

M1 MPRI : Automates et Applications

Lecture 3 Logic for trees

`kn@lri.fr`

January 23, 2023

1 Introduction

2 Logic for trees

3 XML and its languages

We have introduced a new computation model, *tree automata*.
Recall that these can be:

- Deterministic or non-deterministic, as in the word case
- Bottom-up or top-down

Furthermore : $TD_{\text{det}} \subsetneq TD_{\text{nondet}} = BU_{\text{nondet}} = BU_{\text{det}}$

What is a high-level language for trees?

The notion of regular expression for words is standard and pervasive: operating system shells, programming languages, languages theory. For trees, the situation is a bit different. We study two formalisms:

- 1 Defining languages as sets that verify a property given in a particular logic
- 2 Concrete languages from the W3C

1 Introduction

2 Logic for trees

3 XML and its languages

Problems with automata

By definition, regular tree languages are those that can be recognized by a BU TA. We have also seen on examples how to encode some properties as structural properties on trees (e.g. the set of true statements of the prop. logic).

The problem is that TA are very low-level. If a property must hold at for arbitrarily deep paths, one must:

- describe explicitly how to go from the leaves to such interesting nodes
- change state to indicate that we have seen such an interesting node
- continue up to the root

For uninteresting subtrees, must be ignored explicitly.

Can we do it differently?

Recal the example of right combs on $\Sigma = \{f^2, a\}$. We would like to write:

“the set of trees which are either a or f -labeled nodes which have an a left subtree and a tree of the set as a right subtree”

without bothering with loops, refusing nodes (sink states).

Definition

A first-order logical formula is a finite production of the following grammar:

$\phi ::=$	$U(x)$	<i>unary predicate</i>
	$B(x, y)$	<i>binary predicate</i>
	$\neg\phi$	<i>negation</i>
	$\phi \vee \phi$	<i>disjunction</i>
	$\exists x.\phi$	<i>existential quantification</i>

We can use some syntactic sugar to define well known operations:

$$\phi_1 \wedge \phi_2 \equiv \neg(\neg\phi_1 \vee \neg\phi_2) \quad \text{conjunction}$$

$$\phi_1 \Rightarrow \phi_2 \equiv \phi_1 \vee \neg\phi_2 \quad \text{implication}$$

$$\phi_1 \Leftrightarrow \phi_2 \equiv (\phi_1 \Rightarrow \phi_2) \wedge (\phi_2 \Rightarrow \phi_1) \quad \text{equivalence}$$

$$\phi_1 \oplus \phi_2 \equiv \neg(\phi_1 \Leftrightarrow \phi_2) \quad \text{X-or}$$

$$\forall x\phi \equiv \neg\exists x.\neg\phi \quad \text{universal quantification}$$

For a *fixed* ranked alphabet Σ , we define the following predicates:

- $\text{lab}_a(x)$, for all $a \in \Sigma$
- $\text{child}_i(x, y)$, for all $1 \leq i \leq \max\{n \in \mathbb{N} \mid n = |f|, \forall f \in \Sigma\}$
- $x \preceq y$
- $x = y$

Definition

Let Σ be a ranked and $t \in \mathcal{T}(\Sigma)$. Let $\pi, \pi_1, \pi_2 \in \text{dom}(t)$, The truth value of predicates is defined as:

$$\text{lab}_a(\pi) \quad \equiv \quad t(\pi) = a$$

$$\text{child}_i(\pi_1, \pi_2) \quad \equiv \quad \pi_1 i = \pi_2$$

$$\pi_1 \preceq \pi_2 \quad \equiv \quad \pi_1 \leq_{\text{lex}} \pi_2$$

Predicate $x \preceq y$ is called *document order* and it means that x is visited before y during a depth-first, left-to-right traversal of the tree.

Henceforth, when we say first-order logic, we mean first-order logic with this fixed set of predicates (for a given Σ).

The semantics of formulæ is given by the judgement $t, \gamma \vdash \phi$ which means that for a tree t and a valuation γ , the formula ϕ is true. A valuation is a mapping from free variables of the formula to $\text{dom}(t)$.

$$t, \gamma \vdash U(x) \quad \Leftrightarrow \quad U(\gamma(x))$$

$$t, \gamma \vdash B(x, y) \quad \Leftrightarrow \quad B(\gamma(x), \gamma(y))$$

$$t, \gamma \vdash \phi_1 \vee \phi_2 \quad \Leftrightarrow \quad t, \gamma \vdash \phi_1 \text{ OR } t, \gamma \vdash \phi_2$$

$$t, \gamma \vdash \neg\phi \quad \Leftrightarrow \quad t, \gamma \not\vdash \phi$$

$$t, \gamma \vdash \exists x.\phi \quad \Leftrightarrow \quad \exists \pi \in \text{dom}(t), \quad t, \gamma \cup \{x \mapsto \pi\} \vdash \phi$$

Let's revisit the right combs on $\Sigma = \{f^2, a\}$, we can verify that for all such tree, the following formula holds:

$$\forall x. (\text{lab}_a(x) \vee (\text{lab}_f(x) \wedge (\exists y. \text{child}_1(x, y) \wedge \text{lab}_a(y))))$$

Importantly, there is no explicit recursion. We define a local property, and we state that this property holds for all nodes.

Definition

We can define the language L_ϕ of a formula ϕ as:

$$L_\phi = \{t \mid \exists \gamma \text{ such that } t, \gamma \vdash \phi\}$$

Are we limited to languages?

Using logical formulæ allows us to rephrase some problems:

- A closed formula defines a language (as in the example)
- A formula with a free variable x correspond to a *predicate* or a *query*

Indeed, given a tree t :

$$\{\pi \in \text{dom}(t) \mid t, \{x \mapsto \pi\} \vdash \phi\}$$

represents the set of paths in the tree for which the formula holds. This can be generalized to an arbitrary number of variables. That allows us to *select* paths of interest in the tree.

These are just the consequence of the semantics:

$$L_{\phi_1} \cap L_{\phi_2} \Leftrightarrow L_{\phi_1 \wedge \phi_2}$$

$$L_{\phi_1} \cup L_{\phi_2} \Leftrightarrow L_{\phi_1 \vee \phi_2}$$

$$\bar{L}_{\phi} \Leftrightarrow L_{\neg \phi}$$

- For a fixed tree, is a formula ϕ true ? Polynomial.
- The language of a formula is empty \equiv the formula is satisfiable?
 \Rightarrow non-elementary...

$$2^{2^{\dots^{2^c}}} \} |\phi|$$

In some ways this language is too powerful...

...but surprisingly

First order logic cannot express all regular tree languages!

Ex: The set of trees on $\Sigma = \{f^2, a, b\}$ with an even number of a .

Easily recognized by $\mathcal{A} = (\{q_0, q_1\}, \Sigma, \delta, \{q_0\})$:

$$\begin{array}{l} \delta : \quad a \mapsto \{q_1\} \quad f(q_0, q_0) \mapsto \{q_0\} \\ \quad \quad b \mapsto \{q_0\} \quad f(q_1, q_0) \mapsto \{q_1\} \\ \quad \quad f(q_1, q_1) \mapsto \{q_0\} \quad f(q_0, q_1) \mapsto \{q_1\} \end{array}$$

q_0 represent a subtree with even number of a et q_1 odd numbers of a (note: this language is not TD deterministic).

FO on trees?

To summarize:

- compact: one can express properties on path without explicit recursion
- some properties cannot be expressed (FO cannot count, even modulo two)
- crazy complexity

WHAT IF WE TRIED
MORE POWER?



©XKCD

Definition (Monadic Second-order Logic)

A formula of the MSO is a finite production of the following grammar:

$$\begin{array}{ll}
 \phi ::= & \dots \quad \text{formula of FO} \\
 | & x \in X \quad \text{set membership} \\
 | & \exists X.\phi \quad \text{second-order quantification}
 \end{array}$$

X denotes *sets of paths*. We have the syntactic sugar:

$$\forall X.\phi \equiv \neg \exists X.\neg \phi.$$

We naturally express the semantics of MSO with a judgement

$$t, \gamma, \Gamma \vdash \phi$$

which means that given t, γ (from free variables to paths) and Γ (from free Variables to sets of paths), formula ϕ is true.

$$t, \gamma, \Gamma \vdash x \in X \quad \Leftrightarrow \quad \gamma(x) \in \Gamma(X)$$

$$t, \gamma \vdash \exists X. \phi \quad \Leftrightarrow \quad \exists P \subseteq \text{dom}(t), t, \gamma, \Gamma \cup \{X \mapsto P\} \vdash \phi$$

What can we express?

We can characterize the *descendants* of a path. The descendants of a path x is the set of paths Y such that:

$$\forall y. y \in Y \Leftrightarrow x = y \vee \exists z. (z \in Y \wedge (\text{child}_1(z, y) \vee \dots \vee \text{child}_k(z, y)))$$

where k is the (fixed) maximal arity of a symbol in Σ .

What can we express?

We can characterize the set of trees with an even number of a (over $\Sigma = \{f^2, a, b\}$):

What can we express?

We can characterize the set of trees with an even number of a (over $\Sigma = \{f^2, a, b\}$):

$$\begin{aligned} & \exists E. \exists O. \forall x. (x \in E \oplus x \in O) \\ & \quad \wedge \text{lab}_a(x) \Rightarrow x \in O \\ & \quad \wedge \text{lab}_b(x) \Rightarrow x \in E \\ \wedge & \text{lab}_f(x) \Rightarrow (\exists y. \exists z. \text{child}_1(x, y) \wedge \text{child}_2(x, z) \wedge \\ & \quad ((y \in E) \wedge (z \in E)) \vee ((y \in O) \wedge (z \in O))) \end{aligned}$$

Theorem (Thatcher, Wright 68)

The set of regular tree languages is exactly the set of MSO-definable languages

Proof (very rough sketch):

- put the formula in a particular canonical form
- give automata for each basic formula ($x \in X$, $\text{lab}_a(x)$, ...)
- build inductively automata for logical connective by connecting sub-automata (like for Thompson and regex, but there is a massive explosion due to quantification).

We can also go backward:

- label each transition of an automaton with a formula, initially only the formula $lab_a(ax)$
- use an algorithm similar to state elimination (merge transitions then states and repeat)

On a donc un langage :

- compact: can express complex properties without explicit recursion
- exactly equivalent to regular tree languages
- still, non-elementary complexity

1 Introduction

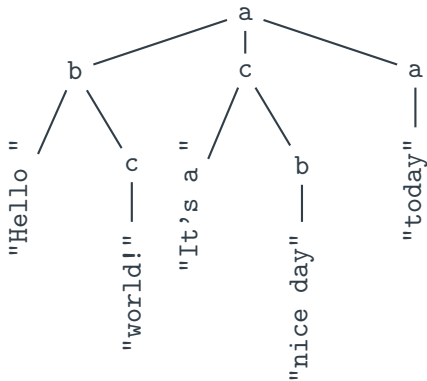
2 Logic for trees

3 XML and its languages

XML, a standard for documents

XML is a W3C standard that defines a format for storing and exchanging tree structured data. Essentially, it defines a *markup* with `<opening>` tags or `</closing>`, that must be balanced:

```
<a>  
  <b>Hello <c>world!</c></b>  
  <c>It's a <b>nice day</b></c>  
  <a>today</a>  
</a>
```



For our overview:

- we ignore many intricacies of the standard (character encodings, attributes, namespaces, ...)
- we assume all text path are replaced by leaves \$.
- we will also assumes that everything applies to HTML

Documents represent *unranked* trees

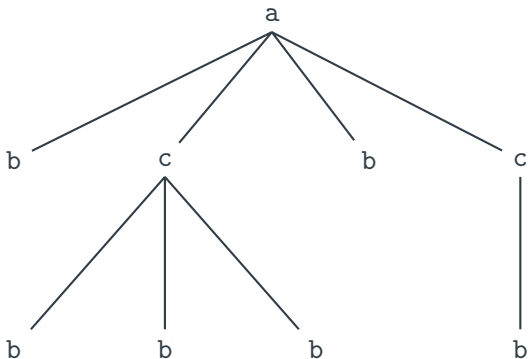
- a path may have a finite, arbitrary number of children
- the labels do not indicate the number of children

Still, we would like to make use of automata or MSO

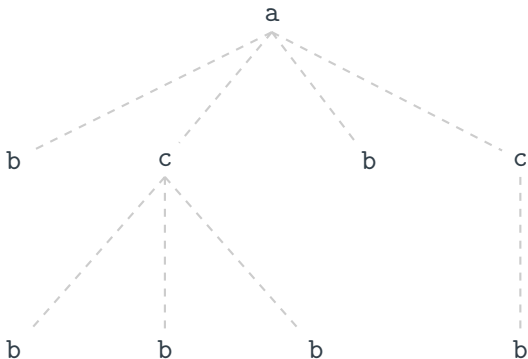
We use a trick to encode any unranked tree into a binary tree:

- the first child of a path (in the n -ary tree) is the first child in the binary tree
- the right sibling of a path (in the n -ary tree) is the second child in the binary tree

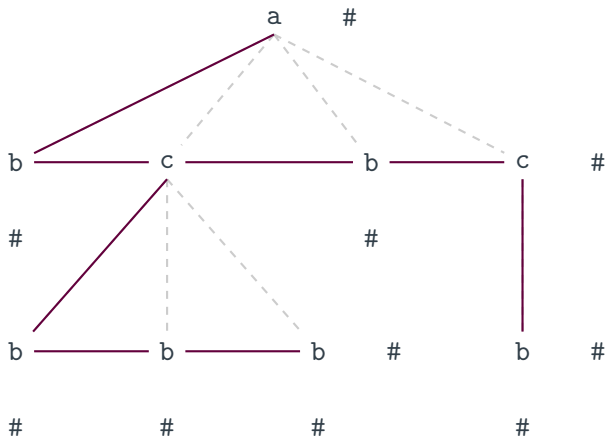
Example



Example



Example



This encoding corresponds to the usual representation of a tree in memory (for instance in C):

- a node stores the label and a pointer to the linked list of children
- the first element of the list is directly accessible, the others require a traversal
- the NULL pointer corresponds to the # leaves

Structure de données (2)

```
struct node;
struct list {
    struct list *next; //right edge
    struct node *node;
};
struct node {
    char * label;
    struct list *children; //down edge
};
```

Document Type Definitions (DTD) allows one to define schemas, that is tree languages:

```
<!ELEMENT a ( (b|c)* ) >
```

```
<!ELEMENT b ( EMPTY ) >
```

```
<!ELEMENT c ( (b|c)* ) >
```

For each tag, we give its content as a regular expression over its children. A tag determines uniquely its content (we cannot have two alternative definitions for the same tag) \Rightarrow TD déterministic

The fact that DTD are TD deterministic means that we can validate a document in *streaming*, that is with memory bounded by the height of the document.

Why is it a good thing?

The fact that DTD are TD deterministic means that we can validate a document in *streaming*, that is with memory bounded by the height of the document.

Why is it a good thing?

We can validate the document while reading the file.

What limitations ?

The fact that DTD are TD deterministic means that we can validate a document in *streaming*, that is with memory bounded by the height of the document.

Why is it a good thing?

We can validate the document while reading the file.

What limitations ?

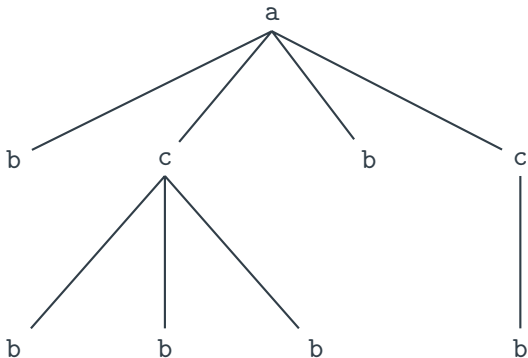
We cannot specify documents where there is a dependency between children.

The XPath language is a query language for trees. It allows one to select nodes. In its simplest form, (navigational Core XPath), queries have the form:

$$path ::= axe_1 :: test_1 [pred_1] / \dots / axe_n :: test_n [pred_n]$$
$$axe ::= child \mid descendant \mid parent \mid ancestor \mid \dots$$
$$test ::= tag \mid *$$
$$pred ::= path \mid pred_1 \text{ or } pred_2 \mid pred_1 \text{ and } pred_2 \mid \text{not}(pred_1)$$

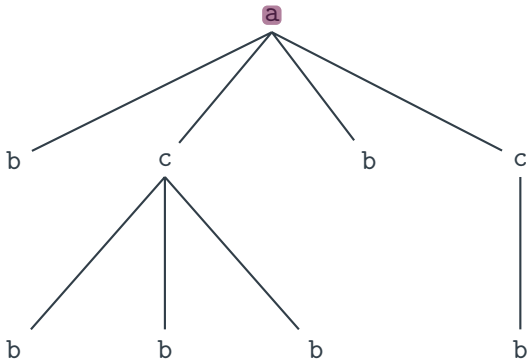
Exemple

/descendant::b[parent::c]/parent::* /child::*



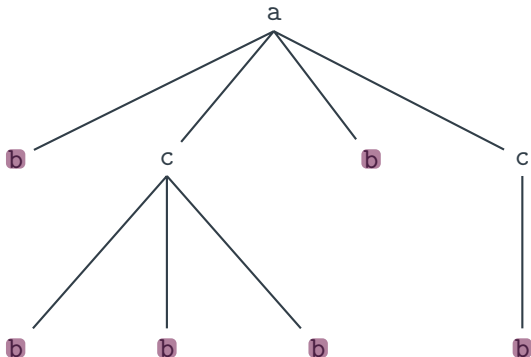
Exemple

`/descendant::b[parent::c]/parent::* /child::*`

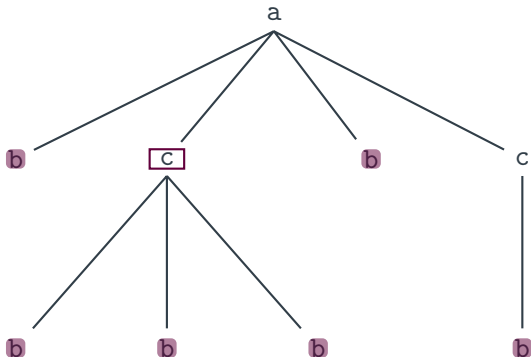


Exemple

`/descendant::b[parent::c]/parent::* /child::*`

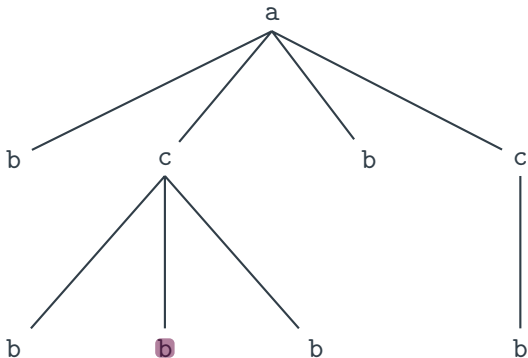


/descendant::b[parent::c]/parent::* /child::*



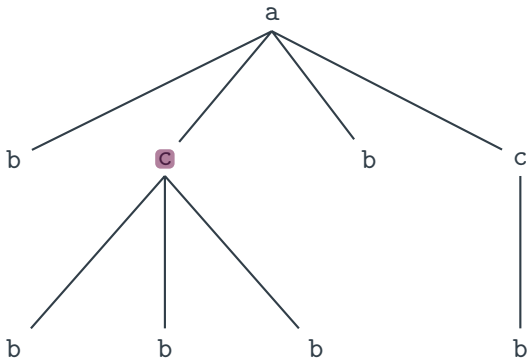
Exemple

`/descendant::b[parent::c]/parent::* /child::*`



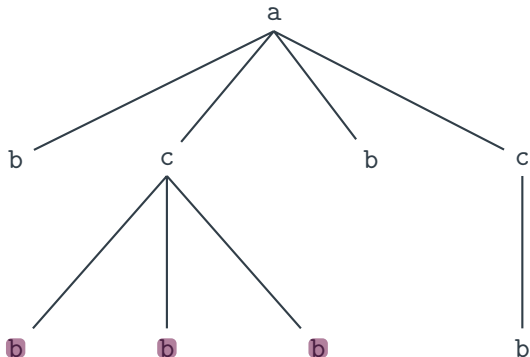
Exemple

`/descendant::b[parent::c]/parent::*`



Exemple

`/descendant::b[parent::c]/parent::*`
`child::*`



Le fragment présenté est moins expressif que FO. Par exemple, on ne peut pas exprimer : Renvoyer tous les nœuds a descendants d'un b tel qu'il n'y a pas de c entre le a et le b .

Quels problèmes souhaite-t-on résoudre ?

- 1 Étant donné un document D , est-il valide par rapport à une DTD T ?

Quels problèmes souhaite-t-on résoudre ?

- 1 Étant donné un document D , est-il valide par rapport à une DTD T ? $O(|D| + |T|)$
- 2 Étant donné un document D , quels sont les nœuds renvoyés par une requête XPath Q ?

- 1 Étant donné un document D , est-il valide par rapport à une DTD T ? $O(|D| + |T|)$
- 2 Étant donné un document D , quels sont les nœuds renvoyés par une requête XPath Q ? $O(|D| \times |Q|)$
- 3 Étant donné une DTD T et une requête XPath Q , la requête est-elle satisfiable ?

Quels problèmes souhaite-t-on résoudre ?

- 1 Étant donné un document D , est-il valide par rapport à une DTD T ? $O(|D| + |T|)$
- 2 Étant donné un document D , quels sont les nœuds renvoyés par une requête XPath Q ? $O(|D| \times |Q|)$
- 3 Étant donné une DTD T et une requête XPath Q , la requête est-elle satisfiable ?EXPTIME

Proposer un algorithme pour (2) basé sur des automates sera le sujet du devoir.