

LEARNING REGULAR LANGUAGES OVER LARGE ORDERED ALPHABETS

IRINI-ELEFThERIA MENS AND ODED MALER

VERIMAG, CNRS and University of Grenoble, France
e-mail address: {irini-eleftheria.mens,oded.maler}@imag.fr

ABSTRACT. This work is concerned with regular languages defined over large alphabets, either infinite or just too large to be expressed enumeratively. We define a generic model where transitions are labeled by elements of a finite partition of the alphabet. We then extend Angluin’s L^* algorithm for learning regular languages from examples for such automata. We have implemented this algorithm and we demonstrate its behavior where the alphabet is a subset of the natural or real numbers. We sketch the extension of the algorithm to a class of languages over partially ordered alphabets.

INTRODUCTION

The main contribution of this paper is a generic algorithm for learning regular languages defined over a large alphabet Σ . Such an alphabet can be infinite, like \mathbb{N} or \mathbb{R} or just so large, like \mathbb{B}^n for very large n or large subsets of \mathbb{N} , so that it is impossible or impractical to treat it in an enumerative way, that is, to write down the entries of the transition function $\delta(q, a)$ for every $a \in \Sigma$. The obvious solution is to use a *symbolic* representation where transitions are labeled by predicates which are applicable to the alphabet in question. Learning algorithms infer an automaton from a finite set of words (the *sample*) for which membership is known. Over small alphabets, the sample should include the set S of all the shortest words that lead to each state (access sequences) and, in addition, the set $S \cdot \Sigma$ of all their Σ -continuations. Over large alphabets this is not a practical option and as an alternative we develop a symbolic learning algorithm over *symbolic words* which are only partially backed up by the sample. In a sense, our algorithm is a combination of automaton learning and learning of non-temporal predicates. Before getting technical, let us discuss briefly some motivation.

Finite automata are among the corner stones of Computer Science. From a practical point of view they are used routinely in various domains ranging from syntactic analysis, design of user interfaces or administrative procedures to implementation of digital hardware and verification of software and hardware protocols. Regular languages admit a very

2012 ACM CCS: [Theory of computation]: Formal languages and automata theory; Theory and algorithms for application domains—Machine learning theory.

Key words and phrases: symbolic automata, active learning.

This paper is an extended version of [MM14].

nice, clean and comprehensive theory where different formalisms such as automata, logic, regular expressions, semigroups and grammars are shown to be equivalent. The problem of learning automata from examples was introduced already in 1956 by Moore [Moo56]. This problem, like the problem of automaton minimization, is closely related to the Nerode right-congruence relation over words associated with every language or sequential function [Ner58]. This relation declares two *input histories* as equivalent if they lead to the same *future continuations*, thus providing a crisp characterization of what a *state* in a dynamical system is in terms of observable input-output behavior. All algorithms for learning automata from examples, starting with the seminal work of Gold [Gol72] and culminating in the well-known L^* algorithm of Angluin [Ang87] are based on this concept [DIH10].

One weakness, however, of the classical theory of regular languages is that it is rather “thin” and “flat”. In other words, the alphabet is often considered as a small set devoid of any additional structure. On such alphabets, classical automata are good for expressing and exploring the temporal (sequential, monoidal) dimension embodied by the concatenation operations, but less good in expressing “horizontal” relationships. To make this statement more concrete, consider the verification of a system consisting of n automata running in parallel, making independent as well as synchronized transitions. To express the set of joint behaviors of this product of automata as a formal language, classical theory will force you to use the exponential alphabet of global states and indeed, a large part of verification is concerned with fighting this explosion using constructs such as BDDs and other logical forms that exploit the sparse interaction among components. This is done, however, without a real interaction with classical formal language theory (one exception is the theory of *traces* [DR95] which attempts to treat this issue but in a very restricted context).

These and other considerations led us to use *symbolic automata* as a generic framework for recognizing languages over large alphabets where transitions outgoing from a state are labeled, semantically speaking, by *subsets* of the alphabet. These subsets are expressed syntactically according to the specific alphabet used: Boolean formulae when $\Sigma = \mathbb{B}^n$ or by some classes of inequalities when $\Sigma \subseteq \mathbb{N}$ or $\Sigma \subseteq \mathbb{R}$. Determinism and completeness of the transition relation, which are crucial for learning and minimization, can be enforced by requiring that the subsets of Σ that label the transitions outgoing from a given state form a *partition* of the alphabet. Such symbolic automata have been used in the past for Boolean vectors [HJJ⁺95] and have been studied extensively in recent years as acceptors and transducers where transitions are guarded by predicates of various theories [HV11, VHL⁺12].

Readers working on program verification or hybrid automata are, of course, aware of automata with symbolic transition guards but it should be noted that in the model that we use, *no auxiliary variables* are added to the automaton. Let us stress this point by looking at a popular extension of automata to infinite alphabets, initiated in [KF94] using *register automata* to accept *data languages* (see [BLP10] for a good exposition of theoretical properties and [HSJC12] for learning algorithms). In that framework, the automaton is augmented with additional registers that can store some input letters. The registers can then be compared with newly-read letters and influence transitions. With register automata one can express, for example, the requirement that the password at login is the same as the password at sign-up. This very restricted use of memory makes register automata much simpler than more notorious automata with variables whose emptiness problem is typically undecidable. The downside is that beyond *equality* they do not really exploit the potential richness of the alphabets and their corresponding theories.

Our approach is different: we do allow the *values* of the input symbols to influence transitions via predicates, possibly of a restricted complexity. These predicates involve domain *constants* and they partition the alphabet into finitely many classes. For example, over the integers a state may have transitions labeled by conditions of the form $c_1 \leq x \leq c_2$ which give real (but of limited resolution) access to the input domain. On the other hand, we insist on a finite (and small) memory so that the exact value of x *cannot* be registered and has no future influence beyond the transition it has triggered. Many control systems, artificial (sequential machines working on quantized numerical inputs) as well as natural (central nervous system, the cell), are believed to operate in this manner. The automata that we use, like the symbolic automata and transducers studied in [HV11, VHL⁺12, VB12], are geared toward languages recognized by automata having a large alphabet and a relatively-small state space.

We then develop a symbolic version of Angluin’s L^* algorithm for learning regular sets from queries and counter-examples whose output is a symbolic automaton. The main difference relative to the concrete algorithm is that in the latter, every transition $\delta(q, a)$ in a conjectured automaton has at least one word in the sample that exercises it. In the symbolic case, a transition $\delta(q, \mathbf{a})$ where \mathbf{a} stands for a *set* of concrete symbols, will be backed up in the sample only by a *subset* of \mathbf{a} . Thus, unlike concrete algorithms where a counter-example always leads to a discovery of one or more new states, in our algorithm it may sometimes only modify the boundaries between partition blocks without creating new states. There are some similarities between our work and another recent adaptation of the L^* algorithm to symbolic automata, the Σ^* algorithm of [BB13]. This work is incomparable to ours as they use a richer model of transducers and more general predicates on inputs and outputs. Consequently their termination result is weaker and is relative to the termination of the counter-example guided abstraction refinement procedure.

The rest of the paper is organized as follows. In Section 1 we provide a quick summary of learning algorithms over small alphabets. In Section 2 we define symbolic automata and then extend the structure which underlies all automaton learning algorithms, namely the *observation table*, to be symbolic, where symbolic letters represent sets, and where entries in the table are supported only by partial evidence. In Section 4 we write down a symbolic learning algorithm, an adaptation of L^* for totally ordered alphabets such as \mathbb{R} or \mathbb{N} and illustrate the behavior of a prototype implementation. The algorithm is then extended to languages over partially ordered alphabets such as \mathbb{N}^d and \mathbb{R}^d where in each state, the labels of outgoing transition from a monotone partition of the alphabet are represented by finitely many points. We conclude by a discussion of past and future work.

1. LEARNING REGULAR SETS

We briefly survey Angluin’s L^* algorithm [Ang87] for learning regular sets from membership queries and counter-examples, with slightly modified definitions to accommodate for its symbolic extension. Let Σ be a finite alphabet and let Σ^* be the set of sequences (words) over Σ . Any order relation $<$ over Σ can be naturally lifted to a lexicographic order over Σ^* . With a language $L \subseteq \Sigma^*$ we associate a *characteristic function* $f : \Sigma^* \rightarrow \{+, -\}$, where $f(w) = +$ if the word $w \in \Sigma^*$ belongs to L and $f(w) = -$, otherwise.

A *deterministic finite automaton* over Σ is a tuple $\mathcal{A} = (\Sigma, Q, \delta, q_0, F)$, where Q is a non-empty finite set of *states*, $q_0 \in Q$ is the *initial* state, $\delta : Q \times \Sigma \rightarrow Q$ is the *transition function*, and $F \subseteq Q$ is the set of *final* or *accepting* states. The transition function δ can

be extended to $\delta : Q \times \Sigma^* \rightarrow Q$, where $\delta(q, \epsilon) = q$, and $\delta(q, u \cdot a) = \delta(\delta(q, u), a)$ for $q \in Q$, $a \in \Sigma$ and $u \in \Sigma^*$. A word $w \in \Sigma^*$ is *accepted* by \mathcal{A} if $\delta(q_0, w) \in F$, otherwise w is *rejected*. The language recognized by \mathcal{A} is the set of all accepted words and is denoted by $L(\mathcal{A})$.

Learning algorithms, represented by the *learner*, are designed to infer an unknown regular language L (the *target language*). The learner aims to construct a finite automaton that recognizes L by gathering information from the *teacher*. The *teacher* knows L and can provide information about it. It can answer two types of queries: *membership queries*, i.e., whether a given word belongs to the target language, and *equivalence queries*, i.e., whether a conjectured automaton suggested by the learner is the right one. If this automaton fails to accept L the teacher responds to the equivalence query by a *counter-example*, a word miss-classified by the conjectured automaton.

In the L^* algorithm, the learner starts by asking membership queries. All information provided is suitably gathered in a table structure, the *observation table*. Then, when the information is sufficient, the learner constructs a *hypothesis automaton* and poses an equivalence query to the teacher. If the answer is positive then the algorithm terminates and returns the conjectured automaton. Otherwise the learner accommodates the information provided by the counter-example into the table, asks additional membership queries until it can suggest a new hypothesis and so on, until termination.

A prefix-closed set $S \uplus R \subset \Sigma^*$ is a *balanced Σ -tree* if $\forall a \in \Sigma$: 1) For every $s \in S$ $s \cdot a \in S \cup R$, and 2) For every $r \in R$, $r \cdot a \notin S \cup R$. Elements of R are called *boundary elements* or *leaves*.¹

Definition 1.1 (Observation Table). An *observation table* is a tuple $T = (\Sigma, S, R, E, f)$ such that Σ is an alphabet, $S \cup R$ is a balanced Σ -tree, E is a subset of Σ^* and $f : (S \cup R) \cdot E \rightarrow \{-, +\}$ is the classification function, a restriction of the characteristic function of the target language L .

The set $(S \cup R) \cdot E$ is the *sample* associated with the table, that is, the set of words whose membership is known. The elements of S admit a tree structure isomorphic to a *spanning tree* of the transition graph rooted in the initial state. Each $s \in S$ corresponds to a state q of the automaton for which s is an *access sequence*, one of the shortest words that lead from the initial state to q . The elements of R should tell us about the back- and cross-edges in the automaton and the elements of E are “experiments” that should be sufficient to distinguish between states. This works by associating with every $s \in S \cup R$ a specialized classification function $f_s : E \rightarrow \{-, +\}$, defined as $f_s(e) = f(s \cdot e)$, which characterizes the row of the observation table labeled by s . To build an automaton from a table it should satisfy certain conditions.

Definition 1.2 (Closed, Reduced and Consistent Tables). An observation table T is:

- Closed if for every $r \in R$, there exists an $s \in S$, such that $f_r = f_s$;
- Reduced if for every $s, s' \in S$ $f_s \neq f_{s'}$;
- Consistent if for every $s, s' \in S$, $f_s = f_{s'}$ implies $f_{s \cdot a} = f_{s' \cdot a}, \forall a \in \Sigma$.

Note that a reduced table is trivially consistent and that for a closed and reduced table we can define a function $g : R \rightarrow S$ mapping every $r \in R$ to the unique $s \in S$ such that $f_s = f_r$. From such an observation table $T = (\Sigma, S, R, E, f)$ one can construct an automaton

¹We use \uplus for disjoint union.

$\mathcal{A}_T = (\Sigma, Q, q_0, \delta, F)$ where $Q = S$, $q_0 = \epsilon$, $F = \{s \in S : f_s(\epsilon) = +\}$ and

$$\delta(s, a) = \begin{cases} s \cdot a & \text{when } s \cdot a \in S \\ g(s \cdot a) & \text{when } s \cdot a \in R \end{cases}$$

The learner attempts to keep the table closed at all times. The table is not closed when there is some $r \in R$ such that f_r is different from f_s for all $s \in S$. To close the table, the learner moves r from R to S and adds the Σ -successors of r , i.e., all words $r \cdot a$ for $a \in \Sigma$, to R . The extended table is then filled up by asking membership queries until it becomes closed.

Variants of the L^* algorithm differ in the way they treat counter-examples, as described in more detail in [BR04]. The original algorithm [Ang87] adds all the *prefixes* of the counter-example to S and thus possibly creating inconsistency that should be fixed. The version proposed in [MP95] for learning ω -regular languages adds all the *suffixes* of the counter-example to E . The advantage of this approach is that the table always remains consistent and reduced with S corresponding exactly to the set of states. A disadvantage is the possible introduction of redundant columns that do not contribute to further discrimination between states. The symbolic algorithm that we develop in this paper is based on an intermediate variant, referred to in [BR04] as the *reduced observation algorithm*, where some prefixes of the counter-example are added to S and some suffixes are added to E .

Example 1.3. We illustrate the behavior of the L^* algorithm while learning a language L over $\Sigma = \{1, 2, 3, 4, 5\}$. We use the tuple $(w, +)$ to indicate a counter-example $w \in L$ rejected by the conjectured automaton, and $(w, -)$ for the opposite case. Initially, the observation table is $T_0 = (\Sigma, S, R, E, f)$ with $S = E = \{\epsilon\}$ and $R = \Sigma$ and we ask membership queries for all words in $(S \cup R) \cdot E$ to obtain table T_0 , shown in Fig. 1. The table is not closed so we move word 1 to S , add its continuations, $1 \cdot \Sigma$ to R and ask membership queries to obtain table T_1 which is now closed. We construct an hypothesis \mathcal{A}_1 (Fig. 2) from this table, and pose an equivalence query for which the teacher returns counter-example $(3 \cdot 1, -)$. We add $3 \cdot 1$ and its prefix 3 to set S and add all their continuations to the boundary of the table resulting table T_2 of Fig. 1. This table is not consistent: two elements ϵ and 3 in S are equivalent but their successors 1 and $3 \cdot 1$ are not. In order to distinguish the two strings we add to E the suffix 1 and end up with a closed and consistent table T_3 . The new hypothesis for this table is \mathcal{A}_3 , shown in Fig. 2. Once more the equivalence query will return a counter-example, $(1 \cdot 3 \cdot 3, -)$. We again add the counter-example and prefixes to the table, ask membership queries to fill in the table and solve the inconsistency that appears for 1 and $1 \cdot 3$ by adding suffix 3 to the table. The table corresponds now to the correct hypothesis \mathcal{A}_5 , and the algorithm terminates. \square

2. SYMBOLIC AUTOMATA

In this section we introduce the variant of *symbolic automata* that we use. Symbolic automata [HV11, VB12] give a more succinct representation for languages over large finite alphabets and can also represent languages over infinite alphabets such as \mathbb{N} , \mathbb{R} , or \mathbb{R}^n . The size of a standard automaton for a language grows linearly with the size of the alphabet and so does the complexity of learning algorithms such as L^* . As we shall see, symbolic automata admit a variant of the L^* algorithm whose complexity is independent of the alphabet size.

T_0	T_1	T_2	T_3	T_4	T_5																																																																																																																																																																																																		
<table border="1" style="display: inline-table; border-collapse: collapse;"><tr><td></td><td>ϵ</td></tr><tr><td>ϵ</td><td>-</td></tr><tr><td>1</td><td>+</td></tr><tr><td>2</td><td>+</td></tr><tr><td>3</td><td>-</td></tr><tr><td>4</td><td>-</td></tr><tr><td>5</td><td>-</td></tr></table>		ϵ	ϵ	-	1	+	2	+	3	-	4	-	5	-	<table border="1" style="display: inline-table; border-collapse: collapse;"><tr><td></td><td>ϵ</td></tr><tr><td>ϵ</td><td>-</td></tr><tr><td>1</td><td>+</td></tr><tr><td>2</td><td>+</td></tr><tr><td>3</td><td>-</td></tr><tr><td>4</td><td>-</td></tr><tr><td>5</td><td>-</td></tr><tr><td>1 · 1</td><td>-</td></tr><tr><td>1 · 2</td><td>-</td></tr><tr><td>1 · 3</td><td>+</td></tr><tr><td>1 · 4</td><td>-</td></tr><tr><td>1 · 5</td><td>-</td></tr></table>		ϵ	ϵ	-	1	+	2	+	3	-	4	-	5	-	1 · 1	-	1 · 2	-	1 · 3	+	1 · 4	-	1 · 5	-	<table border="1" style="display: inline-table; border-collapse: collapse;"><tr><td></td><td>ϵ</td></tr><tr><td>ϵ</td><td>-</td></tr><tr><td>1</td><td>+</td></tr><tr><td>3</td><td>-</td></tr><tr><td>3 · 1</td><td>-</td></tr><tr><td>2</td><td>+</td></tr><tr><td>4</td><td>-</td></tr><tr><td>5</td><td>-</td></tr><tr><td>1 · 1</td><td>-</td></tr><tr><td>1 · 2</td><td>-</td></tr><tr><td>1 · 3</td><td>+</td></tr><tr><td>1 · 4</td><td>-</td></tr><tr><td>⋮</td><td></td></tr></table>		ϵ	ϵ	-	1	+	3	-	3 · 1	-	2	+	4	-	5	-	1 · 1	-	1 · 2	-	1 · 3	+	1 · 4	-	⋮		<table border="1" style="display: inline-table; border-collapse: collapse;"><tr><td></td><td>ϵ</td><td>1</td></tr><tr><td>ϵ</td><td>-</td><td>+</td></tr><tr><td>1</td><td>+</td><td>-</td></tr><tr><td>3</td><td>-</td><td>-</td></tr><tr><td>3 · 1</td><td>-</td><td>-</td></tr><tr><td>2</td><td>+</td><td>-</td></tr><tr><td>4</td><td>-</td><td>-</td></tr><tr><td>5</td><td>-</td><td>-</td></tr><tr><td>1 · 1</td><td>-</td><td>-</td></tr><tr><td>1 · 2</td><td>-</td><td>-</td></tr><tr><td>1 · 3</td><td>+</td><td>-</td></tr><tr><td>1 · 4</td><td>-</td><td>-</td></tr><tr><td>⋮</td><td></td><td></td></tr></table>		ϵ	1	ϵ	-	+	1	+	-	3	-	-	3 · 1	-	-	2	+	-	4	-	-	5	-	-	1 · 1	-	-	1 · 2	-	-	1 · 3	+	-	1 · 4	-	-	⋮			<table border="1" style="display: inline-table; border-collapse: collapse;"><tr><td></td><td>ϵ</td><td>1</td></tr><tr><td>ϵ</td><td>-</td><td>+</td></tr><tr><td>1</td><td>+</td><td>-</td></tr><tr><td>3</td><td>-</td><td>-</td></tr><tr><td>1 · 3</td><td>+</td><td>-</td></tr><tr><td>3 · 1</td><td>-</td><td>+</td></tr><tr><td>1 · 3 · 3</td><td>-</td><td>-</td></tr><tr><td>2</td><td>+</td><td>-</td></tr><tr><td>4</td><td>-</td><td>-</td></tr><tr><td>5</td><td>-</td><td>-</td></tr><tr><td>1 · 1</td><td>-</td><td>-</td></tr><tr><td>1 · 2</td><td>-</td><td>-</td></tr><tr><td>⋮</td><td></td><td></td></tr></table>		ϵ	1	ϵ	-	+	1	+	-	3	-	-	1 · 3	+	-	3 · 1	-	+	1 · 3 · 3	-	-	2	+	-	4	-	-	5	-	-	1 · 1	-	-	1 · 2	-	-	⋮			<table border="1" style="display: inline-table; border-collapse: collapse;"><tr><td></td><td>ϵ</td><td>1</td><td>3</td></tr><tr><td>ϵ</td><td>-</td><td>+</td><td>-</td></tr><tr><td>1</td><td>+</td><td>-</td><td>+</td></tr><tr><td>3</td><td>-</td><td>-</td><td>-</td></tr><tr><td>1 · 3</td><td>+</td><td>-</td><td>-</td></tr><tr><td>3 · 1</td><td>-</td><td>+</td><td>-</td></tr><tr><td>1 · 3 · 3</td><td>-</td><td>-</td><td>-</td></tr><tr><td>2</td><td>+</td><td>-</td><td>+</td></tr><tr><td>4</td><td>-</td><td>-</td><td>-</td></tr><tr><td>5</td><td>-</td><td>-</td><td>-</td></tr><tr><td>1 · 1</td><td>-</td><td>-</td><td>-</td></tr><tr><td>1 · 2</td><td>-</td><td>-</td><td>-</td></tr><tr><td>⋮</td><td></td><td></td><td></td></tr></table>		ϵ	1	3	ϵ	-	+	-	1	+	-	+	3	-	-	-	1 · 3	+	-	-	3 · 1	-	+	-	1 · 3 · 3	-	-	-	2	+	-	+	4	-	-	-	5	-	-	-	1 · 1	-	-	-	1 · 2	-	-	-	⋮			
	ϵ																																																																																																																																																																																																						
ϵ	-																																																																																																																																																																																																						
1	+																																																																																																																																																																																																						
2	+																																																																																																																																																																																																						
3	-																																																																																																																																																																																																						
4	-																																																																																																																																																																																																						
5	-																																																																																																																																																																																																						
	ϵ																																																																																																																																																																																																						
ϵ	-																																																																																																																																																																																																						
1	+																																																																																																																																																																																																						
2	+																																																																																																																																																																																																						
3	-																																																																																																																																																																																																						
4	-																																																																																																																																																																																																						
5	-																																																																																																																																																																																																						
1 · 1	-																																																																																																																																																																																																						
1 · 2	-																																																																																																																																																																																																						
1 · 3	+																																																																																																																																																																																																						
1 · 4	-																																																																																																																																																																																																						
1 · 5	-																																																																																																																																																																																																						
	ϵ																																																																																																																																																																																																						
ϵ	-																																																																																																																																																																																																						
1	+																																																																																																																																																																																																						
3	-																																																																																																																																																																																																						
3 · 1	-																																																																																																																																																																																																						
2	+																																																																																																																																																																																																						
4	-																																																																																																																																																																																																						
5	-																																																																																																																																																																																																						
1 · 1	-																																																																																																																																																																																																						
1 · 2	-																																																																																																																																																																																																						
1 · 3	+																																																																																																																																																																																																						
1 · 4	-																																																																																																																																																																																																						
⋮																																																																																																																																																																																																							
	ϵ	1																																																																																																																																																																																																					
ϵ	-	+																																																																																																																																																																																																					
1	+	-																																																																																																																																																																																																					
3	-	-																																																																																																																																																																																																					
3 · 1	-	-																																																																																																																																																																																																					
2	+	-																																																																																																																																																																																																					
4	-	-																																																																																																																																																																																																					
5	-	-																																																																																																																																																																																																					
1 · 1	-	-																																																																																																																																																																																																					
1 · 2	-	-																																																																																																																																																																																																					
1 · 3	+	-																																																																																																																																																																																																					
1 · 4	-	-																																																																																																																																																																																																					
⋮																																																																																																																																																																																																							
	ϵ	1																																																																																																																																																																																																					
ϵ	-	+																																																																																																																																																																																																					
1	+	-																																																																																																																																																																																																					
3	-	-																																																																																																																																																																																																					
1 · 3	+	-																																																																																																																																																																																																					
3 · 1	-	+																																																																																																																																																																																																					
1 · 3 · 3	-	-																																																																																																																																																																																																					
2	+	-																																																																																																																																																																																																					
4	-	-																																																																																																																																																																																																					
5	-	-																																																																																																																																																																																																					
1 · 1	-	-																																																																																																																																																																																																					
1 · 2	-	-																																																																																																																																																																																																					
⋮																																																																																																																																																																																																							
	ϵ	1	3																																																																																																																																																																																																				
ϵ	-	+	-																																																																																																																																																																																																				
1	+	-	+																																																																																																																																																																																																				
3	-	-	-																																																																																																																																																																																																				
1 · 3	+	-	-																																																																																																																																																																																																				
3 · 1	-	+	-																																																																																																																																																																																																				
1 · 3 · 3	-	-	-																																																																																																																																																																																																				
2	+	-	+																																																																																																																																																																																																				
4	-	-	-																																																																																																																																																																																																				
5	-	-	-																																																																																																																																																																																																				
1 · 1	-	-	-																																																																																																																																																																																																				
1 · 2	-	-	-																																																																																																																																																																																																				
⋮																																																																																																																																																																																																							

FIGURE 1. Observation tables for Example 1.3.

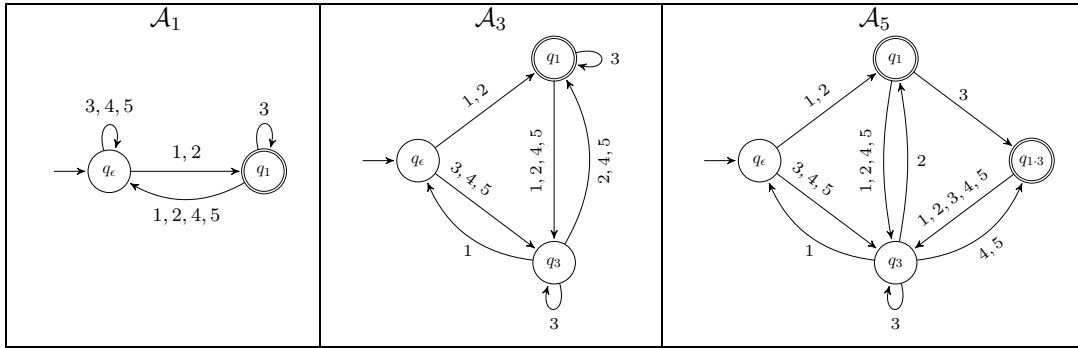


FIGURE 2. Hypotheses for Example 1.3

Let Σ be a large, possibly infinite, alphabet, to which we will refer from now on as the *concrete* alphabet. We define a symbolic automaton to be an automaton over Σ where each state has a small number of outgoing transitions labeled by symbols that represent subsets of Σ . For every state, these subsets form a (possibly different) *partition* of Σ and hence the automaton is complete and deterministic. We start with an arbitrary alphabet viewed as an unstructured set and present the concept in purely semantic manner before we move to ordered sets and inequalities in subsequent sections.

Let $\mathbf{\Sigma}$ be a finite alphabet, that we call the *symbolic alphabet* and its elements *symbolic letters* or *symbols*. Let $\psi : \Sigma \rightarrow \mathbf{\Sigma}$ map concrete letters into symbolic ones. The Σ -semantics of a *symbolic letter* $\mathbf{a} \in \mathbf{\Sigma}$ is defined as $[\mathbf{a}]_\psi = \{a \in \Sigma : \psi(a) = \mathbf{a}\}$ and the set $\{[\mathbf{a}]_\psi : \mathbf{a} \in \mathbf{\Sigma}\}$ forms a *partition* of Σ . We will often omit ψ from the notation and use $[\mathbf{a}]$ when ψ , which is always present, is clear from the context. The Σ -semantics can be extended to symbolic words of the form $\mathbf{w} = \mathbf{a}_1 \cdot \mathbf{a}_2 \cdots \mathbf{a}_k \in \mathbf{\Sigma}^*$ as the concatenation of the concrete one-letter languages associated with the respective symbolic letters or, recursively speaking, $[\epsilon] = \{\epsilon\}$ and $[\mathbf{w} \cdot \mathbf{a}] = [\mathbf{w}] \cdot [\mathbf{a}]$ for $\mathbf{w} \in \mathbf{\Sigma}^*$, $\mathbf{a} \in \mathbf{\Sigma}$.

Definition 2.1 (Symbolic Automaton). A *deterministic symbolic automaton* is a tuple $\mathcal{A} = (\Sigma, \mathbf{\Sigma}, \psi, Q, \delta, \delta, q_0, F)$, where

- Σ is the input alphabet,

- Σ is a finite alphabet, decomposable into $\Sigma = \bigsqcup_{q \in Q} \Sigma_q$,
- $\psi = \{\psi_q : q \in Q\}$ is a family of surjective functions $\psi_q : \Sigma \rightarrow \Sigma_q$,
- Q is a finite set of states,
- $\delta : Q \times \Sigma \rightarrow Q$ and $\delta : Q \times \Sigma \rightarrow Q$ are the concrete and symbolic transition functions respectively, such that $\delta(q, a) = \delta(q, \psi_q(a))$,
- q_0 is the initial state and F is a set of accepting states.

The transition function is extended to words as in the concrete case and the symbolic automaton can be viewed as an acceptor of a concrete language. When at q and reading a concrete letter a , the automaton will take the transition $\delta(q, \mathbf{a})$ where \mathbf{a} is the *unique* element of Σ_q satisfying $a \in [\mathbf{a}]$. Hence $L(\mathcal{A})$ consists of all concrete words whose run leads from q_0 to a state in F . A language L over alphabet Σ is symbolic recognizable if there exists a symbolic automaton \mathcal{A} such that $L = L(\mathcal{A})$.

Remark: The association of a *symbolic language* with a symbolic automaton is more subtle because we allow different partitions of Σ and hence different symbolic input alphabets at different states. The transition to be taken while being in a state q and reading a symbol $\mathbf{a} \notin \Sigma_q$ is well defined only when $[\mathbf{a}] \subseteq [\mathbf{a}']$ for some $\mathbf{a}' \in \Sigma_q$. Such a model can be transformed into an automaton which is complete over a symbolic alphabet which is common to all states as follows. Let

$$\Sigma' = \prod_{q \in Q} \Sigma_q, \text{ with the } \Sigma\text{-semantics } [(\mathbf{a}_1, \dots, \mathbf{a}_n)] = [\mathbf{a}_1] \cap \dots \cap [\mathbf{a}_n],$$

and let $\tilde{\Sigma} = \{\mathbf{b} \in \Sigma' : [\mathbf{b}] \neq \emptyset\}$. Then we define $\tilde{\mathcal{A}} = (\tilde{\Sigma}, Q, \tilde{\delta}, q_0, F)$ where, by construction, for every $\mathbf{b} \in \tilde{\Sigma}$ and every $q \in Q$, there is a unique $\mathbf{a} \in \Sigma_q$ such that $[\mathbf{b}] \subseteq [\mathbf{a}]$ and hence one can define the transition function as $\tilde{\delta}(q, \mathbf{b}) = \delta(q, \mathbf{a})$. This model is more comfortable for language-theoretic studies but in the learning context it introduces an unnecessary blow-up in the alphabet size and the number of queries for every state. For this reason we stick in this paper to the Definition 2.1 which is more economical. A similar approach of state-local abstraction has been taken in [IHS13] for learning parameterized language. The construction of Σ' is similar to the minterm construction of [DV14] used to create a common alphabet in order to apply the minimization algorithm of Hopcroft to symbolic automata. Anyway, in our learning framework symbolic automata are used to read concrete and not symbolic words. \square

It is straightforward that for a finite concrete alphabet Σ the set of languages accepted by symbolic automata coincides with the set of recognizable regular languages over Σ . Moreover, even when the alphabet is infinite, closure under Boolean operations is preserved.

Proposition 2.2 (Closure under Boolean Operations). *Languages accepted by deterministic symbolic automata are effectively closed under Boolean operations.*

Proof. Closure under complement is immediate by complementing the set of accepting states. For intersection the standard product construction is adapted as follows. Let L_1, L_2 be languages recognized by the symbolic automata $\mathcal{A}_1 = (\Sigma, \Sigma_1, \psi_1, Q_1, \delta_1, \delta_1, q_{01}, F_1)$, and $\mathcal{A}_2 = (\Sigma, \Sigma_2, \psi_2, Q_2, \delta_2, \delta_2, q_{02}, F_2)$, respectively. Let $\mathcal{A} = (\Sigma, \Sigma, \psi, Q, \delta, \delta, q_0, F)$, where

- $Q = Q_1 \times Q_2$, $q_0 = (q_{01}, q_{02})$, $F = F_1 \times F_2$,
- For every $(q_1, q_2) \in Q$
 - $\Sigma_{(q_1, q_2)} = \{(\mathbf{a}_1, \mathbf{a}_2) \in \Sigma_1 \times \Sigma_2 \mid [\mathbf{a}_1] \cap [\mathbf{a}_2] \neq \emptyset\}$
 - $\psi_{(q_1, q_2)}(a) = (\psi_{1, q_1}(a), \psi_{2, q_2}(a))$, $\forall a \in \Sigma$

$$- \delta((q_1, q_2), (\mathbf{a}_1, \mathbf{a}_2)) = (\delta_1(q_1, \mathbf{a}_1), \delta_2(q_2, \mathbf{a}_2)), \forall (\mathbf{a}_1, \mathbf{a}_2) \in \Sigma_{(q_1, q_2)}$$

It is sufficient to observe that the corresponding implied concrete automata \mathcal{A}_1 , \mathcal{A}_2 and \mathcal{A} satisfy $\delta((q_1, q_2), a) = (\delta_1(q_1, a), \delta_2(q_2, a))$ and the standard proof that $L(\mathcal{A}) = L(\mathcal{A}_1) \cap L(\mathcal{A}_2)$ follows. Closure under union and set difference is then evident. \square

The above product construction is used to implement equivalence queries where both the target language and the current conjecture are represented by symbolic automata. A counter-example is found by looking for a shortest path in the product automaton from the initial state to a state in $F_1 \times (Q_2 - F_2) \cup (Q_1 - F_1) \times F_2$ and selecting a lexicographically minimal concrete word along that path.

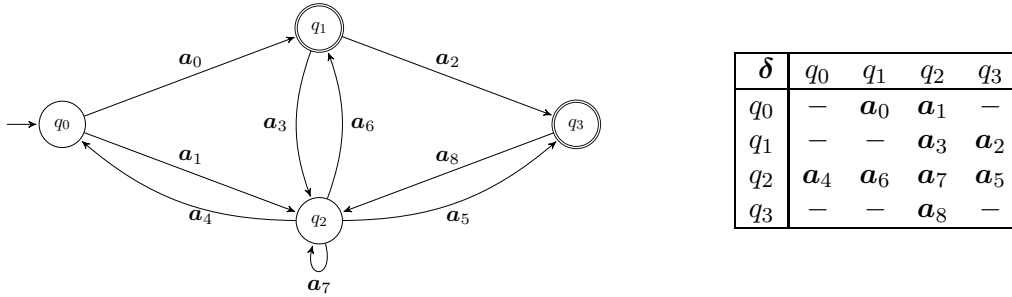


FIGURE 3. A symbolic automaton \mathcal{A} with its symbolic transition function.

Example 2.3. Figure 3 shows a symbolic automaton equivalent to automaton \mathcal{A}_5 of Figure 2. The symbolic alphabets for the states are $\Sigma_{q_0} = \{\mathbf{a}_0, \mathbf{a}_1\}$, $\Sigma_{q_1} = \{\mathbf{a}_2, \mathbf{a}_3\}$, $\Sigma_{q_2} = \{\mathbf{a}_4, \mathbf{a}_5, \mathbf{a}_6, \mathbf{a}_7\}$, $\Sigma_{q_3} = \{\mathbf{a}_8\}$, and the Σ -semantics for the symbols is $[\mathbf{a}_0] = \{1, 2\}$, $[\mathbf{a}_1] = \{3, 4, 5\}$, $[\mathbf{a}_2] = \{3\}$, $[\mathbf{a}_3] = \{1, 2, 4, 5\}$, etc.. The same automaton can accept a language over the uncountable alphabet $\Sigma = [0, 100) \subset \mathbb{R}$, defining ψ as shown in Figure 4.

ψ	0	20	30	50	80	100
Σ_{q_0}		\mathbf{a}_0			\mathbf{a}_1	
Σ_{q_1}	\mathbf{a}_2			\mathbf{a}_3		
Σ_{q_2}	\mathbf{a}_4	\mathbf{a}_5	\mathbf{a}_6	\mathbf{a}_7		
Σ_{q_3}		\mathbf{a}_8				

FIGURE 4. The concrete semantics of the symbols of automaton \mathcal{A} of Fig. 3, when defined over $\Sigma = [0, 100) \subseteq \mathbb{R}$.

3. SYMBOLIC OBSERVATION TABLES

In this section we adapt observation tables to the symbolic setting. They are similar to the concrete case with the additional notions of evidences and evidence compatibility.

Definition 3.1 (Balanced Symbolic Σ -Tree). A *balanced symbolic Σ -tree* is a tuple $(\Sigma, \mathbf{S}, \mathbf{R}, \psi)$ where

- $\mathbf{S} \uplus \mathbf{R}$ is a prefix-closed subset of Σ^*
- $\Sigma = \bigsqcup_{s \in \mathbf{S}} \Sigma_s$ is a symbolic alphabet
- $\psi = \{\psi_s\}_{s \in \mathbf{S}}$ is a family of total surjective functions of the form $\psi_s : \Sigma \rightarrow \Sigma_s$.

It is required that for every $s \in \mathbf{S}$ and $\mathbf{a} \in \Sigma_s$, $s \cdot \mathbf{a} \in \mathbf{S} \cup \mathbf{R}$ and for any $r \in \mathbf{R}$ and $\mathbf{a} \in \Sigma$, $r \cdot \mathbf{a} \notin \mathbf{S} \cup \mathbf{R}$. Elements of \mathbf{R} are called boundary elements of the tree.

We will use observation tables whose rows are symbolic words and hence an entry in the table will constitute a statement about the inclusion or exclusion of a large *set* of concrete words in the language. We will not ask membership queries concerning all those concrete words, but only for a small representative subset that we call *evidence*.

Definition 3.2 (Symbolic Observation Table). A *symbolic observation table* is a tuple $T = (\Sigma, \Sigma, \mathbf{S}, \mathbf{R}, \psi, E, \mathbf{f}, \mu)$ such that

- Σ is an alphabet,
- $(\Sigma, \mathbf{S}, \mathbf{R}, \psi)$ is a balanced symbolic Σ -tree (with \mathbf{R} being its *boundary*),
- E is a subset of Σ^* ,
- $\mathbf{f} : (\mathbf{S} \cup \mathbf{R}) \cdot E \rightarrow \{-, +\}$ is the symbolic classification function
- $\mu : (\mathbf{S} \cup \mathbf{R}) \cdot E \rightarrow 2^{\Sigma^*} - \{\emptyset\}$ is an evidence function satisfying $\mu(\mathbf{w}) \subseteq [\mathbf{w}]$. The image of the evidence function is prefix-closed: $w \cdot a \in \mu(\mathbf{w} \cdot \mathbf{a}) \Rightarrow w \in \mu(\mathbf{w})$.

As for the concrete case we use $\mathbf{f}_s : E \rightarrow \{-, +\}$ to denote the partial evaluation of \mathbf{f} to some symbolic word $s \in \mathbf{S} \cup \mathbf{R}$, such that, $\mathbf{f}_s(e) = \mathbf{f}(s \cdot e)$. Note that the set E consists of *concrete* words but this poses no problem because elements of E are used only to distinguish between states and do not participate in the derivation of the symbolic automaton from the table. Concatenation of a symbolic word and a concrete one follows concatenation of symbolic words as defined above where each concrete letter a is considered as a symbolic letter \mathbf{a} with $[\mathbf{a}] = \{a\}$ and $\mu(\mathbf{a}) = a$. The notions of closed, consistent and reduced table are similar to the concrete case.

The set $\mathbf{M}_T = (\mathbf{S} \cup \mathbf{R}) \cdot E$ is called the *symbolic sample* associated with T . We require that for each word $\mathbf{w} \in \mathbf{M}_T$ there is at least one concrete $w \in \mu(\mathbf{w})$ whose membership in L , denoted by $f(w)$, is known. The set of such words is called the *concrete sample* and is defined as $M_T = \{s \cdot e : s \in \mu(\mathbf{s}), \mathbf{s} \in \mathbf{S} \cup \mathbf{R}, e \in E\}$. A table where all evidences of the same symbolic word admit the same classification is called *evidence-compatible*.

Definition 3.3 (Table Conditions). A table $T = (\Sigma, \Sigma, \mathbf{S}, \mathbf{R}, \psi, E, \mathbf{f}, \mu)$ is

- Closed if $\forall r \in \mathbf{R}, \exists s = g(r) \in \mathbf{S}, \mathbf{f}_r = \mathbf{f}_s$,
- Reduced if $\forall s, s' \in \mathbf{S}, \mathbf{f}_s \neq \mathbf{f}_{s'}$,
- Consistent if $\forall s, s' \in \mathbf{S}, \mathbf{f}_s = \mathbf{f}_{s'}$ implies $\mathbf{f}_{s \cdot \mathbf{a}} = \mathbf{f}_{s' \cdot \mathbf{a}}, \forall \mathbf{a} \in \Sigma_s$.
- Evidence compatible if $\forall \mathbf{w} \in \mathbf{M}_T, \forall w_1, w_2 \in \mu(\mathbf{w}), f(w_1) = f(w_2)$.

When a table T is evidence compatible the symbolic classification function \mathbf{f} can be defined for every $s \in (\mathbf{S} \cup \mathbf{R})$ and $e \in E$ as $\mathbf{f}(s \cdot e) = f(s \cdot e)$, $s \in \mu(\mathbf{s})$.

Theorem 3.4 (Automaton from Table). *From a closed, reduced and evidence compatible table one can construct a deterministic symbolic automaton compatible with the concrete sample.*

Proof. The proof is similar to the concrete case. Let $\mathbf{T} = (\Sigma, \mathbf{\Sigma}, \mathbf{S}, \mathbf{R}, \psi, E, \mathbf{f}, \mu)$ be such a table, which is reduced and closed and thus a function $g : \mathbf{R} \rightarrow \mathbf{S}$ such that $g(\mathbf{r}) = \mathbf{s}$ iff $\mathbf{f}_{\mathbf{r}} = \mathbf{f}_{\mathbf{s}}$ is well defined. The automaton derived from the table is then $\mathcal{A}_{\mathbf{T}} = (\Sigma, \mathbf{\Sigma}, \psi, Q, \delta, q_0, F)$ where:

- $Q = \mathbf{S}$, $q_0 = \epsilon$
- $F = \{\mathbf{s} \in \mathbf{S} \mid \mathbf{f}_{\mathbf{s}}(\epsilon) = +\}$
- $\delta : Q \times \Sigma \rightarrow Q$ is defined as $\delta(\mathbf{s}, \mathbf{a}) = \begin{cases} \mathbf{s} \cdot \mathbf{a} & \text{when } \mathbf{s} \cdot \mathbf{a} \in \mathbf{S} \\ g(\mathbf{s} \cdot \mathbf{a}) & \text{when } \mathbf{s} \cdot \mathbf{a} \in \mathbf{R} \end{cases}$

By construction and like the L^* algorithm, $\mathcal{A}_{\mathbf{T}}$ classifies correctly the symbolic sample and, due to evidence compatibility, this holds also for the concrete sample. \square

4. LEARNING LANGUAGES OVER ORDERED ALPHABETS

In this section we present a symbolic learning algorithm starting with an intuitive verbal description. The algorithmic scheme is similar to the concrete L^* algorithm but differs in the treatment of counter-examples and the new concept of evidence compatibility. Whenever the table is not closed, $\mathbf{S} \cup \mathbf{R}$ is extended until closure. Then a conjectured automaton $\mathcal{A}_{\mathbf{T}}$ is constructed and an equivalence query is posed. If the answer is positive we are done. Otherwise, the teacher provides a counter-example leading to the extension of $\mathbf{S} \cup \mathbf{R}$ and/or E . Whenever such an extension occurs, additional membership queries are posed to fill the table. The table is always kept evidence compatible and reduced except temporarily during the processing of counter-examples.

From now on we assume Σ to be a *totally ordered* alphabet with a minimal element a_0 and restrict ourselves to symbolic automata where the concrete semantics for every symbolic letter is an interval. In the case of a dense order like in \mathbb{R} , we assume the intervals to be left-closed and right-open. The order on the alphabet can be extended naturally to a lexicographic order on Σ^* . Our algorithm also assumes that the teacher provides a counter-example of minimal length which is minimal with respect to the lexicographic order. This strong assumption improves the performance of the algorithm and its relaxation is discussed in Section 7.

The rows of the observation table consist of symbolic words because we want to group together all concrete letters and words that are assumed to induce the same behavior in the automaton. New symbolic letters are introduced in two occasions: when a new state is discovered or when a partition is modified due to a counter-example. In both cases we set the concrete semantics $[\mathbf{a}]$ to the largest possible subset of Σ , given the current evidence (in the first case it will be Σ). As an evidence we always select the smallest possible $a \in [\mathbf{a}]$ (a_0 when $[\mathbf{a}] = \Sigma$). The choice of the right evidences is a key point for the performance of the algorithm as we want to keep the concrete sample as small as possible and avoid posing unnecessary queries. For infinite concrete alphabets this choice of evidence guarantees termination.

The initial symbolic table is $\mathbf{T} = (\Sigma, \mathbf{\Sigma}, \mathbf{S}, \mathbf{R}, \psi, E, \mathbf{f}, \mu)$, where $\mathbf{\Sigma} = \{\mathbf{a}_0\}$, $[\mathbf{a}_0] = \Sigma$, $\mathbf{S} = \{\epsilon\}$, $\mathbf{R} = \{\mathbf{a}_0\}$, $E = \{\epsilon\}$, and $\mu(\mathbf{a}_0) = \{a_0\}$. The table is filled by membership queries concerning ϵ and a_0 . Whenever \mathbf{T} is not closed, there is some $\mathbf{r} \in \mathbf{R}$ such that $\mathbf{f}_{\mathbf{r}} \neq \mathbf{f}_{\mathbf{s}}$ for every $\mathbf{s} \in \mathbf{S}$. To close the table we move \mathbf{r} from \mathbf{R} to \mathbf{S} , recognizing it as a new state, and checking the behavior of its continuation. To this end we add to \mathbf{R} the word $\mathbf{r}' = \mathbf{r} \cdot \mathbf{a}$,

Algorithm 1 The symbolic algorithm

```

1: procedure SYMBOLIC
2:   learned = false
3:   Initialize the table  $\mathbf{T} = (\Sigma, \mathbf{\Sigma}, \mathbf{S}, \mathbf{R}, \psi, E, \mathbf{f}, \mu)$ 
4:    $\mathbf{\Sigma} = \{\mathbf{a}\}; \psi_\epsilon(a) = \mathbf{a}, \forall a \in \Sigma$ 
5:    $\mathbf{S} = \{\epsilon\}; \mathbf{R} = \{\mathbf{a}\}; E = \{\epsilon\}$ 
6:    $\mu(\mathbf{a}) = \{a_0\}$ 
7:   Ask MQ on  $\epsilon$  and  $a_0$  to fill  $\mathbf{f}$ 
8:   if  $\mathbf{T}$  is not closed then
9:     CLOSE
10:  end if
11:  repeat
12:    if  $EQ(\mathcal{A}_{\mathbf{T}})$  then ▷  $\mathcal{A}_{\mathbf{T}}$  is correct
13:      learned = true
14:    else ▷ A counter-example  $w$  is provided
15:       $M = M \cup \{w\}$ 
16:      COUNTER-EX( $w$ ) ▷ Process counter-example
17:    end if
18:  until learned
19: end procedure

```

Procedure 2 Close the table

```

1: procedure CLOSE
2:   while  $\exists \mathbf{r} \in \mathbf{R}$  such that  $\forall \mathbf{s} \in \mathbf{S}, \mathbf{f}_{\mathbf{r}} \neq \mathbf{f}_{\mathbf{s}}$  do
3:      $\mathbf{S}' = \mathbf{S} \cup \{\mathbf{r}\}$  ▷  $\mathbf{r}$  becomes a new state
4:      $\mathbf{\Sigma}' = \mathbf{\Sigma} \cup \{\mathbf{a}_{\text{new}}\}$ 
5:      $\psi' = \psi \cup \{\psi_{\mathbf{r}}\}$  with  $\psi_{\mathbf{r}}(a) = \mathbf{a}_{\text{new}}, \forall a \in \Sigma$ 
6:      $\mathbf{R}' = (\mathbf{R} - \{\mathbf{r}\}) \cup \{\mathbf{r} \cdot \mathbf{a}_{\text{new}}\}$ 
7:      $\mu(\mathbf{r} \cdot \mathbf{a}_{\text{new}}) = \mu(\mathbf{r}) \cdot a_0$ 
8:     Ask MQ for all words in  $\{\mu(\mathbf{r} \cdot \mathbf{a}_{\text{new}}) \cdot e : e \in E\}$ 
9:      $\mathbf{T} = (\Sigma, \mathbf{\Sigma}', \mathbf{S}', \mathbf{R}', \psi', E, \mathbf{f}', \mu')$ 
10:  end while
11: end procedure

```

where \mathbf{a} is a new symbolic letter with $[\mathbf{a}] = \Sigma$. We extend the evidence function by letting $\mu(\mathbf{r}') = \mu(\mathbf{r}) \cdot a_0$, assuming that all elements of Σ behave as a_0 from \mathbf{r} . Once \mathbf{T} is closed we construct a hypothesis automaton as described in the proof of Theorem 3.4.

When a counter-example w is presented, it is of course not part of the concrete sample. A miss-classified word in the conjectured automaton means that somewhere a wrong transition is taken. Hence w admits a factorization $w = u \cdot b \cdot v$ where $u \in \Sigma^*$ and $b \in \Sigma$ is where the first wrong transition is taken. Obviously we do not know u and b in advance but know that this happens in the following two cases. Either b leads to an undiscovered state in the automaton of the target language, or letter b does not belong to the interval it was assumed to belong in the conjectured automaton. The latter case happens only when b does not belong to the evidence function. Since counter-example w is minimal, it admits

Procedure 3 Process counter-example

```

1: procedure COUNTER-EX( $w$ )
2:   Find a factorization  $w = u \cdot b \cdot v$ ,  $b \in \Sigma$ ,  $u, v \in \Sigma^*$  such that
3:      $\exists \mathbf{u} \in \mathbf{M}_T$ ,  $u \in \mu(\mathbf{u})$  and  $\forall \mathbf{u}' \in \mathbf{M}_T$ ,  $u \cdot b \notin \mu(\mathbf{u}')$ 
4:   if  $\mathbf{u} \in \mathbf{S}$  then ▷  $\mathbf{u}$  is already a state
5:     Find  $\mathbf{a} \in \Sigma_{\mathbf{u}}$  such that  $b \in [\mathbf{a}]$  ▷ refine  $[\mathbf{a}]$ 
6:      $\Sigma' = \Sigma \cup \{\mathbf{a}_{\text{new}}\}$ 
7:      $\mathbf{R}' = \mathbf{R} \cup \{\mathbf{u} \cdot \mathbf{a}_{\text{new}}\}$ 
8:      $\mu(\mathbf{u} \cdot \mathbf{a}_{\text{new}}) = \mu(\mathbf{u}) \cdot b$ 
9:     Ask MQ for all words in  $\{\mu(\mathbf{u} \cdot \mathbf{a}_{\text{new}}) \cdot e : e \in E\}$ 
10:     $\psi'_{\mathbf{u}}(a) = \begin{cases} \psi_{\mathbf{u}}(a) & \text{if } a \notin [\mathbf{a}] \\ \mathbf{a}_{\text{new}} & \text{if } a \in [\mathbf{a}] \text{ and } a \geq b \\ \mathbf{a} & \text{otherwise} \end{cases}$ 
11:     $T = (\Sigma, \Sigma', \mathbf{S}, \mathbf{R}', \psi', E, \mathbf{f}', \mu')$ 
12:  else ▷  $\mathbf{u}$  is in the boundary
13:     $\mathbf{S}' = \mathbf{S} \cup \{\mathbf{u}\}$  ▷ and becomes a state
14:    if  $b = a_0$  then
15:       $\Sigma' = \Sigma \cup \{\mathbf{a}_{\text{new}}\}$ 
16:       $\psi' = \psi \cup \{\psi_{\mathbf{u}}\}$ , with  $\psi_{\mathbf{u}}(a) = \mathbf{a}_{\text{new}}, \forall a \in \Sigma$ 
17:       $\mathbf{R}' = (\mathbf{R} - \{\mathbf{u}\}) \cup \{\mathbf{u} \cdot \mathbf{a}_{\text{new}}\}$ 
18:       $E' = E \cup \{\text{suffixes of } b \cdot v\}$ 
19:       $\mu(\mathbf{u} \cdot \mathbf{a}_{\text{new}}) = \mu(\mathbf{u}) \cdot a_0$ 
20:      Ask MQ for all words in  $\{\mu(\mathbf{u} \cdot \mathbf{a}_{\text{new}}) \cdot e : e \in E'\}$ 
21:    else
22:       $\Sigma' = \Sigma \cup \{\mathbf{a}_{\text{new}}, \mathbf{a}'_{\text{new}}\}$ 
23:       $\psi' = \psi \cup \{\psi_{\mathbf{u}}\}$ , with  $\psi_{\mathbf{u}}(a) = \begin{cases} \mathbf{a}'_{\text{new}} & \text{if } a \geq b \\ \mathbf{a}_{\text{new}} & \text{otherwise} \end{cases}$ 
24:       $\mathbf{R}' = (\mathbf{R} - \{\mathbf{u}\}) \cup \{\mathbf{u} \cdot \mathbf{a}_{\text{new}}, \mathbf{u} \cdot \mathbf{a}'_{\text{new}}\}$ 
25:       $E' = E \cup \{\text{suffixes of } b \cdot v\}$ 
26:       $\mu(\mathbf{u} \cdot \mathbf{a}_{\text{new}}) = \mu(\mathbf{u}) \cdot a_0$ ;  $\mu(\mathbf{u} \cdot \mathbf{a}'_{\text{new}}) = \mu(\mathbf{u}) \cdot b$ 
27:      Ask MQ for all words in  $\{(\mu(\mathbf{u} \cdot \mathbf{a}_{\text{new}}) \cup \mu(\mathbf{u} \cdot \mathbf{a}'_{\text{new}})) \cdot e : e \in E'\}$ 
28:    end if
29:     $T = (\Sigma, \Sigma', \mathbf{S}', \mathbf{R}', \psi', E', \mathbf{f}', \mu')$ 
30:  end if
31:  if  $T$  is not closed then
32:    CLOSE
33:  end if
34: end procedure

```

a factorization $w = u \cdot b \cdot v$, where u is the largest prefix of w such that $u \in \mu(\mathbf{u})$ for some $\mathbf{u} \in \mathbf{S} \cup \mathbf{R}$ but $s \cdot b \notin \mu(\mathbf{u}')$ for any word \mathbf{u}' in the symbolic sample. We consider two cases, $\mathbf{u} \in \mathbf{S}$ and $\mathbf{u} \in \mathbf{R}$.

In the first case, when $\mathbf{u} \in \mathbf{S}$, \mathbf{u} is already a state in the hypothesis but b indicates that the partition boundaries are not correctly defined and need refinement. That is, $u \cdot b$ was wrongly considered to be part of $[\mathbf{u} \cdot \mathbf{a}]$ for some $\mathbf{a} \in \Sigma_{\mathbf{u}}$, and thus b was wrongly considered

to be part of $[a]$. Due to minimality, all letters in $[a]$ less than letter b behave like $\mu(a)$. We assume that all remaining letters in $[a]$ behave like b and map them to a new symbol \mathbf{a}_{new} that we add to $\Sigma_{\mathbf{u}}$. We then update $\psi_{\mathbf{u}}$ such that $\psi'_{\mathbf{u}}(a) = \mathbf{a}_{\text{new}}$ for all $a \in [a], a \geq b$, and $\psi'_{\mathbf{u}}(a) = \psi_{\mathbf{u}}(a)$, otherwise. The evidence function is updated by letting $\mu(\mathbf{u} \cdot \mathbf{a}_{\text{new}}) = \mu(\mathbf{u}) \cdot b$ and $\mathbf{u} \cdot \mathbf{a}_{\text{new}}$ is added to \mathbf{R} .

In the second case, the symbolic word \mathbf{u} is part of the boundary. From the counterexample we deduce that \mathbf{u} is not equivalent to any of the existing states in the hypothesis and should form a new state. Specifically, we find the prefix \mathbf{s} that was considered to be equivalent to \mathbf{u} , that is $g(\mathbf{u}) = \mathbf{s} \in \mathbf{S}$. Since the table is reduced $\mathbf{f}_{\mathbf{u}} \neq \mathbf{f}_{\mathbf{s}'}$ for any other $\mathbf{s}' \in \mathbf{S}$. Because w is the shortest counter-example, the classification of $\mathbf{s} \cdot b \cdot v$ in the automaton is correct (otherwise $s \cdot b \cdot v$, for some $s \in [\mathbf{s}]$ would constitute a shorter counter-example) and different from that of $\mathbf{u} \cdot b \cdot v$. We conclude that \mathbf{u} is a new state, which is added to \mathbf{S} . To distinguish between \mathbf{u} and \mathbf{s} we add to E the word $b \cdot v$, possibly with some of its suffixes (see [BR04] for a more detailed discussion of counter-example processing).

As \mathbf{u} is a new state we need to add its continuations to \mathbf{R} . We distinguish two subcases depending on b . If $b = a_0$, the smallest element of Σ , then a new symbolic letter \mathbf{a}_{new} is added to Σ , with $[\mathbf{a}_{\text{new}}] = \Sigma$ and $\mu(\mathbf{u} \cdot \mathbf{a}_{\text{new}}) = \mu(\mathbf{u}) \cdot a_0$, and the symbolic word $\mathbf{u} \cdot \mathbf{a}_{\text{new}}$ is added to \mathbf{R} . If $b \neq a_0$ then *two* new symbolic letters, \mathbf{a}_{new} and \mathbf{a}'_{new} , are added to Σ with $[\mathbf{a}_{\text{new}}] = \{a : a < b\}$, $[\mathbf{a}'_{\text{new}}] = \{a : a \geq b\}$, $\mu(\mathbf{u} \cdot \mathbf{a}_{\text{new}}) = \mu(\mathbf{u}) \cdot a_0$ and $\mu(\mathbf{u} \cdot \mathbf{a}'_{\text{new}}) = \mu(\mathbf{u}) \cdot b$. The words $\mathbf{u} \cdot \mathbf{a}_{\text{new}}$ and $\mathbf{u} \cdot \mathbf{a}'_{\text{new}}$ are added to \mathbf{R} .

A detailed description of the algorithm is given in Algorithm 1 and its major procedures, table closing and counter-example treatment are described in Procedures 2 and 3 respectively. A statement of the form $\Sigma' = \Sigma \cup \{\mathbf{a}\}$ indicates the introduction of a new symbolic letter $\mathbf{a} \notin \Sigma$. We use MQ and EQ as shorthands for membership and equivalence queries, respectively. In the following we illustrate the symbolic algorithm as applied to a language over an infinite alphabet.

Example 4.1. Let $\Sigma = [0, 100) \subset \mathbb{R}$ with the usual order and let $L \subseteq \Sigma^*$ be a target language. Fig. 5 shows the evolution of the symbolic observation tables and Fig. 6 depicts the corresponding automata and the concrete semantics of the symbolic alphabets.

We initialize the table with $\mathbf{S} = \{\epsilon\}$, $\mathbf{R} = \{a_0\}$, $\mu(a_0) = \{0\}$ and $E = \{\epsilon\}$ and ask membership queries for ϵ (rejected) and 0 (accepted). The obtained table, \mathbf{T}_0 is not closed so we move a_0 to \mathbf{S} , introduce $\Sigma_{a_0} = \{a_1\}$, where a_1 is a new symbol, and add $a_0 \cdot a_1$ to \mathbf{R} with $\mu(a_0 \cdot a_1) = 0 \cdot 0$. Asking membership queries we obtain the closed table \mathbf{T}_1 and its automaton \mathcal{A}_1 . We pose an equivalence query and obtain $(50, -)$ as a (minimal) counter-example which implies that all words smaller than 50 are correctly classified. We add a new symbol a_2 to Σ_{ϵ} and redefine the concrete semantics to $[a_0] = \{a < 50\}$ and $[a_2] = \{a \geq 50\}$. As evidence we select the smallest possible letter, $\mu(a_2) = 50$, ask membership queries to obtain the closed table \mathbf{T}_2 and automaton \mathcal{A}_2 .

For this hypothesis we get a counter-example $(0 \cdot 30, -)$ whose prefix 0 is already in the sample, hence the misclassification occurs in the second transition. We refine the alphabet partition for state a_0 by introducing a new symbol a_3 and letting $[a_1] = \{a < 30\}$ and $[a_3] = \{a \geq 30\}$. Table \mathbf{T}_3 is closed but automaton \mathcal{A}_3 is still incorrect and a counter-example $(50 \cdot 0, -)$ is provided. The prefix 50 belongs to the evidence of a_2 and is moved from the boundary to become a new state and its successor $a_2 \cdot a_4$, for a new symbol a_4 , is added to \mathbf{R} . To distinguish a_2 from ϵ , the suffix 0 of the counter-example is added to E resulting in \mathbf{T}_4 which is not closed. The newly discovered state $a_0 \cdot a_1$ is added to \mathbf{S} , the filled table \mathbf{T}_5 is closed and the conjectured automaton \mathcal{A}_5 has two additional states.

T_0			T_1			T_2			T_3			T_4		
	ϵ			ϵ			ϵ			ϵ			ϵ	0
ϵ	-		ϵ	-		ϵ	-		ϵ	-		ϵ	-	+
\mathbf{a}_0	+		\mathbf{a}_0	+		\mathbf{a}_0	+		\mathbf{a}_0	+		\mathbf{a}_0	+	+
			$\mathbf{a}_0 \cdot \mathbf{a}_1$	+		$\mathbf{a}_0 \cdot \mathbf{a}_1$	+		$\mathbf{a}_0 \cdot \mathbf{a}_1$	+		$\mathbf{a}_0 \cdot \mathbf{a}_1$	+	-
			\mathbf{a}_2	-		\mathbf{a}_2	-		\mathbf{a}_2	-		\mathbf{a}_2	-	-
			$\mathbf{a}_0 \cdot \mathbf{a}_3$	-		$\mathbf{a}_0 \cdot \mathbf{a}_3$	-		$\mathbf{a}_0 \cdot \mathbf{a}_3$	-		$\mathbf{a}_0 \cdot \mathbf{a}_3$	-	-
			$\mathbf{a}_2 \cdot \mathbf{a}_4$	-		$\mathbf{a}_2 \cdot \mathbf{a}_4$	-		$\mathbf{a}_2 \cdot \mathbf{a}_4$	-		$\mathbf{a}_2 \cdot \mathbf{a}_4$	-	-
			$\mathbf{a}_0 \cdot \mathbf{a}_1 \cdot \mathbf{a}_5$	-		$\mathbf{a}_2 \cdot \mathbf{a}_6$	+		$\mathbf{a}_0 \cdot \mathbf{a}_3$	-		$\mathbf{a}_0 \cdot \mathbf{a}_1$	+	-
						$\mathbf{a}_0 \cdot \mathbf{a}_3$	-		$\mathbf{a}_2 \cdot \mathbf{a}_4$	-		$\mathbf{a}_0 \cdot \mathbf{a}_1 \cdot \mathbf{a}_5$	-	-
						$\mathbf{a}_2 \cdot \mathbf{a}_4$	-		$\mathbf{a}_0 \cdot \mathbf{a}_1 \cdot \mathbf{a}_5$	-		$\mathbf{a}_2 \cdot \mathbf{a}_6$	+	-
						$\mathbf{a}_0 \cdot \mathbf{a}_1 \cdot \mathbf{a}_5$	-		$\mathbf{a}_2 \cdot \mathbf{a}_6$	+		$\mathbf{a}_2 \cdot \mathbf{a}_7$	-	-
						$\mathbf{a}_2 \cdot \mathbf{a}_6$	+		$\mathbf{a}_2 \cdot \mathbf{a}_7$	-		$\mathbf{a}_2 \cdot \mathbf{a}_8$	+	+
						$\mathbf{a}_2 \cdot \mathbf{a}_7$	-		$\mathbf{a}_2 \cdot \mathbf{a}_8$	+				

FIGURE 5. Observation tables for Example 4.1.

Subsequent equivalence queries result counter-examples $(50 \cdot 20, +)$, $(50 \cdot 80, -)$ and $(50 \cdot 50 \cdot 0, +)$ which are used to refine the alphabet partition at state \mathbf{a}_2 and modify its outgoing transitions progressively as seen in automata \mathcal{A}_6 , \mathcal{A}_7 and \mathcal{A}_8 , respectively. Automaton \mathcal{A}_8 accepts the target language and the algorithm terminates. \square

Note that for the language in Example 1.3, the symbolic algorithm needs around 30 queries instead of the 80 queries required by L^* . If we choose to learn a language as the one described in Example 4.1, restricting the concrete alphabet to the finite alphabet $\Sigma = \{1, \dots, 100\}$, then L^* requires around 1000 queries compared to 17 queries required by our symbolic algorithm. As we shall see in Section 6, the complexity of the symbolic algorithm does not depend on the size of the concrete alphabet, only on the number of transitions.

5. LEARNING LANGUAGES OVER PARTIALLY-ORDERED ALPHABETS

In this section we sketch the extension of the results of this paper to partially-ordered alphabets of the form $\Sigma = X^d$ where X is a totally-ordered set such as an interval $[0, k) \subseteq \mathbb{R}$. Letters of Σ are d -tuples of the form $\mathbf{x} = (x_1, \dots, x_d)$ and the minimal element is $\mathbf{0} = (0, \dots, 0)$. The usual partial order on this set is defined as $\mathbf{x} \leq \mathbf{y}$ if and only if $x_i \leq y_i$ for all $i = 1, \dots, d$. When $\mathbf{x} \leq \mathbf{y}$ and $x_i \neq y_i$ for some i the inequality is strict, denoted by $\mathbf{x} < \mathbf{y}$, and we say then that \mathbf{x} *dominates* \mathbf{y} . Two elements are *incomparable*, denoted by $\mathbf{x} \parallel \mathbf{y}$, if $x_i < y_i$ and $x_j > y_j$ for some i and j .

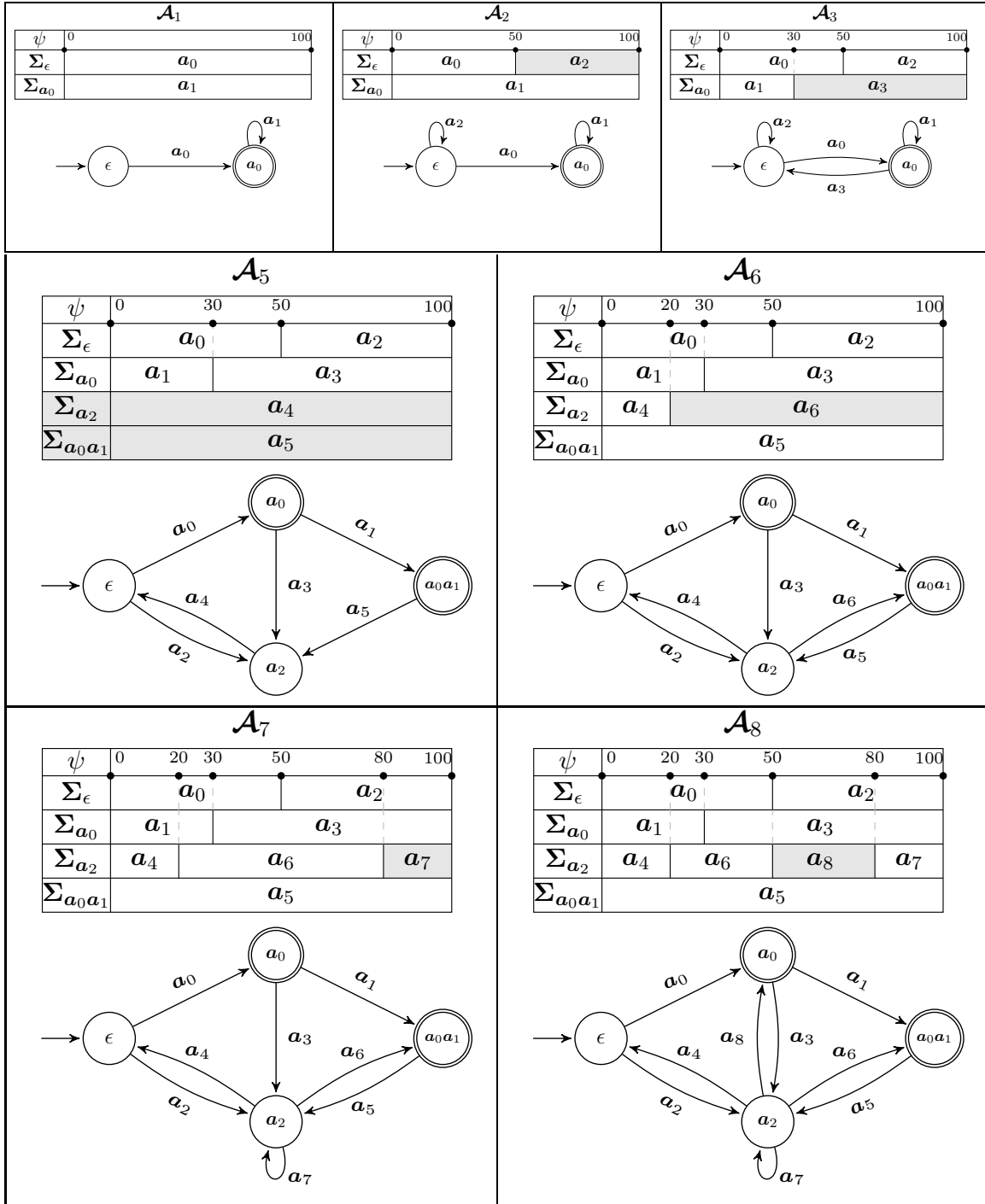


FIGURE 6. Hypotheses and Σ -semantics for Example 4.1

For partially-ordered sets, a natural extension of the partition of an ordered set into intervals is a *monotone* partition, where for each partition block P there are no three points

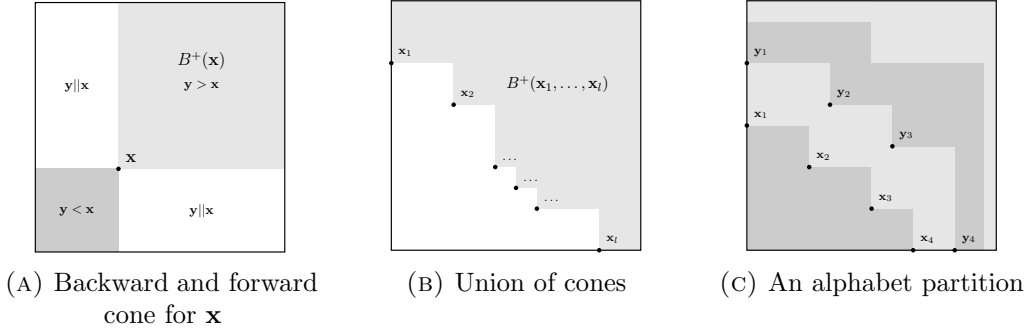
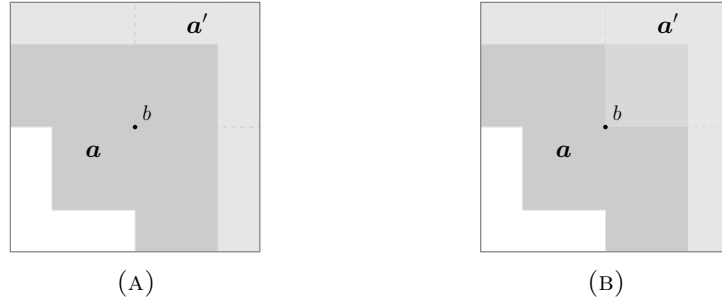


FIGURE 7

FIGURE 8. Modifying the alphabet partition for state \mathbf{u} after receiving $\mathbf{u} \cdot b \cdot v$ as counter-example. Letters above b are moved from $[\mathbf{a}]$ to $[\mathbf{a}']$.

such that $\mathbf{x} < \mathbf{y} < \mathbf{z}$, $\mathbf{x}, \mathbf{z} \in P$, and $\mathbf{y} \notin P$. We define in the following such partitions represented by a finite set of points.

A *forward cone* $B^+(\mathbf{x}) \subset \Sigma$ is the set of all points dominated by a point $\mathbf{x} \in \Sigma$ (see Fig. 7a). Let $F = \{\mathbf{x}_1, \dots, \mathbf{x}_l\}$ be a set of points, then $B^+(F) = B^+(\mathbf{x}_1) \cup \dots \cup B^+(\mathbf{x}_l)$ as shown in Fig. 7b. From a family of sets of points $\mathcal{F} = \{F_0, \dots, F_{m-1}\}$, such that $F_0 = \{\mathbf{0}\}$ satisfying for every i : 1) $\forall \mathbf{y} \in F_i, \exists \mathbf{x} \in F_{i-1}$ such that $\mathbf{x} < \mathbf{y}$, and 2) $\forall \mathbf{y} \in F_i, \forall \mathbf{x} \in F_{i-1}, \mathbf{y} \not< \mathbf{x}$, we can define a monotone partition of the form $\mathcal{P} = \{P_1, \dots, P_{m-1}\}$, where $P_i = B^+(F_{i-1}) - B^+(F_i)$, see Fig. 7c.

A subset P of Σ , as defined above, may have several mutually-incomparable minimal elements, none of which being dominated by any other element of P . One can thus apply the symbolic learning algorithm but without the presence of unique minimal evidence and minimal counter-example. For this reason a symbolic word may have more than one evidence. Evidence compatibility is preserved though due to the nature of the partition.

The teacher is assumed to return a counter-example chosen from a set of incomparable minimal counter-examples. Like in the algorithm for totally ordered alphabet, every counter-example either discovers a new state or refines a partition. The learning algorithm for partially-ordered alphabets is similar to Algorithm 1 and can be applied with only a minor modification in the treatment of the counterexamples and specifically in the refinement procedure. Lines 6-8 of Procedure 3 should be ignored in the case where there exists a symbolic letter \mathbf{a}' , as illustrated in Fig. 8a, such that $\mathbf{f}(\mathbf{u} \cdot b \cdot e) = \mathbf{f}(\mathbf{u} \cdot \mathbf{a}' \cdot e)$ for all $e \in E$. In such a case, function ψ is updated as in line 9 by replacing \mathbf{a}_{new} by \mathbf{a}' and b should

be added to $\mu(\mathbf{a}')$. In Fig. 8b, one can see the partition after refinement, where all letters above b have been moved from $[\mathbf{a}]$ to $[\mathbf{a}']$.

<p>T_0</p> <table border="1" style="margin: auto;"> <tr><td></td><td>ϵ</td></tr> <tr><td>ϵ</td><td>-</td></tr> <tr><td>\mathbf{a}_0</td><td>+</td></tr> <tr><td>$\mathbf{a}_0 \cdot \mathbf{a}_1$</td><td>+</td></tr> </table>		ϵ	ϵ	-	\mathbf{a}_0	+	$\mathbf{a}_0 \cdot \mathbf{a}_1$	+	<p>T_{1-3}</p> <table border="1" style="margin: auto;"> <tr><td></td><td>ϵ</td></tr> <tr><td>ϵ</td><td>-</td></tr> <tr><td>\mathbf{a}_0</td><td>+</td></tr> <tr><td>$\mathbf{a}_0 \cdot \mathbf{a}_1$</td><td>+</td></tr> <tr><td>\mathbf{a}_2</td><td>-</td></tr> </table>		ϵ	ϵ	-	\mathbf{a}_0	+	$\mathbf{a}_0 \cdot \mathbf{a}_1$	+	\mathbf{a}_2	-	<p>T_{4-7}</p> <table border="1" style="margin: auto;"> <tr><td></td><td>ϵ</td></tr> <tr><td>ϵ</td><td>-</td></tr> <tr><td>\mathbf{a}_0</td><td>+</td></tr> <tr><td>$\mathbf{a}_0 \cdot \mathbf{a}_1$</td><td>+</td></tr> <tr><td>\mathbf{a}_2</td><td>-</td></tr> <tr><td>$\mathbf{a}_0 \cdot \mathbf{a}_3$</td><td>-</td></tr> </table>		ϵ	ϵ	-	\mathbf{a}_0	+	$\mathbf{a}_0 \cdot \mathbf{a}_1$	+	\mathbf{a}_2	-	$\mathbf{a}_0 \cdot \mathbf{a}_3$	-	<p>T_8</p> <table border="1" style="margin: auto;"> <tr><td></td><td>ϵ</td><td>$\binom{0}{0}$</td></tr> <tr><td>ϵ</td><td>-</td><td>+</td></tr> <tr><td>\mathbf{a}_0</td><td>+</td><td>+</td></tr> <tr><td>\mathbf{a}_2</td><td>-</td><td>-</td></tr> <tr><td>$\mathbf{a}_0 \cdot \mathbf{a}_1$</td><td>+</td><td>-</td></tr> <tr><td>$\mathbf{a}_0 \cdot \mathbf{a}_3$</td><td>-</td><td>-</td></tr> <tr><td>$\mathbf{a}_2 \cdot \mathbf{a}_4$</td><td>-</td><td>+</td></tr> </table>		ϵ	$\binom{0}{0}$	ϵ	-	+	\mathbf{a}_0	+	+	\mathbf{a}_2	-	-	$\mathbf{a}_0 \cdot \mathbf{a}_1$	+	-	$\mathbf{a}_0 \cdot \mathbf{a}_3$	-	-	$\mathbf{a}_2 \cdot \mathbf{a}_4$	-	+																																																															
	ϵ																																																																																																																				
ϵ	-																																																																																																																				
\mathbf{a}_0	+																																																																																																																				
$\mathbf{a}_0 \cdot \mathbf{a}_1$	+																																																																																																																				
	ϵ																																																																																																																				
ϵ	-																																																																																																																				
\mathbf{a}_0	+																																																																																																																				
$\mathbf{a}_0 \cdot \mathbf{a}_1$	+																																																																																																																				
\mathbf{a}_2	-																																																																																																																				
	ϵ																																																																																																																				
ϵ	-																																																																																																																				
\mathbf{a}_0	+																																																																																																																				
$\mathbf{a}_0 \cdot \mathbf{a}_1$	+																																																																																																																				
\mathbf{a}_2	-																																																																																																																				
$\mathbf{a}_0 \cdot \mathbf{a}_3$	-																																																																																																																				
	ϵ	$\binom{0}{0}$																																																																																																																			
ϵ	-	+																																																																																																																			
\mathbf{a}_0	+	+																																																																																																																			
\mathbf{a}_2	-	-																																																																																																																			
$\mathbf{a}_0 \cdot \mathbf{a}_1$	+	-																																																																																																																			
$\mathbf{a}_0 \cdot \mathbf{a}_3$	-	-																																																																																																																			
$\mathbf{a}_2 \cdot \mathbf{a}_4$	-	+																																																																																																																			
<p>T_9</p> <table border="1" style="margin: auto;"> <tr><td></td><td>ϵ</td><td>$\binom{0}{0}$</td></tr> <tr><td>ϵ</td><td>-</td><td>+</td></tr> <tr><td>\mathbf{a}_0</td><td>+</td><td>+</td></tr> <tr><td>\mathbf{a}_2</td><td>-</td><td>-</td></tr> <tr><td>$\mathbf{a}_0 \cdot \mathbf{a}_1$</td><td>+</td><td>-</td></tr> <tr><td>$\mathbf{a}_0 \cdot \mathbf{a}_3$</td><td>-</td><td>-</td></tr> <tr><td>$\mathbf{a}_2 \cdot \mathbf{a}_4$</td><td>-</td><td>+</td></tr> <tr><td>$\mathbf{a}_0 \cdot \mathbf{a}_1 \cdot \mathbf{a}_5$</td><td>-</td><td>-</td></tr> </table>		ϵ	$\binom{0}{0}$	ϵ	-	+	\mathbf{a}_0	+	+	\mathbf{a}_2	-	-	$\mathbf{a}_0 \cdot \mathbf{a}_1$	+	-	$\mathbf{a}_0 \cdot \mathbf{a}_3$	-	-	$\mathbf{a}_2 \cdot \mathbf{a}_4$	-	+	$\mathbf{a}_0 \cdot \mathbf{a}_1 \cdot \mathbf{a}_5$	-	-	<p>T_{10-11}</p> <table border="1" style="margin: auto;"> <tr><td></td><td>ϵ</td><td>$\binom{0}{0}$</td></tr> <tr><td>ϵ</td><td>-</td><td>+</td></tr> <tr><td>\mathbf{a}_0</td><td>+</td><td>+</td></tr> <tr><td>\mathbf{a}_2</td><td>-</td><td>-</td></tr> <tr><td>$\mathbf{a}_0 \cdot \mathbf{a}_1$</td><td>+</td><td>-</td></tr> <tr><td>$\mathbf{a}_0 \cdot \mathbf{a}_3$</td><td>-</td><td>-</td></tr> <tr><td>$\mathbf{a}_2 \cdot \mathbf{a}_4$</td><td>-</td><td>+</td></tr> <tr><td>$\mathbf{a}_0 \cdot \mathbf{a}_1 \cdot \mathbf{a}_5$</td><td>-</td><td>-</td></tr> <tr><td>$\mathbf{a}_2 \cdot \mathbf{a}_6$</td><td>+</td><td>-</td></tr> </table>		ϵ	$\binom{0}{0}$	ϵ	-	+	\mathbf{a}_0	+	+	\mathbf{a}_2	-	-	$\mathbf{a}_0 \cdot \mathbf{a}_1$	+	-	$\mathbf{a}_0 \cdot \mathbf{a}_3$	-	-	$\mathbf{a}_2 \cdot \mathbf{a}_4$	-	+	$\mathbf{a}_0 \cdot \mathbf{a}_1 \cdot \mathbf{a}_5$	-	-	$\mathbf{a}_2 \cdot \mathbf{a}_6$	+	-	<p>T_{12-15}</p> <table border="1" style="margin: auto;"> <tr><td></td><td>ϵ</td><td>$\binom{0}{0}$</td></tr> <tr><td>ϵ</td><td>-</td><td>+</td></tr> <tr><td>\mathbf{a}_0</td><td>+</td><td>+</td></tr> <tr><td>\mathbf{a}_2</td><td>-</td><td>-</td></tr> <tr><td>$\mathbf{a}_0 \cdot \mathbf{a}_1$</td><td>+</td><td>-</td></tr> <tr><td>$\mathbf{a}_0 \cdot \mathbf{a}_3$</td><td>-</td><td>-</td></tr> <tr><td>$\mathbf{a}_2 \cdot \mathbf{a}_4$</td><td>-</td><td>+</td></tr> <tr><td>$\mathbf{a}_0 \cdot \mathbf{a}_1 \cdot \mathbf{a}_5$</td><td>-</td><td>-</td></tr> <tr><td>$\mathbf{a}_2 \cdot \mathbf{a}_6$</td><td>+</td><td>-</td></tr> <tr><td>$\mathbf{a}_2 \cdot \mathbf{a}_7$</td><td>-</td><td>-</td></tr> </table>		ϵ	$\binom{0}{0}$	ϵ	-	+	\mathbf{a}_0	+	+	\mathbf{a}_2	-	-	$\mathbf{a}_0 \cdot \mathbf{a}_1$	+	-	$\mathbf{a}_0 \cdot \mathbf{a}_3$	-	-	$\mathbf{a}_2 \cdot \mathbf{a}_4$	-	+	$\mathbf{a}_0 \cdot \mathbf{a}_1 \cdot \mathbf{a}_5$	-	-	$\mathbf{a}_2 \cdot \mathbf{a}_6$	+	-	$\mathbf{a}_2 \cdot \mathbf{a}_7$	-	-	<p>T_{16-18}</p> <table border="1" style="margin: auto;"> <tr><td></td><td>ϵ</td><td>$\binom{0}{0}$</td></tr> <tr><td>ϵ</td><td>-</td><td>+</td></tr> <tr><td>\mathbf{a}_0</td><td>+</td><td>+</td></tr> <tr><td>\mathbf{a}_2</td><td>-</td><td>-</td></tr> <tr><td>$\mathbf{a}_0 \cdot \mathbf{a}_1$</td><td>+</td><td>-</td></tr> <tr><td>$\mathbf{a}_0 \cdot \mathbf{a}_3$</td><td>-</td><td>-</td></tr> <tr><td>$\mathbf{a}_2 \cdot \mathbf{a}_4$</td><td>-</td><td>+</td></tr> <tr><td>$\mathbf{a}_0 \cdot \mathbf{a}_1 \cdot \mathbf{a}_5$</td><td>-</td><td>-</td></tr> <tr><td>$\mathbf{a}_2 \cdot \mathbf{a}_6$</td><td>+</td><td>-</td></tr> <tr><td>$\mathbf{a}_2 \cdot \mathbf{a}_7$</td><td>-</td><td>-</td></tr> <tr><td>$\mathbf{a}_2 \cdot \mathbf{a}_8$</td><td>+</td><td>+</td></tr> </table>		ϵ	$\binom{0}{0}$	ϵ	-	+	\mathbf{a}_0	+	+	\mathbf{a}_2	-	-	$\mathbf{a}_0 \cdot \mathbf{a}_1$	+	-	$\mathbf{a}_0 \cdot \mathbf{a}_3$	-	-	$\mathbf{a}_2 \cdot \mathbf{a}_4$	-	+	$\mathbf{a}_0 \cdot \mathbf{a}_1 \cdot \mathbf{a}_5$	-	-	$\mathbf{a}_2 \cdot \mathbf{a}_6$	+	-	$\mathbf{a}_2 \cdot \mathbf{a}_7$	-	-	$\mathbf{a}_2 \cdot \mathbf{a}_8$	+	+
	ϵ	$\binom{0}{0}$																																																																																																																			
ϵ	-	+																																																																																																																			
\mathbf{a}_0	+	+																																																																																																																			
\mathbf{a}_2	-	-																																																																																																																			
$\mathbf{a}_0 \cdot \mathbf{a}_1$	+	-																																																																																																																			
$\mathbf{a}_0 \cdot \mathbf{a}_3$	-	-																																																																																																																			
$\mathbf{a}_2 \cdot \mathbf{a}_4$	-	+																																																																																																																			
$\mathbf{a}_0 \cdot \mathbf{a}_1 \cdot \mathbf{a}_5$	-	-																																																																																																																			
	ϵ	$\binom{0}{0}$																																																																																																																			
ϵ	-	+																																																																																																																			
\mathbf{a}_0	+	+																																																																																																																			
\mathbf{a}_2	-	-																																																																																																																			
$\mathbf{a}_0 \cdot \mathbf{a}_1$	+	-																																																																																																																			
$\mathbf{a}_0 \cdot \mathbf{a}_3$	-	-																																																																																																																			
$\mathbf{a}_2 \cdot \mathbf{a}_4$	-	+																																																																																																																			
$\mathbf{a}_0 \cdot \mathbf{a}_1 \cdot \mathbf{a}_5$	-	-																																																																																																																			
$\mathbf{a}_2 \cdot \mathbf{a}_6$	+	-																																																																																																																			
	ϵ	$\binom{0}{0}$																																																																																																																			
ϵ	-	+																																																																																																																			
\mathbf{a}_0	+	+																																																																																																																			
\mathbf{a}_2	-	-																																																																																																																			
$\mathbf{a}_0 \cdot \mathbf{a}_1$	+	-																																																																																																																			
$\mathbf{a}_0 \cdot \mathbf{a}_3$	-	-																																																																																																																			
$\mathbf{a}_2 \cdot \mathbf{a}_4$	-	+																																																																																																																			
$\mathbf{a}_0 \cdot \mathbf{a}_1 \cdot \mathbf{a}_5$	-	-																																																																																																																			
$\mathbf{a}_2 \cdot \mathbf{a}_6$	+	-																																																																																																																			
$\mathbf{a}_2 \cdot \mathbf{a}_7$	-	-																																																																																																																			
	ϵ	$\binom{0}{0}$																																																																																																																			
ϵ	-	+																																																																																																																			
\mathbf{a}_0	+	+																																																																																																																			
\mathbf{a}_2	-	-																																																																																																																			
$\mathbf{a}_0 \cdot \mathbf{a}_1$	+	-																																																																																																																			
$\mathbf{a}_0 \cdot \mathbf{a}_3$	-	-																																																																																																																			
$\mathbf{a}_2 \cdot \mathbf{a}_4$	-	+																																																																																																																			
$\mathbf{a}_0 \cdot \mathbf{a}_1 \cdot \mathbf{a}_5$	-	-																																																																																																																			
$\mathbf{a}_2 \cdot \mathbf{a}_6$	+	-																																																																																																																			
$\mathbf{a}_2 \cdot \mathbf{a}_7$	-	-																																																																																																																			
$\mathbf{a}_2 \cdot \mathbf{a}_8$	+	+																																																																																																																			

FIGURE 9. Observation tables for Example 5.1

Example 5.1. Let us illustrate the learning process for a target language L defined over $\Sigma = [0, 100]^2$. All tables, hypotheses automata and alphabet partitions for this example are shown in Figures 9, 10, and 11, respectively.

The learner starts asking MQs for the empty word. A symbolic letter \mathbf{a}_0 is chosen to represent its continuations with the minimal element of Σ as evidence, i.e., $\mu(\mathbf{a}_0) = \binom{0}{0}$. The symbolic word \mathbf{a}_0 is moved to \mathbf{S} for the table T_0 to be closed. The symbolic letter \mathbf{a}_1 is added to the alphabet of state \mathbf{a}_0 , and the learner asks a MQ for $\binom{0}{0} \binom{0}{0}$, the evidence of the symbolic word $\mathbf{a}_0 \mathbf{a}_1$. The first hypothesis automaton is \mathcal{A}_0 with Σ -semantics $[\mathbf{a}_0] = [\mathbf{a}_1] = \Sigma$. The counter-example $(\binom{45}{50}, -)$ refines the partition for the initial state. The symbolic alphabet is extended to $\Sigma_\epsilon = \{\mathbf{a}_0, \mathbf{a}_2\}$ with $[\mathbf{a}_2] = \{x \succ \binom{45}{50}\}$, $[\mathbf{a}_0] = \Sigma - [\mathbf{a}_2]$, and $\mu(\mathbf{a}_2) = \binom{45}{50}$. The new observation table and hypothesis are T_1 and \mathcal{A}_1 . Two more counter-examples will come to refine the partition for the initial state, $(\binom{0}{70}, -)$ and $(\binom{0}{70}, -)$, that will modify the partition for the initial state, moving all letters greater than $\binom{60}{0}$ and $\binom{0}{70}$ to the Σ -semantics of \mathbf{a}_2 as can be seen in ψ_2 and ψ_3 respectively.

After the hypothesis \mathcal{A}_3 , the counter-example $(\binom{0}{0} \binom{0}{80}, -)$ adds a new symbol \mathbf{a}_3 and a new transition in the hypothesis automaton. The counter-examples that follow, namely, $(\binom{0}{0} \binom{80}{0}, -)$, $(\binom{0}{0} \binom{40}{15}, -)$, and $(\binom{0}{0} \binom{30}{30}, -)$ refine the Σ -semantics for symbols in $\Sigma_{\mathbf{a}_0}$ as shown in ψ_{4-7} .

Then counter-example $(\binom{45}{50} \binom{0}{0}, +)$ is presented. As we can see, the prefix $\binom{45}{50}$ exist already in $\mu(\mathbf{a}_2)$ and $\mathbf{a}_2 \in \mathbf{R}$ which means \mathbf{a}_2 becomes a state, and to distinguish it from

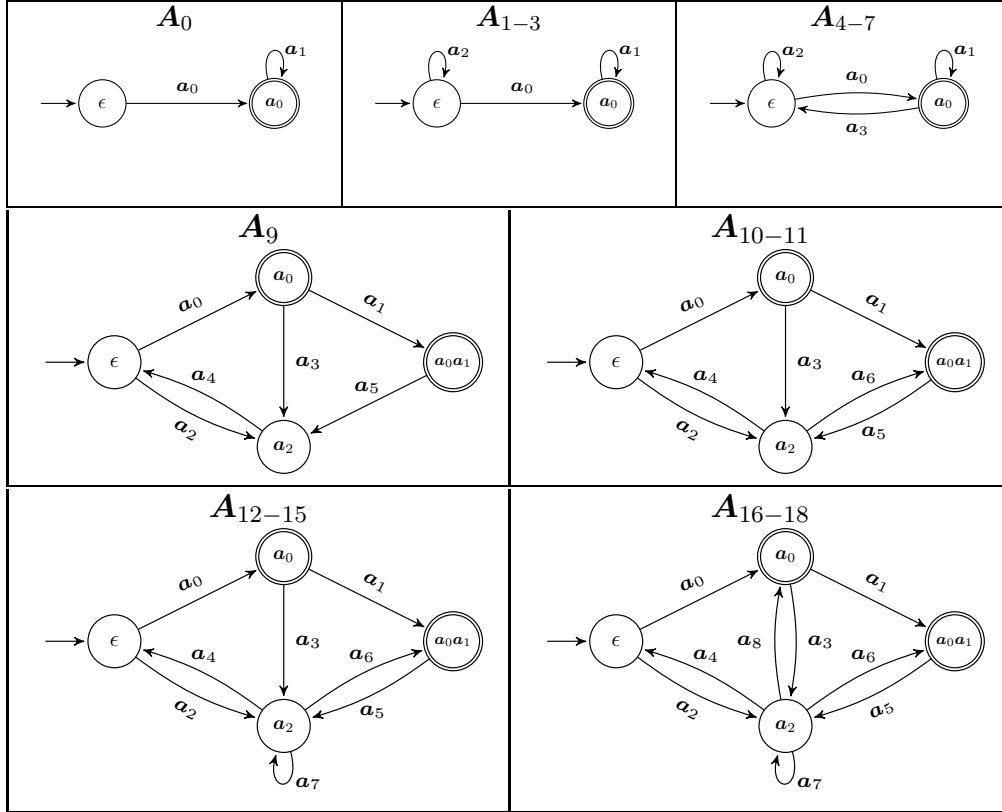


FIGURE 10. Hypothesis automata for Example 5.1

the state represented by the empty word the learner adds to E the suffix of the counter-example $\binom{0}{0}$. The resulting table \mathbf{T}_8 is not closed and $\mathbf{a}_0\mathbf{a}_1$ is moved to \mathbf{S} . The new table \mathbf{T}_9 is closed and evidence compatible. The hypothesis \mathcal{A}_9 has now four states and the symbolic alphabet and Σ -semantics for each state can be seen in ψ_9 . The counter-examples that follow will refine the partition at state \mathbf{a}_2 . The new transitions discovered and all refinements are shown in \mathcal{A}_{10-18} and $\psi_{10} - \psi_{18}$. The language was learned using 20 membership queries and 17 counter-examples. \square

6. ON COMPLEXITY

The complexity of the symbolic algorithm is influenced not by the size of the alphabet but by the resolution (partition size) with which we observe it. Let $L \subset \Sigma$ be the target language and let \mathcal{A} be the minimal symbolic automaton recognizing this language with state set Q of size n and a symbolic alphabet $\Sigma = \bigsqcup_q \Sigma_q$ such that $|\Sigma_q| \leq m$ for every q .

Each counter-example improves the hypothesis in one out of two ways. Either a new state is discovered or a partition gets refined. Hence, at most $n - 1$ equivalence queries of the first type can be asked and $n(m - 1)$ of the second, resulting in $\mathcal{O}(mn)$ equivalence queries.

Concerning the size of the table, the set of prefixes \mathbf{S} is monotonically increasing and reaches the size of exactly n elements. Since the table, by construction, is always kept

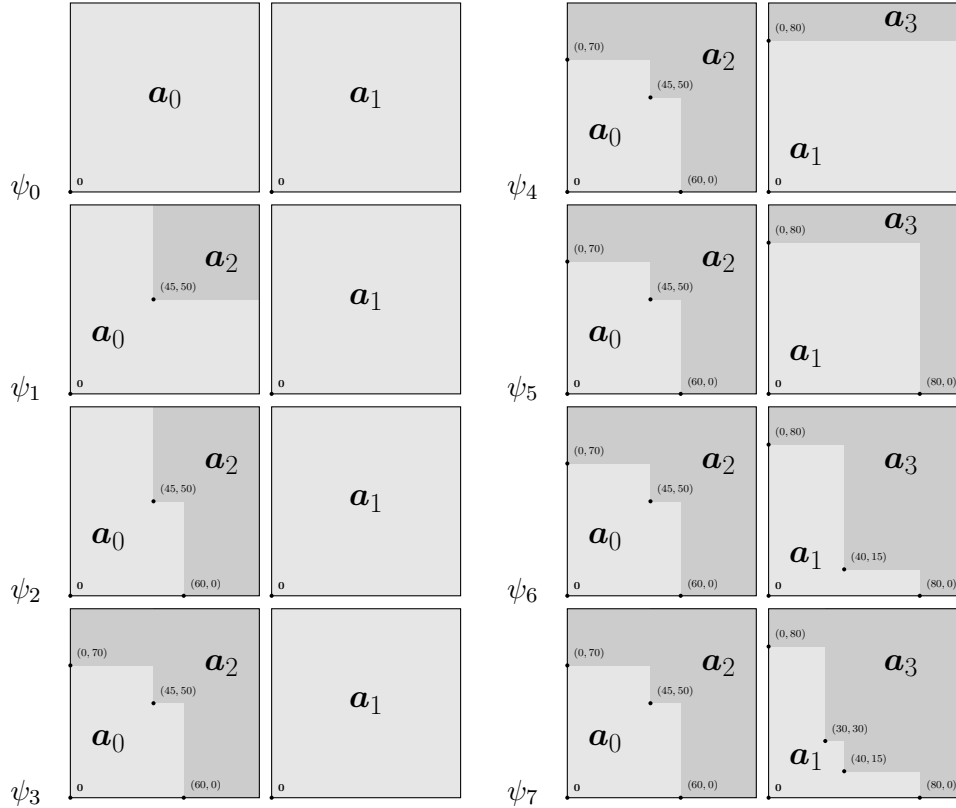


FIGURE 11. Alphabet partition for Example 5.1 (part 1)

reduced, the elements in \mathbf{S} represent exactly the states of the automaton. The size of the boundary is always smaller than the total number of transitions in the automaton, that is $mn - n + 1$. The number of suffixes in E , playing a distinguishing role for the states of the automaton, range between $\log_2 n$ and n . Hence, the size of the table ranges between $(n + m) \log_2 n$ and $n(mn + 1)$.

For a totally ordered alphabet the size of the concrete sample coincides with the size of the symbolic sample associated with the table and hence the number of membership queries asked is $\mathcal{O}(mn^2)$. For a partially ordered alphabet with each F_i defined by at most l points, some additional queries are asked. For every row in \mathbf{S} , at most $n(m - 1)(l - 1)$ additional words are added to the concrete sample, hence more membership queries might need to be asked. Furthermore, at most $l - 1$ more counter-examples are given to refine a partition. To conclude, the number of queries in total asked to learn language L is $\mathcal{O}(mn^2)$ if $l < n$ and $\mathcal{O}(lmn)$ otherwise.

7. CONCLUSION

We have defined a generic algorithmic scheme for automaton learning, targeting languages over large alphabets that can be recognized by finite symbolic automata having a modest number of states and transitions. Some ideas similar to ours have been proposed for the particular case of parametric languages [BJR06] and recently in a more general setting

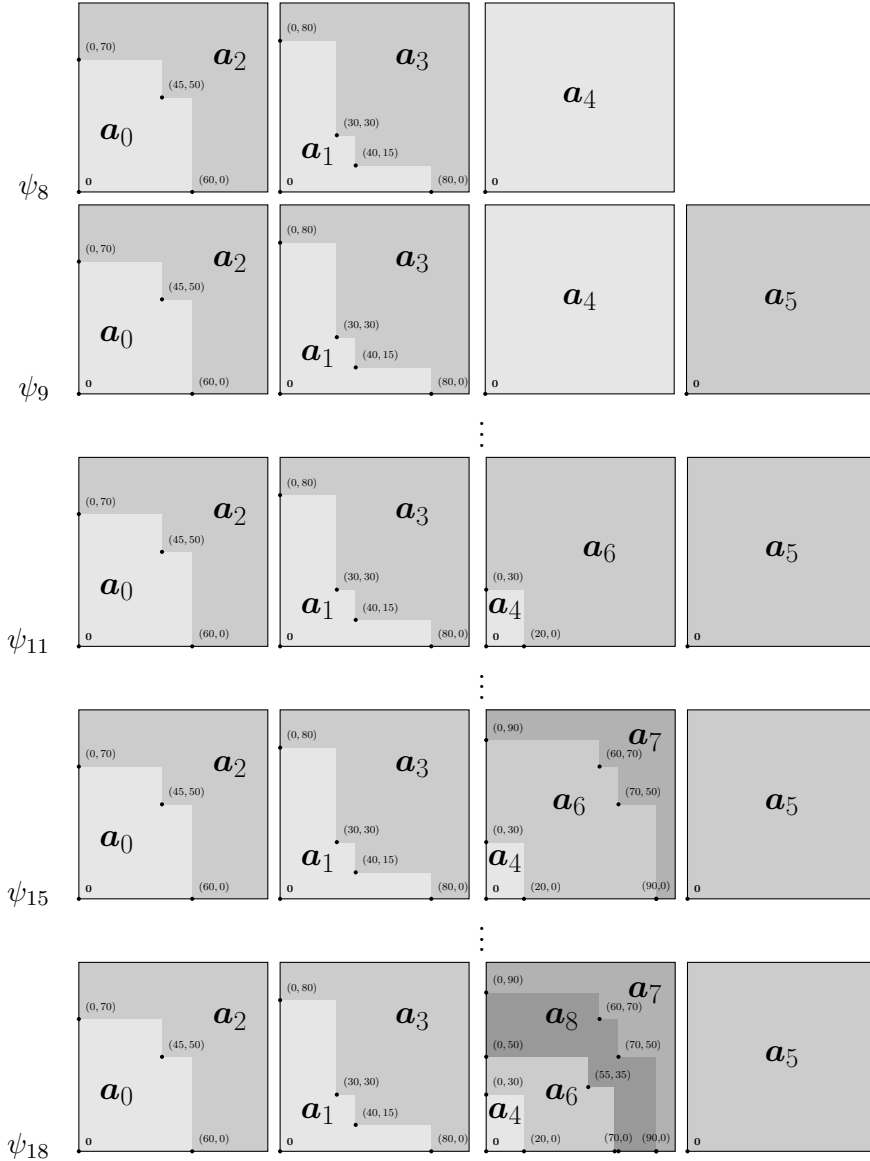


FIGURE 11. Alphabet partition for Example 5.1 (part 2)

[HSM11, IHS13, BB13] including partial evidential support and alphabet refinement during the learning process.

The genericity of our algorithm is due to a semantic approach (alphabet partitions) but of course, each and every domain will have its own semantic and syntactic specialization in terms of the size and shape of the alphabet partitions. In this work we have implemented an instantiation of this scheme for alphabets such as (\mathbb{N}, \leq) and (\mathbb{R}, \leq) . When dealing with numbers, the partition into a finite number of intervals (and monotone sets in higher dimensions) is very natural and used in many application domains ranging from quantization of sensor readings to income tax regulations. It will be interesting to compare the expressive

power and succinctness of symbolic automata with other approaches for representing numerical time series and to compare our algorithm with other inductive inference techniques for sequences of numbers.

As a first excursion into the domain, we have made quite strong assumptions on the nature of the equivalence oracle, which, already for small alphabets, is a bit too strong and pedagogical to be realistic. We assumed that it provides the shortest counter-example and also that it chooses always the minimal available concrete symbol. We can relax the latter (or both) and even omit this oracle altogether and replace it by random sampling, as already proposed in [Ang87] for concrete learning. Over large alphabets, it might be even more appropriate to employ probabilistic convergence criteria a-la *PAC learning* [Val84] and be content with a correct classification of a large fraction of the words, thus tolerating imprecise tracing of boundaries in the alphabet partitions. This topic is subject to ongoing work. Another challenging research direction is the adaptation of our framework to languages over Boolean vectors.

ACKNOWLEDGEMENT

This work was supported by the French project EQINOCS (ANR-11-BS02-004). We thank Peter Habermehl, Eugene Asarin and anonymous referees for useful comments and pointers to the literature.

REFERENCES

- [Ang87] Dana Angluin. Learning regular sets from queries and counterexamples. *Information and Computation*, 75(2):87–106, 1987.
- [BB13] Matko Botinčan and Domagoj Babić. Sigma*: Symbolic learning of Input-Output specifications. In *POPL*, pages 443–456. ACM, 2013.
- [BJR06] Therese Berg, Bengt Jonsson, and Harald Raffelt. Regular inference for state machines with parameters. In *FASE*, volume 3922 of *LNCS*, pages 107–121. Springer, 2006.
- [BLP10] Michael Benedikt, Clemens Ley, and Gabriele Puppis. What you must remember when processing data words. In *AMW*, volume 619 of *CEUR Workshop Proceedings*, 2010.
- [BR04] Therese Berg and Harald Raffelt. Model checking. In *Model-Based Testing of Reactive Systems*, volume 3472 of *LNCS*, pages 557–603. Springer, 2004.
- [DIH10] Colin De la Higuera. *Grammatical inference: learning automata and grammars*. Cambridge University Press, 2010.
- [DR95] Volker Diekert and Grzegorz Rozenberg. *The Book of Traces*. World Scientific, 1995.
- [DV14] Loris D’Antoni and Margus Veanes. Minimization of symbolic automata. In *POPL*, pages 541–554. ACM, 2014.
- [Gol72] E. Mark Gold. System identification via state characterization. *Automatica*, 8(5):621–636, 1972.
- [HJJ⁺95] Jesper G. Henriksen, Ole J.L. Jensen, Michael E. Jrgensen, Nils Klarlund, Robert Paige, Theis Rauhe, and Anders B. Sandholm. Mona: Monadic second-order logic in practice. In *TACAS*, volume 1019 of *LNCS*, pages 80–110. Springer, 1995.
- [HSJC12] Falk Howar, Bernhard Steffen, Bengt Jonsson, and Sofia Cassel. Inferring canonical register automata. In *VMCAI*, volume 7148 of *LNCS*, pages 251–266. Springer, 2012.
- [HSM11] Falk Howar, Bernhard Steffen, and Maik Merten. Automata learning with automated alphabet abstraction refinement. In *VMCAI*, volume 6538 of *LNCS*, pages 263–277. Springer, 2011.
- [HV11] Pieter Hooimeijer and Margus Veanes. An evaluation of automata algorithms for string analysis. In *VMCAI*, volume 6538 of *LNCS*, pages 248–262. Springer, 2011.
- [IHS13] Malte Isberner, Falk Howar, and Bernhard Steffen. Inferring automata with state-local alphabet abstractions. In *NASA Formal Methods*, volume 7871 of *LNCS*, pages 124–138. Springer, 2013.

- [KF94] Michael Kaminski and Nissim Francez. Finite-memory automata. *Theoretical Computer Science*, 134(2):329–363, 1994.
- [MM14] Oded Maler and Irini-Eleftheria Mens. Learning regular languages over large alphabets. In *TACAS*, volume 8413 of *LNCS*, pages 485–499. Springer, 2014.
- [Moo56] Edward F Moore. Gedanken-experiments on sequential machines. In *Automata studies*, volume 34 of *Annals of Mathematical Studies*, pages 129–153. Princeton, 1956.
- [MP95] Oded Maler and Amir Pnueli. On the learnability of infinitary regular sets. *Information and Computation*, 118(2):316–326, 1995.
- [Ner58] Anil Nerode. Linear automaton transformations. *Proceedings of the American Mathematical Society*, 9(4):541–544, 1958.
- [Val84] Leslie G. Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984.
- [VB12] Margus Veanes and Nikolaž Björner. Symbolic automata: The toolkit. In *TACAS*, volume 7214 of *LNCS*, pages 472–477. Springer, 2012.
- [VHL⁺12] Margus Veanes, Pieter Hooimeijer, Benjamin Livshits, David Molnar, and Nikolaž Björner. Symbolic finite state transducers: algorithms and applications. In *POPL*, pages 137–150. ACM, 2012.