = \Pr_{\mathcal{D}_1}(w_3) = \frac{1}{2}$ and the distribution \mathcal{D}_2 s.t. $\Pr_{\mathcal{D}_2}(w_2) = \Pr_{\mathcal{D}_2}(w_3) = \frac{1}{2}$. Basically, if one learns either L_1 from positive examples drawn according to \mathcal{D}_2 , or L_2 from positive examples drawn according to \mathcal{D}_1 , only the string w_3 will be used. However, the error will be $\geq \frac{1}{2}$. \square

Theorem 3. *(1) $\mathcal{GB}(\Sigma)$ and (2) $\mathcal{DFA}(\Sigma)$ are not polynomially PAC-learnable from positive examples only.*

Proof. An immediate consequence of Lemma 3 with $L_1 = B_1(a)$, $L_2 = B_1(b)$, $w_1 = aa$, $w_2 = bb$ and $w_3 = ab$. \square

4 Query Learning

Learning from queries involves the Learner (he) being able to interrogate the Oracle (she) using queries from a given set [8]. The goal of the Learner is to identify a grammar of an unknown language L . The Oracle knows L and properly answers to the queries (*i.e.*, she does not lie). Below, we will use three kinds of queries. With the Membership Queries (MQ), the Learner submits a string w to the Oracle and she answers YES if $w \in L$, NO otherwise. With the Equivalence Queries (EQ), he submits (the grammar of) a language K and she answers YES if $K = L$, and a string belonging to $K \oplus L$ otherwise. With the Correction Queries based on the Edit Distance (CQ_{EDIT}), he submits a string w and she answers YES if $w \in L$, and any correction $z \in L$ at minimum edit distance of w otherwise.

Definition 3. *A class \mathcal{G} is polynomially identifiable from queries if_{def} there is an algorithm \mathfrak{A} able to identify every $G \in \mathcal{G}$ s.t. at any call of a query, the total number of queries and of time used up to that point by \mathfrak{A} is polynomial both in $\|G\|$ and in the size of the information presented up to that point by the Oracle.*

In the case of good balls, we have shown:

Theorem 4 ([20]). *(1) $\mathcal{GB}(\Sigma)$ is not polynomially identifiable from MQ and EQ. (2) $\mathcal{GB}(\Sigma)$ is polynomially identifiable from CQ_{EDIT} .*

Notice however that if the Learner is given one string from a good ball, then he can learn using a polynomial number of MQ only.

Concerning the class of the DFA, we get:

Theorem 5. (1) $\mathcal{DFA}(\Sigma)$ is polynomially identifiable from MQ and EQ [21], but is not from (2) MQ only [8], nor (3) EQ only [10], nor (4) CQ_{EDIT} only.

Proof. (4) Let \mathcal{A}_w denote the DFA that recognizes $\Sigma^* \setminus \{w\}$. Let $n \in \mathbb{N}$ and $\mathcal{DFA}_{\leq n} = \{\mathcal{A}_w : w \in \Sigma^{\leq n}\}$. Following [10], we describe an Adversary that maintains a set X of all the possible DFA. At the beginning, $X = \mathcal{DFA}_{\leq n}$. Each time the correction of any string w is demanded, the Adversary answers YES and eliminates only one DFA of X : \mathcal{A}_w . As there is $\Omega(|\Sigma|^n)$ DFA in $\mathcal{DFA}_{\leq n}$, identifying one of them will require $\Omega(|\Sigma|^n)$ queries in the worst case. \square

5 Polynomial Identification in the Limit

Identification in the limit [7] is standard: a Learner receives an infinite sequence of information (presentation) that should help him to find the grammar $G \in \mathcal{G}$ of an unknown target language $L \in \mathcal{L}$. The set of admissible presentations is denoted by **Pres**, each presentation being a function $\mathbb{N} \rightarrow X$ where X is any set. Given $\mathbf{f} \in \mathbf{Pres}$, we will write \mathbf{f}_m for the $m + 1$ first elements of \mathbf{f} , and $\mathbf{f}(n)$ for its n^{th} element. Below, we will consider two sorts of presentations. When **Pres**=TEXT, all the strings in L are presented: $\mathbf{f}(\mathbb{N}) = \mathbb{L}(G)$. When **Pres**=INFORMANT, a presentation is of labelled pairs (w, l) where $(w \in L \Rightarrow l = +)$ and $(w \notin L \Rightarrow l = -)$: $\mathbf{f}(\mathbb{N}) = \mathbb{L}(G) \times \{+\} \cup \overline{\mathbb{L}(G)} \times \{-\}$; we will write **Pres** = PRESENTATION for all the results that concern both TEXT and INFORMANT.

Definition 4. We say that \mathcal{G} is identifiable in the limit from **Pres** if_{def} there exists an algorithm \mathfrak{A} s.t. for all $G \in \mathcal{G}$ and for any presentation \mathbf{f} of $\mathbb{L}(G)$, there exists a rank n s.t. for all $m \geq n$, $\mathfrak{A}(\mathbf{f}_m) = \mathfrak{A}(\mathbf{f}_n)$ and $\mathbb{L}(\mathfrak{A}(\mathbf{f}_n)) = \mathbb{L}(G)$.

This definition yields a number of learnability results. However, the absence of efficiency constraints often leads to unusable algorithms. Firstly, it seems reasonable that the amount of time an algorithm has to learn should be bounded:

Definition 5 (Polynomial Update Time). An algorithm \mathfrak{A} is said to have polynomial update time if_{def} there is a polynomial $p()$ s.t., for every presentation \mathbf{f} and every integer n , computing $\mathfrak{A}(\mathbf{f}_n)$ requires $\mathcal{O}(p(\|\mathbf{f}_n\|))$ time.

It is known that polynomial update time is not sufficient [11]: a Learner could receive an exponential number of examples without doing anything but wait, and then use the amount of time he saved to solve any \mathcal{NP} -hard problem... Polynomiality should also concern the minimum amount of data that any Learner requires:

Definition 6 (Polynomial Characteristic Sample). We say that \mathcal{G} admits polynomial characteristic samples if_{def} there exist an algorithm \mathfrak{A} and a polynomial $p()$ s.t. for all $G \in \mathcal{G}$, there exists $\text{Cs} \subseteq X$ s.t. (1) $\|\text{Cs}\| \leq p(\|G\|)$, (2) $\mathbb{L}(\mathfrak{A}(\text{Cs})) = \mathbb{L}(G)$ and (3) for all $\mathbf{f} \in \mathbf{Pres}$, for all $n \geq 0$, if $\text{Cs} \subseteq \mathbf{f}_n$ then $\mathfrak{A}(\mathbf{f}_n) = \mathfrak{A}(\text{Cs})$. Such a set Cs is called a characteristic sample of G for \mathfrak{A} . If \mathfrak{A} exists, we say that \mathcal{G} is identifiable in the limit in Cs polynomial time.

Lastly, polynomiality may concern either the number of implicit prediction errors [11] or the number of mind changes (MC) [28] done by the learner:

Definition 7 (Implicit Prediction Errors). *We say that an algorithm \mathfrak{A} makes an implicit prediction error (IPE) at time n of a presentation \mathbf{f} if_{def} $\mathfrak{A}(\mathbf{f}_{n-1}) \not\vdash \mathbf{f}(n)$. \mathfrak{A} is called consistent if_{def} it changes its mind each time a prediction error is detected with the new presented element.*

\mathfrak{A} identifies \mathcal{G} in the limit in IPE polynomial time if_{def} (1) \mathfrak{A} identifies \mathcal{G} in the limit, (2) \mathfrak{A} has polynomial update time and (3) \mathfrak{A} makes a polynomial number of implicit prediction errors: let $\#\text{IPE}(\mathbf{f}) = |\{k \in \mathbb{N} : \mathfrak{A}(\mathbf{f}_k) \not\vdash \mathbf{f}(k+1)\}|$; there exists a polynomial $p()$ s.t. for each $G \in \mathcal{G}$ and each presentation \mathbf{f} of $\mathbb{L}(G)$, $\#\text{IPE}(\mathbf{f}) \leq p(\|\mathbf{f}\|)$.

Definition 8 (Mind Changes). *We say that an algorithm \mathfrak{A} changes its mind (MC) at time n of presentation \mathbf{f} if_{def} $\mathfrak{A}(\mathbf{f}_n) \neq \mathfrak{A}(\mathbf{f}_{n-1})$. \mathfrak{A} is called conservative if_{def} it never changes its mind when the current hypothesis is consistent with the new presented element.*

\mathfrak{A} identifies \mathcal{G} in the limit in MC polynomial time if_{def} (1) \mathfrak{A} identifies \mathcal{G} in the limit, (2) \mathfrak{A} has polynomial update time and (3) \mathfrak{A} makes a polynomial number of mind changes: Let $\#\text{MC}(\mathbf{f}) = |\{k \in \mathbb{N} : \mathfrak{A}(\mathbf{f}_k) \neq \mathfrak{A}(\mathbf{f}_{k+1})\}|$; there exists a polynomial $p()$ s.t. for each $G \in \mathcal{G}$ and each presentation \mathbf{f} of $\mathbb{L}(G)$, $\#\text{MC}(\mathbf{f}) \leq p(\|\mathbf{f}\|)$.

Concerning both last notions, one can notice that if an algorithm \mathfrak{A} is consistent then $\#\text{IPE}(\mathbf{f}) \leq \#\text{MC}(\mathbf{f})$ for every presentation \mathbf{f} . Likewise, if \mathfrak{A} is conservative then $\#\text{MC}(\mathbf{f}) \leq \#\text{IPE}(\mathbf{f})$. So we deduce the following lemma:

Lemma 4. *If \mathfrak{A} identifies the class \mathcal{G} in MC polynomial time and is consistent, then \mathfrak{A} identifies \mathcal{G} in IPE polynomial time. Conversely, if \mathfrak{A} identifies \mathcal{G} in IPE polynomial time and is conservative, then \mathfrak{A} identifies \mathcal{G} in MC polynomial time.*

5.1 Polynomial Identification from Text

The aim of this section is to show the following result:

Theorem 6. $\mathcal{GB}(\Sigma)$ is identifiable in the limit from TEXT in (1) MC polynomial time, (2) IPE polynomial time and (3) CS polynomial time.

Notice that as the DFA recognize a *superfinite* class of languages (*i.e.*, containing all the finite languages and at least one infinite language), it is impossible to identify the class in the limit from positive examples only:

Theorem 7 ([7]). $\mathcal{DFA}(\Sigma)$ is not identifiable in the limit from TEXT.

In order to prove Theo. 6, we will need to build the minimum consistent good ball containing a set $X = \{x_1, \dots, x_n\}$ of strings (sample). This problem is \mathcal{NP} -hard but some instances are efficiently solvable. Let X^{\max} (*resp.* X^{\min}) denote the set of all longest (*resp.* shortest) strings of X . A *minimality fingerprint* is a subset $\{u, v, w\} \subseteq X$ s.t. (1) $u, v \in X^{\max}$, (2) $w \in X^{\min}$, (3) $|u| - |w| = 2r$ for

some $r \in \mathbb{N}$, (4) u and v have only one longest common subsequence, that is, $lcs(u, v) = \{o\}$ for some $o \in \Sigma^*$, (5) $|o| = |u| - r$ and (6) $X \subseteq B_r(o)$.

Checking if X contains a minimality fingerprint, and computing o and r can be achieved in polynomial time (in $\|X\|$). Indeed, the only critical point is that the cardinal of $lcs(u, v)$ may be $> 1.442^n$ [29] (where $n = |u| = |v|$); nevertheless, a data structure such as the LCS-graph [30] allows one to conclude polynomially. Moreover, the minimality fingerprints are meaningful. Indeed, only the properties of the edit distance are needed to show that if X contains a minimality fingerprint $\{u, v, w\}$ for the ball $B_r(o)$ and $X \subseteq B_{r'}(o')$, then either $r' > r$, or ($r' = r$ and $o' = o$). In other words, $B_r(o)$ is the smallest ball containing X *w.r.t.* the radius.

We can now consider Algo. 1. This algorithm does identify $\mathcal{GB}(\Sigma)$ in the limit since if $B_r(o)$ denotes the target ball, then at some point, the algorithm will meet the strings $u = a^r o$, $v = b^r o$, and some w of length $|o| - r$ that constitute a minimality fingerprint for $B_r(o)$. Moreover, it obviously has a polynomial update time. Finally, it makes a polynomial number of MC. Indeed, it only changes its hypothesis in favour of a valid ball if the ball has a larger radius than all the valid balls it has ever conjecture, that may happen $\leq r$ times. And it only changes its hypothesis in favour of a junk ball if either it must abandon a valid ball, or if the actual junk ball does not contain all the examples, that may happen $\leq r + 2r$ times. So the total number of MC is $\leq 4r$. So Claim (1) holds.

Concerning Claim (2), note that Algo. 1 is consistent (thanks to the use of the junk balls), thus Claim (2) holds by Lemma 4. Lastly, every minimality fingerprint is a characteristic set that makes Algo. 1 converge, so Claim (3) holds.

Algorithm 1. Identification of good balls from text.

```

Data: A text  $f = \{x_1, x_2, \dots\}$ 
read( $x_1$ );  $c \leftarrow x_1$ ; output ( $x_1, 0$ );
while true do
  | read( $x_i$ );
  | if  $f_i$  is a minimality fingerprint for  $B_r(o)$  then
  |   | output ( $o, r$ ) (* valid ball *)
  | else
  |   | if  $c \notin f_i^{nax}$  then  $c \leftarrow$  any string in  $f_i^{nax}$ ;
  |   | output ( $c, |c|$ ) (* junk ball *)
  | end
end

```

5.2 Polynomial Identification from Informant

Theorem 8. (1) $\mathcal{GB}(\Sigma)$ is not identifiable from INFORMANT in IPE polynomial time, but is identifiable in (2) MC polynomial time and (3) CS polynomial time.

Proof. (1) Similar proof as that of [11] for the DFA: if $\mathcal{GB}(\Sigma)$ was identifiable in IPE polynomial time from INFORMANT, then $\mathcal{GB}(\Sigma)$ would be polynomially identifiable from EQ, that contradicts Theo. 4. (2) As the hypotheses are not

necessarily consistent with the data, one can use Algo. 1, ignoring the negative examples. (3) Same characteristic sets as those of Theo. 6, Claim (3). \square

Theorem 9. (1) $\mathcal{DFA}(\Sigma)$ is not identifiable from INFORMANT in IPE polynomial time [11], but is identifiable in (2) MC polynomial time and (3) CS polynomial time [31,32].

Let us prove Claim (2) with minimality fingerprints again. We say that $X = \langle X_+, X_- \rangle$ contains a *minimality fingerprint* *if_{def}* the following conditions hold: (1) let $A = \langle \Sigma, Q, q_0, F, \delta \rangle$ be the DFA computed by RPNI [32] on X possibly completed with one hole state; (2) for all $q \in Q$, the smallest string w_q *w.r.t.* the hierarchical order \preceq *s.t.* $\delta(q_0, w_q) = q$ belongs to either X_+ if $q \in F$, or X_- if $q \notin F$; (3) for all $q \in Q, a \in \Sigma$, $w_q a$ belongs to either X_+ if $\delta(q, a) \in F$, or X_- if $\delta(q, a) \notin F$; (4) for all $p, q, r \in Q, a \in \Sigma$ *s.t.* $\delta(p, a) = q \neq r$, there exists $f \in \Sigma^*$ *s.t.* either $(w_p a f \in X_+, w_r f \in X_-)$, or $(w_p a f \in X_-, w_r f \in X_+)$.

Notice that not all the DFA have minimality fingerprints. Moreover, every fingerprint contains a characteristic sample of A for RPNI [32], plus new information: actually, *all* the states and the transitions of A are determined by the fingerprint, so any other complete DFA A' compatible with X necessarily has more states than A . A is thus the unique complete minimal DFA compatible with X . Lastly, checking if X contains a minimality fingerprint, and computing A is achievable in polynomial time (in $\|X\|$).

We can now define a Learner \mathfrak{A} . At each step, \mathfrak{A} tests if \mathbf{f}_i contains a minimality fingerprint. If yes, \mathfrak{A} changes its mind in favour of the DFA A_i returned by RPNI on \mathbf{f}_i . If no, \mathfrak{A} returns the previous hypothesis (that may not be consistent). Clearly, the number of states of A_i strictly increases (thanks to the fingerprints). As this number is bounded by the number of states of the target DFA A , we get $\#\text{MC}(\mathbf{f}) \leq \|A\|$. Moreover, at some point, a fingerprint (thus a characteristic set of A) will appear in the data, and then \mathfrak{A} will converge.

6 Conclusion

In this paper, we have performed a systematic study of two classes of languages whose definitions is based on very different principles. Table 1 summarize our results. Those marked with a \dagger were proved in this article.

Clearly, the goal of this work was not to show that any paradigm is equivalent or better than any other: comparing two classes is not sufficient. Nevertheless, we have shown that several hints were wrong. For instance, it is wrong to think that the identification in MC polynomial time implies the identification in IPE polynomial time. It is also wrong to think that it is easier to learn from positive and negative examples (INFORMANT) than from positive examples only (TEXT) (because in some paradigm, misclassifying negative examples is expensive in terms of complexity).

In Table 1, we also show (without proof, due to the lack of space), the results that concern all the balls including those that are not good. Let us recall that a ball $B_r(o)$ is good *if_{def}* $r \leq |o|$, so for a bad ball, it is possible that $r \gg 2^{|o|}$. In

Table 1. A synthetic view of the results presented in this paper. [†] marks the theorems proved above. Due to the lack of space, we just claim the results concerning the general balls of $\mathcal{BALC}(\Sigma)$.

Criterion	$\mathcal{GB}(\Sigma)$		$\mathcal{DFA}(\Sigma)$		$\mathcal{BALC}(\Sigma)$
PAC INFORM.	NO [†]	Theo. 1	NO	Theo. 2	No
PAC TEXT	NO [†]	Theo. 3 (1)	NO [†]	Theo. 3 (2)	No
IPE INFORM.	NO [†]	Theo. 8 (1)	NO	Theo. 9 (1)	No
IPE TEXT	YES [†]	Theo. 6 (2)	NO	Theo. 7	No
MC INFORM.	YES [†]	Theo. 8 (2)	YES [†]	Theo. 9 (2)	YES
MC TEXT	YES [†]	Theo. 6 (1)	NO	Theo. 7	No
CS INFORM.	YES [†]	Theo. 8 (3)	YES	Theo. 9 (3)	No
CS TEXT	YES [†]	Theo. 6 (3)	NO	Theo. 7	No
MQ (or EQ)	NO	Theo. 4 (1)	NO	Theo. 5 (2,3)	No
MQ and EQ	NO	Theo. 4 (1)	YES	Theo. 5 (1)	No
CQ _{EDIT}	YES	Theo. 4 (2)	NO	Theo. 5 (4)	No

this case, the longest strings of $B_r(o)$, whose length is $|o| + r$, which delimitate the upper boarder of $B_r(o)$, are not polynomially related to the size of the ball ($|o| + \log r$). So the picture is the same as that of the *non deterministic automata* for which one has to consider strings of exponential length in order to distinguish two states [17,33]. Hence, studying the learnability of all the balls is a way to explore the limits of the paradigms, reached when an algorithm cannot get round of exponential strings anymore.

Finally, one can note that the good balls are not learnable from a polynomial number of MQ and EQ, that is the case of the DFA. As the balls are finite languages, they are recognizable with DFA. Thus a subclass of a learnable class could be non learnable! We conjecture, following [5,2], that this is because the size of the minimal DFA recognizing a ball is exponential in the size of the ball.

References

1. Levenshtein, V.I.: Binary codes capable of correcting deletions, insertions, and reversals. Doklady Akademii Nauk SSSR 163(4), 845–848 (1965)
2. Navarro, G.: A guided tour to approximate string matching. ACM computing surveys 33(1), 31–88 (2001)
3. Chávez, E., Navarro, G., Baeza-Yates, R.A., Marroquín, J.L.: Searching in metric spaces. ACM Computing Survey 33(3), 273–321 (2001)
4. Kohonen, T.: Median strings. Pattern Recognition Letters 3, 309–313 (1985)
5. Schulz, K.U., Mihov, S.: Fast string correction with Levenshtein automata. Int. Journal on Document Analysis and Recognition 5(1), 67–85 (2002)
6. Sagot, M.F., Wakabayashi, Y.: Pattern inference under many guises. In: Recent Advances in Algorithms and Combinatorics, pp. 245–287. Springer, Heidelberg (2003)

7. Gold, E.M.: Language identification in the limit. *Information and Control* 10(5), 447–474 (1967)
8. Angluin, D.: Queries and concept learning. *Machine Learning Journal* 2, 319–342 (1987)
9. Valiant, L.G.: A theory of the learnable. *Communications of the ACM* 27(11), 1134–1142 (1984)
10. Angluin, D.: Negative results for equivalence queries. *Machine Learning Journal* 5, 121–150 (1990)
11. Pitt, L.: Inductive inference, DFA’s, and computational complexity. In: Jantke, K.P. (ed.) AII 1989. LNCS, vol. 397, pp. 18–44. Springer, Heidelberg (1989)
12. Li, M., Vitanyi, P.: Learning simple concepts under simple distributions. *Siam Journal of Computing* 20, 911–935 (1991)
13. Denis, F.: Learning regular languages from simple positive examples. *Machine Learning Journal* 44(1), 37–66 (2001)
14. Parekh, R.J., Honavar, V.: On the relationship between models for learning in helpful environments. In: Oliveira, A.L. (ed.) ICGI 2000. LNCS (LNAI), vol. 1891, pp. 207–220. Springer, Heidelberg (2000)
15. Haussler, D., Kearns, M.J., Littlestone, N., Warmuth, M.K.: Equivalence of models for polynomial learnability. *Information and Computation* 95(2), 129–161 (1991)
16. Kearns, M., Valiant, L.: Cryptographic limitations on learning boolean formulae and finite automata. In: 21st ACM Symposium on Theory of Computing (STOC 1989), pp. 433–444 (1989)
17. de la Higuera, C.: Characteristic sets for polynomial grammatical inference. *Machine Learning Journal* 27, 125–138 (1997)
18. Wagner, R., Fisher, M.: The string-to-string correction problem. *Journal of the ACM* 21, 168–178 (1974)
19. Papadimitriou, C.M.: *Computational Complexity*. Addison Wesley, New York (1994)
20. Becerra-Bonache, L., de la Higuera, C., Janodet, J.C., Tantini, F.: Learning balls of strings with correction queries. In: Kok, J.N., Koronacki, J., Lopez de Mantaras, R., Matwin, S., Mladenić, D., Skowron, A. (eds.) ECML 2007. LNCS (LNAI), vol. 4701, pp. 18–29. Springer, Heidelberg (2007)
21. Angluin, D.: Learning regular sets from queries and counterexamples. *Information and Control* 39, 337–350 (1987)
22. Warmuth, M.: Towards representation independence in PAC-learning. In: Jantke, K.P. (ed.) AII 1989. LNCS, vol. 397, pp. 78–103. Springer, Heidelberg (1989)
23. Kearns, M., Vazirani, U.: *An Introduction to Computational Learning Theory*. MIT Press, Cambridge (1994)
24. Pitt, L., Valiant, L.G.: Computational limitations on learning from examples. *Journal of the ACM* 35(4), 965–984 (1988)
25. Maier, D.: The complexity of some problems on subsequences and supersequences. *Journal of the ACM* 25, 322–336 (1977)
26. de la Higuera, C., Casacuberta, F.: Topology of strings: Median string is NP-complete. *Theoretical Computer Science* 230, 39–48 (2000)
27. Pitt, L., Warmuth, M.: The minimum consistent DFA problem cannot be approximated within any polynomial. *Journal of the ACM* 40(1), 95–142 (1993)
28. Angluin, D., Smith, C.: Inductive inference: theory and methods. *ACM computing surveys* 15(3), 237–269 (1983)
29. Greenberg, R.I.: Bounds on the number of longest common subsequences. Technical report, Loyola University (2003), <http://arXiv.org/abs/cs/0301030v2>

30. Greenberg, R.I.: Fast and simple computation of all longest common subsequences. Technical report, Loyola University (2002), <http://arXiv.org/abs/cs.DS/0211001>
31. Gold, E.M.: Complexity of automaton identification from given data. *Information and Control* 37, 302–320 (1978)
32. Oncina, J., García, P.: Identifying regular languages in polynomial time. In: *Advances in Structural and Syntactic Pattern Recognition. Series in Machine Perception and Artificial Intelligence*, vol. 5, pp. 99–108. World Scientific, Singapore (1992)
33. Denis, F., Lemay, A., Terlutte, A.: Learning regular languages using RFSA. *Theoretical Computer Science* 313(2), 267–294 (2004)