

MAT Learners for Tree Series – an Abstract Data Type and Two Realizations

Frank Drewes¹, Johanna Högberg¹, and Andreas Maletti^{2,*}

¹ Department of Computing Science, Umeå University
S-901 87 Umeå, Sweden, {drewes, johanna}@cs.umu.se

² Universitat Rovira i Virgili, Departament de Filologies Romàniques
Av. Catalunya 35, 43002 Tarragona, Spain, andreas.maletti@urv.cat

Abstract. We propose abstract observation tables, an abstract data type for learning deterministic weighted tree automata in Angluin’s *minimal adequate teacher* (MAT) model, and show that every correct implementation of abstract observation tables yields a correct MAT learner. Besides the “classical” observation table, we show that abstract observation tables can also be implemented by observation trees. The advantage of the latter is that they often require fewer queries to the teacher.

1 Introduction

We study the problem of learning deterministic weighted (bottom-up) tree automata (dwta) [5, 7] over a semifield \mathbb{S} . A weighted tree automaton computes a recognizable tree series [4], i.e., a function that maps trees to values in \mathbb{S} . This is accomplished by assigning a weight in \mathbb{S} to every transition. The weight of a computation (also called a run) is the product of all its transitions, multiplied with an additional final weight that is associated with the final state reached. Finally, the weight of a tree is the sum of the weights of all its runs. In other words, the addition of the semifield is used to handle nondeterminism. Ordinary bottom-up tree automata correspond to the special case obtained by choosing the Boolean semifield as \mathbb{S} . We recommend [14] for an introduction to recognizable tree series.

As mentioned above, this article is devoted to the problem of algorithmically learning dwta, also called grammatical inference. The general setting considered in grammatical inference is characterized by a series ψ (or, in the traditional case, a language) for which, a priori, no explicit representation is available. The learning algorithm (henceforth called the learner) only has access to some restricted information. Its goal is to derive, from this information, an automaton computing ψ .

There exist various learning models that make this general picture more precise. Most of them fall into one of the following categories: Gold’s *learning from examples with identification in the limit* [16], Valiant’s *probably approximately correct (PAC) learning* [23], and Angluin’s *query learning* [2]. In this article, we focus on the most prominent type of query learning, proposed in [1], in which

* Supported by the *Ministerio de Educación y Ciencia* (MEC) grant JDCI-2007-760.

the learner has access to an oracle called minimal adequate teacher (MAT). Angluin’s original model was designed for learning ordinary deterministic finite automata (dfa). Its first generalization to bottom-up tree automata was proposed by Sakakibara in [21] and later improved in [12].

In the weighted case, if ψ is the target tree series, then a MAT is an oracle able to answer two types of queries. A *coefficient query* passes a tree t to the MAT and receives $\psi(t)$ as the answer. An *equivalence query* passes a dwta A , the hypothesis, to the MAT. If A computes ψ , a special token is returned. Otherwise, the answer is a counterexample, i.e., a tree t such that $A(t) \neq \psi(t)$. One of the first extensions of the MAT learner for dfa to stochastic automata appears in [9]. This algorithm is extended in [20] to certain cancellative semirings [15], which makes the algorithm applicable to string transducers. The first extension to dwta is found in [13]. This learner generalizes the learner in [12] from bottom-up tree automata to ‘all-accepting’ dwta, and was in turn extended to general dwta in [19]. In both cases, the weight structure is a semifield (i.e., a semiring, in which multiplicative inverses exist). All those algorithms learn deterministic devices. In [17], the first MAT learner for nondeterministic weighted tree automata over fields was presented. An overview of MAT learners for weighted and unweighted tree automata can be found in [10].

Typically, MAT learnability results are based on MYHILL-NERODE-like characterizations of recognizability. In essence, the learner learns (representatives of) equivalence classes, which it refines in the light of the counterexamples provided by the teacher. This process terminates when the equivalence has converged to the MYHILL-NERODE equivalence. The learners in [12, 13, 19] are all very similar. Following Angluin’s original approach, they maintain an observation table, whose rows are indexed by trees representing the states and transitions of an automaton. The columns are indexed by contexts (trees with a “hole”) whose purpose is to separate the discovered equivalence classes from each other. From the table, the learner repeatedly constructs a dwta consistent with the information in the table and asks an equivalence query. If the MAT accepts the dwta, then the learning process has converged. Otherwise, the counterexample received is inspected by a technique known as contradiction backtracking [22]. This will either reveal a tree not yet in the table, that represents a new transition, or a context that separates previously equivalent trees in the table, thus making the learner discover a new state.

In [18, Chapter 8], Kearns and Vazirani give a rough description of Angluin’s original learner (for regular string languages), in which the observation table has been replaced by a tree-like data structure. As we shall see in Section 5, a generalized version of these observation trees can be used even in the case of dwta, thus yielding an alternative version of the learner in [19]. In Section 6, we shall see that this learner often asks fewer coefficient queries, and thus promises advantages in practical applications.

Looking at all these different but closely related learners and their correctness proofs, one cannot but realize that the same algorithms and arguments are repeated over and over again, with only slight differences. Thus, rather than

cloning [19] to come up with a version using observation trees, we propose an abstract data type, called abstract observation table (aot), in Section 3 of this paper. This captures in a formally precise way the abstract properties of observation tables that are needed to obtain a correct MAT learner. In this way, we hope to be able to contribute to a better understanding of the formal basis of MAT learning and the degrees of freedom one has when implementing it. To obtain a concrete realization of abstract observation tables, one needs to implement an interface consisting of a few abstract routines that have to satisfy certain conditions. The learning algorithm itself is hidden within the aot, and its correctness is established once and for all (depending on the correct implementation of the interface, of course). In this way, the common parts of the correctness proofs of different realizations of the learner are encapsulated. What remains to be considered in each individual case is what is specific to the particular realization at hand.

We provide two such realizations. Firstly, the learner in [19] can easily be seen to be an instance of the aot. Secondly, in Section 5, we instantiate aots as observation trees and give a correctness proof (which becomes easy thanks to the results of Section 3). As mentioned above, we expect observation trees to have the advantage of requiring fewer coefficient queries than observation tables, even though the theoretical best and worst cases coincide. To confirm this, we have implemented both learners and have conducted a set of experiments whose results are reported in Section 6. Finally, we show in Section 6 how to optimize the equivalence test derived from [6, Corollary 5.6] to run in linear time in the product of the number of transitions of the two dwta. The straightforward algorithm derived from [6, Corollary 5.6] runs in linear time in the product of the size of the transition tables, which could be exponentially larger than the number of actual transitions.

In summary, the structure of this paper is the following. In the next section, the necessary preliminaries around trees, tree series, and MAT learning are gathered. In Section 3, we present the aot and prove its correctness. This is the first major contribution of this paper. In Section 4, we show briefly that the learner in [19] is an instance of the aot. Section 5 introduces the instance based on observation trees. Together with Section 6, which confirms that observation trees often use fewer coefficient queries than observation tables, this is the second major contribution of the paper. In addition, Section 6 presents the algorithm for deciding the equivalence of dwta over semifields that we have used in our experiments for the purpose of implementing the teacher. This represents the final major contribution.

2 Preliminaries

The set of natural numbers (including 0) is denoted by \mathbb{N} . For every $k \in \mathbb{N}$ the set $\{1, \dots, k\}$ is denoted by $[k]$. For a set S , the set of all finite sequences (or strings) over S , including the empty sequence ϵ , is denoted by S^* . The cardinality of S and the length of $w \in S^*$ are denoted by $|S|$ and $|w|$, respectively.

The *index* of an equivalence relation \equiv on a set A is the number of equivalence classes induced by \equiv ; i.e., the cardinality of the quotient set A/\equiv .

Given a function $f: A \rightarrow B$ and pairs $(a_1, b_1), \dots, (a_n, b_n)$, where the a_i are pairwise distinct, we let $f\langle a_1 := b_1, \dots, a_n := b_n \rangle$ denote the function $g: A' \rightarrow B'$ with $A' = A \cup \{a_1, \dots, a_n\}$ and $B' = B \cup \{b_1, \dots, b_n\}$ such that $g(a_i) = b_i$ for $i \in [n]$, and $g(a) = f(a)$ for all $a \in A \setminus \{a_1, \dots, a_n\}$. If f is a function with empty domain, then we also write $\langle a_1 := b_1, \dots, a_n := b_n \rangle$ instead of $f\langle a_1 := b_1, \dots, a_n := b_n \rangle$.

A (commutative) semiring $\mathbf{S} = \langle \mathbb{S}, +, \cdot, 0, 1 \rangle$ consists of a set \mathbb{S} together with binary addition and multiplication operations $+$ and \cdot , respectively, as well as distinct elements $0, 1 \in \mathbb{S}$ such that (i) $\langle \mathbb{S}, +, 0 \rangle$ and $\langle \mathbb{S}, \cdot, 1 \rangle$ are commutative monoids, (ii) multiplication distributes over addition, and (iii) 0 is absorbing with respect to multiplication. The product $a_1 \cdot \dots \cdot a_k$ of finitely many elements $a_1, \dots, a_k \in \mathbb{S}$ is denoted by $\prod_{i=1}^k a_i$ or equivalently $\prod_{i \in [k]} a_i$ (note that the order is irrelevant since we consider only semirings with commutative multiplication). We say that \mathbf{S} is a semifield if every element $a \in \mathbb{S} \setminus \{0\}$ has a multiplicative inverse, which is denoted by a^{-1} (i.e., $a \cdot a^{-1} = 1$). Throughout the rest of this paper, we let \mathbf{S} be a semifield, which we simply denote by its domain \mathbb{S} . In fact, since we will only consider deterministic weighted tree automata, we will not need the addition of \mathbb{S} . For sequences $\alpha, \beta \in \mathbb{S}^*$, we write $\alpha \sim \beta$ if β is a multiple of α ; i.e., if there exists $a \in \mathbb{S} \setminus \{0\}$ such that $a \cdot \alpha = \beta$ (where multiplication is extended to a function on $\mathbb{S} \times \mathbb{S}^* \rightarrow \mathbb{S}^*$ in the obvious way). Note that \sim is an equivalence relation due to the fact that \mathbb{S} has multiplicative inverses.

Alphabets, Trees, and Contexts

A *ranked set* is a set $\Sigma = \bigcup_{k \in \mathbb{N}} \Sigma_k$ of labels (also called symbols) consisting of (not necessarily disjoint) subsets Σ_k . The labels in Σ_k are said to have rank k . We write $f^{(k)}$ to indicate that $f \in \Sigma_k$. The set T_Σ of all trees over Σ consists of all mappings $t: N \rightarrow \Sigma$ (called trees) with the following properties:

- The set N of *nodes* of t is a finite and non-empty prefix-closed subset of \mathbb{N}^* . Thus, for every $vw \in N$ with $v, w \in \mathbb{N}^*$, it holds that $v \in N$.
- For all $v \in N$, there is $k \in \mathbb{N}$ such that $t(v) \in \Sigma_k$ and $\{i \in \mathbb{N} \mid vi \in N\} = [k]$.

For a tree $t: N \rightarrow \Sigma$, we also write $\text{nodes}(t)$ for N . The *height* of t is denoted by $\text{hg}(t)$; i.e., $\text{hg}(t) = \max\{|v| \mid v \in \text{nodes}(t)\}$. Given a tree t and a node $v \in \text{nodes}(t)$, the *subtree of t rooted at v* is denoted by t/v . It is defined by $\text{nodes}(t/v) = \{w \in \mathbb{N}^* \mid vw \in \text{nodes}(t)\}$ and $(t/v)(w) = t(vw)$ for all $w \in \text{nodes}(t/v)$. The set $\text{leaves}(t)$ of all *leaves of t* is the set

$$\text{leaves}(t) = \{v \in \text{nodes}(t) \mid v1 \notin N\} .$$

We shall denote a tree t as $f[t_1, \dots, t_k]$ if $t(\epsilon) = f^{(k)}$ and $t/i = t_i$ for all $i \in [k]$. In the special case where $k = 0$ (i.e., $\text{nodes}(t) = \{\epsilon\}$), the brackets may be omitted, thus denoting t as f . For a set T of trees, the set of all trees of the form $f[t_1, \dots, t_k]$ such that $k \in \mathbb{N}$, $f \in \Sigma_k$, and $t_1, \dots, t_k \in T$ is denoted by $\Sigma(T)$.

We will frequently decompose a tree into a context and a subtree. For this purpose, let us reserve the special symbol \square of rank 0 (and no other rank) that does not occur in Σ . A tree $c \in \mathsf{T}_{\Sigma \cup \{\square\}}$ in which \square occurs exactly once is called a *context* (over Σ). The set of all contexts over Σ is denoted by C_Σ . Given a context $c \in \mathsf{C}_\Sigma$ and a tree t , we denote by ct the concatenation of c and t at \square , which is obtained by substituting t for the unique leaf labelled \square in c . More precisely, if $v \in \text{nodes}(c)$ is the unique node such that $c(v) = \square$, then $\text{nodes}(ct) = \text{nodes}(c) \cup \{vw \mid w \in \text{nodes}(t)\}$ with $ct(w) = c(w)$ for all $w \in \text{nodes}(c) \setminus \{v\}$, and $ct/v = t$.

For every $k \in \mathbb{N}$, $t, t_1, \dots, t_k \in \mathsf{T}_\Sigma$, and $v_1, \dots, v_k \in \text{nodes}(t)$ such that v_i is not a prefix of v_j for any $i, j \in [k]$, we denote by $t[v_1 \leftarrow t_1, \dots, v_k \leftarrow t_k]$ the tree obtained by replacing t/v_i by t_i for every $i \in [k]$.

Tree Series and Weighted Tree Automata

A *tree series* over \mathbb{S} is a mapping $\psi: \mathsf{T}_\Sigma \rightarrow \mathbb{S}$. Its support is

$$\text{supp}(\psi) = \{t \in \mathsf{T}_\Sigma \mid \psi(t) \neq 0\} .$$

A tree $t \in \mathsf{T}_\Sigma$ is called *live* (in ψ) if there exists $c \in \mathsf{C}_\Sigma$ such that $ct \in \text{supp}(\psi)$. Such a context c is also called a *sign of life* for t (in ψ).

We now recall the definition of deterministic weighted (bottom-up finite-state) tree automata [4, 5, 7] (dwta, for short) over \mathbb{S} . Such a dwta is a tuple $A = (\Sigma, Q, \delta, \lambda)$ where

- Σ is the finite input ranked set,
- Q is the finite set of *states*, considered as symbols of rank 0 and no other rank,
- $\delta: \Sigma(Q) \rightarrow Q \times (\mathbb{S} \setminus \{0\})$ is the partial *transition function*, and
- $\lambda: Q \rightarrow \mathbb{S}$ is the *final weight mapping*.

Intuitively, $\delta(f[q_1, \dots, q_k]) = (q, a)$ determines the behaviour of the dwta in the situation, in which it processes an occurrence of the symbol f of rank k and the k subtrees t_1, \dots, t_k of f have already been processed in states q_1, \dots, q_k , respectively. Then the dwta continues in state q and charges the weight a .

The transition function δ extends to trees in a straightforward way, which yields the partial function $\hat{\delta}: \mathsf{T}_{\Sigma \cup Q} \rightarrow Q \times (\mathbb{S} \setminus \{0\})$ such that $\hat{\delta}(q) = (q, 1)$ for every $q \in Q$ and, for every $f[t_1, \dots, t_k]$ with $f \in \Sigma_k$ and $t_1, \dots, t_k \in \mathsf{T}_{\Sigma \cup Q}$:

- If $\hat{\delta}(t_i) = (q_i, a_i)$ for every $i \in [k]$ and $\delta(f[q_1, \dots, q_k]) = (q, a)$ are all defined, then $\hat{\delta}(f[t_1, \dots, t_k]) = (q, a \cdot \prod_{i=1}^k a_i)$.
- Otherwise, $\hat{\delta}(f[t_1, \dots, t_k])$ is undefined.

The tree series computed by A is given as follows. For all $t \in \mathsf{T}_\Sigma$, if $\hat{\delta}(t) = (q, a)$, then $A(t) = a \cdot \lambda(q)$. Otherwise, $A(t) = 0$.

Example 1. The following dwta will be used as a running example throughout the paper. Consider the *tropical semifield* $\mathbb{T} = (\mathbb{Z} \cup \{-\infty\}, \min, +, -\infty, 0)$, where ‘min’ and ‘+’ are ordinary minimum and addition, resp., on \mathbb{Z} extended by $-\infty$. Note that ‘+’ is the multiplicative operation in this semifield. Moreover, as mentioned above, the additive operation ‘min’ will not be needed, so we are actually using the group $(\mathbb{Z}, +, 0)$ enriched by an absorbing “error element” $-\infty$.

Let $\Delta = \{f^{(2)}, g^{(1)}, e^{(0)}\}$. To save space in later discussions of this example, we shall denote trees and contexts without brackets. Thus, for instance, $f[e, g[g[e]]]$ will simply be denoted as $f e g g e$.

We define the tree series $\chi: \mathbb{T}_\Delta \rightarrow \mathbb{T}$ as follows. For every $t \in \mathbb{T}_\Delta$, we have $t \in \text{supp}(\chi)$ [i.e., $(\chi, t) \neq -\infty$] if and only if

- (a) $t \neq e$,
- (b) no node labeled g has an f -labeled node as its child, and
- (c) the children of each f -labeled node have equal root labels.

On all trees t that fulfill these conditions, $\chi(t)$ is equal to the number of nodes labeled f whose first subtree equals $g e$.

The states in a dwta $A_{\text{count}} = (\Delta, Q, \delta, \lambda)$ recognizing χ will have to distinguish between the tree e , the tree $g e$, the remaining monadic trees, and the trees having the root f . In other words, we need four states q_e, q_g, q_m, q_f . The transition function reflects the meaning of these states. It charges a weight of 0 except when an f -labeled node is processed whose first subtree is $g e$, in which case the weight 1 is charged:

$$\begin{aligned} \delta(e) &= (q_e, 0) \\ \delta(g q_e) &= (q_g, 0) & \delta(f q_e q_e) &= (q_f, 0) & \delta(f q_g q_m) &= (q_f, 1) \\ \delta(g q_g) &= (q_m, 0) & \delta(f q_g q_g) &= (q_f, 1) & \delta(f q_m q_g) &= (q_f, 0) \\ \delta(g q_m) &= (q_m, 0) & \delta(f q_f q_f) &= (q_f, 0) & \delta(f q_m q_m) &= (q_f, 0) . \end{aligned}$$

The final weight of the states q_g, q_m , and q_f is 0, while $\lambda(q_e) = -\infty$ to make sure that $A_{\text{count}}(e) = -\infty$. The dwta is illustrated in Figure 1. Note also that $\hat{\delta}$ is undefined on all trees violating condition (b) or (c). Hence, only trees satisfying all three conditions are in the support. *To be continued.*

Every dwta can be made *total* without affecting the tree series computed, where a dwta A is total if δ is a total function. This is achieved by adding a non-final (i.e., final weight 0) sink state that is the target of all missing transitions (using, for example, weight 1). For the dwta of Example 1, the corresponding total dwta would have $1 + 5 + 5^2 = 31$ instead of 10 transitions. In general, exponentially many (in the number of states or existing transitions) additional transitions might be required. Thus, making a dwta total can yield an exponential blow-up.

Throughout the remainder of this paper, let Σ be a finite ranked set, $\psi: \mathbb{T}_\Sigma \rightarrow \mathbb{S}$ be a tree series, and $\mathbf{S} = \langle \mathbb{S}, +, \cdot, 0, 1 \rangle$ be a commutative semifield.

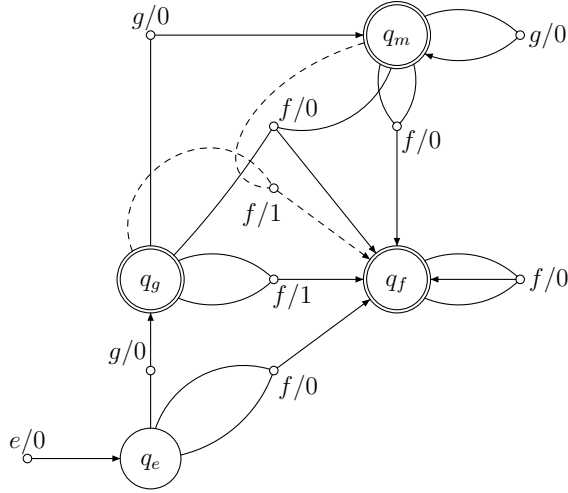


Fig. 1. Illustration of the dwta of Example 1. Lines (unbroken or dashed) to a transition node are ordered counter-clockwise from the arrow leading to the target state. The dashed lines are only used to increase readability and carry no special significance.

Let $s, t \in T_\Sigma$. For $C \subseteq C_\Sigma$, we write $s \equiv_C t$ if there exists $a \in \mathbb{S} \setminus \{0\}$ such that $\psi(cs) = a \cdot \psi(ct)$ for every $c \in C$. Note that \equiv_C is an equivalence relation on T_Σ . The relation \equiv_{C_Σ} is simply denoted by \equiv .

Theorem 2 (see [5, Theorems 2 and 3]). *The tree series ψ can be computed by some dwta if and only if \equiv has finite index. Moreover, any dwta computing ψ has at least $|L/\equiv|$ states where $L = \{t \in T_\Sigma \mid t \text{ live in } \psi\}$.*

An easy but important observation for the learning algorithms considered in this paper concerns the question how one can establish that $t \not\equiv t'$ for trees $t, t' \in T_\Sigma$ having a common sign of life c . Let $\psi(ct') = a \cdot \psi(ct)$. By definition, $t \equiv t'$ is equivalent to saying that $\psi(c't') = a \cdot \psi(c't)$ for all $c' \in C_\Sigma$. In other words, we have the following lemma, that we will usually make use of without explicitly mentioning this fact.

Lemma 3. *Let $t, t' \in T_\Sigma$ have a common sign of life c . Then $t \equiv t'$ if and only if, for all $c' \in C_\Sigma$,*

$$\frac{\psi(c't)}{\psi(ct)} = \frac{\psi(c't')}{\psi(ct')} .$$

Example 4 (cont'd). In our running example, one of the (infinitely many) common signs of life for the trees e and ge is the context $c = g\Box$. Taking $c' = \Box$, we find that $\chi(c'e) - \chi(ce) = -\infty - 0 \neq 0 - 0 = \chi(c'ge) - \chi(cge)$, which shows that $e \not\equiv ge$. *To be continued.*

The Minimal Adequate Teacher

In the following, we will consider grammatical inference of ψ . The aim is to build a dwta computing ψ (if such a dwta exists) using an appropriately extended version of Angluin’s *minimal adequate teacher* (MAT) as the source of information about ψ (cf. [1, 2, 3]).

A learning algorithm that infers a dwta computing ψ will henceforth be called a *learner*. Such a learner may ask *coefficient* and *equivalence queries* to the MAT (and the MAT will answer them correctly):

Coefficient: Given a tree $t \in T_\Sigma$ (provided by the learner), what is $\psi(t)$?

Equivalence: Given a dwta A (provided by the learner), does A compute ψ ?

If so, the teacher returns the token \perp indicating that $A = \psi$. Otherwise, a counterexample is returned; i.e., a tree $t \in T_\Sigma$ such that $A(t) \neq \psi(t)$.

The algorithms presented in this paper are all supposed to have access to a MAT. In particular, this implies that these algorithms can obtain $\psi(t)$ for $t \in T_\Sigma$ by asking a coefficient query. Thus, the reader should bear in mind that every mention of $\psi(t)$ in our algorithms indicates that a coefficient query is asked unless $\psi(t)$ has already been provided earlier, in which case a reasonable implementation would have memorized the value.

3 An abstract data type for MAT learners

In this section, we will develop an abstract data type specification, called *abstract observation table*, for MAT learners [1, 2, 3]. The commonly used ‘observation table’ [1, 12, 13, 19] will be an instance of this specification, but in the next section we will present another data type, called *observation tree*, that avoids (in our experiments) some of the coefficient queries typically asked when the learner fills the observation table.

In essence, our abstract data type manages the two sets S and T of trees, which are related by $S \subseteq T \subseteq \Sigma(S)$. Eventually, S and T will respectively correspond to the states and the transitions of the deterministic weighted tree automaton that we construct. Moreover, our abstract data type maintains two mappings, $\text{sol}: S \rightarrow C_\Sigma$ and $\rho: T \rightarrow S$, which assign a sign of life (in ψ) for each tree in S and a tree of S to each tree $t \in T$, respectively. This tree $\rho(t)$ is called the representative of t and eventually encodes the transition target $\delta(t) = \rho(t)$, where δ is the transition function of the constructed dwta. The name ‘representative’ is justified by the fact that only $\rho(t)$ of all trees in S can be \equiv -equivalent to t .

To simplify the notation (and to avoid some parentheses), we write $t\rho$ instead of $\rho(t)$ for every $t \in T$. Let us present a mathematical definition of the universal invariant of an abstract observation table.

Definition 5 (cf. [19, Definition 8]). *Let $S \subseteq T \subseteq \Sigma(S)$ be finite subsets of T_Σ , $\text{sol}: S \rightarrow C_\Sigma$, and $\rho: T \rightarrow S$. Then (S, T, sol, ρ) is an abstract observation table (with respect to ψ) if*

1. $\text{sol}(t\rho) = \square$ for every $t \in T \cap \text{supp}(\psi)$ *(trivial sign of life for support)*
2. $\text{sol}(t\rho)t \in \text{supp}(\psi)$ for every $t \in T$ *(sol($t\rho$) sign of life for t in ψ)*
3. $t \not\equiv s$ for every $t \in T$ and $s \in S \setminus \{t\rho\}$ *(t distinct from other states)*

In the following, we will abbreviate $\text{sol}(t\rho)$ by $\text{sol}(t)$ for $t \in T_\Sigma$ whenever appropriate. Thus, with this convention in mind, the first condition becomes $\text{sol}(t) = \square$ for every $t \in T \cap \text{supp}(\psi)$, and the second becomes $\text{sol}(t)t \in \text{supp}(\psi)$ for every $t \in T$.

Note that the third condition yields $s\rho = s$ for every $s \in S$ because trivially $s \equiv s$. In addition, this implies that $s_1 \not\equiv s_2$ for all distinct $s_1, s_2 \in S$. Finally, for every $t \in T$, $t \in \text{supp}(\psi)$ if and only if $t\rho \in \text{supp}(\psi)$, which can be seen as follows. If $t \in \text{supp}(\psi)$ or $t\rho \in \text{supp}(\psi)$, then $\text{sol}(t) = \square$ by the first condition and thus $t \in \text{supp}(\psi)$ and $t\rho \in \text{supp}(\psi)$ by the second condition.

Conceptually, we note that all known variants of Angluin's original algorithm (at least those we are aware of) maintain a set of contexts for the purpose of distinguishing between the equivalence classes of trees in T . In our abstract version, the only contexts that are explicitly required to be maintained are the signs of life for the trees in T . Thus, it might, in principle, be possible to implement abstract observation tables without managing additional contexts. The difficulty is, of course, to make sure that the third condition is satisfied, because this is usually what separating contexts are used for (cf. the definition of \equiv).

In the following, we will use the following interface to manipulate an abstract observation table. To simplify the description of the semantic properties, let (S, T, sol, ρ) and $(S', T', \text{sol}', \rho')$ be the abstract observation table before and after execution, respectively.

- INITIALIZE *(constructor)*
 - Post-condition: $S' = T' = \emptyset$
- ADDTRANSITION(t, c) with $t \in \Sigma(S)$ and $c \in C_\Sigma$ *(add new transition)*
 - Pre-conditions: $t \notin T$ and $ct \in \text{supp}(\psi)$
 - Post-conditions: $S \subseteq S'$ and $T \cup \{t\} \subseteq T'$
- ADDSTATE(t, c) with $t \in T$ and $c \in C_\Sigma$ *(add new state)*
 - Pre-condition: $\frac{\psi(ct)}{\psi(\text{sol}(t)t)} \neq \frac{\psi(c(t\rho))}{\psi(\text{sol}(t)(t\rho))}$
 - Post-conditions: $S \cup \{t\} \subseteq S'$ and $T \subseteq T'$

In the next definition, it is shown how to construct a dwta $A_{\text{aot}} = (\Sigma, S, \delta, \lambda)$ from an abstract observation table (S, T, sol, ρ) . The motivation for the pre-conditions of ADDTRANSITION and ADDSTATE are closely related to this definition. As we shall see, the trees in S will be turned into the states of A_{aot} , whereas the trees in T will give rise to its transitions. As a consequence, $\delta(t)$ and $\hat{\delta}(t)$ are undefined for trees $t \in \Sigma(S) \setminus T$. In other words, a tree t that fulfills the pre-conditions of ADDTRANSITION is not live in A_{aot} , whereas the second part of the pre-condition states that it, in fact, is live in ψ . Thus, the tree t must be added (as a transition) to T . The pre-condition of ADDSTATE is motivated by the fact that A_{aot} maps t and $t\rho$ to the same state (namely $t\rho$), while the pre-condition expresses that $\text{sol}(t)$ and c separate $t\rho$ from t (by Lemma 3). Hence, t should rather be mapped to a new state.

Definition 6 (cf. [5, Definition 4 and p. 9]). Let $\text{aot} = (S, T, \text{sol}, \rho)$ be an abstract observation table. Let $\psi_{\text{aot}}: T \rightarrow \mathbb{S}$ be such that for every $t \in T$

$$\psi_{\text{aot}}(t) = \frac{\psi(\text{sol}(t)t)}{\psi(\text{sol}(t)(t\rho))} .$$

We construct the deterministic weighted tree automaton $A_{\text{aot}} = (\Sigma, S, \delta, \lambda)$ with

- $\lambda(s) = \psi(s)$ for every $s \in S$,
- $\delta(t) = (t\rho, \psi_{\text{aot}}(t))$ for every $t \in T$, and
- $\delta(t)$ is undefined for all $t \in \Sigma(S) \setminus T$.

Example 7 (cont'd). Let $S = \{e, ge\}$ and $T = \{e, ge, fee, fgege\}$, where

$s \in S$	$\text{sol}(s)$	$s\rho$	$t \in T \setminus S$	$t\rho$
e	$g\Box$	e	fee	ge
ge	\Box	ge	$fgege$	ge

Then $\text{aot} = (S, T, \text{sol}, \rho)$ is an abstract observation table. The dwta A_{aot} has the states e and ge , where

$$\delta(e) = (e, 0) \quad \delta(ge) = (ge, 0) \quad \delta(fee) = (ge, 0) \quad \delta(fgege) = (ge, 1) .$$

The final weights are $\lambda(e) = -\infty$ and $\lambda(ge) = 0$. Figure 2 illustrates A_{aot} . This dwta still computes some wrong weights. For example, it computes weight 2 for $ffgegege$, but it computes the correct weight for all trees $t \in T$ (namely 1 if $t = fgege$ and 0 otherwise). As we shall see next, this is not a coincidence.

To be continued.

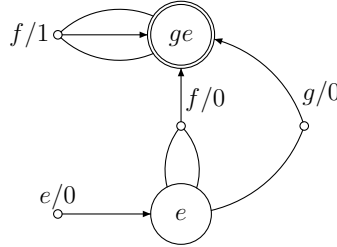


Fig. 2. Illustration of the dwta of Example 7. As before, lines to a transition node are ordered counter-clockwise from the arrow leading to the target state.

We observe some easy properties of the automaton A_{aot} . First, note that $\psi_{\text{aot}}(s) = 1$ for every $s \in S$ (as $s\rho = s$). Second, A_{aot} computes ψ on all trees of T , which we prove in the next lemma.

Lemma 8 (see [19, Lemma 12]). Let $\text{aot} = (S, T, \text{sol}, \rho)$ be an abstract observation table. Then $A_{\text{aot}}(t) = \psi(t)$ for every $t \in T$.

Algorithm 1 Learn a minimal deterministic weighted tree automaton for ψ

Post-conditions: returned automaton computes ψ

```

aot.INITIALIZE // initialize data structure
2: loop
     $t \leftarrow \text{EQUAL?}(A_{\text{aot}})$  // ask equivalence query
4: if  $t = \perp$  then
    return  $A_{\text{aot}}$  // return the approved automaton
6:  $\text{aot} \leftarrow \text{EXTEND}(\text{aot}, t)$  // extend the data structure

```

Proof. Let $A_{\text{aot}} = (\Sigma, S, \delta, \lambda)$. First, we claim that $\hat{\delta}(t) = (t\rho, \psi_{\text{aot}}(t))$ for every $t \in T$ by induction on t . Since $S \subseteq T \subseteq \Sigma(S)$, this induction is possible and we have $t = f[s_1, \dots, s_k]$ for some $k \in \mathbb{N}$, $f \in \Sigma_k$, and $s_1, \dots, s_k \in S$. The induction base corresponds to the case $k = 0$ and is contained in the induction step. The induction hypothesis is $\hat{\delta}(s_i) = (s_i, 1)$ because $s_i\rho = s_i$ (since $s_i \in S$) and $\psi_{\text{aot}}(s_i) = 1$ for every $i \in [k]$. Since $\delta(t) = (t\rho, \psi_{\text{aot}}(t))$, we obtain

$$\hat{\delta}(t) = \left(t\rho, \psi_{\text{aot}}(t) \cdot \prod_{i=1}^k 1 \right) = (t\rho, \psi_{\text{aot}}(t)) .$$

This proves the claim. Now we can prove the statement as follows.

$$\begin{aligned} A_{\text{aot}}(t) &= \psi_{\text{aot}}(t) \cdot \lambda(t\rho) = \frac{\psi(\text{sol}(t)t)}{\psi(\text{sol}(t)(t\rho))} \cdot \psi(t\rho) \\ &= \begin{cases} \frac{\psi(t)}{\psi(t\rho)} \cdot \psi(t\rho) & \text{if } t \in \text{supp}(\psi) \\ 0 & \text{otherwise} \end{cases} \\ &= \psi(t) , \end{aligned}$$

where the penultimate equality uses the first condition of Definition 5. \square

The principal structure of the MAT learner [12, 13, 19] is shown in Algorithm 1. Note that we only adapted it to work with our abstract observation table. We start with the initial empty data structure ‘aot’ and iteratively query the teacher for counterexamples to our current hypothesis (the current deterministic weighted tree automaton A_{aot}), which is constructed from the current abstract observation table (see Definition 6). We update our abstract data structure with the returned information (using EXTEND), and if the teacher eventually approves our dwta, then we simply return it.

Example 9 (cont’d). Let us consider a run of Algorithm 1 using an unknown but correct implementation of abstract observation tables. After initialization, an equivalence query is asked, passing the empty dwta as a parameter. The teacher may return the counterexample $f\text{ge}g\text{ge}$. In EXTEND this yields a call to either ADDTRANSITION($e, f\text{g}\square\text{g}\text{ge}$) or ADDTRANSITION($e, f\text{ge}g\text{g}\square$). Thus, we may obtain the new abstract observation table

$$(\{\mathbf{e}\}, f\text{ge}g\text{g}\square) ,$$

where we denote an abstract observation table as the set of pairs $(\rho^{-1}(s), \text{sol}(s))$ with $s \in S$ and the element $s \in \rho^{-1}(s)$ is printed in boldface. In addition, we drop the outermost set braces.

Let us assume that the teacher returns the same counterexample as long as it is not treated correctly.³ Then three more calls of `ADDTRANSITION` lead to the addition of ge , gge , and $fgege$ to the table, all of which have the sign of life \square . Hence, the table may become

$$(\{\mathbf{e}\}, fgegg\square) \quad , \quad (\{\mathbf{ge}, gge, fgege\}, \square) \quad .$$

Note that there are some other possible outcomes as well, such as the table $\{(\{\mathbf{e}\}, fgegg\square), (\{\mathbf{ge}, gge\}, \square), (\{\mathbf{fgege}\}, \square)\}$. The result depends on the particular implementation of the used abstract observation table.

Since a transition for processing $f ee$ is missing, the teacher may return a counterexample containing this subtree, which may lead to the table

$$(\{\mathbf{e}\}, fgegg\square) \quad , \quad (\{\mathbf{ge}, gge, f ee, fgege\}, \square) \quad .$$

However, the resulting dwta assigns the weight 0 to $gf ee$, so this tree, which is outside $\text{supp}(\chi)$, may be the next counterexample. This results in the first call of `ADDSTATE`, namely `ADDSTATE(f ee, g\square)`, where $g\square$ is the context that separates $f ee$ from ge . As a consequence, the table may become

$$(\{\mathbf{e}\}, fgegg\square) \quad , \quad (\{\mathbf{ge}, gge\}, \square) \quad , \quad (\{\mathbf{f ee}, fgege\}, \square) \quad .$$

The next problem is the counterexample $f g g e ge$, to which the dwta assigns the weight 1 rather than 0. Processing this counterexample leads to the call `ADDSTATE(gge, f\square ge)` and subsequently to a call of `ADDTRANSITION`. The result is

$$(\{\mathbf{e}\}, fgegg\square) \quad , \quad (\{\mathbf{ge}\}, \square) \quad , \quad (\{\mathbf{gge}, ggge\}, \square) \quad , \quad (\{\mathbf{f ee}, fgege\}, \square) \quad .$$

It remains to discover the four missing transitions, all of which are revealed by the counterexample $f f f g g e ge f g g e g e f ge g e$. Repeatedly processing this counterexample leads to the final table

$$\begin{aligned} &(\{\mathbf{e}\}, fgegg\square) \quad , \quad (\{\mathbf{ge}\}, \square) \quad , \quad (\{\mathbf{gge}, ggge\}, \square) \quad , \\ &(\{\mathbf{f ee}, fgege, f g g e ge, f g g e g e, f ge g e, f f e e f e e\}, \square) \quad . \end{aligned}$$

To be continued.

We say that an algorithm works *correctly* if whenever the pre-conditions are met at the beginning of the algorithm, then (i) the algorithm terminates and (ii) the post-conditions hold at the point of return. For the next statements, we additionally assume that a correct implementation of our abstract observation table is used.

³ An optimized learner can check this without asking equivalence queries; see [11].

Theorem 10 (see [19, Theorem 13]). *If EXTEND works correctly (see the pre- and post-conditions given in Algorithm 2) and ψ can be computed by some dwta, then Algorithm 1 terminates and returns a minimal dwta computing ψ .*

Proof. Suppose that ψ can be computed by some dwta. Then \equiv has finite index by Theorem 2. Let $n = |\mathbb{T}_\Sigma/\equiv|$. Clearly, by the third condition in Definition 5, the set S of an abstract observation table (S, T, sol, ρ) contains at most n elements. Trivially, EXTEND is always called with a counterexample, because the counterexample is provided by the teacher. Since $|S|$ and $|T|$ are uniformly bounded and each call to EXTEND increases $|S| + |T|$, there can only be finitely many calls to EXTEND, which yields that Algorithm 1 terminates. Moreover, the returned dwta A_{aot} was approved by the teacher, so A_{aot} trivially computes ψ . By the construction of A_{aot} (see Definition 6), we know that it has at most n states. Since all states (recall that they are trees) of A_{aot} are live in A_{aot} , this shows that it is a minimal dwta computing ψ by Theorem 2. \square

Finally, let us discuss the function EXTEND, which is displayed in Algorithm 2. Given the counterexample t , it searches for a minimal subtree of t that is still a counterexample, using a technique called *contradiction backtracking* [22]. Let $\text{aot} = (S, T, \text{sol}, \rho)$ be the abstract observation table and $t \in \mathbb{T}_\Sigma$ be the counterexample; i.e., a tree t such that $A_{\text{aot}}(t) \neq \psi(t)$. We decompose t into a context $c \in C_\Sigma$ and a tree u that is itself not in S but whose direct subtrees are all in S . In some sense, this is a minimal subtree that could possibly be offending, because A_{aot} computes the correct coefficient on all trees in T by Lemma 8. Moreover, such a subtree must exist, because $t \notin S$ (since t is a counterexample).

Now, we distinguish two cases. If u was already seen (i.e., $u \in T$), then by Lemma 8, A_{aot} returns $\psi(u)$ if applied to u . Thus an error is made when processing the context c . To this end, we test whether the context c separates u and $u\rho$; the latter is the state that represents u . Provided that c does not distinguish between u and $u\rho$, then we continue our search for an error with the simplified counterexample $c(u\rho)$. In the other cases, either u and $u\rho$ could be separated or u was not seen before. Consequently, we either add u as a new state (in the former case) or as a new transition (in the latter case). Overall, the post-condition of the algorithm is trivially met.

It is clear that the pre-conditions of aot.ADDSTATE and aot.ADDTRANSITION are met as well. It remains to prove that the recursive call of EXTEND meets the pre-conditions of EXTEND. To this end, we need to prove that $c(u\rho)$ is also a counterexample in line 6. This is shown in the next lemma.

Lemma 11 (cf. [19, Lemma 16]). *Let $\text{aot} = (S, T, \text{sol}, \rho)$ be an abstract observation table, $t \in T$, and $c \in C_\Sigma$ such that*

$$\frac{\psi(ct)}{\psi(\text{sol}(t)t)} = \frac{\psi(c(t\rho))}{\psi(\text{sol}(t)(t\rho))} . \quad (1)$$

If $A_{\text{aot}}(ct) \neq \psi(ct)$, then also $A_{\text{aot}}(c(t\rho)) \neq \psi(c(t\rho))$.

Algorithm 2 Function EXTEND(t) for aot = (S, T, sol, ρ)

Pre-conditions: $t \in T_\Sigma$ with $A_{\text{aot}}(t) \neq \psi(t)$

Post-conditions: return an abstract observation table aot' = $(S', T', \text{sol}', \rho')$ such that $S \subseteq S'$ and $T \subseteq T'$ and one inclusion is strict

Decompose t into $t = cu$ where $c \in C_\Sigma$ and $u \in \Sigma(S) \setminus S$

2: **if** $u \notin T$ **then**
 return aot.ADDTRANSITION(u, c) // u not reachable so far; add transition

4: **if** $\frac{\psi(cu)}{\psi(\text{sol}(u)u)} \neq \frac{\psi(c(u\rho))}{\psi(\text{sol}(u)(u\rho))}$ **then**
 return aot.ADDSTATE(u, c) // add new state u

6: **return** EXTEND(aot, $c(u\rho)$) // normalize and continue

Proof. Let $A_{\text{aot}} = (\Sigma, S, \delta, \lambda)$. By the claim in the proof of Lemma 8 it follows that $\hat{\delta}(t) = (t\rho, \psi_{\text{aot}}(t))$ and $\hat{\delta}(t\rho) = (t\rho, 1)$ because $t \in T$. Trivially, the former yields that $\hat{\delta}(ct)$ is defined if and only if $\hat{\delta}(c(t\rho))$ is defined. Moreover, if they are defined, then there exist $s \in S$ and $a \in \mathbb{S} \setminus \{0\}$ such that $\hat{\delta}(ct) = (s, a \cdot \psi_{\text{aot}}(t))$ and $\hat{\delta}(c(t\rho)) = (s, a)$. Now we distinguish three cases:

- First, let $ct \notin \text{supp}(A_{\text{aot}})$. Then clearly also $c(t\rho) \notin \text{supp}(A_{\text{aot}})$, because $A_{\text{aot}}(t) = \psi(t) \neq 0$. Since ct is a counterexample, we have $ct \in \text{supp}(\psi)$ and thus also $c(t\rho) \in \text{supp}(\psi)$ by (1), which proves that $A_{\text{aot}}(c(t\rho)) \neq \psi(c(t\rho))$.
- Second, let $ct \notin \text{supp}(\psi)$. Since ct is a counterexample by assumption, we obtain that $ct \in \text{supp}(A_{\text{aot}})$ and thus also $c(t\rho) \in \text{supp}(A_{\text{aot}})$. Moreover, equation (1) yields that $c(t\rho) \notin \text{supp}(\psi)$, which again proves that $A_{\text{aot}}(c(t\rho)) \neq \psi(c(t\rho))$.
- Third, let $ct \in \text{supp}(A_{\text{aot}}) \cap \text{supp}(\psi)$. By the same reasoning as in the previous cases, this yields that $c(t\rho) \in \text{supp}(A_{\text{aot}}) \cap \text{supp}(\psi)$. The observation $\hat{\delta}(ct) = (s, a \cdot \psi_{\text{aot}}(t))$ stated in the first paragraph allows us to conclude that

$$\begin{aligned} A_{\text{aot}}(ct) &= a \cdot \psi_{\text{aot}}(t) \cdot \lambda(s) = \frac{\psi(\text{sol}(t)t)}{\psi(\text{sol}(t)(t\rho))} \cdot A_{\text{aot}}(c(t\rho)) \\ &= \frac{\psi(ct)}{\psi(c(t\rho))} \cdot A_{\text{aot}}(c(t\rho)) \end{aligned}$$

by (1). Since all factors are nonzero, we obtain

$$\frac{A_{\text{aot}}(ct)}{\psi(ct)} = \frac{A_{\text{aot}}(c(t\rho))}{\psi(c(t\rho))} .$$

By assumption, the left-hand side is different from 1, which proves that $A_{\text{aot}}(c(t\rho)) \neq \psi(c(t\rho))$. \square

Consequently, the recursive call of EXTEND is correct. An easy size argument (counting the subtrees of t that are not in T) can be used to show that the recursion terminates (see [13, Lemma 5.3]). Thus we obtain the main statement of this section.

Corollary 12 (of Theorem 10). *If ψ can be computed by some dwta, then Algorithm 1 terminates and returns a minimal dwta computing ψ .*

4 Observation Tables

Let us briefly present the “classical” implementation of our abstract observation table: the ‘observation table’. Several similar implementations exist; the one presented here corresponds to the one in [19].

Definition 13 (see [19, Definition 8]). *Let $T \subseteq \Sigma(T)$ and $C \supseteq \{\square\}$ be finite subsets of T_Σ and C_Σ , respectively. An observation table is a $(T \times C)$ -matrix \mathbf{P} with $\mathbf{P}(t, c) = \psi(ct)$ for every $t \in T$ and $c \in C$ such that, for every $t \in T$, there exists $c \in C$ with $\mathbf{P}(t, c) \neq 0$.*

Given a set $S \subseteq T$, the pair (S, \mathbf{P}) is an S -observation table, if

- $S \subseteq \Sigma(S)$,
- $s_1 \not\equiv_C s_2$ for all $s_1, s_2 \in S$, and (no \equiv_C -equivalent rows in S)
- for every $t \in T$ there exists $s \in S$ such that $t \equiv_C s$. (no new rows in T)

Note that $t \equiv_C t'$ for trees $t, t' \in T$ means that the row indexed by t is a multiple of the one indexed by t' (by a nonzero factor). It has essentially been shown in [12, 13, 19] that every observation table \mathbf{P} can be turned into an S -observation table by choosing an appropriate set $S \subseteq T$, and that the operations of our abstract observation table can be implemented with the help of observation tables. Let us quickly show how this works.

Lemma 14. *Observation tables implement abstract observation tables.*

Proof. Let $(S, \mathbf{P}: T \times C \rightarrow \mathbb{S})$ be an S -observation table. The *abstract observation table* $(S, T', \text{sol}, \rho)$ represented by (S, \mathbf{P}) is given follows:

- $T' = T \cap \Sigma(S)$,
- for every $s \in S$,

$$\text{sol}(s) = \begin{cases} \square & \text{if } s \in \text{supp}(\psi) \\ c & \text{otherwise, for some } c \in C \text{ such that } cs \in \text{supp}(\psi) \end{cases}$$

(by Definition 13, such an element exists for every $s \in S$), and

- $t\rho = s$ where $s \in S$ is such that $t \equiv_C s$ (by the second and third condition for S -observation trees s exists and is unique).

To verify the conditions of our abstract observation table, let $t \in T'$.

1. If $t \in \text{supp}(\psi)$ then $t\rho \in \text{supp}(\psi)$ because $t \equiv_C t\rho$. Consequently, $\text{sol}(t) = \square$.
2. By definition, $\text{sol}(t)(t\rho) \in \text{supp}(\psi)$. Again $t \equiv_C t\rho$ and since $\text{sol}(t) \in C$, we obtain $\text{sol}(t)t \in \text{supp}(\psi)$.
3. Let $s \in S \setminus \{t\rho\}$. By the definition of ρ and the first condition of Definition 13, $t \equiv_C t\rho$ and $t\rho \not\equiv_C s$. Since \equiv_C is an equivalence relation and $\equiv \subseteq \equiv_C$ (see remarks below [19, Definition 7]), this yields $t \not\equiv_C s$ and $t \neq s$.

It remains to define INITIALIZE, ADDTRANSITION, and ADDSTATE. Of course, INITIALIZE returns $(\emptyset, \mathbf{P}_\epsilon)$, where \mathbf{P}_ϵ is the empty matrix (i.e., $T = C = \emptyset$). The function ADDTRANSITION simply adds t to T and c to C (and extends \mathbf{P} by means of coefficient queries).⁴ If necessary, it *completes* S by adding elements of $T \cup \{t\}$ to it until the third condition of Definition 13 is fulfilled. Similarly, ADDSTATE adds t to S and c to C , updates \mathbf{P} , and completes S . For both ADDTRANSITION and ADDSTATE, it is straightforward to check that the resulting pair (S', \mathbf{P}') is an S' -observation table, and that the abstract observation table it represents fulfills the post-condition of the respective function. \square

Example 15 (cont'd). Re-using the first counterexample in Example 9 as long as possible may lead to the following observation table:

	\square	$f\ ge\ gg\square$	$f\ ge\ g\square$	$f\ ge\ \square$
\mathbf{e}	$-\infty$	1	1	$-\infty$
\mathbf{ge}	0	1	1	1
gge	0	1	1	1
$\mathbf{f\ ge\ ge}$	1	$-\infty$	$-\infty$	$-\infty$

The tree $f\ ge\ ge$ has been separated from the others “by chance”. Hence, the table corresponds to

$$(\{\mathbf{e}\}, f\ ge\ gg\square) \quad , \quad (\{\mathbf{ge}, gge\}, \square) \quad , \quad (\{\mathbf{f\ ge\ ge}\}, \square)$$

rather than to $\{(\{\mathbf{e}\}, f\ ge\ gg\square), (\{\mathbf{ge}, gge, f\ ge\ ge\}, \square)\}$. *To be continued.*

5 Observation Trees

We are now going to show that the abstract data type proposed in the previous section can alternatively be implemented by an *observation tree*. For MAT learning of regular string languages, this idea has roughly been described earlier by Kearns and Vazirani in [18, Chapter 8]. The expected advantage of observation trees over observation tables is that they require a smaller number of coefficient queries to be asked. This is important if we want to make practical use of MAT learners for recognizable tree series, because such practical use normally requires an (exact or approximate) simulation of the teacher, which means that the complexity of answering coefficient and equivalence queries cannot be neglected.

To understand the idea behind observation trees, it is useful to have a look at Definition 13 and the proof of Lemma 14. Intuitively, the major purpose of the matrix \mathbf{P} is to be able to guarantee that condition (iii) of Definition 5 holds. In other words, the collected contexts provide explicit evidence that trees t, t' with $t\rho \neq t'\rho$ belong to distinct congruence classes. Suppose that, at some stage of the algorithm, there are trees $t, t' \in T$ such that $t\rho = s = t'\rho$, but the teacher provides the learner with a counterexample that, via EXTEND, reveals a

⁴ In fact, C can be left unchanged if it already contains a sign of life for t .

Algorithm 3 Function $\text{ADDTRANSITION}(t, c)$ for an abstract observation table (S, T, sol, ρ) represented by (S, τ)

Pre-conditions: $t \in \Sigma(S) \setminus T$ and $c \in C_\Sigma$ with $ct \in \text{supp}(\psi)$

Post-conditions: return an S' -observation tree (S', τ') such that $S \subseteq S'$ and $T \cup \{t\} \subseteq T(\tau')$, representing an abstract observation table $\text{aot}' = (S', T', \text{sol}', \rho')$

```

    v ← nodτ(t)
2: if τ = ∅ then
    return ({t}, □[c[{t}]]) // aot' = ({t}, {t}, ⟨t := c'⟩, ⟨t := t⟩), c' ∈ {c, □}
4: if v ∈ leaves(τ) then
    return (S, τ[v ← τ(v) ∪ {t}]) // aot' = (S, T ∪ {t}, sol, ρ⟨t := s⟩)
    // where S ∩ τ(v) = {s}
6: let τ/v = c'[τ1, ..., τk] and u = c[{t}]
    if {ψ(dt) | d ∈ Cτ(v)} ≠ {0} then
8:   c ← d; u ← {t} // sign of life c not needed
    return (S ∪ {t}, τ[v ← c'[τ1, ..., τk, u]]) // aot' = (S ∪ {t}, T ∪ {t},
    // sol⟨t := c⟩, ρ⟨t := t⟩)

```

separating context c . The addition of c to the table divides the set $T' = \rho^{-1}(s)$ into subsets T'_1, \dots, T'_k with $k \geq 2$. The addition of c may also subdivide some of the other sets $\rho^{-1}(s')$ with $s' \in S \setminus \{s\}$ as a side effect. A similar effect may occur when a transition is added; see Example 15. Although these effects are welcome (because they speed up convergence), they have the disadvantage of forcing us to query the teacher for all the coefficients $\psi(ct')$ with $t' \in T$. To avoid the latter, we may organize our data in a tree, where the internal nodes are contexts and the leaves are the sets in $T/\ker(\rho)$.⁵ In the situation considered above, when the new context c has been discovered, the leaf T' would be replaced with $c[T'_1, \dots, T'_k]$. Then, only the coefficients $\psi(ct')$ for all $t' \in T'$ need to be asked for.

Formally, let Ω be the infinite ranked set such that $\Omega_0 = \text{fin}(T_\Sigma)$ and $\Omega_k = C_\Sigma$ for every $k \geq 1$. For a tree $\tau \in T_\Omega$ and $v \in \text{nodes}(\tau) \setminus \text{leaves}(\tau)$, we let $C_\tau(v)$ denote the set of contexts on the path from the root of τ to v , including the latter. In other words, if $v_1 = \epsilon, \dots, v_n = v$ are the prefixes of v , then $C_\tau(v) = \{\tau(v_1), \dots, \tau(v_n)\}$. Below, we also use the notation $T(\tau)$ to designate the set $\bigcup_{v \in \text{leaves}(\tau)} \tau(v)$.

Definition 16. A tree $\tau \in T_\Omega$ is an observation tree if

1. $\tau(\epsilon) \in \{\square, \emptyset\}$,
2. for all $v \in \text{nodes}(\tau) \setminus \text{leaves}(\tau)$, if $\tau/v = c[\tau_1, \dots, \tau_k]$, then

$$T(\tau/v)/_{C_\tau(v)} = \{T(\tau_1), \dots, T(\tau_k)\} \text{ , and}$$

3. for all $v \in \text{leaves}(\tau)$ and $t \in \tau(v)$, $C_\tau(v)$ contains a sign of life for t .

⁵ $T/\ker(\rho)$ denotes the quotient of T under the equivalence $\{(t, t') \in T^2 \mid t\rho = t'\rho\}$. For technical convenience, we let $T/\ker(\rho) = \{\emptyset\}$ in the special case where $T = \emptyset$.

Algorithm 4 Function $\text{ADDSTATE}(t, c)$ for an abstract observation table (S, T, sol, ρ) represented by (S, τ)

Pre-conditions: $t \in T$ and $c \in C_\Sigma$ with $\frac{\psi(ct)}{\psi(\text{sol}(t)t)} \neq \frac{\psi(c(t\rho))}{\psi(\text{sol}(t)(t\rho))}$

Post-conditions: return an S' -observation tree (S', τ') such that $S \cup \{t\} \subseteq S'$ and $T \subseteq T(\tau')$, representing an abstract observation table $\text{aot}' = (S', T', \text{sol}', \rho')$

let $\text{nod}_\tau(t) = v = ui$ with $i \in \mathbb{N}$ and $\{T_1, \dots, T_k\} = \tau(v) / \equiv_{C_\tau(u) \cup \{c\}}$

2: choose $s_1 \in T_1, \dots, s_k \in T_k$, such that $\{t, t\rho\} \subseteq \{s_1, \dots, s_k\}$

return $(S \cup \{s_1, \dots, s_k\}, \tau[v \leftarrow c[T_1, \dots, T_k]])$

// $\text{aot}' = (S \cup \{s_1, \dots, s_k\}, T, \text{sol}', \rho')$, where

// $\text{sol}' = \text{sol}(s_1 := \text{sol}(t), \dots, s_k := \text{sol}(t))$,

// $\rho'(t_i) = s_i$ for $i \in [k]$ and $t_i \in T_i$, and

// $\rho'(u) = u\rho$ for $u \in T \setminus \tau(v)$

Given a set S such that $S \subseteq T(\tau) \subseteq \Sigma(S)$, the pair (S, τ) is an S -observation tree if $\tau = \emptyset$ or $|\tau(v) \cap S| = 1$ for all $v \in \text{leaves}(\tau)$.

Let us now see how observation trees can implement abstract observation tables. For this, let (S, τ) be an S -observation tree. We define the *abstract observation table* (S, T, sol, ρ) represented by (S, τ) , as follows. The set T is given by $T(\tau)$. By the second condition, for every tree $t \in T$, there is a unique leaf u of τ such that $t \in \tau(u)$. Henceforth, we denote u by $\text{nod}_\tau(t)$. Now, for every $s \in S$, define $\text{sol}(s) = \tau(v)$ where v is the shortest prefix of $\text{nod}_\tau(s)$ such that $\tau(v)$ is a sign of life for s . By the third condition, v exists, and by the first condition it is equal to ϵ (yielding $\text{sol}(s) = \square$) if $s \in \text{supp}(\psi)$. Finally, the definition of ρ is straightforward: $t\rho$ is the unique element of $\tau(\text{nod}_\tau(t)) \cap S$ for every $t \in T$.

It should be clear that the tuple (S, T, sol, ρ) constructed in this way fulfils the conditions of Definition 5. It remains to give implementations of INITIALIZE , ADDTRANSITION , and ADDSTATE . Unsurprisingly, INITIALIZE returns (S, τ) with $S = \emptyset$ and $\tau = \emptyset$. The functions ADDTRANSITION and ADDSTATE are given in Algorithms 3 and 4, respectively. In their definitions, we use the following extension of nod_τ . For a tree $t \in T_\Sigma$, let $\text{nod}_\tau(t)$ be the maximal node $v \in \text{nodes}(\tau)$ (with respect to $|v|$) such that $t \equiv_{C_\tau(u)} t'$ for all $t' \in T(\tau/v)$ and all proper prefixes u of v . Note that v is uniquely determined, and that the requirement “ $t \equiv_{C_\tau(u)} t'$ for all $t' \in T(\tau/v)$ ” is equivalent to “ $t \equiv_{C_\tau(u)} t'$ for a $t' \in T(\tau/v)$ ” (both by the second condition). The latter makes it possible to find $\text{nod}_\tau(t)$ efficiently. The reader should also notice that the extension of nod_τ is consistent with the earlier definition of $\text{nod}_\tau(t)$ for $t \in T$.

The last case distinction in ADDTRANSITION is needed only for efficiency reasons; i.e., to keep the observation tree small. If efficiency is not a concern, then $(S, \tau[v \leftarrow c'[\tau_1, \dots, \tau_k, c[\{t\}]]])$ can be returned in either case. In fact, a similar case distinction could be made in line 3, because c is not needed if $t \in \text{supp}(\psi)$.

Example 17 (cont'd). If we use observation trees when processing the counterexamples of Example 9 (again, always re-using counterexamples as long as possible), we get the sequence of observation trees depicted in Figure 3.

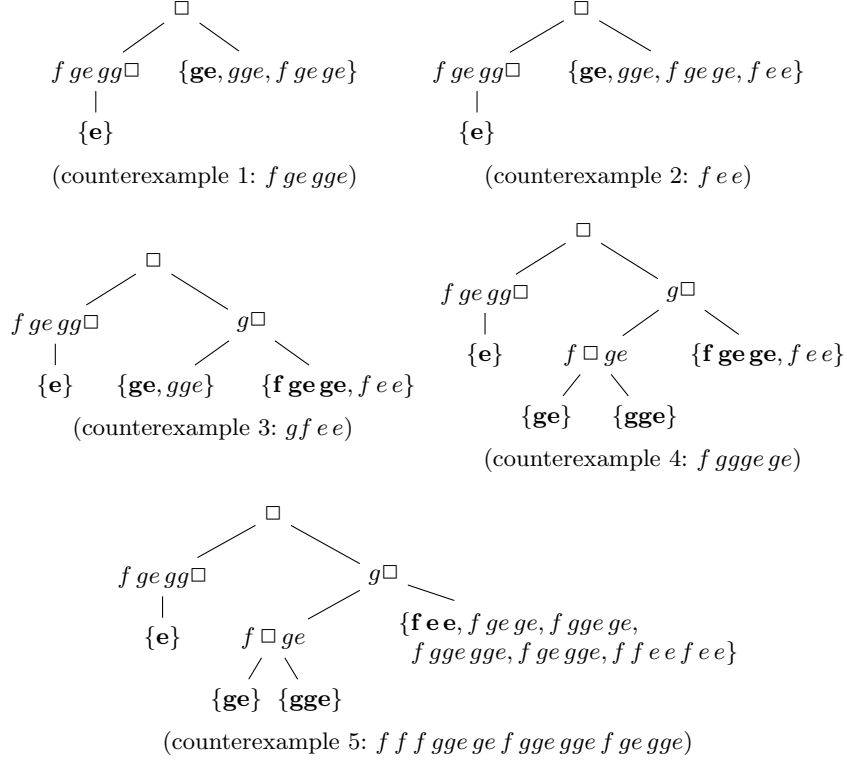


Fig. 3. Observation trees encountered when using the counterexamples of Example 9.

Lemma 18. *Observation trees implement abstract observation tables.*

Proof. It suffices to show that `ADDTRANSITION` and `ADDSTATE` are correct; i.e., that they return S' -observation trees (S', τ') with $S \subseteq S'$ and $T \cup \{t\} \subseteq T(\tau')$ in case of `ADDTRANSITION`, and $S \cup \{t\} \subseteq S'$ and $T \subseteq T(\tau')$ in case of `ADDSTATE`.

Correctness of ADDTRANSITION. The return statement in line 3 is obviously correct, because c is a sign of life for t . If the condition in line 4 holds, and u is the parent node of v , then $t \equiv_{C_\tau(u)} t'$ for all $t' \in \tau(v)$. Hence, it follows that the addition of t to $\tau(v)$ does not violate any of the requirements imposed on S -observation trees. Finally, consider the third case. By the definition of $\text{nod}_\tau(t)$, for all $t' \in \tau(v)$, it holds that $t \not\equiv_{C_\tau(v)} t'$ but $t \equiv_{C_\tau(u)} t'$ for all proper prefixes u of v . Consequently, $(S \cup \{t\}, \tau[v \leftarrow c'[\tau_1, \dots, \tau_k, \{t\}]])$ satisfies all conditions imposed on S -observation trees, with the possible exception of condition 3. If condition 3 is violated, then $(S \cup \{t\}, \tau[v \leftarrow c'[\tau_1, \dots, \tau_k, c[\{t\}]])$ satisfies it, since c is a sign of life for t . Clearly, the remaining conditions are not affected by the insertion of c . Hence, the return statement in line 9 is correct.

Correctness of ADDSTATE. Let (S', τ') be the pair returned by the algorithm. Concerning line 1, notice first that $v \neq \epsilon$, because $T \neq \emptyset$ and, thus, $\tau(\epsilon) = \square$. The pre-condition ensures that $t \not\equiv_{C_\tau(u) \cup \{c\}} t\rho$; i.e., $t \in T_i$ and $t\rho \in T_j$

for distinct $i, j \in [k]$ in line 1. Hence, s_1, \dots, s_k can be chosen as required in line 2, which means that (S', τ') with $S' = S \cup \{s_1, \dots, s_k\}$ and $\tau' = \tau[v \leftarrow c[T_1, \dots, T_k]]$ satisfies condition 2 of Definition 16. Further, condition 3 is satisfied since τ satisfies it, $T(\tau') = T = T(\tau)$, and $C_{\text{nod}_\tau(t')} \subseteq C_{\text{nod}_{\tau'}(t')}$ for all $t' \in T$. Hence, (S', τ') is an S' -observation tree. \square

Let us roughly compare the size of an observation table \mathbf{P} , and the number of coefficient queries required to build it, with the corresponding numbers for an observation tree τ . Clearly, the number of rows of \mathbf{P} is equal to the cardinality of $T(\tau)$, because both are equal to $|T|$. For non-trivial cases, the number K of columns of \mathbf{P} lies between 2 and $|S|$. These bounds are sharp. On the one hand, two contexts may separate any number of equivalence classes from each other. On the other hand, $|S|$ contexts may be needed to separate the $|S|$ equivalence classes from each other. Thus, \mathbf{P} has between $2|T|$ and $|S| \cdot |T|$ cells, requiring as many coefficient queries.

When using an observation tree τ , the number of coefficient queries required to build it is determined by the depth d at which the trees of T reside in τ . More precisely, let $d(t) = |\text{nod}_\tau(t)|$ for $t \in T$. Since a coefficient query has to be asked for each node v such that v is a proper prefix of $\text{nod}_\tau(t)$, the overall number of coefficient queries used to build τ is $D(\tau) = \sum_{t \in T} d(t)$. From the observation that $d(t) \leq |S| + 1$ for all $t \in T$, we obtain the worst-case estimation $D(\tau) \leq (|S| + 1) \cdot |T|$, which is essentially the same as above. In the best case, $d(t) = 2$ for all $t \in T$, again yielding the same estimation as above.

So, why should τ have an advantage over \mathbf{P} ? The reason is that, in most cases, one may expect the average of all $d(t)$ to be considerably smaller than the number K of columns of \mathbf{P} . This is because the contexts indexing the columns of \mathbf{P} must simultaneously separate *all* trees in S from each other, whereas the contexts in $C_\tau(u)$, for a leaf $v = ui$ of τ , only need to separate the one tree of $\text{nod}_\tau(v) \subseteq S$ from the remaining ones. In other words, we expect $d_{\text{avg}} = \frac{1}{|T|} \sum_{t \in T} d(t)$ to be considerably smaller than K , and thus, $D(\tau) = d_{\text{avg}}|T|$ to be considerably smaller than $K|T|$.

Of course, in concrete cases, there are many factors that can affect d_{avg} . A thorough study of d_{avg} and its relation with K is beyond the scope of this article. Such a study should take into account the properties of the tree series ψ to be learned and suitable probabilistic assumptions regarding the behaviour of the MAT.

6 Experiments

As established in Section 5, observation trees and observation tables are both proper realizations of the abstract observation table (aot) of Section 3. The termination and correctness of the learner are thus guaranteed when instantiated with either data structure. As argued in Section 5, observation trees are expected to have an advantage over observation tables as the former should usually require fewer coefficient queries, but only slightly more equivalence queries, than the latter. To confirm this expectation, we implemented the relevant data structures

and algorithms in Java and conducted a series of experiments.⁶ In particular, the aot is implemented as an abstract class, the learner as an algorithm instantiated with an aot, and the observation tree and observation table as data structures realizing the aot. The teacher is implemented for various restricted and unrestricted weighted tree automata over semifields. From here on, we refer to the learner as L_*^{table} when instantiated with an observation table, and as L_*^{tree} when instantiated with an observation tree.

In our experiments, we record both the number of coefficient and equivalence queries — the latter because one may suspect that L_*^{table} , if it by chance receives contexts separating many trees from each other, may use fewer equivalence queries than L_*^{tree} .

6.1 Results and discussion

We investigate the performance of L_*^{table} and L_*^{tree} with respect to several families of tree series, each parameterized by a natural number. The formal definitions of the families are given below. To verify that the maximal rank of the underlying alphabet has little influence on the relative performance of the learners, the last two families in the list are defined over variable-rank alphabets.

- The tree series SIZE_n , $n \in \mathbb{N}$, is defined over the arctic (or max-plus) semiring $(\mathbb{N} \cup \{-\infty\}, \max, +, -\infty, 0)$ and the ranked alphabet $\Sigma = \{a^{(0)}, f^{(2)}\}$. For every $t \in \mathbb{T}_\Sigma$,

$$\text{SIZE}_n(t) = \begin{cases} n & \text{if } n = |\text{nodes}(t)|, \text{ and} \\ -\infty & \text{otherwise .} \end{cases}$$

- The tree series NUMBER_n , $n \in \mathbb{N}$, over \mathbb{R} and $\Sigma = \Sigma_0 \cup \Sigma_2$, where $\Sigma_0 = [n]$ and $\Sigma_2 = \{f\}$, is such that for every $t \in \mathbb{T}_\Sigma$,

$$\text{NUMBER}_n(t) = \begin{cases} \prod_{v \in \text{leaves}(t)} t(v) & \text{if } \text{hg}(t) \geq 1, \text{ and} \\ 0 & \text{otherwise .} \end{cases}$$

- The tree series QUARTER_n , $n \in \mathbb{N}$, over the field \mathbb{R} of reals and $\Sigma = \Sigma_0 \cup \Sigma_2$, where $\Sigma_0 = [\lfloor n/4 \rfloor]$ and $\Sigma_2 = \{f\}$, is such that for every $t \in \mathbb{T}_\Sigma$,

$$\text{QUARTER}_n(t) = \begin{cases} 1 & \text{if } t \in \mathbb{T}_{\{f, i\}}, i \in \Sigma_0, \text{ and } \text{hg}(t) \equiv 1 \pmod{n}, \text{ and} \\ 0 & \text{otherwise .} \end{cases}$$

- The tree series POWER_n , $n \in \mathbb{N}$, over \mathbb{R} and $\Sigma = \Sigma_0 \cup \Sigma_2$, where $\Sigma_0 = 2^{[n]}$ and $\Sigma_2 = \{f\}$, is such that for every $t \in \mathbb{T}_\Sigma$,

$$\text{POWER}_n(t) = \begin{cases} 1 & \text{if } \bigcup_{v \in \text{leaves}(t)} t(v) = [n], \text{ and} \\ 0 & \text{otherwise .} \end{cases}$$

⁶ The source files can be downloaded from <http://www.cs.umu.se/~johanna/adt/>.

- The tree series GIRD_n , $n \in \mathbb{N}$, over \mathbb{R} and $\Sigma = \Sigma_0 \cup \Sigma_n$, where $\Sigma_0 = \{a, b\}$ and $\Sigma_n = \{f\}$, is such that for every $t \in \mathsf{T}_\Sigma$,

$$\text{GIRD}_n(t) = \begin{cases} 1 & \text{if } \text{hg}(t) = 5 \text{ and } |t|_b = 1, \text{ and} \\ 0 & \text{otherwise.} \end{cases}$$

- The tree series ORDER_n , $n \in \mathbb{N}$, over the arctic semiring and $\Sigma = \Sigma_0 \cup \Sigma_n$, where $\Sigma_0 = \Sigma_n = \{f, g\}$, is such that for every $t \in \mathsf{T}_\Sigma$,

$$\text{ORDER}_n(t) = \begin{cases} |t| & \text{if } \nexists v \in \text{nodes}(t) \text{ and } i, j \in \mathbb{N} \text{ such that} \\ & i < j, t(vi) = g, \text{ and } t(vj) = f, \text{ and} \\ 0 & \text{otherwise.} \end{cases}$$

The number of coefficient and equivalence queries posed by L_*^{table} and L_*^{tree} when inferring the tree series are plotted in Figures 4–6. In all experiments, L_*^{tree} requires fewer coefficient queries than L_*^{table} . The only example in which it makes more equivalence queries than L_*^{table} is given by the tree series **POWER**. To learn **POWER**, L_*^{tree} needed 3–4% more equivalence queries. However, the overall computation time (including the time consumed by the teacher) was still faster than L_*^{table} . In our implementation, the overhead of using observation trees is larger than for observation tables, so although L_*^{tree} is typically faster than L_*^{table} , this was not the case for the tree series **NUMBER**. Here, the savings in terms of coefficient queries was too small to make up for the additional overhead, so L_*^{tree} ran slightly slower than L_*^{table} .

Let us, finally, discuss how we implemented the teacher. In our implementation, it is initialized with a dwta A computing the target tree series. To answer a coefficient query for a tree t , the teacher simply runs A on t . To answer an equivalence query; i.e., deciding whether A is equivalent to some dwta B , the teacher searches for a tree on which A and B disagree. It is easy to see that if A and B are dwta over the Boolean semiring, or are all-accepting dwta (meaning that every final weight is non-zero) over a semifield, then such a tree can, if it exists, be found in time $\mathcal{O}(n^r m^r)$ where r is the maximum rank of the ranked set, and n and m are the number of states of A and B , respectively. It was shown in [8] that the equivalence problem for probabilistic string automata over fields is in time $\mathcal{O}(|\Sigma|(n+m)^3)$. When the weights are taken from a semifield, the problem can, as we shall see in the subsequent section, be solved in time $\mathcal{O}(r^{nm})$ by an algorithm based on the pumping lemma of [6, Corollary 5.6].

6.2 Deciding the equivalence problem for dwta over semifields

Let us, finally, discuss how to decide the equivalence problem for dwta over semifields (and how to obtain a counterexample if the answer is negative), since such an algorithm is obviously needed to implement the teacher.

In the following, for a function $f: A \rightarrow B_1 \times B_2$, we let f_i ($i \in [2]$) denote the composition of f with the projection onto B_i . In particular, for a total dwta

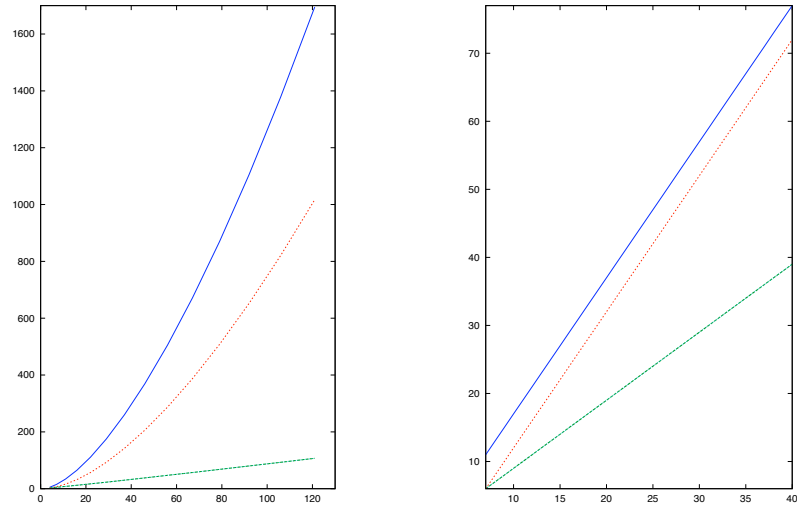


Fig. 4. The outcome of applying L_*^{table} and L_*^{tree} to the tree series SIZE (left plot) and NUMBER (right plot). The x -axis is labeled with the size of the target dwta; the y -axis with the number of queries posed. The curves are in turn (from above to below): the number of coefficient queries posed by L_*^{table} ; the number of coefficient queries posed by L_*^{tree} ; and the number of equivalence queries posed by L_*^{table} and L_*^{tree} (which coincide).

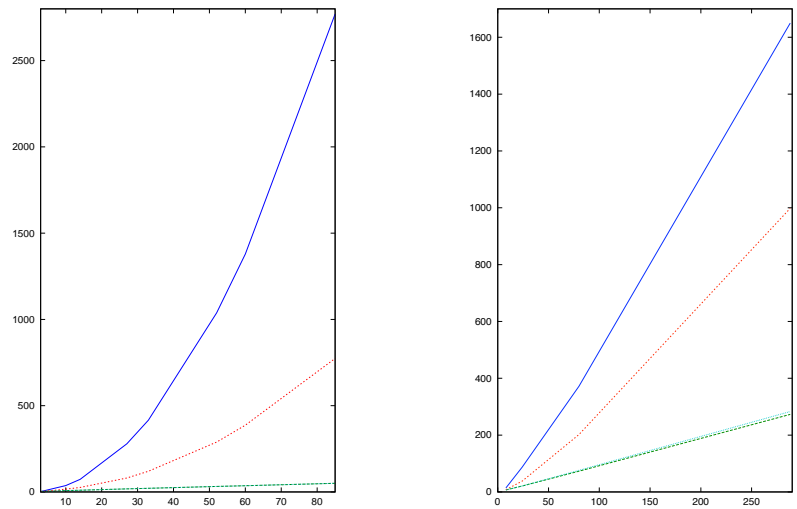


Fig. 5. The outcome of applying L_*^{table} and L_*^{tree} to the tree series QUARTER (left plot) and POWER (right plot). Curves and labels are as in Figure 4.

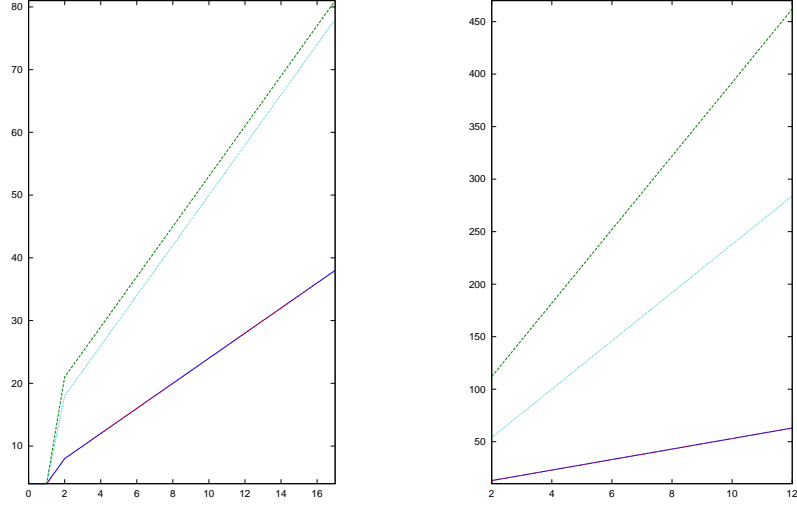


Fig. 6. The outcome of applying L_*^{table} and L_*^{tree} to the tree series RANK (left plot) and GIRD (right plot). Curves and labels are as in Figure 4.

$(\Sigma, P, \delta, \lambda)$, $\delta_1: \Sigma(P) \rightarrow P$ and $\delta_2: \Sigma(P) \rightarrow \mathbb{S} \setminus \{0\}$ are the (total) functions such that $\delta(u) = (\delta_1(u), \delta_2(u))$ for every $u \in \Sigma(P)$.

Let $A = (\Sigma, P, \delta, \lambda)$ and $B = (\Sigma, Q, \eta, \nu)$ be total dwta. By [7, Theorem 6.1.6] we can suppose, without loss of generality, that A and B have Boolean final weights (i.e., $\lambda: P \rightarrow \{0, 1\}$ and $\nu: Q \rightarrow \{0, 1\}$). We first construct the direct product total dwta $A \cdot B^{-1} = (\Sigma, P \times Q, \delta', \lambda')$ by

$$\delta'(f[\langle p_1, q_1 \rangle, \dots, \langle p_k, q_k \rangle]) = (\langle \delta_1(f[p_1, \dots, p_k]), \eta_1(f[q_1, \dots, q_k]) \rangle, \delta_2(f[p_1, \dots, p_k]) \cdot \eta_2(f[q_1, \dots, q_k])^{-1})$$

and

$$\lambda'(\langle p, q \rangle) = \begin{cases} 1 & \text{if } \lambda(p) = 1 = \nu(q) \\ 0 & \text{otherwise.} \end{cases}$$

This construction (without the inverses) is taken from [6, Definition 3.7]. Next, we observe that if any state $\langle p, q \rangle$ with $\lambda(p) \neq \nu(q)$ is reachable (i.e., $\hat{\delta}'_1(t) = \langle p, q \rangle$ for some $t \in T_\Sigma$), then clearly A and B are not equivalent because $A(t) \neq B(t)$ for any tree t that reaches this state. If such a pair does not exist, we can eliminate all states $\langle p, q \rangle$ with $\lambda(p) \neq \nu(q)$ from $A \cdot B^{-1}$ without changing the tree series computed. Moreover, in this case, for every $t \in T_\Sigma$

$$(A \cdot B^{-1})(t) = \begin{cases} 0 & \text{if } A(t) = 0 = B(t) \\ A(t) \cdot B^{-1}(t) & \text{otherwise} \end{cases}$$

(essentially by [6, Lemma 3.8]). Thus, to decide whether A and B are equivalent, it is sufficient to decide whether this reduced total dwta computes a Boolean tree

series (i.e., all coefficients are either 0 or 1). This can (effectively) be decided by [6, Corollary 6.9].

In practice, we often find ourselves working with partial automata, and in these cases it is, of course, sensible to use a decision algorithm that avoids processing dead states in the product automaton (and their associated transitions), because the transition table of a partial dwta may be exponentially smaller than the corresponding total dwta. We therefore conclude with an algorithm that operates on the same principle, but iterates over transitions rather than trees to take advantage of sparsity in the transition table.

In the following, $A = (\Sigma, P, \delta, \lambda)$ and $B = (\Sigma, Q, \eta, \nu)$ are *partial* dwta with Boolean final weights. It is computationally easy to decide if A and B have the same support. Since we operate in a semifield, the support is regular, so it suffices to test the equality of two regular languages. Moreover, the involved automata already yield deterministic unweighted tree automata for the support by setting every nonzero weight to one, so the size of the deterministic unweighted automata for the support is the same as the size of the input automata. Thus, the equality of the supports can be decided by the classical equivalence test for deterministic tree automata, running in time $\mathcal{O}(|A| \cdot |B|)$. Thus, we may henceforth assume that A and B have the same support, and that A and B have been purged of useless (i.e., unreachable or dead) states.

Lemma 19. *Let A and B be partial dwta that contain no useless states, and are such that $\text{supp}(A) = \text{supp}(B)$. The automata A and B are equivalent iff for every $p \in P$ and $q \in Q$ there is a constant $a_{p,q} \in \mathbb{S}$, such that, for all $t \in \hat{\delta}_1^{-1}(q) \cap \hat{\eta}_1^{-1}(p)$,*

- (i) $a_{p,q} = \frac{\hat{\delta}_2(t)}{\hat{\eta}_2(t)}$, and
- (ii) if p (and thus q) is final, then $a_{p,q} = 1$.

Proof. For the “if” direction, we combine conditions (i) and (ii), and thus obtain that if p and q are final states, then $\hat{\delta}_2(t) = \hat{\eta}_2(t)$ for all $t \in \hat{\delta}_1^{-1}(p) \cap \hat{\eta}_1^{-1}(q)$.

For the opposite direction, let A and B be equivalent. Consider $p \in P$, $q \in Q$, and $t, u \in \hat{\delta}_1^{-1}(p) \cap \hat{\eta}_1^{-1}(q)$. To establish (i), we have to show that $\frac{\hat{\delta}_2(t)}{\hat{\eta}_2(t)} = \frac{\hat{\delta}_2(u)}{\hat{\eta}_2(u)}$. Since A contains no dead states and $\text{supp}(A) = \text{supp}(B)$, there is a sign of life $c \in C_\Sigma$ for t and u in both A and B . Moreover,

$$\frac{\hat{\delta}_2(ct)}{\hat{\delta}_2(t)} = \hat{\delta}_2(cp) = \frac{\hat{\delta}_2(cu)}{\hat{\delta}_2(u)} \quad \text{and} \quad \frac{\hat{\eta}_2(ct)}{\hat{\eta}_2(t)} = \hat{\eta}_2(cq) = \frac{\hat{\eta}_2(cu)}{\hat{\eta}_2(u)},$$

because $\hat{\delta}_1(t) = p = \hat{\delta}_1(u)$ and $\hat{\eta}_1(t) = q = \hat{\eta}_1(u)$. We can now compute as follows:

$$\begin{aligned} \frac{\hat{\delta}_2(t)}{\hat{\eta}_2(t)} &= \frac{\hat{\delta}_2(ct)}{\hat{\delta}_2(cp)} \cdot \frac{\hat{\eta}_2(cq)}{\hat{\eta}_2(ct)} &&= \frac{\hat{\delta}_2(ct) \cdot \hat{\delta}_2(u)}{\hat{\delta}_2(cu)} \cdot \frac{\hat{\eta}_2(cu)}{\hat{\eta}_2(ct) \cdot \hat{\eta}_2(u)} \\ &= \frac{\hat{\delta}_2(ct) \cdot \hat{\delta}_2(u) \cdot \hat{\eta}_2(cu)}{\hat{\eta}_2(ct) \cdot \hat{\eta}_2(u) \cdot \hat{\delta}_2(cu)} = \frac{\hat{\delta}_2(u)}{\hat{\eta}_2(u)} \end{aligned}$$

because $\frac{\hat{\delta}_2(ct)}{\hat{\eta}_2(ct)} = 1 = \frac{\hat{\eta}_2(cu)}{\hat{\delta}_2(cu)}$ since A and B are equivalent and have Boolean final weights. This last observation also shows that the second condition holds. \square

Algorithm 5 traverses the transitions of the product automaton $A \cdot B^{-1}$ to compute a constant $a_{p,q} \in \mathbb{S}$ for each pair of states $\langle p, q \rangle \in P \times Q$ that is reachable; i.e., each pair with $\hat{\delta}_1^{-1}(p) \cap \hat{\eta}_1^{-1}(q) \neq \emptyset$. (Note that the constants $a_{p,q}$ for unreachable pairs of states are irrelevant, as they can be chosen arbitrarily.) For this purpose, the algorithm maintains a partial mapping τ that assigns to each pair of states $\langle p, q \rangle \in P \times Q$ found reachable a constant in \mathbb{S} that is the current candidate for $a_{p,q}$. In the algorithm, the domain of τ is denoted by $\text{dom}(\tau)$. The algorithm starts with the totally undefined mapping \perp and terminates when it discovers a violation of Lemma 19 or has reached a stable state.

Algorithm 5 Decide if A and B are equivalent.

Pre-conditions: $A = (\Sigma, P, \delta, \lambda)$ and $B = (\Sigma, Q, \eta, \nu)$ are partial dwta with Boolean final weights, contain no useless states, and are such that $\text{supp}(A) = \text{supp}(B)$.

```

     $\tau \leftarrow \perp$ 
2: repeat
    unchanged  $\leftarrow true$ 
4:   for all  $\langle p, q \rangle \in P \times Q$  such that  $\exists f \in \Sigma_k$  and  $\langle p_1, q_1 \rangle, \dots, \langle p_k, q_k \rangle \in \text{dom}(\tau)$  with
        $\delta_1(f[p_1, \dots, p_k]) = p$  and  $\eta_1(f[q_1, \dots, q_k]) = q$  do
        $a_{p,q} \leftarrow \delta_2(f[p_1, \dots, p_k]) \cdot \eta_2(f[q_1, \dots, q_k])^{-1} \cdot \tau(\langle p_1, q_1 \rangle) \cdot \dots \cdot \tau(\langle p_k, q_k \rangle)$ 
6:   if  $p$  is final and  $a_{p,q} \neq 1$  then
       return false
8:   if  $\langle p, q \rangle \notin \text{dom}(\tau)$  then
        $\tau(\langle p, q \rangle) \leftarrow a_{p,q}$ 
10:  unchanged  $\leftarrow false$ 
       else if  $\tau(\langle p, q \rangle) \neq a_{p,q}$  then
12:    return false
    until unchanged
14: return true

```

Lemma 20. *Algorithm 5 decides if A and B are equivalent.*

Proof. Termination is obvious, because every execution of the main loop except the last enlarges $\text{dom}(\tau)$. We show that Algorithm 5 returns *true* if and only if there is, for every $p \in P$ and $q \in Q$, a constant $a_{p,q} \in \mathbb{S}$ that fulfills Conditions (i) and (ii) of Lemma 19.

Suppose that Algorithm 5 returns *false*. This can happen in two cases.

- In the first case (lines 11 and 12), the algorithm has reached the same pair of states $\langle p, q \rangle$ on distinct trees s and t , such that $\frac{\hat{\delta}_2(t)}{\hat{\eta}_2(t)} \neq \frac{\hat{\delta}_2(s)}{\hat{\eta}_2(s)}$. This violates Condition (i) of Lemma 19.

- In the second case (lines 6 and 7), the algorithm has discovered that a pair of final states $\langle p, q \rangle$ are reachable on a tree t such that $\frac{\hat{\delta}_2(t)}{\hat{\eta}_2(t)} \neq 1$, which violates Condition (ii) of Lemma 19.

For the other direction, suppose that the algorithm returns *true* after some iterations of the main loop. For a contradiction, assume that there is a minimal tree $t = f[t_1, \dots, t_k] \in \hat{\delta}_1^{-1}(p) \cap \hat{\eta}_1^{-1}(q)$ such that $\frac{\hat{\delta}_2(t)}{\hat{\eta}_2(t)} \neq \tau(\langle p, q \rangle)$, which includes the possibility that $\tau(\langle p, q \rangle)$ is undefined. Let $p_i = \hat{\delta}_1(t_i)$ and $q_i = \hat{\eta}_1(t_i)$ for all $i \in [k]$. By the minimality assumption, we have $\frac{\hat{\delta}_2(t_i)}{\hat{\eta}_2(t_i)} = \tau(\langle p_i, q_i \rangle)$. Now, consider the last execution of the body of the main loop. At some point during that execution, the body of the *for all* loop will be executed with the given choice of $\langle p, q \rangle$, f , and $\langle p_1, q_1 \rangle, \dots, \langle p_k, q_k \rangle$. Since neither the condition in line 6 nor the one in line 8 is fulfilled (the latter because $\text{dom}(\tau)$ does not change, according to the termination condition of the main loop), line 11 is reached. However, since line 12 is not reached, this means that

$$\begin{aligned}
&= a_{p,q} \\
&= w_A \cdot w_B^{-1} \cdot \tau(\langle p_1, q_1 \rangle) \cdot \dots \cdot \tau(\langle p_k, q_k \rangle) \\
&= (w_A \cdot \hat{\delta}_2(t_1) \cdot \dots \cdot \hat{\delta}_2(t_k)) \cdot (w_B \cdot \hat{\eta}_2(t_1) \cdot \dots \cdot \hat{\eta}_2(t_k))^{-1} \\
&= \frac{\hat{\delta}_2(t)}{\hat{\eta}_2(t)},
\end{aligned}$$

contradicting the assumption that $\frac{\hat{\delta}_2(t)}{\hat{\eta}_2(t)} = \tau(\langle p, q \rangle)$. □

Lemma 21. *Algorithm 5 executes in time $\mathcal{O}(|A||B|)$.*

Proof. For an efficient implementation, we begin by calculating a sign of life for every state in the smaller of the two automata. This is done by computing representative trees for each state, at a cost of $\mathcal{O}(\min(|A|, |B|))$ operations, and then searching for the shortest path from each state to an accepting state, consuming another $\mathcal{O}(\min(|P|^2, |Q|^2)) \leq \mathcal{O}(\min(|A|^2, |B|^2))$ operations. In the worst case, the main loop must traverse every transition in $A \cdot B^{-1}$, but since this is sufficient, the complexity of this loop, and of the entire algorithm, is $\mathcal{O}(|A \cdot B^{-1}|) = \mathcal{O}(|A| \cdot |B|)$. □

We finally note that Algorithm 5 can easily be extended to return a counterexample whenever A and B are found to be different. For this, it suffices to store, along with each of the values $\tau(\langle p, q \rangle)$, a corresponding tree $t_{p,q}$ such that $\frac{\hat{\delta}_2(t_{p,q})}{\hat{\eta}_2(t_{p,q})} = \tau(\langle p, q \rangle)$. Then the two *return* statements in lines 7 and 12 can be adapted in the obvious way to return a counterexample.

Acknowledgment We thank the referees for their useful comments.

References

- [1] Dana Angluin. Learning regular sets from queries and counterexamples. *Inform. and Comput.*, 75(2):87–106, 1987.
- [2] Dana Angluin. Queries and concept learning. *Machine Learning*, 2(4):319–342, 1987.
- [3] Dana Angluin. Queries revisited. In *Proc. 12th Int. Conf. Algorithmic Learning Theory*, volume 2225 of LNCS, pages 12–31. Springer, 2001.
- [4] Jean Berstel and Christophe Reutenauer. Recognizable formal power series on trees. *Theoret. Comput. Sci.*, 18(2):115–148, 1982.
- [5] Björn Borchardt. The Myhill-Nerode theorem for recognizable tree series. In *Proc. 7th Int. Conf. Developments in Language Theory*, volume 2710 of LNCS, pages 146–158. Springer, 2003.
- [6] Björn Borchardt. A pumping lemma and decidability problems for recognizable tree series. *Acta Cybernet.*, 16(4):509–544, 2004.
- [7] Björn Borchardt. *The Theory of Recognizable Tree Series*. PhD thesis, Technische Universität Dresden, 2005.
- [8] Corinna Cortes, Mehryar Mohri, and Ashish Rastogi. Lp distance and equivalence of probabilistic automata. *J. Comput. System Sci.*, 18(4):761–779, 2007.
- [9] Colin de la Higuera and José Oncina. Learning stochastic finite automata. In *Proc. 7th Int. Coll. Grammatical Inference*, volume 3264 of LNCS, pages 175–186. Springer, 2004.
- [10] Frank Drewes. MAT learners for recognizable tree languages and tree series. *Acta Cybernet.*, 19(2):249–274, 2009.
- [11] Frank Drewes and Johanna Högberg. Extensions of a MAT learner for regular tree languages. In *Proc. 23rd Annual Workshop of the Swedish Artificial Intelligence Society*, pages 35–44, 2006.
- [12] Frank Drewes and Johanna Högberg. Query learning of regular tree languages: How to avoid dead states. *Theory of Comput. Syst.*, 40(2):163–185, 2007.
- [13] Frank Drewes and Heiko Vogler. Learning deterministically recognizable tree series. *J. Automata, Languages and Combinatorics*, 12(3):332–354, 2007.
- [14] Zoltán Fülöp and Heiko Vogler. Weighted tree automata and tree transducers. In Werner Kuich, Manfred Droste, and Heiko Vogler, editors, *Handbook of Weighted Automata*, chapter 9, pages 313–403. Springer, 2009.
- [15] Jonathan S. Golan. *Semirings and their Applications*. Kluwer Academic, Dordrecht, 1999.
- [16] E. Mark Gold. Language identification in the limit. *Information and Control*, 10(5):447–474, 1967.
- [17] Amaury Habrard and Jose Oncina. Learning multiplicity tree automata. In *Proc. 8th Int. Coll. Grammatical Inference*, volume 4201 of LNCS, pages 268–280. Springer, 2006.
- [18] Michael J. Kearns and Umesh V. Vazirani. *An Introduction to Computational Learning Theory*. MIT Press, 1994.
- [19] Andreas Maletti. Learning deterministically recognizable tree series — revisited. In *Proc. 2nd Int. Conf. Algebraic Informatics*, volume 4728 of LNCS, pages 218–235. Springer, 2007.
- [20] José Oncina. Using multiplicity automata to identify transducer relations from membership and equivalence queries. In *Proc. 9th Int. Coll. Grammatical Inference*, volume 5278 of LNCS, pages 154–162. Springer, 2008.

- [21] Yasubumi Sakakibara. Learning context-free grammars from structural data in polynomial time. *Theoret. Comput. Sci.*, 76(2–3):223–242, 1990.
- [22] Ehud Y. Shapiro. *Algorithmic Program Debugging*. ACM Distinguished Dissertation. MIT Press, 1983.
- [23] Leslie G. Valiant. A theory of the learnable. *Commun. ACM*, 27(11):1134–1142, 1984.