
SAT based learning Algorithms

Author:
Rajarshi Roy

Supervisor:
Dr. Daniel Neider

A thesis submitted for the degree of

*Masters degree in Computer Science
at Chennai Mathematical Institute*

June 15, 2019



MAX PLANCK INSTITUTE
FOR SOFTWARE SYSTEMS



CHENNAI
MATHEMATICAL
INSTITUTE

Abstract

In this thesis, we present algorithms for learning regular-expressions, omega-regular expressions, LTL formulas and PSL formulas, from a given sample of positive and negative examples. The precise learning problem studied in this thesis—is to find a model consistent with a given sample, that is, a model that accepts the positive examples and rejects the negative examples in the sample. SAT-solver based algorithms have been efficient in solving these problems in certain models such as Linear Temporal Logic (LTL).

Property Specification Language (PSL) is an extension of LTL, with increased expressive power. Now, PSL relies heavily on regular expressions. Hence, it is necessary to look at the learning problem for regular expressions and ω -regular expressions for solving the learning problem of PSL. While it is straightforward to apply the existing techniques for solving the problem for regular expressions, it is challenging to translate the same ideas to ω -regular expression. This is because the sample given in the case of ω -regular expressions contains infinite words. The same can be said about PSL formulas as well.

Acknowledgements

I would like to thank my supervisor, Daniel Neider, for presenting me with such a wonderful problem for my master thesis and being a great mentor while I worked on it. I would also, express my gratitude to my professors at CMI for teaching me courses, which made me competent enough to tackle a problem of this kind. Last but not the least, I would like to thank all my colleagues and friends at MPI and CMI for their constant assistance in helping me complete my thesis.

Contents

1 Introduction	3
1.1 Possible Applications	4
2 Preliminaries	5
3 Learning Regular Expressions	7
3.1 Regular Expressions	7
3.2 The Learning Problem	9
3.2.1 Proof of correctness	10
4 Learning ω-Regular Expressions	13
4.1 ω -Regular expression	13
4.2 The Learning Problem	16
5 Learning of Linear Temporal Language	19
5.1 Linear Temporal Logic	19
5.2 The Learning Problem	20
6 Learning Program Specification Language	22
6.1 Property Specification Language	22
6.2 The Learning Problem	23
7 Conclusion	25
A Proofs from Chapter 4	26
B Proof of Claim 2:	27
Bibliography	31

Chapter 1

Introduction

Constructing a model of system behavior from observation traces, in a form that is understandable and meaningful to a human, is central to human interpretability of complex systems. Naturally, ability to learn human-interpretable models is quite crucial in order to understand general behaviour of a system. But, various scenarios require us to come up with different models that suitably represent it. To this end, we have presented learning techniques for a variety of models, which could find applications in diverse spheres.

The techniques used in this thesis have been motivated by ideas from the work of Neider and Gavran[6], where they present novel algorithms for learning formulas in Linear Temporal Logic (LTL)[7]. LTL has a natural syntax useful in specifying properties of programs for verification. But, there are certain program behaviour which cannot be characterized by LTL formulas. For instance, the property "Event e_1 happens after even number of steps of e_2 " cannot be expressed in LTL[8]. This led us to look at Property Specification Language (PSL)[4], which extends the expressive power of LTL and captures a wider range of program behaviour. PSL is the IEEE standard for temporal logic (IEEE). In terms of expressiveness, PSL is as powerful as ω -regular expressions while LTL can express only star-free ω -regular expressions. Consequently, formulating learning algorithms for PSL formulas is the natural next step.

The increased expressive power of PSL is due to the incorporation of regular expressions in its syntax. Thus, in an attempt to solve the aforementioned problem, it is necessary to understand how the same problem could be solved when the model being learnt is a regular expression. Moreover, since PSL formula argue about infinite strings, it is also useful to devise algorithms for learning ω -regular expression as well.

The main goal of this thesis is to devise algorithms for learning high-level descriptions of system behaviour provided in the form of a sample consisting of some classified examples. Moreover, the description needs to be succints, it is not quite challenging to construct arbitrary large models which adhere to the given examples. To be precise, given a sample S consisting of two finite sets of positive and negative examples, we want to learn the 'smallest' model \mathcal{M} that is consistent with S —in the sense that all positive examples satisfy \mathcal{M} , whereas all negative examples violate \mathcal{M} . The model being learnt could be—a regular expression (Chapter 3), ω -regular expressions (Chapter 4), LTL fomulas (Chapter 5) or PSL formulas (Chapter 6).

The general strategy that has been employed for solving the learning problems is encoding the problem in a series of satisfiability formulas and then using highly optimized SAT-solver to search for a solution.

1.1 Possible Applications

Most of the learning problems for various models emerge out of some software verification problem. But, by extending the ideas to solve learning problems for regular expressions and ω -regular expressions, there could be potential applications in diverse fields. This is because patterns and sequences are abundant in nature and identifying them could prove to be quite crucial. Here, we present some possible applications.

- *Bug Detection:* From a buggy program, it is possible to procure the data of desirable and undesirable program behaviour. In order to have a deeper understanding of the root of the bug, it could be helpful to apply the learning algorithms using the data of the program traces as the sample. High level description of the bug, could benefit bug-detection. Since, there exists learning algorithms for several models, the model which captures the bug more appropriately, can be learned.
- *Classifying biological sequences:* Consider the following table which represents biological sequences that have been associated with a group of diseases called amyloidosis.

String	Class
stviil	positive
ktvive	negative
stviie	positive
st piie	negative

Now, from this data it can be deduced that any sequence satisfying the regular expression $\Sigma^*(stv)\Sigma^*$ could potentially be harmful. Hence, identifying such a pattern could be beneficial in this case. This problem has been discussed in [1].

Chapter 2

Preliminaries

Words and Languages An alphabet Σ is a nonempty, finite set of symbols. A *finite word* $u = a_0 \cdots a_{n-1}$ is a finite sequence of symbols $a_i \in \Sigma$ for $i \in \{0, \dots, n-1\}$. The empty sequence is called the empty word and is denoted by ε . The length of a word u , denoted by $|u|$, is the number of symbols in u . For two words $u = a_0 \cdots a_{m-1}$ and $v = b_0 \cdots b_{n-1}$, the concatenation of u and v is the word $u \circ v = uv = a_0 \cdots a_{m-1} b_0 \cdots b_{n-1}$. The set of all finite words over an alphabet Σ is denoted by Σ^* . Let $w[i, j]$ refer to the subword of w starting at position i and ending at position $j - 1$. Also, $w[i, i) = \varepsilon$ in this definition. Assume, here that the indexing of word starts from position 0.

An *infinite word* $\alpha = a_0 a_1 \dots$ is an infinite sequence of symbols $a_i \in \Sigma$ for each $i \geq 0$. Given a word $v \in \Sigma^+$, the infinite repetition of v is the infinite word $v^\omega = vv \dots \in \Sigma^\omega$. We say that a word $\alpha \in \Sigma^\omega$ is *ultimately periodic* if it can be written as wv^ω with $w \in \Sigma^*$ and $v \in \Sigma^+$. Moreover, let $\alpha[i, j) = a_i a_{i+1} \cdots a_{j-1}$ be the finite infix of the infinite word $\alpha = a_0 a_1 \dots \in \Sigma^\omega$. Similarly, let $\alpha[i, \infty)$ be the infinite suffix $a_i a_{i+1} \dots \in \Sigma^\omega$.

A subset $L \subseteq \Sigma^\omega$ of infinite words is called a ω -language. Concatenation can be also be extended to account for ω -languages as well and has been referred to as \circ_ω in this thesis. But, $L \circ_\omega L'$ is well-defined only when L is a language of finite words and L' is a ω -language. Also, the ω -iteration of a language $L \subseteq \Sigma^*$ is the ω -language $L^\omega = \{w_1 w_2 w_3 \cdots \mid w_i \in L \setminus \{\varepsilon\}\}$.

Deterministic Finite Automata A *Deterministic Finite Automaton* (DFA) is a tuple $\mathcal{A} = (\Sigma, Q, q_0, \delta, F)$, where Σ is a finite alphabet, Q is a non-empty, finite set of states, $q_0 \in Q$ is the initial state, $\delta : Q \times \Sigma \rightarrow Q$ is the transition function, and $F \subseteq Q$ is the set of final (or accepting) states. We can extend δ to a function $\delta : Q \times \Sigma^* \rightarrow Q$ by $\delta(q, \varepsilon) = q$ and $\delta(q, wa) = \delta(\delta(q, w), a)$ for all $a \in \Sigma$ and $w \in \Sigma^*$. The language accepted by \mathcal{A} is the set $\mathcal{L}(\mathcal{A}) = \{w \in \Sigma^* \mid \delta(q_0, w) \in F\}$. The run of a DFA on a word $w = a_0 \cdots a_n$ is defined as a sequence of states of the DFA $q_0, q_1, \dots, q_n, q_{n+1}$, such that $\delta(q_i, a_{i+1}) = q_{i+1}$ for all $i = 0 \cdots n$. A word w is accepted if the run of the DFA on w terminates in a final state. A language L is regular iff there exists a DFA \mathcal{A} such that $\mathcal{L}(\mathcal{A}) = L$.

Propositional Boolean Logic: Let AP be a set of propositional variables, which take Boolean values from $\mathbb{B} = \{0, 1\}$ (0 representing false and 1 representing true). Formulas in propositional (Boolean) logic have the following grammar:

$$\varphi ::= \mathbf{p} \mid \neg\varphi \mid \varphi \vee \varphi$$

, where $\mathbf{p} \in \text{AP}$ Moreover, we add syntactic sugar and allow the formulas *true*, *false*, $\varphi \wedge \psi$, $\varphi \implies$

$\psi, \varphi \iff \psi$. A propositional valuation is a mapping $v : AP \rightarrow \mathbb{B}$, which maps propositional variables to Boolean values.

Semantics: The semantics of propositional logic is given by a satisfaction relation \models that is inductively defined as follows:

$$v \models \mathbf{p} \text{ iff } v(\mathbf{p}) = 1; v \models \neg\varphi \text{ iff } v \not\models \varphi; v \models \varphi \vee \psi \text{ iff } v \models \varphi \text{ or } v \models \psi$$

In the case that $v \models \varphi$, we say that v satisfies φ and call it a model of φ . A propositional formula φ is satisfiable if there exists a model v of φ . The size of a formula is the number of its subformulas (defined in the usual way).

Chapter 3

Learning Regular Expressions

3.1 Regular Expressions

Regular expression consist of a finite alphabet Σ along with three operators $+$ (union), \circ (concatenation) $*$ (Kleene star) with the following grammar

$$r ::= \varepsilon \mid a \in \Sigma \mid r + r \mid r \circ r \mid r^*$$

The set of all regular expressions over alphabet set Σ is referred to as \mathcal{R}_Σ . With slight abuse of notation, another commonly used expression is r^+ which refers to the expression $r \circ r^*$.

A regular expression can also be represented as a *syntax tree* where the nodes are labelled by elements from Λ_R , where $\Lambda_R = \Sigma \cup \{\varepsilon, +, \circ, *\}$ -a combination of alphabets and operators. Leaf nodes are labelled by only alphabets or ε whereas the internal nodes are labelled by operator with the operands being their children.

A *syntax DAG* can also be created by merging the common subtrees in different branches of a syntax tree. A syntax DAG essentially eliminates redundant subexpressions from a syntax tree.

The *size* of a regular expression is the number of unique subexpression present. It coincides with the number of nodes in the syntax DAG representation of the expression.

The nodes of syntax DAG with n nodes can be indexed by natural numbers $1, \dots, n$. The indexing that we follow in this thesis, satisfies the property that children of a node are always indexed by numbers smaller than the node itself and only a leaf node can have the index 1. Note that the indexing need not be unique. Fig. 3.1 shows an example of syntax tree and a syntax DAG for the regular expression $(a \circ b) + b^*$.

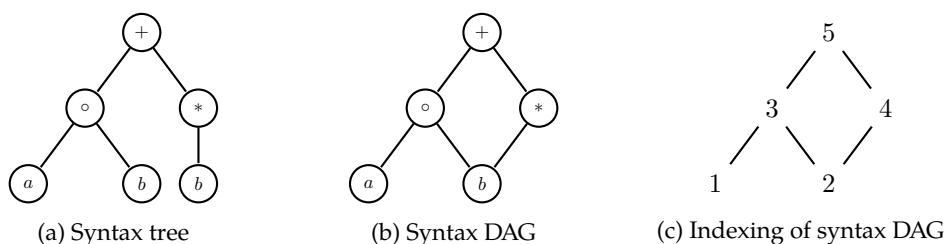


Figure 3.1: Syntax tree, syntax DAG with indexing of node, of the regular expression $(a \circ b) + b^*$

Semantics The semantics of regular expression is generally defined in terms of the language they define:

$$\llbracket \epsilon \rrbracket = \{\epsilon\}; \llbracket a \rrbracket = \{a\}; \llbracket r_1 + r_2 \rrbracket = \llbracket r_1 \rrbracket \cup \llbracket r_2 \rrbracket; \llbracket r_1 \circ r_2 \rrbracket = \llbracket r_1 \rrbracket \circ \llbracket r_2 \rrbracket; \llbracket r^* \rrbracket = \llbracket r \rrbracket^*$$

Definition 1. Matching relation $\vdash \subseteq \Sigma^* \times \mathcal{R}_\Sigma$ is defined inductively on the structure of a regular expression r in the following manner:

$$\begin{aligned} w[i, j] \vdash \epsilon &\iff i = j \\ w[i, j] \vdash a \text{ where } a \in \Sigma &\iff w[i, j] = a \\ w[i, j] \vdash r_1 + r_2 &\iff w[i, j] \vdash r_1 \text{ or } w[i, j] \vdash r_2 \\ w[i, j] \vdash r_1 \circ r_2 &\iff \exists k \in \mathbb{N} \text{ such that } i \leq k \leq j \text{ and } w[i, k] \vdash r_1 \text{ and } w[k+1, j] \vdash r_2 \\ w[i, j] \vdash r^* &\iff \begin{cases} i = j \text{ or} \\ \exists k \in \mathbb{N} \text{ such that } i < k \leq j \text{ and } w[i, k] \vdash r \text{ and } w[k, j] \vdash r^* \end{cases} \end{aligned}$$

The matching relation \vdash holds true when a word belongs to the language defined by the regular expression. More precisely, we have that $w[i, j] \vdash r \iff w[i, j] \in \llbracket r \rrbracket$. This statement can be proved using a simple induction on the structure of r .

Motivated by definition of the matching relation, it is possible to design a *dynamic programming* algorithm for checking whether a given word w matches a regular expression r . Let T_w be a three dimensional table with $|w| + 1$ number of rows and columns both indexed using $0 \dots n$ and the third index is used to denote the subexpression of the expression r for which we are computing the matching relation. The entries of the table have the following meaning.

$$T_w(i, j, r) = 1 \iff w[i, j] \vdash r$$

Given, a regular expression r , the table T_w is constructed in a iterative manner starting from the portion of the table which refers to the leaf elements of the syntax DAG. These table entries are used for propagating the results to table entries for the internal nodes.

- For the regular expression ϵ : $T_w(i, j, \epsilon) = 1$ iff $i = j$
- For a leaf node say a , $T_w(i, j, a) = 1$ iff $w[i, j] = a$
- For + (union) node, with children r_1 and r_2 :
 $T_w(i, j, r_1 + r_2) = 1$ if $T_w(i, j, r_1) = 1$ or $T_w(i, j, r_2) = 1$
- For \circ (concat) node, with children r_1 and r_2 :
 $T_w(i, j, r_1 \circ r_2) = 1$ iff $\exists k, i \leq k \leq j$ such that $T_w(i, k, r_1) = 1$ and $T_w(k, j, r_2) = 1$
- For * (Kleene star) node, with a child r :
 $T_w(i, j, r^*) = 1$ iff $i = j$ or $\exists k, i \leq k \leq j$ such that $T_w(i, k, r) = 1$ and $T_w(k, j, r^*) = 1$
 In this case, T_w is filled starting from (n, n) -th entry up along the columns.

3.2 The Learning Problem

We assume that the data to learn from is given as a pair $\mathcal{S} = (\mathcal{P}, \mathcal{N})$ consisting of two finite, disjoint sets $(\mathcal{P}, \mathcal{N}) \subset \Sigma^*$ of finite words such that $(\mathcal{P} \cap \mathcal{N}) \neq \emptyset$. We call this pair a *sample*. Moreover, we say that a regular expression r over Σ is *consistent* with a sample $\mathcal{S} = (\mathcal{P}, \mathcal{N})$ if

1. $r \vdash u$ for each $u \in \mathcal{P}$; and
2. $r \not\vdash u$ for each $u \in \mathcal{N}$.

SAT solver based solution: A possible solution could be to use SAT-solvers to find the regular expression consistent with the given sample. The formula $\Phi_n^{\mathcal{S}}$ that we feed into the solver has the following structure:

$$\Phi_n^{\mathcal{S}} = \Phi_n^{\text{structure}} \wedge \Phi_n^{\text{consistency}} \quad (3.1)$$

Here, satisfying assignment of $\Phi_n^{\text{structure}}$ provides encoding of a syntax DAG which describes expression of size n , while the $\Phi_n^{\text{consistency}}$ checks whether the guessed expression is consistent with the given sample. $\Phi_n^{\text{structure}}$ uses the following variables:

- $x_{p,\lambda}$ where $p \in \{1, \dots, n\}$ and $\lambda \in \Lambda_R$
- $l_{p,q}$ where $p \in \{2, \dots, n\}$ and $q \in \{1, \dots, p-1\}$
- $r_{p,q}$ where $p \in \{2, \dots, n\}$ and $q \in \{1, \dots, p-1\}$

Now, this is how the formula looks like. It is formed by taking conjunction of the following subformulas.

$$\left[\bigwedge_{1 \leq p \leq n} \bigvee_{\lambda \in \Lambda_R} x_{p,\lambda} \right] \wedge \left[\bigwedge_{1 \leq p \leq n} \bigwedge_{\lambda \neq \lambda' \in \Lambda_R} \neg x_{p,\lambda} \vee \neg x_{p,\lambda'} \right] \quad (3.2)$$

$$\left[\bigwedge_{2 \leq p \leq n} \bigvee_{1 \leq q \leq p} l_{p,q} \right] \wedge \left[\bigwedge_{2 \leq p \leq n} \bigwedge_{1 \leq q \leq q' \leq n} \neg l_{p,q} \vee \neg l_{p,q'} \right] \quad (3.3)$$

$$\left[\bigwedge_{2 \leq p \leq n} \bigvee_{1 \leq q \leq p} r_{p,q} \right] \wedge \left[\bigwedge_{2 \leq p \leq n} \bigwedge_{1 \leq q \leq q' \leq n} \neg r_{p,q} \vee \neg r_{p,q'} \right] \quad (3.4)$$

$$x_{1,\epsilon} \vee \bigvee_{a \in \Sigma} x_{1,a} \quad (3.5)$$

Subformula (3.2) says that each node would be uniquely labelled by an element from Λ_R . Subformulas (3.3) and (3.4) say that each node would have a unique left and right child respectively. Last but not the least, Subformula (3.5) says that the first node can have only ϵ or an alphabet.

Now, it needs to be ensured that the regular expression guessed by the first part $\Phi_n^{\text{structure}}$ is consistent with the given sample. For each word w in the given sample, consider a three dimensional table whose entries are basically variables $y_{i,j,p}^w$, where the indices i and j refers to the substring of the word w and p refers to the sub-expression we are looking at. For simplicity, we assume $m = |w|$. Now, Φ_n^w is a conjunction of the following subformulas:

$$\bigwedge_{1 \leq p \leq n} x_{p,\varepsilon} \rightarrow \left[\bigwedge_{0 \leq i \leq j \leq m} y_{i,j,p}^w \leftrightarrow [i = j] \right] \quad (3.6)$$

$$\bigwedge_{1 \leq p \leq n} \bigwedge_{a \in \Sigma} x_{p,a} \rightarrow \left[\bigwedge_{0 \leq i \leq j \leq m} \begin{cases} y_{i,j,p}^w & \text{if } w[i, j] = a \\ \neg y_{i,j,p}^w & \text{if } w[i, j] \neq a \end{cases} \right] \quad (3.7)$$

$$\bigwedge_{\substack{1 \leq p \leq n \\ 1 \leq q, q' \leq p}} x_{p,+} \wedge l_{p,q} \wedge r_{p,q'} \rightarrow \left[\bigwedge_{0 \leq i \leq j \leq m} \left[y_{i,j,p}^w \leftrightarrow y_{i,j,q}^w \vee y_{i,j,q'}^w \right] \right] \quad (3.8)$$

$$\bigwedge_{\substack{1 \leq p \leq n \\ 1 \leq q, q' \leq p}} x_{p,\circ} \wedge l_{p,q} \wedge r_{p,q'} \rightarrow \left[\bigwedge_{0 \leq i \leq j \leq m} \left[y_{i,j,p}^w \leftrightarrow \bigvee_{i \leq k \leq j} y_{i,k,q}^w \wedge y_{k,j,q'}^w \right] \right] \quad (3.9)$$

$$\bigwedge_{\substack{1 \leq p \leq n \\ 1 \leq q, q' \leq p}} x_{p,*} \wedge l_{p,q} \rightarrow \left[\bigwedge_{0 \leq i \leq j \leq m} \left[y_{i,j,p}^w \leftrightarrow [i = j] \vee \bigvee_{i < k \leq j} y_{i,k,q}^w \wedge y_{k,j,p}^w \right] \right] \quad (3.10)$$

The idea for each of these subformulas is motivated by the procedure for filling out the table T_w , which in turn essentially uses the way the matching relation has been defined. Using the formula Φ_n^w for various w , we construct $\Phi_n^{\text{consistency}}$ in the following manner. Notice, that variable $y_{0,m,n}^w$ encodes whether word w matches the guessed regular expression. Hence, $y_{0,m,n}^w$ should be *true* for positive words and should be *false* for negative words.

$$\Phi_n^{\text{consistency}} = \left[\bigwedge_{w \in \mathcal{P}} \Phi_n^w \wedge y_{0,m,n}^w \right] \wedge \left[\bigwedge_{w \in \mathcal{N}} \Phi_n^w \wedge \neg y_{0,m,n}^w \right] \quad (3.11)$$

Given that we have devised SAT formula for this problem, we utilise it to come up with a simple algorithm using SAT-solvers to find out the smallest regular expression consistent with the sample.

Algorithm 1: SAT-based learning algorithm for regular expressions

Input: A sample \mathcal{S}

- 1 $n \leftarrow 0$
 - 2 **repeat**
 - 3 $n \leftarrow n + 1$
 - 4 Construct and solve the formula $\Phi_n^{\mathcal{S}}$ using a SAT solver
 - 5 **until** $\Phi_n^{\mathcal{S}}$ is satisfiable, say with model v
 - 6 $R^v \leftarrow \text{Extract}(\Phi_n^{\mathcal{S}}, v)$
 - 7 **return** the regular expression R^v
-

In the above algorithm, if we can find a model v of $\Phi_n^{\mathcal{S}}$, we can extract a regular expression R^v using the procedure $\text{Extract}(\Phi_n^{\mathcal{S}}, v)$. The first part $\Phi_n^{\text{structure}}$ of the formula $\Phi_n^{\mathcal{S}}$ encodes the structure of the regular expression as a syntax DAG. The model v assigns valuations to the variables $x_{p,\lambda}$, $l_{p,q}$ and $r_{p,q}$ which determines the shape of the syntax DAG. Hence, the procedure $\text{Extract}(\Phi_n^{\mathcal{S}}, v)$ uses this information to derive the expression R^v from the syntax tree.

3.2.1 Proof of correctness

Claim 1. Let $\mathcal{S} = (\mathcal{P}, \mathcal{N})$ be a sample, $n \in \mathbb{N} \setminus \{0\}$, and $\Phi_n^{\mathcal{S}}$ be the propositional formula defined above. Then, the following holds:

1. If there exists a regular expression of size n , R^S that is consistent with \mathcal{S} , then the propositional formula Φ_n^S is satisfiable.
2. If $v \models \Phi_n^S$, then R^v is a regular expression of size n that is consistent with \mathcal{S} .

Proof. Using the syntax DAG of the expression R^S which is indexed using $1 \cdots n$, we formulate a valuation v for the propositional variables in Φ_n^S . We use R_p^S to refer to the regular expression rooted at the p^{th} node.

- We set $v(x_{p,\lambda}) = 1$ iff the p^{th} node is labelled by λ .
- We set $v(l_{p,q}) = 1$ iff q^{th} node is the left child of the p^{th} node and similarly, $v(r_{p,q}) = 1$ iff q^{th} node is the right child of the p^{th} node.
- We set $v(y_{i,j,p}^w) = 1$ iff $w[i, j] \vdash R_p^S$

Firstly, it is easy to see that $v \models \Phi_n^{\text{structure}}$, since the formulated valuation ensures the uniqueness of the labels of the nodes as well as that of the left and right children. Moreover, $v \models \Phi_n^w$ for $w \in \mathcal{P} \cup \mathcal{N}$, since, the values of the variables $y_{i,j,p}^w$ correspond exactly to the valuation of each subexpression R_i^S on the subword $w[i, j]$. Finally, the fact that R^S is consistent with \mathcal{S} implies $v \models y_{i,j,p}^w$ for each $w \in \mathcal{P}$ and $v \models \neg y_{i,j,p}^w$ for each $w \in \mathcal{N}$. Thus, $v \models \Phi_n^S$, which proves that Φ_n^S is satisfiable.

For the second statement, since, $v \models \Phi_n^S$, we have $v \models \Phi_n^{\text{structure}}$ as well. Hence, the valuation of the variables $x_{p,\lambda}$, $l_{p,q}$, and $r_{p,q}$ encode a syntax DAG from which we obtain the regular expression R^v . Again here, we consider R_p^v to be the subexpression of R^v rooted at the node p . Now, it needs to be shown that R^v is indeed consistent with the sample \mathcal{S} . For that we show, $v(y_{i,j,p}^w) = 1 \iff w[i, j] \vdash R_p^v$ for any $i, j \in \{0, \dots, m\}$ where $m = |w|$ and $p \in \{0, \dots, n\}$. This can be done using induction on the structure of R^v .

- *Base case $R_p^v = \varepsilon$:* In this case, $v(x_{p,\varepsilon}) = 1$, which leads to the following implications:

$$\begin{aligned} v(y_{i,j,p}^w) = 1 &\iff \mathbf{p} \in w^\omega[i, j] && \text{[Using Subformula (3.6)]} \\ &\iff w[i, j] \vdash \varepsilon && \text{[Using Def. 1]} \end{aligned}$$

- *Base case $R_p^v = \mathbf{p}$:* In this case, $v(x_{p,a}) = 1$ and this leads to the following implications:

$$\begin{aligned} v(y_{i,j,p}^w) = 1 &\iff w[i, j] = a && \text{[Using Subformula (3.7)]} \\ &\iff w[i, j] \vdash a && \text{[Using Def. 1]} \end{aligned}$$

- *Case $R_p^v = R_q^v + R_{q'}^v$:* In this case, $v(x_{p,+})$, $v(l_{p,q})$, $v(r_{p,q'})$ are all set to 1, and this leads to the following implications:

$$\begin{aligned} v(y_{i,j,p}^w) = 1 &\iff v(y_{i,j,q}^w) = 1 \text{ or } v(y_{i,j,q'}^w) = 1 && \text{[Using Subformula (3.8)]} \\ &\iff w[i, j] \vdash R_q^v \text{ or } w[i, j] \vdash R_{q'}^v && \text{[Using ind. hypothesis]} \\ &\iff w[i, j] \vdash R_q^v + R_{q'}^v && \text{[Using Def. 1]} \end{aligned}$$

- *Case $R_p^v = R_q^v \circ R_{q'}^v$:* In this case, $v(x_{p,\circ})$, $v(l_{p,q})$, $v(r_{p,q'})$ are all set to 1 and this leads to the

following implications:

$$\begin{aligned}
v(y_{i,j,p}^w) = 1 &\iff \exists k \in \mathbb{N}, 1 \leq k \leq m, v(y_{i,k,q}^w) = 1 \text{ and } v(y_{k,j,q'}^w) = 1 && \text{[Using Subformula (3.9)]} \\
&\iff \exists k \in \mathbb{N}, 1 \leq k \leq m, w[i,k] \vdash R_q^v \text{ and } w[k,j] \vdash R_{q'}^v && \text{[Using ind. hypothesis]} \\
&\iff w[i,j] \vdash R_q^v \circ R_{q'}^v && \text{[Using Def. 1]}
\end{aligned}$$

– *Case $R_p^v = (R_q^v)^*$* : In this case, $v(x_{p,*})$ and $v(l_{p,q})$ are set to 1 and this leads to the following implications:

$$\begin{aligned}
v(y_{i,j,p}^w) = 1 &\iff \begin{cases} i = j \text{ or} \\ \exists k \in \mathbb{N}, 1 \leq k \leq m, v(y_{i,k,q}^w) = 1 \text{ and } v(y_{k,j,p}^w) = 1 \end{cases} && \text{[Using Subformula (3.10)]} \\
&\iff \begin{cases} i = j \text{ or} \\ \exists k \in \mathbb{N}, 1 \leq k \leq m, w[i,k] \vdash R_q^v \text{ and } w[k,j] \vdash (R_q^v)^* \end{cases} && \text{[Using ind. hypothesis]} \\
&\iff w[i,j] \vdash (R_q^v)^* && \text{[Using Def. 1]}
\end{aligned}$$

Here, we assumed that $v(y_{k,j,p}^w) = 1$ iff $w[k,j] \vdash R_p^v$ in the second step, which cannot be deduced from the induction on the structure of R_p^v . This we obtained using another level of induction, which is on the length of the subword that matches R_p^v . More precisely, we induct on k which ranges from j to $i + 1$, where by induction hypothesis, it is assumed that

$$v(y_{k,j,p}^w) = 1 \iff w[k,j] \vdash R_p^v \quad \forall k, i < k \leq j$$

The base case occurs when $k = j$, and then, we have $v(y_{k,j,p}^w) = 1 \iff k = j \iff w[k,j] \vdash R_p^v$.

□

Chapter 4

Learning ω –Regular Expressions

4.1 ω –Regular expression

An ω -regular expression r_ω is an expression with the following grammar:

$$r_\omega ::= r^\omega \mid r \circ_\omega r_\omega \mid r_\omega +_\omega r_\omega$$

where, r is a regular expression as described in Section 3.1. Let the set of ω -regular expressions be denoted by $\mathcal{R}_\Sigma^\omega$. Again here, an ω -regular expression can be represented in the form of a syntax tree or a syntax DAG, similar to ones for regular expressions. The only difference here is that the nodes of the syntax tree and the syntax DAG are labelled by elements from $\Lambda_\omega = \Lambda_R \cup \{\omega, +_\omega, \circ_\omega\}$.

Semantics The semantics of an ω -regular expression is as usual defined in terms of the language they define; obtained in the following manner:

$$\llbracket r^\omega \rrbracket = \llbracket r \rrbracket^\omega; \llbracket r \circ_\omega r_\omega \rrbracket = \llbracket r \rrbracket \circ_\omega \llbracket r_\omega \rrbracket; \llbracket r_\omega +_\omega r'_\omega \rrbracket = \llbracket r_\omega \rrbracket \cup \llbracket r'_\omega \rrbracket$$

It should be noted that r^ω is well defined only when $\epsilon \notin \llbracket r \rrbracket$, since otherwise, $\llbracket r^\omega \rrbracket$ might contain invalid infinite words.

Now, let us look at some results regarding membership of ultimately periodic words in languages defined by ω -regular expressions. In the following results, we assume $u \in \Sigma^*$ and $v \in \Sigma^+$ and r is a regular expression in \mathcal{R}_Σ , such that $\epsilon \notin \llbracket r \rrbracket$. Moreover, assume that the minimal DFA of $\llbracket r \rrbracket$ has size n .

Proposition 1. $uv^\omega \in \llbracket r^\omega \rrbracket \iff \exists i, j \in \mathbb{N}, |u| < i < j, uv^\omega[0, i) \in \llbracket r^+ \rrbracket$ and $uv^\omega[i, j) \in \llbracket r^+ \rrbracket$ where $j \equiv i \pmod{|v|}$.

Proof. (\Rightarrow) From the semantics of ω -regular expressions, we have $uv^\omega \in \llbracket r^\omega \rrbracket \implies uv^\omega \in \llbracket r \rrbracket^\omega$. Hence, there exists an infinite sequence of natural numbers i_1, i_2, i_3, \dots satisfying $|u| < i_1 < i_2 < \dots$ such that $uv^\omega[0, i_1) \in \llbracket r^+ \rrbracket$ and from then on $uv^\omega[i_1, i_2) \in \llbracket r \rrbracket, uv^\omega[i_2, i_3) \in \llbracket r \rrbracket$ and so on. Since, i_0, i_1, i_2, \dots is an infinite sequence, it is possible to find numbers i and j in the sequence, such that $j \equiv i \pmod{|v|}$ using pigeonhole principle. Now, clearly, $uv^\omega[0, i) \in \llbracket r^+ \rrbracket$ and $uv^\omega[i, j) \in \llbracket r^+ \rrbracket$, since i and j belong to the sequence.

(\Leftarrow) Let $uv^\omega[0, i) = v_1$ and $uv^\omega[i, j) = v_2$, where both $uv^\omega[0, i)$ and $uv^\omega[i, j)$ satisfy the properties mentioned in the lemma. It is an easy observation that $v_1 v_2^\omega = uv^\omega$. As we already know $v_1 \in \llbracket r^+ \rrbracket$ and $v_2 \in \llbracket r^+ \rrbracket$, we can conclude that $v_1 v_2^\omega \in \llbracket r^+ \rrbracket \circ \llbracket r^+ \rrbracket^\omega \implies v_1 v_2^\omega \in \llbracket r^\omega \rrbracket \implies uv^\omega \in \llbracket r^\omega \rrbracket$. \square

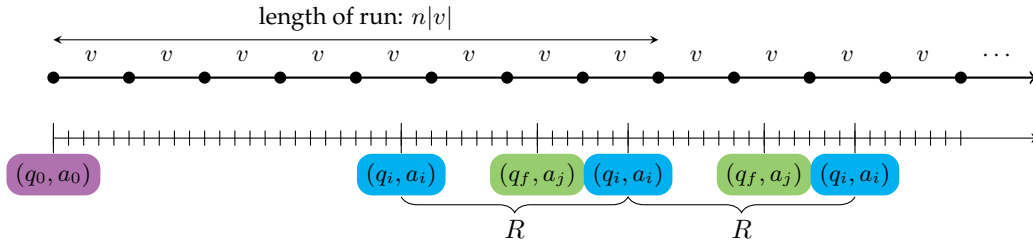


Figure 4.1: The run of the DFA of $\llbracket r \rrbracket$ on $v^\omega[i, j]$, where R is the repeating run starting at (i, j) . The first occurrence of R happens within $n|v| + 1$ steps and hence, the first occurrence of state q_f also happens within that many steps.

Lemma 1. $v^\omega \in \llbracket r^\omega \rrbracket \iff \exists i \in \mathbb{N}, v^\omega[0, i] \in \llbracket r^{|\mathbb{N}|+1} \rrbracket$.

Proof. (\Rightarrow) Since, $v^\omega \in \llbracket r \rrbracket^\omega$, for any natural number n , there is a prefix $v^\omega[0, i]$ of v^ω such that $v^\omega[0, i] \in \llbracket r^n \rrbracket$. Therefore, it holds true for $n = |\mathbb{N}| + 1$ as well.

(\Leftarrow) Here, we have $v^\omega[0, i] \in \llbracket r^{|\mathbb{N}|+1} \rrbracket$ for some $i \in \mathbb{N}$. Therefore, we can find a finite sequence of distinct natural numbers $i_1, i_2, \dots, i_{|\mathbb{N}|+1}$, where $i_{|\mathbb{N}|+1} = i$ and also, $v^\omega[0, i_1] \in \llbracket r \rrbracket$, $v^\omega[i_1, i_2] \in \llbracket r \rrbracket$ and so on. Now, notice that it is possible to find numbers i', j' with $i' < j'$, in the sequence such that $j' \equiv i' \pmod{|\mathbb{N}|}$ using pigeonhole principle (Since, the sequence contains more than $|\mathbb{N}|$ elements). Consequently, we have $v^\omega[0, i'] \in \llbracket r^+ \rrbracket$, $v^\omega[i', j'] \in \llbracket r^+ \rrbracket$ and $j' \equiv i' \pmod{|\mathbb{N}|}$, which gives us $v^\omega \in \llbracket r^\omega \rrbracket$ straightaway using Prop. 1. \square

Lemma 2. Let $uv^\omega \in \llbracket r^\omega \rrbracket$ such that for some $i > |u|$, $uv^\omega[0, i] \in \llbracket r \rrbracket$ and $uv^\omega[i, \infty) \in \llbracket r^\omega \rrbracket$. Then, $\exists j \in \mathbb{N}, j > i$, such that $uv^\omega[0, j] \in \llbracket r^{|\mathbb{N}|+1} \rrbracket$

The proof for this proceeds exactly the same way as Lemma 1. The next lemma shows how the length of each match can be bounded.

Lemma 3. $v^\omega[i, j] \in \llbracket r \rrbracket \implies \exists k \in \mathbb{N}, k - i \leq n|v|$ such that $v^\omega[i, k] \in \llbracket r \rrbracket$ and $k \equiv j \pmod{|\mathbb{N}|}$, where n is the size of the minimal DFA for $\llbracket r \rrbracket$.

Proof. If $j - i \leq n|v|$, we are done, since we simply take $k = j$. On the other hand, if $j - i > n|v|$, finding the suitable k is a bit more involved. Firstly, since, $v^\omega[i, j] \in \llbracket r \rrbracket$, there is an accepting run of the DFA \mathcal{A} for $\llbracket r \rrbracket$ on $v^\omega[i, j]$, which is of length greater than $n|v|$. Now, we look at this run of \mathcal{A} on $v^\omega[i, j]$. Fig 4.1 provides a pictorial depiction of the run. Here we consider the run to be a sequence of tuples of the form (q, m) , where the first entry refers to the state of the automaton at that instant and the second entry refers to the position in v which will be read next by the automaton. Hence, it is easy to observe that any run longer than $n|v|$ would mean that there exists a tuple which repeats itself during the run, due to pigeonhole principle. Let (q, l) be the tuple which repeats itself and let the run starting at the first occurrence of (q, l) to the next, be referred to as R . Notice that due to the deterministic nature of the automaton, R repeats itself during the rest of the run. Hence, if a final state q_f occurs after $n|v|$ steps, there must be a tuple (q_f, m) which belongs to the run R . Clearly, (q_f, m) must have been also visited during the first occurrence of R , which happened within $n|v|$ steps. Thus, we get a prefix $v^\omega[i, k] \in \llbracket r \rrbracket$, which has length less than $n|v|$. Moreover, $v^\omega[i, j]$ and $v^\omega[i, k]$ terminate at the same position in v , meaning $k \equiv j \pmod{|\mathbb{N}|}$. \square

Lemma 4. $uv^\omega[i, j] \in \llbracket r \rrbracket$, where $i < |u| < j \implies \exists k \in \mathbb{N}, |u| < k \leq |u| + n|v|$ such that $uv^\omega[i, k] \in \llbracket r \rrbracket$ and $k \equiv j \pmod{|\mathbb{N}|}$.

The proof for this statement goes along the same lines as Lemma 3. Finally, having these results, it is possible to verify whether an infinite word (ultimately periodic word in this case) actually belongs to the language of regular expression by checking only a finite portion of the word. The following theorems come in handy when designing algorithms for checking matching of infinite word with ω -regular expressions, by bounding the length on which checking is done.

Theorem 1. $v^\omega \in \llbracket r^\omega \rrbracket \iff \exists i, j \in \mathbb{N}, 0 < i < j \leq n|v| + n|v|^2$ $v^\omega[0, i) \in \llbracket r^+ \rrbracket$ and $v^\omega[i, j) \in \llbracket r^+ \rrbracket$ where $j \equiv i \pmod{|v|}$

Proof. (\Leftarrow) This can be seen directly from Prop. 1, assuming $u = \epsilon$.

(\Rightarrow) Using Lemma 1, there exists some i' such that $v^\omega[0, i') \in \llbracket r^{+|v|+1} \rrbracket$. As a result, we can generate a sequence of distinct natural numbers, $i_1, \dots, i_{|v|+1}$, with $i_{|v|+1} = i'$ such that $v^\omega[0, i_1) \in \llbracket r \rrbracket$, $v^\omega[i_1, i_2) \in \llbracket r \rrbracket$ and so on. Now, we can find another sequence $j_1, \dots, j_{|v|+1}$ with $j_{|v|+1} = i'$, which satisfy all the properties that the sequence of i 's satisfy and additionally have $j_1 \leq n|v|$, $j_2 - j_1 \leq n|v|$, so on and also, $j_1 \equiv i_1 \pmod{|v|}$, $j_2 \equiv i_2 \pmod{|v|}$, so on. This is a direct consequence of Lemma 3. It is easy to see that $j_{|v|+1} \leq n|v| + n|v|^2$. Moreover, we can again find some i and j in this sequence such that $j \equiv i \pmod{|v|}$ and clearly, $v^\omega[0, i) \in \llbracket r^+ \rrbracket$ and $v^\omega[i, j) \in \llbracket r^+ \rrbracket$. \square

Theorem 2. $uv^\omega \in \llbracket r^\omega \rrbracket \iff \exists i, j \in \mathbb{N}, |u| < i < j \leq |u| + n|v| + n|v|^2$, $uv^\omega[0, i) \in \llbracket r^+ \rrbracket$ and $uv^\omega[i, j) \in \llbracket r^+ \rrbracket$ where $j \equiv i \pmod{|v|}$

Proof. (\Leftarrow) This can be seen directly from Prop. 1.

(\Rightarrow) When $uv^\omega \in \llbracket r^\omega \rrbracket$, there could be two possible cases arising from how uv^ω matches r^ω .

- First case would be one in which $u \in \llbracket r^+ \rrbracket$ and $v^\omega \in \llbracket r^\omega \rrbracket$. In this case, Theorem 1 can be applied directly to obtain the result.
- In the second case, $\exists i', j' \in \mathbb{N}$ with $i' < |u| < j'$, such that, $uv^\omega[0, i') \in \llbracket r^+ \rrbracket$, $uv^\omega[i', j') \in \llbracket r \rrbracket$, and $uv^\omega[j', \infty) \in \llbracket r^\omega \rrbracket$. Now, it can be seen that infix $uv^\omega[i', \infty)$ satisfies the conditions of Lemma 2. Hence, there is some $k' > j'$, such that $uv^\omega[i', k') \in \llbracket r^{+|v|+1} \rrbracket$. As a result, we can generate a sequence of distinct natural numbers $i_1, \dots, i_{|v|+1}$, with $i_1 = j'$ and $i_{|v|+1} = k'$ such that $uv^\omega[i', i_1) \in \llbracket r \rrbracket$, $uv^\omega[i_1, i_2) \in \llbracket r \rrbracket$ and so on. Combining results of Lemma 4 and Lemma 3, we can find another sequence $j_1, j_2 \dots j_{|v|+1}$ which satisfy the properties of the sequence of i 's and additionally have $j_1 \leq |u| + n|v|$, $j_2 - j_1 \leq n|v|$, $j_3 - j_2 \leq n|v|$, so on and also, $j_1 \equiv i_1 \pmod{|v|}$, $j_2 \equiv i_2 \pmod{|v|}$, so on. It is easy to see that $j_{|v|+1} \leq |u| + n|v| + n|v|^2$. Moreover, in this sequence we can find some i and j such that $j \equiv i \pmod{|v|}$ and clearly, $uv^\omega[0, i) \in \llbracket r^+ \rrbracket$ and $uv^\omega[i, j) \in \llbracket r^+ \rrbracket$. \square

Motivated by the above results, we define the matching relation for infinite words which would help us to algorithmically check the matching of a ultimately periodic word with an ω -regular expression.

Definition 2. The matching relation \vdash for finite words with regular expressions, could be extended to account for matching of ultimately periodic words with ω -regular expressions as well. Here, n refers to the size of the minimal DFA of $\llbracket r \rrbracket$ in all the cases.

$$\begin{aligned}
uv^\omega[i, \infty) \vdash r^\omega &\iff \exists j, k \in \mathbb{N}, \max(i, |u|) < j < k \leq i + |u| + |v| + n|v| + n|v|^2, \\
&\quad uv^\omega[i, j) \vdash r^+ \text{ and } uv^\omega[j, k) \vdash r^+ \text{ where } k \equiv j \pmod{|v|} \\
uv^\omega[i, \infty) \vdash r \circ r_\omega &\iff \exists j \in \mathbb{N}, i \leq j \leq |u| + |v| + n|v|, uv^\omega[i, j) \vdash r \text{ and } uv^\omega[j, \infty) \vdash r_\omega \\
uv^\omega[i, \infty) \vdash r_\omega + r'_\omega &\iff uv^\omega[i, \infty) \vdash r_\omega \text{ or } uv^\omega[i, \infty) \vdash r'_\omega
\end{aligned}$$

It is necessary to check whether the matching relation satisfies $uv^\omega[i, \infty) \vdash r_\omega \iff v^\omega[i, \infty) \in \llbracket r_\omega \rrbracket$. This check proceeds via structural induction on the expression r_ω .

Observation 4.1: Let $uv^\omega \in \Sigma^\omega$. Then, $uv^\omega[|u| + i, \infty) = uv^\omega[|u| + j, \infty)$ for $j \equiv i \pmod{|v|}$. Because of this observation, we can say $r_\omega \vdash uv^\omega[|u| + i, \infty) \iff r_\omega \vdash uv^\omega[|u| + j, \infty)$. Hence, the matching relation has been defined only for suffix $uv^\omega[i, \infty)$, where $i < |uv|$.

Unlike regular expression, the construction of the tables for algorithmically checking matching of an ω -regular expression with an ω -word needs a slightly different approach. There are two types of tables that need to be maintained for each node. The first table keeps track of which infinite suffixes of the ω -word matches which ω -regular subexpressions. The other table maintains record of matching of finite infixes of the ω -word with regular subexpressions. The entries of both types of tables are propagated to tables of higher indexed nodes using the matching relation.

4.2 The Learning Problem

We assume that the data to learn from is given as a pair $\mathcal{S} = (\mathcal{P}, \mathcal{N})$ consisting of two finite, disjoint sets $\mathcal{P}, \mathcal{N} \subset \Sigma^*$ of ultimately periodic words such that $\mathcal{P} \cap \mathcal{N} \neq \emptyset$.

Just like regular expressions, we say that an ω -regular expression r_ω over Σ is *consistent* with a sample $\mathcal{S} = (\mathcal{P}, \mathcal{N})$ if

1. $r_\omega \vdash u$ for each $u \in \mathcal{P}$; and
2. $r_\omega \not\vdash u$ for each $u \in \mathcal{N}$.

A solution similar to the one used for regular expression can be used for this problem as well. In fact, the Algorithm 1 would work here and would return a ω -regular expression. The formula $\Phi_n^{\mathcal{S}}$ used in the algorithm is designed slightly differently for this case. The formula has two parts as described in Equation 3.1. The first part $\Phi_n^{\text{structure}}$ which encodes a syntax DAG of an ω -regular expression has the following variables.

- $x_{p,\lambda}$ where $p \in \{1, \dots, n\}$ and $\lambda \in \Lambda_\omega$
- $l_{p,q}, r_{p,q}$ where $p \in \{2, \dots, n\}$ and $q \in \{1, \dots, p-1\}$

Now, the variables should satisfy certain constraints in order to form a valid syntax DAG. But, these constraints are precisely the ones that needed to be satisfied by the syntax DAG of regular expressions. Thus, the exact description of the formulas have been skipped.

While the first part of the formula guesses an ω -regular expression r , the second part $\Phi_n^{\text{consistency}}$ checks whether words in the sample matches r . This formula makes use of the following variables. Here, we assume $b = n + n|v| + 1$ and $c = n + 1$.

- $y_{i,p}^{u,v}$ where $0 \leq i \leq |uv| - 1$
- $z_{i,j,p'}^{uv^b}$ where $0 \leq i \leq j \leq |uv^b|$.

- $\bar{z}_{i,j,p}^{uv^b}$, where $0 \leq i \leq j \leq |uv^b|$

The variable $y_{i,p}^{u,v}$ basically checks the matching of the infinite suffix $uv^\omega[i, \infty)$ with the subexpression of the guessed expression r rooted at an infinite operator indexed by p . The variable $z_{i,j,p}^{uv^b}$ checks the matching of finite infix $uv^b[i, j]$ with the subexpression of r rooted at a finite operator indexed by p . $\bar{z}_{i,j,p}^{uv^b}$ is an auxiliary variable which captures that whether $uv^b[i, j]$ matches \bar{r}^+ if \bar{r} is the subexpression rooted under the p^{th} node. Φ_n^w is a disjunction of several formulas as mentioned below.

$$\bigwedge_{1 \leq p \leq n} x_{p,\varepsilon} \rightarrow \left[\bigwedge_{0 \leq i \leq j \leq |uv^b|} z_{i,j,p}^{uv^b} \leftrightarrow [i = j] \right] \quad (4.1)$$

$$\bigwedge_{1 \leq p \leq n} \bigwedge_{a \in \Sigma} x_{p,a} \rightarrow \left[\bigwedge_{0 \leq i \leq j \leq |uv^b|} \begin{cases} z_{i,j,p}^{uv^b} & \text{if } w[i, j] = a \\ \neg z_{i,j,p}^{uv^b} & \text{if } w[i, j] \neq a \end{cases} \right] \quad (4.2)$$

$$\bigwedge_{\substack{1 \leq p \leq n \\ 1 \leq q, q' \leq p}} x_{p,+} \wedge l_{p,q} \wedge r_{p,q'} \rightarrow \left[\bigwedge_{0 \leq i \leq j \leq |uv^b|} [z_{i,j,p}^{uv^b} \leftrightarrow z_{i,j,q}^{uv^b} \vee z_{i,j,q'}^{uv^b}] \right] \quad (4.3)$$

$$\bigwedge_{\substack{1 \leq p \leq n \\ 1 \leq q, q' \leq p}} x_{p,\circ} \wedge l_{p,q} \wedge r_{p,q'} \rightarrow \left[\bigwedge_{0 \leq i \leq j \leq |uv^b|} [z_{i,j,p}^{uv^b} \leftrightarrow \bigvee_{i \leq k \leq j} z_{i,k,q}^{uv^b} \wedge z_{k,j,q'}^{uv^b}] \right] \quad (4.4)$$

$$\bigwedge_{\substack{1 \leq p \leq n \\ 1 \leq q, q' \leq p}} x_{p,*} \wedge l_{p,q} \rightarrow \left[\bigwedge_{0 \leq i \leq j \leq |uv^b|} [z_{i,j,p}^{uv^b} \leftrightarrow [i = j] \vee \bigvee_{i < k \leq j} z_{i,k,q}^{uv^b} \wedge z_{k,j,p}^{uv^b}] \right] \quad (4.5)$$

The above formulas are for the finite operators in the guessed expression. They are motivated by Formulas 3.6 to 3.10 for regular expressions. The formulas for the infinite operators are as follows.

$$\bigwedge_{\substack{1 \leq p \leq n \\ 1 \leq q \leq p}} x_{p,\omega} \wedge l_{p,q} \rightarrow \left[\bigwedge_{0 \leq i < |uv|} [y_{i,p}^{u,v} \leftrightarrow \bigvee_{\substack{\max(i, |u|) < j < k \leq |uv^b| \\ k \equiv j \pmod{|v|}}} [\bar{z}_{i,j,q}^{uv^b} \wedge \bar{z}_{j,k,q}^{u,v}]] \right] \\ \wedge \left[\bigwedge_{0 \leq i \leq j \leq |uv^b|} [\bar{z}_{i,j,q}^{u,v} \leftrightarrow \bigvee_{i < k \leq j} z_{i,k,q}^{uv^b} \wedge \bar{z}_{k,j,q}^{uv^b}] \right] \quad (4.6)$$

$$\bigwedge_{\substack{1 \leq p \leq n \\ 1 \leq q, q' \leq p}} x_{p,\circ_\omega} \wedge l_{p,q} \wedge r_{p,q'} \rightarrow \left[\bigwedge_{0 \leq i < |uv|} [y_{i,p}^{u,v} \leftrightarrow \bigvee_{i \leq j \leq |uv^c|} [z_{i,j,q}^{uv^b} \wedge y_{M(j),q'}^{u,v}]] \right] \quad (4.7)$$

$$\bigwedge_{\substack{1 \leq p \leq n \\ 1 \leq q, q' \leq p}} x_{p,+_\omega} \wedge l_{p,q} \wedge r_{p,q'} \rightarrow \left[\bigwedge_{0 \leq i < |uv|} y_{i,p}^{u,v} \leftrightarrow y_{i,q}^{u,v} \vee y_{i,q'}^{u,v} \right] \quad (4.8)$$

The formulas for the infinite operators are motivated by how the matching relation is defined. In the Formula 4.6, the second disjunct computes the matching relation of \bar{r}^+ for the subexpression \bar{r} rooted at the q^{th} node. $M(j)$ used in Formula 4.7 for \circ_ω is defined as follows:

$$M(j) = \begin{cases} j & \text{for } j < |uv| \\ |u| + ((j - |u|) \% |v|) & \text{for } j \geq |uv| \end{cases} \quad (4.9)$$

$a \% b = r$ if r is the remainder when a is divided by b . $\%$ is a 'remainder' operation commonly used

in programming languages such as C . Finally, all the Formulas 4.1 to 4.8 are put together in Φ_n^w . $\Phi_n^{\text{consistency}}$ is constructed in the following manner:

$$\Phi_n^{\text{consistency}} = \left[\bigwedge_{w \in \mathcal{P}} \Phi_n^w \wedge y_{0,n}^w \right] \wedge \left[\bigwedge_{w \in \mathcal{N}} \Phi_n^w \wedge \neg y_{0,n}^w \right] \quad (4.10)$$

Claim 2. Let $\mathcal{S} = (\mathcal{P}, \mathcal{N})$ be a sample, $n \in \mathbb{N} \setminus \{0\}$, and Φ_n^S be the propositional formula defined above. Then, the following holds:

1. If there exists a ω -regular expression of size n , R_ω^S that is consistent with \mathcal{S} , then the propositional formula Φ_n^S is satisfiable.
2. If $v \models \Phi_n^S$, then R^v is a ω -regular expression of size n that is consistent with \mathcal{S} .

The proof for this is similar to the proof for Claim 1. Appendix B has the complete proof.

Chapter 5

Learning of Linear Temporal Language

The learning algorithm in this chapter has been taken from the paper by Neider and Gavran[6].

5.1 Linear Temporal Logic

Linear Temporal Logic (LTL)[3] is an extension of propositional Boolean logic with modalities that allow expressing temporal properties. The syntax of LTL is given by the following grammar:

$$\varphi ::= p \mid \varphi \vee \psi \mid \neg\varphi \mid \mathbf{X}\varphi \mid \varphi\mathbf{U}\psi$$

where $p \in AP$, where AP be a set of atomic propositions. The set of all LTL formulas over AP is referred to as \mathcal{F}_{LTL}

LTL is allowed to have formulas commonly used in propositional logic such as *true*, *false*, $\varphi \wedge \psi$, $\psi \rightarrow \varphi$ which are defined as usual. Moreover, additional temporal formulas are also allowed— $\mathbf{F}\varphi := true\mathbf{U}\varphi$ (“finally”) and $\mathbf{G}\varphi := \neg\mathbf{F}\neg\varphi$ (“globally”). For sake of the learning problem, it is enough to consider the primary operators. Hence, let $\Lambda_L = AP \cup \{\vee, \neg, \mathbf{X}, \mathbf{U}\}$

Syntax trees and syntax DAGs are appropriate representation for LTL formulas as well. Here, the labels of the syntax tree or syntax DAG come from the set Λ_L . The size of an LTL formula is naturally defined as the number of unique subformulas in it, which is equivalent to number of nodes in the syntax DAG representation of the formula.

Semantics: Given a word $w \in (2^{AP})^\omega$ and a LTL formula φ , we define the relation $w \models_{LTL} \varphi$ (read “ w satisfies φ ”).

$$\begin{aligned} w \models_{LTL} p &\text{ iff } p \in w[0, 0) \\ w \models_{LTL} \neg\varphi &\text{ iff } w \not\models_{LTL} \varphi \\ w \models_{LTL} \varphi \vee \psi &\text{ iff } w \models_{LTL} \varphi \text{ or } w \models_{LTL} \psi \\ w \models_{LTL} \mathbf{X}\varphi &\text{ iff } w[1, \infty) \models_{LTL} \varphi \\ w \models_{LTL} \varphi\mathbf{U}\psi &\text{ iff } \exists i \geq 0 \text{ such that, } w[i, \infty) \models_{LTL} \psi \text{ and } \forall k \ 0 < k < i, w[k, \infty) \models \varphi \end{aligned}$$

Just as in the case of regular expressions and ω -regular expressions, a satisfaction relation needs to be defined for LTL formulas, for checking whether an infinite word over 2^{AP} satisfies an LTL

formula φ . Here also, we restrict infinite words to only ultimately periodic words, because, they have certain structural properties, which aid the process of checking satisfaction.

Due to observation 4.1, it is enough to define the satisfaction relation for only those infinite suffix $uv^\omega[i, \infty)$ of uv^ω , where $0 \leq i \leq |uv| - 1$.

Definition 3. The satisfaction relation $\models_{\subseteq} (2^{\text{AP}})^\omega \times \mathcal{F}_{\text{LTL}}$ is defined inductively on the structure of an LTL formula φ in the following manner:

$$uv^\omega[i, \infty) \models p \iff p \in uv^\omega[i, i)$$

$$uv^\omega[i, \infty) \models \neg\varphi \iff uv^\omega[i, \infty) \not\models \varphi$$

$$uv^\omega[i, \infty) \models \varphi \vee \psi \iff uv^\omega[i, \infty) \models \varphi \text{ or } uv^\omega[i, \infty) \models \psi$$

$$uv^\omega[i, \infty) \models \mathbf{X}\varphi \iff \begin{cases} uv^\omega[i+1, \infty) \models \varphi, & \text{for } 0 \leq i < |uv| - 1 \\ uv^\omega[|u|, \infty) \models \varphi, & \text{for } i = |uv| - 1 \end{cases}$$

$$uv^\omega[i, \infty) \models \varphi \mathbf{U} \psi \iff \begin{cases} \exists j, i \leq j \leq |uv| - 1, uv^\omega[j, \infty) \models \psi \text{ and } \forall k, i < k < j, uv^\omega[k, \infty) \models \varphi & \text{for } i < |u| \\ \exists j, |u| \leq j \leq |uv| - 1, uv^\omega[j, \infty) \models \psi \text{ and } \forall k \in \mathcal{I}_{u,v}(i, j), uv^\omega[k, \infty) \models \varphi & \text{for } i \geq |u| \end{cases}$$

The set $\mathcal{I}^{u,v}(i, j)$ used in defining the satisfaction relation is defined as follows:

$$\mathcal{I}^{u,v}(i, j) = \begin{cases} \{i, \dots, j-1\} & \text{for } i \leq j \\ \{|u|, \dots, j-1\} \cup \{i, \dots, |uv| - 1\} & \text{for } i > j \end{cases}$$

Now, it needs to be shown that the satisfaction relation so defined is actually consistent with the semantics of LTL. Precisely, we need to prove, $uv^\omega[i, \infty) \models \varphi \iff uv^\omega[i, \infty) \models_{\text{LTL}} \varphi$. But, this is not difficult to see and requires an induction on the structure of φ .

5.2 The Learning Problem

The problem statement for this case remains almost the same as in ω -regular expressions. Given a sample $\mathcal{S} = (\mathcal{P}, \mathcal{N})$ consisting of two finite, disjoint sets $\mathcal{P}, \mathcal{N} \subset \Sigma^*$ of ultimately periodic words such that $\mathcal{P} \cup \mathcal{N} \neq \emptyset$, we need to find a LTL formula φ over AP consistent with a sample \mathcal{S} .

The solution for this learning problem is not much different from the ones described so far in this thesis. In fact, the same Algorithm 1 can be used to learn a LTL formula from the given sample. The only difference is in how the formula $\Phi_n^{\mathcal{S}}$ is formulated. The formula has two parts as usual.

$$\Phi_n^{\mathcal{S}} = \Phi_n^{\text{structure}} \wedge \Phi_n^{\text{consistency}}$$

The first part $\Phi_n^{\text{structure}}$ which encodes a syntax DAG of a LTL formula has the following variables.

- $x_{p,\lambda}$ where $p \in \{1, \dots, n\}$ and $\lambda \in \Lambda_L$
- $l_{p,q}, r_{p,q}$ where $p \in \{2, \dots, n\}$ and $q \in \{1, \dots, p-1\}$

Now, the variables should satisfy certain constraints in order to form a valid syntax DAG. But, these constraints are precisely the ones that needed to be satisfied by the syntax DAG of regular expressions. Thus, the exact description of the formulas have been skipped.

The second part of the formula checks whether the guessed syntax DAG of the LTL formula is consistent with the sample. Φ_n^w is a disjunction of following formulas.

$$\bigwedge_{1 \leq p \leq n} \bigwedge_{p \in \text{AP}} x_{p,p} \rightarrow \left[\bigwedge_{0 \leq i < |uv|} \begin{cases} y_{i,p}^{u,v} & \text{if } p \in uv[i, i) \\ \neg y_{i,p}^{u,v} & \text{if } p \notin uv[i, i) \end{cases} \right] \quad (5.1)$$

$$\bigwedge_{\substack{1 \leq p \leq n \\ 1 \leq q, q' \leq p}} x_{p,\vee} \wedge l_{p,q} \wedge r_{p,q'} \rightarrow \left[\bigwedge_{0 \leq i < |uv|} \left[y_{i,p}^{u,v} \leftrightarrow y_{i,q}^{u,v} \vee y_{i,j,q'}^{u,v} \right] \right] \quad (5.2)$$

$$\bigwedge_{\substack{1 \leq p \leq n \\ 1 \leq q \leq p}} x_{p,\times} \wedge l_{p,q} \rightarrow \left[\bigwedge_{0 \leq i < |uv|-1} \left[y_{i,p}^{u,v} \leftrightarrow y_{i+1,q}^{u,v} \right] \wedge \left[y_{|uv|-1,p}^{u,v} \leftrightarrow y_{|u|,q}^{u,v} \right] \right] \quad (5.3)$$

$$\bigwedge_{\substack{1 \leq p \leq n \\ 1 \leq q, q' \leq p}} x_{p,\cup} \wedge l_{p,q} \wedge r_{p,q'} \rightarrow \left[\bigwedge_{0 \leq i \leq |u|} \left[y_{i,p}^{u,v} \leftrightarrow \bigvee_{i \leq j < |uv|} \left[y_{j,q}^{u,v} \wedge \bigwedge_{i \leq k < j} y_{k,q}^{u,v} \right] \right] \right. \\ \left. \wedge \left[\bigwedge_{|u| \leq i < |uv|} \left[y_{i,p}^{u,v} \leftrightarrow \bigvee_{|u| \leq j < |uv|} \left[y_{j,q}^{u,v} \wedge \bigwedge_{k \in \mathcal{I}^{u,v}(i,j)} y_{k,q}^{u,v} \right] \right] \right] \right] \quad (5.4)$$

These formulas originate from how the satisfaction relation could be used to check whether word uv^w satisfies the guessed formula. The final step is to collate all the formulas and ensure consistency with the sample \mathcal{S} .

$$\Phi_n^{\text{consistency}} = \left[\bigwedge_{w \in \mathcal{P}} \Phi_n^w \wedge y_{0,n}^w \right] \wedge \left[\bigwedge_{w \in \mathcal{N}} \Phi_n^w \wedge \neg y_{0,n}^w \right] \quad (5.5)$$

The proof of correctness for the learning algorithm for LTL formulas goes along the same lines as the proof of Claim 1. The complete proof can be found in [6].

Chapter 6

Learning Program Specification Language

6.1 Property Specification Language

Property Specification Language (PSL)[4] is a language for the formal specification of hardware. It is used to describe properties that are required to hold in the design under verification. PSL subsumes LTL and has an increased expressive power. While the expressive power of LTL is that of star-free ω -regular expressions, the expressive power of PSL is the same as the full class of ω -regular expressions[2]. As PSL is a widely used standard it contains a large amount of derived operators; we concentrate here only on a subset of operators that suffice to obtain the expressive power of ω -regular expressions. We call this subset corePSL. The syntax of corePSL is given by the following grammar:

$$\varphi ::= \mathbf{p} \mid \varphi \vee \varphi \mid \neg\varphi \mid \mathbf{X}\varphi \mid \varphi_1 \mathbf{U}\varphi_2 \mid r \mapsto \varphi$$

where $\mathbf{p} \in \text{AP}$, where Let AP be a set of atomic propositions and r is a regular expression over 2^{AP} . The operator \mapsto is called *suffix implication* or *trigger*. The set of all PSL formulas over AP is referred to as \mathcal{F}_{PSL} . PSL uses operators from both regular expressions and from LTL. Hence, $\Lambda_P = \Lambda_R \cup \Lambda_L\{\mapsto\}$, where $\Sigma = 2^{\text{AP}}$, that is the alphabet for the regular expressions in PSL is the subsets of the propositions in AP.

Abbreviations We mention only the main abbreviations related to the suffix implication operator.

$$\begin{aligned} r \Vdash \varphi &::= r \mapsto (\text{true} \cdot \varphi) \\ r \diamond \rightarrow \varphi &::= \neg(r \mapsto \neg\varphi) \\ r &::= r \diamond \rightarrow \text{true} \end{aligned}$$

Semantics Given a word $w \in (2^{\text{AP}})^\omega$ and a corePSL formula φ , we define the relation $w \models_{\text{PSL}} \varphi$ (read “ w satisfies φ ”) similarly to the way it is defined for LTL. We provide the semantics only for the new operator, the suffix implication operator.

$$w \models_{\text{PSL}} r \mapsto \varphi \text{ iff } \forall j \in \mathbb{N}, (w[0, j] \in \llbracket r \rrbracket \implies w[j - 1, \infty) \models_{\text{PSL}} \varphi)$$

That is, w satisfies $r \mapsto \varphi$ if for every prefix u of w that matches the regular expression r , the suffix of w starting where u ends (with one letter overlap) satisfies the formula φ .

Definition 4. The satisfaction relation $\models_{\subseteq} (2^{\text{AP}})^{\omega} \times \mathcal{F}_{\text{PSL}}$ is defined inductively on the structure of an PSL formula φ . Since, the satisfaction relation for the LTL operators in PSL remain the same—satisfaction relation is defined only for the \mapsto operator.

$$w^{\omega}[i, \infty) \models r \mapsto \varphi \iff \forall j \in \mathbb{N}, j < |u| + |v| + n|v|, (w^{\omega}[i, j) \vdash r \implies w^{\omega}[j - 1, \infty) \models \varphi)$$

6.2 The Learning Problem

The learning problem in this case is same as LTL. Hence, the approach taken here is same as the one for LTL. The problem is solved by encoding it as a series of satisfiability problem and then using a SAT solver. The formula designed has two parts. The first part encodes the structure of the PSL formula and is similar to the one for LTL, with the added PSL operators and operators for regular expression.

The interesting part is to design $\Phi_n^{\text{consistency}}$, which essentially depends on how the satisfaction relation is defined. There are two type of variables used in the formula.

- $y_{i,p}^{u,v}$ where $0 \leq i \leq |uv| - 1$
- $z_{i,j,p}^{uv^b}$ where $0 \leq i \leq j \leq |uv^b|$.

Let, $b = n + 1$ for simplicity. The variables $y_{i,p}^{u,v}$ ensure satisfaction of the guessed PSL formula with a word w^{ω} present in sample while the variables $z_{i,j,p}^{uv^b}$ ensure matching of the regular expressions in the formula with certain finite infixes of the word w^{ω} . Φ_n^w is a disjunction of the formulas described below.

$$\bigwedge_{\substack{1 \leq p \leq n \\ 1 \leq q, q' \leq p}} x_{p, \mapsto} \wedge l_{p,q} \wedge r_{p,q'} \rightarrow \left[\bigwedge_{0 \leq i < |uv|} \left[y_{i,p}^{u,v} \leftrightarrow \bigwedge_{i \leq j \leq |uv^b|} \left[z_{i,j,q}^{uv^b} \rightarrow y_{M(j-1), q'}^{u,v} \right] \right] \right] \quad (6.1)$$

$$\bigwedge_{1 \leq p \leq n} x_{p, \varepsilon} \rightarrow \left[\bigwedge_{0 \leq i \leq j \leq |uv^b|} z_{i,j,p}^{uv^b} \leftrightarrow [i = j] \right] \quad (6.2)$$

$$\bigwedge_{1 \leq p \leq n} \bigwedge_{a \in 2^{\text{AP}}} x_{p,a} \rightarrow \left[\bigwedge_{0 \leq i \leq j \leq |uv^b|} \begin{cases} z_{i,j,p}^{uv^b} & \text{if } w[i, j) = a \\ \neg z_{i,j,p}^{uv^b} & \text{if } w[i, j) \neq a \end{cases} \right] \quad (6.3)$$

$$\bigwedge_{\substack{1 \leq p \leq n \\ 1 \leq q, q' \leq p}} x_{p,+} \wedge l_{p,q} \wedge r_{p,q'} \rightarrow \left[\bigwedge_{0 \leq i \leq j \leq |uv^b|} \left[z_{i,j,p}^{uv^b} \leftrightarrow z_{i,j,q}^{uv^b} \vee z_{i,j,q'}^{uv^b} \right] \right] \quad (6.4)$$

$$\bigwedge_{\substack{1 \leq p \leq n \\ 1 \leq q, q' \leq p}} x_{p,o} \wedge l_{p,q} \wedge r_{p,q'} \rightarrow \left[\bigwedge_{0 \leq i \leq j \leq |uv^b|} \left[z_{i,j,p}^{uv^b} \leftrightarrow \bigvee_{i \leq k \leq j} z_{i,k,q}^{uv^b} \wedge z_{k,j,p}^{uv^b} \right] \right] \quad (6.5)$$

$$\bigwedge_{\substack{1 \leq p \leq n \\ 1 \leq q, q' \leq p}} x_{p,*} \wedge l_{p,q} \rightarrow \left[\bigwedge_{0 \leq i \leq j \leq |uv^b|} \left[z_{i,j,p}^{uv^b} \leftrightarrow [i = j] \vee \bigvee_{i < k \leq j} z_{i,k,q}^{uv^b} \wedge z_{k,j,p}^{uv^b} \right] \right] \quad (6.6)$$

It can be observed that for operators of regular expressions, we have used formulas which are motivated by Formulas 3.6 to 3.10. Also, $M(j)$ has been defined in 4.9. Finally $\Phi_S^{\text{consistency}}$ is defined as 5.5, as for LTL. In other words, given a sample \mathcal{S} , we construct a formula $\Phi_n^{\mathcal{S}}$

Claim 3. Let $\mathcal{S} = (\mathcal{P}, \mathcal{N})$ be a sample, $n \in \mathbb{N} \setminus \{0\}$, and $\Phi_n^{\mathcal{S}}$ be the propositional formula defined above. Then, the following holds:

1. If there exists a PSL formula of size n , $\varphi^{\mathcal{S}}$ that is consistent with \mathcal{S} , then the propositional formula $\Phi_n^{\mathcal{S}}$ is satisfiable.
2. If $v \models \Phi_n^{\mathcal{S}}$, then φ^v is a PSL formula of size n that is consistent with \mathcal{S} .

The proof has been presented in the Appendix B

Chapter 7

Conclusion

The main problem talked about in this thesis is the following: given a sample \mathcal{S} consisting of two finite sets of positive and negative examples, learn a model \mathcal{M} that is consistent with \mathcal{S} . The problem has been looked at for the following models: regular expressions, ω -regular expressions, Linear Temporal Logic, Property Specification Language.

There are certain interesting directions of future work which emerge from the work presented in thesis. An implementation of the theoretical results presented in the thesis using SAT solvers such as Z3 could be of great use. Checking whether the bounds mentioned in the lemmas and theorems of Chapter 4 are actually tight—could have theoretical relevance. Moreover, the learning techniques used in this thesis could potentially be applied to many more models including timed regular expressions[5].

Appendix A

Proofs from Chapter 4

Lemma 5. $uv^\omega[i, \infty) \vdash r_\omega \iff uv^\omega[i, \infty) \in \llbracket r_\omega \rrbracket$

Proof. The proof goes via induction on the structure of the the ω -regular expression.

- Case $r_\omega = r^\omega$: When $i < |u|$, let $u[i, |u|) = u'$. Now, we need to prove $u'v \vdash r^\omega$. When $i \geq |u|$, we use Theorem 1 for the proof.
- Case $r_\omega = r \circ r'_\omega$: Consider the following implications:

$$\begin{aligned} uv^\omega[i, \infty) \vdash r \circ r'_\omega &\iff \exists i \leq j \leq |u| + |v| + n|v|, uv^\omega[i, j) \vdash r \text{ and } uv^\omega[j, \infty) \vdash r'_\omega \\ &\iff \exists i \leq j \leq |u| + n|v|, uv^\omega[i, j) \in \llbracket r \rrbracket \text{ and } uv^\omega[j, \infty) \in \llbracket r'_\omega \rrbracket && \text{[Using ind. hypothesis]} \\ &\iff \exists i \leq j, uv^\omega[i, j) \in \llbracket r \rrbracket \text{ and } uv^\omega[j, \infty) \in \llbracket r'_\omega \rrbracket && \text{[Using Lemma 4]} \\ &\iff v^\omega[i, \infty) \in \llbracket r \circ r'_\omega \rrbracket \end{aligned}$$

An extra term $|v|$ gets added in the inequality in lines 1 and 2, because $|u| \leq i < |uv|$ in which case the inequality in Lemma 3 gets shifted by $|v|$ factor.

- Case $r_\omega = r'_\omega + r''_\omega$: Consider the following implications:

$$\begin{aligned} uv^\omega[i, \infty) \vdash r'_\omega + r''_\omega &\iff v^\omega[i, \infty) \vdash r'_\omega \text{ or } uv^\omega[i, \infty) \vdash r''_\omega \\ &\iff v^\omega[i, \infty) \in \llbracket r'_\omega \rrbracket \text{ or } uv^\omega[i, \infty) \in \llbracket r''_\omega \rrbracket && \text{[Using induction hypothesis]} \\ &\iff uv^\omega[i, \infty) \in \llbracket r'_\omega + r''_\omega \rrbracket \end{aligned}$$

□

Appendix B

Proof of Claim 2:

Proof. Using the syntax DAG of the expression R_ω^S which is indexed using $1 \cdots n$, we formulate a valuation v for the propositional variables in Φ_n^S . We use R_p^S to refer to the expression rooted at the p^{th} node.

- We set $v(x_{p,\lambda}) = 1$ iff the p^{th} node is labelled by λ .
- We set $v(l_{p,q}) = 1$ iff q^{th} node is the left child of the p^{th} node and similarly, $v(r_{p,q}) = 1$ iff q^{th} node is the right child of the p^{th} node.
- We set $v(y_{i,p}^{u,v}) = 1$ iff $uv^\omega[i, j] \models R_p^S$, when label at p is an infinite operator.
- We set $v(z_{i,j,p}^{uv^b}) = 1$ iff $uv^b[i, j] \models R_p^S$, when label at p is a finite operator or an alphabet.
- We set $v(\bar{z}_{i,j,p}^{uv^b}) = 1$ iff $uv^b[i, j] \models (R_q^S)^+$, when label at p is an ω operator and q^{th} node is the left child of the p^{th} node.

Firstly, it can be seen that $v \models \Phi_n^{\text{structure}}$, since the formulated valuation ensures the uniqueness of the labels of the nodes as well as that of the left and right children. Moreover, $v \models \Phi_n^w$ for $w \in \mathcal{P} \cup \mathcal{N}$, since, the values of the variables $z_{i,j,p}^{uv^b}$ correspond exactly to the valuation of each subexpression R_p^S on the subword $w[i, j]$ for finite operators and alphabets; and the values of the variables $y_{i,p}^{u,v}$ correspond exactly to the valuation of each subexpression R_p^S on the subword $w[i, \infty)$ for infinite operators. Finally, the fact that R_ω^S is consistent with \mathcal{S} implies $v \models y_{0,n}^{u,v}$ for each $w = uv^\omega \in \mathcal{P}$ and $v \models \neg y_{0,n}^{u,v}$ for each $w = uv^\omega \in \mathcal{N}$. Thus, $v \models \Phi_n^S$, which proves that Φ_n^S is satisfiable.

For the second statement, since, $v \models \Phi_n^S$, we have $v \models \Phi_n^{\text{structure}}$ as well. Hence, the valuation of the variables $x_{p,\lambda}$, $l_{p,q}$, and $r_{p,q}$ encode a syntax DAG from which we obtain the ω -regular expression R^v . Again here, we consider R_p^v to be the subexpression of R^v rooted at the p^{th} node. Now, it needs to be shown that R^v is indeed consistent with the sample \mathcal{S} . For that we show, $v(y_{i,p}^{u,v}) = 1 \iff uv^\omega[i, \infty) \vdash R_p^v$ for any $i \in \{0, \dots, |uv| - 1\}$ for the p^{th} node which is labelled by an infinite operators. For this proof, we consider only those nodes which are labelled by infinite operators. For the finite operators, we need to show $v(z_{i,j,p}^{uv^b}) = 1 \iff uv^b[i, j] \vdash R_p^v$. But, this is exactly done in proof of Claim 1. This can be done using induction on the structure of R^v .

- *Case $R_p^v = (R_q^v)^\omega$:* In this case, $v(x_{p,\omega})$ and $v(l_{p,q})$ are set to 1. There could be two cases depending on value of i . In, both the cases we consider that $k \equiv j \pmod{|v|}$.

– When $i < |u|$

$$\begin{aligned} v(y_{i,p}^{u,v}) = 1 &\iff \exists j, k, |u| \leq j \leq k \leq |uv^b|, v(\bar{z}_{i,j,q}^{uv^b}) = 1 \text{ and } v(\bar{z}_{j,k,q}^{uv^b}) = 1 && \text{[Using Subform. 4.6]} \\ &\iff \exists j, k, |u| \leq j \leq k \leq |uv^b|, uv^b[i, j] \vdash (R_q^v)^+ \text{ and } uv^b[j, k] \vdash (R_q^v)^+ && \text{[Using ind. hypo.]} \\ &\iff (R_q^v)^\omega && \text{[Using Theorem 2]} \end{aligned}$$

– When $i > |u|$

$$\begin{aligned} v(y_{i,p}^{u,v}) = 1 &\iff \exists j, k, i \leq j \leq k \leq |uv^b|, v(\bar{z}_{i,j,q}^{uv^b}) = 1 \text{ and } v(\bar{z}_{j,k,q}^{uv^b}) = 1 && \text{[Using Subform. 4.6]} \\ &\iff \exists j, k, i \leq j \leq k \leq |uv^b|, uv^b[i, j] \vdash (R_q^v)^+ \text{ and } uv^b[j, k] \vdash (R_q^v)^+ && \text{[Using ind. hypo.]} \\ &\iff (R_q^v)^\omega && \text{[Using Theorem 1]} \end{aligned}$$

In both the cases above, we assumed that $v(\bar{z}_{i,j,q}^{uv^b}) = 1$ iff $w[i, j] \vdash (R_q^v)^+$ in the second step, which cannot be deduced from the induction on the structure of R_p^v . This is obtained in the following manner:

$$\begin{aligned} v(\bar{z}_{i,j,p}^{uv^b}) = 1 &\iff \exists k \in \mathbb{N}, i \leq k \leq j, v(z_{i,k,q}^{uv^b}) = 1 \text{ and } v(\bar{z}_{k,j,p}^{uv^b}) = 1 && \text{[Using Subformula (4.6)]} \\ &\iff \exists k \in \mathbb{N}, i \leq k \leq m, w[i, k] \vdash R_q^v \text{ and } w[k, j] \vdash (R_q^v)^+ && \text{[Using ind. hypothesis]} \\ &\iff w[i, j] \vdash (R_q^v)^+ && \text{[Using Def. 1]} \end{aligned}$$

Here, $v(\bar{z}_{k,j,p}^{uv^b}) = 1$ iff $w[k, j] \vdash (R_q^v)^+$ requires another level of induction, similar to the last case of the proof of Claim 1.

– Case $R_p^v = R_q^v \circ_\omega R_{q'}^v$: In this case, $v(x_{p,\circ_\omega}), v(l_{p,q}), v(r_{p,q'})$ are all set to 1 and this leads to the following implications:

$$\begin{aligned} v(y_{i,p}^{u,v}) = 1 &\iff \exists k \in \mathbb{N}, 1 \leq j \leq |uv^c|, v(z_{i,k,q}^{u,v}) = 1 \text{ and } v(y_{M(j),q'}^{u,v}) = 1 && \text{[Using Subformula 4.7]} \\ &\iff \exists k \in \mathbb{N}, 1 \leq k \leq m, w[i, j] \vdash R_q^v \text{ and } w[M(j), \infty] \vdash R_{q'}^v && \text{[Using ind. hypothesis]} \\ &\iff \exists k \in \mathbb{N}, 1 \leq k \leq m, w[i, j] \vdash R_q^v \text{ and } w[j, \infty] \vdash R_{q'}^v && \text{[Using Observation 4.1]} \\ &\iff w[i, j] \vdash R_q^v \circ_\omega R_{q'}^v && \text{[Using Def. 2]} \end{aligned}$$

– Case $R_p^v = R_q^v +_\omega R_{q'}^v$: In this case, $v(x_{p,+_\omega}), v(l_{p,q}), v(r_{p,q'})$ are all set to 1, and this leads to the following implications:

$$\begin{aligned} v(y_{i,p}^{u,v}) = 1 &\iff v(y_{i,q}^{u,v}) = 1 \text{ or } v(y_{i,q'}^{u,v}) = 1 && \text{[Using Subformula 4.8]} \\ &\iff w[i, \infty] \vdash R_q^v \text{ or } w[i, \infty] \vdash R_{q'}^v && \text{[Using ind. hypothesis]} \\ &\iff w[i, j] \vdash R_q^v +_\omega R_{q'}^v && \text{[Using Def. 2]} \end{aligned}$$

□

Proof of Claim 3

Proof. Using the syntax DAG of the expression φ^S which is indexed using $1 \cdots n$, we formulate a valuation v for the propositional variables in Φ_n^S . We use φ_p^S to refer to the formula rooted at the p^{th} node.

- We set $v(x_{p,\lambda}) = 1$ iff the p^{th} node is labelled by λ .
- We set $v(l_{p,q}) = 1$ iff q^{th} node is the left child of the p^{th} node and similarly, $v(r_{p,q}) = 1$ iff q^{th} node is the right child of the p^{th} node.
- We set $v(y_{i,p}^{u,v}) = 1$ iff $uv^\omega[i, \infty) \vdash \varphi_p^S$, when label at p is an LTL operator or a triggers(\mapsto) operator or a proposition in AP.
- We set $v(z_{i,j,p}^{uv^b}) = 1$ iff $uv^b[i, j) \vdash \varphi_p^S$, when label at p is a operator from regular expressions or a letter of the alphabet 2^{AP} .

Firstly, it can be seen that $v \models \Phi_n^{\text{structure}}$, since the formulated valuation ensures the uniqueness of the labels of the nodes as well as that of the left and right children. Moreover, $v \models \Phi_n^w$ for $w \in \mathcal{P} \cup \mathcal{N}$, since, the values of the variables $z_{i,j,p}^{uv^b}$ correspond exactly to the valuation of each subformula φ_p^S on the subword $uv^b[i, j)$ for \mapsto , LTL operators and letters of alphabet 2^{AP} ; and the values of the variables $y_{i,p}^{u,v}$ correspond exactly to the valuation of each subformula φ_p^S on the subword $uv^\omega[i, \infty)$ for \mapsto , LTL operators or a proposition in AP. Finally, the fact that φ^S is consistent with \mathcal{S} implies $v \models y_{0,n}^{u,v}$ for each $w = uv^\omega \in \mathcal{P}$ and $v \models \neg y_{0,n}^{u,v}$ for each $w = uv^\omega \in \mathcal{N}$. Thus, $v \models \Phi_n^S$, which proves that Φ_n^S is satisfiable.

For the second statement, since, $v \models \Phi_n^S$, we have $v \models \Phi_n^{\text{structure}}$ as well. Hence, the valuation of the variables $x_{p,\lambda}$, $l_{p,q}$, and $r_{p,q}$ encode a syntax DAG from which we obtain the PSL formula φ^v . Again here, we consider φ_p^v to be the subformula of φ^v rooted at the p^{th} node. Now, it needs to be shown that φ^v is indeed consistent with the sample \mathcal{S} . For that we show, $v(y_{i,p}^{u,v}) = 1 \iff uv^\omega[i, \infty) \vdash \varphi_p^v$ for any $i \in \{0, \dots, |uv| - 1\}$ for the p^{th} node which is labelled by \mapsto , LTL operators or a proposition in AP. For the regular operators and its corresponding alphabet, we need to show $v(z_{i,j,p}^{uv^b}) = 1 \iff uv^b[i, j) \vdash \varphi_p^v$. But, this is exactly done in proof of Claim 1. For this proof, we consider only those nodes which are labelled by only the \mapsto operator. This can be done using induction on the structure of φ^v .

- *Base case* $\varphi_p^v = \mathbf{p}$: In this case, $v(x_{p,\mathbf{p}}) = 1$, which leads to the following implications:

$$\begin{aligned} v(y_{i,p}^{u,v}) = 1 &\iff \mathbf{p} \in uv^\omega[i, i) && \text{[Using Subformula (5.1)]} \\ &\iff uv^\omega[i, j) \models \mathbf{p} && \text{[Using Def. 3]} \end{aligned}$$

- *Case* $\varphi_p^v = \varphi_q^v \vee \varphi_{q'}^v$: In this case, $v(x_{p,\vee})$, $v(l_{p,q})$, $v(r_{p,q'})$ are all set to 1, and this leads to the following implications:

$$\begin{aligned} v(y_{i,p}^{u,v}) = 1 &\iff v(y_{i,q}^{u,v}) = 1 \text{ or } v(y_{i,q'}^{u,v}) = 1 && \text{[Using Subformula (5.2)]} \\ &\iff uv^\omega[i, \infty) \models \varphi_q^v \text{ or } uv^\omega[i, \infty) \models \varphi_{q'}^v && \text{[Using ind. hypothesis]} \\ &\iff uv^\omega[i, \infty) \models \varphi_q^v \vee \varphi_{q'}^v && \text{[Using Def. 3]} \end{aligned}$$

- *Case* $\varphi_p^v = \mathbf{X}(\varphi_q^v)$: In this case, $v(x_{p,\mathbf{X}})$ and $v(l_{p,q})$ are set to 1 and this leads to the following

implications:

$$\begin{aligned}
v(y_{i,p}^{u,v}) = 1 &\iff \begin{cases} v(y_{i+1,q}^{u,v}) = 1 \text{ for } 0 \leq i < |uv| - 1 \\ v(y_{|u|,q}^{u,v}) = 1 \text{ for } i = |uv| - 1 \end{cases} && \text{[Using Subformula (5.3)]} \\
&\iff \begin{cases} uv^\omega[i+1, \infty) \models \varphi_q^v \text{ for } 0 \leq i < |uv| - 1 \\ uv^\omega[|u|, \infty) \models \varphi_q^v \text{ for } i = |uv| - 1 \end{cases} && \text{[Using ind. hypothesis]} \\
&\iff uv^\omega[i, \infty) \models \mathbf{X}(\varphi_q^v) && \text{[Using Def. 3]}
\end{aligned}$$

- Case $\varphi_p^v = \varphi_q^v \mathbf{U} \varphi_{q'}^v$: In this case, $v(x_{p,\cup}), v(l_{p,q}), v(r_{p,q'})$ are all set to 1, and this leads to the following implications:

$$\begin{aligned}
v(y_{i,p}^{u,v}) = 1 &\iff \begin{cases} \exists j, i \leq j \leq |uv| - 1, v(y_{j,q'}^{u,v}) = 1 \text{ and } \forall k, i \leq k < j, v(y_{k,q}^{u,v}) = 1 \text{ for } i < |u| & \text{[Using} \\ \exists j, |u| \leq j < |uv|, v(y_{j,q'}^{u,v}) = 1 \text{ and } \forall k, k \in \mathcal{I}^{u,v}(i, j), v(y_{k,q}^{u,v}) = 1 \text{ for } i \geq |u| & \text{Subform. (5.4)]} \end{cases} \\
&\iff \begin{cases} \exists j, i \leq j \leq |uv| - 1, uv^\omega[j, \infty) \models \varphi_{q'}^v \text{ and } \forall k, i \leq k < j, uv^\omega[k, \infty) \models \varphi_q^v \text{ for } i < |u| & \text{[Using} \\ \exists j, |u| \leq j < |uv|, uv^\omega[j, \infty) \models \varphi_{q'}^v \text{ and } \forall k, k \in \mathcal{I}^{u,v}(i, j), uv^\omega[k, \infty) \models \varphi_q^v \text{ for } i \geq |u| & \text{ind. hyp} \end{cases} \\
&\iff uv^\omega[i, j) \models \varphi_q^v \mathbf{U} \varphi_{q'}^v
\end{aligned}$$

- Case $\varphi_p^v = \varphi_q^v \mapsto \varphi_{q'}^v$: In this case, $v(x_{p,\rightarrow}), v(l_{p,q}), v(r_{p,q'})$ are all set to 1, and this leads to the following implications:

$$\begin{aligned}
v(y_{i,p}^{u,v}) = 1 &\iff \forall j, i \leq j \leq |uv^b|, v(z_{i,j,q}^{uv^b}) = 1 \implies v(y_{M(j-1),q'}^{u,v}) = 1 && \text{[Using Subformula (6)} \\
&\iff \forall j, i \leq j \leq |u| + |v| + n|v|, uv^b[i, j) \vdash \varphi_q^v \implies uv^\omega[M(j-1), q') \models \varphi_{q'}^v && \text{[Using ind. hypothesis]} \\
&\iff \forall j, i \leq j \leq |u| + |v| + n|v|, uv^b[i, j) \vdash \varphi_q^v \implies uv^\omega[j-1, q') \models \varphi_{q'}^v && \text{[Using Observation 4.]} \\
&\iff uv^\omega[i, \infty) \models \varphi_q^v \mapsto \varphi_{q'}^v && \text{[Using Def. 3]}
\end{aligned}$$

Here, we assumed that $v(y_{k,j,p}^w) = 1$ iff $w[k, j) \models \varphi_p^v$ in the second step, which cannot be deduced from the induction on the structure of φ_p^v . This we obtained using another level of induction, which is on the length of the subword that matches φ_p^v . More precisely, we induct on k which ranges from j to $i+1$, where by induction hypothesis, it is assumed that

$$v(y_{k,j,p}^w) = 1 \iff w[k, j) \models \varphi_p^v \quad \forall k, i < k \leq j$$

The base case occurs when $k = j$, and then, we have $v(y_{k,j,p}^w) = 1 \iff k = j \iff w[k, j) \models \varphi_p^v$.

□

Bibliography

- [1] ALVES, T., TERESA MARTINS, A., AND MARTINS FERREIRA, F. On finding a first-order sentence consistent with a sample of strings. vol. 277, pp. 220–234.
- [2] ARMONI, R., FIX, L., FLAISHER, A., GERTH, R., GINSBURG, B., KANZA, T., LANDVER, A., MADOR-HAIM, S., SINGERMAN, E., TIEMEYER, A., VARDI, M. Y., AND ZBAR, Y. The forspec temporal logic: A new temporal property-specification language. In *TACAS (2002)*.
- [3] BAIER, C., AND KATOEN, J.-P. *Principles of Model Checking*, vol. 26202649. 01 2008.
- [4] EISNER, C., AND FISMAN, D. *A Practical Introduction to PSL*. 07 2006.
- [5] NARAYAN, A., CUTULENCO, G., JOSHI, Y., AND FISCHMEISTER, S. Mining timed regular specifications from system traces. *ACM Trans. Embed. Comput. Syst.* 17, 2 (Jan. 2018), 46:1–46:21.
- [6] NEIDER, D., AND GAVRAN, I. Learning linear temporal properties. In *2018 Formal Methods in Computer Aided Design, FMCAD 2018, Austin, TX, USA, October 30 - November 2, 2018 (2018)*, pp. 1–10.
- [7] PNUELI, A. The temporal logic of programs. In *Proceedings of the 18th Annual Symposium on Foundations of Computer Science (Washington, DC, USA, 1977), SFCS '77, IEEE Computer Society*, pp. 46–57.
- [8] WOLPER, P. Temporal logic can be more expressive. In *Proceedings of the 22Nd Annual Symposium on Foundations of Computer Science (Washington, DC, USA, 1981), SFCS '81, IEEE Computer Society*, pp. 340–348.