

Document downloaded from:

<http://hdl.handle.net/10251/37983>

This paper must be cited as:

García Gómez, P.; López Rodríguez, D.; Vázquez-De-Parga Andrade, M. (2012). Polynomial characteristic sets for DFA identification. *Theoretical Computer Science*. 448:41-46. doi:10.1016/j.tcs.2012.04.042.



The final publication is available at

<http://dx.doi.org/doi:10.1016/j.tcs.2012.04.042>

Copyright Elsevier

Polynomial characteristic sets for DFA identification*

Pedro García, Damián López and Manuel Vázquez de Parga

Departamento de Sistemas Informáticos y Computación.
Universidad Politécnica de Valencia. Valencia (Spain).
email: {pgarcia,dlopez,mvazquez}@dsic.upv.es

Abstract

We study the order in *Grammatical Inference* algorithms, and its influence on the polynomial (with respect to the data) identification of languages. This work is motivated by recent results on the polynomial convergence of data-driven grammatical inference algorithms. In this paper, we prove a sufficient condition that assures the existence of a characteristic sample whose size is polynomial with respect to the minimum DFA of the target language.

Keywords: DFA identification; Grammatical inference.

1 Introduction

A *Grammatical Inference* (GI) algorithm is a method that tries to obtain a representation of a target language L from some information about L [1, 2, 3]. In this work we study the inference of deterministic automata for regular string languages using complete presentation (samples that belong or not to the target language) [4, 5, 6, 7, 8]. We note here that there are other results that, using the same data presentation, tackle the inference of non-deterministic automata [9, 10] as well as results that study the identification of languages using queries [11, 12, 13].

A common approach to the GI of regular string languages takes into account an initial machine that represents the input data. Some states of this machine are then merged in order to obtain some generalization. It is worth to be noted that the generalization obtained depends both on the data supplied as well as on the order in which the states of the machine are traversed.

An important issue to determine whether a given GI algorithm has good behaviour or not, is the amount of information needed to identify the target

*Work partially supported by the Spanish Ministerio de Economía y Competitividad under research project TIN2011-28260-C03-01

language. In that way it is important the concept of *characteristic sample* of a target language for an inference algorithm. Given a class of languages H , the characteristic sample for an algorithm A and a language $L \in H$ is defined as a set of words in L (usually denoted by D_+) and a set of words not in L (denoted by D_-) such that: whenever the algorithm A is run with input (D_+, D_-) the algorithm outputs a correct representation of L ; this representation does not change even though more words are added to the input.

The model of learning called identification in the limit by Gold, establishes that an algorithm identifies a class of languages H *in the limit* if and only if every language in the class has associated a characteristic set for that algorithm [14]. Taking into account Gold's work, as well as results by Pitt [15] and Angluin[16], de la Higuera [17] defines *polynomial time and data identifiability*, as an extension of Gold's definitions which consider characteristic sets of polynomial size.

As mentioned above, the most referred results on inference of regular languages from complete presentation propose algorithms based on the merging of states: the *RPNI* algorithm proposed by Oncina and García [7] and the *Blue-Fringe* algorithm by Lang et al. [8]. In both approaches, from an initial representation of the input information (usually a Moore machine of the training set), the algorithms consider a fixed order to traverse of the states of the machine, usually the canonical order. In [18], de la Higuera et al. propose an algorithmic scheme able to consider a broader set of orderings (the chosen order is an input parameter of the method), including any *data-driven* one. In that work, the order in which the states of the initial representation are promoted (considered as states of the automaton to be output) is consequence of the order in which the states are merged. In fact, the authors do not distinguish among both orders. In that article, it is proved that, whenever the (merging) order does not consider the input data, then, for that algorithm, there exists a polynomial characteristic sample for any language. The authors also show that, given any size of automata, there is at least one data-driven order for which it is possible to find automata with non-polynomial characteristic set. This seems to indicate that any *interesting* data-driven order would imply an exponential characteristic sample.

All these results led to the GI community to think that *Blue-Fringe* algorithm (which was somewhat inspired by the results in [18]) to have no polynomial characteristic set, because *Blue-Fringe* merging order is in fact data-driven. Despite this assumption, *Blue-Fringe* became the *state of art* algorithm because its experimental behaviour in practical tasks, which outperformed the results of previous approaches. In fact, the merging order used by *Blue-Fringe* led this algorithm to be more data-efficient with respect other GI algorithms, even when the training set does not include a characteristic set for the target language.

Nevertheless, it has been recently proved that *Blue-Fringe* algorithm has a polynomial characteristic set [19]. The proof of this, takes into account that the *Blue-Fringe* promotion order does not depend on the merging order. This last result motivates the study of the order influence in the polynomial convergence of GI algorithms.

2 Definitions and notation

Let Σ be a finite alphabet and let Σ^* be the set of possible strings over Σ . Let also λ denote the empty string. A *language* L over Σ is a subset of Σ^* . For any given set q , we will denote the cardinality of q with $|q|$. Given $x \in \Sigma^*$, if $x = uv$ with $u, v \in \Sigma^*$, then u (resp. v) is called *prefix* (resp. *suffix*) of x . Let us denote with $Pr(L)$ the set of prefixes of L .

A *Deterministic Finite Automaton (DFA)* is a 5-tuple $A = (Q, \Sigma, \delta, q_0, F)$, where Q is a finite set of states, Σ is an alphabet, $q_0 \in Q$ is the initial state, $F \subseteq Q$ is the set of final states and $\delta : Q \times \Sigma \rightarrow Q$ is the transition function. The language accepted by an automaton A is denoted $L(A)$.

Given any *DFA* A , any two states p and q of A are usually said to be equivalent if and only if, for any $x \in \Sigma^*$, $\delta(p, x) \in F$ if and only if $\delta(q, x) \in F$. This relation allows to obtain the *minimal DFA* for the language $L(A)$ (the deterministic automaton with the smallest set of states that accept the language). We note that this minimal automaton is unique up to isomorphism.

A Moore machine is a 6-tuple $M = (Q, \Sigma, \Gamma, \delta, q_0, \Phi)$, where Σ (resp. Γ) is the input (resp. output) alphabet, δ is a partial function that maps $Q \times \Sigma$ in Q and Φ is a function that maps Q in Γ called *output function*. The function δ can be extended in a natural way to consider strings over Σ .

Throughout this paper, the behavior of M will be given by the partial function $t_M : \Sigma^* \rightarrow \Gamma$ defined as $t_M(x) = \Phi(\delta(q_0, x))$ for every $x \in \Sigma^*$ such that $\delta(q_0, x)$ is defined.

A *DFA* $A = (Q, \Sigma, \delta, q_0, F)$ can be simulated by a Moore machine $M = (Q, \Sigma, \{+, -\}, \delta, q_0, \Phi)$, where $\Phi(q) = +$ if $q \in F$ and $\Phi(q) = -$ otherwise. Then, the language defined by M is $L(M) = \{x \in \Sigma^* : \Phi(\delta(q_0, x)) = +\}$.

Given two disjoint finite sets of strings D_+ and D_- , we define the (D_+, D_-) -*prefix tree Moore machine (PTMM(D_+, D_-))* as the Moore machine having $\Gamma = \{+, -, ?\}$, $Q = Pr(D_+ \cup D_-)$, $q_0 = \lambda$ and $\delta(u, a) = ua$ if $u, ua \in Q$ and $a \in \Sigma$. For every state u , the value of the output function associated to u is $+$, $-$ or $?$ (undefined) depending whether u belongs to D_+ , to D_- or to $Q - (D_+ \cup D_-)$ respectively.

A Moore machine $M = (Q, \Sigma, \{+, -, ?\}, \delta, q_0, \Phi)$ is *consistent* with (D_+, D_-) if $\forall x \in D_+$ we have $t_M(x) = +$ and $\forall x \in D_-$ we have $t_M(x) = -$.

Algorithm 3.1 Inference from complete presentation. A general scheme.

Input: $M = PTMM(D_+, D_-) = (Q, \Sigma, \{0, 1, ?\}, \delta, q_0, \Phi)$;

Input: An order among the states of the input $PTMM$

Output: A Moore Machine consistent with respect to the input data

```
1: Method
2:  $red = \{\lambda\}$ 
3:  $blue = \{q \in Q : q = \delta(p, a), p \in red \wedge a \in \Sigma\} - red$ 
4: while  $blue \neq \emptyset$  do
5:    $NonEqStatesList = MergibleStatesList = \emptyset$ 
6:   for all  $q \in blue$  /* traversed following the given order */ do
7:      $merged = False$ 
8:     for all  $p \in red$  /* traversed following the given order */ do
9:       if  $(p, q)$  are mergible then
10:         $AppendTo(MergibleStatesList, (p, q))$ 
11:         $merged = True$ 
12:       end if
13:     end for
14:     if not merged then  $AppendTo(NonEqStatesList, q)$  end if
15:   end for
16:   Set option value among  $\{merge, promote\}$ 
17:   if option = promote then
18:      $red = red \cup \{q\} : q = First(NonEqStatesList)$ 
19:   else
20:     Let  $(p, q) \in MergibleStatesList$  be chosen following whichever cri-
      terion, and deterministically merge them
21:   end if
22:    $blue = \{q \in Q : q = \delta(p, a), p \in red \wedge a \in \Sigma\} - red$ 
23: end while
24: Return( $M$ );
25: End Method.
```

3 Convergence of GI algorithms using polynomial data

In this section we will prove a condition that, when fulfilled by a GI algorithm, assures the existence of a polynomial characteristic set.

In order to do so, we first propose a general framework that unifies all the previous results. This framework will be also useful to prove the main result

3.1 A general framework

Algorithm 3.1 summarizes the most relevant GI algorithms. Please note that this algorithm puts aside efficiency. We now show that this algorithm can be considered as a general framework able to implement any previous result.

Note that, in order to implement *RPNI* algorithm with the proposed scheme, it is enough to take into account which list (*NonEqStatesList* or *MergibleStatesList*) has been first updated in order to choose among promotion or merge (line 16). Note also that, in order to implement *Blue-Fringe* algorithm, it is enough to consider the merge of states whenever the list of promotable states is empty.

Let us recall here the algorithm scheme proposed in [18]. The fact that this scheme does not use a *PTMM* to represent the training data can be considered secondary. In order to avoid extra notation, we will describe it in terms of a *PTMM* representation of the training set.

In any given iteration of that scheme, a score for every $p \in red$ and $q \in blue$ is obtained, no matter whether the merge is possible or not. This set of scores drives the traverse of the pairs of *red* and *blue* states. In the traverse of the pairs of states, it is checked whether or not the merge of the states is possible. If so, the pair is merged and the *blue* set is updated. If the merge is not possible, the algorithm checks if the state q has been considered to be merged with all states in the *red* set. In that case, the state q is promoted to *red*, otherwise, the *blue* set is updated and the set of scores is reseted. The algorithm ends when the *blue* set is empty.

Note that, in order to use Algorithm 3.1 to implement the algorithm by de la Higuera et al., it is necessary to compute the score among the *red* and *blue* states before the loop that traverses the *blue* states (line 6). The scores obtained guide the traverse of both *blue* and *red* sets, and therefore, the loops of lines 6 and 8 are also reduced to just one loop. The update of the *NonEqStatesList* (line 14) has also to be modified (to check whether or not all the possible merges have been considered), as well as included into the loop. Finally, the *option* flag is set to *merge* or *promote* (line 6) taking into account which list is updated first (in the same way the *RPNI* algorithm was previously adapted to our scheme).

3.2 A stronger result

We now will consider the proposed framework to prove a sufficient condition that assure the existence of a polynomial characteristic set for any given regular language. We first show that, for any regular language L , and whichever the order Algorithm 3.1 considers, it is possible to obtain a polynomial characteristic sample that identifies L .

The usual way to compute the characteristic set for any given language

L is based on the definition of the *minimal set of test states*. Thus, given $A = (Q, \Sigma, \delta, q_0, F)$ the minimum DFA for L , the set $S \subset \Sigma^*$ is a minimal set of test states if for every $q \in Q$ there exists only one string $x \in S$ such that $\delta(q_0, x) = q$.

Usually, for each state q , the set S contains the first string in canonical order that reaches q (note that, so defined, S is minimal ($\text{Card}(S) = \text{Card}(Q)$) and prefix closed). We note here that, for any order used by Algorithm 3.1, no state is considered before any of its prefixes. This follows from the fact that Algorithm 3.1 chooses a state among those in the *blue* set. In the following, we will consider only this kind of *effective* orders.

Whichever the effective order chosen, it can be used to obtain a prefix closed minimal set of test states.

Example 1 *Let the automaton in Figure 1. Let us also consider the alphabetical order.*

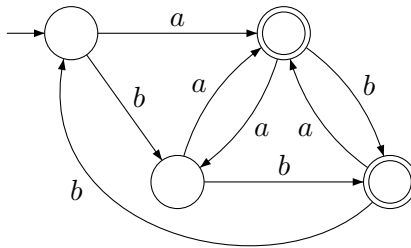


Figure 1: Automaton example.

In order to obtain the minimal prefix closed test states set it is necessary to find, for each state q , the first string that reaches q such that it does not visit a state twice. The alphabetic order does not give priority to shorter strings, thus, the minimal prefix closed set of test states for this example is $S = \{\lambda, a, aa, aab\}$. Note that the set of test states obtained according the canonic order would have been $S = \{\lambda, a, b, ab\}$.

Let us also note that, for instance, a set of test states using the alphabetical order is $S = \{\lambda, a, aa, aaab\}$, but it is not prefix closed. Note also that, to make this set closed under prefixes implies that it would be non-minimal.

Taking into account any prefix closed minimal set of test states, Algorithm 3.2 shows the way to obtain two sets $D_+(S)$ and $D_-(S)$. A rough bound for the size of $D_+(S) \cup D_-(S)$ is easily seen to be quadratic in the size of Q .

In this algorithm, the condition of the loop in line 6 refers to two *undistinguished* states, that is, two elements u and v in S such that there is no $w \in E$ such that just one of the strings uw or vw belongs to the language L . In order to find two such states, it is possible to use the matrix T and

Algorithm 3.2 Algorithm to obtain the characteristic set for a language L .

Input: The minimal DFA A for L

Output: The polynomial characteristic set for L

1: **Method**
2: Let S be the minimal set of test states for A
3: $E = \{\lambda\}$
4: Let $S' = S\Sigma - S$
5: Let T be a matrix indexed by the strings $u \in S \cup S'$ and $e \in E$ that stores the membership of the string ue to the language
6: **while** there exist two undistinguished $u, v \in S$ and a symbol $a \in \Sigma$ such that $T[ua, e] \neq T[va, e]$ for some $e \in E$ **do**
7: $E = E \cup \{ae\}$
8: **end while**
9: $(D_+, D_-) = \text{Data in } T$
10: Return(D_+, D_-);
11: **End Method.**

look for two identical rows indexed by elements in S . Example 2 illustrates this procedure.

Example 2 Let us also consider the automaton in Figure 1 and the prefix closed minimal set of test states $S = \{\lambda, a, aa, aab\}$.

Following table summarizes the process of obtaining the characteristic set. For the sake of clarity, we represent separately the elements in S and those in $S\Sigma - S$. Initially the only column available is the one with label λ . The 1 and 0 entries in the table represent if the strings obtained by concatenation of the strings that label the row and column belong or not to the language L .

	λ	b
λ	0	0
a	1	1
aa	0	1
aab	1	0
b	0	1
ab	1	0
aaa	1	1
$aaba$	1	1
$aabb$	0	0

Note that, taking into account just the column labelled λ , the undistinguished elements in S are $\{\lambda, aa\}$ and $\{a, aab\}$. It is possible to distinguish the first one using the suffix b . Once the table is filled in, all the elements in

S are distinguished, therefore, the sets $D_+(S)$ and $D_-(S)$ for the language are the following:

$$\begin{aligned} D_+(S) &= \{a, ab, bb, aaa, aab, aaab, aaba, aabab\} \\ D_-(S) &= \{\lambda, b, aa, abb, aabb, aabbb\} \end{aligned}$$

We now prove that, if the minimal set of test states was built using the same order Algorithm 3.1 uses, then, the sets $D_+(S)$ and $D_-(S)$ are a characteristic sample that identifies L .

Theorem 3 *Any GI algorithm, such that the promotion of states is independent from the input set, has a polynomial characteristic set, no matter the order followed to carry out the merge of states.*

Proof. *Let us consider any order over a finite subset of Σ^* as defined previously. Let $S \subset \Sigma^*$ be a prefix-closed minimal set of test states obtained according the defined order. Let also $D_+(S)$ and $D_-(S)$ be the positive and negative sets of strings obtained from S .*

Let us assume first that promotion has priority over merge. We first will prove that $\text{red} \subseteq S$ always hold.

Initially, $\text{red} = \{\lambda\}$ and $\text{blue} = \Sigma$. Let us consider any given iteration such that $\text{red} \subsetneq S$, then $\text{blue} = \text{red}\Sigma - \text{red} \subset (S\Sigma \cup S)$. Note that there is at least one state in $\text{blue} \cap S$ which, by construction of the characteristic sample, can be distinguished from any state in red . Let q denote the first of those states, q appears also the first in NonEqStateList and is promoted to red (line 18 in Algorithm 3.1). Eventually, all the elements in S will be promoted and thus $\text{red} = S$ and $\text{blue} = S\Sigma - S$. At that moment, again by the construction of the characteristic sample, each $q \in \text{blue}$ can be distinguished from any element in S but just one. Therefore, the criterion followed to merge the remaining states in blue is irrelevant.

Let us assume now that promotion has no priority over merge. We prove now that the order in which the merges are carried out affects only when the training set is not characteristic. Under this conditions, whenever $\text{red} \subset S$ it is fulfilled that $\text{blue} \subseteq S\Sigma \cup S$. By construction of the characteristic sample, for any state $q \in \text{blue}$ there is only one red state p such that the pair (p, q) is in $\text{MergibleStatesList}$. Once the merge of the pair of states (p, q) has been carried out, the red set does not change and the blue set continues being included into $S\Sigma \cup S$. \square

Example 4 *Let the automaton A in Figure 1 and the characteristic sample for the language $L(A)$ obtained in Example 2. The $\text{PTMM}(D_+, D_-)$ is the one shown in Figure 2.*

The symbol inside each state represents the output value for that state. Note also the figure shows also the numbering of the states according the alphabetical order. The red and blue sets are initialized respectively to $\{1\}$ and $\{2, 13\}$. The first state to analyze is state 2 and added to the

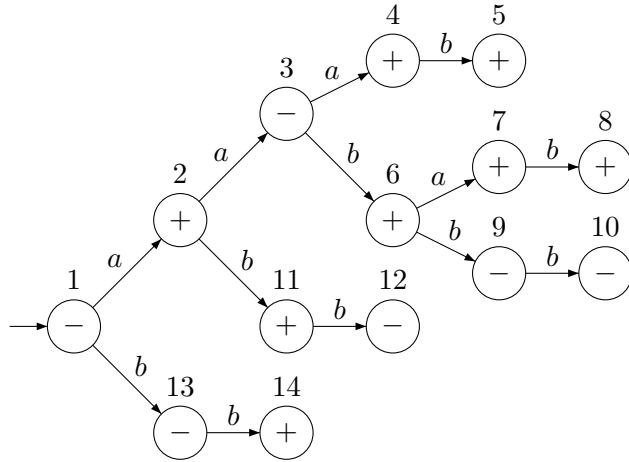


Figure 2: *PTMM* for the example (characteristic) training set.

NonEqStatesList as well as state 13 which is analyzed afterwards. Thus, the state 2 is promoted to the red set, and therefore, $\text{red} = \{1, 2\}$ and $\text{blue} = \{3, 11, 13\}$.

The second iteration analyses first the state 3 which is added to the list of *NonEqStatesList*, as well as state 11 and state 13. Thus, the algorithm only option is to promote the first state in *NonEqStatesList*, and thus, $\text{red} = \{1, 2, 3\}$ and $\text{blue} = \{4, 6, 11, 13\}$.

The next iteration starts analysing state 4 which is found to be mergible with state 2, and therefore it is added to the *MergibleStatesList*. State 6 is the next state took into account and it is added to *NonEqStatesList* because it is not equivalent to any state in red (as well as state 11). Last state analysed is state 13 which is found to be equivalent to state 3 and therefore added to *MergibleStatesList*. First, let us note that the two possible merges in this stage consider the pairs of states (2, 4) and (3, 13) (which will be of interest at the end of this example). In this run, we will choose first to promote rather than to merge. Therefore, the first state in *MergibleStatesList* is added to the red set, and thus $\text{red} = \{1, 2, 3, 6\}$ and $\text{blue} = \{4, 7, 9, 11, 13, \}$.

The analysis of the blue states carried out in the next iteration detects that the pair of states (2, 4) can be merged, as well as the pairs (2, 7), (1, 9), (6, 11) and (3, 13). It is worth to be noted here that: first, each blue state can only be merged with one red state; second that the possible merges detected in previous iterations are also considered in this last iteration, and therefore, when characteristic sample is used, it does not matter which choice is done in previous iterations because the output of the algorithm is always the same.

4 Conclusions

The experimental behaviour of *Blue-Fringe* algorithm proves that the merging order in a GI algorithm is important to obtain good results in applied tasks. This is mainly due to the fact that a guided order can take profit from *evidences* in the training set. In general, the consideration of a guided orders in a GI algorithm lead to a more data-efficient method.

The proof that *Blue-Fringe* algorithm has a polynomial characteristic set, which was assumed not to exist by the GI community, motivates this work. In this paper we prove a sufficient condition for GI algorithms to have polynomial characteristic sample. The result allows the consideration of any interesting data-driven criterion to establish the merging order of a GI algorithm. Thus, the use of *ad-hoc* orders in the application of GI algorithms to real tasks, under some conditions, could lead on the one hand to very efficient algorithms with respect to the data, and on the other hand, does not threaten the polynomial convergence which can be achieved.

References

- [1] Y. Sakakibara. Recent advances of grammatical inference. *Theoretical Computer Science*, 185:15–45, 1997.
- [2] C. de la Higuera. A bibliographical study of grammatical inference. *Pattern Recognition*, 38:1332–1348, 2005.
- [3] C. de la Higuera. *Grammatical Inference. Learning Automata and Grammars*. Cambridge University Press, 2010.
- [4] E. M. Gold. Complexity of automaton identification from given data. *Information and Control*, 37:302–320, 1978.
- [5] B.A. Trakhtenbrot and Ya. M. Barzdin. *Finite automata. Behavior and Synthesis*. North-Holland Pub. Co., 1973.
- [6] K.J. Lang. Random DFA’s can be approximately learned from sparse uniform examples. In *Proceedings of the fifth annual workshop on Computational learning theory*, pages 45–52, 1992.
- [7] J. Oncina and P. García. *Pattern recognition and image analysis*, volume 1, chapter Inferring regular languages in polynomial updated time, pages 49–61. World Scientific, 1992.
- [8] K. J. Lang, B. A. Pearlmutter, and R. A. Price. Results of the ab-badingo one dfa learning competition and a new evidence-driven state merging algorithm. *LNAI*, 1433:1–12, 1998. 4th International Colloquium, ICGI-98.

- [9] F. Denis, A. Lemay, and A. Terlutte. Learning regular languages using rfsa. *Theoretical Computer Science*, 313(2):267–294, 2004.
- [10] P. García, M. Vázquez de Parga, G. I. Álvarez, and J. Ruiz. Universal automata and NFA learning. *Theoretical Computer Science*, 407(1-3):192–202, 2008.
- [11] D. Angluin. Learning regular sets from queries and counterexamples. *Information and Computation*, 75:87–106, 1987.
- [12] D. Angluin. Queries and Concept Learning. *Machine Learning*, 2(4):319–342, 1988.
- [13] M. J. Kearns and U. V. Vazirani. *An introduction to computational learning theory*. The MIT Press, 1994.
- [14] E. M. Gold. Language identification in the limit. *Information and Control*, 10:447–474, 1967.
- [15] L. Pitt. Inductive inference, dfas, and computational complexity. *LNAI*, 397:18–44, 1989. Proc. 2nd Workshop on Analogical and Inductive Inference.
- [16] D. Angluin and C.H. Smith. Inductive inference: Theory and Methods. *Computing Surveys*, 15(3):237–269, 1983.
- [17] C. de la Higuera. Characteristic sets for polynomial grammatical inference. *Machine Learning*, 27:125–138, 1997.
- [18] C. de la Higuera, J. Oncina, and E. Vidal. Identification of dfa: data-dependent vs data-independent algorithms. *LNAI*, 1147:313–325, 1996. 3rd International Colloquium, ICGI-96.
- [19] P. García, M. Vázquez de Parga, D. López, and J. Ruiz. Learning automata teams. *LNAI*, 6339:52–65, 2010. 10th International Colloquium, ICGI-10.