

Apprentissage des langages réguliers d'arbres à l'aide d'un expert

Jérôme Besombes*, Jean-Yves Marion*

LORIA-INRIA Lorraine
615, rue du jardin botanique
54602 Villers-lès-Nancy, France

et

INPL-Ecole des Mines de Nancy
Parc de Saurupt,
54042 Nancy Cedex, France

{Jerome.Besombes, Jean-Yves.Marion}@loria.fr

Résumé : Dans le but de modéliser l'apprentissage des langues naturelles, nous nous intéressons à l'apprentissage des langages réguliers d'arbres à partir d'exemples positifs et en interaction avec un expert. L'expert est un oracle qui connaît le langage appris et qui répond à des questions d'appartenance. Nous donnons un algorithme efficace d'apprentissage dans ce paradigme, nous montrons sa correction et l'illustrons sur un exemple détaillé.

Mots clef : langages réguliers d'arbres, apprentissage, questions d'appartenance

1 Introduction

1.1 Des motivations linguistiques

Depuis les années 60, Chomsky et son école ont mis en évidence l'existence d'une grammaire universelle innée dans laquelle chaque langue peut être décrite. Cette capacité à gérer les données permet de découvrir la structure des phrases et les catégories associées aux mots. Une conséquence est notre faculté à décider si une phrase est ou n'est pas grammaticalement correcte sans connaître nécessairement le sens de chaque mot qui la compose. L'information apportée par le sens d'un mot n'est pas directement utile, mais permet juste de désambigüiser plus rapidement une phrase. Nous voyons que le constituant de base de l'apprentissage est la phrase et non le mot. Le sens d'un mot pris comme suite de phonèmes est arbitraire. L'apprentissage d'une langue passe obligatoirement par l'apprentissage de son lexique. Cependant, une phrase n'est pas un amalgame de mots qui serait générée par un automate d'états finis piochant dans un lexique. Au contraire, nous ne connaissons pas toutes les phrases, par contre nous pouvons les générer. L'apprentissage du langage consiste à déterminer les structures des phrases.

Pour cela, il faut penser à une phrase comme un arbre dans lequel les mots dépendent les uns des autres. Cet arbre est construit à partir de la spécialisation de la grammaire universelle de la langue étudiée. Ainsi Christophe (Christophe, 2000) montre le rôle capital de l'intonation, ce que Pinker (Pinker, 1994) appelle les leçons de grammaire implicite infligées aux bébés en « Motherese ». L'intonation structure les mots de la phrase. Un bébé a accès non à une phrase linéaire mais bien à une information semi-structurée faite d'une suite de phonèmes séparés par des frontières marquées. A partir d'un vocabulaire, l'enfant commence à produire des phrases qu'il n'a jamais entendues. Cette généralisation est la conséquence de la construction de règles cohérentes avec les exemples. Ces règles mettent en relation des mots et permettent la production de structures similaires aux structures des phrases entendues. La phrase est comprise comme un objet semi-structuré, i.e. un arbre comme cela est illustré par la figure 1, et que les phrases sont construites par calcul.

Un autre aspect important doit être mentionné : l'environnement. Un enfant apprend une langue par interaction avec ses proches. Ce dialogue fournit des informations à l'enfant. Nous pourrions entamer une longue discussion sur la nature exacte de l'information apportée par l'environnement. Dans l'hypothèse la plus basse, celle-ci se réduit à savoir si oui ou non une phrase produite est correcte. Ainsi, il paraît raisonnable de penser que l'enfant peut conclure du fait qu'une phrase qu'il vient de prononcer n'est pas comprise, que celle-ci n'est pas correcte.

1.2 Inférence des langages réguliers d'arbres

D'un point de vue formel, nous nous plaçons dans le cadre général des grammaires régulières d'arbres. Bien que ces langages soient très utilisés en informatique, ils semblent délaissés pour le traitement automatique des langues (TAL). Pourtant, les grammaires pour le TAL manipulent très souvent à un certain niveau des langages réguliers d'arbres. Par exemple, les ensembles d'arbres d'analyse des grammaires algébriques et catégorielles sont réguliers. L'apprentissage des langages réguliers d'arbres a été étudiée d'un point de vue stochastique dans (Carrasco *et al.*, 1998). Notre cadre d'inférence grammaticale est celui de l'identification exacte à partir d'un ensemble fini d'exemples et de requêtes d'appartenance. Nous disposons d'un ensemble fini d'arbres d'un langage régulier L et nous cherchons un automate qui reconnaît L en posant des requêtes d'appartenance à un oracle. Pour l'apprentissage du langage, les exemples sont des phrases formées des mots et de leurs dépendances (i.e. des arbres de dépendances), et les requêtes d'appartenance sont obtenues par « l'oracle » formé par les parents de l'enfant.

1.3 Inférence à partir de requêtes d'appartenance

Le paradigme que nous considérons, est bien connu de la communauté puisqu'il s'agit d'apprentissage à partir de requêtes. Le lecteur pourra se référer à différents états de l'art comme celui de Angluin (Angluin, 1988). Différents types de requêtes sont pris en compte, à savoir l'appartenance, l'équivalence, le test de sous-ensemble, ...

De nombreux travaux ont porté sur l'inférence grammaticale à partir de tests d'appartenance et d'équivalence. Angluin (Angluin, 1987) a démontré que les langages

La phrase « *Ertol est un chat* » peut être mise sous forme d'un arbre de dépendance

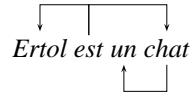


FIG. 1 – Exemples semi-structurés

réguliers de mots sont ainsi identifiables. Sakakibara (Sakakibara, 1988) a étudié les langages d'arbres de dérivation désétiqueté (*skeleton tree*) de grammaires algébriques.

A notre connaissance, le seul travail qui ne prend en compte que des requêtes d'appartenance pour l'inférence grammaticale est celui de Angluin (Angluin, 1981) au sujet des langages réguliers de mots. Dans un contexte linguistique, une requête d'équivalence n'a pas de sens. Un enfant ne peut pas demander à ses parents si la grammaire qu'il a élaboré est bien équivalente à celle du français.

Dans cet article, nous proposons un algorithme efficace, polynômial dans la taille des exemples, qui apprend les langages réguliers d'arbres. La convergence est assurée si l'ensemble d'exemples contient les exemples caractéristiques du langage cible. Typiquement, un exemple caractéristique est un exemple qui excite une transition de l'automate d'arbres. L'algorithme est proche de ceux décrits par Angluin dans (Angluin, 1981) et (Angluin, 1987). Mais il comporte aussi des différences notoires. Tout d'abord, la signature des arbres entraîne, si l'on n'y prend garde, une explosion combinatoire. En effet, Angluin teste toutes les transitions possibles d'un automate de mots. Or, le coût de cette énumération serait exponentiel en l'arité des symboles de la signature dans le cas des arbres. De ce fait, nous remplaçons cette approche par la définition d'un ensemble suffisant de contextes pour démarrer avec une partition correcte. L'autre différence est la construction d'un contexte quand l'algorithme constate que l'automate est non-déterministe. Dans le cas des mots, il suffit de rajouter une lettre. Pour les arbres, il faut calculer un ensemble de contexte pour déterminer quel est le sous arbre coupable.

2 Langages réguliers d'arbres

Pour tout ensemble de symboles \mathcal{F} et tout ensemble de variables¹ \mathcal{X} , les termes sont définis par : un symbole de \mathcal{F} est un terme, une variable de \mathcal{X} est un terme et si f est un symbole de \mathcal{F} et t_1, \dots, t_n sont des termes, alors $f(t_1, \dots, t_n)$ est un terme. L'ensemble des termes est noté $\mathcal{T}(\mathcal{F}, \mathcal{X})$ et l'ensemble $\mathcal{T}(\mathcal{F}) = \mathcal{T}(\mathcal{F}, \emptyset)$ est l'ensemble des termes clos. Ainsi, les arbres ordonnés et étiquetés se représentent par des termes. Pour tout terme t , l'ensemble des sous-termes de t est l'ensemble $\mathcal{S}(t)$ défini par : $t \in \mathcal{S}(t)$ et si $f(t_1, \dots, t_n) \in \mathcal{S}(t)$ alors $t_1, \dots, t_n \in \mathcal{S}(t)$. Pour tout ensemble de termes \mathcal{E} , $\mathcal{S}(\mathcal{E})$ est la réunion des ensembles des sous-termes des éléments de \mathcal{E} . Pour tout terme t et tout sous-terme s de t , la profondeur $\text{depth}(t, s)$ de

1. Les variables sont considérées comme des symboles spéciaux d'arité 0.

s dans t est égale à 0 si $t = s$ et égale à $1 + \text{depth}(t, f(t_1, \dots, t_i, s, t_{i+1}, \dots, t_n))$ si $f(t_1, \dots, t_i, s, t_{i+1}, \dots, t_n)$ est un sous-terme de t .

Un *contexte* est un terme $c[\diamond]$ qui contient une variable particulière \diamond à une seule occurrence. La variable \diamond marque un emplacement vide dans le terme. La substitution de \diamond par un terme s se note $c[s]$. On pourra utiliser cette notation pour indiquer que le terme s est un sous-terme de $c[s]$. Pour un ensemble donné de termes \mathcal{E} , $C[\mathcal{E}]$ est l'ensemble des contextes obtenus à partir de tout élément t de \mathcal{E} , en remplaçant exactement un sous-terme de t par \diamond . En particulier pour tout ensemble de termes \mathcal{E} , $C[\mathcal{E}]$ contient \diamond .

Un *automate d'arbres d'états fini* (NFTA) est un quadruplet $\mathcal{A} = \langle \mathcal{F}, \mathcal{Q}, \mathcal{Q}_F, \rightarrow_{\mathcal{A}} \rangle$ où \mathcal{Q} est un ensemble fini d'états, $\mathcal{Q}_F \subseteq \mathcal{Q}$ est l'ensemble des états finaux et $\rightarrow_{\mathcal{A}}$ est l'ensemble des transitions. Une transition est une règle de la forme $f(q_1, \dots, q_n) \rightarrow_{\mathcal{A}} q$ où q et q_1, \dots, q_n sont des états de \mathcal{Q} et f est un symbole de \mathcal{F} . En particulier, une transition peut être de la forme $a \rightarrow_{\mathcal{A}} q$ où le symbole a a une arité nulle.

Un automate d'arbres d'états fini est *déterministe* (DFTA) s'il n'existe pas deux transitions distinctes possédant des parties gauches identiques. Pour ces automates, la relation $\rightarrow_{\mathcal{A}}$ est le graphe d'une fonction notée $\delta_{\mathcal{A}}$. Comme dans le cas des automates de mots, il a été montré que pour tout langage reconnu par un automate non-déterministe, il existe un automate déterministe qui le reconnaît (Comon *et al.*, 1997).

La dérivation simple $\rightarrow_{\mathcal{A}}$ est définie par $t \rightarrow_{\mathcal{A}} s$ si et seulement si il existe une transition $f(q_1, \dots, q_n) \rightarrow_{\mathcal{A}} q$ telle que $t = u[f(q_1, \dots, q_n)]$ et $s = u[q]$. On notera que $\rightarrow_{\mathcal{A}} \subseteq \mathcal{T}(\mathcal{F}, \mathcal{Q}) \times \mathcal{T}(\mathcal{F}, \mathcal{Q})$, et que les états de \mathcal{Q} sont vus comme des variables.

La relation de dérivation $\rightarrow_{*,\mathcal{A}}$ est la clôture réflexive et transitive de $\rightarrow_{\mathcal{A}}$. Le langage d'arbres *reconnu* par \mathcal{A} est l'ensemble $\mathcal{L}_{\mathcal{A}} = \{t \in \mathcal{T}(\mathcal{F}) : t \rightarrow_{*,\mathcal{A}} q_f \text{ et } q_f \in \mathcal{Q}_F\}$.

Soit un automate d'arbres $\mathcal{A} = \langle \mathcal{F}, \mathcal{Q}, \mathcal{Q}_F, \rightarrow_{\mathcal{A}} \rangle$ et une fonction surjective ϕ de \mathcal{Q} dans \mathcal{Q}' , où $\mathcal{Q}' = \phi(\mathcal{Q})$ est un ensemble fini quelconque. ϕ définit un automate $\mathcal{A}' = \phi(\mathcal{A})$, avec $\mathcal{A}' = \langle \mathcal{F}, \phi(\mathcal{Q}), \phi(\mathcal{Q}_F), \phi(\rightarrow_{\mathcal{A}}) \rangle$ et $\phi(\rightarrow_{\mathcal{A}})$ est l'ensemble des règles de la forme $f(\phi(q_1), \dots, \phi(q_n)) \rightarrow_{\mathcal{A}'} \phi(q)$, pour toutes règles $f(q_1, \dots, q_n) \rightarrow_{\mathcal{A}} q$ dans $\rightarrow_{\mathcal{A}}$. Pour tout arbre t , si $t \rightarrow_{*,\mathcal{A}} q$ alors $t \rightarrow_{*,\mathcal{A}'} \phi(q)$, ce qui implique que $\mathcal{L}_{\mathcal{A}} \subseteq \mathcal{L}_{\mathcal{A}'}$. Si ϕ est bijective, on dit que ϕ est un renommage des états et l'on a : $\mathcal{L}_{\mathcal{A}} = \mathcal{L}_{\mathcal{A}'}$.

D'après le théorème de Myhill-Nerode (Comon *et al.*, 1997), tout langage \mathcal{L} , définit une relation d'équivalence $\equiv_{\mathcal{L}}$ par : pour tout t et tout s dans $\mathcal{S}(\mathcal{L})$, $t \equiv_{\mathcal{L}} s$ si et seulement si pour tout contexte $c[\diamond]$ dans $C[\mathcal{T}(\mathcal{F})]$, $c[t] \in \mathcal{L} \Leftrightarrow c[s] \in \mathcal{L}$. Le théorème de Myhill-Nerode implique que pour tout langage \mathcal{L} , il existe un automate DFTA \mathcal{A} , nommé *automate minimal* de \mathcal{L} tel que :

- $\mathcal{L} = \mathcal{L}_{\mathcal{A}}$,
- pour tout automate \mathcal{A}' tel que $\mathcal{L} = \mathcal{L}_{\mathcal{A}'}$, le nombre d'états de \mathcal{A}' est supérieur ou égal au nombre d'états de \mathcal{A} .

De plus, cet automate est unique à un renommage des états près, et pour tout arbre t et tout arbre t' on a : $\delta_{\mathcal{A}}(t) = \delta_{\mathcal{A}}(t')$ si et seulement si $t \equiv_{\mathcal{L}} t'$.

3 Ensemble caractéristique

Si \mathcal{L} est un langage régulier d'arbres quelconque et \mathcal{C} un sous-ensemble de \mathcal{L} , on dira que \mathcal{C} est un *ensemble caractéristique* de \mathcal{L} si pour un automate minimal \mathcal{A} de \mathcal{L} , pour toute règle $f(q_1, \dots, q_n) \rightarrow_{\mathcal{A}} q$ de \mathcal{A} , il existe un arbre $t = f(t_1, \dots, t_n)$ dans $\mathcal{S}(\mathcal{L})$ tel que $\forall 1 \leq i \leq n, \delta_{\mathcal{A}}(t_i) = q_i$ ². Puisque l'automate minimal d'un langage est unique à un renommage des états près, la définition d'ensemble caractéristique est indépendante du choix d'un automate minimal particulier.

De manière informelle, un ensemble caractéristique d'un langage \mathcal{L} est un sous ensemble fini \mathcal{C} tel que pour automate minimal de \mathcal{L} , chaque règle a été utilisée pour produire une fois au moins un élément de \mathcal{C} . Enfin tout ensemble fini contenant un ensemble caractéristique est aussi caractéristique.

Exemple 1 Soit l'automate minimal \mathcal{A}_1 défini par les règles :

$$\begin{array}{lcl} c & \rightarrow_{\mathcal{A}_1} & q \\ b & \rightarrow_{\mathcal{A}_1} & q_1 \\ a(q_1, q) & \rightarrow_{\mathcal{A}_1} & q_0 \\ a(q_0, q) & \rightarrow_{\mathcal{A}_1} & q_0 \end{array}$$

où q_0 est le seul état final. On a $\mathcal{L}_{\mathcal{A}_1} = \{ \underbrace{a(\dots(a(b, c)) \dots, c)}_n : n > 0 \}$ et pour tout $n > 1$, le singleton $\{ \underbrace{a(\dots(a(b, c)) \dots, c)}_n \}$ est un ensemble caractéristique de $\mathcal{L}_{\mathcal{A}_1}$. Soit pour tout $m > 0$, l'automate \mathcal{A}_m suivant :

$$\begin{array}{lcl} c & \rightarrow_{\mathcal{A}_m} & q \\ b & \rightarrow_{\mathcal{A}_m} & q_m \\ a(q_m, q) & \rightarrow_{\mathcal{A}_m} & q_{m-1} \\ & \vdots & \\ a(q_1, q) & \rightarrow_{\mathcal{A}_m} & q_0 \\ a(q_0, q) & \rightarrow_{\mathcal{A}_m} & q_0 \end{array}$$

où q_0 est le seul état final. Alors $\mathcal{L}_{\mathcal{A}_m} = \{ \underbrace{a(\dots(a(b, c)) \dots, c)}_n : n > m \}$. Pour tout n , le singleton $\{ \underbrace{a(\dots(a(b, c)) \dots, c)}_n \}$ est un ensemble caractéristique de tout langage $\mathcal{L}_{\mathcal{A}_m}$, pour $0 < m < n$.

4 Apprentissage

Nous nous intéressons à l'identification de la classe des langages réguliers d'arbres à partir d'un ensemble caractéristique. On dira qu'un algorithme I est un algorithme

2. Puisque \mathcal{A} est déterministe, on a $\delta_{\mathcal{A}}(t) = q$.

d'apprentissage pour la classe des langages réguliers d'arbres si : pour tout langage \mathcal{L} , si I reçoit un ensemble caractéristique en entrée, alors A donne un automate minimal \mathcal{A} en sortie avec $\mathcal{L} = \mathcal{L}_{\mathcal{A}}$. D'après l'exemple qui précède, il ne peut exister d'algorithme d'apprentissage ainsi défini pour la classe des langages réguliers d'arbres. En effet, si un tel algorithme existe, il doit être capable de trouver un unique automate pour l'entrée $\{a(\dots(a(b,c))\dots,c)\}$, où n est un entier quelconque. Or pour tout

$n > 1$, $\{a(\dots(a(b,c))\dots,c)\}$ est un ensemble caractéristique de plusieurs langages

(l'ensemble des $\mathcal{L}_{\mathcal{A}_m}$, pour $0 < m < n$). Pour résoudre le problème d'apprentissage, I doit pouvoir bénéficier d'information supplémentaire concernant le langage à identifier. Cette information supplémentaire lui est fournie grâce à une interaction avec un oracle O . L'oracle O connaît le langage cherché et il est capable de répondre à des questions d'appartenance d'un terme proposé par I à ce langage. Ainsi I prend en entrée un ensemble caractéristique d'un langage inconnu \mathcal{L} , il pose une série de questions d'appartenance de termes t qu'il choisit à un oracle O et l'oracle lui répond *vrai* ou *faux* suivant que t appartient à \mathcal{L} ou pas. Le processus s'arrête quand I n'a plus de questions à soumettre à O et I fournit un automate \mathcal{A} en sortie.

4.1 Réduction du problème

Classiquement, nous allons montrer que le problème d'identification d'un langage régulier à partir d'un ensemble caractéristique \mathcal{C} peut se réduire à un problème de partitionnement de \mathcal{C} . Soit \mathcal{L} un langage régulier d'arbres quelconque et \mathcal{C} un ensemble caractéristique de \mathcal{L} . Si \equiv est une relation d'équivalence sur \mathcal{C} alors \equiv définit un automate $\mathcal{A}_{\mathcal{C},\equiv}$ par : $\mathcal{A}_{\mathcal{C},\equiv} = \langle \mathcal{F}_{\mathcal{C},\equiv}, \mathcal{Q}_{\mathcal{C},\equiv}, \mathcal{Q}_{F,\mathcal{C},\equiv}, \rightarrow_{\mathcal{A}_{\mathcal{C},\equiv}} \rangle$, où

- $\mathcal{F}_{\mathcal{C},\equiv}$ est l'ensemble des symboles étiquetant les arbres de \mathcal{C} ,
- $\mathcal{Q}_{\mathcal{C},\equiv} = \{[t], t \in \mathcal{S}(\mathcal{C})\}$, où $[t]$ est la classe d'équivalence de t ,
- $\mathcal{Q}_{F,\mathcal{C},\equiv} = \{[t], t \in \mathcal{C}\}$,
- $\rightarrow_{\mathcal{A}_{\mathcal{C},\equiv}} = \{f([t_1], \dots, [t_n]) \rightarrow_{\mathcal{A}_{\mathcal{C},\equiv}} [f(t_1, \dots, t_n)], f(t_1, \dots, t_n) \in \mathcal{S}(\mathcal{C})\}$.

On notera que dans le cas général, $\mathcal{A}_{\mathcal{C},\equiv}$ n'a aucune raison particulière d'être déterministe.

Lemme 2

Pour tout langage \mathcal{L} , pour tout ensemble caractéristique \mathcal{C} de \mathcal{L} , l'automate $\mathcal{A}_{\mathcal{C},\equiv_{\mathcal{L},\mathcal{C}}}$ est un automate minimal de \mathcal{L} , où $\equiv_{\mathcal{L},\mathcal{C}}$ est la restriction de $\equiv_{\mathcal{L}}$ sur $\mathcal{S}(\mathcal{C})$.

Démonstration. Considérons un automate minimal $\mathcal{A} = \langle \mathcal{F}, \mathcal{Q}, \mathcal{Q}_F, \rightarrow_{\mathcal{A}} \rangle$ de \mathcal{L} et la fonction $\psi : \mathcal{Q} \mapsto \{[t], t \in \mathcal{S}(\mathcal{C})\}$ définie par $\psi(\delta_{\mathcal{A}}(t)) = [t]$, pour tout t dans $\mathcal{S}(\mathcal{C})$. Puisque \mathcal{C} est un ensemble caractéristique, à chaque état q de \mathcal{Q} , correspond un terme t dans $\mathcal{S}(\mathcal{C})$ tel que $\delta_{\mathcal{A}}(t) = q$ (chaque état q apparaît dans la partie droite d'une règle de $\rightarrow_{\mathcal{A}}$). Alors ψ est bien définie et on vérifie facilement que ψ est un renommage des états de \mathcal{A} .

D'après le lemme précédent, qui est une conséquence directe du théorème de Myhill-Nerode, le problème de construction d'un automate minimal pour un langage \mathcal{L} à partir

d'un ensemble caractéristique \mathcal{C} peut se réduire au problème de partitionnement de \mathcal{C} suivant la relation $\equiv_{\mathcal{L},\mathcal{C}}$. En effet le théorème de Myhill-Nerode caractérise tout langage régulier d'arbres (résultat équivalent sur les mots) en terme d'ensemble fini de classes d'équivalence. Chaque classe correspondant à un état d'un automate minimal pour le langage, reconstruire cet automate revient à reconstruire la relation d'équivalence sur un ensemble caractéristique.

4.2 Tableau d'observation

Comme dans (Angluin, 1987), l'information traitée par l'algorithme d'apprentissage sera stockée dans un tableau nommé *tableau d'observation*.

Soit \mathcal{L} un langage quelconque, \mathcal{E} un ensemble quelconque de termes fini et F un ensemble fini quelconque de contextes. Le *tableau d'observation* $T_{\mathcal{L}}(\mathcal{E}, F)$ relatif à \mathcal{L} , \mathcal{E} and F , est défini par :

- chaque ligne correspond à un élément de \mathcal{E} ,
- chaque colonne correspond à un élément de F ,
- pour chaque t dans \mathcal{E} et chaque $c[\diamond]$ dans F , $T_{\mathcal{L}}(t, c[\diamond]) = \text{vrai}$ si $c[t]$ est dans \mathcal{L} et $T_{\mathcal{L}}(t, c[\diamond]) = \text{faux}$ sinon.

Un tableau d'observation $T_{\mathcal{L}}(\mathcal{E}, F)$ définit une relation d'équivalence $\equiv_{T_{\mathcal{L}}(\mathcal{E}, F)}$ sur \mathcal{E} par : pour tout t et tout t' dans \mathcal{E} , $t \equiv_{T_{\mathcal{L}}(\mathcal{E}, F)} t'$ si et seulement si les mots *ligne*(t) et *ligne*(t') construits sur l'alphabet $\{\text{vrai}, \text{faux}\}^*$ et correspondant respectivement aux lignes de t et t' dans $T_{\mathcal{L}}(\mathcal{E}, F)$ sont égaux.

Remarque 3 Pour tous termes t et t' de \mathcal{E} , si $t \not\equiv_{T_{\mathcal{L}}(\mathcal{E}, F)} t'$, alors $t \not\equiv_{\mathcal{L}, F} t'$. En d'autres termes, la partition $\equiv_{\mathcal{L}, F}$ raffine toujours la partition $\equiv_{T_{\mathcal{L}}(\mathcal{E}, F)}$.

Un tableau d'observation $T_{\mathcal{L}}(\mathcal{E}, F)$ est dit *cohérent* lorsque pour tous termes $f(t_1, \dots, t_n)$ et $f(t'_1, \dots, t'_n)$ de $\mathcal{S}(\mathcal{E})$, si pour tout $1 \leq j \leq n$ on a :

$$t_j \equiv_{T_{\mathcal{L}}(\mathcal{E}, F)} t'_j$$

alors

$$f(t_1, \dots, t_n) \equiv_{T_{\mathcal{L}}(\mathcal{E}, F)} f(t'_1, \dots, t'_n).$$

Nous allons maintenant définir notre algorithme d'apprentissage. Celui-ci prend un ensemble caractéristique d'un langage inconnu en entrée, puis construit une série de tableaux d'observation à l'aide de l'oracle. L'algorithme termine lorsque le dernier tableau construit est cohérent.

4.3 L'algorithme IRTL

Initialisation : $F = C[\mathcal{C}]$;

Construction du tableau $T_{\mathcal{L}}(\mathcal{S}(\mathcal{C}), F)$ à l'aide de questions d'appartenance à O ;

TANTQUE

$T_{\mathcal{L}}(\mathcal{S}(\mathcal{C}), F)$ il existe $f(t_1, \dots, t_n)$ et $f(t'_1, \dots, t'_n)$ dans $\mathcal{S}(\mathcal{E})$ tels que $f(t_1, \dots, t_n) \not\equiv_{T_{\mathcal{L}}(\mathcal{S}(\mathcal{C}), F)} f(t'_1, \dots, t'_n)$ et pour tout $1 \leq i \leq n$,

$$t_i \equiv_{T_{\mathcal{L}}(\mathcal{S}(\mathcal{C}), \mathcal{F})} t'_i$$

FAIRE

Trouver un contexte $c[\diamond]$ dans \mathcal{F} tel que $c[f(t_1, \dots, t_n)] \in \mathcal{L}$

et $c[f(t'_1, \dots, t'_n)] \notin \mathcal{L}$;

$\mathcal{F} = \mathcal{F} \cup \{c[f(t_1, \dots, t_{i-1}, \diamond, t_{i+1}, \dots, t_n)], 1 \leq i \leq n\}$;

Construction de $T_{\mathcal{L}}(\mathcal{S}(\mathcal{C}), \mathcal{F})$ à l'aide de questions d'appartenance à \mathcal{O} ;

FINTANTQUE

Retourner l'automate $\mathcal{A}_{\equiv_{T_{\mathcal{L}}(\mathcal{S}(\mathcal{C}), \mathcal{F})}}$.

4.4 Preuve de correction

Nous allons montrer que la construction itérative du tableau et la vérification de la cohérence à chaque étape permet de raffiner la partition jusqu'à l'obtention de $\equiv_{\mathcal{L}, \mathcal{C}}$.

Lemme 4

Soient \mathcal{L} un langage régulier d'arbres, \mathcal{E} un ensemble de termes quelconque, \mathcal{C} un ensemble caractéristique pour \mathcal{L} , \mathcal{F} un ensemble de contextes contenant $\mathcal{C}[\mathcal{C}]$ et \mathcal{A} un automate minimal pour \mathcal{L} . Pour toute paire de termes t et t' de \mathcal{E} , si $t \equiv_{T_{\mathcal{L}}(\mathcal{E}, \mathcal{F})} t'$, alors pour toute règle de la forme $f(q_1, \dots, q_{i-1}, \delta_{\mathcal{A}}(t), q_{i+1}, \dots, q_n) \rightarrow_{\mathcal{A}} r$ dans \mathcal{A} , il existe également une règle de la forme $f(q_1, \dots, q_{i-1}, \delta_{\mathcal{A}}(t'), q_{i+1}, \dots, q_n) \rightarrow_{\mathcal{A}} r'$ dans \mathcal{A} .

Démonstration. Supposons $t \equiv_{T_{\mathcal{L}}(\mathcal{E}, \mathcal{F})} t'$ dans \mathcal{E} et $f(q_1, \dots, q_{i-1}, \delta_{\mathcal{A}}(t), q_{i+1}, \dots, q_n) \rightarrow_{\mathcal{A}} r$ une règle de \mathcal{A} . Puisque \mathcal{C} est caractéristique et \mathcal{F} contient $\mathcal{C}[\mathcal{C}]$, il existe un contexte de la forme $c[f(t_1, \dots, t_{i-1}, \diamond, t_{i+1}, \dots, t_n)]$ dans \mathcal{F} , où $\delta_{\mathcal{A}}(t_j) = q_j$, $\forall 1 \leq j \neq i \leq n$. D'après la définition de $T_{\mathcal{L}}(\mathcal{E}, \mathcal{F})$, on a $T_{\mathcal{L}}(t, c[f(t_1, \dots, t_{i-1}, \diamond, t_{i+1}, \dots, t_n)]) = \text{vrai}$. De plus, $t \equiv_{T_{\mathcal{L}}(\mathcal{E}, \mathcal{F})} t'$ implique que $T_{\mathcal{L}}(t', c[f(t_1, \dots, t_{i-1}, \diamond, t_{i+1}, \dots, t_n)]) = \text{vrai}$. Puisque $c[f(t_1, \dots, t_{i-1}, t', t_{i+1}, \dots, t_n)]$ est dans \mathcal{L} , il doit exister une règle $f(q_1, \dots, q_{i-1}, \delta_{\mathcal{A}}(t'), q_{i+1}, \dots, q_n) \rightarrow_{\mathcal{A}} r'$ dans \mathcal{A} .

Théorème 5

Soit \mathcal{L} un langage régulier d'arbres, \mathcal{C} un ensemble caractéristique de \mathcal{L} et \mathcal{F} un ensemble de contextes contenant $\mathcal{C}[\mathcal{C}]$. Si $\equiv_{T_{\mathcal{L}}(\mathcal{C}, \mathcal{F})}$ n'est pas égale à $\equiv_{\mathcal{L}, \mathcal{C}}$, il existe deux termes $f(t_1, \dots, t_{i-1}, t, t_{i+1}, \dots, t_n)$ et $f(t'_1, \dots, t'_{i-1}, t', t'_{i+1}, \dots, t'_n)$ dans $\mathcal{S}(\mathcal{C})$ tels que :

- $f(t_1, \dots, t_{i-1}, t, t_{i+1}, \dots, t_n) \not\equiv_{T_{\mathcal{L}}(\mathcal{C}, \mathcal{F})} f(t'_1, \dots, t'_{i-1}, t', t'_{i+1}, \dots, t'_n)$,
- $t_j \equiv_{\mathcal{L}, \mathcal{C}} t'_j$, pour tout $1 \leq j \neq i \leq n$,
- $t \equiv_{T_{\mathcal{L}}(\mathcal{C}, \mathcal{F})} t'$.

Démonstration. Si $\equiv_{T_{\mathcal{L}}(\mathcal{C}, \mathcal{F})}$ n'est pas égale à $\equiv_{\mathcal{L}, \mathcal{C}}$, il existe deux termes t et t' dans $\mathcal{S}(\mathcal{C})$ tels que $t \equiv_{T_{\mathcal{L}}(\mathcal{C}, \mathcal{F})} t'$ et $t \not\equiv_{\mathcal{L}, \mathcal{C}} t'$. $t \not\equiv_{\mathcal{L}, \mathcal{C}} t'$ implique l'existence d'un contexte $c[\diamond]$ tel que $c[t] \in \mathcal{L}$ et $c[t'] \notin \mathcal{L}$.

Définissons le contexte $c[\diamond]$ par :

- il existe deux termes s et s' dans $\mathcal{S}(\mathcal{C})$ tels que
 - $s \equiv_{T_{\mathcal{L}}(\mathcal{C}, \mathcal{F})} s'$

- $c[s] \in \mathcal{L}$ et $c[s'] \notin \mathcal{L}$
- $\text{depth}(c[\diamond]) = \min\{\text{depth}(c'[\diamond]), \exists u, u' \in \mathcal{S}(\mathcal{C}), u \equiv_{T_{\mathcal{L}}(\mathcal{C}, \mathcal{F})} u', c'[u] \in \mathcal{L}$ et $c'[u'] \notin \mathcal{L}\}$

Montrons tout d'abord que $\text{depth}(c[\diamond]) > 0$. Si $\text{depth}(c[\diamond]) = 0$, on a $c[\diamond] = \diamond$. Or \diamond est un contexte de \mathcal{F} , donc si $s \in \mathcal{L}$ et $s' \notin \mathcal{L}$, alors $s \not\equiv_{T_{\mathcal{L}}(\mathcal{C}, \mathcal{F})} s'$ ce qui montre une contradiction. Par conséquent $c[\diamond]$ peut s'écrire $d[\mathfrak{f}(s_1, \dots, s_{i-1}, \diamond, s_i, \dots, s_n)]$, où $s_1, \dots, s_{i-1}, s_{i+1}, \dots, s_n$ sont des termes et $d[\diamond]$ un contexte tel que $\text{depth}(c[\diamond]) = \text{depth}(d[\diamond]) + 1$. Puisque $d[\mathfrak{f}(s_1, \dots, s_{i-1}, s, s_{i+1}, \dots, s_n)]$ est dans \mathcal{L} , dans tout automate \mathcal{A} minimal pour \mathcal{L} , il existe une règle $\mathfrak{f}(q_1, \dots, q_{i-1}, \delta_{\mathcal{A}}(s), q_{i+1}, \dots, q_n) \rightarrow_{\mathcal{A}} r$, avec $q_j = \delta_{\mathcal{A}}(s_j)$, pour tout $1 \leq j \neq i \leq n$. D'après le lemme 4, il existe aussi une règle $\mathfrak{f}(q_1, \dots, q_{i-1}, \delta_{\mathcal{A}}(s'), q_{i+1}, \dots, q_n) \rightarrow_{\mathcal{A}} r'$ dans \mathcal{A} , avec $q_j = \delta_{\mathcal{A}}(s_j)$, pour tout $1 \leq j \neq i \leq n$. Si $r = r'$, alors on a $\delta_{\mathcal{A}}(\mathfrak{f}(s_1, \dots, s_{i-1}, s, s_{i+1}, \dots, s_n)) = \delta_{\mathcal{A}}(\mathfrak{f}(s_1, \dots, s_{i-1}, s', s_{i+1}, \dots, s_n))$ et $d[\mathfrak{f}(s_1, \dots, s_{i-1}, s, s_{i+1}, \dots, s_n)]$ est dans \mathcal{L} implique que $d[\mathfrak{f}(s_1, \dots, s_{i-1}, s', s_{i+1}, \dots, s_n)]$ est dans \mathcal{L} ce qui contredit l'hypothèse. \mathcal{C} est un ensemble caractéristique, alors aux règles $\mathfrak{f}(q_1, \dots, q_{i-1}, \delta_{\mathcal{A}}(s), q_{i+1}, \dots, q_n) \rightarrow_{\mathcal{A}} r$ et $\mathfrak{f}(q_1, \dots, q_{i-1}, \delta_{\mathcal{A}}(s'), q_{i+1}, \dots, q_n) \rightarrow_{\mathcal{A}} r'$ correspondent deux termes $\mathfrak{f}(t_1, \dots, t_{i-1}, t, t_{i+1}, \dots, t_n)$, $\mathfrak{f}(t'_1, \dots, t'_{i-1}, t', t'_{i+1}, \dots, t'_n)$ dans $\mathcal{S}(\mathcal{C})$, avec $\delta_{\mathcal{A}}(t_j) = \delta_{\mathcal{A}}(t'_j) = q_j$, pour tout $1 \leq j \neq i \leq n$ et $\delta_{\mathcal{A}}(\mathfrak{f}(t_1, \dots, t_{i-1}, t, t_{i+1}, \dots, t_n)) = r$ et $\delta_{\mathcal{A}}(\mathfrak{f}(t'_1, \dots, t'_{i-1}, t', t'_{i+1}, \dots, t'_n)) = r'$. Finalement supposons que $\mathfrak{f}(t_1, \dots, t_{i-1}, t, t_{i+1}, \dots, t_n) \equiv_{T_{\mathcal{L}}(\mathcal{C}, \mathcal{F})} \mathfrak{f}(t'_1, \dots, t'_{i-1}, t', t'_{i+1}, \dots, t'_n)$ et montrons que le contexte $d[\diamond]$ contredit la minimalité de $\text{depth}(c[\diamond])$. En effet :

- $\mathfrak{f}(t_1, \dots, t_{i-1}, t, t_{i+1}, \dots, t_n) \equiv_{T_{\mathcal{L}}(\mathcal{C}, \mathcal{F})} \mathfrak{f}(t'_1, \dots, t'_{i-1}, t', t'_{i+1}, \dots, t'_n)$
- $d[\mathfrak{f}(t_1, \dots, t_{i-1}, t, t_{i+1}, \dots, t_n)] \in \mathcal{L}$ et $d[\mathfrak{f}(t'_1, \dots, t'_{i-1}, t', t'_{i+1}, \dots, t'_n)] \notin \mathcal{L}$
- $\text{depth}(c[\diamond]) > \text{depth}(d[\diamond])$

On a donc bien $\mathfrak{f}(t_1, \dots, t_{i-1}, t, t_{i+1}, \dots, t_n) \not\equiv_{T_{\mathcal{L}}(\mathcal{C}, \mathcal{F})} \mathfrak{f}(t'_1, \dots, t'_{i-1}, t', t'_{i+1}, \dots, t'_n)$ ce qui permet de conclure.

Remarque 6 Si $\equiv_{T_{\mathcal{L}}(\mathcal{C}, \mathcal{F})}$ n'est pas égale à $\equiv_{\mathcal{L}, \mathcal{C}}$, le théorème précédent implique que $T_{\mathcal{L}}(\mathcal{C}, \mathcal{F})$ n'est pas cohérent. En d'autres termes, si l'automate $T_{\mathcal{L}}(\mathcal{C}, \mathcal{F})$ est cohérent, le langage est appris et l'algorithme s'arrête.

Corollaire 7

L'algorithme *IRTL* apprend la classe des langages réguliers d'arbres.

Démonstration. L'algorithme commence par la construction de $T_{\mathcal{L}}(\mathcal{E}, \mathcal{F}) = T_{\mathcal{L}}(\mathcal{S}(\mathcal{C}), \mathcal{C}[\mathcal{C}])$ puis entre dans une boucle. D'après le théorème 5, si le programme sort de cette boucle, le tableau est cohérent et l'on a $\equiv_{T_{\mathcal{L}}(\mathcal{C}, \mathcal{F})}$ est égale à $\equiv_{\mathcal{L}, \mathcal{C}}$, et d'après le lemme 2, l'automate $\mathcal{A}_{\mathcal{C}, \equiv_{\mathcal{L}, \mathcal{C}}}$ est un automate minimal pour \mathcal{L} . Montrons alors que l'algorithme termine. La boucle collecte tous les couples de termes $\mathfrak{f}(t_1, \dots, t_n)$ et $\mathfrak{f}(t'_1, \dots, t'_n)$ dans \mathcal{C} , tels qu'il existe un contexte $c[\diamond]$ dans \mathcal{F} avec pour tout $1 \leq i \leq n$, $t_i \equiv_{T_{\mathcal{L}}(\mathcal{C}, \mathcal{F})} t'_i$, $c[\mathfrak{f}(t_1, \dots, t_n)] \in \mathcal{L}$ et $c[\mathfrak{f}(t'_1, \dots, t'_n)] \notin \mathcal{L}$. D'après le théorème 5, parmi ces couples, il en est un $\mathfrak{f}(t_1, \dots, t_{i-1}, t, t_{i+1}, \dots, t_n)$ et $\mathfrak{f}(t'_1, \dots, t'_{i-1}, t', t'_{i+1}, \dots, t'_n)$, tel que $t_j \equiv_{\mathcal{L}, \mathcal{C}} t'_j$, pour tout $1 \leq j \neq i \leq n$. $c[\mathfrak{f}(t'_1, \dots, t'_{i-1}, t', t'_{i+1}, \dots, t'_n)] \notin \mathcal{L}$ est alors équivalent à $c[\mathfrak{f}(t_1, \dots, t_{i-1}, t, t_{i+1}, \dots, t_n)] \notin \mathcal{L}$. Par définition d'*IRTL*, lorsque la boucle se termine, le contexte $c[\mathfrak{f}(t_1, \dots, t_{i-1}, \diamond, t_{i+1}, \dots, t_n)]$ a été ajouté à \mathcal{F} , ce qui

différencie les colonnes de t et de t' dans $T_{\mathcal{L}}(\mathcal{C}, F \cup \{c[\mathfrak{f}(t_1, \dots, t_{i-1}, \diamond, t_{i+1}, \dots, t_n)]\})$. Comme ces colonnes étaient précédemment semblables et comme deux colonnes différentes dans $T_{\mathcal{L}}(\mathcal{C}, F)$ le restent dans $T_{\mathcal{L}}(\mathcal{C}, F \cup \{c[\mathfrak{f}(t_1, \dots, t_{i-1}, \diamond, t_{i+1}, \dots, t_n)]\})$, à chaque passage dans la boucle, la relation d'équivalence $\equiv_{T_{\mathcal{L}}(\mathcal{C}, F)}$ se raffine (augmentation du nombre de classes). Or $\equiv_{\mathcal{L}, \mathcal{C}}$ est toujours plus fine que $\equiv_{T_{\mathcal{L}}(\mathcal{C}, F)}$, et comme le nombre de classes de $\equiv_{\mathcal{L}, \mathcal{C}}$ est fini (égal au nombre d'état d'un automate minimal pour \mathcal{L}), l'algorithme *IRTL* ne peut exécuter cette boucle qu'un nombre fini de fois.

4.5 Complexité

La complexité en temps de l'algorithme dépend de la taille des données d'entrée ainsi que de la taille de l'automate à déterminer. Le premier tableau construit est de taille n^2 , où n est le nombre total de nœuds des arbres de l'ensemble caractéristique \mathcal{C} donné. En effet, chaque ligne du tableau correspond à un élément de $\mathcal{S}(\mathcal{C})$ et chaque colonne à un élément de $\mathcal{C}[\mathcal{C}]$ et $|\mathcal{S}(\mathcal{C})| = |\mathcal{C}[\mathcal{C}]| = n$. Ensuite, le nombre de lignes n'évolue plus et c'est le nombre de colonnes qui augmente de la manière suivante : chaque fois qu'une itération ajoute k contextes à F , deux lignes du tableau initialement égales sont différenciées. Or le nombre de lignes différentes correspond au nombre de classes d'équivalence et est toujours plus petit que le nombre m d'états de l'automate minimal. De plus, le nombre k de contextes ajoutés à chaque itération est inférieur à n . Le nombre de contextes ajoutés est donc inférieur à $m \times n$ et le nombre final de colonnes est inférieur à $m \times n + n$, soit un tableau final de $n \times (m \times n + n)$ cases remplies par des questions d'appartenance. La complexité est donc de $n \times (m \times n + n)$ multiplié par une constante correspondant au temps de réponse de l'oracle. A cela nous devons ajouter le temps de vérification du tableau qui est lui aussi polynomial dans sa taille.

4.6 Comparaison

Il est ici intéressant de comparer notre algorithme à celui de Angluin dans (Angluin, 1981). Cet algorithme identifie la classe des langages réguliers (de mots) à partir d'un ensemble *live-complete* et de questions d'appartenance. Un ensemble *live-complete* est un ensemble fini d'exemples positifs vérifiant la propriété que, pour produire cet exemple, chaque état d'un automate minimal a été utilisé. Cette définition est moins restrictive que celle d'ensemble caractéristique puisque tout ensemble caractéristique est *live-complete*. Il est alors pertinent de voir si l'algorithme fonctionnant à partir d'un ensemble *live-complete* est adaptable à la classe des langages réguliers d'arbres. Il apparaît que l'utilisation de chaque transition dans les exemples est nécessaire à l'apprentissage. Pour obtenir la propriété d'un ensemble caractéristique à partir d'un ensemble *live-complete*, Angluin commence par considérer chaque mot $\omega\alpha$ formé d'un préfixe ω d'un mot appartenant à son ensemble *live-complete* et d'un suffixe α qui est une lettre de l'alphabet. Ainsi, puisqu'à chaque état d'un automate minimal correspond au moins un préfixe ω d'un mot de l'ensemble *live-complete*, alors l'ensemble $\{\omega\alpha\}$ obtenu est caractéristique. Pour appliquer cette technique sur les arbres, il nous faut trouver pour chaque sous-arbre d'un ensemble *live-complete* donné lc (avec une définition équiva-

lente à celle des langages de mots), tout contexte susceptible de correspondre à une transition de l'automate minimal cherché. Ainsi pour chaque sous-arbre t il nous faut construire tous les termes de la forme $f(t_1, \dots, t_i, t, t_{i+1}, \dots, t_n)$ pour les ajouter aux lignes du tableau, f étant un élément de l'alphabet de taille k d'arité $n(k)$ choisi parmi tous les k possibles et les t_j pour $1 \leq j \leq n$ étant choisis parmi tous les sous-arbres de lc . Cette opération nécessiterait donc $\sum_k n(k)^{|lc|}$ étapes de calcul, ce qui en accroît très nettement la complexité. Notre choix de considérer les ensembles caractéristiques permet donc d'obtenir un algorithme efficace avec une propriété un peu plus restrictive sur l'ensemble d'exemples positifs donné.

5 Un exemple d'exécution

Soit le langage \mathcal{L} défini par l'automate minimal suivant \mathcal{A}^3 :

$$\begin{array}{llll} a(q_1) \rightarrow_{\mathcal{A}} q_0 & d(q_3, q_5) \rightarrow_{\mathcal{A}} q_1 & e(q_7) \rightarrow_{\mathcal{A}} q_3 & g \rightarrow_{\mathcal{A}} q_7 \\ a(q_2) \rightarrow_{\mathcal{A}} q_0 & d(q_4, q_5) \rightarrow_{\mathcal{A}} q_1 & e(q_8) \rightarrow_{\mathcal{A}} q_4 & h \rightarrow_{\mathcal{A}} q_8 \\ b(q_1) \rightarrow_{\mathcal{A}} q_0 & d(q_3, q_6) \rightarrow_{\mathcal{A}} q_1 & f(q_9) \rightarrow_{\mathcal{A}} q_5 & i \rightarrow_{\mathcal{A}} q_9 \\ c \rightarrow_{\mathcal{A}} q_1 & d(q_4, q_6) \rightarrow_{\mathcal{A}} q_2 & f(q_{10}) \rightarrow_{\mathcal{A}} q_6 & j \rightarrow_{\mathcal{A}} q_{10} \end{array}$$

où q_0 est l'unique état final. D'après cette définition, on a :

$$\begin{aligned} \mathcal{L} = \{ & a(c), b(c), a(d(e(g), f(i))), a(d(e(h), f(i))), a(d(e(g), f(j))), \\ & a(d(e(h), f(j))), b(d(e(g), f(i))), b(d(e(h), f(i))), b(d(e(g), f(j))) \}. \end{aligned}$$

Supposons alors donné l'ensemble caractéristique suivant :

$$C = \{ b(c), a(d(e(g), f(i))), a(d(e(h), f(i))), a(d(e(g), f(j))), a(d(e(h), f(j))) \}.$$

IRTL commence par la construction de $F = C[C]$ et de $T_{\mathcal{L}}(\mathcal{S}(C), F)$ à l'aide de questions d'appartenance (tableau 1).

IRTL remarque les trois couples de termes :

$$d(e(h), f(i)) \text{ et } d(e(h), f(j))$$

$$d(e(g), f(i)) \text{ et } d(e(h), f(j))$$

$$d(e(g), f(j)) \text{ et } d(e(h), f(j))$$

En effet,

$$e(g) \equiv_{T_{\mathcal{L}}(\mathcal{S}(C), F)} e(h)$$

et

$$f(i) \equiv_{T_{\mathcal{L}}(\mathcal{S}(C), F)} f(j)$$

3. D'après le lemme 4, la partition définie par un tableau d'observation n'est pas correcte si l'automate minimal du langage appris contient des règles particulières : des règles identiques à un état près dans leur partie gauche. Si un automate ne contient pas de telles règles, le premier tableau construit est cohérent et l'algorithme s'arrête. Afin d'illustrer au mieux une exécution, nous proposons donc un exemple de langage dont l'automate minimal possède ce type de règles.

TAB. 1 – $T_L(S(C), C[C])$

$S(C)$	F	\diamond	$a(d(e(\diamond), f(i)))$	$a(d(e(g), f(\diamond)))$	$a(d(\diamond, f(i)))$	$a(d(e(g), \diamond))$	$a(\diamond)$	$a(d(e(h), f(\diamond)))$	$a(d(e(h), \diamond))$	$a(d(e(\diamond), f(j)))$	$a(d(\diamond, f(j)))$	$b(\diamond)$
$a(d(e(g), f(i)))$	<i>vrai</i>	<i>faux</i>	<i>faux</i>	<i>faux</i>	<i>faux</i>	<i>faux</i>	<i>faux</i>	<i>faux</i>	<i>faux</i>	<i>faux</i>	<i>faux</i>	<i>faux</i>
$e(g)$	<i>faux</i>	<i>faux</i>	<i>faux</i>	<i>faux</i>	<i>vrai</i>	<i>faux</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>vrai</i>	<i>faut</i>
$f(i)$	<i>faux</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>
g	<i>faut</i>	<i>vrai</i>	<i>faut</i>	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>
i	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>
$a(d(e(h), f(i)))$	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>
$d(e(h), f(i))$	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>vrai</i>
$e(h)$	<i>faut</i>	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>vrai</i>	<i>faut</i>
h	<i>faut</i>	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>
$a(d(e(g), f(j)))$	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>
$d(e(g), f(j))$	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>vrai</i>
$f(j)$	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>
j	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>
$a(d(e(h), f(i)))$	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>
$d(e(h), f(j))$	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>
$b(c)$	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>
c	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>vrai</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>faut</i>	<i>vrai</i>

mais

$$b(d(e(g), f(i))) \in \mathcal{L},$$

$$b(d(e(g), f(j))) \in \mathcal{L},$$

$$b(d(e(h), f(i))) \in \mathcal{L}$$

et

$$b(d(e(h), f(j))) \notin \mathcal{L}.$$

IRTL ajoute les contextes

$$b(d(e(g), \diamond)),$$

$$b(d(e(h), \diamond)),$$

$$b(d(\diamond, f(i)))$$

et

$$b(d(\diamond, f(j)))$$

à F et complète le tableau d'observation (tableau 2).

IRTL remarque maintenant que

$$b(d(e(h), f(i))) \in \mathcal{L},$$

$$b(d(e(h), f(j))) \notin \mathcal{L}$$

et

$$i \equiv_{T_{\mathcal{L}}(S(\mathcal{C}), F)} j$$

et que

$$b(d(e(g), f(j))) \in \mathcal{L},$$

$$b(d(e(h), f(j))) \notin \mathcal{L}$$

et

$$g \equiv_{T_{\mathcal{L}}(S(\mathcal{C}), F)} h.$$

Il ajoute donc les nouveaux contextes

$$b(d(e(\diamond), f(j)))$$

et

$$b(d(e(h), f(\diamond)))$$

à F et complète une dernière fois le tableau (tableau 3).

Le tableau est maintenant cohérent et l'algorithme donne l'automate $\mathcal{A}_{\equiv_{T_{\mathcal{L}}(S(\mathcal{C}), F)}}$ en sortie. On vérifie aisément que cet automate est un renommage de \mathcal{A} et donc un automate minimal pour le langage appris.

Références

- ANGLUIN D. (1981). A note on the number of queries needed to identify regular languages. *Information and Control*, **51**, 76–87.
- ANGLUIN D. (1987). Learning regular sets from queries and counter examples. *Information and Control*, **75**, 87–106.
- ANGLUIN D. (1988). Queries and concept learning. *Machine learning*, **2**, 319–342.
- CARRASCO R., ONCINA J. & CALERA-RUBIO J. (1998). Stochastic inference of regular tree languages. In *Grammatical inference: Algorithms and Applications (ICGI)*, volume 1433 of *Lecture Notes in Computer Science*, p. 187–198.
- CHRISTOPHE A. (2000). L'apprentissage du langage. In *Université de tous les savoirs*, volume 2, p. 41–51: Odile Jacob.
- COMON H., DAUCHET M., GILLERON R., JACQUEMARD F., LUGIEZ D., TISON S. & TOMMASI M. (1997). Tree automata techniques and applications. Disponible à : <http://www.grappa.univ-lille3.fr/tata>.
- PINKER S. (1994). *The language instinct*. Harper.
- SAKAKIBARA Y. (1988). Learning context-free grammars from structural data in polynomial time. In *Proceedings of the first annual workshop on Computational learning theory*, p. 330–344: Morgan Kaufmann Publishers Inc.