# Communauté UNIVERSITÉ Grenoble Alpes

## THÈSE

Pour obtenir le grade de

### DOCTEUR DE LA COMMUNAUTÉ UNIVERSITÉ GRENOBLE ALPES

Spécialité : Mathématiques et Informatique

Arrêté ministérial : 25 mai 2016

Présentée par

## Irini-Eleftheria Mens

Thèse dirigée par Oded Maler,

préparée au sein du Laboratoire VERIMAG dans l'Ecole Doctorale Mathématiques, Sciences et Technologies de l'Information, Informatique

# Learning Regular Languages over Large Alphabets

# Apprentissage de langages réguliers sur des alphabets de grandes tailles

Thèse soutenue publiquement le **10 October 2017**, devant le jury composé de :

#### Monsieur Oded Maler

DIRECTEUR DE RECHERCHE, CNRS DÉLÉGATION ALPES, Directeur de thèse

Madame Dana Angluin PROFESSEURE, UNIVERSITÉ YALE A NEW HAVEN - USA, Rapporteur Monsieur Peter Habermehl MAÎTRE DE CONFÉRENCES, UNIVERSITÉ PARIS 7, Rapporteur Monsieur Eric Gaussier PROFESSEUR, UNIVERSITÉ GRENOBLE ALPES, Président Monsieur Frits W. Vaandrager PROFESSEUR, UNIV. RADBOUD DE NIMEGUE - PAYS-BAS, Examinateur Monsieur Laurent Fribourg DIRECTEUR DE RECHERCHE, CNRS DÉLÉGATION PARIS-VILLEJUIF, Examinateur



## Abstract

Learning regular languages is a branch of machine learning, a domain which has proved useful in many areas, including artificial intelligence, neural networks, data mining, verification, etc. In addition, interest in languages defined over large and infinite alphabets has increased in recent years. Although many theories and properties generalize well from the finite case, learning such languages is not an easy task. As the existing methods for learning regular languages depend on the size of the alphabet, a straightforward generalization in this context is not possible.

In this thesis, we present a generic algorithmic scheme that can be used for learning languages defined over large or infinite alphabets, such as bounded subsets of  $\mathbb{N}$  or  $\mathbb{R}$  or Boolean vectors of high dimensions. We restrict ourselves to the class of languages accepted by deterministic symbolic automata that use predicates to label transitions, forming a finite partition of the alphabet for every state.

Our learning algorithm, an adaptation of Angluin's  $L^*$ , combines standard automaton learning by state characterization, with the learning of the static predicates that define the alphabet partitions. We use the online learning scheme, where two types of queries provide the necessary information about the target language. The first type, membership queries, answer whether a given word belongs or not to the target. The second, equivalence queries, check whether a conjectured automaton accepts the target language and provide a counter-example otherwise.

We study language learning over large or infinite alphabets within a general framework but our aim is to provide solutions for particular concrete instances. For this, we focus on the two main aspects of the problem. Initially, we assume that equivalence queries always provide a counter-example which is minimal in the length-lexicographic order when the conjecture automaton is incorrect. Then, we drop this "strong" equivalence oracle and replace it by a more realistic assumption, where equivalence is approximated by testing queries, which use sampling on the set of words. Such queries are not guaranteed to find counter-examples and certainly not minimal ones. In this case, we obtain the weaker notion of PAC (probably approximately correct) learnability and learn an approximation of the target language. All proposed algorithms have been implemented and their performance, as a function of automaton and alphabet size, has been empirically evaluated.

## Résumé

L'apprentissage de langages réguliers est une branche de l'apprentissage automatique qui s'est révélée utile dans de nombreux domaines tels que l'intelligence artificielle, les réseaux de neurones, l'exploration de données, la vérification, etc. De plus, l'intérêt dans les langages définis sur des alphabets infinis ou de grande taille s'est accru au fil des années. Même si plusieurs propriétés et théories se généralisent à partir du cas fini, l'apprentissage de tels langages est une tâche difficile. En effet, dans ce contexte, l'application naïve des algorithmes d'apprentissage traditionnel n'est pas possible.

Dans cette thèse, nous présentons un schéma algorithmique général pour l'apprentissage de langages définis sur des alphabets infinis ou de grande taille, comme par exemple des sous-ensembles bornés de  $\mathbb{N}$  or  $\mathbb{R}$  ou des vecteurs booléens de grandes dimensions. Nous nous restreignons aux classes de langages qui sont acceptés par des automates déterministes symboliques utilisant des prédicats pour définir les transitions, construisant ainsi une partition finie de l'alphabet pour chaque état.

Notre algorithme d'apprentissage, qui est une adaptation du  $L^*$  d'Angluin, combine l'apprentissage classique d'un automate par la caractérisation de ses états, avec l'apprentissage de prédicats statiques définissant les partitions de l'alphabet. Nous utilisons l'apprentissage incrémental avec la propriété que deux types de requêtes fournissent une information suffisante sur le langage cible. Les requêtes du premier type sont les requêtes d'appartenance, qui permettent de savoir si un mot proposé appartient ou non au langage cible. Les requêtes du second type sont les requêtes d'équivalence, qui vérifient si un automate proposé accepte le langage cible; dans le cas contraire, un contre-exemple est renvoyé.

Nous étudions l'apprentissage de langages définis sur des alphabets infinis ou de grande tailles dans un cadre théorique et général, mais notre objectif est de proposer des solutions concrètes pour un certain nombre de cas particuliers. Ensuite, nous nous intéressons aux deux principaux aspects du problème. Dans un premier temps, nous supposerons que les requêtes d'équivalence renvoient toujours un contre-exemple minimal pour un ordre de longueur-lexicographique quand l'automate proposé est incorrect. Puis dans un second temps, nous relâchons cette hypothèse forte d'un oracle d'équivalence, et nous la remplaçons avec une hy-

pothèse plus réaliste où l'équivalence est approchée par un test sur les requêtes qui utilisent un échantillonnage sur l'ensemble des mots. Dans ce dernier cas, ce type de requêtes ne garantit pas l'obtention de contre-exemples, et par conséquent de contre-exemples minimaux. Nous obtenons alors une notion plus faible d'apprentissage PAC (*Probably Approximately Correct*), permettant l'apprentissage d'une approximation du langage cible. Tout les algorithmes ont été implémentés, et leurs performances, en terme de construction d'automate et de taille d'alphabet, ont été évaluées empiriquement.

# Contents

Abstract						
Résumé iii						
No	otatio	n	ix			
1	Introduction					
	1.1	Synopsis	1			
	1.2	Motivation	3			
	1.3	Outline	4			
2	2 Preliminaries					
	2.1	Regular Languages and Automata	7			
	2.2	Sets and Partitions	10			
	2.3 Learning Partitions					
		2.3.1 Learning Binary Decision Trees	15			
3	Language Identification					
	3.1	Connection with Machine Learning	21			
	3.2	Learning Languages	23			
	3.3	The $L^*$ Learning Algorithm	24			
		3.3.1 Observation table	25			
		3.3.2 The Learning Algorithm	26			
4	4 Symbolic Automata		33			
	4.1	State of the Art	33			
	4.2	Definition	34			
	4.3	Operations on Symbolic Automata	36			

4.4	Alphabets and Partitions	39		
	4.4.1 Interval Automata	39		
	4.4.2 Automata over Partially-ordered Alphabets	41		
	4.4.3 Boolean Vectors	41		
Learning Symbolic Automata				
5.1	Definitions	44		
5.2	Comparison to Related Work	47		
5.3	The Symbolic Learning Algorithm	49		
Lea	rning with a Helpful Teacher	57		
6.1	Learning Languages over Ordered Alphabets	57		
6.2	Learning over Partially-ordered Alphabets	63		
Lea	rning without a Helpful Teacher	69		
7.1	Approximating the Equivalence Query	70		
7.2	Learning Languages over $\mathbb{N}, \mathbb{R}$	71		
7.3	Learning Languages over $\mathbb{B}^n$	76		
Theoretical Analysis				
8.1	Updating the Hypothesis: Counter-Examples	83		
8.2	Hypothesis Error	84		
	8.2.1 A Probability Distribution on $\Sigma^*$	85		
	8.2.2 Computing the Relative Volumes	86		
8.3	Complexity and Termination	88		
	8.3.1 Using a Helpful Teacher (Minimal Counter-Examples)	88		
	8.3.2 Equivalence using Random Tests	89		
Emp	pirical Results	93		
9.1	General Comments on the Implementation	93		
9.2	On the Behavior of the Symbolic Learning Algorithm	95		
9.3	Comparison with Other Algorithms	97		
9.4	Learning Passwords	102		
9.5	Learning over the Booleans	104		
9.6	Comparing Boolean Vectors to Numerical Alphabets	107		
9.7	Conclusions	111		
Con	clusions and Future Work	113		
Bibliography				
	<ul> <li>4.4</li> <li>Lean</li> <li>5.1</li> <li>5.2</li> <li>5.3</li> <li>Lean</li> <li>6.1</li> <li>6.2</li> <li>Lean</li> <li>7.1</li> <li>7.2</li> <li>7.3</li> <li>The</li> <li>8.1</li> <li>8.2</li> <li>8.3</li> <li>Emp</li> <li>9.1</li> <li>9.2</li> <li>9.3</li> <li>9.4</li> <li>9.5</li> <li>9.6</li> <li>9.7</li> <li>Con</li> </ul>	4.4       Alphabets and Partitions       4.4.1         4.4.1       Interval Automata       4.4.2         Automata over Partially-ordered Alphabets       4.4.3         Boolean Vectors       4.4.3         Boolean Vectors       5.1         Learning Symbolic Automata       5.1         5.1       Definitions       5.1         5.2       Comparison to Related Work       5.3         5.3       The Symbolic Learning Algorithm       5.1         Learning with a Helpful Teacher       6.1       Learning Languages over Ordered Alphabets         6.2       Learning over Partially-ordered Alphabets       6.2         Learning without a Helpful Teacher       7.1       Approximating the Equivalence Query       7.2         7.3       Learning Languages over $\mathbb{B}^n$ 7.3       Learning Languages over $\mathbb{B}^n$ 7.3         7.3       Learning Languages over $\mathbb{B}^n$ 8.1       Updating the Hypothesis: Counter-Examples       8.2         8.1       Updating the Hypothesis: Counter-Examples       8.2       8.2.1       A Probability Distribution on $\Sigma^*$ 8.2.2       Computing the Relative Volumes       8.3.2       Equivalence using Random Tests       8.3.2       Equivalence using Random Tests       8.3.2       Equivalence using Random Tests       8.3.2		

# **List of Algorithms**

1	Grow Tree	16
2	Main Learning Algorithm	27
3	Table Closing	28
4	Make Table Consistent	28
5	Counter-Example Treatment: $L^*$	28
6	Counter-Example Treatment: Adding Suffixes	30
7	Counter-Example Treatment: Breakpoint	30
8	Symbolic Learning Algorithm	50
9	Table Initialization	50
10	Symbol Initialization	51
11	Table Closing	51
12	Make Evidence Compatible	52
13	Counter-Example Treatment	55
14	Counter-Example Treatment (with Helpful Teacher) - $\mathbb{R}$	59
15	Counter-Example Treatment (with Helpful Teacher) - $\mathbb{R}^n$	64
16	Testing Oracle	70
17	Make Evidence Compatible (without Helpful Teacher) - $\mathbb{R}$	73
18	Make Evidence Compatible (without Helpful Teacher) - $\mathbb{B}^n$	77

# Notation

General notation	
N	the set of <i>natural</i> numbers
Q	the set of <i>rational</i> numbers
$\mathbb{R}$	the set of <i>real</i> numbers
$\mathbb{B}^n$	the set $\{0,1\}^n$ , where $n \in \mathbb{N}$
A	the cardinality of a finite set $A$
$\max(A)$	the <i>maximum</i> value in a set A
$\min(A)$	the <i>minimum</i> value in a set A
$\arg\max_x\{f(x)\}$	the value $x$ which maximizes the function $f$
$A\oplus B$	the symmetric difference of two sets
$A \uplus B$	the union of pairwise disjoint sets
Notation used in t	his thesis
$\Sigma$	concrete alphabet (finite or infinite)
$\Sigma$	symbolic alphabet, a finite set of symbols
$\boldsymbol{a}, \boldsymbol{b}, \dots$	symbolic letters (or symbols)
$\llbracket a \rrbracket$	semantics of symbol <i>a</i>
$\mathcal{V}(oldsymbol{a})$	set of symbols with adjacent semantics to $[\![a]\!]$
Abbreviations	
DFA	Deterministic Finite Automaton
BDT	Binary Decision Tree
DNF	Disjunctive Normal Form
PBNF	Pseudo-Boolean Normal Form
MQ	Membership Query
EQ	Equivalence Query

# **Introduction**

#### 1.1 Synopsis

In its most general and abstract form, the problem of machine learning can be phrased as follows. Given a function  $f : X \to Y$  and a sample S consisting of (x, f(x)) pairs, construct a function f' that agrees, exactly or approximately, with f on the sample. In many cases, where X is a large domain and Y is a small one, the function f is viewed as a classifier, i.e., a rule that classifies objects in X into one of the finitely many categories in Y. In other words, X is partitioned into finitely many subsets that share some property. In the case where  $Y = \mathbb{B} = \{0, 1\}$ , f is the characteristic function of a set  $L \subseteq X$ , and the sample consists of positive (f(x) = 1) and negative (f(x) = 0) examples.

Formal languages are subsets of  $\Sigma^*$ , that is, the (infinite) set of all finite sequences (words) over an alphabet  $\Sigma$ . The major topic of this thesis concerns the problem of language identification, an instance of machine learning where one wants to infer the characteristic function of a target language  $L \subseteq \Sigma^*$ . More precisely, we work within the framework of active learning, that is, the learning algorithm can select, in an adaptive manner, the words for which membership in L is queried. Our starting point is the  $L^*$  algorithm of Angluin [Ang87], a well-known algorithm for learning regular languages based on membership queries, which also makes use of counter-examples provided by a helpful teacher.

The essence of automaton learning algorithms, such as  $L^*$ , is to organize the information on word membership in a two-dimensional observation table whose rows indicate words (prefixes) that lead to states in the minimal automaton for L and its columns represent suffixes whose acceptance status after the prefix characterizes the states based on the Nerode right-congruence relation [Ner58]. Depending on the learning protocol, the number of queries needed to learn a language L is polynomial or exponential in the number of states of the minimal automaton accepting L. It is linear in the size of the alphabet  $\Sigma$  [BR05].

We extend this algorithm to work in situations where the size of the alphabet is prohibitively large or even infinite, as in the case of  $\mathbb{B}^n$  or  $\mathbb{R}$ , and the accepting automaton cannot be represented traditionally by writing down all the letters that label its transitions. We represent languages on such large alphabets using symbolic automata with a modest number of transitions, each labeled by a predicate on the input alphabet. In each state, these predicates define a partition of the alphabet resulting in a deterministic automaton.

We adapt the automaton learning framework to this setting in a generic way, which does not make strong assumptions on the type of the large alphabet. The idea is to associate with each state a symbolic alphabet whose letters label the transitions outgoing from that state. Symbolic sequences serve as access sequences to states (rows in the observation table) in the same manner as concrete sequences do for the  $L^*$  algorithm. The concrete semantics of the symbolic letters are defined via predicates that map them to sets of concrete letters.

The crucial difference from the small alphabet setting is that, rather than the entire alphabet, only a small sample (evidence) of the presumed semantics of each symbolic letter is tested with membership queries. As a consequence, this semantics is not guaranteed to be correct and it might be that the partition boundaries are imprecise, or even that some symbolic letter (and its corresponding transition) has not been discovered yet.

To explain this with a concrete example, let the alphabet be [0, 10) and suppose that at some point in time we associate two symbolic letters  $\boldsymbol{a}$  and  $\boldsymbol{b}$  with a state  $q_1$ . Consider their respective semantics to be  $[\![\boldsymbol{a}]\!] = [0, 5)$  and  $[\![\boldsymbol{b}]\!] = [5, 10)$  with transitions  $\delta(q_1, \boldsymbol{a}) = q_2$  and  $\delta(q_1, \boldsymbol{b}) = q_3$ , respectively. At some later point in time, we may discover that some concrete letters in [4, 6) behave differently. We introduce then a new symbol c, modify the semantics to  $[\![\boldsymbol{a}]\!] = [0, 4)$ ,  $[\![\boldsymbol{c}]\!] = [4, 6)$ , and  $[\![\boldsymbol{b}]\!] = [6, 10)$ , and add the transition  $\delta(q_1, \boldsymbol{c}) = q_1$ . As we will see, the process of learning over large alphabets interleaves two activities: the discovery of new states, as in standard automaton learning, and the modification of the partition boundaries in each state. The latter part is done for each state using sampling and static machine learning techniques, depending on the nature of the alphabet.

Algorithm  $L^*$ , in addition to membership queries, has a helpful teacher that, whenever the learner conjectures an automaton, either confirms it or presents a counter-example. Such a counter-example simplifies the discovery of new states. One version of our algorithm uses such a helpful teacher and, moreover, assumes the counter-example to be minimal in the lexicographic order (or partial order) induced by the domain. This assumption simplifies the modification of the partition boundaries due to counter-examples, allowing exact learning of the partitions. In another version of the algorithm, we replace this teacher by a more realistic equivalence checking technique that uses random sampling and membership queries. This may or may not lead to counter-examples. We develop a method for analyzing these (non-minimal) counter-examples and see whether they lead to a discovery of a new state or to a modification in the alphabet partition for an existing state. In this setting, a certain convergence to the exact target language cannot be guaranteed and for this, it is replaced by the probably approximately correct (PAC) condition.

We consider, then, the following types of alphabets and partitions/predicates.

- Most of the work is done for one-dimensional numerical domains, such as bounded sub-intervals of ℝ or ℕ, partitioned into finitely many intervals. To this domain we apply algorithms with and without a helpful teacher that provides counter-examples to conjectures;
- For the domain B<sup>n</sup> of Boolean vectors, we consider partitions expressed using decision trees and apply a teacher-free algorithm;
- 3. For multi-dimensional domains, such as bounded subsets of  $\mathbb{R}^n$  or  $\mathbb{N}^n$ , we use monotone partitions that generalize intervals to partially-ordered sets and apply an algorithm that uses a helpful teacher.

#### 1.2 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 nice, clean and comprehensive theory where different formalisms such as automata, logic, regular expressions, semigroups, and grammars are shown to be equivalent.

One weakness, however, of the classical theory of regular languages and automata is that it mainly deals with alphabets that are rather "thin" and "flat", that is, sets that are small and do not admit any additional structure. In many applications, however, alphabets are large and structured. In hardware verification, for example, behaviors are sequences of states and inputs ranging over valuations of Boolean state variables that give rise to exponentially large alphabets, treated symbolically using BDDs and other logical formalisms. As another motivation, one can consider the verification of continues or hybrid systems against specifications written in a formalism such as signal temporal logic (STL) [MN04, MNP08]. Automata over numerical alphabets, admitting an order or partially-order relation, can define the semantics of such requirements.

Symbolic automata [VBDM10], like the ones we use here, 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  $\mathbb{R}^n$ . Other applications that benefit from the succinct representation of symbolic automata and the usage of large alphabets include, among other, natural language processing [VNG01], regular expressions in the context of static and dynamic program analysis [VDHT10], constraint solving and symbolic analysis with SMT solvers [VBDM10] or security analysis of string sanitizers [Vea13].

Finally, learning automata from examples provides a crisp characterization of what a *state* in a dynamical system is in terms of observable input-output behavior. In contrast to classical machine learning, language inference aims to build an automaton model that will be used later for further analysis rather than simple characterization and prediction of some unobserved behavior of a future input. State system identification and language inference admit a variety of applications ranging from linguistics and biology to computer science. It is also widely used in software analysis and verification. Within software analysis, it is applied more specifically, in program analysis [WBAH08], software testing [AJU10, MN10], security testing [CSS<sup>+</sup>10], dynamic testing [RMSM09], integration testing [Sha08], and black-box components testing [Nie03, SL07, Sha08].

#### 1.3 Outline

The rest of the thesis is organized as follows.

#### Preliminaries.

Chapter 2 provides the basic definitions and notations for the theory of regular languages and automata. Moreover, it presents different types of domain partitioning as well as static machine learning algorithms that can be used to learn such partitions from given sample points.

Chapter 3 provides a summary of learning algorithms over small alphabets. The state of the art is followed by a presentation of the  $L^*$  algorithm and a discussion on the treatment of counter-examples.

In Chapter 4 one can find a full presentation of symbolic automata in general, accompanied with special cases of domain-specific symbolic automata.

#### **Contributions.**

Chapter 5 first provides the basic elements that are used in the symbolic learning algorithm and then presents the algorithm without making any assumptions neither on the nature of the teacher or on the alphabet. Moreover in this chapter, one can find a comparison to related work. Chapter 6 provides an adaptation of the more general symbolic algorithm presented in Chapter 5 to the case of a helpful teacher that always provides a counterexample when there is one, and in addition, this is a minimal one in the case of ordered alphabets. We study the cases of numerical alphabets such as  $\mathbb{N}$ ,  $\mathbb{R}$  and  $\mathbb{R}^n$ [MM14, MM15].

Chapter 7 describes the symbolic algorithm in the context of a teacher who can only answer membership queries and thus counter-examples when provided, are not necessarily minimal. We focus on the cases of one-dimensional numerical domains and Boolean alphabets, where the partitions are inferred using static classification learning algorithms [MM17].

Chapter 8 provides some theoretical evaluation of the algorithm in terms of termination and the trade-off between accuracy and complexity.

Chapter 9 presents case studies that we have tested to validate the algorithms, evaluate their performance and compare them, when possible, with concrete learning algorithms.

#### **Conclusion.**

Chapter 10 provides the conclusion and future directions.

# Preliminaries

In this chapter, we provide the basic definitions and notations that are used throughout this thesis. Section 2.1 succinctly provides the definitions and notations concerning regular languages and finite state automata; for a full reference on the topic, the reader may consult the books [LP97, HMU06, Sip06]. In Section 2.2, we explain what a domain partition into classes is, and present some specific cases that are used in the subsequent chapters. Finally in Section 2.3, we discuss some learning methods that can be applied to identify an unknown partition.

#### 2.1 Regular Languages and Automata

An *alphabet*  $\Sigma$  is a non-empty set whose elements are called *letters*. A *word* (or *string*) is a finite sequence of letters, i.e.,  $w = a_1 \dots a_n$ , chosen from the alphabet  $\Sigma$ ; the number of letters in this sequence determines the *length* of the word, i.e., |w| = n. The sequence of zero length is called the *empty word* and it is denoted by  $\varepsilon$ . The set of all words over an alphabet  $\Sigma$  is denoted by  $\Sigma^*$ .

Let  $w = a_1 \dots a_n$  and  $u = b_1 \dots b_m$  be two words. The two words w and u are said to be *equal* when they have equal length and the same letter at each position, i.e., w = u, if and only if n = m and  $a_i = b_i$  for all  $i = 1, \dots, n$ . The *concatenation* of two words w and u is the sequence  $w \cdot u = a_1 \dots a_n b_1 \dots b_m$ . A word u is a *prefix* (*resp. suffix*) of w, if there exists  $v \in \Sigma^*$  such that  $w = u \cdot v$  (*resp.*  $w = v \cdot u$ ).

Let  $\leq$  be an order relation over  $\Sigma$ . This can be naturally lifted to an order relation over  $\Sigma^*$ . Some of the most interesting order relations defined on the set of words are the following.

- 1. The *prefix order*,  $(\Sigma^*, \leq_{pr})$ , where  $u \leq_{pr} w$  if and only if u is a prefix of w.
- 2. The *lexicographic order*,  $(\Sigma^*, \leq_{lex})$ , where  $u \leq_{lex} w$  if and only if either

 $u \leq_{pr} w$ , or  $u = v \cdot a \cdot v'$  and  $w = v \cdot b \cdot v''$  for some  $v, v', v'' \in \Sigma^*$ ,  $a, b \in \Sigma$ , where  $a \leq b$  and  $a \neq b$ .

3. The *length-lexicographic order*,  $(\Sigma^*, \leq_{ll})$ , where  $u \leq_{ll} w$  if and only if either |u| < |w|, or |u| = |w| and  $u \leq_{lex} w$ .

Note that, if  $(\Sigma, \leq)$  defines a total order, then the orders  $(\Sigma^*, \leq_{lex})$  and  $(\Sigma^*, \leq_{ll})$  are also total.

Any subset of words L chosen from  $\Sigma^*$  is called a *language* over  $\Sigma$ . To each language L we can associate a *characteristic function*  $f : \Sigma^* \to \{+, -\}^1$  such that f(w) = + if the word w belongs to L and f(w) = -, otherwise. Moreover, with every  $s \in \Sigma^*$  we can associate a *residual characteristic function* defined as  $f_s(w) = f(s \cdot w)$ .

Several operations can be defined on languages. The *Boolean operations* include the *union*, *intersection*, *complementation*, and *difference*. The *symmetric difference* is also commonly used. Formally, let L and L' be two languages over  $\Sigma$ , these operations are defined, respectively, as

$$L \cup L' = \{ w \in \Sigma^* : w \in L \text{ or } w \in L' \}$$
  

$$L \cap L' = \{ w \in \Sigma^* : w \in L \text{ and } w \in L' \}$$
  

$$\bar{L} = \{ w \in \Sigma^* : w \notin L \}$$
  

$$L \setminus L' = \{ w \in \Sigma^* : w \in L \text{ and } w \notin L' \}$$
  

$$L \oplus L' = (L \setminus L') \cup (L' \setminus L)$$

The *concatenation* (or *product*) of two languages L and L' is the set of words that can be formed by taking any word in L and concatenating it with any word in L', that is,

$$L \cdot L' = \{ w \cdot w' : w \in L \text{ and } w' \in L' \}.$$

The *power* of a language L is defined as  $L^{k+1} = L^k \cdot L = L \cdot L^k$  for all  $k \ge 1$  where  $L^0 = \{\varepsilon\}$ . The *star* operation is the union of all the powers of a language L and it is defined as

$$L^* = \bigcup_{k \in \mathbb{N}} L^k.$$

The left quotient (resp. right quotient) of L is the language

$$u^{-1}L = \{v \in \Sigma^* : uv \in L\}$$
 (resp.  $Lu^{-1} = \{v \in \Sigma^* : vu \in L\}$ ).

A language  $L \subseteq \Sigma^*$  is *prefix-closed* (*resp. suffix-closed*) if and only if for all  $w, u \in \Sigma^*, w \cdot u \in L$  implies  $w \in L$  (*resp.*  $u \in L$ ).

<sup>&</sup>lt;sup>1</sup>The notation  $\{1, 0\}$  or  $\{True, False\}$  is also commonly used.

A language  $L \subseteq \Sigma^*$  is a *regular language* if it can be recursively constructed by applying the union, concatenation and star operations finitely many times. By default, the empty language, the language  $\{\varepsilon\}$ , and the one letter languages  $\{a\}$  for all  $a \in \Sigma$ , are regular languages. Regular languages over finite alphabets can be represented by deterministic finite automata.

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

The transition function  $\delta$  can be extended to  $\delta: Q \times \Sigma^* \to Q$ , where  $\delta(q, \varepsilon) = q$ 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* (or *accepted*) by  $\mathcal{A}$ , denoted by  $L(\mathcal{A})$ , consists of the set of all accepted words.

A language  $L \subseteq \Sigma^*$  is *recognizable* if there exists an automaton  $\mathcal{A}$  such that  $L = L(\mathcal{A})$ . Let  $REC(\Sigma)$  denote the set of all recognizable languages over the alphabet  $\Sigma$ ; this set coincides with the set of regular languages.

Two automata that recognize the same language are said to be *equivalent*. Among all automata that are equivalent there exists a *minimal* automaton, also known as the *canonical* automaton, which is minimal in the number of states and unique up to isomorphism.

Let  $L \subseteq \Sigma^*$  be a language. This induces an equivalence relation  $\sim_L$  on  $\Sigma^*$ , the *Nerode equivalence relation*, where

$$s \sim_L r$$
 if and only if  $f_s = f_r$  (2.1)

Intuitively, two words are Nerode equivalent whenever by appending the same string to s and r, the resulting two strings are either both in L or both not in L. The relation  $\sim_L$  is a right congruence satisfying  $s \sim_L r \rightarrow s \cdot a \sim_L r \cdot a$  and is known as the *syntactic congruence* associated with L. Its equivalence classes correspond to the states of the minimal automaton that accepts L. The following theorem is part of what is known as the Myhill-Nerode Theorem

**Theorem 2.1** (Nerode [Ner58]). Let  $L \subseteq \Sigma^*$  be a language, and let  $\sim_L$  be an equivalence relation as defined in (2.1). Then, L is a regular language if and only if  $\sim_L$  is of finite index.

As a result, for any regular language  $L \subseteq \Sigma^*$ , one can find an automaton that is minimal using the syntactic congruence  $\sim_L$ . Based on Theorem 2.1, this automaton can be defined as  $\mathcal{A}_M = (\Sigma, Q, q_0, \delta, F)$ , where the set of states Q is the set of equivalence classes, i.e,  $Q = \Sigma^* / \sim_L$ ; the initial state is the class that contains the empty word, i.e.,  $q_0 = [\varepsilon]$ ; the transition function is defined as  $\delta([u], a) = [u \cdot a]$  for all  $a \in \Sigma$ ; and, the final states are all those classes that contain accepted words,  $F = \{[u] : u \cdot \varepsilon \in L\}.$ 

According to this, for an automaton to be considered minimal, all pairs of states in it should be distinguishable and thus not be equivalent. Two states p, q are distinguishable if there exists a word w such that  $\delta(p, w) \in F$  and  $\delta(q, w) \notin F$  (or  $\delta(p, w) \notin F$  and  $\delta(q, w) \in F$ ).

As the minimal representation is useful in many areas, several minimization algorithms have been proposed. One can refer to [Moo56, Brz62, Hop71] for the most well known automaton minimization algorithms. The Nerode congruence and the identification of minimal automata play a crucial role in most automata learning algorithms since [Gol72]. We discuss further this topic in Chapter 3.

#### 2.2 Sets and Partitions

Let  $\Sigma$  be a finite or an infinite alphabet. A representation of  $\Sigma$  as the finite union of a family of pairwise disjoint subsets is called a *partition of*  $\Sigma$  *into classes*. This representation plays an important role in many problems [KF70]. Such partition is usually based on certain criteria that allow the assignment of elements from  $\Sigma$  to one class or another.

Let  $C = \{a_1, \dots, a_k\}$  be a finite set of *classes* (or *labels*). A partition of  $\Sigma$  into classes can be defined as a mapping  $\psi : \Sigma \to C$  that associates each letter in  $\Sigma$  with a class in C. Such a mapping is known as a *classification function*, and the tuple  $(\Sigma, \psi)$  is often called a *labeled alphabet*. The  $\Sigma$ -semantics of a label  $a \in C$  denotes the set of elements from  $\Sigma$  that are associated to the same class a, i.e.,  $[\![a]\!] = \{a \in \Sigma : \psi(a) = a\}$ .

In the following, we provide some examples of such partitions. For readability, the alphabet  $\Sigma$  is presented in these examples as a subset of  $\mathbb{N}$ ,  $\mathbb{R}$ ,  $\mathbb{R}^2$  or  $\mathbb{B}^4$ , but one should consider that  $\Sigma$  can be any subset of  $\mathbb{N}^n$ ,  $\mathbb{R}^n$  and  $\mathbb{B}^n$ , where  $n \in \mathbb{N}$ .

#### Intervals

Let  $\Sigma$  be a subset of the reals that admits the usual order relation  $\leq$ . We consider as an *interval* any connected segment<sup>2</sup> of  $\mathbb{R}$ . Let  $a, b \in \mathbb{R}$  be the endpoints of an interval, we can determine the following interval types:

- open interval:  $(a, b) = \{x \in \mathbb{R} : a < x < b\},\$
- closed interval:  $[a, b] = \{x \in \mathbb{R} : a \le x \le b\},\$
- *closed-open* interval:  $[a, b) = \{x \in \mathbb{R} : a \le x < b\}$ , and

<sup>&</sup>lt;sup>2</sup>A connected segment of  $\mathbb{R}$  is a non-empty subset X of  $\mathbb{R}$  such that for all  $a, b \in X$  and  $c \in \mathbb{R}$ ,  $a \leq c \leq b$  implies  $c \in X$ .



Figure 2.1: Partition of  $\Sigma = [0, 10] \subseteq \mathbb{R}$  ( or  $\mathbb{N}$ ) into three intervals.

- open-closed interval:  $(a, b] = \{x \in \mathbb{R} : a < x \le b\}.$ 

An interval is empty when b < a (or  $b \le a$  when the interval is not closed). By abuse of notation, we use the interval notation also for other numerical domains such as  $\mathbb{N}$ , e.g.,  $[a, b) = \{x \in \mathbb{N} : a \le x < b\}$ . An example of a partition of  $\Sigma =$ [0, 10] into intervals is shown in Figure 2.1. In this example,  $\Sigma$  is partitioned into three intervals, each designated to a class from  $\mathcal{C} = \{a, b, c\}$ . The classification function is given by

$$\psi(x) = \begin{cases} a, & \text{for } x \in [0, 4) \\ b, & \text{for } x \in [4, 6) \\ c, & \text{for } x \in [6, 10] \end{cases}$$

For a given partition of  $\Sigma$ , distinct classes with adjacent semantics can be considered as neighbors. Formally, a class  $b \in C$  is a *right (resp. left) adjacent* to class  $a \in C$  ( $a \neq b$ ), if  $a \in [\![a] \!] \land c \in [\![b] \!]$  (*resp.*  $a \in [\![b] \!] \land c \in [\![a] \!]$ ) implies  $b \in [\![a] \!] \cup [\![b] \!]$ , for any three letters  $a, b, c \in \Sigma$  such that a < b < c. The set of all right and left adjacent classes to  $a \in \Sigma$  forms the *neighborhood* of a, denoted by  $\mathcal{V}(a)$ . Intuitively, the neighborhood of a consists of all those classes  $b \in C \setminus \{a\}$  with which a shares the same boundaries in their semantics. Referring to the example of Figure 2.1, class b is a right adjacent to a and a left adjacent to c, but a and c are not adjacent. The neighborhoods of a and b are  $\mathcal{V}(a) = \{b\}$  and  $\mathcal{V}(b) = \{a, c\}$ , respectively.

#### **Monotone Partitions**

Let  $\Sigma$  be a partially-ordered alphabet of the form  $\Sigma = X^n$ , where X is a totallyordered set such as an interval  $[0, k) \subseteq \mathbb{R}$ . Letters of  $\Sigma$  are *n*-tuples of the form  $\mathbf{x} = (x_1, \ldots, x_n)$  with minimal element  $\mathbf{0} = (0, \ldots, 0)$ . 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  $i = 1, \ldots, n$ . When  $\mathbf{x} \leq \mathbf{y}$  and  $x_i < y_i$  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  $\mathbf{x} ||\mathbf{y}$ , if  $x_i < y_i$  and  $x_j > y_j$  for some *i* and *j*.

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



Figure 2.2: (a) Backward and forward cone for x, (b) union of cones, (c) an alphabet partition into three classes  $C = \{a_1, a_2, a_3\}$ .

no three points such that  $\mathbf{x} < \mathbf{y} < \mathbf{z}$ ,  $\mathbf{x}, \mathbf{z} \in P$ , and  $\mathbf{y} \notin P$ . We will see in the following that a monotone partition with a finite number of partition blocks can be represented by a finite set of points.

A forward cone  $B^+(\mathbf{x}) \subseteq \Sigma$  is the set of all points dominated by a point  $\mathbf{x} \in \Sigma$ , see Figure 2.2a. For a set of points  $F = {\mathbf{x}_1, \dots, \mathbf{x}_l}$ , this cone can be extended to  $B^+(F) = B^+(\mathbf{x}_1) \cup \dots \cup B^+(\mathbf{x}_l)$ , as shown in Figure 2.2b. When  $\mathbf{x}_1, \dots, \mathbf{x}_l$ are mutually incomparable, they constitute the set of minimal elements of F.

From a family of sets of points  $\mathcal{F} = \{F_0, \ldots, F_{m-1}\}$ , which satisfies  $F_0 = \{\mathbf{0}\}$  and 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)  $\forall \mathbf{y} \in F_i$ ,  $\forall \mathbf{x} \in F_{i-1}$ ,  $\mathbf{y} \not\leq \mathbf{x}$ ; one can define a monotone partition of the form  $\mathcal{P} = \{P_1, \ldots, P_{m-1}\}$ , where  $P_i = B^+(F_{i-1}) - B^+(F_i)$ .

An example of such a partition appears in Figure 2.2c. The partition is defined by the family of sets of incomparable points  $\mathcal{F} = \{F_0, F_1, F_2\}$ , where  $F_0 = \{\mathbf{0}\}$ ,  $F_1 = \{x_1, \ldots, x_4\}$ , and  $F_2 = \{y_1, \ldots, y_4\}$ . The classification function can be given by  $\psi(\mathbf{x}) = \mathbf{a}_i$  for all  $\mathbf{x} \in P_i$ .

The extension of the notions of adjacent classes and neighbors to the case of partially-ordered sets is straightforward. Hence in the example, the label  $a_1$  is a left adjacent to  $a_2$ , but is not an adjacent to  $a_3$ , and for the labels  $a_1$  and  $a_2$  the neighborhoods are  $\mathcal{V}(a_1) = \{a_2\}$  and  $\mathcal{V}(a_2) = \{a_1, a_3\}$ , respectively.

#### **Partitions of the Boolean Cube**

Let  $\Sigma = \mathbb{B}^n$  be the *n*-dimension Boolean hyper-cube, where letters are boolean vectors and can be accessed by the Boolean variables  $X = \{x_1, \ldots, x_n\}$ . A literal is either a Boolean variable x, or its negation  $\bar{x}$ . A sub-cube (or term) of  $\mathbb{B}^n$  is a conjunction of literals. We denote a sub-cube by  $\phi_k = \prod_{i \in A_k} x_i \cdot \prod_{j \in B_k} \bar{x}_j$  where  $A_k, B_k$  are disjoint sets of indices from  $\{1, \ldots, n\}$ .

A possible way to partition a Boolean cube into classes is to partition it into sub-cubes, not necessarily of equal size, and associate each sub-cube with a class. A classification function for the Boolean cube is thus a pseudo-boolean function<sup>3</sup>  $\psi : \mathbb{B}^n \to \mathcal{C}$ . Such a function can be represented in several ways; we either use a pseudo-boolean normal form or a binary decision tree to represent it, but in the case where the number of variables n is small, we often use a Karnaugh map representation as it offers a better visualization of the partition.

A *pseudo-boolean normal form* (PBNF), an analogous of a disjunctive normal form (DNF) for Boolean functions, is defined as:

$$\psi = \sum_{k=1}^{m} a_k (\prod_{i \in A_k} x_i) (\prod_{j \in B_k} \bar{x}_j) = \sum_{k=1}^{m} a_k \cdot \phi_k$$
(2.2)

where the coefficients  $a_1, \ldots, a_m$  are labels from C and  $A_k, B_k$  are sets of indices such that i)  $A_k, B_k \subseteq \{1, \ldots, n\}$ , ii)  $A_k \cap B_k = \emptyset$ , iii)  $A_k \cup B_k \neq \emptyset$ , and iv)  $(A_k \cap B_l) \cup (A_l \cap B_k) \neq \emptyset$  for  $k, l = 1, 2, \ldots, m, k \neq l$ .

Each one of the four conditions is required for the function  $\psi$  to be well defined and to form a partition of  $\Sigma$ . More precisely, we let  $A_k \cup B_k \neq \emptyset$  to avert redundant terms in  $\psi$ ; the function  $\sum_k (\prod_{i \in A_k} x_i) (\prod_{j \in B_k} \bar{x}_j)$  is a DNF only when  $A_k \cap B_k = \emptyset$ ; and moreover, it is orthogonal when no two terms intersect, i.e.,  $(A_k \cap B_l) \cup$  $(A_l \cap B_k) \neq \emptyset$  for  $k \neq l$ . Orthogonality guarantees that  $\psi$  partitions  $\Sigma$ . The set of all letters  $a \in \Sigma$  that model  $\phi_k$  form the  $\Sigma$ -semantics of the class labeled by  $a_k$ , i.e.,  $\|a_k\| = \|\phi_k\| = \{a \in \Sigma : \phi_k(a)\}.$ 

A binary decision tree (BDT) (or classification tree) takes a letter as an input and, based on its attributes, assigns it to a class. A BDT is a tree T which at each node either tests some variable of the given input, or provides a class label. A BDT can hence be seen as a binary tree<sup>4</sup> T equipped with a mapping  $d : \mathcal{P}os(T) \rightarrow X \uplus \mathcal{C}$ , where  $\mathcal{P}os(T)$  is the set of nodes in the tree, X is the set of attributes, and  $\mathcal{C}$  is the set of classes. Let leaves(T) denote the set of all leaf nodes in T. A mapping d interprets a BDT if and only if  $d(t) \in \mathcal{C}$  when  $t \in leaves(T)$ , and  $d(t) \in X$  otherwise.

To assign a class to a letter  $a \in \Sigma$ , one needs to follow a path in T. This path is a sequence of nodes, which starts at the root of the tree and traverses it until reaching a leaf node. Formally, a path in T as a sequence  $path_T(a) : t_0, \ldots, t_m$ , such that  $t_0 = \varepsilon$ ,  $t_j \in \{t_{j-1}0, t_{j-1}1\}$  for all  $j = 1, \ldots, m$ , and  $t_m \in leaves(T)$ . The path is computed as follows: start at the root; when at inner node  $t_j$ , which is mapped to the variable  $x_i$ , make a test on a; proceed to the left child node  $t_j0$  if  $x_i = 0$ , proceed to the right child node  $t_j1$ , otherwise; repeat until reaching a leaf

<sup>&</sup>lt;sup>3</sup>A *pseudo-boolean function* is the analogous of a Boolean function, with the difference that the range of the function is any set other than  $\mathbb{B}$ .

<sup>&</sup>lt;sup>4</sup> A binary tree T is a prefix closed set of positions  $\mathcal{P}os(T) \subseteq \{0,1\}^*$  with  $\varepsilon$  representing its root, such that (i) for every inner node  $t \in \mathcal{P}os(T)$ , nodes t0 and t1 are the left and right *child nodes*, respectively, and (ii) a node t is a leaf node when  $t0 \notin \mathcal{P}os(T)$  and  $t1 \notin \mathcal{P}os(T)$ .



Figure 2.3: The (a) hypercube of four dimensions  $\mathbb{B}^4$ , and (b) a partition of it into three sub-cubes  $\{\bar{x}_3, \bar{x}_1x_3, x_1x_3\}$ , each one assigned to one class from  $\mathcal{C} = \{a_1, a_2, a_3\}$ . A representation of the classification function can be given as (c) a Karnaugh map, (d) a BDT, or (e) a pseudo-boolean function.

node. For each letter  $a \in \Sigma$ , there exists a unique path in T. We use T(a) to denote the label of the leaf node reached by  $path_T(a)$ .

Each node  $t \in \mathcal{P}os(T)$  is naturally associated with a sub-cube  $\llbracket t \rrbracket$  of  $\Sigma$ , which consists of all letters  $a \in \Sigma$  whose path in T contains the node t. Consequently, for every BDT T there exists an equivalent PBNF  $\psi$  that represents the same classification function. To find this PBNF, we work as follows. Let  $\phi_t$  denote the cube associated to node t. This can be defined recursively as  $\phi_{\varepsilon} = 1$  and  $\phi_{t \cdot a} = \phi_t \cdot x_i^{\kappa}$ , where  $x_i$  is the variable carried by node t and  $x^{\kappa}$  denotes x and  $\bar{x}$  when  $\kappa = 1$  and  $\kappa = 0$ , respectively. Let T a BDT, we define its equivalent PBNF as

$$\psi = \sum_{t \in leaves(T)} d(t) \cdot \phi_t$$

where d(t) is the label associated to the leaf node t. To emphasize the relation between a PBNF  $\psi$  and its equivalent BDT we denote the latter by  $T_{\psi}$ .

A *Karnaugh map* is a classical representation of Boolean functions that offers a better visualization of functions involving a small number of variables (up to 5 or 6). A Karnaugh map is given in a matrix structure, where rows and columns are indexed by the values of some variables, and each cell contains the value of the function in the corresponding Boolean vector. For instance, in Figure 2.3c, rows and columns are indexed by the values of the pairs of variables  $(x_1, x_2)$  and  $(x_3, x_4)$ , respectively, and the cell in the second row and third column of the map contains the value (label)  $a_2$  since  $\psi(0111) = a_2$ .

**Example 2.2.** Let the alphabet be the Boolean cube  $\Sigma = \mathbb{B}^4$ , see Figure 2.3a. A possible partition of this cube into three classes  $\mathcal{C} = \{a_1, a_2, a_3\}$  is shown in Figure 2.3b. This partition can be equivalently represented by either the BDT in Figure 2.3d, or the Karnaugh map in Figure 2.3c, or the PBNF  $\psi = a_1 \cdot \bar{x_3} + a_2 \cdot \bar{x_1} \cdot x_3 + a_3 \cdot x_1 \cdot x_3$ . Using the classification tree representation, the letter a = (0, 0, 1, 0) follows the path  $path_T(a) = \varepsilon$ , 1, 10 in the tree, and thus it is assigned to the class  $a_2$ , i.e.,  $\psi(a) = a_2$ .

#### 2.3 Learning Partitions

In the previous section, we have defined classification functions and have shown how a domain, an alphabet  $\Sigma$ , can be partitioned into classes. Moreover, we have presented various classification functions  $\psi : \Sigma \to C$ , where C denotes a set of classes, typically finite and small.

However, a typical problem is to find such a representation, a classification function, when only a small set of classified elements from  $\Sigma$  is given. This problem is commonly encountered in machine learning and belongs in the context of *predictive* or *supervised learning*. To formalize the problem, let S = $\{(a^1, y^1), \ldots, (a^r, y^r)\} \subseteq \Sigma \times C$  be a *labeled sample*. We assume the existence of an underlying classification function  $\psi$ , which is unknown outside S, such that  $y^i = \psi(a^i)$  for all  $i = 1, \ldots, r$ . Our aim is to find a function  $\psi'$  which approximates well  $\psi$  on the labeled data from S, and, moreover, it generalizes well to make predictions on novel unseen inputs from  $\Sigma$ .

Depending on the set of labels C, the problem is known as *classification* or *pattern recognition* when C is finite, or as *regression* when C is the set of reals. In this thesis, we restrict ourselves to a finite number of classes, and in this section, we describe classification learning techniques that are commonly used to learn the partitions on Boolean cubes presented in Section 2.2.

#### 2.3.1 Learning Binary Decision Trees

Let S be a labeled sample from an unknown classification function  $\psi$ . In this section, we discuss the problem of learning a BDT to represent  $\psi$  that is compatible with the given sample S. As many such trees can be found, we prefer an optimal tree. Optimality usually refers to the tree size, such as its depth; the node purity, like the number of mismatching points in the sample; or other criteria such as

Ockham's razor. Since finding a globally optimal tree is NP-hard [BFSO84], a greedy splitting algorithm is commonly used to build the tree.

A greedy splitting algorithm works roughly as follows. It starts with a tree that consists of a single root node; all sample points are associated with this node. A node is said to be *pure* if all its sample points have the same label. For each impure node, two descendants are created and the sample is split among them based on the value of some selected variable  $x_i$ . The variable is chosen according to some purity measure, such as information gain, that characterizes the quality of the split based on each variable. The selection is greedy and does not depend on the optimality of following splits. The algorithm terminates when the tree becomes sample compatible and sends each sample point to a pure leaf node.

In the following, we describe more formally this learning algorithm. First, we present the case where the labeled sample is assumed to be given in advance. This algorithmic scheme is used by the most popular implementations for learning BDTs such as CART [BFSO84], ID3 [Qui86] and C4.5 [Qui93]. Then, we discuss the problem of learning a BDT when the sample is not completely provided before hand, but is updated during the learning process by adding new sample points. Popular incremental algorithms of this kind are, among others, the ID4 [SF86], ID5 [E88], and ID5R [Utg89] algorithms.

#### Using a Fixed Labeled Sample

Let  $S = \{(a^1, y^1), \ldots, (a^r, y^r)\}$  be a labeled sample, where  $y^i = \psi(a^i)$  and  $\psi : \Sigma \to C$  is an unknown classification function. We want to learn a BDT T such that  $T \simeq T_{\psi}$ . Each node  $t \in \mathcal{P}os(T)$  associates to a set  $S_{|t} = \{(x, y) \in S : t \in path(x)\}$  that consists of all sample points reaching node t. Note that  $S_{|\varepsilon} = S$  for any labeled sample  $S \subseteq \Sigma \times C$ . A node  $t \in leaves(T)$  is said to be *consistent* with S whenever T(x) = y for all  $(x, y) \in S_{|t}$ . A BDT is consistent with a labeled sample S if all its leaf nodes are consistent with S. A node  $t \in \mathcal{P}os(T)$  is pure when all points in  $S_{|t}$  carry the same label, i.e.,  $|\{y : (x, y) \in S_{|t}\}| = 1$ .

Algorithm 1 Grow Tree

1: **initialize** (T, d), such that  $\mathcal{P}os(T) = \{\varepsilon\}$  and  $d(\varepsilon) = label(S)$ 2: **while**  $\exists t \in leaves(T)$  such that t is not pure **do** 3: Let  $x_i = best\_attribute(S_{|t}, X)$ 4:  $\mathcal{P}os(T) = \mathcal{P}os(T) \cup \{t0, t1\}$ 5:  $d(t) = x_i; d(t0) = label(S_{|t0}); d(t1) = label(S_{|t1})$ 6: **end while** 

The greedy splitting algorithm is shown in Algorithm 1. The algorithm uses a quality measure to do splitting; the most common quality measures are presented

in detail below. We denote by label(S) the most common label encountered in a sample S, that is,  $label(S) = \arg \max_{c \in C} |\{(x, y) \in S : y = c\}|$ ; we let label(S) be equal to a random label from C when S is empty. The  $best\_attribute(S, X)$  denotes the Boolean variable that performs the best split of a sample S into two subsets. The algorithm performs as follows. After being initialized to a one node tree, the BDT T is extended by splitting its leaf nodes until it becomes consistent with the sample. Whenever a leaf node t in T is not pure it becomes an inner node and two child leaf nodes t0 and t1 are added to  $\mathcal{P}os(T)$  (line 4). The inner node t updates and maps to the attribute  $x_i \in X$  that performs the best split according to the chosen quality measure (line 3); note that this variable does not appear so far in the path reaching the node t. The two new leaf nodes map to the most common label encountered in the sub-sample associated to them (line 5).

#### **Splitting Quality Measures**

To measure the spitting quality of an attribute and decide whether a split is better than another, we need to compare the gain achieved when splitting the sample *S* using each attribute. We compare on the gain on the node impurity and for this we use impurity measures. Other measures, as for instance on the size of the tree, are not considered here. As a further reference, one can consult [BFSO84, Mur12].

The set  $\Delta^k \subseteq \mathbb{R}^k$ , known as the standard k-simplex, is the set of all real tuples  $(t_1, \dots, t_k) \in \Delta^k$  such that  $\sum_{i=1}^k t_i = 1$  and  $t_i \ge 0$  for all i. A function  $\mathcal{I} : \Delta^k \to \mathbb{R}$  is an *impurity measure* if i) it is symmetric with regard to its arguments, ii) it has minima at the vertices of  $\Delta^k$ , and iii) it has one maximum at  $(\frac{1}{k}, \dots, \frac{1}{k})$ . One can measure the *impurity of a labeled sample* S, with labels chosen from  $\mathcal{C} = \{c_1, \dots, c_k\}$ , by letting

$$\mathcal{I}(S) = \mathcal{I}(p_1, \dots, p_k) \text{ and } p_i = \frac{|\{(x, y) \in S : y = c_i\}|}{|S|}.$$

The impurity measures, most widely used in learning BDT, are the *misclassification rate*, the *entropy measure*, and the *Gini index*. These are defined on a tuple  $(p_1, \ldots, p_k)$  as:

Misclassification rate:	$\mathcal{I}_{error}(p_1,\ldots,p_k) = 1 - \max_i p_i$
Entropy measure:	$\mathcal{H}(p_1,\ldots,p_k) = -\sum_{i=1}^k p_i \log_2 p_i$
Gini index:	$\mathcal{I}_G(p_1,\ldots,p_k) = 1 - \sum_{i=1}^k p_i^2$

Now we can measure the impurity gain when S is split using a specific attribute, and see whether a split is better than another. Let  $x_i \in X$  be an attribute that splits S into two sub-samples  $S_{x_i=0}$  and  $S_{x_i=1}$ , each containing the sample points that satisfy  $x_i = 0$  and  $x_i = 1$ , respectively. The impurity gain achieved by splitting



Figure 2.4: The hypercube and Karnaugh map of  $\mathbb{B}^4$  with the sample points indicated by their labels red (r), green (g), and blue (b). On the right is the BDT learned from S using Algorithm 1.

using attribute  $x_i$  is given by

$$Gain_{\mathcal{I}}(S, x_i) = \mathcal{I}(S) - \sum_{j=0,1} \frac{|S_{x_i=j}|}{|S|} \cdot \mathcal{I}(S_{x_i=j}),$$

where  $\mathcal{I}$  is the impurity measure used. Let us demonstrate the learning algorithm with a simple example.

**Example 2.3.** Let the alphabet  $\Sigma$  be the Boolean hypercube  $\mathbb{B}^4$  and let  $\mathcal{C} = \{r, g, b\}$  be the set of labels, which correspond to the color names red (r), green (g), and blue (b), respectively. We use the learning algorithm described above to build a BDT. We let the labeled sample be  $S = \{(0001, r), (0111, b), (1011, g), (1100, r), (1110, g)\}$ . Figure 2.4a shows a 4-dimensional cube whose vertices denote the Boolean vectors and the sample points are denoted by label colors. The Karnaugh map in Figure 2.4b provides a more intuitive representation. As a split quality measure we use information gain (entropy measure). All the calculations needed for this example are shown in Table 2.1.

First, we initialize the BDT to (T, d) such that  $\mathcal{P}os(T) = \{\varepsilon\}$  and  $d(\varepsilon) = label(S) = r$ . The node  $\varepsilon$  is not pure, i.e.,  $\mathcal{H}(S_{|\varepsilon}) > 0$ . For this reason, the root node needs to be split. Comparing all possible splits, we find that the maximum information gain is achieved by splitting over  $x_3$ . Hence, the tree updates to  $\mathcal{P}os(T) = \{\varepsilon, 0, 1\}, d(\varepsilon) = x_3, d(0) = label(S_{|0}) = r$ , and  $d(1) = label(S_{|1}) = g$ . The leaf node t = 0 is pure, i.e., has entropy 0, and does not need to split further. The leaf node t = 1, on the other hand, is not pure. The best split is given for  $x_1$ , and the tree updates to  $\mathcal{P}os(T) = \{\varepsilon, 0, 1, 00, 01\}, d(1) = x_1, d(00) = b$ , and d(01) = g. Now, all leaf nodes are pure and the algorithm terminates. The final tree can be seen in Figure 2.4c.

S	(r,g,b)	$\mathcal{H}(S)$	$Gain_{\mathcal{H}}(S, x_1)$	$Gain_{\mathcal{H}}(S, x_2)$	$Gain_{\mathcal{H}}(S, x_3)$	$Gain_{\mathcal{H}}(S, x_4)$
$S_{ \varepsilon }$	(2, 2, 1)	1.52193	0.570951	0.170951	0.970951	0.170951
$S_{ 0 }$	(2, 0, 0)	0	-	-	-	-
$S_{ 1}$	(0, 2, 1)	0.918296	0.918296	0.251629	0	0.251629
$S_{ 00 }$	(0, 0, 1)	0	-	-	-	-
$S_{ 01}$	(0, 2, 0)	0	-	-	-	-

Table 2.1: Information gain for all samples and sub-samples of Example 2.3

#### **Incremental Learning Algorithm**

In this section, we change the setting of the learning problem and let the labeled sample be a stream of sample points received during the learning process. That is, the tree is learned and made compatible with the current sample when another new sample point arrives. The new observation incorporates to the sample and the tree updates according to the new information that is now available.

Formally, let  $S^i = \{(a^1, y^1), \dots, (a^i, y^i)\}$  denote the sample that consists of i sample points, and let  $(T^i, d^i)$  be a BDT tree compatible to  $S^i$ . When a new observation  $(a^{i+1}, y^{i+1})$  arrives, we want to update  $T^i$  according to this new information. Several algorithms have been proposed in the literature to handle this problem. In the following, we explain some of these algorithms and discuss on their efficiency and complexity results.

One naive algorithm, which also appears in [SF86] for comparison purposes, uses Algorithm 1 as an incremental algorithm. Whenever a new observation or a set of new observations arrives, this is added to the sample. Then, the Algorithm 1 is called and outputs a new tree which is constructed from scratch based on the updated sample.

It is observed that a large part of the tree need not change at all, making the previous method unnecessarily costly. In [SF86], the authors propose an incremental algorithm where only a sub-tree needs to be rebuilt. The existing tree  $T^i$  is revised based on the new sample point  $(a^{i+1}, y^{i+1})$ , which traverses the nodes of the tree starting from the root. The quality measures are updated including the new sample point. If for a node  $t \in \mathcal{P}os(T^i)$  the selected attribute is not the best anymore, then the sub-tree rooted at t is deleted and rebuilt. The sub-tree with root t is reconstructed using Algorithm 1 on the sample  $S^i_{|t|} \cup \{(a^{i+1}, y^{i+1})\}$ .

To avoid preserving the whole set of sample points that have been seen so far, the ID4 algorithm [SF86] uses a series of tables, located at each node, that carry all information needed to calculate the quality measures when revising the tree.

Finally, another algorithm that is proposed is the ID5 algorithm [E88], which is followed by two improved versions, the ID5R [Utg89] and ITI [UBC97] algorithms. These algorithms use an alternative revising mechanism. That is, whenever

a sub-tree of  $T^i$  is found such that the splitting attribute at its root is not the best split, then this is only restructured instead of reconstructed as in ID4. To restructure the sub-tree, the algorithm uses a *pull-up* method in a way that all attributes at the nodes are the same.

A worst-case analysis and comparison of these incremental algorithms can be found in [Utg89]. It is shown that the ID4 algorithm cannot incrementally learn the parity function when the sample is given in a certain order, because it has always to update the whole tree. The cost of calculating the gain function and choosing the best attribute is significantly lower in the case of ID5R.

In this section, we restricted our discussion to algorithms where a) the order in which sample points arrive does not impact the final result, and b) the final BDT, which is constructed incrementally, is the same as if the tree is built using the non-incremental algorithm discussed in Section 2.3.1.

In the literature, one can find many extensions and variations of the algorithms discussed above that handle the same or slightly different settings of the learning BDT problem. For instance, when the sample is noisy or incomplete, most algorithms include a pruning phase after the tree has been fully constructed in order to avoid overfitting of the data [BFSO84, Qui93]. During this phase the stopping criteria are relaxed and some branches of the tree are pruned back whenever they do not add any significant gain to the final tree.

# **Language Identification**

This chapter introduces in more detail the problem of learning formal languages, that is a subset of  $\Sigma^*$ , also known as language identification. We start, in Section 3.1, by connecting language identification to the more general machine learning context. Then, in Section 3.2, we provide the state of the art in the field and explain the notions of identification in the limit and query learning. The presentation of the active learning algorithm  $L^*$  is given in Sections 3.3, and it is followed by a description of some well-known variations. A full reference on the topics presented in this chapter can be found in [Ang87, BR05, DlH10].

#### 3.1 Connection with Machine Learning

Machine learning has gained popularity in the last decade with a variety of applications being associated with the field. Recognizing patterns after having observed a relatively small amount of data is the main task of almost all machine learning techniques.

A general and abstract description of a problem in this context is the following. Let  $f : \Sigma \to C$  be a function, where  $\Sigma$  is a set of objects and C is a set of labels. Given a labeled sample S that consists of pairs  $(a, f(a)) \in \Sigma \times C$ , construct a function f' that exactly or approximately agrees with f on S, and behaves well on unseen data. Depending on whether the labels are taken from a finite or infinite set C, the problem is known as *classification* or *regression* problem, respectively.

Machine learning problems, in the setting described above, are called *supervised*. In an *unsupervised* setting, the values of C are not provided, and the objects from S are grouped, clustered and classified, according to similarity measures.

Another important factor in learning problems is whether the learning is done *online* or *offline*. In online learning, the sample S gets updated during the learning

process and the classification or regression function can be updated to accommodate the new information that arrives. On the other hand, in offline problems the learning process starts after all sample data is given.

The learning algorithm can have an *active* role in the process by choosing the sample points that are believed to provide more information. In such a case the algorithm is called *active*. The algorithm is called *passive* otherwise.

Even though the primary goal of machine learning is to learn a *model* (or *func-tion*) and then use it for prediction, a conjectured model can be useful also in other contexts. For instance, data may come from a system that needs to be analyzed or checked for correctness. Such a system might be a complex program or a cyber-physical system with multiple components and embedded sub-systems. Often, a formal model of such a system is not available or particular components of the system are not accessible. As a result, the system can be viewed as a black box, or sometimes a gray box, where learning is used to identify the underlying formal model.

One main characteristic that differentiates learning models as system representations from conventional machine learning problems is that typically the inputs are *sequences* of events in contrast to the *static* objects that are more commonly used in machine learning. Typically, instances stand for behaviors of some dynamical system, a continuous one or an automaton.

Formally, the problem of learning an unknown system or model from observations is known as *inductive* or *model inference* or, more generally, *system iden-tification*. To characterize the type of the inductive inference problem, one needs to define three main aspects of the problem. First, the type of model that is to be learned; second, the way the observations are given; and third, the rules used to measure success. For instance, data may be given in an on-line fashion or in a batch; data may contain observations that are labeled or not; the underlying model is assumed to be of a certain type.

When the target model is assumed to be a language, the problem of inductive inference is called *language identification*. In this case, the learned model is usually a grammar or an automaton. Language identification can be found in the literature under the names of *grammatical inference*, *grammar induction*, *automata learning*, or *automata inference*. Language identification is concerned with learning language representations from information, which can be text, examples and counter-examples, or anything that can provide an insight into the elements of the languages being acquired.

Although, the terms induction and inference are interchangeably used, their meaning is not identical [DlH10]. In particular, when one has data and tries to find a representation that better explains them, e.g., a grammar or an automaton, the term *induction* is used; the derived model is basically used for prediction. On

the other hand, the term *inference* refers to problems where we assume that there exists an underlying target language, from which the data emerge, and one tries to discover it; the learning process, which is being examined and measured, is more important in this case.

Language identification is considered a sub-field that lies at the intersection of many fields such as: artificial intelligence, algorithmic learning theory, machine learning, statistics, formal languages and automata theory.

#### 3.2 Learning Languages

The problem of learning finite state systems was first suggested by Moore [Moo56], who described the situation where you find a black box with buttons and lights and try to infer the internal structure (identification) and solve other problems based on pushing the buttons and observing the lights. One of the first positive results on automata was given by Gold in the context of inductive inference, in an online and passive setting, where one observes a sequence of classified words and at each point in time maintains an automaton compatible with the sample seen so far. He showed that an algorithm based on the enumeration of all finite automata, will identify regular languages in the limit [Gol67, OG92]. This means that given any presentation of any language in the class the learner will produce only a finite number of wrong representations, and therefore converge on the correct representation in a finite number of steps, without however necessarily being able to announce its correctness since a counter-example to that representation could appear as an element arbitrarily long after.

Most 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 the concept of Nerode right-congruence relation [Ner58], which declares two *input histories* as equivalent if they lead to the same *future continuations*. In passive learning problems, given a sample of labeled examples, learning the smallest deterministic finite state automaton representation of the target regular language is NP-complete [Gol78]. Additionally, it cannot be approximated within any polynomial [PW93].

We are interested in active learning algorithms where the learner can interact with the teacher and select the words in the membership queries. The  $L^*$  algorithm [Ang87], on which this thesis is based, is one such algorithm. It overcomes the negative complexity results by introducing a minimally-adequate teacher, who can answer also equivalence queries that either confirm a conjectured automaton or provide a counter-example. In this setting, regular languages have been shown to be learnable using a polynomial number of queries.



Figure 3.1: The Teacher-Learner model consists of a learner (the learning algorithm) and a teacher. The teacher can answer two types of queries about the target language, membership queries (MQ), and equivalence queries (EQ). The learner asks MQ's until being able to construct a hypothesis H. Then, the validity of an hypothesis is checked by an EQ. If H is correct the algorithm terminates returning H, otherwise, a counter-example (cex) is returned and the learner improves the hypothesis.

#### **3.3** The *L*\* Learning Algorithm

Represented by the *learner*, the learning algorithm is designed to infer an unknown regular language L, the *target language*. The learner aims to construct a deterministic finite state automaton that recognizes the target language by gathering information from a *teacher*. The *teacher*, who knows the target language, can provide information about it by answering queries.

The algorithmic scheme used in  $L^*$ , is shown in Figure 3.1. The learner starts by asking *membership queries*, i.e., whether a word belongs to the target language or not. All information provided is suitably gathered in a table structure, the *observation table*. When the information is sufficient, the learner constructs a *hypothesis automaton* and poses an *equivalence query* to the teacher. If the answer is positive, i.e., the hypothesis suggested by the learner is the right one, the algorithm terminates and returns the conjectured automaton. Otherwise, the teacher responds to the equivalence query with a *counter-example*, a word misclassified by the conjectured automaton. Then, the learner incorporates the information provided by the counter-example into the table and repeats the procedure until a correct hypothesis is constructed.
#### **3.3.1** Observation table

An active learner organizes all information about the target language in a table structure, called the *observation table*. To construct a deterministic automaton, the learner needs to determine a set of states, with a subset of those determined as accepting, and a transition function. Hence, the observation table consists of three main parts: a set of candidate states and transitions, that is, a set of prefixes  $S \cup R \subseteq \Sigma^*$  that label the rows of the table; a set of distinguishing suffixes  $E \subseteq \Sigma^*$  that label the columns of the table and help to distinguish between states; and, the main body of the table, which consists of the classification of sequences whose membership is known. The rows in the table induce a Nerode right-congruence which is used to construct the hypothesis automaton.

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, that is, one of the shortest words that lead from the initial state to q. The elements of R hold information about the back- and cross-edges in the automaton.

Formally, let  $\Sigma$  be a finite alphabet. A prefix-closed set  $S \uplus R \subset \Sigma^*$  is a finite *balanced*  $\Sigma$ -*tree* if for all letters  $a \in \Sigma$ : 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*).

**Definition 3.1** (Observation Table). An *observation table* is a tuple  $T = (\Sigma, S, R, E, f)$  such that  $\Sigma$  is a finite alphabet;  $S \uplus R$  is a finite balanced  $\Sigma$ -tree, with R being its boundary; E is a suffix-closed subset of  $\Sigma^*$ ; and  $f : (S \cup R) \cdot E \to \{-,+\}$  is a classification function.

Let  $L \subseteq \Sigma^*$  be the target language. The classification function f of the table is a restriction of the characteristic function of L. The set  $M = (S \cup R) \cdot E$  consists the *sample* associated to the table. For all  $s \cdot e \in M$ , the element  $f(s \cdot e)$  is placed in the (s, e) position of the table. We say that a table has no holes when membership is known for all words in M.

With every prefix  $s \in S \cup R$ , we associate the residual  $f_s : E \to \{-,+\}$ , defined as  $f_s(e) = f(s \cdot e)$ , which characterizes the row of the observation table labeled by s. Two sequences s and r are considered as Nerode equivalent with respect to L, i.e,  $s \sim_L r$ , if  $f_s = f_r$ . Based on this right congruence relation, the learner can build an automaton from the table which, moreover, is minimal. However, the table should satisfy certain conditions.

**Definition 3.2.** An observation table  $T = (\Sigma, S, R, E, f)$  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  $f_{s \cdot a} = f_{s' \cdot a}, \forall a \in \Sigma$ .

**Theorem 3.3.** From a closed and reduced table, one can construct a finite deterministic automaton that is compatible with the sample.

*Proof.* Let  $T = (\Sigma, S, R, E, f)$  be a closed and reduced observation table. For a closed and reduced table we let  $g : R \to S$  be a function that maps every  $r \in R$  to the unique  $s \in S$  such that  $f_s = f_r$ .

We can define the automaton  $\mathcal{A}_T = (\Sigma, Q, q_0, \delta, F)$  such that  $Q = S, q_0 = \varepsilon$ ,  $F = \{s \in S : f_s(\varepsilon) = +\}$  and

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

С		
L		
L		

Note that a reduced table is trivially consistent. Moreover, every consistent non reduced table has a reduced form.

To explain this better, a table is not reduced when there exist  $s_1$  and  $s_2 \in S$  such that  $f_{s_1} = f_{s_2}$ . To reduce the table, we select one out of the two prefixes, preferably the longest one, and move it from S to R. This step may repeat several times until the table becomes reduced and S consists of rows that distinct. Note that, after converting a non reduced table to a reduced one, the set R is not necessarily equal to  $S \cdot \Sigma$ , and it may contain redundant elements.

To construct a finite deterministic automaton of a closed and consistent observation table that is compatible with the sample, we let the set of states of the automaton, which is a subset of S, to be the set of all the prefixes  $s \in S$  with distinguishable residuals  $f_s$ . The transition function is constructed similarly as for a reduced table.

#### 3.3.2 The Learning Algorithm

The learning algorithm, shown in Algorithm 2, proceeds as follows. The learner initializes the table to  $T = (\Sigma, S, R, E, f)$  such that  $S = E = \{\varepsilon\}$  and  $R = \Sigma$ . The table is then filled in, i.e., the classification f is determined for all words in the sample M that are missing, by asking membership queries to the teacher.

The learner attempts to keep the table closed and consistent at all times. The table T is not closed when there exists a prefix  $r \in R$  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 to R the  $\Sigma$ -successors of r, i.e., all words  $r \cdot a$  for  $a \in \Sigma$ , see Procedure 3. The table T

A	lgorithm	2 N	Main	Learni	ing	Al	lgori	tl	hr	n
---	----------	-----	------	--------	-----	----	-------	----	----	---

```
1: learned = false
 2: Let T = (\Sigma, S, R, E, f) such that S = E = \{\varepsilon\} and R = \Sigma \triangleright initialize table
 3: Ask MQ's to fill in T
 4: repeat
 5:
         while T is not closed or not consistent do
             CLOSE
 6:
             CONSISTENT
 7:
         end while
 8:
        Construct \mathcal{A}_T
 9:
         if EQ(A_T) then
10:
                                                                            \triangleright check hypothesis
             learned = true
11:
         else
                                                        \triangleright a counter-example w is provided
12:
              COUNTEREX(\mathcal{A}_T, w)
                                                                  \triangleright process counter-example
13:
         end if
14:
15: until learned
16: return \mathcal{A}_T
```

is inconsistent when there exist  $s_1, s_2 \in S$ ,  $e \in E$ , and  $a \in \Sigma$  such that  $f_{s_1} = f_{s_2}$ but  $f_{s_1}(a \cdot e) \neq f_{s_2}(a \cdot e)$ . To fix the inconsistency, the string  $a \cdot e$  is added to E as a new distinguishing experiment, see Procedure 4. In both cases the extended table is filled in by asking membership queries.

When the observation table  $T = (\Sigma, S, R, E, f)$  becomes both closed and consistent, following Theorem 3.3, the learner makes a hypothesis automaton  $\mathcal{A}_T = (\Sigma, Q, \varepsilon, \delta, F)$ , where  $Q \subseteq S$  is the set of states,  $F \subseteq Q$  is the set of final states, and  $\delta$  is the transition function. The set of states consists of all elements  $s \in S$ with distinct residual, i.e.,  $s, s' \in Q$  if and only if  $f_s \neq f_{s'}$ , and moreover, for all  $s \in S$  there exists  $s' \in Q$  such that  $f_s = f_{s'}$ . When the table is reduced, then Q = S. A state  $s \in Q$  is final when  $f_s(\varepsilon) = +$ . The transition function is defined as  $\delta(s, a) = s'$  such that  $f_{s'} = f_{s \cdot a}$ , for all  $s \in Q$  and  $a \in \Sigma$ .

When a counter-example is presented, the learner incorporates this into the sample and tries to construct another hypothesis. Variants of the  $L^*$  algorithm differ in the way they treat counter-examples. Subsequently, we present how a counter-example is treated in the algorithm  $L^*$ , followed by two variations, proposed in [MP95] and [RS93], respectively.

In any case, the table is extended and new states are added to the hypothesis. The learner asks new membership queries to fill in the table, closes the table and makes it consistent when needed, before making the new conjecture. This repeats until the teacher responds positively to an equivalence query.

#### Procedure 3 Table Closing

1: procedure CLOSE 2: Given  $r \in R$  such that  $\forall s \in S, f_r \neq f_s$ 3:  $S = S \cup \{r\}$ 4:  $R = (R - \{r\}) \cup r \cdot \Sigma$ 5: Ask MQ's to fill in T6: end procedure

Procedure 4 Make Table Consistent

1:	procedure CONSISTENT
2:	Given $s, s' \in S, a \in \Sigma$ and $e \in E$ such that
3:	$f_s = f_{s'}$ and $f_{s \cdot a}(e) \neq f_{s' \cdot a}(e)$
4:	$E = E \cup \{a \cdot e\}$
5:	Ask MQ's to fill in $T$
6:	end procedure

#### **Counter-Example Treatment**

**The classic**  $L^*$  **algorithm.** The original algorithm [Ang87] adds all the *prefixes* of a counter-example w to S, see Procedure 5. Then all continuations of the new prefixes are added to R, such that  $S \cdot \Sigma \subseteq S \cup R$  holds. The learner fills in the table by asking membership queries. The counter-example treatment creates an inconsistency or renders the table non closed.

Notice that this treatment of the counter-example adds many redundant rows to the table and only a subset of S forms the set of states in the conjectured automaton.

```
Procedure 5 Counter-Example Treatment: L^*
 1: procedure COUNTEREX
 2:
         let w \in \Sigma^* be a counter-example
         for u \in \Sigma^* such that w = u \cdot v, v \in \Sigma^* do
 3:
             S = S \cup \{u\}
                                                           \triangleright add w and all its prefixes to S
 4:
             R = R \cup u \cdot \Sigma
 5:
         end for
 6٠
         Ask MQ's to fill in T
 7:
 8: end procedure
```

Adding Suffixes. This version is proposed in [MP95] and was initially used for learning  $\omega$ -regular languages. The learner adds all the *suffixes* of a counter-example

w to E, see Procedure 6. Then the learner fills in the table.

The advantage of this approach is that the table always remains reduced and thus consistent, with S corresponding exactly to the set of states. A disadvantage, however, is the possible introduction of redundant columns that do not contribute to further discrimination between states.

**Breakpoint.** Another variation of the counter-example treatment in this algorithm is due to [RS93] and is referred in [BR05] as the *reduced algorithm*. In this variation, the counter-example is first processed by the so-called *breakpoint* method, and then only the suffixes necessary to detect a new state are added to E.

The breakpoint method, see Procedure 7, searches for a cut point in the counterexample, the *breakpoint*. This provides a factorization of the counter-example whose suffix detects a new state. Formally, let  $w = a_1 \cdots a_{|w|}$  be a counterexample to a hypothesis  $\mathcal{A}_T$  associated with a table T. An *i*-factorization of w is  $w = u_i \cdot a_i \cdot v_i$  such that  $u_i = a_1 \cdots a_{i-1}$  and  $v_i = a_{i+1} \cdots a_{|w|}$ . Let  $s_i = \delta(\varepsilon, u_i \cdot a_i)$ be the state of  $\mathcal{A}_T$  (element of S in table T) reached after reading  $u_i \cdot a_i$  and let  $s_0$ be the initial state.

**Proposition 3.4** (Breakpoints [RS93]). If w is a counter-example to  $A_T$  then there is an *i*-factorization of w such that  $f(s_{i-1} \cdot a_i \cdot v_i) \neq f(s_i \cdot v_i)$ .

In other words, adding  $v_i$  to E distinguishes  $s_{i-1} \cdot a_i$  from  $s_i$ , makes the table not closed and leads to the identification of  $s_{i-1} \cdot a_i$  as a new state. The breakpoint method iterates over i values between 1 and |w| to find new states and experiments that distinguish them from old states. Note that a counter-example may admit more than one breakpoint.

We let i take values in a monotonically decreasing order. The reason for doing this is to favor the detection of smaller distinguishing strings that are added as suffixes to the table. The iteration in an increasing order can be used as well, making the breakpoint method more effective when we need to add prefixes to the table. Alternatively, we can proceed by choosing the values i in a binary search manner and be sure to find a breakpoint in logarithmic time, or in an exponential search manner when of favor to short suffixes. These different iteration methods are experimentally tested in [IS14].

**Example 3.5.** We illustrate the behavior of the  $L^*$  algorithm while learning the target language  $L = a \cdot \Sigma^*$  over  $\Sigma = \{a, b\}$ . We use +w to indicate a counterexample  $w \in L$  rejected by the conjectured automaton, and -w otherwise.

Initially, the observation table is  $T_0 = (\Sigma, S, R, E, f)$  with  $S = E = \{\varepsilon\}$  and  $R = \Sigma$ . We ask membership queries for all words in  $(S \cup R) \cdot E = \{\varepsilon, a, b\}$  to obtain table  $T_0$ , shown in Figure 3.2. The table is not closed so we move a to S, add its continuations,  $a \cdot a$  and  $a \cdot b$  to R, and ask membership queries to obtain the

Procedure 6 Counter-Example Treatment: Adding Suffixes

- 1: procedure COUNTEREX
- 2: let  $w \in \Sigma^*$  be a counter-example
- 3: for each  $v \in \Sigma^*$  such that  $w = u \cdot v, u \in \Sigma^*$  do
- 4:  $E = E \cup \{v\}$   $\triangleright$  add w and all its suffixes to E
- 5: end for
- 6: Ask MQ's to fill in T
- 7: end procedure

Procedure 7 Counter-Example Treatment: Breakpoint

1: procedure COUNTEREX let  $w \in \Sigma^*$  be a counter-example 2: for i = |w|, ..., 1 do 3: let  $u_i \cdot a_i \cdot v_i$  be the *i*-factorization of w 4: if  $f(s_{i-1} \cdot a_i \cdot v_i) \neq f(s_i \cdot v_i)$  then 5:  $E = E \cup \{v_i\}$ 6: break 7: end if 8: 9: end for Ask MQ's to fill in T10: 11: end procedure



Figure 3.2: Observation tables of Example 3.5 when learning  $a \cdot \Sigma^*$  using the counterexample treatment method that is used in the classic  $L^*$  algorithm.



Figure 3.3: Hypothesis automata of Example 3.5 when learning  $a \cdot \Sigma^*$ .

closed table  $T_1$ . From  $T_1$  we construct the hypothesis automaton  $A_1$  of Figure 3.3. In response to the equivalence query for  $A_1$ , counter-example  $-b \cdot a$  is presented.

First, let us use the classical counter-example treatment used by the  $L^*$  algorithm. In this case, the learner adds all prefixes of the counter-example  $b \cdot a$  to the set of prefixes, that is, it adds the prefixes b and ba to S. All successors are added to R, resulting in table  $T_2$  of Figure 3.2. This table is not consistent since the two elements  $\varepsilon$  and b in S are considered equivalent but their a-successors a and  $b \cdot a$  are not. Adding a to E and asking membership queries yields a closed and consistent table  $T_3$ . The derived automaton  $A_3$  is the desired automaton that recognizes the target language L and the algorithm terminates.

In the following, we will use the same example, but we will apply another method for treating the counter-example  $-b \cdot a$ . We have chosen an example where the hypotheses remain the same, but this is frequently not the case. We focus here on the different observation tables  $T_2$  and  $T_3$  that are produced for each method.

Let us use the suffixes method for treating  $-b \cdot a$ , we add all suffixes, i.e., a and  $b \cdot a$ , to E. The resulting table is thus  $T_2^m$ , see Figure 3.4. We recall here, that with this method the learner need not check for consistency, as the table is always kept consistent. After closing the table  $T_2^m$  updates  $T_3^m$ . The hypothesis associated to this table is  $A_3$ , shown in Figure 3.3. The algorithm terminates as  $A_3$  is correct.

Finally, we analyze the counter-example  $-b \cdot a$  using the breakpoint method. This results in finding one distinguishing string that is enough to be added to the table and to distinguish some states. The *i*-factorization for i = 1 yields such a suffix, since  $f(b \cdot a) \neq f(\varepsilon \cdot a)$ , as Procedure 7 indicates. As a result, the suffix  $v_1 = a$  is added to E and the new table is  $T_2^b$  is shown in Figure 3.5. After making it closed, the table updates to  $T_3^b$  and the hypothesis  $A_3$ , see Figure 3.3, is conjectured. The algorithm terminates with this hypothesis being correct.

		~~~				$T_{z}$	$m_{3}$	
$T_2^m$						ε	a	$b \cdot a$
	ε	a	$b \cdot a$	]	ε	_	+	_
ε	—	+	_		a	+	+	+
a	+	+	+		b	-	_	_
b	_	_	_	]	$a \cdot a$	+	+	+
$a \cdot a$	+	+	+		$a \cdot b$	+	+	+
$a \cdot b$	+	+	+		$b \cdot a$	-	—	—
					$b \cdot b$	-	_	_

Figure 3.4: Observation tables of Example 3.5 when learning  $a \cdot \Sigma^*$  using the counterexample treatment of the suffixes method.

	<b>T</b> h			$T_3^b$	
-	$T_{2}^{0}$			ε	a
	ε	a	ε	_	+
ε	—	+	a	+	+
a	+	+	b	_	_
b	_	—	$a \cdot a$	+	+
$a \cdot a$	+	+	$a \cdot b$	+	+
$a \cdot b$	+	+	$b \cdot a$	-	_
			$b \cdot b$	_	—

Figure 3.5: Observation tables of Example 3.5 while learning  $a \cdot \Sigma^*$  and using the counterexample treatment of the breakpoint method.

# Symbolic Automata

#### 4.1 State of the Art

In recent years, interest in languages defined over large and infinite alphabets has increased both for theoretical and practical consideration [Seg06]. When regular languages are defined over a large or infinite domain it is impractical or unfeasible to use standard finite automata to represent them. Several authors consider automata models that use infinite alphabets, namely, *finite memory automata*, also known as *register automata*, [KF94, NSV04]; *pebble automata* [NSV04]; *data automata* [BMS<sup>+</sup>06]; and *variable automata* [GKS10]. These are finite state automata that use a finite memory, e.g., registers, pebbles, or variables, to handle the infinite nature of the alphabet.

Other automata models over infinite alphabets are finite state automata, where transitions are labeled by predicates applicable to the alphabet in question. Hence, transitions are associated with formulas representing sets of letters rather than individual letters resulting in automata with finitely many transitions and thus with a finite representation. No memory is used and predicates are evaluated only on the current input. The idea of such automata was mentioned already in [Wat96]. However, it was formally introduced and studied later, in [VNG01], under the name of *predicate-augmented finite state recognizers*. In a similar flavor, *lattice automata* [LGJ07] use elements of an atomic lattice to label transitions rather than general predicates. *Interval automata*, which are automata where a move is realized whenever the input letter belongs to a predefined interval, can be viewed as a special case of both lattice automata and predicate-augmented finite state recognizers. Automata with predicates were further studied in [VDHT10, Vea13], under the name *symbolic automata* and predicates are chosen from a Boolean Algebra. Finally, one can find a special variant of symbolic automata in [HJJ<sup>+</sup>95, BKR96, YBC108],

where transition relations are represented as Multi-terminal Binary Decision Diagrams.

Hereafter, we use the automata, as defined in [VNG01, Vea13], and refer to them as *symbolic automata*. The input, which comes from a concrete alphabet that may be large or infinite, is treated in a symbolic way. Transitions are labeled by symbols taken from a finite symbolic alphabet. Each symbol stands for a predicate and represents a non-empty set of concrete letters. Symbolic automata are shown to be closed under union, complement, concatenation, and star operations [VNG01, YBC108]. Using the product construction one can show that they are also closed under intersection and difference [VNG01, VBDM10, HV11]. Moreover, symbolic automata are determinizable, they can be made complete, and  $\varepsilon$ transitions can be eliminated [VNG01, VDHT10]. Minimization algorithms are adapted to the symbolic setting and studied in detail in [DV14]. Equivalence checking is decidable and counter-examples can be found in the case where two symbolic automata are not equivalent [VHL<sup>+</sup>12, DV14].

Their succinct representation makes symbolic automata suitable for language representation not only when the input alphabet is infinite, but also when the alphabet is finite but large, e.g., the Unicode alphabet UTF-16. The smaller size of the automaton allows faster implementation of various operations [VNG01]. The necessity of introducing and studying symbolic automata derives from several applications which have their roots in quite diverse areas such as natural language processing [VNG01]; support for regular expressions in the context of static and dynamic program analysis [VDHT10]; constraint solving and symbolic analysis with SMT solvers [VBDM10]; support for some specific regular expressions, e.g., LIKE-patterns, used in the context of symbolic analysis of database queries [VGDHT09, VTDH10]; and security analysis of string sanitizers [Vea13].

In this chapter, we first introduce symbolic automata (Section 4.2) and then discuss their closure properties and canonical form (Section 4.3). We end the chapter with some examples of symbolic automata defined over specific domains, such as numerical alphabets, e.g.,  $\mathbb{N}$ ,  $\mathbb{R}$ , or sets of vectors, e.g.,  $\mathbb{B}^n$ .

#### 4.2 Definition

Let  $\Sigma$  be a large, possibly infinite, alphabet, to which we will refer from now on as the *concrete* alphabet. A *symbolic automaton* over  $\Sigma$  is a finite automaton where each state has a small number of outgoing transitions labeled by *symbolic letters* (or *symbols*) taken from a finite alphabet  $\Sigma$ ; we call this the *symbolic alphabet*. Symbols represent subsets of  $\Sigma$ , which may differ from state to state.

Formally, let  $\Sigma$  be a disjoint union of finite symbolic alphabets of the form  $\Sigma_q$ , each associated with a state q of the automaton. We choose a state-depended

alphabet as being more intuitive and efficient for the purposes of the current thesis. Nevertheless, for language-theoretic studies, one may find that a globally defined symbolic alphabet, where the alphabet is the same for all states, is a more suitable option. A state local alphabet can be easily transformed into a global one; the procedure is presented in detail at the end of this section.

Concrete letters are mapped to symbols through a family of mappings  $\psi$  decomposable into state-specific mappings  $\psi_q : \Sigma \to \Sigma_q$ , where q is a state in the automaton. The  $\Sigma$ -semantics of a symbolic letter  $a \in \Sigma$  denotes the subset of concrete letters from  $\Sigma$  that take the transition labeled by a in the automaton. This can be seen as a mapping  $\llbracket \cdot \rrbracket : \Sigma \to 2^{\Sigma}$  defined as  $\llbracket a \rrbracket = \{a \in \Sigma : \psi_q(a) = a\}$  for any symbol  $a \in \Sigma_q$ . We often omit  $\psi$  and  $\psi_q$  from the notation and use  $\llbracket a \rrbracket$  when  $\psi$ , which is always present, is clear from the context.

The  $\Sigma$ -semantics is extended to symbolic words of the form  $w = a_1 \cdots a_{|w|} \in \Sigma^*$  as the concatenation of the concrete one-letter languages associated with the respective symbolic letters or, recursively speaking, by letting  $[\![\varepsilon]\!] = \{\varepsilon\}$  and  $[\![w \cdot a]\!] = [\![w]\!] \cdot [\![a]\!]$  for  $w \in \Sigma^*$ ,  $a \in \Sigma$ .

The symbolic automaton is *deterministic* when the family  $\psi$  consists of functions; and, it is *complete* when the mappings in  $\psi$  are total. Moreover, if the functions in  $\psi$  are surjective, there are no redundant symbols in  $\Sigma$ , i.e., symbols whose semantics is empty. As a result, a total surjective function  $\psi_q$ , where q is a state in the symbolic automaton, implies that the set of  $\Sigma$ -semantics of all symbols in  $\Sigma_q$ , {[[a]] :  $a \in \Sigma_q$ }, forms a *partition* of  $\Sigma$ . This should hold for all states in the automaton for it to be complete and deterministic. Formally,

**Definition 4.1** (Symbolic Automaton). A *deterministic symbolic automaton* is a tuple  $\mathcal{A} = (\Sigma, \Sigma, \psi, Q, \delta, q_0, F)$ , where

- $-\Sigma$  is the input alphabet,
- $\Sigma$  is a finite alphabet, decomposable into  $\Sigma = \biguplus_{q \in Q} \Sigma_q$ ,
- $-\psi = \{\psi_q : q \in Q\}$  is a family of total surjective functions  $\psi_q : \Sigma \to \Sigma_q$ ,
- -Q is a finite set of states,
- $-q_0$  is the initial state,
- -F is the set of accepting states, and
- $-\delta: Q \times \Sigma \to Q$  is a partial transition function decomposable into a family of total functions of the form  $\delta_q: \{q\} \times \Sigma_q \to Q$ .

Let  $w = a_1 \cdots a_n \in \Sigma^*$  be a concrete word. A run of w on  $\mathcal{A}$  produces a path  $P(w) : q_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} \ldots \xrightarrow{a_n} q_n$  that corresponds to the symbolic path  $P(w) : q_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} \ldots \xrightarrow{a_n} q_n$ , where  $q_i \in Q$ ,  $a_i \in \Sigma_{q_{i-1}}$ ,  $a_i \in [a_i]$ , and  $\delta(q_{i-1}, a_i) = q_i$  for all  $i \in \{1, \cdots, n\}$ .

Hence, a symbolic automaton  $\mathcal{A}$  can be viewed as an alternative representation of a concrete deterministic automaton whose concrete transition function  $\delta : Q \times$  $\Sigma \to Q$  is defined as  $\delta(q, a) = \delta(q, \psi_q(a))$ . That is, when at q and reading a concrete letter a, the automaton takes the transition  $\delta(q, a)$ , where a is the *unique* element of  $\Sigma_q$  satisfying  $a \in [\![a]\!]$ . The transition function extends to words as in the concrete case. The language recognized by the symbolic automaton, denoted by  $L(\mathcal{A})$ , consists of all concrete words whose run leads from  $q_0$  to a state in F.

A language  $L \subseteq \Sigma$  is symbolic recognizable if there exists a symbolic automaton  $\mathcal{A}$  such that  $L = L(\mathcal{A})$ . Let us denote by  $SREC(\Sigma)$  the set of all symbolic recognizable languages defined over the large alphabet  $\Sigma$ . One can easily verify that when  $\Sigma$  is a finite alphabet, then  $REC(\Sigma) \equiv SREC(\Sigma)$ . As for the classical case, two symbolic automata are *equivalent* when they accept the same language.

As mentioned before, the association of a symbolic language with a symbolic automaton is more subtle because we allow different partitions of  $\Sigma$  and hence different symbolic alphabets at different states, rendering the transition function *partial* with respect to  $\Sigma$ . When in a state q and reading a symbol  $a \notin \Sigma_q$ , the transition to be taken is well defined only when  $[\![a]\!] \subseteq [\![a']\!]$  for some  $a' \in \Sigma_q$ . The model can, nevertheless, be made deterministic and complete over a refinement of the symbolic alphabet. For this, let

$$\Sigma' = \prod_{q \in Q} \Sigma_q$$
, with the  $\Sigma$ -semantics  $\llbracket (\boldsymbol{a}_1, \dots, \boldsymbol{a}_n) \rrbracket = \llbracket \boldsymbol{a}_1 \rrbracket \cap \dots \cap \llbracket \boldsymbol{a}_n \rrbracket$ 

and let  $\tilde{\Sigma} = \{ \boldsymbol{b} \in \Sigma' : [\boldsymbol{b}] \neq \emptyset \}$ . An ordinary automaton  $\tilde{\boldsymbol{\mathcal{A}}} = (\tilde{\Sigma}, Q, \tilde{\delta}, q_0, F)$  can be defined, where, by construction, for every  $\boldsymbol{b} \in \tilde{\Sigma}$  and every  $q \in Q$ , there is  $\boldsymbol{a} \in \Sigma_q$  such that  $[\![\boldsymbol{b}]\!] \subseteq [\![\boldsymbol{a}]\!]$  and hence one can define the transition function as  $\tilde{\delta}(q, \boldsymbol{b}) = \delta(q, \boldsymbol{a})$ . This model is more comfortable for language-theoretic studies but we stick hereafter to Definition 4.1 as it is more efficient.

#### 4.3 **Operations on Symbolic Automata**

After defining and understanding symbolic automata and symbolic recognizable languages, it is important to discuss closure properties as well as decidability of classical problems of this class of languages. In the following, one can find that the set of symbolic recognizable languages is closed under the *Boolean* and *regular* operations. Moreover, decidability problems such as *membership*, *emptiness*, *determinization*, and *equivalence* do not depend on the decidability of the predicate theory that is used and thus they are easily adaptable in the case of symbolic automata [VNG01, HV11, DV14]. Finally, symbolic automata admit a canonical representation which is minimal in the number of states. This is an important

property, especially when inferring symbolic recognizable languages.

**Proposition 4.2** (Closure under Boolean Operations). *Languages accepted by deterministic symbolic automata are closed under the Boolean operations.* 

*Proof.* Closure under complementation is immediate by complementing the set of accepting states. For the intersection, we adapt the standard product construction. Then, closure under union follows.

Let  $L_1, L_2$  be two languages recognized, respectively, by the symbolic automata  $\mathcal{A}_1 = (\Sigma, \Sigma^1, \psi^1, Q^1, \delta^1, q_0^1, F^1)$  and  $\mathcal{A}_2 = (\Sigma, \Sigma^2, \psi^2, Q^2, \delta^2, q_0^2, F^2)$ . The intersection of the two languages is then recognized by the symbolic automatom  $\mathcal{A} = (\Sigma, \Sigma, \psi, Q, \delta, q_0, F)$ , where

 $\begin{aligned} - & Q = Q^1 \times Q^2, q_0 = (q_0^1, q_0^2), F = F^1 \times F^2 \\ - & \text{For every } (q^1, q^2) \in Q \\ & - & \sum_{(q^1, q^2)} = \{(a_1, a_2) \in \Sigma^1 \times \Sigma^2 : \llbracket a_1 \rrbracket \cap \llbracket a_2 \rrbracket \neq \emptyset \} \\ & - & \psi_{(q^1, q^2)}(a) = (\psi_{q^1}^1(a), \psi_{q^2}^2(a)) \ \forall a \in \Sigma \\ & - & \delta((q^1, q^2), (a_1, a_2)) = (\delta^1(q^1, a_1), \delta^2(q^2, a_2)), \forall (a_1, a_2) \in \Sigma_{(q^1, q^2)} \end{aligned}$ 

It is sufficient to observe that the corresponding implied concrete automata  $A_1$ ,  $A_2$ and A satisfy  $\delta((q^1, q^2), a) = (\delta^1(q^1, a), \delta^2(q^2, a))$  and the standard proof that  $L(A) = L(A_1) \cap L(A_2)$  follows.

One can use the same construction as for the intersection, just by changing the set of final states, to construct a symbolic automaton that recognizes the difference or symmetric difference of two languages. Algorithms for the product and difference constructions for symbolic automata appear in [HV11], as well as experimental evaluation on their performance. Moreover, it is shown that symbolic recognizable languages are closed under the regular operations of concatenations and star [VNG01], following the corresponding constructions for classical finite automata.

A symbolic automaton with  $\varepsilon$ -transitions is a symbolic automaton that contains transitions of the form  $\delta(q, \varepsilon) = q'$ , where  $\varepsilon$ , which is a symbol that represents the empty word, is not a symbol of  $\Sigma_q$ . An  $\varepsilon$ -free symbolic automaton is an automaton that does not contain such transitions. Any symbolic automaton can be transformed to an equivalent  $\varepsilon$ -free symbolic automaton. To eliminate  $\varepsilon$ -transitions from a symbolic automaton  $\mathcal{A}$  the  $\varepsilon$ -closure method [HMU06] is used, where transitions need to adapt to the symbolic setting [VDHT10].

A symbolic automaton can be made complete simply by adding a sink state  $q_{sink}$ . Then, whenever  $\bigcup_{a \in \Sigma_q} \llbracket a \rrbracket \neq \Sigma$ , add a new symbol  $a_{sink}$  to  $\Sigma_q$ , and let  $\delta(q, a_{sink}) = q_{sink}$  and  $\llbracket a_{sink} \rrbracket = \Sigma \setminus \bigcup_{a \in \Sigma_q} \llbracket a \rrbracket$ .

Determinization of symbolic automata also follows the same idea as for the concrete case. The set of states in the deterministic automaton is a subset of  $2^Q$ . A small but nevertheless important difference is that when we handle symbolic automata the transitions should be defined more carefully not to overlap. To explain this, let us assume non determinism at a state q such that there are two transitions  $(q, a, q_1)$  and  $(q, b, q_2)$ , where  $a, b \in \Sigma_q$  and  $[\![a]\!] \cap [\![b]\!] \neq \emptyset$ . To make the state deterministic we replace the involved states by new states, and replace the transitions by the new ones  $(\{q\}, a_1, \{q_1\}), (\{q\}, a_{12}, \{q_1, q_2\}),$  and  $(\{q\}, a_2, \{q_2\})$ , where  $[\![a_1]\!] = [\![a]\!] \setminus [\![b]\!], [\![a_{12}]\!] = [\![a]\!] \cap [\![b]\!]$ , and  $[\![a_2]\!] = [\![b]\!] \setminus [\![a]\!]$ , as long as the semantics remain non empty.

Finally, symbolic automata admit a canonical representation. For any symbolic recognizable language, there exists a unique minimal automaton recognizing it, up to isomorphism on the states. As in the classical case, we can define a Nerode equivalence relation and then construct the minimal automaton.

The left quotient can be extended to a symbolic prefix  $u \in \Sigma^*$  as follows:

$$\boldsymbol{u}^{-1}L = \{ v \in \Sigma^* : uv \in L, u \in \llbracket \boldsymbol{u} \rrbracket \}.$$

A Nerode equivalence relation, induced by the language  $L \subseteq \Sigma^*$ , can be lifted to symbolic words as follows

$$\boldsymbol{u} \equiv_L \boldsymbol{v}$$
 if and only if  $\boldsymbol{u}^{-1}L = \boldsymbol{v}^{-1}L$ .

Based on the Nerode Theorem (Theorem 2.1) we build a minimal symbolic automaton  $\mathcal{A}$  using the syntactic congruence. We define this automaton as  $\mathcal{A} = (\Sigma, \Sigma, Q, q_0, \delta, F)$ , where Q is the set of left quotients of L, i.e.,  $Q = \Sigma^* / \equiv_L$ ; the initial state is the class that contains the empty word, i.e.,  $q_0 = [\varepsilon]$ ; and final states are all those classes that contain accepted words, i.e.,  $F = \{[u] : u \cdot \varepsilon \in L\}$ . Defining the symbolic letters is a bit more tricky and only symbols with non-empty semantics are allowed. Let p and q be two states from Q that correspond to the classes [u] and [v] respectively, we define  $\Sigma_p$  as the set of letters  $a_{pq}$  whose semantics is defined as  $[\![a_{pq}]\!] = \{a \in \Sigma : [u \cdot a] \equiv_L [v]\}$  and it is not empty. The symbolic transition function is then given by  $\delta(p, a_{pq}) = q$ . By construction, this automaton has no redundant symbols, moreover, it is complete and deterministic. Therefore,  $\mathcal{A}$  is minimal [DV14].

Moreover,  $\mathcal{A}$  is unique up to renaming of states and up to the equivalence of the  $\Sigma$ -semantics of the symbols. That is, if p, q are two states in  $\mathcal{A}$ , and p', q' are two states in  $\mathcal{B}$ , which is another minimal automaton that recognizes L, such that  $p \sim p'$  and  $q \sim q'$ , then  $[\![a_{pq}]\!] = [\![b_{p'q'}]\!]$  where  $a_{pq}$  and  $b_{p'q'}$  are symbolic letters in  $\mathcal{A}$  and  $\mathcal{B}$ , respectively.

Algorithms for minimization of symbolic automata can be found in [DV14]. In this paper, the authors present an adaptation of the classical minimization algorithms, such as Moore's and Hopcroft's algorithms, to the symbolic setting as well as a symbolic minimization algorithm.

#### 4.4 Alphabets and Partitions

Symbolic automata, as defined in the two previous sections, can be defined over any type of alphabet and the predicates can be of arbitrary form. However, we restrict ourselves to predicates that form a partition over the alphabet at each state, and thus to deterministic symbolic automata. This section presents some examples of symbolic automata over specific alphabet domains and partitions. Details on sets and partitions can be found in Section 2.2

#### 4.4.1 Interval Automata

In this section, we assume an ordered alphabet that is a finite or infinite subset of a numerical set such as  $\mathbb{N}$ ,  $\mathbb{Q}$ , or  $\mathbb{R}$ . We let the  $\Sigma$ -semantics of each symbol form closed-open intervals. To preserve convexity<sup>1</sup>, unions of intervals are not allowed in the semantics. This practice is preferable due to simplicity, even though it may cause the symbolic automaton to not be minimal with respect to the number of transitions, as defined in Section 4.3. Assuming a small number of intervals in each partition, the symbolic automaton can still be well defined and have a finite symbolic alphabet. In the following, we provide some concrete examples.

**Example 4.3.** Let  $\mathcal{A}$  be an interval symbolic automaton as the one shown in Figure 4.1. The symbolic transition function  $\delta$  is finite and can be given as a table. The symbolic alphabet  $\Sigma = \{a_0, \ldots, a_7\}$  can also be viewed as the union of the state-depended symbolic alphabets  $\Sigma_{q_0} = \{a_0, a_1\}, \Sigma_{q_1} = \{a_2, a_3\}, \Sigma_{q_2} = \{a_4, a_5, a_6\}, \text{ and } \Sigma_{q_3} = \{a_7\}.$  The concrete input alphabet of  $\mathcal{A}$  is  $\Sigma = [0, 100) \subseteq \mathbb{R}$ . The  $\Sigma$ -semantics for the symbols are defined as  $[\![a_0]\!] = [0, 50), [\![a_1]\!] = [50, 100), [\![a_2]\!] = [0, 30)$ , etc., see Figure 4.2.

A similar symbolic automaton can be defined by letting the concrete input alphabet be a subset of the natural numbers  $\mathbb{N}$ , that is  $\Sigma = \{0, 1, \dots, 99\}$ . This automaton has the same structure as  $\mathcal{A}$  and differs only in its  $\Sigma$ -semantics, which should be intersected with  $\mathbb{N}$  for all  $a_i \in \Sigma$ , e.g.,  $[\![a_0]\!] = \{0, \dots, 49\}, [\![a_1]\!] = \{50, \dots, 99\}, [\![a_2]\!] = \{0, \dots, 29\}$ , etc.

<sup>&</sup>lt;sup>1</sup>A subset  $\Sigma$  of  $\mathbb{R}^n$  is said to be *convex* if it contains all the line segments connecting any pair of its points. Every interval is, by definition, a convex set.



Figure 4.1: A symbolic automaton  $\mathcal{A}$  with its symbolic transition function  $\delta$ .

$\psi$	0 2	0 3	0 5	0	8	0 100
$\mathbf{\Sigma}_{q_0}$		$oldsymbol{a}_0$			$a_1$	
$\mathbf{\Sigma}_{q_1}$	$a_2$			$oldsymbol{a}_3$		
$\mathbf{\Sigma}_{q_2}$	$a_4$		a	5		$oldsymbol{a}_6$
$\mathbf{\Sigma}_{q_3}$	$a_7$					

Figure 4.2: The concrete semantics of the symbols of automaton  $\mathcal{A}$  of Figure 4.1. The input alphabet is  $\Sigma = [0, 100) \subseteq \mathbb{R}$ .



Figure 4.3: The monotone partitions over the concrete alphabet  $\Sigma = [0, 100)^2$  as defined by  $\psi_q$  at each state q of the automaton  $\mathcal{A}$  (Figure 4.1).

#### 4.4.2 Automata over Partially-ordered Alphabets

In this section we assume the input alphabet to be a partially-ordered set of the form  $\Sigma = X^d$ , where X is a totally-ordered set such as an interval  $[0, k) \subseteq \mathbb{R}$ . We can define a symbolic automaton  $\mathcal{A}$  where at each state q the function  $\psi_q$  forms a monotone partition of the concrete alphabet defined as a finite union of cones, see Section 2.2. Such a partition is of the form  $\mathcal{P} = \{P_1, \ldots, P_{m-1}\}$ , where  $P_i = B^+(F_{i-1}) - B^+(F_i)$  and  $\mathcal{F} = \{F_0, \ldots, F_{m-1}\}$  is a family of sets of minimal incomparable points in  $\Sigma$ . Each symbol  $a_i \in \Sigma_q$  associates to the partition block  $P_i$ , such that  $[\![a_i]\!] = P_i$ . We assume that partition blocks are simply-connected and the  $\Sigma$ -semantics does not form unions of non-connected subsets of  $\Sigma$ .

**Example 4.4.** Let the alphabet be  $\Sigma = [0, 100)^2$ . A symbolic automaton  $\mathcal{A}$  with a monotone partition at each state is shown in Figure 4.1. The semantics function  $\psi$ , which is shown in Figure 4.3, forms at each state a monotone partition of  $\Sigma$ . The partitions are defined from the following families of points:

$$\begin{split} P_{q_0} &= \big\{\{(0,0)\}, \{(0,70), (45,50), (60,0)\}\big\}, \\ P_{q_1} &= \big\{\{(0,0)\}, \{(0,80), (30,30), (40,15), (80,0)\}\big\}, \\ P_{q_2} &= \big\{\{(0,0)\}, \{(0,30), (20,0)\}, \{(0,90), (60,70), (70,50), (90,0)\}\big\}, \text{and} \\ P_{q_3} &= \big\{\{(0,0)\}\big\}. \end{split}$$

#### 4.4.3 Boolean Vectors

Let us define a symbolic automaton  $\mathcal{A}$  over the concrete alphabet  $\Sigma = \mathbb{B}^n$ , that is, the Boolean hyper-cube of dimension  $n \in \mathbb{N}$  where letters are vectors accessed by the Boolean variables  $X = \{x_1, \ldots, x_n\}$ . We let the  $\Sigma$ -semantics of the symbolic letters form sub-cubes of  $\Sigma$ , represented by a family of functions  $\{\psi_q\}_{q \in Q}$  as usual. Hence, for each state q of a symbolic automaton, the function  $\psi_q$  is a *pseudoboolean function*  $\psi_q : \mathbb{B}^n \to \Sigma_q$ . There are multiple ways to represent such a function. We choose the representations of a *binary decision tree* (BDT), a *pseudoboolean normal form* (PBNF), or Karnaugh maps when the number of variables nis small. We refer the reader to Section 2.2 for a brief description of partitions over a Boolean Cube and its different representations.

**Example 4.5.** Let the symbolic automaton  $\mathcal{A}$  have the same structure as the one of Figure 4.1. Transitions are taken based on the same transition function but the semantics are determined by a family  $\psi = \{\psi_q\}_{q \in Q}$  of pseudo-boolean functions.



Figure 4.4: The partitions over the concrete alphabet  $\Sigma = \mathbb{B}^4$  as defined by  $\psi_q$  at each state q of the automaton  $\mathcal{A}$  (Figure 4.5).

The  $\Sigma$ -semantics, which are shown in Figure 4.4 as Karnaugh maps, are defined as

$$\psi_{q_0} = a_0 \cdot \bar{x_1} + a_1 \cdot x_1,$$
  
 $\psi_{q_1} = a_2 \cdot \bar{x_2} + a_3 \cdot x_2,$   
 $\psi_{q_2} = a_4 \cdot x_1 + a_5 \cdot \bar{x_1}x_3 + a_6 \cdot \bar{x_1}\bar{x_3},$  and  
 $\psi_{q_3} = a_7.$ 

The symbolic automaton  $\mathcal{A}$  can be seen in Figure 4.5 where the symbols are replaced by Boolean guards.



Figure 4.5: The symbolic automaton  $\mathcal{A}$  of Figure 4.1 where a symbolic transition is taken when a guard, which is a Boolean function, is satisfied.

### Learning Symbolic Automata

Learning algorithms, as presented in Chapter 3, 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 of the shortest words that lead to each state and, in addition, the set of all their  $\Sigma$ -continuations, where  $\Sigma$  is the input alphabet. Over large alphabets, these algorithms do not scale well. 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.

In this chapter, we introduce the *symbolic learning algorithm*, an algorithm for learning languages over large alphabets that can be represented by symbolic automata models (see Chapter 4). Our algorithm can be seen as a symbolic adaptation of the  $L^*$  algorithm. We assume an oracle, the *teacher*, that knows the target language and can answer *membership* (MQ) and *equivalence queries* (EQ). The *learner* first asks MQ's and collects all the information in a table structure, the *symbolic observation table*, and then, when the information in the symbolic observation table is sufficient, the learner constructs a hypothesis. The hypothesis is tested by posing an EQ; whenever the teacher returns a counter-example, the learner revises the hypothesis, it terminates otherwise.

The entries of the symbolic observation table constitute statements about the inclusion or exclusion of a large set of concrete words in the language. To prevent asking MQ's concerning all concrete words, only a small subset of words is chosen, the *evidences*; to avoid exponential growth of the table, only one of them, the *representative*, will be used for subsequent queries.

Handling multiple evidences to determine partitions is a crucial feature in learning automata over large alphabets. Evidences of the same partition should behave the same; if they do not, this causes evidence incompatibility. Finding and resolving evidence incompatibility is a major novel feature that is used in symbolic learning algorithms, as presented in this chapter.

However, parts of the algorithm can be adapted to the particularities of a specific learning problem, for the algorithm to become more efficient or result in a better hypothesis. For instance, the inferred model cannot be precise when the equivalence query is approximated. Furthermore, the input alphabet of the target language, or knowledge of the nature of the partitions, affects the way we learn transitions. Specialized methods to select evidences may result in more precise conjectures.

We are not the first to study the application of automaton learning to large alphabets. In order to make a comparison between related work and ours meaning-ful, we postpone it to Section 5.2, after the relevant notions have been introduced. In Section 5.1, we give formal definitions of symbolic observation tables and their properties, define evidences and representatives, and explain the notion of evidence incompatibility. The main algorithm and procedures follow and are presented in a general form in Section 5.3. The instantiation of the algorithm to specific alphabets and equivalence checking methods are discussed separately in Chapters 6 and 7, respectively.

#### 5.1 Definitions

To identify states and transitions and finally conjecture a hypothesis, the learner organizes all the examples and information that is available in a data-structure that is called the *symbolic observation table*. The rows of such a table correspond to symbolic words, which represent access sequences to the states; the columns consist of concrete words that are used to distinguish the states; and the entries of the table are instances of the symbolic characteristic function of the target language, which are incrementally filled in after posing MQ's to the teacher. Individual symbols represent transitions and need not be the same in each state. Each symbol has its own semantics, evidences and representatives.

Learning an automaton is like growing a spanning tree, where nodes stand for states and edges are transitions. With the initial state placed at the root of the tree, it is sufficient to grow the tree until all leaf nodes correspond to states that have been already reached via another path in the tree. In the following, we present such trees, which we call *balanced trees*, and relate them to the symbolic observation tables.

The structure of a balanced tree appears in Figure 5.1b together with its corresponding automaton in Figure 5.1c. The underlying intuition is that elements of S, also known as access sequences, correspond to a spanning tree of the transition graph of the automaton to be learned, while elements of the boundary R correspond to back- and cross-edges relative to this spanning tree.



Figure 5.1: (a) A symbolic observation table, (b) its balanced symbolic  $\Sigma$ -tree, and (c) the conjectured automaton.

Let  $\Sigma$  be a large input alphabet; let  $S \uplus R$  be the set of access sequences, a prefix-closed subset of  $\Sigma^*$ , where  $\Sigma$  is a symbolic alphabet such that  $\Sigma = \bigcup_{s \in S} \Sigma_s$ ; and, let  $\psi = \{\psi_s\}_{s \in S}$  be a family of total surjective functions of the form  $\psi_s : \Sigma \to \Sigma_s$ , that defines the semantics of the symbols. A *balanced* symbolic  $\Sigma$ -tree is a tuple  $(\Sigma, S, R, \psi)$ , where for every  $s \in S$  and  $a \in \Sigma_s$ ,  $s \cdot a \in S \cup R$ , and for any  $r \in R$  and  $a \in \Sigma$ ,  $r \cdot a \notin S \cup R$ . Elements of R are called *boundary elements* of the tree.

**Definition 5.1** (Symbolic Observation Table). A symbolic observation table is a tuple  $T = (\Sigma, \Sigma, S, R, \psi, E, f, \mu, \hat{\mu})$  such that

- $-\Sigma$  is an alphabet,
- $(\boldsymbol{\Sigma}, \boldsymbol{S}, \boldsymbol{R}, \psi)$  is a balanced symbolic  $\Sigma$ -tree,
- $E \subseteq \Sigma^*$  is a set of *distinguishing words*,
- $f: (S \cup R) \times E \rightarrow \{-,+\}$  is the symbolic classification function,
- $-\mu: \Sigma \to 2^{\Sigma} \{\emptyset\}$  is the evidence function, where  $\mu(a) \subseteq [a]$  for all  $a \in \Sigma$ ,
- $-\hat{\mu}: \Sigma \to \Sigma$  is the *representative function*, where  $\hat{\mu}(a) \in \mu(a)$  for all  $a \in \Sigma$ .

The evidence and representative functions are extended to symbolic words in  $S \cup R$  as follows:

$$\mu(\varepsilon) = \{\varepsilon\} \qquad \hat{\mu}(\varepsilon) = \varepsilon \mu(s \cdot a) = \hat{\mu}(s) \cdot \mu(a) \qquad \hat{\mu}(s \cdot a) = \hat{\mu}(s) \cdot \hat{\mu}(a).$$
(5.1)

The symbolic characteristic function values are based on the representative of the symbolic prefix rather than the set of all evidences, i.e., to fill the (s, e) entry in the table we let  $f(s, e) = f(\hat{\mu}(s) \cdot e)$  where f is the characteristic function of the target language. By this, an undesirable growth in the size of the queries and size

of the table is avoided. With every  $s \in S \cup R$  we associate a residual classification function  $f_s : E \to \{-,+\}$  defined as  $f_s(e) = f(s,e), e \in E$ .

The symbolic sample associated with T is the set  $M_T = (S \cup R) \cdot E$  while the concrete sample, the set of all words whose classification is known, is  $M_T = \mu(S \cup R) \cdot E$ , and coincides with the domain of f. Note that for each word  $w \in M_T$ there is at least one concrete  $w \in \mu(w)$  whose membership in L is known.

Let  $\mu_s = \bigcup_{a \in \Sigma_s} \mu(a)$  be the set of all evidences for a state  $s \in S$ . Evidences of the same symbol should behave the same and when this is not the case, that is, when two concrete letters in the evidence of a symbol lead to different residual functions, we call this a manifestation of *evidence incompatibility*. Evidence incompatibility can be characterized and measured as follows.

**Definition 5.2** (Incompatibility Instance). A state  $s \in S$  has an *incompatibility instance* at evidence  $a \in \mu_s$  when  $f_{\hat{\mu}(s) \cdot a} \neq f_{\hat{\mu}(s \cdot a)}$ , and  $\psi_s(a) = a$ ; this fact is denoted by INC(s, a), where

$$INC(\boldsymbol{s}, a) = \begin{cases} 1, & \text{if } f_{\hat{\mu}(\boldsymbol{s}) \cdot a} \neq f_{\hat{\mu}(\boldsymbol{s}) \cdot \hat{\mu}(\psi_{\boldsymbol{s}}(a))} \\ 0, & \text{otherwise} \end{cases}$$

**Definition 5.3** (Incompatibility Degree). The *evidence incompatibility degree* associated to a state  $s \in S$  is the total number of incompatibility instances that the state s has. We denote this by M(s), where  $M(s) = |\{a \in \mu_s : INC(s, a)\}|$ .

Naturally, the evidence incompatibility of a table T is the sum of the incompatibility degree over all its states  $s \in S$ , i.e.,  $M(T) = \sum_{s \in S} M(s)$ . Note that the incompatibility degree is bounded by the number of states times the maximum number of evidences per state.

**Definition 5.4** (Table Properties). A table  $T = (\Sigma, \Sigma, S, R, \psi, E, f, \mu, \hat{\mu})$  is

- Closed if  $\forall r \in R, \exists s \in S, f_r = f_s$ ,
- Reduced if  $\forall s, s' \in S, f_s \neq f_{s'}$ , and
- Evidence compatible if  $M(s) = 0, \forall s \in S$ .

The properties of the table, as defined above, are crucial to hold when one wants to conjecture a target language and build a symbolic automaton compatible with the concrete sample. A reduced table does not include any redundant prefixes and thus restricts the set of states to be exactly S. A closed table leads to a well defined transition function. Finally, evidences should be compatible and behave as their representatives for the automaton to be compatible with  $M_T$ .

The following result is the natural generalization of the derivation of an automaton from an observation table to the symbolic setting, see Theorem 3.3. Figure 5.1 illustrates the connection of the observation table to an automaton.

**Theorem 5.5** (Automaton from Table). *From a closed, reduced and evidence compatible table one can construct a deterministic symbolic automaton compatible with the concrete sample.* 

*Proof.* Let  $T = (\Sigma, \Sigma, S, R, \psi, E, f, \mu, \hat{\mu})$  be a symbolic observation table that is closed, reduced, and evidence compatible. We define a function  $g : R \to S$ , such that g(r) = s if and only if  $f_r = f_s$ . The automaton derived from the table is then  $\mathcal{A}_T = (\Sigma, \Sigma, \psi, Q, \delta, q_0, F)$  where:

$$-Q = S, q_0 = \varepsilon$$

$$-F = \{s \in S : f_s(\varepsilon) = +\}$$

$$-\delta : Q \times \Sigma \to Q \text{ is defined as } \delta(s, a) = \begin{cases} s \cdot a, & \text{when } s \cdot a \in S \\ g(s \cdot a), & \text{when } s \cdot a \in R \end{cases}$$

By construction and like in [Ang87],  $A_T$  classifies correctly via f the symbolic sample and, due to evidence compatibility, this classification agrees with the characteristic function f on the concrete sample.

#### 5.2 Comparison to Related Work

To start with, our work on symbolic automata should not be confused with work dealing with register automata, a popular extension of automata to infinite alphabets [KF94, BLP10, HSJC12]. These are automata augmented with additional variables that can store in registers some input letters encountered while reading a word. Newly-read letters can be compared with the registers but typically not with constants in the domain. With register automata one can express, for example, the requirement that the password at login is the same as the password at sign-up. In the most recent work on learning register automata [CHJS16], a strong tree oracle is used. Given a concrete prefix and a symbolic prefix, the teacher returns a special type of a register automaton that has a tree structure. This fills in the entries of the observation table and provides the information about the registers and guards in the automaton. This algorithm is efficient only in the presence of shortest counter-examples and, in addition, when applied to a theory of inequalities and extended to use constants, these constants should be known in advance.

Our approach, which uses symbolic automata, is different. First, no auxiliary variables are added to the automaton. The values of the input symbols 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  $c_1 \le x \le c_2$ , which give real (but of limited resolution) access to the input domain. On the other hand, we insist on 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.

Ideas similar to ours have been suggested and explored in a series of papers [BJR06, HSM11, IHS13] that also adapt automaton learning and Angluin's algorithm to large alphabets and are the closest conceptually to our work. These papers have some design decisions which are similar, for example, the usage of distinct symbolic alphabets at every state [IHS13], the notion of evidence or the adaptation of the breakpoint method [HSM11].

In contrast to our learning algorithm, the scheme presented in [BJR06] is always based on alphabet refinement and has the potential of generating new symbols indefinitely. On the other hand, the algorithms in [HSM11, IHS13] result in a partially defined hypothesis, where the transition function is not defined outside the observed evidence.

The main distinguishing feature in our framework is that it is more rigorous and comprehensive in the way it treats evidence incompatibility and the modification of partition boundaries. We do not consider each modification as a partition refinement, but rather try first to modify the boundaries and only add a new symbol when necessary. As a result, we have the following properties, whenever we conclude the treatment of evidence incompatibility. First, the number of symbolic letters is never larger than the minimal number needed; and second, the mapping of concrete letters to symbols is always sample-compatible and it is well-defined for the whole concrete alphabet.

Recently, new results presented in [DD17] in the context of learning symbolic automata give a more general justification for a learning scheme like ours by proving that learnability is closed under product and disjoint union.

Moreover, some similarities can also be expected between our symbolic algorithm and the  $S^*$  algorithm in [BB13]. However, this work seems to be incomparable to ours as they use a richer model of transducers as well as more general predicates on inputs and outputs. Their termination result is weaker and is closely related to the termination of the counter-example guided abstraction refinement procedure.

To conclude, we believe that our algorithm for learning languages over large alphabets is unique in employing all the following features:

1. It is based on a clean and general definition of the relation between the concrete and symbolic alphabets;

- 2. It treats the modification of alphabet partitions in a rigorous way which guarantees that no superfluous symbols are introduced;
- It accommodates for a teacher free setting, where counter-examples need not be minimal, and PAC learnability can be achieved;
- 4. It employs an adaptation of the breakpoint method to analyze in an efficient way the information provided by counter-examples;
- 5. It is modular, separating the general aspects from those that are alphabet specific, thus providing for a relatively easy adaptation to different alphabets.

#### 5.3 The Symbolic Learning Algorithm

In this section, we present the symbolic learning algorithm in a general and abstract way that does not assume any structure on the alphabet. Moreover, the type of equivalence checking is not specified. As one might notice, some of the procedures and methods cannot be fully specified in this general setting. The input alphabet, the nature of the partitions, and the type of equivalence oracle are all important elements that can affect the way transitions are learned, the efficiency of the algorithm, and sometimes, even its termination.

The main sampling-based symbolic learning algorithm, summarized in Algorithm 8, alternates between two phases. In the first phase it attempts to make the table closed (line 5) and evidence compatible (line 6) so as to construct a symbolic automaton (line 8). In the second phase, after formulating an equivalence query (line 9), it processes the provided counter-example (line 12), which renders the table not closed or evidence incompatible. These phases alternate until no counter-example is found. Note that the table, by construction, is always kept reduced. The procedures that appear in Algorithm 8 are explained in detail hereafter.

Whenever MQ is invoked for word w, then w is added to the concrete sample and the characteristic function f is updated such that f(w) = MQ(w).

Table Initialization (Procedure 9). The algorithm builds an initial observation table T, with  $\Sigma_{\varepsilon} = \{a\}, S = \{\varepsilon\}, R = \{a\}, E = \{\varepsilon\}$ . The newly introduced symbol a is initialized with concrete semantics, evidences and a representative, via the procedure INITSYMBOL which is invoked each time a new state is introduced. Then membership queries are posed to update f and fill the table.

Symbol Initialization (Procedure 10). For a new symbolic letter a, we let  $\llbracket a \rrbracket = \Sigma$  and as an evidence  $\mu(a)$  we take a set of k concrete letters, denoted by  $sample(\Sigma, k)$ . One element of the evidence, denoted by  $select(\mu(a))$ , is chosen as a representative and will be used to fill table entries for all rows where a appears.

Algorithm 8	Symbolic	Learning Algorithm
-------------	----------	--------------------

1:	learned = false	
2:	INITTABLE( $T$ )	
3:	repeat	
4:	while $T$ is not closed or not e	vidence compatible <b>do</b>
5:	CLOSE	
6:	EvComp	
7:	end while	
8:	Construct $\mathcal{A}_T$	
9:	if $EQ(\mathcal{A}_{T})$ then	$\triangleright$ check hypothesis $\mathcal{A}_{T}$
10:	learned = true	
11:	else	$\triangleright$ a counter-example w is provided
12:	$COUNTEREX(\mathcal{A}_T, w)$	▷ process counter-example
13:	end if	
14:	until learned	
15:	return $\mathcal{A}_T$	
Pro	cedure 9 Table Initialization	

1: **procedure** INITTABLE(T) 2:  $\Sigma_{\varepsilon} = \{a\}; S = \{\varepsilon\}; R = \{a\}; E = \{\varepsilon\}$   $\triangleright a$  is a new symbol 3: INITSYMBOL(a) 4: Ask MQ(u) for all  $u \in \mu(a) \cup \{\varepsilon\}$ 5:  $f(\varepsilon) = f(\varepsilon); f(a) = f(\hat{\mu}(a))$ 6:  $T = (\Sigma, \Sigma, S, R, \psi, E, f, \mu, \hat{\mu})$ 7: end procedure

The  $sample(\Sigma, k)$  or *select* methods can use a uniform distribution, but one can think of evidence or representative selection as a more adaptive process that may depend on the outcome of membership queries or other assumptions we have made, e.g., the minimality on the counter-examples, a case explained in more detail in Section 6.1, allows us to select as evidence and representative always the minimal element of a partition.

**Table Closing (Procedure 11).** A table is not closed when there exists some  $r \in R$  without any equivalent element  $s \in S$  such that  $f_r = f_s$ . To render the table closed, r should be considered as a new state. To this end, r is moved from R to S with alphabet  $\Sigma_r = \{a\}$ , where a is a new symbol which is initialized. To balance the table the new word  $r \cdot a$  is added to R, its evidence  $\mu(r \cdot a)$  and representative  $\hat{\mu}(r \cdot a)$  are computed following (5.1) and membership queries are

#### Procedure 10 Symbol Initialization

```
1: procedure INITSYMBOL(a)

2: \llbracket a \rrbracket = \Sigma

3: \mu(a) = sample(\Sigma, k)

4: \hat{\mu}(a) = select(\mu(a))
```

5: end procedure

#### Procedure 11 Table Closing

1:	procedure CLOSE	
2:	Given $m{r}\inm{R}$ such that $orallm{s}\inm{S},m{f_r} eqm{f_s}$	
3:	$oldsymbol{S} = oldsymbol{S} \cup \{oldsymbol{r}\}$	$\triangleright$ declare $r$ a new state
4:	${oldsymbol{\Sigma}}_{oldsymbol{r}}=\{a\}$	$\triangleright$ introduce a new symbol $a$
5:	INITSYMBOL $(a)$	
6:	$oldsymbol{R} = (oldsymbol{R} - \{oldsymbol{r}\}) \cup \{oldsymbol{r} \cdot oldsymbol{a}\}$	▷ add new boundary element
7:	Ask $MQ(u)$ for all $u \in \mu(\boldsymbol{r} \cdot \boldsymbol{a}) \cdot E$	
8:	$oldsymbol{f_{r\cdot a}} = f_{\hat{\mu}(oldsymbol{r\cdot a})}$	
9:	end procedure	

posed to update f and fill the table.

Fixing Evidence Incompatibility (Procedure 12). A table is not evidence compatible when the evidence incompatibility degree of a state  $s \in S$  is greater than zero. This happens when there exist a letter  $a \in \mu_s$  and a suffix  $e \in E$  such that  $f(\hat{\mu}(s) \cdot \hat{\mu}(\psi(a)) \cdot e) \neq f(\hat{\mu}(s) \cdot a \cdot e)$ . Evidence incompatibility can occur already at symbol initialization or as a result of adding a new evidence due to a counter-example treatment that is described in the sequel. Some elements of  $\mu(a)$  may behave differently from the representative and this will flag an evidence incompatibility condition to be treated subsequently.

To fix evidence incompatibility, the evidence function  $\mu$  and the semantics function  $\psi$  should be updated such that the total evidence incompatibility degree becomes zero. This is done incrementally and for each state *s* separately by consecutive calls to EVCOMP until the total incompatibility degree of the observation table becomes zero.

For a state s, an incompatibility instance at a indicates either that the partition boundary is imprecise or that a transition (and its corresponding symbol) is missing. Let  $s \in S$  be a state; let a be an evidence of the symbol  $a \in \Sigma_s$  which is incompatible, i.e.,  $a \in \mu(a)$  and  $INC(s, a) \ge 1$ . The algorithm searches for a symbol  $b \in \Sigma_s$  for which  $f(\hat{\mu}(s) \cdot \hat{\mu}(b) \cdot e) = f(\hat{\mu}(s) \cdot a \cdot e)$ . In the favorable case, where such a symbol b exists and belongs to the neighborhood of a, the partition

Procedure 12 Make Evidence Comp	oatible	
---------------------------------	---------	--

1:	procedure EvComp	
2:	Let $s \in S$ be a state with $\mathrm{M}(s)$	> 0
3:	repeat	
4:	Let $a \in \mu_s$ such that INC $(s,$	$(a)=1  ext{ and } \psi_{oldsymbol{s}}(a)=oldsymbol{a}$
5:	if there exists $oldsymbol{b}\in oldsymbol{\Sigma}_{oldsymbol{s}}$ such the	hat $f_{\hat{\mu}(\boldsymbol{s})\cdot\hat{\mu}(\boldsymbol{b})} = f_{\hat{\mu}(\boldsymbol{s})\cdot\boldsymbol{a}}$ then
6:	$\mu(oldsymbol{a})=\mu(oldsymbol{a})-\{a\};\mu(oldsymbol{b})$	$(oldsymbol{b})=\mu(oldsymbol{b})\cup\{a\}$
7:	else	
8:	$oldsymbol{\Sigma}_{oldsymbol{s}} = oldsymbol{\Sigma}_{oldsymbol{s}} \cup \{oldsymbol{b}\}$	$\triangleright$ introduce new symbol $b$
9:	$oldsymbol{R} = oldsymbol{R} \cup \{oldsymbol{s} \cdot oldsymbol{b}\}$	▷ add new boundary element
10:	$\mu(oldsymbol{b})=\{a\};\hat{\mu}(oldsymbol{b})=a$	
11:	end if	
12:	Let $\psi(a) = \boldsymbol{b}$	
13:	until $\mathbf{M}(\boldsymbol{s}) = 0$	
14:	Update $\psi_{s}$	
15:	end procedure	

and boundaries can be modified while preserving convexity of the partition blocks when defining the semantics; the letter a is moved from  $\mu(a)$  to  $\mu(b)$ . Otherwise, a new symbol b is introduced; the letter a becomes both evidence and representative of b and the incompatibility degree decreases by one. This is repeated for the next incompatible evidence, until M(s) = 0 for all  $s \in S$ .

Finally, the semantics are updated according to the new evidence function. This may depend on the type of the alphabet but a general approach is to define the semantics according to the closest evidence point. That is  $\psi(a) = a$ , where a is the symbol whose evidences are closest to a. In other words, we let  $\psi(a) = \arg \min_{a \in \Sigma_s} \min_{b \in \mu(a)} d(a, b)$ , where d denotes some distance measure over the concrete letters.

When the table becomes both closed and evidence compatible we construct from it the hypothesis automaton and pose an equivalence query.

**Processing Counter-Examples (Procedure 13).** A counter-example is a word w misclassified by the current hypothesis. The conjectured automaton should be modified to classify w correctly while remaining compatible with the evidence accumulated so far. These modifications can be of two major types that we call *vertical* and *horizontal*. The first type, which is the only possible modification in concrete learning, involves the discovery of a new state  $s \cdot a$ . A counter-example which demonstrates that some letter a took a wrong transition  $\delta(s, a)$  has an horizontal effect that fixes a transition or adds a new one. The procedure described in



(a) vertical expansion

(b) horizontal expansion

Figure 5.2: A counter-example expands a hypothesis either (a) vertically, discovering a new state; or (b) horizontally, modifying the alphabet partition in a state.

the sequel reacts to the counter-example by adding a to the evidence of s and thus modifying the table, which should then be made closed and evidence compatible before the learner continues with a new hypothesis.

Counter-examples are treated using a symbolic variant of the breakpoint method, which is described hereafter (for the classical breakpoint method, one can refer to Section 3.3.2). As a matter of fact, a counter-example w may be misclassified also in the new updated hypothesis due to multiple breakpoints. Hence, w can be used again as a counter-example until finding a hypothesis that classifies it correctly.

Let  $\mathcal{A}_T$  be a symbolic automaton derived from a symbolic table T, and let  $w = a_1 \cdots a_{|w|}$  be a counter-example whose symbolic image is  $a_1 \cdots a_{|w|}$ . An *i-factorization* of w is  $w = u_i \cdot a_i \cdot v_i$  such that  $u_i = a_1 \cdots a_{i-1}$  and  $v_i = a_{i+1} \cdots a_{|w|}$ . For every *i*-factorization of w, let  $u_i$  be the symbolic image of  $u_i$ , and  $s_i = \delta(\varepsilon, u_i \cdot a_i)$  be the symbolic state (an element of S) reached in  $\mathcal{A}_T$  after reading  $u_i \cdot a_i$ .

**Proposition 5.6** (Symbolic Breakpoint). If w is a counter-example to  $A_T$  then there exists an *i*-factorization of w such that either

$$f(\hat{\mu}(\boldsymbol{s}_{i-1}) \cdot \boldsymbol{a}_i \cdot \boldsymbol{v}_i) \neq f(\hat{\mu}(\boldsymbol{s}_{i-1}) \cdot \hat{\mu}(\boldsymbol{a}_i) \cdot \boldsymbol{v}_i)$$
(5.2)

or

$$f(\hat{\mu}(\boldsymbol{s}_{i-1} \cdot \boldsymbol{a}_i) \cdot v_i) \neq f(\hat{\mu}(\boldsymbol{s}_i) \cdot v_i)$$
(5.3)

*Proof.* Condition (5.2) states that  $a_i$  is not well represented by  $\hat{\mu}(a_i)$  while Condition (5.3) implies  $s_{i-1} \cdot a_i$  is a new state different from  $s_i$ , see Figure 5.2. We prove the proposition assuming that none of the inequalities above hold for any *i*-factorization of w. Then, by alternatively using the negations of (5.2) and (5.3) for all values of *i*, we conclude that  $f(\hat{\mu}(s_0) \cdot a_1 \cdot v_1) = f(\hat{\mu}(s_{|w|}))$ , where  $\hat{\mu}(s_0) \cdot a_1 \cdot v_1$  is the counter-example and  $s_{|w|}$  is the state reached in  $\mathcal{A}_T$  after reading w. Hence, w is not a counter-example, which is not true.

This argument is illustrated in more detail below. The actual counter-example  $w = a_1 \dots a_m$  is shown as an upper-script, such that the factorization of w is clear at each step.

$$f(w) = f(\hat{\mu}(\hat{s}_0) \cdot \hat{a}_1 \cdot \hat{v}_1^{a_2...a_m}) \qquad \stackrel{(5.2)}{=} f(\hat{\mu}(\hat{s}_0) \cdot \hat{\mu}(\hat{a}_1) \cdot \hat{v}_1^{a_2...a_m}) \qquad \stackrel{(5.3)}{=}$$

$$= f(\hat{\mu}(\mathbf{s}_{1}) \cdot \overset{a_{2}}{a_{2}} \cdot \overset{a_{3}...a_{m}}{v_{2}}) \qquad \stackrel{(5.2)}{=} f(\hat{\mu}(\mathbf{s}_{1}) \cdot \hat{\mu}(\mathbf{a}_{2}) \cdot \overset{a_{3}...a_{m}}{v_{2}}) \qquad \stackrel{(5.3)}{=}$$
  
$$\vdots$$

$$= f(\hat{\mu}(\mathbf{s}_{i-1}) \cdot \hat{a}_{i} \cdot \hat{v}_{i}) \stackrel{(5.2)}{=} f(\hat{\mu}(\mathbf{s}_{i-1}) \cdot \hat{\mu}(\mathbf{a}_{i}) \cdot \hat{v}_{i}) \stackrel{(5.3)}{=}$$

$$= f(\hat{\mu}(\boldsymbol{s}_{m-1}) \cdot \hat{a}_{m}^{a} \cdot \hat{v}_{m}^{\varepsilon}) \qquad \stackrel{(5.2)}{=} f(\hat{\mu}(\boldsymbol{s}_{m-1}) \cdot \hat{\mu}(\boldsymbol{a}_{m}) \cdot \hat{v}_{m}) \qquad \stackrel{(5.3)}{=} \\ = f(\hat{\mu}(\boldsymbol{s}_{m}))$$

Let us consider a hypothesis automaton  $\mathcal{A}_T$  for which  $w \in \Sigma^*$  is given as a counterexample. We use the symbolic breakpoint method to treat this counter-example and update the symbolic table accordingly. As we can see in Procedure 13, the learner iterates over *i* values and checks whether one of the conditions (5.2) and (5.3) holds for some *i*-factorization of the counter-example. We can choose *i* in any order, as described in Section 3.3.2, either to retain smaller prefixes or suffixes, or to find a breakpoint faster. We let *i* take values in a monotonically descending order and keep the suffixes as short as possible. In this case, it suffices to compare  $f(\hat{\mu}(s_{i-1} \cdot a_i) \cdot v_i)$  and  $f(\hat{\mu}(s_{i-1}) \cdot a_i \cdot v_i)$  with the classification of the counterexample, which is kept in a flag variable. When we iterate (line 3) in a descending order, the flag value (line 2) should be set to flag = f(w), that is, the classification of the counter-example by the Teacher.

The two main conditions of the breakpoint method are checked in lines 5 and 9. Specifically, the learner checks first condition (5.3) in line 5, where its satisfaction implies that a new distinguishing string has been found. Adding  $v_i$  to E will distinguish between states  $s_{i-1} \cdot a_i$  and  $s_i$ , resulting a non closed table. Otherwise, if condition (5.2) holds, a check made in line 9, it means that the letter  $a_i$  does not

1:	<b>procedure</b> COUNTEREX( $\mathcal{A}_{T}, w$ )	
2:	$ ext{flag} = f(\hat{\mu}(oldsymbol{\delta}(oldsymbol{arepsilon},oldsymbol{w}))) \hspace{1.5cm}  riangle  ext{flag} =$	$f(w)$ when iterating over $1, \ldots,  w $
3:	for $i= w ,\ldots,1$ do	
4:	For an <i>i</i> -factorization $w = u_i \cdot a_i \cdot v_i$	
5:	if $f(\hat{\mu}(\boldsymbol{s}_{i-1}\cdot\boldsymbol{a}_i)\cdot v_i) \neq flag$ then	$\triangleright$ check (5.3)
6:	$E = E \cup \{v_i\}$	> add $v_i$ as a new distinguishing word
7:	Ask mq $(u)$ for all $u \in \mu(oldsymbol{S} \cup oldsymbol{R}) \cdot v_i$	
8:	break	
9:	else if $f(\hat{\mu}(\boldsymbol{s}_{i-1}) \cdot a_i \cdot v_i) \neq \text{flag}$ then	⊳ check (5.2)
10:	$\mu(\boldsymbol{a}_i) = \mu(\boldsymbol{a}_i) \cup \{a_i\}$	$\triangleright$ add $a_i$ as new evidence
11:	Ask MQ $(u)$ for all $u \in \hat{\mu}(\boldsymbol{s}_{i-1}) \cdot a_i$ .	E
12:	if $\mathbf{M}(\boldsymbol{s}_i) = 0$ then	
13:	$E = E \cup \{v_i\} \qquad \triangleright z$	add distinguishing word when needed
14:	Ask MQ $(u)$ for all $u \in \mu(oldsymbol{S} \cup oldsymbol{R})$	$) \cdot v_i$
15:	end if	
16:	break	
17:	end if	
18:	end for	
19:	end procedure	

behave as its representative in the hypothesis. As a result, the letter  $a_i$  is added to the evidence of  $a_i$  and new membership queries are posed to fill in the table. These queries will render the table evidence incompatible and will lead to refining of  $[a_i]$ . If no evidence incompatibility occurs after introducing the new evidence and filling in the table, this is due to a missing distinguishing string, which is  $v_i$  that is not yet part of the table. In such a case, also the suffix  $v_i$  is added to E, as it is the only witness for the incompatibility, see lines 12-15. The procedure terminates when either condition (5.2) or condition (5.3) is satisfied (line 9 and line 5).

Note here that checking the conditions involves supplementary membership queries that are based on the suffix of the counter-example w where the prefix  $u_i$  of w is replaced by  $\hat{\mu}(s_{i-1})$ , the representative of its shortest equivalent symbolic word in the table. Both cases will lead to a new conjectured automaton that might still not classify w correctly. Hence, if w remains a counter-example for the new conjecture, Procedure 13 should be invoked with the same counter-example until the new hypothesis  $\mathcal{A}_T$  classifies w correctly.

## Learning with a Helpful Teacher

In this chapter, we present an instantiation of the learning algorithm that uses the strong assumption of a helpful teacher. Such a teacher responds positively to an equivalence query  $EQ(\mathcal{A})$ , where  $\mathcal{A}$  is an automaton conjectured by the learning algorithm, only if  $L(\mathcal{A})$  is indeed equivalent to the target language; otherwise, it returns a *minimal counter-example* which is a counter-example of minimal length which is also minimal with respect to a lexicographic order. A *minimal counterexample* helps the learner to localize the modification site. This strong assumption improves the performance of the algorithm resulting in an exact conjecture; its relaxation is discussed in Chapter 7. In the following sections, we discuss the learning problem for the cases of totally and partially-ordered input alphabets, respectively.

#### 6.1 Learning Languages over Ordered Alphabets

In this section, we assume  $\Sigma$  to be a *totally-ordered* alphabet with minimum and restrict ourselves to interval symbolic automata (see Section 4.4.1). In these automata the concrete alphabet is any left-closed subset of  $\mathbb{N}$ ,  $\mathbb{Q}$  or  $\mathbb{R}$ , and the semantics, in the case of a dense order as in  $\mathbb{R}$ , form closed-open intervals.

The order on the alphabet is extended naturally to a lexicographic order on  $\Sigma^*$ . The minimality of the counter-example, in length and lexicographically, improves the learning procedure. The treatment of counter-examples becomes cleaner and learning the semantics is accurate. Minimality allows us to use fewer evidences, keep the concrete sample as small as possible, and avoid posing unnecessary membership queries.

In Section 5.3 we have seen that new symbolic letters are introduced on 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  $[\![a]\!]$  to the largest possible subset of  $\Sigma$ , given the current evidence.

The assumption of a minimal counter-example allows us to specialize some of the procedures of Chapter 5, which results in fewer evidences and a cleaner and cheaper counter-examples treatment.

For each symbol a in  $\Sigma$  we choose exactly one letter as evidence, that is the smallest possible  $a \in [\![a]\!]$ . As this minimal element is the only evidence, we choose it also as the representative. Let us denote by  $a_0$  the minimal element of the input alphabet,  $a_0 = \min \Sigma$ . We adapt Procedure 10, the initialization of a new symbol, to this setting. We replace  $sample(\Sigma, k)$  by  $\{a_0\}$ ; the evidence, representative, and semantics functions are updated such that, for each symbol  $a \in \Sigma$ , the following holds:

$$\mu(\boldsymbol{a}) = \{\hat{\mu}(\boldsymbol{a})\} = \{\min[\boldsymbol{a}]\}.$$
(6.1)

A misclassified word in the conjectured automaton is due to a wrong transition that has been taken. This happens in two cases: either some evidence is missing, or the target state has not been discovered yet. Unlike the general case described in Section 5.3, minimality assures that the counter-example uses the representatives, which are the smallest elements of each transition, to reach states in the automaton through a loop-free path. After detecting the state at which the hypothesis is problematic, the missing state or a transition is added appropriately. Let w be a counter-example, since this is minimal it admits an *i*-factorization  $w = u_i \cdot a_i \cdot v_i$ , with  $u_i$  being the largest prefix of w such that  $u_i \in \mu(u_i)$  for some  $u_i \in S \cup R$ but  $u_i \cdot a_i \notin \mu(u')$  for any prefix u' in  $S \cup R$ . We consider two cases,  $u_i \in S$  and  $u_i \in R$ .

In the first case, when  $u_i$  is already a state in the hypothesis, i.e.,  $u_i \in S$ , the letter  $a_i$  indicates that the partition boundaries are not correctly defined and need refinement. Letter  $a_i$  is wrongly considered to be part of  $[\![a_i]\!]$ ; a new symbol b is added to  $\Sigma_{u_i}$ , with  $\hat{\mu}(b) = a_i$ , and solves the evidence incompatibility. Due to minimality of the counter-example, all letters in  $[\![a_i]\!]$  less than letter  $a_i$  behave like  $\hat{\mu}(a_i)$ . We assume that all remaining letters in  $[\![a_i]\!]$  behave like  $a_i$  and map them to b. The semantics function  $\psi_{u_i}$  is updated such that  $\psi_{u_i}(a) = b$  for all  $a \in [\![a_i]\!]$ ,  $a \ge a_i$ , and  $\psi_{u_i}(a)$  stays unchanged otherwise. Finally, the prefix  $u_i \cdot b$  is added to R.

In the second case, when the symbolic word  $u_i$  is part of the boundary,  $u_i \in \mathbf{R}$ ; we infer that  $u_i$  is not equivalent to any of the existing states in the hypothesis. Let  $s \in S$  be the prefix that was considered to be equivalent to  $u_i$  for which  $f_{u_i} = f_s$ . Since the table is reduced s is unique. Because w is minimal, the classification of words in  $[\![s]\!] \cdot a_i \cdot v_i$  in the hypothesis automaton is correct; otherwise, there exists some  $s \in [\![s]\!]$  such that  $s \cdot a_i \cdot v_i$  constitutes a shorter counter-example. From this we

<b>Procedure 14</b> Counter-Example Treatment (with Helpful Teacher) - $\mathbb{R}$			
1:	procedure COUNTEREX(w)		
2:	Find an <i>i</i> -factorization $w = u_i \cdot a_i \cdot v_i, a_i \in \Sigma, u_i, v_i \in \Sigma^*$ such that		
3:	$\exists oldsymbol{u}_i \in oldsymbol{S} \cup oldsymbol{R},  u_i \in \mu(oldsymbol{u}_i)  ext{ and } u_i \cdot a_i \notin \mu(oldsymbol{u'}), orall oldsymbol{u'} \in oldsymbol{S} \cup oldsymbol{R}$		
4:	if $oldsymbol{u}_i\inoldsymbol{S}$ then		
5:	Let $a_i \in \Sigma_{u_i}$ such that $a_i \in [\![a_i]\!]$ $\triangleright$ refinement of $[\![a_i]\!]$		
6:	$\Sigma_{u_i} = \Sigma_{u_i} \cup \{b\}$ bintroduce new symbol b		
7:	$\mu(oldsymbol{b})=\{a_i\}\ ; \hat{\mu}(oldsymbol{b})=a_i$		
8:	$\llbracket \boldsymbol{b} \rrbracket = \{ a \in \llbracket \boldsymbol{a}_i \rrbracket : a \geq a_i \}$		
9:	$\llbracket \boldsymbol{a}_i \rrbracket = \llbracket \boldsymbol{a}_i \rrbracket \setminus \llbracket \boldsymbol{b} \rrbracket$		
10:	$\boldsymbol{R} = \boldsymbol{R} \cup \{\boldsymbol{u}_i \cdot \boldsymbol{b}\}$		
11:	Ask $MQ(u)$ for all $u \in \mu(\boldsymbol{u}_i \cdot \boldsymbol{b}) \cdot E$		
12:	else $ ho$ if $oldsymbol{u}_i \in oldsymbol{R}$		
13:	$S = S \cup \{u_i\}; E = E \cup \{a_i \cdot v_i\}$ $\triangleright u_i$ becomes a new state		
14:	if $a_i = a_0$ then		
15:	$oldsymbol{\Sigma}_{oldsymbol{u}_i} = \{oldsymbol{b}\}$		
16:	$\mu(oldsymbol{b}) = \{a_0\}; \hat{\mu}(oldsymbol{b}) = a_0; \llbracket oldsymbol{b}  rbracket = \Sigma$		
17:	else		
18:	$oldsymbol{\Sigma}_{oldsymbol{u}_i} = \{oldsymbol{b},oldsymbol{b}'\}$		
19:	$\mu(m{b}) = \{a_0\}; \hat{\mu}(m{b}) = a_0 \ ; \llbracket m{b}  rbracket = \{a \in \Sigma : a < a_i\}$		
20:	$\mu(\boldsymbol{b}') = \{a_i\}; \hat{\mu}(\boldsymbol{b}') = a_i; \llbracket \boldsymbol{b}' \rrbracket = \{a \in \Sigma : a \ge a_i\}$		
21:	end if		
22:	$oldsymbol{R} = (oldsymbol{R} - \{oldsymbol{u}_i\}) \cup oldsymbol{u}_i \cdot oldsymbol{\Sigma}_{oldsymbol{u}_i}$		
23:	Ask MQ $(u)$ for all $u \in \mu(oldsymbol{u}_i \cdot oldsymbol{\Sigma}_{oldsymbol{u}_i}) \cdot E$		
24:	end if		
25:	25: end procedure		

conclude that  $u_i$  is a new state and the hypothesis automaton should be expanded vertically moving  $u_i$  to S. To distinguish between  $u_i$  and s, the suffix  $a_i \cdot v_i$  is added to E. The letter  $a_i$  is an evidence for the new state  $u_i$  and, for 5.1 to hold, we distinguish the following two sub-cases. If  $a_i$  is the smallest element of  $\Sigma$ , i.e.,  $a_i = a_0$ , one symbolic letter b is initialized and added to  $\Sigma_{u_i}$ . Otherwise, if  $a_i$  is not the smallest and  $a_i > a_0$ , two new symbolic letters, b and b', are added to  $\Sigma_{u_i}$ with  $[\![b]\!] = \{a : a < a_i\}, [\![b']\!] = \{a : a \ge a_i\}, \hat{\mu}(b) = a_0$ , and  $\hat{\mu}(b') = a_i$ . The prefixes  $u \cdot \Sigma_{u_i}$  are added to R.

The counter-example treatment for ordered alphabets and minimal counterexamples is summarized in Procedure 14. This procedure significantly differs from the general Procedure 13 discussed in the previous section. More precisely, evidence incompatibility appears only during the counter-example treatment where it is instantly solved. Hence, when the Procedure 13 is replaced by Procedure 14,



Figure 6.1: Minimal automaton accepting the target language used in Example 6.1.

then EVCOMP in Algorithm 8 is never invoked.

However, one can use Procedure 13 instead of Procedure 14, with one minor modification in line 2, where the iteration order should be changed to increasing, and the flag should be initialized to f(w). In such a case, the procedure EVCOMP should be called whenever necessary. Although the final conjecture will be the same, the algorithm becomes less efficient in the sense that more repetitions in the main algorithm are needed to reach the same result. Moreover, the algorithm calls COUNTEREX (Procedure 13) more times, and a number of additional MQ's are posed in order to perform the tests in lines 5 and 9. Note that MQ's are not required at all when the counter-example is treated using Procedure 14.

**Example 6.1** (Learning Interval Symbolic Automata with Minimal Counter-Examples). Let the concrete alphabet be  $\Sigma = [0, 100) \subset \mathbb{R}$  equipped with the usual order, and let  $L \subseteq \Sigma^*$  be a target language whose automaton is shown in Figure 6.1. Figure 6.2 shows the evolution of the symbolic observation tables and Figure 6.3 depicts the corresponding automata and the concrete semantics of the symbolic alphabets after each update.

The learner initializes the table by letting  $S = \{\varepsilon\}$ ,  $\Sigma = \{a_0\}$ ,  $R = \{a_0\}$ ,  $\mu(a_0) = \{0\}$ , and  $E = \{\varepsilon\}$ . The initial observation table  $T_0$  is filled in after membership queries are posed for  $\varepsilon$ , which is rejected, and for 0, which is accepted. To close the table, the learner makes  $a_0$  a new state and introduces  $\Sigma_{a_0} = \{a_1\}$ , where  $a_1$  is a new symbol with  $\mu(a_1) = \{0\}$ . The symbolic word  $a_0 \cdot a_1$  is then added to R. A membership query MQ( $0 \cdot 0$ ) fills the table. The new observation table is  $T_1$  and it is closed. The symbolic automaton  $\mathcal{A}_1$  derived from  $T_1$  has two states, one accepting and one rejecting.

The learner poses an equivalence query for  $A_1$  and obtains (50, -) as a minimal counter-example. The counter-example's minimality implies that all one-letter words smaller than 50 are correctly classified, in contrast to words greater than


Figure 6.2: Observation tables used in Example 6.1.

50. For this reason, the symbol  $a_0$  is not a good representation for all  $\Sigma$ , and a new symbol  $a_2$  is added to  $\Sigma_{\varepsilon}$ . The concrete semantics for this state is updated to  $[\![a_0]\!] = \{a \in \Sigma : a < 50\}$  and  $[\![a_2]\!] = \{a \in \Sigma : a \ge 50\}$ . As an evidence the smallest possible letter is selected, i.e.,  $\mu(a_2) = \{50\}$ . The learner asks membership queries, fills in the table, which is closed and appears in Figure 6.2 as  $T_2$ , and constructs the new hypothesis  $\mathcal{A}_2$ .

As the symbolic automaton  $\mathcal{A}_2$  is not the correct one, the teacher returns the counter-example  $(0 \cdot 30, -)$ . The prefix 0 already exists in the sample, which means that the misclassification must occur in the second transition, that is, at  $a_1$ . Therefore, the alphabet partition for state  $a_0$  is refined by introducing a new symbol  $a_3$  and by letting  $[a_1]] = \{a \in \Sigma : a < 30\}$  and  $[a_3]] = \{a \in \Sigma : a \ge 30\}$ . The new table  $T_3$  is closed, and the new hypothesis automaton is  $\mathcal{A}_3$ . Again, this hypothesis automaton does not recognize the target language, and  $(50 \cdot 0, -)$  is a new minimal counter-example. The prefix 50 belongs to the evidence of  $a_2$ ; this signifies that  $a_2$  goes to a state that is different from  $\varepsilon$ . The suffix 0 of the counter-example is added to E to distinguish the two states. The learner moves  $a_2$  to the set of states and introduces the set  $\Sigma_{a_2} = \{a_4\}$ , where  $a_4$  is a new symbol with



Figure 6.3: Hypotheses and  $\Sigma$ -semantics as learned in Example 6.1.

 $\mu(a_4) = \{0\}$ . Then,  $a_2 \cdot a_4$  is added to R. After all updates, the new table  $T_4$  is made closed, resulting in  $T_5$  and the conjectured automaton  $A_5$ , which has two new states.

In the following, three more counter-examples  $(50 \cdot 20, +)$ ,  $(50 \cdot 80, -)$  and  $(50 \cdot 50 \cdot 0, +)$  are provided. Progressively, each of them expands horizontally

the hypothesis by adding a new transition. The automata, after each update, are  $\mathcal{A}_6$ ,  $\mathcal{A}_7$  and  $\mathcal{A}_8$ , respectively. The learner terminates with the last automaton  $\mathcal{A}_8$  accepting the target language.

At this point, it is important to note that the language L that is defined over a subset of  $\mathbb{R}$  cannot be learned by the classical  $L^*$  algorithm. To be able to compare it with this algorithm, we restrict the problem to a finite discrete alphabet, such that  $\Sigma = \{1, \ldots, 100\}$ , and let the target language be recognized from the same symbolic automaton. The  $L^*$  algorithm requires approximately 1000 queries to learn L, in contrast to a total of 17 queries that are required by the symbolic algorithm.  $\Box$ 

# 6.2 Learning over Partially-ordered Alphabets

Let  $\Sigma$  be a partially-ordered alphabet such as a bounded sub-rectangle of  $\mathbb{R}^n$ , i.e.,  $\Sigma = X^d$ , where X is a totally-ordered set such as a sub-interval  $[0, k) \subseteq \mathbb{R}$ . The target language L is accepted by a canonical symbolic automaton where partitions are monotone, as described in Section 4.4.2. Monotone partitions of this kind are defined using finite unions of forward cones, and equivalently, a partition block can be determined by a finite set of minimum mutually incomparable letters from  $\Sigma$ . We denote a monotone partition by  $\mathcal{P} = \{P_1, \ldots, P_{m-1}\}$ , where  $P_i = B^+(F_{i-1}) - B^+(F_i)$  and  $B^+(F)$  is a forward cone defined by a set of incomparable points  $F \subseteq \Sigma$ .

The partial order on the alphabet extends to a partial lexicographic order over  $\Sigma^*$ . The teacher is assumed to return a minimal counter-example, which in this case is chosen from a finite set of incomparable minimal counter-examples. As for totally-ordered alphabets, the minimality of the counter-examples improves the learning algorithm, resulting in fewer membership queries for counter-example treatment, and to a precise final hypothesis. Evidences derived from counter-examples denote the minimal points of a partition block. Hence, multiple evidences are used to determine each partition block. In the following, we present the symbolic learning algorithm adapted to this special nature of the partitions. We mainly describe the modifications that should be considered to handle the evidences and solve incompatibility whenever this occurs.

The main algorithm (Algorithm 8) is applied and the procedure of table initialization, symbol initializations, and table closure (see Procedures 9, 10 and 11) remain unchanged.

Counter-examples either discover new states or refine partitions. When refinement occurs, minimality of the counter-examples constrains the teacher to provide all the counter-examples needed to find the minimal points that completely define **Procedure 15** Counter-Example Treatment (with Helpful Teacher) -  $\mathbb{R}^n$ 

1:	procedure COUNTEREX(w)
2:	Find a <i>i</i> -factorization $w = u_i \cdot a_i \cdot v_i, a_i \in \Sigma, u_i, v_i \in \Sigma^*$ such that
3:	$\exists m{u}_i \in m{S} \cup m{R},  u_i \in \mu(m{u}_i)  ext{ and } u_i \cdot a_i \notin \mu(m{u'}), orall m{u'} \in m{S} \cup m{R}$
4:	if $oldsymbol{u}_i\inoldsymbol{S}$ then
5:	Let $\boldsymbol{a}_i \in \boldsymbol{\Sigma}_{\boldsymbol{u}_i}$ such that $a_i \in \llbracket \boldsymbol{a}_i  rbracket$
6:	Let $a \in \Sigma_{u_i}$ such that $a > a_i$ adjacent to $a_i \qquad \qquad \triangleright a \in \mathcal{V}(a_i)$
7:	if $a_i    a$ for all $a \in \mu(a)$ and $f_{u_i \cdot a_i} = f_{u_i \cdot \hat{\mu}(a)}$ then
8:	$\mu(\boldsymbol{a}) = \mu(\boldsymbol{a}) \cup \{a_i\}$ $\triangleright$ change boundary
9:	$\llbracket \boldsymbol{a} \rrbracket = \llbracket \boldsymbol{a} \rrbracket \cup \{ a \in \llbracket \boldsymbol{a}_i \rrbracket : a \geq a_i \}$
10:	$\llbracket oldsymbol{a}_i  rbracket = \llbracket oldsymbol{a}_i  rbracket \setminus \{a \in \Sigma: a \geq a_i\}$
11:	else > partition refinement
12:	$\Sigma_{\boldsymbol{u}_i} = \Sigma_{\boldsymbol{u}_i} \cup \{\boldsymbol{b}\}$ $\triangleright$ introduce new symbol $\boldsymbol{b}$
13:	$\mu(oldsymbol{b})=\{a_i\}$ ; $\hat{\mu}(oldsymbol{b})=a_i$
14:	$\llbracket \boldsymbol{b} \rrbracket = \{ a \in \llbracket \boldsymbol{a}_i \rrbracket : a \geq a_i \}$
15:	$[\![\boldsymbol{a}_i]\!] = [\![\boldsymbol{a}_i]\!] \setminus [\![\boldsymbol{b}]\!]$
16:	$\boldsymbol{R} = \boldsymbol{R} \cup \{\boldsymbol{u}_i \cdot \boldsymbol{b}\}$
17:	Ask $MQ(u)$ for all $u \in \mu(\boldsymbol{u}_i \cdot \boldsymbol{b}) \cdot E$
18:	end if
19:	else $\triangleright$ if $oldsymbol{u}_i \in oldsymbol{R}$
20:	$S = S \cup \{u_i\}; E = E \cup \{a_i \cdot v_i\}$ $\triangleright u_i$ becomes a new state
21:	if $a_i = a_0$ then
22:	$\Sigma_{u_i} = \{b\}$ $\triangleright$ introduce new symbol $b$
23:	$\mu(oldsymbol{b}) = \{a_0\}; \hat{\mu}(oldsymbol{b}) = a_0; \llbracket oldsymbol{b}  rbracket = \Sigma$
24:	else
25:	$\Sigma_{\boldsymbol{u}_i} = \{\boldsymbol{b}, \boldsymbol{b}'\}$ $\triangleright$ introduce new symbols $\boldsymbol{b}$ and $\boldsymbol{b}'$
26:	$\mu(oldsymbol{b}) = \{a_i\}; \hat{\mu}(oldsymbol{b}) = a_i; \llbracket oldsymbol{b}  rbracket = \{a \in \Sigma : a \ge a_i\}$
27:	$\mu(oldsymbol{b}') = \{a_0\}; \hat{\mu}(oldsymbol{b}') = a_0; \llbracketoldsymbol{b}' rbracket = \Sigma \setminus \llbracketoldsymbol{b} rbracket$
28:	end if
29:	$oldsymbol{R} = (oldsymbol{R} - \{oldsymbol{u}_i\}) \cup oldsymbol{u}_i \cdot oldsymbol{\Sigma}_{oldsymbol{u}_i}$
30:	Ask $MQ(u)$ for all $u \in \mu(\boldsymbol{u}_i \cdot \boldsymbol{\Sigma}_{\boldsymbol{u}_i}) \cdot E$
31:	end if
32:	end procedure

the new partition. Even though the counter-examples are provided in an arbitrary order, defined by the teacher, this does not affect the behavior of the algorithm.

Let  $s \in S$  be a state in the observation table T, and let a be a symbol in  $\Sigma_s$ . The set of evidences of a is the set of the mutually-incomparable minimum elements belonging to its semantics. These evidences are sufficient to define each partition block as the set difference between two unions of cones, each cone defined by a letter  $a \in \mu(a)$ . When s is a new state, we let  $\Sigma_s = \{a\}$ , where a is a new



Figure 6.4: Modifying the alphabet partition for state  $u_i$  after receiving  $u_i \cdot a_i \cdot v_i$  as counter-example. Letters greater than  $a_i$  are moved from  $[\![a_i]\!]$  to  $[\![a]\!]$ .

symbol which is initialized; the concrete semantics  $[\![a]\!]$  is set to  $\Sigma$  and the minimal element  $a_0 = 0$  serves as evidence and representative. Evidences are added later, given by counter-examples. New evidences refine the partition either by adding a new symbol, or by modifying the set of cones for an existing partition block. Symbols may have many evidences, but keep as representative the first evidence admitted.

Procedure 14, that treats the counter-examples for ordered alphabets, should be modified in the case where the counter-example only modifies a partition boundary, without adding new states or transitions. For this an additional condition should be checked. The prefix  $u_i$  belongs to a state, but the new letter  $a_i \in [\![a_i]\!]$ , behaves like the adjacent symbol  $a > a_i$ . There is no need of a new symbol and the evidence incompatibility is solved by moving  $a_i$  from  $[\![a_i]\!]$  to  $[\![a]\!]$  along with all letters greater or equal to it. The letter  $a_i$  is then added to the evidences of a. Figure 6.4 shows this update of the semantics. The full modification of the procedure appears in Procedure 15. In the lines 8-10, the learning handles the case of boundary modification. The rest of the procedure remains unchanged, and one or two new symbols are introduced to refine the existing partition block that is evidence incompatible.

**Example 6.2.** Let us illustrate the working of the algorithm on a target language L defined over  $\Sigma = [0, 100]^2$ . We assume partitions to be monotone. All resulting tables, hypotheses automata and alphabet partitions for this example are shown in Figures 6.5, 6.6, and 6.7, respectively.

The learner starts by asking a membership query about the empty word. A symbolic letter  $a_0$  is chosen to represent the continuations from the initial state, initially placed in the boundary. The learner chooses the minimal element of  $\Sigma$  as evidence, i.e.,  $\Sigma_{\varepsilon} = \{a_0\}, \mu(a_0) = \{\binom{0}{0}\}, \text{ and } [\![a_0]\!] = \Sigma$ . The table is not closed, and for this, the learner declares  $a_0$  a state, introduces  $\Sigma_{a_0} = \{a_1\}$  with  $\mu(a_1) = \{\binom{0}{0}\}, [\![a_1]\!] = \Sigma$ , and  $R = \{a_0 \cdot a_1\}$ . The resulting table  $T_0$  is now closed; the first hypothesis  $\mathcal{A}_0$  is constructed.



Figure 6.5: Observation tables for Example 6.2.

The counter-example  $\binom{45}{50}, -)$  arrives, to refine the partition at the initial state by introducing a new symbol, and thus to add a new transition to the automaton. The symbolic alphabet is extended to  $\Sigma_{\varepsilon} = \{a_0, a_2\}$  with  $\llbracket a_2 \rrbracket = \{a \in \Sigma : a \ge \binom{45}{50}\}$ ,  $\llbracket a_0 \rrbracket = \Sigma - \llbracket a_2 \rrbracket$ , and  $\mu(a_2) = \{\binom{45}{50}\}$ . The new observation table and hypothesis are  $T_1$  and  $\mathcal{A}_1$ , respectively. Two more counter-examples, namely  $\binom{60}{0}, -)$  and  $\binom{0}{70}, -)$ , are provided and they refine the partition at the initial state causing a boundary modification. No new state or symbol is added to the automaton, which leaves the structure of the automaton unchanged, see Figure 6.6. However, the semantics function for the initial state updates to  $\psi_1, \psi_2$  and  $\psi_3$  (Figure 6.7), where all letters greater than  $\binom{60}{0}$  and  $\binom{0}{70}$  are moved to the  $\Sigma$ -semantics of  $a_2$ . An equivalence query on hypothesis  $\mathcal{A}_3$  will result in the counter-example  $\binom{0}{0}\binom{0}{80}, -)$ , which, in turn, adds a new symbol  $a_3$  to  $\Sigma_{a_0}$  and a new transition in the hypothesis.

The three 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  $\Sigma$ -semantics for the symbols in  $\Sigma_{a_0}$  as shown in  $\psi_4$ though  $\psi_7$ . The next counter-example  $\binom{45}{50}\binom{0}{0}$ , +) discovers a new state. Its prefix  $\binom{45}{50}$  already exists in  $\mu(a_2)$  and  $a_2 \in \mathbf{R}$ , indicating that the prefix  $a_2$  is a distinct state. To distinguish it from the initial state, the learner adds the suffix  $\binom{0}{0}$  to E. The resulting table  $T_8$  is not closed. To render it closed, the prefix  $a_0a_1$  moves to S. From the resulting table  $T_9$ , which is both closed and evidence compatible,



Figure 6.6: Hypothesis automata for Example 6.2.

we conclude the hypothesis  $\mathcal{A}_9$ . This automaton has four states. The updated  $\Sigma$ semantics for each state can be seen in  $\psi_9$ . Partitions at state  $a_2$  are refined further
due to some more counter-examples, see  $\mathcal{A}_{10-18}$  and  $\psi_{10} - \psi_{18}$  for the updates.
The algorithm terminates with  $\mathcal{A}_{18}$  as its final hypothesis, after using a total of 20
queries and treating 17 counter-examples.



Figure 6.7:  $\Sigma$ -semantics as updated in Example 6.2. The rows correspond to states  $\varepsilon$ ,  $a_0$ ,  $a_2$  and  $a_0 \cdot a_1$ , respectively. The function  $\psi_i$  refers to the hypothesis  $\mathcal{A}_i$ .

# Learning without a Helpful Teacher

In this chapter, we relax the strong assumption of a helpful teacher. In the new relaxed setting, equivalence queries are approximated by *testing queries*: a call to EQ yields membership queries for a set of randomly selected words; when all of them agree with the hypothesis, the algorithm terminates with a non-zero probability of misclassification; otherwise, we have a counter-example to process. The number of such queries may depend on what we assume about the distribution over  $\Sigma^*$  and what we want to prove about the algorithm, for example, PAC learnability, which is further discussed in Section 8.

The main learning algorithm used in this chapter is Algorithm 8 as it appears in Section 5.3 and counter-examples are treated by applying the symbolic breakpoint method (Procedure 13). In the following, we discuss how we specialize the treatment of evidence incompatibility for alphabets having a particular structure. Specifically, we show alternative methods and convenient modifications that we can apply on the procedures  $sample(\Sigma, k)$ , select and EvCOMP.

To improve the learning procedure and make it more efficient, we use multiple evidences per state. However, we let only one evidence, the representative, take part in the construction of the observation table and the hypothesis. By doing this we avoid the exponential growth of the concrete sample. The sampling and selection criteria we use to define the evidences and representatives vary and rely upon the type of the alphabet and the partitions. We discuss different types of alphabets and semantics such as interval subsets of countable or uncountable sets and sets of Boolean vectors.

In the previous chapter, the symbolic learning algorithm manages to exactly learn an automaton representation for the target language. Nevertheless, the key to learn such languages, without error, is the presence of a teacher able to answer to equivalence queries and return counter-examples that are minimal. Minimality is what allows the learner to precisely determine the boundaries of each partition. In this chapter, we assume a teacher who is not that helpful and is able to answer only to membership queries. As we shall see, the symbolic algorithm is able to learn a symbolic automaton to represent the target language, but this time as an approximation. Since minimality of the counter-examples cannot be guaranteed, we choose to use a PAC learning criterion for termination. The set of symbolic recognizable languages is shown to be PAC learnable resulting, almost always, in an automaton that accepts the target language approximately well.

# 7.1 Approximating the Equivalence Query

We approximate equivalence queries by testing as follows. First, a word  $w \in \Sigma^*$  is generated according to some probability distribution  $\mathcal{P}$  over  $\Sigma^*$  (line 4). Then a membership query MQ(w) is posed and the learner checks whether the outcome matches the conjectured hypothesis. If not, the testing procedure terminates returning the word w as a counter-example. If no counter-example has been found after a sufficiently large number of tested words, the equivalence query terminates and the hypothesis is considered correct with some probability of error. The probability distribution that is used to generate words need not be known to the learner.

Let  $\varepsilon$  and  $\delta$  be the accuracy and confidence parameters. A hypothesis automaton  $\mathcal{A}$  is an  $\varepsilon$ -approximation of the target L with a probability higher than  $1 - \delta$ . In such a case we say that the hypothesis is *probably approximately correct* (PAC). For more details on the PAC learnability see Chapter 8.3.2. To guarantee a PAC hypothesis, the testing method should successfully test at least  $r_i$  random words, where  $r_i = \frac{1}{\varepsilon} \left( \log \frac{1}{\delta} + (i+1) \log 2 \right)$ , and i is the number of hypotheses that have already been tested.

Pro	Procedure 16 Testing Oracle							
1:	procedure $\text{TEST}(\mathcal{A}^i, \varepsilon, \delta)$	$\triangleright \varepsilon$ accuracy, $\delta$ confidence						
2:	counter = 1							
3:	while $counter < r_i  \mathbf{do}$							
4:	$w = random\_word_{\mathcal{P}}(\Sigma^*)$	$\triangleright \mathcal{P}$ is a prob. distribution over $\Sigma^*$						
5:	if $ extsf{MQ}(w)  eq oldsymbol{\mathcal{A}}^i(w)$ then							
6:	return w	▷ counter-example found						
7:	end if							
8:	counter = counter + 1							
9:	end while							
10:	return True							
11:	end procedure							

# **7.2** Learning Languages over $\mathbb{N}, \mathbb{R}$

As in Section 6.1, we consider bounded subsets of  $\mathbb{N}$  or  $\mathbb{R}$  as the input alphabet  $\Sigma$ , such that the semantics of each symbol is an interval. We disallow disconnected partition blocks, for example, two subsets of even and odd numbers, respectively. Thus, if two disconnected intervals take the same transition, two symbolic letters will be considered. In this setting, the endpoints of an interval associated with a symbolic letter are such that all evidence points between them have the same residual function, while the nearest points outside the interval have different residuals. In an interval setting, semantics can be interpreted as conjunctions of inequalities where the partitioning points reside between two consecutive evidences with different residual functions.

To infer the partition for a state s and define the outgoing transitions in the conjectured automaton we use evidences. Having no prior knowledge of the behavior of a new symbol, when initializing it, we use multiple evidences to get a better approximation. For this reason, when s is a new state, the partition is formed from a sample set of evidences. However, a new evidence may arrive later, introduced by some counter-example, in order to fix an incorrectly conjectured partition.

The initial set of evidences can be any set of k concrete letters denoted by  $sample(\Sigma, k)$ , see line 3 of Procedure 10. This set can be selected randomly or be the result of a fixed step quantization of  $\Sigma$ , that is,  $sample(\Sigma, k) = \{a + i \cdot \Delta : i = 0, ..., k - 1\}$  where  $a = \min \Sigma$  and  $\Delta = |\Sigma|/(k - 1)$ . Moreover, it can be the same each time INITSYMBOL is evoked. One can think of evidence selection by a more adaptive process that depends on the outcome of membership queries. Each symbol admits one element from the set of evidences, usually randomly chosen, as a representative.

Evidence incompatibility appears either right after symbol initialization or after a counter-example treatment. Evidence incompatibility is solved by consecutive calls to EvCOMP until the total incompatibility degree of the observation table becomes zero. Note that partition blocks form convex intervals and, as a consequence, several symbols in the symbolic alphabet may have the same behavior, that is, label a transition that leads to the same target state in the automaton.

For a state s, an incompatibility instance at  $a \in [\![a]\!]$  indicates either that the partition boundary is imprecise or that we missed a symbol and its corresponding transition. In the first case, the incompatible evidence a appears next to the boundary of the interval and its classification matches the classification of a neighboring symbol  $a' \in \mathcal{V}(a)$ . In this situation, modifying the boundary so that a is moved to  $[\![a']\!]$  resolves the incompatibility. On the other hand, when the evidence a is in the interior of an interval, or does not behave as a neighboring symbol in  $\mathcal{V}(a)$ , the incompatibility is resolved by adding a new symbol and refining the existing



Figure 7.1: Evidence incompatibility solved by boundary modification. The incompatible evidence  $a^i$  matches the classification of its neighboring symbol, and a boundary update solves the incompatibility without adding a new symbol.



Figure 7.2: Evidence incompatibility solved by introducing new symbols. More symbols may need to be added calling EvCOMP multiple times until the incompatibility is eliminated.

partition. These two cases are illustrated in Figures 7.1 and 7.2, respectively.

Formally, let  $s \in S$  be a state with positive incompatibility degree M(s) > 0, and let  $\mu_s = \{a^1, \ldots, a^k\} \subset S$  be the set of evidences, ordered such that  $a^{i-1} < a^i$ for all *i*. To simplify notation,  $f^i$  denotes the residual  $f_{\hat{\mu}(s) \cdot a^i}$  when state *s* is understood from the context.

Let  $a^j$  and  $a^{j+1}$  denote symbols in  $\Sigma_s$  with adjacent semantics, and let  $a^{i-1}, a^i \in \mu_s$  be two evidences from the same partition block that behave differently,  $f^{i-1} \neq f^i$ . Moreover, let  $a^j \in \Sigma_s$  be the symbol such that  $a^{i-1}, a^i \in \mu(a^j)$  where  $[\![a^j]\!] = [c, c')$ . There exists at least one partitioning point  $p \in (a^{i-1}, a^i)$ , and it is given by a split method  $p = split(a^{i-1}, a^i)$ . We let split return the middle point, that is split(a, a') = (a + a')/2; other more sophisticated methods can be applied instead, for instance, a split method which asks membership queries in a binary search fashion.

The remedy to evidence incompatibility is to separate  $a^{i-1}$  and  $a^i$  and map them to different symbols. Procedure 17 describes how evidence incompatibility is solved and how evidences and semantics are updated.

The way this separation is realized, with or without introducing a new symbol,

**Procedure 17** Make Evidence Compatible (without Helpful Teacher) -  $\mathbb{R}$ 1: procedure EVCOMP Let  $s \in S$ , for which M(s) > 0, where 2:  $\mu_{s} = \{a^{1}, \dots, a^{k}\}$  such that  $a^{i-1} < a^{i}, \forall i = 2, \dots, k$ 3: Let  $\mathbf{a}^j \in \mathbf{\Sigma}_{\mathbf{s}}, [\![\mathbf{a}^j]\!] = [c, c')$ , such that  $\exists i : f^{i-1} \neq f^i$  for  $a^{i-1}, a^i \in \mu(\mathbf{a}^j)$ 4:  $p = split(a^{i-1}, a^i)$ ▷ new partitioning point 5: if  $f^{i} = f^{i+1} = \cdots = f^{i+l+1}$ , where  $a^{i}, \ldots, a^{i+l} \in \mu(a^{j}), a^{i+l+1} \in \mu(a^{j+1})$  then 6:  $\llbracket a^{j} \rrbracket = [c, p); \llbracket a^{j+1} \rrbracket = [p, c') \cup \llbracket a^{j+1} \rrbracket$ ▷ change right frontier 7:  $\mu(a^{j+1}) = (\mu(a^{j+1}) \cup \mu(a^j)) \cap \llbracket a^{j+1} \rrbracket$ 8:  $\mu(\boldsymbol{a}^j) = \mu(\boldsymbol{a}^j) \cap \llbracket \boldsymbol{a}^j \rrbracket$ 9: else if  $f^{i-1} = \cdots = f^{i-l}$ , where  $a^{i-1}, \ldots, a^{i-l+1} \in \mu(a^j), a^{i-l} \in \mu(a^{j-1})$  then 10:  $\llbracket a^{j-1} \rrbracket = \llbracket a^{j-1} \rrbracket \cup [c, p); \llbracket a^{j} \rrbracket = [p, c')$ ▷ change left frontier 11:  $\mu(a^{j-1}) = (\mu(a^{j-1}) \cup \mu(a^j)) \cap \llbracket a^{j-1} \rrbracket$ 12:  $\mu(\boldsymbol{a}^j) = \mu(\boldsymbol{a}^j) \cap \llbracket \boldsymbol{a}^j \rrbracket$ 13: else 14:  $\Sigma_s = \Sigma_s \cup \{b\}$ ▷ introduce a new symbol 15:  $R = R \cup \{s \cdot b\}$ 16: if  $\hat{\mu}(\boldsymbol{a}^j) < p$  then 17:  $[\![a^j]\!] = [c, p); [\![b]\!] = [p, c')$ 18: 19: else  $[\![b]\!] = [c, p); [\![a^j]\!] = [p, c')$ 20: 21: end if  $\mu(\boldsymbol{b}) = \mu(\boldsymbol{a}^j) \cap \llbracket \boldsymbol{b} \rrbracket; \mu(\boldsymbol{a}^j) = \mu(\boldsymbol{a}^j) \cap \llbracket \boldsymbol{a}^j \rrbracket$ 22:  $\hat{\mu}(\boldsymbol{b}) = select(\mu(\boldsymbol{b}))$ 23: 24:  $f_{\boldsymbol{s}\cdot\boldsymbol{b}} = f_{\hat{\mu}(\boldsymbol{s}\cdot\boldsymbol{b})}$ 25: end if 26: end procedure

depends on the positions of  $a^{i-1}$  and  $a^i$  in the set of evidences and the residual values of their neighboring partition blocks. In the following,  $a^{j-1}$ ,  $a^j$  and  $a^{j+1}$  denote symbols in  $\Sigma_s$  with adjacent semantics.

Boundary modification. Suppose the incompatibility instance is at a<sup>i</sup> ∈ μ(a<sup>j</sup>), and that all other evidences μ(a<sup>j</sup>) to the right of a<sup>i</sup> behave like min μ(a<sup>j+1</sup>). By changing the partition boundaries and moving a<sup>i</sup> from [[a<sup>j</sup>]] to [[a<sup>j+1</sup>]], the incompatibility instance at a<sup>i</sup> is eliminated. The new boundary between these two intervals is set to p, see Figure 7.1. The semantics and evidence functions are updated accordingly. The symmetric case, where the incompatibility occurs at a<sup>i-1</sup> ∈ a<sup>j</sup> with all other evidences of μ(a<sup>j</sup>) on its left behaving like max μ(a<sup>j-1</sup>), is treated similarly.

Symbol introduction. When the above condition does not hold and boundary modification cannot be applied, which happens either when the evidences are internal to [[a<sup>j</sup>]], or when their residuals do not match the residuals of any neighboring symbol, the incompatibility is solved by refining the partition. The semantics [[a<sup>j</sup>]] is split into two intervals [c, p) and [p, c'). A new symbol b is introduced and the interval not containing µ(a<sup>j</sup>) is moved from [[a<sup>j</sup>]] to [[b]] along with the evidences it contains, see Figure 7.2.

Counter-examples are treated by calling Procedure 13 (see Chapter 5) where each counter-example expands the hypothesis automaton either vertically or horizon-tally. Let us demonstrate the working of the algorithm in learning a target language over a subset of  $\mathbb{R}$ .



Figure 7.3: Minimal automaton accepting the target language used in Example 7.1.

**Example 7.1** (Learning Interval Automata without Minimal Counter-Examples). Let  $\Sigma = [0, 100) \subseteq \mathbb{R}$  be the concrete alphabet and a target language  $L \subseteq \Sigma^*$ , shown in Figure 7.3. We recall that equivalence is approximated by testing and that the teacher answers only membership queries. The observation tables, semantics functions and hypotheses used in this example are shown in Figures 7.4 and 7.5.

The table is initialized with  $S = \{\varepsilon\}$  and  $E = \{\varepsilon\}$ . To determine the alphabet partition at the initial state  $\varepsilon$ , the learner asks membership queries for the randomly selected one-letter words  $\{13, 42, 68, 78, 92\}$ . All words in this set except 13 are rejected. Consequently, there are at least two distinct intervals and we take split(13, 42) = 27 as their boundary. Each interval is represented by a symbolic letter resulting in  $\Sigma_{\varepsilon} = \{a_1, a_2\}, \ \mu(a_1) = \{13\}, \ \hat{\mu}(a_1) = 13, \ \mu(a_2) = \{42, 68, 78, 92\}, \text{ and } \hat{\mu}(a_2) = 68$ . The representatives are randomly chosen from the set of evidences. The semantics,  $\psi$  maps all letters smaller than 27 to  $a_1$ , and maps the rest to  $a_2$ , that is,  $[\![a_1]\!] = [0, 27)$  and  $[\![a_2]\!] = [27, 100)$ . The table boundary updates to  $\mathbf{R} = \{a_1, a_2\}$  and the observation table is  $T_0$ , shown in Figure 7.4.

Table  $T_0$  is not closed and in order to fix this, the learner moves  $a_1$  to the set of states S. To find the possible partitions of  $\Sigma$  at this new state  $a_1$ , the learner randomly chooses a sample  $\{2, 18, 26, 46, 54\}$  of letters and asks membership queries

			$T_3$	$T_{4-5}$
$T_0$	$T_1$	$T_2$	$\varepsilon$ 11	$\varepsilon$ 11 $\varepsilon$ - +
ε	$\varepsilon$	$\varepsilon$ 11 $\varepsilon$ - +	$egin{array}{c c} arepsilon & arepsilon \ a_1 & arepsilon & arepsilon \ a_1 & arepsilon & arepsilon \ \end{array}$	$\begin{vmatrix} a_1 \end{vmatrix} + \begin{vmatrix} - \end{vmatrix}$
$egin{array}{c c c c c c c c c c c c c c c c c c c $	$a_1 + a_1$	$a_1 + -$	$a_2$	$egin{array}{c c c c c c c c c c c c c c c c c c c $
$ a_2 $ –	$egin{array}{c c} oldsymbol{a}_2 &= & \ oldsymbol{a}_1 oldsymbol{a}_3 &= & \ oldsymbol{a}_1 oldsymbol{a}_3 &= & \ \end{array}$	$egin{array}{c c} a_2 &=& - \ a_1a_3 &-& + \end{array}$	$egin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{vmatrix} a_1 a_6 \\ a_2 a_4 \end{vmatrix} + - \end{vmatrix}$
			$ a_2 a_5  + -$	$\begin{vmatrix} a_2 a_4 \\ a_2 a_5 \end{vmatrix} + - \end{vmatrix}$

Figure 7.4: Observation tables used in Example 7.1.

concerning the words in  $\{13 \cdot 2, 13 \cdot 18, 13 \cdot 26, 13 \cdot 46, 13 \cdot 54\}$ . Note that the prefix used here is the representative of  $a_1$ . The teacher classifies all words as rejected. The new table is  $T_1$  with  $\Sigma_{a_1} = \{a_3\}, \mu(a_3) = \{2, 18, 26, 46, 54\}, \hat{\mu}(a_3) = 18$ , and  $[a_3] = [0, 100)$ . The new table is closed and the first hypothesis  $\mathcal{A}_1$  is conjectured.

The hypothesis is tested on a set of words, randomly chosen from some distribution, typically unknown to the learner. After some successful tests, a word  $35 \cdot 52 \cdot 11$  is found, which is accepted by  $\mathcal{A}_1$  but is outside the target language. The learner takes this word as a counter-example and analyzes it using the symbolic breakpoint method. At iteration i = 2 of Procedure 13, condition (5.3) is violated, in particular  $MQ(\hat{\mu}(\boldsymbol{\varepsilon} \cdot \boldsymbol{a}_2) \cdot 11) = MQ(68 \cdot 11) \neq flag = +$ . Thus, the suffix 11 is added as a distinguishing word to E. The observation table  $T_2$  obtained after adding the new suffix is, as expected, not closed. The table is made closed by letting  $\boldsymbol{a}_2$  be a new state, resulting in table  $T_3$ , where  $\boldsymbol{\Sigma}_{\boldsymbol{a}_2} = \{\boldsymbol{a}_4, \boldsymbol{a}_5\}, \, \mu(\boldsymbol{a}_4) = \{17, 27\}, \, \hat{\mu}(\boldsymbol{a}_4) = 17, \, [\![\boldsymbol{a}_4]\!] = [0, 45), \, \mu(\boldsymbol{a}_5) = \{64, 72, 94\}, \, \hat{\mu}(\boldsymbol{a}_5) = 72$  and  $[\![\boldsymbol{a}_5]\!] = [45, 100)$ . The corresponding new conjecture is  $\mathcal{A}_3$ .

Automaton  $\mathcal{A}_3$  is tested and a counter-example  $12 \cdot 73 \cdot 4$  is provided. The breakpoint method discovers that condition (5.2) is violated because letter 73 is not part of the semantics of  $a_3$ . This letter is added as a new evidence to  $\mu(a_3)$ . The evidence incompatibility is solved by splitting the existing partition into two subintervals. A new symbol  $a_6$  is added to  $\Sigma_{a_1}$ , such that  $\mu(a_6) = \{73\}$  and  $[a_6] = [63, 100)$ . The new observation table and hypothesis automaton are  $T_4$  and  $\mathcal{A}_4$ , respectively.

The next counter-example 52.47, also adds a new evidence, this time to symbol  $a_5$ . The classification of the new evidence matches the classification of  $a_4$ , which is a neighboring symbol. The boundary between  $[\![a_4]\!]$  and  $[\![a_5]\!]$  is moved from 45 to 55, thus resolving the evidence incompatibility. The new hypothesis  $\mathcal{A}_5$  is suc-



Figure 7.5: Symbolic automata and semantics function learned in Example 7.1.

cessfully tested without discovering any other counter-example and the algorithm terminates while returning  $A_5$  as an answer.

# **7.3** Learning Languages over $\mathbb{B}^n$

We demonstrate the versatility of the algorithm, which was first developed for numerical alphabets, by adapting it to languages over the alphabet  $\Sigma = \mathbb{B}^n$  of Boolean vectors accessed by variables  $\{x_1, \ldots, x_n\}$ . All components of the algorithm remain the same except the construction of alphabet partitions and their modification due to evidence incompatibility. These should be adapted to the particular nature of the Boolean hyper-cube. The concrete semantics of the symbolic letters in a state s, which can either form sub-cubes (terms) or finite unions of sub-cubes, will be defined by a function  $\psi_s : \mathbb{B}^n \to \Sigma_s$ .

At any given moment, the raw data for inducing the alphabet partition at s is the

Pro	cedure 18 Make Evidenc	e Compatible (without Helpful Teacher) - $\mathbb{B}^n$
1:	procedure EvComp	
2:	Let $s \in S$ be a state $t$	for which $\mathbf{M}(\boldsymbol{s}) > 0$
3:	$\text{Update}(T_{\boldsymbol{s}})$	▷ make tree consistent with sample
4:	for all $h \in \mathcal{F}_{\boldsymbol{s}}$ do	$\triangleright$ Update $\xi_s$
5:	if $\exists oldsymbol{a} \in oldsymbol{\Sigma}_{oldsymbol{s}}$ s.t. $h$	$= f_{\hat{\mu}(\boldsymbol{s})\cdot\hat{\mu}(\boldsymbol{a})}$ then
6:		$\triangleright h$ is already associated with an existing symbol $a$
7:	$\xi_{s}(h) = a$	
8:	$\mu(oldsymbol{a})=\{a^i\in$	$\{\mu_{m{s}}: f^i = h\}$
9:	else	$\triangleright$ h does not match any pre-existing residual
10:	${old \Sigma}_{oldsymbol{s}}={old \Sigma}_{oldsymbol{s}}\cup ig \{$	b > introduce new symbol $b$
11:	$oldsymbol{R} = oldsymbol{R} \cup \{oldsymbol{s}$	$\cdot b$ }
12:	$\xi_{s}(h) = b$	
13:	$\mu(\boldsymbol{b})=\{a^{i}\in$	$\mu_{\boldsymbol{s}}: f^i = h\}$
14:	$\hat{\mu}(oldsymbol{b}) = selec$	$t(\mu(oldsymbol{b}))$
15:	end if	
16:	end for	
17:	$\psi_{\boldsymbol{s}} = \xi_{\boldsymbol{s}} \circ T_{\boldsymbol{s}}$	
18:	end procedure	

 $\psi_s = \xi_s \circ I_s$ nd procedure

sample  $\{(a^i, f^i) : a^i \in \mu_s\}$  where  $\mu_s$  is the set of all evidences for state s, and for every  $a^i$ ,  $f^i = f_{\hat{\mu}(s) \cdot a^i}$  is the residual associated with  $a^i$ . Let  $\mathcal{F}_s = \{f^i : a^i \in \mu_s\}$ be the set of all observed distinct residuals associated with the one-letter successors of s. We view  $\psi_s$ , associated with a state  $s \in S$ , as a composition of two functions, that is,  $\psi_s = \xi_s \circ T_{\psi_s}$ , where  $T_{\psi_s} : \mathbb{B}^n \to \mathcal{F}_s$  and  $\xi_s : \mathcal{F}_s \to \Sigma_s$ . To simplify the notation, we use  $T_s$  instead of  $T_{\psi_s}$ .

The function  $T_s$  is represented by a binary decision tree (see Section 2.2) whose leaf nodes are labeled by elements of  $\mathcal{F}_s$  and which is considered compatible with the sample if it agrees with it on the elements of  $\mu_s$ . The function  $\xi_s$  can be seen as a relabeling that maps each distinct residual to a symbol from  $\Sigma_s$ . We let this function be a bijection. Had we wanted to follow the "convex" partition approach that we used for numerical alphabets, we should have associated a fresh symbol with each leaf node of the tree, thus letting  $[\![a]\!]$  be a cube/term for every  $a \in \Sigma_q$ . However, unlike numerical alphabets, we prefer here to associate the same symbol with multiple leaf nodes that share the same label (residual), allowing the semantics of a symbol to be a finite union of cubes. This results in  $|\Sigma_s| = |\mathcal{F}_s|$  and the existence of at most one symbol that labels a transition between any pair of states. For the function  $\xi_s$  to be well defined, we require that it sends each representative to the symbol it represents, i.e.,

$$\xi_{\boldsymbol{s}}(f_{\hat{\mu}(\boldsymbol{s}),\hat{\mu}(\boldsymbol{a})}) = \boldsymbol{a} \text{ for all symbols } \boldsymbol{a} \in \boldsymbol{\Sigma}_{\boldsymbol{s}}.$$
(7.1)

To learn the semantics  $\psi_s$  for state s, it suffices to learn the two functions  $\xi_s$  and  $T_s$ . We first build  $T_s$  as a decision tree where all evidences mapped to the same leaf node agree on their residual function. Hence, learning alphabet partitions is an instance of learning decision trees using algorithms such as CART [BFSO84], ID3 [Qui86], or ID5 [Utg89] that construct a tree compatible with a labeled sample. Then  $\xi_s$  may be updated and new symbols may be added to  $\Sigma_s$  if necessary.

Evidence incompatibility in a state s appears when the decision tree  $T_s$  is not compatible with the sample. This may happen in three occasions during the execution of the algorithm, the first being symbol initialization. Recall that when a new state s is introduced, we create a new symbol a and collect evidences for it, which may have different residuals while being associated with the same single root node. The second occasion occurs when new evidence is added to a symbol, making a leaf node in the tree impure. Finally, when some new suffix is added to E, the set  $\mathcal{F}_s$  of distinct residuals (rows in the table) may increase and the labels of existing evidences may change.

The simplest way to fix incompatibility of a decision tree is to split impure leaf nodes until purification. However, this may lead to very deep trees and it is preferable to reconstruct the tree each time the sample is updated in a way that leads to incompatibility. In the simple (second) case where a new evidence is added, we can use an incremental algorithm such as ID5, see Section 2.3, which restructures only parts of the tree that need to be modified, leaving the rest of the tree intact. This algorithm produces the same tree as a non-incremental algorithm would, while performing less computation. In the third case, we build the tree from scratch and this is also what we do after initialization where the incremental and non-incremental algorithms coincide.

Once a tree  $T_s$  is made compatible with the sample, the semantics of the symbolic alphabet, expressed via  $\psi_s$ , is updated. This is nothing but updating  $\xi_s$  to map the possibly new residuals from  $\mathcal{F}_s$  to  $\Sigma_s$ . First, with each symbol a that already exists, we re-associate the leaves of  $T_s$  that agree with the labels of its representative such that (7.1) holds. Note that,  $\mu_s$  may have changed but the representative of an existing symbol remains the same. Then, in the case where the set  $\mathcal{F}_s$  of distinct residuals has increased in number, we introduce a new symbolic letter for each new residual and select its representative.

The whole process is described in Procedure 18. We use UPDATE $(T_s)$  for any tree learning algorithm that can be used to make the tree compatible with the sample. This can be an incremental or non-incremental algorithm as the ones presented



Figure 7.6: Minimal automaton accepting the target language used in Example 7.2.

in Section 2.3 and need not be specified here.

**Example 7.2** (Learning Languages over Boolean Alphabets). In this example, the algorithm is applied to learn the target language L over  $\Sigma = \mathbb{B}^4$ , which is shown in Figure 7.6. We denote by  $X = \{x_1, x_2, x_3, x_4\}$  the set of Boolean variables used to access the attributes of the alphabet. All tables encountered during the execution of the algorithm are shown in Figure 7.7, and the semantics functions  $\psi_s$ , realized by decision trees, appear in Figure 7.9 in the form of Karnaugh maps. All letters that are used as evidences in the example are noted in Figure 7.9, where different colors and shapes indicate different residuals.

The learner starts by initializing the observation table to  $T_0$ . Like every new state, the initial state  $\varepsilon$  admits one outgoing transition. Let us label this transition by the symbol  $a_0$  and let the symbolic alphabet be  $\Sigma_{\varepsilon} = \{a_0\}$ . Since  $a_0$ represents, initially, all concrete letters in  $\Sigma$ , the semantics function  $\psi_{\varepsilon}^0$  is a decision tree consisting of a single node, and  $[a_0] = \Sigma$ . To find the behavior of  $a_0$ , a set of concrete letters is sampled and used as the evidence for it, e.g.,  $\mu(a_0) = \{(0000), (0010), (1011), (1000), (1101)\}$ . The letter  $\hat{\mu}(a_0) = (0000)$  is chosen as a representative. The table is filled in after all evidences are queried.

As not all evidences behave the same, e.g., evidence  $(1101) \in \mu(a_0)$  behaves differently than the representative (0000), the observation table  $T_0$  is not evidence compatible and thus the partition at state  $\varepsilon$  needs refinement. The learner applies the tree induction algorithm CART, which is used throughout this example, to refine the BDT  $\psi_{\varepsilon}$ , and finds that  $\Sigma$  is best split into two blocks based on the values of variable  $x_2$ . A new symbol  $a_1$ , which is added to  $\Sigma_{\varepsilon}$ , is used to denote the new block (the new behavior). With this in mind, all letters for which  $x_2 = 0$  are mapped to  $a_0$  and all letters for which  $x_2 = 1$  are mapped to  $a_1$ . The semantics function for  $\varepsilon$  updates to  $\psi_{\varepsilon}^1$ .

					$T_2$				$T_3$	
$T_0$		$T_1$			ε			ε	0000	
	-		ε		ε	_	1 [	ε	_	_
6 -	_	ε	$\varepsilon$ –		$\boldsymbol{a}_1$	+		$oldsymbol{a}_1$	+	_
		$oldsymbol{a}_0$	_		$\boldsymbol{a}_0$	_		$oldsymbol{a}_0$	-	_
$\boldsymbol{u}_0$		$oldsymbol{a}_1$	+	a	$_1 \boldsymbol{a}_2$	_		$a_1a_2$	-	+
				a	$_1 a_3$	+		$a_1a_3$	+	—
$\overline{}$										
/	$T_{4}$			5-6		[		ε	0000	1110
	-	0000		ε	000	0	ε	_	_	+
	ε	0000	ε	_	_		$oldsymbol{a}_1$	+	—	+
E Q	_	_	$a_1$	+	—		$oldsymbol{a}_1oldsymbol{a}_2$	-	+	_
	Ŧ	_	$oldsymbol{a}_1oldsymbol{a}_2$	_	+		$oldsymbol{a}_0$	-	_	_
			$a_0$	_	_		$oldsymbol{a}_5$	-	+	_
	_		$a_5$	_	+		$oldsymbol{a}_0oldsymbol{a}_7$	-	_	+
			$a_1a_3$	+	_		$oldsymbol{a}_0oldsymbol{a}_8$	-	+	—
	+	_	$ a_1a_2a_4 $	—	—		$oldsymbol{a}_1oldsymbol{a}_3$	+	—	+
<b>u</b> <sub>1</sub> <b>u</b> <sub>2</sub> <b>u</b> <sub>6</sub>	I		$a_1a_2a_6$	+	_		$oldsymbol{a}_1oldsymbol{a}_2oldsymbol{a}_4$	-	—	—
							$a_1a_2a_6$	+	_	+

Figure 7.7: Observation tables generated during the execution of the algorithm on Example 7.2.

The resulting observation table  $T_1$  is not closed. To close it,  $a_1$  becomes a state. The symbolic alphabet for  $a_1$  is set initially to  $\Sigma_{a_1} = \{a_2\}$ . The evidence for the new state/symbol is sampled to  $\mu(a_1) = \{(0000), (0010), (1011), (1000), (1101)\}$  and  $\hat{\mu}(a_1) = (0000)$  is chosen as a representative. For simplicity, we use the same sample for each new state, but this is not a requirement in applying the algorithm. After solving evidence incompatibility that is caused by the symbol initialization, we obtain the symbolic alphabet  $\Sigma_{a_1} = \{a_2, a_3\}$ , the semantics function  $\psi = \{\psi_{\varepsilon}^1, \psi_{a_1}\}$ , and the observation table  $T_2$ . The first conjectured automaton is  $\mathcal{A}_2$ , shown in Figure 7.8.

After testing the hypothesis  $\mathcal{A}_2$ , a counter-example  $w = (1010) \cdot (0000)$  has been found. The learner, applying the breakpoint method, detects the distinguishing word (0000) and adds it to E. After posing the necessary MQ's, the new observation table is  $T_3$ . The table is neither closed nor evidence compatible. We first close  $T_3$  and then solve the incompatibility at  $\varepsilon$  which was caused by changing the residual functions of the evidence, see  $\psi_{\varepsilon}^2$ .

The table is made closed by making  $a_1a_2$  a state. After sampling evidence,



Figure 7.8: Intermediate and final Automata conjectures made in Example 7.2.



Figure 7.9: Semantics functions used in Example 7.2. We show the evolution of  $\psi_{\varepsilon}$  over time, while for the other states we show only the final partition. We use symbols such as  $\{\bullet, \bullet, \bullet\}$  to indicate different residuals.

selecting the representative, and building an evidence compatible partition at  $a_1a_2$ , defined by  $\psi_{a_1a_2}$ , the table updates to  $T_4$ . Now, we return to state  $\varepsilon$  and solve the evidence incompatibility that appeared in  $\psi_{\varepsilon}^2$  when (0000) was added to E. The decision tree  $\psi_{\varepsilon}$  is reconstructed from scratch to become compatible with the updated sample. The symbols are rearranged such that they match the residuals of their representatives. A new symbol  $a_5$  is added to the symbolic alphabet. The new evidence compatible partition is  $\psi_{\varepsilon}^3$ .

The new table is  $T_5$ . The new hypothesis  $\mathcal{A}_5$  is tested for equivalence, providing the counter-example w = (1111). Applying the breakpoint method on the new counter-example shows that the letter (1111) does not behave like  $a_5$  as it was assumed in the hypothesis. The letter (1111) is, thus, added to  $\mu(a_5)$  as a new evidence, causing once more an incompatibility at the initial state. The BDT  $\psi_{\varepsilon}$  is fixed and updated to  $\psi_{\varepsilon}^4$ . Since this incompatibility is due to a new evidence and not to an update of the whole sample, the tree is updated using an incremental algorithm. Observe that this last counter-example fixes the partition by rearranging the sub-cubes that form the partition of  $\mathbb{B}^4$  without adding any new transition. The observation table remains unchanged.

Testing the next hypothesis  $\mathcal{A}_6$ , provides the counter-example  $w = (1000) \cdot (1000) \cdot (0000) \cdot (1110)$ . The treatment of the counter-example adds the new suffix (1110) to E. This distinguishes the initial state from the prefix  $a_0$ , which is now identified as a state. The partition at state  $a_0$ , after refining and making evidence compatible, is given by  $\psi_{a_0}$ , and the observation table  $T_7$  is obtained. The new final hypothesis  $\mathcal{A}_7$  is tested without discovering any counter-example and the algorithm terminates.

# **Theoretical Analysis: Complexity and Termination**

In this chapter, we study the termination and complexity of the algorithm, first under the assumption of a helpful teacher and then in the random testing setting. Most of the analysis resembles that of  $L^*$  except for the special feature of symbolic learning, the modification of partition boundaries without adding new states and transitions, which requires a special treatment. We also propose one way to define probabilities on  $\Sigma^*$  and show how to compute the probability of a misclassification error based on the symmetric difference between the target language and a hypothesis.

# 8.1 Updating the Hypothesis: Counter-Examples

Let  $L \subseteq \Sigma^*$  be the target language and let  $\mathcal{A}$  be the minimal symbolic automaton recognizing L such that  $L = L(\mathcal{A})$ . We assume that  $\mathcal{A}$  has n states and at most m outgoing transitions per state, that is, a symbolic alphabet  $\Sigma = \biguplus_q \Sigma_q$  such that  $|\Sigma_q| \leq m$  for every state q.

A typical run of the algorithm produces a sequence  $\{A_i\}$  of hypotheses with each  $A_i$  different from  $A_{i-1}$ , see Figure 8.1. We say that  $A_i$  is *expansive* if it is different from  $A_{i-1}$  in the transition graph, that is,  $A_i$  has more states and/or transitions than  $A_{i-1}$ . When the structure of the transition graph does not change and only the partition boundaries differ, we call  $A_i$  a *non-expansive* hypothesis. Our goal is to study the behavior of this sequence under different assumptions and in particular to show under what conditions it is finite. Moreover, we show under which conditions exact learnability can be replaced by PAC learnability.

Finite number of expansions. Our learning algorithm produces, by construction,



Figure 8.1: A sequence of hypothesis automata  $\{A_i\}$  produced during a run of the algorithm. We use the symbols  $\{\bullet, \downarrow\}$  to denote expansive and non-expansive hypotheses, respectively.

minimal automata and the size of the hypothesis cannot exceed the size of the symbolic automaton  $\mathcal{A}$ . Consequently, the number of expansive hypotheses is bounded by  $\mathcal{O}(mn)$ , because at most n-1 involve the discovery of a new state and at most n(m-1) involve a new transition. What remains to be shown is that the number of non-expansive hypotheses is bounded. This type of hypothesis is highly influenced by the type of teacher and his ability to answer equivalence queries correctly, as well as the nature of the alphabet. Before moving on to study the special cases of teachers and alphabets, we explain in more detail how a boundary modification affects the hypothesis. Then, it is important to define the error of a hypothesis and how to compute this.

**Boundary modifications.** To infer the partitions, the learner uses partial information and hence, boundaries are often only approximated. As a result, a word may be misclassified just by the fact that it involves a letter falling on the wrong side of the boundary. Such a word is then provided as a counter-example. The misclassified letter is added and yields evidence incompatibility, which is solved by modifying the existing boundary. This procedure may not affect the structure of the hypothesis automaton.

Let us consider a partition at some state of the hypothesis. Figure 8.2 shows such an example with a one dimensional numerical alphabet. Let the real boundary p between two neighboring intervals be approximated by p' = split(a, b), where aand b are evidences with different residuals. At this point the split(a, b) method, which always results a point in [a, b], is not specified. The letters in [p, p'] are wrongly classified as equivalent to a. We call this interval the *boundary error*. This boundary is modified when a new evidence b' is added to [p, p']. The partitioning point p' updates to p'' = split(a, b') in the next hypothesis.

## 8.2 Hypothesis Error

Let  $L \subseteq \Sigma^*$  be the target language and let  $\mathcal{A}$  be a conjectured automaton accepting the language  $L_{\mathcal{A}}$ . The language  $L_{\mathcal{A}}$  is not equivalent to L when some states or transitions are still missing, or there exists a partition that contains some boundary error.



Figure 8.2: A partition boundary where p, the real partition boundary, is approximated by p' = split(a, b), where a and b are evidences with different residuals. This is the case where p < p' and for the case p' < p we get some a' in the counter-example.

The quality of  $\mathcal{A}$  is determined by its distance from L, denoted by  $d(L, L_{\mathcal{A}})$ . One way<sup>1</sup> to define this distance is to measure the probability of error, that is, the probability of choosing a misclassified word  $w \in \Sigma^*$  that belongs to one language but not to the other. This definition includes all possible sources of error, i.e., incomplete structure, boundary errors.

The set of all misclassified words is exactly the symmetric difference of the languages L and  $L_A$ , i.e.,  $L \oplus L_A = L \setminus L_A \uplus L_A \setminus L$ . Hence, we define the distance as

$$d(L, L_{\mathcal{A}}) = \mathcal{P}(L \oplus L_{\mathcal{A}}),$$

where  $\mathcal{P}$  is a probability distribution over  $\Sigma^*$ .

When an automaton representation of the target language is available, one can find the probability of the error by first computing the symmetric difference of the target and conjectured languages, and then computing the probability of this language. To compute the probability we use the *relative volumes* that are described below.

#### **8.2.1** A Probability Distribution on $\Sigma^*$

In this section we define a natural probability distribution  $\mathcal{P}$  over the set of all words  $\Sigma^*$ , where  $\Sigma$  is an input alphabet. First we choose some distribution  $\mathcal{P}_{\ell}$ :  $\mathbb{N} \to [0, 1]$ , which provides the probability of a word w to be of certain length. Then, each letter in the word obeys a uniform distribution over the alphabet, that is  $\mathcal{P}_{\Sigma}(a) = 1/|\Sigma|$  for all  $a \in \Sigma$ . The probability of a given word  $w = a_1 \cdots a_k \in \Sigma^*$ 

<sup>&</sup>lt;sup>1</sup>Other ways to measure the error (or distance of two languages) have been proposed in the literature, see for instance [SVVV14].

can then be defined as

$$\mathcal{P}(a_1 \cdots a_k) = \mathcal{P}_{\ell}(k) \prod_{i=1}^k \mathcal{P}_{\Sigma}(a_i) = \frac{\mathcal{P}_{\ell}(k)}{|\Sigma|^k}$$

As usual, probability distributions  $\mathcal{P}_{\Sigma}$  and  $\mathcal{P}$  extend to sets of letters and words, respectively. Thus,

$$\mathcal{P}_{\Sigma}(I) = \sum_{a \in I} \mathcal{P}_{\Sigma}(a) = |I|/|\Sigma|, \text{ and}$$
  
 $\mathcal{P}(L) = \sum_{w \in L} \mathcal{P}(w),$ 

where  $I \subseteq \Sigma$  and  $L \subseteq \Sigma^*$ .

By partitioning the set of words according to their length, the probability of a language can be written as

$$\mathcal{P}(L) = \sum_{k=0}^{\infty} \mathcal{P}(L \cap \Sigma^{k})$$
  
=  $\sum_{k=0}^{\infty} \mathcal{P}_{\ell}(k) \cdot \mathcal{P}(L \mid \Sigma^{k}),$  (8.1)

where  $\mathcal{P}(L|\Sigma^k)$  is the probability that a word belongs to L given its length is k. We call the last probability the *relative* k-volume of L and denote it as  $D_k(L)$ . In the following we show how to compute these relative k-volumes using stochastic matrices.

#### 8.2.2 Computing the Relative Volumes

The relative k-volume of a language  $L \subseteq \Sigma^*$ , where  $k \ge 0$  is an integer, is defined as  $D_k(L) = \mathcal{P}(L|\Sigma^k) = vol(L_k)/vol(\Sigma^k)$  where  $vol(\Sigma^k) = |\Sigma|^k$  and  $vol(L_k)$  is the volume of  $L_k = L \cap \Sigma^k$ , which is the subset of L that contains only the words of length k. The volume  $vol(L_k)$  is defined and computed recursively by value iteration on the structure of the automaton that recognizes L [ABD15].

Let  $\mathcal{A} = (\Sigma, \Sigma, \psi, Q, \delta, q_0, F)$  be a symbolic automaton such that  $L = L(\mathcal{A})$ . We can define on this automaton, a *volume transition matrix*  $V = (\nu_{ij}) \in M_{|Q|}(\mathbb{R})$ , where  $\nu_{ij}$  denotes the volume of direct transitions from state *i* to state *j* in  $\mathcal{A}$  defined as  $\nu_{ij} = |\{a \in \Sigma : \delta(i, \psi(a)) = j\}|$ . In the symbolic automata setting, this is equivalent to  $\nu_{ij} = |\cup_{a \in \Sigma_i, \delta(i,a) = j} [a]|$ . Intuitively, this indicates the number of possible ways there are for going from state *i* to state *j* in one step. Note that, for



Figure 8.3: A symbolic automaton  $\mathcal{A}$  and its semantics.

a complete and deterministic automaton  $\mathcal{A}$ ,  $\sum_{j=0}^{|Q|} \nu_{ij} = vol(\Sigma) = |\Sigma|$  for every  $i \in Q$ , and hence  $V/|\Sigma|$  forms a right stochastic matrix. We can find all powers of the matrix V, such that  $V^0$  is the identity matrix,  $V^1 = V$ , and  $V^n = V^{n-1}V$  for all n > 1.

The element  $\nu_{ij}^{(n)}$ , which denotes the (i, j) entry of the matrix  $V^n$ , gives the number of possible paths of length n that lead from a state i to j. The volume of  $L_k$  is given by

$$vol(L_k) = \sum_{j \in F} \nu_{q_0 j}^{(k)},$$

where  $q_0$  is the initial state and F is the set of accepting states in A.

We illustrate this computation with an example. Let  $\Sigma = [0, 100) \subseteq \mathbb{R}$  be an alphabet and let  $L \subseteq \Sigma^*$  be a language, accepted by the symbolic automaton  $\mathcal{A}$  shown in Figure 8.3. Each state  $q_i$  in the automaton has two outgoing transitions labeled by  $a_{2i}$  and  $a_{2i+1}$ , respectively. Given that the volume of an interval [a, b) is vol([a, b)) = b - a, the volume transition matrix defined by  $\mathcal{A}$  is

$$V = \begin{pmatrix} 0 & vol(\boldsymbol{a}_0) & vol(\boldsymbol{a}_1) & 0\\ vol(\boldsymbol{a}_2) & 0 & 0 & vol(\boldsymbol{a}_3)\\ 0 & vol(\boldsymbol{a}_4) & 0 & vol(\boldsymbol{a}_5)\\ 0 & 0 & vol(\boldsymbol{a}_6) & vol(\boldsymbol{a}_7) \end{pmatrix} = \begin{pmatrix} 0 & 25 & 75 & 0\\ 40 & 0 & 0 & 60\\ 0 & 60 & 0 & 40\\ 0 & 0 & 90 & 10 \end{pmatrix}$$

To find the volume of  $L_k$  we compute the matrix  $V^k$  and sum up all numbers of the first row of matrix  $V^k$  that correspond to final states. In the example this is the last element of the first row. In Figure 8.4, we present some powers of the volume transition matrix, where the values are already normalized. According to this figure, we see that  $vol(L_1) = 0$  and  $vol(L_2) = 45$ , etc., and that  $D_2(L) =$  $0.45, D_{45}(L) = 0.315$ .

$V^0 = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$	$V^3/ \Sigma^3  =$	$\begin{pmatrix} 0.18 \\ 0.04 \\ 0 \\ 0.216 \end{pmatrix}$	$\begin{array}{c} 0.025 \\ 0.5040 \\ 0.276 \\ 0.054 \end{array}$	$\begin{array}{c} 0.48 \\ 0.054 \\ 0.54 \\ 0.333 \end{array}$	$\begin{array}{c} 0.315 \\ 0.4020 \\ 0.184 \\ 0.3970 \end{array}$
$V^2/ \Sigma^2  = \begin{pmatrix} 0.1 & 0.45 & 0 \\ 0 & 0.1 & 0.84 \\ 0.24 & 0 & 0.36 \\ 0 & 0.54 & 0.09 \end{pmatrix}$	$ \begin{array}{c} 0.45\\ 0.06\\ 0.4\\ 0.37 \end{array} \right)  V^n /  \Sigma^n  = \\ \end{array} $	(0.0945 0.0945 0.0945 0.0945 for a	$\begin{array}{c} 0.2362 \\ 0.2362 \\ 0.2362 \\ 0.2362 \\ 0.2362 \\ \text{ill } n > 30 \end{array}$	0.3543 0.3543 0.3543 0.3543 0.3543 ).	$\begin{array}{c} 0.3150\\ 0.3150\\ 0.3150\\ 0.3150\\ 0.3150 \end{array}$

Figure 8.4: Some powers of the volume transition matrix V defined by  $\mathcal{A}$  of Figure 8.3.

## 8.3 Complexity and Termination

The termination and complexity of the algorithm are highly affected by the nature of the teacher, that is, the way equivalence queries are realized. We considered two types of equivalence checking. The first is a helpful teacher (Chapter 6), which will always provide a counter-example when there exists one, and moreover, this counter-example is minimal. The second (Chapter 7) implements equivalence by random sampling from the same distribution used to compute the error. This implementation may, of course, miss some counter-examples, but some of its runs may also have the opposite effect and produce a series of counter-examples that reduce the error as little as possible. In the case where  $\Sigma$  is finite this may lead to a behavior as bad as that of a concrete algorithm and in the case of alphabet  $\mathbb{R}$ , even to non termination. However, as we will argue, under probabilistic assumptions, the probability of such a scenario is zero.

#### 8.3.1 Using a Helpful Teacher (Minimal Counter-Examples)

First, we examine the case where a helpful teacher is present. Such a teacher, as presented in Chapter 6, returns a minimal counter-example whenever a conjecture is wrong. 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.

Concerning the size of the observation table T, the set of prefixes S monotonically increases until it reaches the size of exactly n elements. Since the table, by construction, is always kept reduced, the elements in S represent exactly the states of the automaton. The size of the boundary cannot exceed the total number of transitions in the automaton, and thus it is bounded by the size 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. Taking all these into account, we can conclude that the size of the table ranges between  $(n + m) \log_2 n$  and n(mn + 1).

Regarding the size of the samples, and in consequence the number of membership queries, we discuss separately totally and partially ordered alphabets. In the case of numerical totally-ordered alphabets, the learner uses Procedure 14 to treat counter-examples. Partitions are defined as closed-open intervals, the number of sufficient evidence, which is needed to define a partition block, is exactly one, the minimal letter belonging to the interval. The evidence incompatibility, caused by a new evidence, is immediately solved by splitting the current partition block and introducing a transition. As only one evidence is used per transition, the size of the concrete sample coincides with the size of the symbolic sample associated with the table, and the number of membership queries needed is  $\mathcal{O}(mn^2)$ .

Moreover, due to minimality of the counter-examples in the length-lexicographic order, choosing as evidence the minimal element of a partition block the boundaries are correctly defined and allow no error. This way, the run of the algorithm does not contain any non-expansive hypotheses. The algorithm terminates after all states and transitions have been discovered and defined. Hence, for termination at most  $\mathcal{O}(mn)$  equivalence queries are used in total, from which at most n-1 of them add a new state to the hypothesis, and at most n(m-1) add a transition.

When the alphabet is partially-ordered, we have seen in Section 6.2 that each partition block  $F_i$  may need multiple evidences in order to be fully defined. For this, some additional queries are asked. We know that when we use monotone partitions, as in Section 6.2, a finite (at most l) number of incomparable minimal points is sufficient to fully define  $F_i$ . Hence, for every row in T, at most (l-1) supplementary words are added to the concrete sample. The number of membership queries is thus  $O(n^2ml)$ .

Furthermore, in the case of partially-ordered alphabets, additional counterexamples may be given in order to refine existing partitions. That is, at most l-1additional evidences are returned per partition block. This implies that the number of non-expansive hypotheses is limited to nm(l-1). We conclude that the number of equivalence queries is bounded by  $\mathcal{O}(mn^2)$  if l < n and  $\mathcal{O}(lmn)$  otherwise.

**Proposition 8.1.** Let  $L \subseteq \Sigma^*$  be a target language. The symbolic learning algorithm, which uses a helpful teacher, terminates with a correct conjecture after asking at most  $O(n^2ml)$  membership queries and at most O(lmn) equivalence queries.

#### 8.3.2 Equivalence using Random Tests

The algorithm implements equivalence checks by comparing membership in L and in  $L_A$  for words randomly selected according to a probability distribution  $\mathcal{P}$  over  $\Sigma^*$ . The procedure that performs equivalence testing is shown in Procedure 16. It takes as input a hypothesis  $\mathcal{A}_i$  and the accuracy and confidence parameters and returns either a True statement or a counter-example. Note that testing does not guarantee the discovery of a counter-example when the conjecture is wrong.

Within this setting, the learner can only guarantee to return an approximation of the target language as there is no way to ensure exact learnability. We use the PAC (probably approximately correct) learnability criterion instead, where the learner terminates with an automaton, which is almost always a good approximation of the target language L. The approximation measure that is used is the distance between the conjecture language and the target, see Section 8.2.

**Definition 8.2** (PAC Learning [Val84]). A learning algorithm learns a language L in a probably-approximately correct (PAC) manner with probability parameters  $\epsilon$  (accuracy) and  $\delta$  (confidence) if its output  $\mathcal{A}$  satisfies

$$Pr(\mathcal{P}(L \oplus L_{\mathcal{A}}) \le \epsilon) \ge 1 - \delta, \tag{8.2}$$

The minimum number of tests sufficient to conclude that a hypothesis  $\mathcal{A}$  satisfies (8.2) depends on the accuracy and confidence parameters as well as the number of previous hypotheses that are made. This number was introduced in [Ang87, Ang88] for the class of automata and applies in a straightforward way to the symbolic case.

**Proposition 8.3.** The symbolic learning algorithm PAC-learns a language L if the *i*-th equivalence query tests  $r_i = \frac{1}{\epsilon} (\ln \frac{1}{\delta} + (i+1) \ln 2)$  random words without finding a counter-example.

*Proof.* Let  $\mathcal{A}$  be the *i*-th hypothesis tested. We show that after testing  $r_i = \frac{1}{\epsilon} (\ln \frac{1}{\delta} + (i+1) \ln 2)$  random words, then (8.2) holds, or equivalently,

$$Pr(\mathcal{P}(L \oplus L_{\mathcal{A}}) \ge \epsilon) \le \delta.$$

Taking into account that a) the probability of choosing a word which does not belong to the symmetric difference is (1 - p) when the probability of finding a counter-example is p, and b) all words are independently drawn from the same probability distribution, then not drawing a counter-example after  $r_i$  draws implies that

$$Pr(\mathcal{P}(L \oplus L_{\mathcal{A}}) \ge \epsilon) \le (1-\epsilon)^{r_i} \le \sum_{i=1}^{\infty} (1-\epsilon)^{r_i} \le \sum_{i=1}^{\infty} e^{-\epsilon r_i} \le \sum_{i=1}^{\infty} \frac{\delta}{2^{i+1}} \le \delta$$

As for termination, it is always guaranteed when the alphabet is finite, but when  $\Sigma$  is a subset of  $\mathbb{R}$  there is a theoretical possibility of an infinite sequence of non-

expansive hypotheses. We give below an informal argument for the improbability of such a behavior of the algorithm.

In concrete automata, a sample whose prefixes contains all words of size up to 2n will exercise all the paths in the automaton and is sufficient to discover all states and transitions. In the symbolic setting, the same holds for valid symbolic words and since all the transition guards are intervals of a positive measure, we will obtain such a sample after a finite number of equivalence queries. Hence we can restrict our discussion to a sequence of non-expansive hypotheses applied to the final automaton.

Secondly, we have the following about what happens when we modify an error interval [p, p'] associated with a boundary between two neighboring intervals.

**Proposition 8.4.** Each time an error interval [p, p'] is modified, the upper bound on its size is divided by at least 2

*Proof.* Let  $a, b \in \Sigma$  be two consecutive evidences with different residuals in the neighboring intervals. Let  $p \in [a, b]$  be the real boundary which is approximated by p' = split(a, b), as shown in Figure 8.2. Without loss of generality, let p < p'. The boundary error satisfies  $|p - p'| \leq \frac{|a-b|}{2}$ . Now, let  $c \in [p, p']$  be an evidence indicating that p' is wrong. It updates the approximation to p'' = split(a, c). The upper bound on the error is reduced at least to half,  $|p - p''| \leq \frac{|a-c|}{2} \leq \frac{|a-b|}{4}$ .

Naturally, this will decrease by two the bound on the probability to get a counter-example involving this interval in a subsequent step. To have an infinite run of the algorithm requires that each equivalence query (which involves  $r_i$  tests) yields a counter-example. Each such counter-example modifies at least one error interval and halves the bounds on its size. It can be shown, based on probability estimates, that the sizes of *all* boundary intervals decrease with some rate which is at least quadratic in the number of tests, and so does the probability of finding a counter-example. This rate of decrease is faster then the linear growth in the number of queries due to the  $r_i$  formula and, hence, the scenario of an infinite sequence of equivalence queries has a probability zero.

\_\_\_\_\_

# Empirical Results

The symbolic learning algorithm has been implemented and applied to several case studies. In this chapter, we first give some insights on the implementation of the algorithms and procedures in Section 9.1. Then in Section 9.2, we present a simple example of learning a language over a subset of the reals, which illustrates the behavior of the algorithm, especially in the setting of PAC learnability. In Section 9.3, we compare the algorithm's performance to the three concrete learning algorithms presented in Chapter 3. For the comparison to be feasible, it is necessary to restrict ourselves to languages whose alphabet is finite. Another case study appears in Section 9.4, where the target languages represent valid types of passwords that are defined over the ASCII characters. Finally, we study the symbolic algorithm for Boolean vectors and compare its performance with the numerical case in Sections 9.5 and 9.6, respectively.

# 9.1 General Comments on the Implementation

All algorithms presented in this thesis have been implemented using Python, in a prototype implementation. Preserving the notation used in the previous chapters, the method  $sample(\Sigma, k)$  in Procedure 10 returns a sample of size k = 5 chosen uniformly from  $\Sigma$ , and the method select(A) chooses one representative from a uniform distribution over a set of evidences A. The split(a, b) method, used in Procedure 17, returns the middle point of the interval (a, b), where  $a, b \in \mathbb{N}$  or  $\mathbb{R}$ .

The teacher is implemented separately and it can be used by all types of learners. A teacher is characterized as *helpful* of not, depending on the way equivalence is checked. In both cases, we let H denote the hypothesis automaton given to the teacher. For a helpful teacher, the target language  $L \subseteq \Sigma^*$  is given as a DFA or a symbolic automaton  $\mathcal{A}$ . Then, the teacher computes the symmetric difference  $L(A) \oplus L(H)$  and checks emptiness. It returns a *true* statement if this language is empty. Otherwise, it finds the shortest accepting paths in the automaton and returns the lexicographically minimal word in these as a counter-example.

On the other hand, a non-helpful teacher checks equivalence through testing. In this case, the target language need not necessarily be represented as an automaton. To check equivalence, the membership of random words is compared for H and L. The sufficient number of tested words is set to  $\frac{1}{\epsilon}(\ln \frac{1}{\delta} + (i+1)\ln 2)$  where  $\epsilon$  and  $\delta$  are parameters given by the user, and i is a counter for the number of hypotheses made so far.



Figure 9.1: The log-normal distribution  $\Lambda(13.53, 7.75^2)$  along with its normalization over the natural numbers. This distribution is used for the experiments presented in Section 9.4.

To perform testing we implemented an oracle that produces random words from  $\Sigma^*$ . Since determining a probability distribution on  $\Sigma^*$  is not a trivial task, we implemented this by combining two distributions (see Section 8.2.1). The first distribution  $\mathcal{P}_{\ell}$ , defined over  $\mathbb{N}$ , determines the length of the random word; while the second distribution  $\mathcal{P}_{\Sigma}$ , defined over the concrete alphabet  $\Sigma$ , selects the letters that fill each position in the word. In order to prevent very long counter-examples, it is preferable to use a right skewed distribution to determine the length of the random words. We use the *log-normal distribution*<sup>1</sup> in our experiments, normalized over the natural numbers; Figure 9.1 shows the log-normal distribution, which is used for the experiments in Section 9.4. In all our experiments, we use a uniform distribution over the alphabet  $\Sigma$ .

It is important to state here that all comparisons that appear in the sequel pose

<sup>&</sup>lt;sup>1</sup> The *log-normal distribution* is given by  $\Lambda(x; \mu, \sigma^2) = \frac{1}{x\sigma\sqrt{2\pi}}e^{-\frac{1}{2}\left(\frac{\ln x-\mu}{\sigma}\right)^2}$ , where the variable x > 0 and the parameters  $\mu$  and  $\sigma > 0$  are real numbers. Intuitively, one can think that if u is distributed according to the normal distribution  $N(\mu, \sigma^2)$  and  $u = \ln x$ , then x is distributed according to the log-normal distribution  $\Lambda(\mu, \sigma^2)$ .



Figure 9.2: The minimal automaton accepting the target language  $L \subseteq \Sigma^*$ , where  $\Sigma = [0, 100) \subseteq \mathbb{R}$ , used for the empirical evaluation in Section 9.2.

queries to the same teacher and they are compared on learning the same target languages. Any distribution over  $\Sigma^*$  that is used during the testing phase does not change until a learner terminates, and the error is measured using this same distribution. The performance is measured in terms of time, number of queries and size of the final observation table and the samples.

# 9.2 On the Behavior of the Symbolic Learning Algorithm

When the equivalence query is realized by random testing, different runs of the same algorithm may produce slightly different hypotheses. In this section, our aim is to show how different counter-examples impact the evolution of the learning process. We observe the final hypothesis as well as the intermediate conjectures made by the learner and discuss the performance of the symbolic algorithm according to the automaton size, the number of queries and the approximation error. We apply the symbolic algorithm to the target language  $L \subseteq \Sigma^*$ , accepted by the automaton in Figure 9.2, where the concrete input alphabet is  $\Sigma = [0, 100) \subset \mathbb{R}$ . We set the accuracy and confidence parameters to  $\epsilon = \delta = 0.05$ . For our discussion, we have selected four representative runs that appear in Figure 9.3.

The symbolic algorithm manages to recover the whole structure of the target automaton in the majority of the runs, allowing a small error in the partition boundaries, see for instance, the runs R2 and R3. However, we also observe runs, such as R4, where the final error is small, less than 0.2% as required, even though there are missing transitions or states in the final hypothesis. Finally, since PAC learnability is used as the termination criterion, the algorithm terminates with probability 95% ( $\delta = 5\%$ ) with a hypothesis with a total error greater than  $\epsilon = 5\%$ . This happens in run R1, where the learner failed to discover all states of the target, resulting in an error 0.0528. Figure 9.4 shows how the error evolves for all hypotheses made along the learning path for each run.

]	R1		R2	R3				R4	
		÷	• (a) - <u>35</u> (a) 30.			(q)			
		(a) <u>'a0</u>	- (1) - <u>11</u> 11. (2) - 11.	(1) (1)	'a0' (q1 ) (q0) (32'	(q2)		a2' (q1) (a1) (a4', (a5') (q2)	
	41 - (a1) - (a2) - ((a2)) - ((a2))	<b>⊷</b> @ <u>∞</u> (		÷.	'a0' (q1  a2',  a/	) 31 42 42	-0-0-		
		-0-		€ -	a) <u>w</u> (1) <u>w</u>	90 95 96			
		-0		<del>.</del> 	40 <u>77</u>	(4) (4) (4) (4) (4) (4) (4) (4) (4) (4)			
				<del>.</del> 					
				÷ →					
	S	$ \Sigma $	$ M_T $	$ M_T $	MQ's	EQ's	tests	error (D)	
Run R1	4	8	36	83	153	10	750	5.28%	
$\operatorname{Run} R2$	5	12	39	116	1052	13	2040	0.17%	
Run $R3$	5	12	52	149	820	13	1346	0.26%	
Run $R4$	5	11	36	105	1653	10	4235	0.17%	

Figure 9.3: The evolution of the conjectured automata in four different runs of the symbolic algorithm. The target language is accepted by the minimal automaton shown in Figure 9.2. Columns |S| and  $|\Sigma|$  denote the number of states and transitions learned. The size of the symbolic and concrete samples appear in columns  $|M_T|$  and  $|M_T|$ , respectively; MQ's denotes the total number of membership queries used to fill the table, to fix evidence incompatibility and to treat counter-examples. The number of equivalence queries, as well as the total number of words tested, are shown in columns EQ's and *tests*, respectively. The last column shows the approximation error of the final hypothesis.


Figure 9.4: The error of each hypothesis made during the runs of the symbolic algorithm described in Section 9.2.

## 9.3 Comparison with Other Algorithms

In this section, we compare the performance of the symbolic learning algorithm (SL) with the three non-symbolic learning algorithms presented and discussed in Chapter 3. Namely, the  $L^*$  (L) algorithm [Ang87], the reduced  $L^*$  (LR) [RS93], and the suffixes variation of  $L^*$  (LM) taken from [MP95].

The four algorithms are tested on learning the same target languages. The input alphabet  $\Sigma$ , which is adapted to fit the needs of the concrete algorithms, is a finite subset of  $\mathbb{N}$ . All algorithms have access to the same non-helpful teacher and use the same membership and equivalence queries. The target languages, provided as minimal symbolic automata, are randomly generated. The accuracy and confidence parameters are kept fixed to  $\epsilon = \delta = 0.05$ .

We compare the performance of the different algorithms based on two main criteria, the size of the input alphabet and the size of target automaton. To this end, we perform two groups of experiments. In the first group, the four learning algorithms are asked to learn target languages that share a fixed structure of a minimal automaton that consists of 15 states and a maximum of 5 outgoing transitions per state. Then the alphabet size ranges from 2 to 200. The partitions of the alphabet at each state are randomly generated. An example of such an automaton appears in Figure 9.5. In the second group, we are interested in testing the effect of the automaton size. For this, we fix the alphabet to  $\Sigma = \{0, \dots, 150\}$  and generate automata of a random structure. We let the number of states, in their minimal representation, to range from 3 to 45 and allow a maximal number of 5 partition blocks per state.

The different algorithms are tested and compared in terms of several size and



Figure 9.5: An example of an automaton used in the comparison of Section 9.3, where the automaton structure is kept fixed while the alphabet is increasing and the partitions are randomly generated.

time aspects. All results are measured and averaged over multiple runs for each target automaton. Specifically, for each different size of alphabet (in the first group) and each size of automaton (in the second group), we generate 10 different target languages and apply each algorithm 10 times.

We measured the total number of equivalence queries needed until reaching a final hypothesis, the number of intermediate conjectures made by the learner, the number of counter-examples found, as well as the total number of tests performed to guarantee the validity of a hypothesis. Moreover, the number of membership queries is measured, counting separately the queries that are used to characterize the sample, to treat the counter-examples, and used during the testing phase. In addition to the final size of the table (prefixes, boundary, suffixes) and the size of the samples, we also compare the number of states and transitions discovered upon termination. The quality of each final hypothesis is measured by mean of the total execution time, the total time spent for testing, and the time spent in treating counter-examples.

In the following, we show the most interesting experimental results and discuss the differences in the behavior of the algorithms. We display the results of the two groups of experiments in Figures 9.6 and 9.8, respectively.

To start with, the main results of the first group of experiments, where the algorithms are compared on different alphabet sizes, appear in Figure 9.6. The figure shows, the performance of all four algorithms with respect to the total number of membership and equivalence queries, the number of states conjectured in the final



Figure 9.6: A comparison on the average number of membership queries posed, the number of equivalence checks, the states learned and the execution time for each of the algorithms: *SL* (the symbolic algorithm), *L* ([Ang87]), *LM* ([MP95]), and *LR* ([RS93]). The four algorithms are compared over several concrete alphabet sizes that ranges from 10 to 200 letters.

hypothesis, and the average execution times.

As expected, the symbolic algorithm admits the most modest growth in the number of membership queries, shown in the first row of Figure 9.6. It outperforms, in this sense, the best of the concrete algorithms, that is, the reduced algorithm (LR). This becomes more evident when the alphabet size grows beyond a certain point. We present the same results in Figure 9.7, where now the num-



Figure 9.7: A box plot of the number of MQ's posed by each algorithm. This plot representation provides the average, variance, range and error points for each algorithm. Each column represents the experimental results where the size of the alphabet is kept fixed to sizes 10, 25, 80, 75, 100, 150 and 200. One can see that, especially when the size of the alphabet is growing, even the (SL) worst case outperforms the best values of all other algorithms.

ber of membership queries is shown in the form box-plots. This representation provides a more clear visualization that contains the standard deviations and min max values in the data, rather than just the average. We observe that the behavior of the symbolic algorithm is better than that of the other algorithms, but what is more interesting is the fact that even the worst case run of the symbolic algorithm outperforms the best run of all other algorithms.

Referring to the number of states conjectured in the final hypothesis, shown in the third plot of Figure 9.6, we observe that the symbolic algorithm (SL) almost always discovers more states than the other algorithms. Taking into account that all algorithms produce automata that are minimal in the number of states, we can say that SL comes closer to the target automaton representation than the other algorithms. This seems to be a consequence of the increased number of equivalence queries that the algorithm asks in contrast to the other algorithms (see the second row in Figure 9.6).

Surprisingly, in terms of execution times, see last plot of Figure 9.6, the concrete reduced  $L^*$  algorithm (LR) has a similar performance to the symbolic algorithm (SL) when the alphabet size is smaller than 150, a result that is probably due



Figure 9.8: A comparison of the average number of membership queries posed, the number of equivalence checks, and the execution time for each of the algorithms: *SL* (the symbolic algorithm), *L* ([Ang87]), *LM* ([MP95]), and *LR* ([RS93]). The four algorithms are compared on variable sizes of target automata where the number of states ranges from 3 to 45.

to the increased number of equivalence queries that the symbolic algorithm needs.

Figure 9.8 shows the results of the second group of experiments. Here the algorithms are compared on different sizes of target automata, where the number of states ranges from 0 to 45. Similarly to the first group of experiments, one can see that beyond 20 states the symbolic algorithm is better than all the non-symbolic algorithms in terms of membership queries, see the first row. The number of equivalence queries, naturally larger for the symbolic algorithm, increases linearly in the number of states in the automaton.

Finally, referring to execution timings, the concrete reduced  $L^*$  algorithm (LR) has a similar performance to the symbolic algorithm (SL) only when the size of the target automata is kept smaller than 25 states.

To conclude this comparison, the symbolic algorithm asks much fewer membership queries in order to fill the table than any of the concrete algorithms, which need to ask queries for every  $a \in \Sigma$ . On the other hand, the symbolic algorithm gets more counter-examples, generates more hypotheses and needs more equiva-

0	NUL	16	DLE	32	SPC	48	0	64	0	80	Р	96		112	р
1	SOH	17	DC1	33	!	49	1	65	Α	81	Q	97	а	113	q
2	STX	18	DC2	34		50	2	66	В	82	R	98	b	114	r
3	ETX	19	DC3	35	#	51	3	67	С	83	S	99	С	115	S
4	EOT	20	DC4	36	\$	52	4	68	D	84	Т	100	d	116	t
5	ENQ	21	NAK	37	%	53	5	69	E	85	U	101	е	117	u
6	АСК	22	SYN	38	&	54	6	70	F	86	V	102	f	118	v
7	BEL	23	ETB	39	1	55	7	71	G	87	W	103	g	119	w
8	BS	24	CAN	40	(	56	8	72	Н	88	Х	104	h	120	х
9	HT	25	EM	41	)	57	9	73	1	89	Y	105	i.	121	у
10	LF	26	SUB	42	*	58	1.1	74	J	90	Ζ	106	j	122	z
11	VT	27	ESC	43	+	59	;	75	K	91	[	107	k	123	{
12	FF	28	FS	44		60	<	76	L	92	1	108	1	124	
13	CR	29	GS	45	-	61	=	77	М	93	]	109	m	125	}
14	SO	30	RS	46		62	>	78	N	94	^	110	n	126	~
15	SI	31	US	47	/	63	?	79	0	95		111	0	127	DEL

Figure 9.9: The table of all ascci characters, split into groups of *control characters*, *alphabetic*, *numerals*, and *punctuation symbols*.

lence queries. Remarkably, this does not affect the performance of the symbolic algorithm, which is still better in the total number of membership queries that are posed, including the queries used for testing.

## 9.4 Learning Passwords

In this section, we compare the same four algorithms in a more realistic setting. We let the target language be a set of valid passwords, where the input alphabet is the set of ASCII characters, represented by the integers  $\Sigma = \{0, ..., 127\}$ . A password is a string over  $\Sigma$ , and it is characterized as valid when certain rules are satisfied. Rules of this kind typically concern the length of a password and groups of characters that it may or may not contain.

The ASCII characters are split into groups, see Figure 9.9. Namely, *control characters*, i.e.,  $\{0, \ldots, 31\} \cup \{127\} \subseteq \Sigma$ , which are non printable characters that provide meta-information; *alphabetic* of two types, upper case and lower case letters, that is,  $\{65, \ldots, 90\}$  and  $\{97, \ldots, 122\}$  respectively; *numerals*, i.e.,  $\{48, \ldots, 57\}$ ; and, finally, the *punctuation symbols*, that is the set which contains all remaining characters.

Passwords are used in several types of accounts and often providers set up rules on how passwords should be formulated in order to guaranty security and/or simplicity to the users. Due to this, we are familiar with PIN codes, typically used for phones, credit cards, etc., but also longer passwords that should contain letters and/or punctuation characters, most commonly used for online accounts such as email accounts, bank accounts, etc. Passwords with characters from more than one



Figure 9.10: A performance comparison of the symbolic algorithm (SL) to the three concrete algorithms (LR), (LM) and (L). The four algorithms are evaluated on learning a set of different types of valid passwords: (A) PIN code, (B) easy password, (C) medium password, (D) medium-strong password, and (E) strong password. The average number of states learned by each algorithm for each password type is shown in (a), while (b) illustrates the average number of membership queries used in total.

or two groups of ASCII characters are considered stronger and more secure.

In the sequel, we define five different target languages, each one accepting a different type of valid passwords. Each language is defined by a set of rules, which are of increasing complexity resulting in languages of varying difficulty degree. For all languages selected here, no password is valid when it contains any non-printable character.

The PIN code passwords (A) contains 4 to 8 numbers. The *easy* passwords (B), with a maximum length of 8, contain any printable ASCII character, i.e., numerical, letter, or punctuation character. The *medium* passwords (C) contain from 6 up to 14 characters and use at least one character from two different sub-alphabets of printable characters; punctuation symbols are not allowed. The *medium-strong* passwords (D), with the same length as the medium passwords, require the presence of both lower-case letters and numerals and punctuation characters are allowed. Finally, the *strong* passwords (E), contain at least one character from all different sub-alphabets.

After applying all four algorithms to the different target languages, we obtain the results of 10 runs for each algorithm. A visual comparison of the results is summarized in Figure 9.10, where one can observe that the symbolic algorithm requires fewer MQ's than any other method in all five password types, see Figure 9.10b. The difference increases as the target becomes more complicated and requires more states and transitions. Remarkable is the fact that the symbolic algorithm manages to discover more states, in almost all different password types, and especially when the target language becomes more complicated and is represented by a larger automaton, see Figure 9.10a. In Figure 9.11, we provide two instances of symbolic automata for the passwords of type (D) and (E), which are inferred using the symbolic algorithm.

## 9.5 Learning over the Booleans

In Chapter 7, the symbolic algorithm has been adapted to the case where the input alphabet is a set of Boolean vectors, i.e.,  $\Sigma = \mathbb{B}^n$ . We have implemented the learning algorithm following the algorithms and procedures presented in Section 7.3. To represent the semantics of the symbols, we use binary decision trees (see Section 2.2), where, in order to learn the partitions at each state, we use an optimized version of the CART algorithm (see Section 2.3), which is available in the *scikit-learn* module of machine learning in Python [PVG<sup>+</sup>11]. As a splitting quality measure we use entropy, and best split is chosen. As a termination criterion we use purity, that is, trees are expanded until all leaves are pure and thus the observation table is evidence compatible.



Figure 9.11: Two instances of symbolic automata for the passwords of type (D) and (E), which are learned using the symbolic algorithm.

We test the learning algorithm on two groups of experiments, as in the case of numerical alphabets in Section 9.3. The two groups of experiments are used to evaluate the algorithm's performance when the size of the automaton or the size of the alphabet increases, while other parameters are kept constant. All results are averaged over the total number of experiments for each case, that is, 10 different target automata for each automaton or alphabet size, where the learning algorithm is applied 10 times to the same target.

For the first group of experiments, we use the same fixed minimal automaton structure that appears in Figure 9.5. Partitions are randomly chosen at each state and the alphabet's size ranges from 8 to 32768 letters, that is, 3 to 15 Boolean variables. In the second group of experiments, we fix the alphabet to  $\Sigma = \mathbb{B}^8$  and generate random automata of increasing size, with the number of states ranging from 2 to 50. Transitions are determined by randomly generated pseudo-boolean functions. They are limited to contain at most k = 4 literals per term, which, intuitively, restricts to BDTs of maximal depth k and corresponds to at most  $2^k$  outgoing transitions per state.

As in Section 9.3, we have measured several size and time parameters. Size related, we consider the size of the final observation table, the concrete and symbolic sample size, as well as the number of membership and equivalence queries posed to the teacher.

The major observation of our experimental results for the Boolean alphabets show that, if the complexity of the alphabet partition remains fixed, the number of membership, equivalence queries and sample size remain constant as the alphabet increases. This can be seen in Figure 9.12. Note that, the number of variables increases, but the depth of the trees that are used to define the partitions are bounded



Figure 9.12: The performance of the symbolic algorithm when applied on the Boolean vector alphabet  $\Sigma = \mathbb{B}^n$  where the number of variables n ranges from 3 to 15. Here we illustrate the average number of membership queries, equivalence queries, the size of the final sample, and the execution time as the alphabet size increases.

to have at most depth 4. That is, the variables that play a crucial role in defining a partition at each state is limited to 4.

A growth in the total execution time seems exponential in the number of variables, which corresponds to linear growth in the alphabet size, see the last row of Figure 9.12. Profiling shows that most of the time is spent in testing, indicating that, in the implementation, this depends on the size of the concrete alphabet.

On the other hand, as we see in Figure 9.13, the number of queries, the size of the symbolic sample, the number of hypotheses conjectured, and the total execution time all appear to increase linearly in the number of states in the target automaton. This result does not come as a surprise since the number of BDTs that need to be



Figure 9.13: The performance of the symbolic algorithm when applied on target automata of variable size, where the number of states in the minimal automaton that represents the target ranges from 3 to 50. Here we illustrate the average number of membership queries, equivalence queries, the size of the final sample, and the execution time as the size of the automaton increases.

learned depends on the number of states.

# 9.6 Comparing Boolean Vectors to Numerical Alphabets

#### **Comparison over Alphabet and Target Size**

In this section, we compare the results that appear in Sections 9.3 and 9.5, which are related to the two symbolic algorithms, i.e., the symbolic learning algorithm on numbers (SL) and on the Booleans (SL Booleans), respectively. All results are illustrated in Figure 9.14.



Figure 9.14: A comparison on the performance of the symbolic algorithm when applied on interval alphabets (SL) and Boolean vector alphabets (SL Booleans). The rows show 1) the average number of membership queries, 2) the average size of the symbolic sample, and 3) the average number of equivalence queries posed in total by each algorithm, respectively.

In the first group of experiments (first column), the same fixed underlying structure is used to generate the target languages while the alphabet size grows. We observe that the Boolean algorithm clearly outperforms the numerical algorithm in terms of membership and equivalence queries. In the second group of experiments (second column), where we compare over different automaton sizes, the numerical algorithms use a fixed alphabet of 150 letters, while the Boolean algorithm uses a 256-letter alphabet. Although the absolute numbers differ for the three algorithms, it is worth noting that the growth of the number of queries, as well as the size of the symbolic sample is linear with respect to the size of the automaton in both symbolic algorithms. We observe again that in the case of Boolean vectors, the number of membership queries is much smaller than for numerical alphabets.

Finally, in the last plot of the first group of experiments (first column, third row of Figure 9.14), we can see that fewer EQ's are needed in general for the case of Boolean vector alphabets. A possible explanation for this superiority can rely on the different restrictions that were applied in defining the partitions in the target languages. To explain this better, we give some insights on the random generator that is used to define the target languages. For interval partitions, a restriction to



Figure 9.15: The performance of the symbolic algorithm on learning languages that represent valid passwords. Two adaptations of the algorithm are compared on two different representations of the input alphabet, the set of ASCII characters. In the first, the alphabet is a set of integers (SL), and in the second, it is a set of Boolean vectors (SLBool).

k intervals implies that their size itself grows linearly with the alphabet size and thus one may expect more errors and counter-examples when the boundaries are approximated. This is often the case, especially when some intervals are small compared to others. On the other hand, for partitions on the Boolean cube, the decision trees are restricted to have at most depth k. This, in turn, restricts the number of variables that interact in a product alphabet, and it gives a lower-bound on the size of the Boolean sub-cubes appearing in the partition. This lower bound grows linearly with the size of the alphabet and such partitions are easier to learn.

#### **Comparing on Learning Passwords**

In Section 9.4, the symbolic algorithm is used to learn languages that represent valid passwords. The set of ASCII characters, in their decimal representation, is used as the input alphabet. The alphabet  $\Sigma = \{0, ..., 127\}$  is partitioned into intervals, where intervals represent groups of characters, e.g., *digits, lowercase* 



Figure 9.16: The BDT representing the partition of the ASCII characters into groups of characters, i.e., *digits, lowercase letters, punctuation symbols*, etc.

#### letters, special symbols, etc.

In this section, we apply the symbolic algorithm on the same set of target languages, but this time, we use as an input alphabet the set  $\Sigma = \mathbb{B}^7$  of Boolean vectors, which corresponds to the set of ASCII characters' binary representation. For instance, the letter 'a' corresponds to the integer 97 and to the Boolean vector (1100001). In contrast to the interval partitions, which appears to be convenient when classifying ASCII into groups of characters, the Boolean vector adaptation requires a complicated and deep BDT to determine all partition blocks in the alphabet, see Figure 9.16.

Figure 9.15 shows the performance of the two algorithms. We compare them on the number of queries they use as well as the size of their final conjectures in terms of number of states and transitions. As we can observe, due to the higher complexity on the partitions for the Boolean case, the algorithm requires more evidences to define the partitions well. Hence, more equivalence and membership queries are used.

The two algorithms terminate with conjectures that contain similar numbers of states. A small advantage is observed in the case of Boolean Vectors as valid passwords become more complicated. Although, this is not surprising when we consider the increased number of equivalence queries the algorithm uses.

Finally, concerning the size of the symbolic alphabets, in the case of intervals more symbols are used. This is due to the fact, that partitions and semantics are determined in a slightly different way in the two adaptations of the algorithm. That is, BDTs use the same symbol to represent multiple leaf nodes allowing at most one transition that leads from a state q to a state q'. In contrast, in the case of the intervals, the semantics are restricted to preserve convexity, and thus, multiple transitions from q to q' are allowed.

We conclude that both algorithms can be used when the alphabet admits both representations resulting equivalent automata. However, the complexity on the partitions should be taken into account, as it affects the number of resources the algorithm uses.

## 9.7 Conclusions

In this chapter, we have experimented with the symbolic algorithm and applied it to several case studies. Moreover, we compared the symbolic algorithm to other known concrete algorithms, restricting the alphabet to be finite. The symbolic algorithm always terminates faster and with better results in the PAC setting. Our algorithm, almost always, discovers more states and since we consider that the resulting automaton is minimal, it means that the conjecture is close to the target's representation. The size of the sample and the number of membership queries are notably lower than any other algorithm, since the words are partially queried and evidences are used to represent larger set of letters. Naturally, this approach increases the number of equivalence queries our algorithm requires, however, this does not seem to influence the time that the algorithm needs until termination.

Comparing the symbolic algorithms on different alphabets, we observe that when applied to Boolean vectors, it performs better, partially due to the different nature of partition complexity restrictions on decision trees and intervals, and also the fact that we use a better method to learn partition boundaries. Learning BDTs allows a smaller boundary error than the simple splitting we use in the interval case. The latter can be ameliorated by using a more precise method, such as binary search, for detecting the boundaries.

# **Conclusions and Future Work**

In this thesis we presented an algorithmic scheme for learning languages over large alphabets. We aimed at languages accepted by symbolic automata with a modest number of states and transitions. The transitions are guarded by simple constraints on the alphabet, which by itself can be arbitrarily large. The constraints form a partition of the alphabet in each state and hence the automaton is deterministic. Our framework can be characterized by the following features:

- 1. It is based on a clean and general definition of the relation between the concrete and symbolic alphabets, via concepts such as alphabet partitions, evidences and representatives;
- 2. It separates the sequential aspects of learning (modifying the automaton structure) and the learning of partition boundaries, which is specific to the alphabet;
- 3. It treats the modification of alphabet partitions in a rigorous way based on the concept of evidence compatibility. This guarantees that no superfluous symbols are introduced;
- It can work with or without a helpful teacher, where in the latter case, equivalence checking is done by random sampling and exact learnability is replaced by PAC;
- 5. It is modular, separating the generic aspects from those that are alphabet specific, thus providing for a relatively easy adaptation to new alphabets.

We have implemented the testing-based algorithm for two classes of alphabets, numerical (sub-intervals of  $\mathbb{R}$  or  $\mathbb{N}$ ) and Boolean ( $\mathbb{B}^n$ ). For the former, symbolic letters were associated with sub-intervals while in the latter we used decision trees to define the alphabet partitions in every state. For both cases we ran some initial experiments on mostly synthetic examples, to see how the empirical performance of the algorithm varies with the size of the automaton and the alphabets. The preliminary results were encouraging and the symbolic algorithm outperforms concrete learning algorithms over the same finite alphabets. There is still much to do in tuning and optimizing our algorithm.

What is still missing is a convincing class of real-world applications that benefits from such algorithms. In the numerical domain, these should be mechanisms where discrete transitions are taken based on threshold crossings of numerical variables without remembering their precise values. In the Boolean domain, one needs to find applications in the formal specification of large complex systems with many components (digital circuit, distributed multi-agent systems). Hopefully such specifications could be expressible by symbolic automata where the complexity can be confined to the alphabet partitions while the number of states does not explode.

A natural future extension of the algorithm is to consider alphabets which are subsets of  $\mathbb{N}^n$  and  $\mathbb{R}^n$ . Preliminary work in this direction has been reported in Section 6.2, but it used a very restricted type of monotone partitions in order to keep the notion of a minimal counter-example meaningful. In such a future algorithm one can use more general partitions, represented by regression trees, a generalization of decision trees to numerical domains.

We feel that the application of our framework for learning over Boolean alphabets is very promising and is worth being explored much further. It combines well-known concepts and algorithms from traditional machine learning (decision trees) with automaton learning. As such, it can suggest an alternative to other machine learning approaches (recurrent neural learning, decision tree learning over variables that correspond to different time indices, auto-correlation) used to learn from data-rich and time-dependent examples. Automata bring some sophistication and economy in representing sets of time series and might bring a genuine practical contribution to machine learning. This will require, however, some changes in the automaton learning approach, most notably to allow noise in the examples and to relax the active learning framework where the learner can fully dictate the set of observed examples.

# Bibliography

[ABD15]	Eugene Asarin, Nicolas Basset, and Aldric Degorre. Entropy of reg- ular timed languages. <i>Information and Computation</i> , 241:142–176, 2015. ( <i>Cited on page 86</i> .)
[AJU10]	Fides Aarts, Bengt Jonsson, and Johan Uijen. Generating models of infinite-state communication protocols using regular inference with abstraction. In <i>IFIP International Conference on Testing Software and Systems</i> , pages 188–204. Springer, 2010. ( <i>Cited on page 4.</i> )
[Ang87]	Dana Angluin. Learning regular sets from queries and counterexamples. <i>Information and Computation</i> , 75(2):87–106, 1987. ( <i>Cited on pages 1, 21, 23, 28, 47, 90, 97, 99, and 101.</i> )
[Ang88]	Dana Angluin. Queries and concept learning. <i>Machine learning</i> , 2(4):319–342, 1988. ( <i>Cited on page 90</i> .)
[BB13]	Matko Botinčan and Domagoj Babić. Sigma*: Symbolic learning of input-output specifications. In <i>POPL</i> , pages 443–456. ACM, 2013. ( <i>Cited on page 48.</i> )
[BFSO84]	Leo Breiman, Jerome Friedman, Charles J Stone, and Richard A Olshen. <i>Classification and regression trees</i> . CRC press, 1984. ( <i>Cited on pages 16, 17, 20, and 78.</i> )
[BJR06]	Therese Berg, Bengt Jonsson, and Harald Raffelt. Regular inference for state machines with parameters. In <i>FASE</i> , volume 3922 of <i>LNCS</i> , pages 107–121. Springer, 2006. ( <i>Cited on page 48</i> .)
[BKR96]	Morten Biehl, Nils Klarlund, and Theis Rauhe. Algorithms for guided tree automata. In <i>International Workshop on Implementing Automata</i> , pages 6–25. Springer, 1996. ( <i>Cited on page 33.</i> )

[BLP10]	Michael Benedikt, Clemens Ley, and Gabriele Puppis. What you must remember when processing data words. In <i>AMW</i> , volume 619 of <i>CEUR Workshop Proceedings</i> , 2010. ( <i>Cited on page 47.</i> )
[BMS <sup>+</sup> 06]	Mikolaj Bojanczyk, Anca Muscholl, Thomas Schwentick, Luc Segoufin, and Claire David. Two-variable logic on words with data. In <i>21st Annual IEEE Symposium on Logic in Computer Science</i> ( <i>LICS'06</i> ), pages 7–16. IEEE, 2006. ( <i>Cited on page 33</i> .)
[BR05]	Therese Berg and Harald Raffelt. 19 Model checking. In <i>Model-Based Testing of Reactive Systems: Advanced Lectures</i> , volume 3472 of <i>LNCS</i> , pages 557–603. Springer, 2005. ( <i>Cited on pages 2, 21, and 29.</i> )
[Brz62]	Janusz A Brzozowski. Canonical regular expressions and minimal state graphs for definite events. <i>Mathematical theory of Automata</i> , 12(6):529–561, 1962. ( <i>Cited on page 10</i> .)
[CHJS16]	Sofia Cassel, Falk Howar, Bengt Jonsson, and Bernhard Steffen. Ac- tive learning for extended finite state machines. <i>Formal Aspects of</i> <i>Computing</i> , 28(2):233–263, 2016. ( <i>Cited on page</i> 47.)
[CSS <sup>+</sup> 10]	Chia Yuan Cho, Eui Chul Richard Shin, Dawn Song, et al. Infer- ence and analysis of formal models of botnet command and control protocols. In <i>Proceedings of the 17th ACM conference on Computer</i> <i>and communications security</i> , pages 426–439. ACM, 2010. ( <i>Cited</i> <i>on page 4.</i> )
[DD17]	Samuel Drews and Loris D'Antoni. Learning symbolic automata. In <i>International Conference on Tools and Algorithms for the Construc-</i> <i>tion and Analysis of Systems</i> , pages 173–189. Springer, 2017. ( <i>Cited on page</i> 48.)
[DlH10]	Colin De la Higuera. <i>Grammatical inference: learning automata and grammars</i> . Cambridge University Press, 2010. ( <i>Cited on pages 21 and 22.</i> )
[DV14]	Loris D'Antoni and Margus Veanes. Minimization of symbolic automata. In <i>POPL</i> , pages 541–554. ACM, 2014. ( <i>Cited on pages 34</i> , <i>36</i> , <i>38</i> , <i>and 39</i> .)

[E88] Utgoff Paul E. ID5: an incremental ID3. In John E. Laird, editor, Machine Learning, Proceedings of the Fifth International Conference on Machine Learning, Ann Arbor, Michigan, USA, June 12-14, 1988, pages 107–120. Morgan Kaufmann, 1988. (*Cited on pages 16 and 19.*)

- [GKS10] Orna Grumberg, Orna Kupferman, and Sarai Sheinvald. Variable automata over infinite alphabets. In *International Conference on Language and Automata Theory and Applications*, pages 561–572. Springer, 2010. (*Cited on page 33.*)
- [Gol67] E Mark Gold. Language identification in the limit. *Information and Control*, 10(5):447–474, 1967. (*Cited on page 23.*)
- [Gol72] E Mark Gold. System identification via state characterization. *Automatica*, 8(5):621–636, 1972. (*Cited on pages 10 and 23*.)
- [Gol78] E Mark Gold. Complexity of automaton identification from given data. *Information and control*, 37(3):302–320, 1978. (*Cited on page 23.*)
- [HJJ<sup>+</sup>95] Jesper G Henriksen, Ole JL 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. (*Cited on page 33*.)
- [HMU06] John E Hopcroft, Rajeev Motwani, and Jeffrey D Ullman. Introduction to Automata Theory, Languages, and Computation (3rd Edition). Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 2006. (Cited on pages 7 and 37.)
- [Hop71] John E Hopcroft. An *n logn* algorithm for minimizing states in a finite automaton. Technical report, DTIC Document, 1971. (*Cited on page 10.*)
- [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. (Cited on page 47.)
- [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. (Cited on page 48.)
- [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. (Cited on pages 34, 36, and 37.)

[IHS13]	Malte Isberner, Falk Howar, and Bernhard Steffen. Inferring automata with state-local alphabet abstractions. In <i>NASA Formal Methods</i> , volume 7871 of <i>LNCS</i> , pages 124–138. Springer, 2013. ( <i>Cited on page 48.</i> )
[IS14]	Malte Isberner and Bernhard Steffen. An abstract framework for counterexample analysis in active automata learning. In <i>ICGI</i> , pages 79–93, 2014. ( <i>Cited on page 29</i> .)
[KF70]	Andreĭ Nikolaevich Kolmogorov and Sergeĭ Vasilevich Fomin. <i>In-</i> <i>troductory real analysis</i> . New York: Dover Publications Inc., 1970. ( <i>Cited on page 10</i> .)
[KF94]	Michael Kaminski and Nissim Francez. Finite-memory automata. <i>Theoretical Computer Science</i> , 134(2):329–363, 1994. ( <i>Cited on pages 33 and 47.</i> )
[LGJ07]	Tristan Le Gall and Bertrand Jeannet. Lattice automata: A repre- sentation for languages on infinite alphabets, and some applications to verification. In <i>International Static Analysis Symposium</i> , pages 52–68. Springer, 2007. ( <i>Cited on page 33</i> .)
[LP97]	Harry R Lewis and Christos H Papadimitriou. <i>Elements of the Theory of Computation</i> . Prentice Hall PTR, 1997. ( <i>Cited on page 7.</i> )
[MM14]	Oded Maler and Irini-Eleftheria Mens. Learning regular languages over large alphabets. In <i>TACAS</i> , volume 8413 of <i>LNCS</i> , pages 485–499. Springer, 2014. ( <i>Cited on page 5.</i> )
[MM15]	Oded Maler and Irini-Eleftheria Mens. Learning regular languages over large ordered alphabets. <i>Logical Methods in Computer Science</i> ( <i>LMCS</i> ), 11(3), 2015. ( <i>Cited on page 5.</i> )
[MM17]	Oded Maler and Irini-Eleftheria Mens. A generic algorithm for learning symbolic automata from membership queries. In <i>Models, Algorithms, Logics and Tools</i> , pages 146–169. Springer, 2017. ( <i>Cited on page 5.</i> )
[MN04]	Oded Maler and Dejan Nickovic. Monitoring temporal properties of continuous signals. In <i>Formal Techniques, Modelling and Analysis of Timed and Fault-Tolerant Systems</i> , pages 152–166. Springer, 2004.

(Cited on page 3.)

[MN10]	Karl Meinke and Fei Niu. A learning-based approach to unit testing of numerical software. In <i>IFIP International Conference on Testing Software and Systems</i> , pages 221–235. Springer, 2010. ( <i>Cited on page 4.</i> )
[MNP08]	Oded Maler, Dejan Nickovic, and Amir Pnueli. Checking temporal properties of discrete, timed and continuous behaviors. In <i>Pillars of computer science</i> , pages 475–505. Springer, 2008. ( <i>Cited on page 3.</i> )
[Moo56]	Edward F Moore. Gedanken-experiments on sequential machines. In <i>Automata studies</i> , volume 34 of <i>Annals of Mathematical Studies</i> , pages 129–153. Princeton, 1956. ( <i>Cited on pages 10 and 23</i> .)
[MP95]	Oded Maler and Amir Pnueli. On the learnability of infinitary regular sets. <i>Information and Computation</i> , 118(2):316–326, 1995. ( <i>Cited on pages 27, 28, 97, 99, and 101.</i> )
[Mur12]	Kevin P Murphy. <i>Machine Learning: A Probabilistic Perspective</i> . The MIT Press, 1 edition, 2012. ( <i>Cited on page 17.</i> )
[Ner58]	Anil Nerode. Linear automaton transformations. <i>Proceedings of the American Mathematical Society</i> , 9(4):541–544, 1958. ( <i>Cited on pages 1, 9, and 23.</i> )
[Nie03]	Oliver Niese. <i>An integrated approach to testing complex systems</i> . PhD thesis, Universität Dortmund, 2003. ( <i>Cited on page 4.</i> )
[NSV04]	Frank Neven, Thomas Schwentick, and Victor Vianu. Finite state machines for strings over infinite alphabets. <i>ACM Transactions on Computational Logic (TOCL)</i> , 5(3):403–435, 2004. ( <i>Cited on page 33.</i> )
[OG92]	Jose Oncina and Pedro Garcia. Identifying regular languages in polynomial time. In <i>Advances in Structural and Syntactic Pattern Recognition</i> , volume volume 5 of <i>Series in Machine Perception and Artificial Intelligence</i> , pages 99–108. World Scientific, 1992. ( <i>Cited on page 23.</i> )
[PVG <sup>+</sup> 11]	F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Van- derplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. <i>Jour- nal of Machine Learning Research</i> , 12:2825–2830, 2011. ( <i>Cited on</i> <i>page 104.</i> )

[PW93]	Leonard Pitt and Manfred K Warmuth. The minimum consistent DFA problem cannot be approximated within any polynomial. <i>Journal of the ACM (JACM)</i> , 40(1):95–142, 1993. ( <i>Cited on page 23.</i> )
[Qui86]	J Ross Quinlan. Induction of decision trees. <i>Machine learning</i> , 1(1):81–106, 1986. ( <i>Cited on pages 16 and 78.</i> )
[Qui93]	J Ross Quinlan. <i>C4.5: Programs for Machine Learning</i> . Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1993. ( <i>Cited on pages 16 and 20.</i> )
[RMSM09]	Harald Raffelt, Maik Merten, Bernhard Steffen, and Tiziana Mar- garia. Dynamic testing via automata learning. <i>International jour-</i> <i>nal on software tools for technology transfer</i> , 11(4):307–324, 2009. ( <i>Cited on page 4.</i> )
[RS93]	Ronald L Rivest and Robert E Schapire. Inference of finite automata using homing sequences. <i>Information and Computation</i> , 103(2):299–347, 1993. ( <i>Cited on pages 27, 29, 97, 99, and 101.</i> )
[Seg06]	Luc Segoufin. Automata and logics for words and trees over an infinite alphabet. In <i>International Workshop on Computer Science Logic</i> , pages 41–57. Springer, 2006. ( <i>Cited on page 33.</i> )
[SF86]	Jeffrey C Schlimmer and Douglas Fisher. A case study of incremen- tal concept induction. In AAAI, pages 496–501, 1986. ( <i>Cited on</i> pages 16 and 19.)
[Sha08]	Muzammil Shahbaz. <i>Reverse Engineering Enhanced State Models of Black Box Components to support Integration Testing.</i> PhD thesis, Grenoble Institute of Technology, 2008. ( <i>Cited on page 4.</i> )
[Sip06]	Michael Sipser. <i>Introduction to the Theory of Computation</i> , volume 2. Thomson Course Technology Boston, 2006. ( <i>Cited on page</i> 7.)
[SL07]	Guoqiang Shu and David Lee. Testing security properties of protocol implementations-a machine learning based approach. In 27th International Conference on Distributed Computing Systems (ICDCS'07), pages 25–25. IEEE, 2007. (Cited on page 4.)
[SVVV14]	Rick Smetsers, Michele Volpato, Frits Vaandrager, and Sicco Ver- wer. Bigger is not always better: on the quality of hypotheses in active automata learning. In <i>International Conference on Grammat</i> -

ical Inference, pages 167-181, 2014. (Cited on page 85.)

[UBC97]	Paul E Utgoff, Neil C Berkman, and Jeffery A Clouse. Decision tree
	induction based on efficient tree restructuring. Machine Learning,
	29(1):5–44, 1997. (Cited on page 19.)

- [Utg89] Paul E Utgoff. Incremental induction of decision trees. *Machine learning*, 4(2):161–186, 1989. (*Cited on pages 16, 19, 20, and 78.*)
- [Val84] Leslie G Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984. (*Cited on page 90.*)
- [VBDM10] Margus Veanes, Nikolaj Bjørner, and Leonardo De Moura. Symbolic automata constraint solving. In International Conference on Logic for Programming Artificial Intelligence and Reasoning, volume 6397, pages 640–654. Springer, 2010. (Cited on pages 4 and 34.)
- [VDHT10] Margus Veanes, Peli De Halleux, and Nikolai Tillmann. Rex: Symbolic regular expression explorer. pages 498–507, 2010. (*Cited on pages 4, 33, 34, and 37.*)
- [Vea13] Margus Veanes. Applications of symbolic finite automata. In International Conference on Implementation and Application of Automata, volume 7982, pages 16–23. Springer, 2013. (Cited on pages 4, 33, and 34.)
- [VGDHT09] Margus Veanes, Pavel Grigorenko, Peli De Halleux, and Nikolai Tillmann. Symbolic query exploration. In *International Conference on Formal Engineering Methods*, pages 49–68. Springer, 2009. (Cited on page 34.)
- [VHL<sup>+</sup>12] Margus Veanes, Pieter Hooimeijer, Benjamin Livshits, David Molnar, and Nikolaj Björner. Symbolic finite state transducers: algorithms and applications. In *POPL*, pages 137–150. ACM, 2012. (*Cited on page 34.*)
- [VNG01] Gertjan Van Noord and Dale Gerdemann. Finite state transducers with predicates and identities. *Grammars*, 4(3):263–286, 2001. (*Cited on pages 4, 33, 34, 36, and 37.*)
- [VTDH10] Margus Veanes, Nikolai Tillmann, and Jonathan De Halleux. Qex: Symbolic SQL query explorer. In International Conference on Logic for Programming Artificial Intelligence and Reasoning, pages 425– 446. Springer, 2010. (Cited on page 34.)

[Wat96]	Bruce W Watson. Implementing and using finite automata toolk-
	its. Natural Language Engineering, 2(04):295-302, 1996. (Cited on
	page 33.)
[WBAH08]	Neil Walkinshaw, Kirill Bogdanov, Shaukat Ali, and Mike Hol-
	combe. Automated discovery of state transitions and their func-
	tions in source code. Software Testing, Verification and Reliability,
	18(2):99–121, 2008. (Cited on page 4.)
[YBCI08]	Fang Yu, Tevfik Bultan, Marco Cova, and Oscar H Ibarra. Symbolic
	string verification: An automata-based approach. In International
	SPIN Workshop on Model Checking of Software, pages 306-324.

Springer, 2008. (Cited on pages 33 and 34.)

# Learning Regular Languages over Large Alphabets

# Irini-Eleftheria Mens

Thèse dirigée par Oded Maler

Learning regular languages is a branch of machine learning, a domain which has proved useful in many areas, including artificial intelligence, neural networks, data mining, verification, etc. In addition, interest in languages defined over large and infinite alphabets has increased in recent years. Although many theories and properties generalize well from the finite case, learning such languages is not an easy task. As the existing methods for learning regular languages depend on the size of the alphabet, a straightforward generalization in this context is not possible.

In this thesis, we present a generic algorithmic scheme that can be used for learning languages defined over large or infinite alphabets, such as bounded subsets of  $\mathbb{N}$  or  $\mathbb{R}$  or Boolean vectors of high dimensions. We restrict ourselves to the class of languages accepted by deterministic symbolic automata that use predicates to label transitions, forming a finite partition of the alphabet for every state.

Our learning algorithm, an adaptation of Angluin's  $L^*$ , combines standard automaton learning by state characterization, with the learning of the static predicates that define the alphabet partitions. We use the online learning scheme, where two types of queries provide the necessary information about the target language. The first type, membership queries, answer whether a given word belongs or not to the target. The second, equivalence queries, check whether a conjectured automaton accepts the target language and provide a counter-example otherwise.

We study language learning over large or infinite alphabets within a general framework but our aim is to provide solutions for particular concrete instances. For this, we focus on the two main aspects of the problem. Initially, we assume that equivalence queries always provide a counter-example which is minimal in the length-lexicographic order when the conjecture automaton is incorrect. Then, we drop this "strong" equivalence oracle and replace it by a more realistic assumption, where equivalence is approximated by testing queries, which use sampling on the set of words. Such queries are not guaranteed to find counter-examples and certainly not minimal ones. In this case, we obtain the weaker notion of PAC (probably approximately correct) learnability and learn an approximation of the target language. All proposed algorithms have been implemented and their performance, as a function of automaton and alphabet size, has been empirically evaluated.

L'apprentissage de langages réguliers est une branche de l'apprentissage automatique qui s'est révélée utile dans de nombreux domaines tels que l'intelli-gence artificielle, les réseaux de neurones, l'exploration de données, la vérification, etc. De plus, l'intérêt dans les langages définis sur des alphabets infinis ou de grande taille s'est accru au fil des années. Même si plusieurs propriétés et théories se généralisent à partir du cas fini, l'apprentissage de tels langages est une tâche difficile. En effet, dans ce contexte, l'application naïve des algorithmes d'apprentissage traditionnel n'est pas possible.

Dans cette thèse, nous présentons un schéma algorithmique général pour l'ap-prentissage de langages définis sur des alphabets infinis ou de grande taille, comme par exemple des sous-ensembles bornés de  $\mathbb{N}$  or  $\mathbb{R}$  ou des vecteurs booléens de grandes dimensions. Nous nous restreignons aux classes de langages qui sont acceptés par des automates déterministes symboliques utilisant des prédicats pour définir les transitions, construisant ainsi une partition finie de l'alphabet pour chaque état. Notre algorithme d'apprentissage, qui est une adaptation du  $L^*$  d'Angluin, combine l'apprentissage

Notre algorithme d'apprentissage, qui est une adaptation du  $L^*$  d'Angluin, combine l'apprentissage classique d'un automate par la caractérisation de ses états, avec l'apprentissage de prédicats statiques définissant les partitions de l'alphabet. Nous utilisons l'apprentissage incrémental avec la propriété que deux types de requêtes fournissent une information suffisante sur le langage cible. Les requêtes du premier type sont les requêtes d'appartenance, qui permettent de savoir si un mot proposé appartient ou non au langage cible. Les requêtes du second type sont les requêtes d'équivalence, qui vérifient si un automate proposé accepte le langage cible; dans le cas contraire, un contre-exemple est renvoyé.

Nous étudions l'apprentissage de langages définis sur des alphabets infinis ou de grande tailles dans un cadre théorique et général, mais notre objectif est de proposer des solutions concrètes pour un certain nombre de cas particuliers. Ensuite, nous nous intéressons aux deux principaux aspects du problème. Dans un premier temps, nous supposerons que les requêtes d'équivalence renvoient toujours un contreexemple minimal pour un ordre de longueur-lexicographique quand l'automate proposé est incorrect. Puis dans un second temps, nous relâchons cette hypothèse forte d'un oracle d'équivalence, et nous la remplaçons avec une hypothèse plus réaliste où l'équivalence est approchée par un test sur les requêtes ne garantit pas l'obtention de contre-exemples, et par conséquent de contre-exemples minimaux. Nous obtenons alors une notion plus faible d'apprent-issage PAC (*Probably Approximately Correct*), permettant l'apprentissage d'une approximation du langage cible. Tout les algorithmes ont été implémentés, et leurs performances, en terme de construction d'automate et de taille d'alphabet, ont été évaluées empiriquement.

Communauté

