

APPRENTISSAGE DES LANGAGES RÉGULIERS D'ARBRES ET APPLICATIONS

Jérôme Besombes, Jean-Yves Marion *

Résumé

Nous nous intéressons à l'apprentissage des langages réguliers d'arbres à partir d'exemples positifs. Notre modèle d'apprentissage est l'identification à la limite de Gold et les exemples sont des arbres. Nous démontrons que les langages réguliers d'arbres réversibles sont identifiables. Nous présentons trois conséquences de ce résultat. (1) Nous montrons que les grammaires algébriques *reset-free* sont identifiables à partir des arbres de dérivation. (2) Nous montrons que les langages d'arbres de dépendances sont identifiables. (3) Nous donnons une nouvelle démonstration de l'apprentissage des grammaires catégorielles rigides.

1 INTRODUCTION

Des linguistes comme Chomsky et Pinker (Pinker, 1994), ont montré que l'acquisition des langages naturels est basée sur l'analyse de structures de phrases correctes. Ceci illustre l'importance de l'identification à la limite à partir d'exemples positifs de Gold. Notre objectif est de contribuer à la formalisation et à la compréhension du processus d'acquisition du langage. Dans ce but, nous nous basons sur les travaux d'Angluin et Sakakibara concernant l'apprentissage. Nous montrons que la classe des langages réguliers d'arbres réversibles est identifiable. Tous ces résultats issus de motivations différentes, partagent un point de vue commun ; ils supposent qu'un langage est produit par un mécanisme fini et réversible et ainsi, qu'il doit être possible d'inverser une production. Angluin (Angluin, 1982) a montré que les langages de mots reconnus par des automates d'états finis réversibles sont identifiables, puis, Sakakibara (Sakakibara, 1992) a proposé une notion de grammaires algébriques réversibles. Pour ce dernier, les données ne sont plus des mots mais des structures d'arbres correspondant à des squelettes d'arbres de dérivation. Ainsi, le problème est réduit à l'identification d'une sous-classe de langages réguliers d'arbres réversibles et les exemples sont des éléments du langage. Nous avons évolué de l'apprentissage de grammaires à l'apprentissage de langages, avec l'avantage d'une plus grande simplicité conceptuelle. Notre résultat peut s'appliquer à différents types de grammaires :

1. les grammaires de dépendances réversibles sont identifiables à partir d'exemples positifs. Voir le théorème 30,

* LORIA, BP 239. 54506 Vandoeuvre-lès-Nancy

2. les grammaires catégorielles rigides sont identifiables à partir des squelettes d'arbres de dérivation. Voir le théorème 30 et la nouvelle démonstration de ce résultat établi par Kanazawa (Kanazawa, 1998).

D'un point de vue pratique, il apparaît que nous avons accès aux éléments d'un langage, mais pas au mécanisme de production de ces éléments. Les langages d'arbres réguliers représentent un cadre formel général pour des domaines aussi variés que les documents XML ou le langage naturel ce qui nous intéresse particulièrement. L'idée de départ étant que les données sont des arbres annotés et que ces annotations représentent des informations sémantiques et/ou syntaxiques qui, une fois prises en compte, doivent faciliter l'apprentissage. De telles données se retrouvent dans un grand nombre de grammaires formelles de linguistique computationnelle. Mentionnons par exemple les LFG, les TAG ou encore les grammaires catégorielles.

Travaux reliés

Un certain nombre d'articles concernent l'apprentissage de grammaires d'arbres. Nous avons déjà mentionné les travaux d'Angluin et de Sakakibara (Angluin, 1982; Sakakibara, 1992); d'autres (Edwards *et al.*, 1976; Levine, 1981; Kamata, 1988; Fukuda & Kamata, 1984) se basent sur l'idée d'inférence *k-tail* des langages réguliers de mots de Biermann et Feldman (Biermann & Feldman, 1972) puis Knuutila et Steinby (Knuutila & Steinby, 1994). Récemment, Fernau (Fernau, 2001) a utilisé ce cadre pour les grammaires XML. Carrasco et Oncina dans (Edwards *et al.*, 1976) ont montré que les langages réguliers d'arbres sont identifiables à partir d'exemples probabilistes. Le lecteur intéressé par l'apprentissage des langages algébriques peut consulter les études (Lee, 1996; Sakakibara, 1997; Mäkinen, 1997).

2 LANGAGES RÉGULIERS D'ARBRES

2.1 Les arbres comme des termes

Le livre (Comon *et al.*, 1997) donne les bases de la théorie des langages réguliers d'arbres. Un alphabet ordonné est un couple $(\mathcal{F}, \text{arite})$ où \mathcal{F} est un ensemble fini de symboles et *arite* est une fonction de \mathcal{F} dans \mathbb{N} qui indique l'arité des symboles. Pour tout ensemble de variables \mathcal{X} , les termes sont définis par : un symbole d'arité 0 est un terme, une variable de \mathcal{X} est un terme et si f est un symbole d'arité n 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 $\text{sub}(t)$ défini par : $t \in \text{sub}(t)$ et si $f(t_1, \dots, t_n) \in \text{sub}(t)$ alors t_1, \dots, t_n sont dans $\text{sub}(t)$. Pour tout ensemble de termes \mathcal{E} , $\text{sub}(\mathcal{E})$ est la réunion des ensembles des sous-termes des éléments de \mathcal{E} . 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. Ainsi, la substitution de \diamond par un terme u se note $c[u]$. On pourra utiliser cette notation pour indiquer que le terme u est un sous-terme de $c[u]$.

2.2 Les automates d'arbres

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

Un automate d'arbres d'états finis est *déterministe* (AEFD) si il n'existe pas deux transitions distinctes possédant des parties gauches identiques. Pour ces automates, $\xrightarrow{\mathcal{A}}$ est une fonction.

La dérivation simple $\xrightarrow{\mathcal{A}}$ est définie par $t \xrightarrow{\mathcal{A}} u$ si et seulement si il existe une transition $f(q_1, \dots, q_n) \xrightarrow{\mathcal{A}} q$ telle que $t = v[f(q_1, \dots, q_n)]$ et $u = v[q]$. On notera que $\xrightarrow{\mathcal{A}} \subseteq \mathcal{T}(\mathcal{F}, \mathcal{Q}) \times \mathcal{T}(\mathcal{F}, \mathcal{Q})$, et que les états de \mathcal{Q} sont des symboles d'arité nulle.

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

2.3 Langages réguliers d'arbres réversibles

Définition 1

Un AEFD \mathcal{A} est *réversible* si et seulement si

1. il n'existe pas deux règles distinctes possédant des parties gauches qui ne diffèrent que d'un unique symbole. C'est à dire qu'il n'existe pas deux transitions $f(p_1, \dots, p_n, q, p_{n+1}, \dots, p_m) \xrightarrow{\mathcal{A}} p$ et $f(p_1, \dots, p_n, q', p_{n+1}, \dots, p_m) \xrightarrow{\mathcal{A}} p$ avec $q \neq q'$,
2. \mathcal{A} possède un unique état final.

Un langage d'arbres est *réversible* si et seulement si il est reconnu par un AEFD réversible.

On remarquera que dans ce cas, un symbole d'arité n et $n - 1$ états définissent au plus une transition possible.

Exemple 2

L'AEFD $\mathcal{A}_1 = \langle \{Xor, Not, true, false\}, \{q_0, q_1\}, q_1, \xrightarrow{\mathcal{A}_1} \rangle$ suivant est réversible.
 $\xrightarrow{\mathcal{A}_1} = \{ true \xrightarrow{\mathcal{A}_1} q_1, false \xrightarrow{\mathcal{A}_1} q_0, Not(q_0) \xrightarrow{\mathcal{A}_1} q_1, Not(q_1) \xrightarrow{\mathcal{A}_1} q_0, \}$

$Xor(q_0, q_0) \xrightarrow{\mathcal{A}_1} q_0, Xor(q_0, q_1) \xrightarrow{\mathcal{A}_1} q_1, Xor(q_1, q_0) \xrightarrow{\mathcal{A}_1} q_1, Xor(q_1, q_1) \xrightarrow{\mathcal{A}_1} q_0$.

Le langage régulier d'arbres réversible $\mathcal{L}_{\mathcal{A}_1}$ est l'ensemble des expressions booléennes vraies.

2.4 Les grammaires régulières d'arbres réversibles

Une *grammaire régulière d'arbres*, notée GRA, est un quadruplet $\Gamma = \langle \mathcal{F}, \mathcal{X}, \xrightarrow{\Gamma}, S \rangle$ où $S \in \mathcal{X}$ est la variable de départ. Toutes les productions sont de la forme $X \xrightarrow{\Gamma} t$ où X est une variable de \mathcal{X} appelée tête et t est un terme de $\mathcal{T}(\mathcal{F}, \mathcal{X})$. Ainsi, les productions de la forme $X \xrightarrow{\Gamma} Y$, où Y est une variable de \mathcal{X} ne sont pas autorisées.

Définissons $t \xrightarrow{\Gamma} u$ par : $t \xrightarrow{\Gamma} u$ si et seulement si il existe une production $X \xrightarrow{\Gamma} v$ telle que $u = t[X \leftarrow v]$, où $t[X \leftarrow v]$ est le terme obtenu par substitution de X dans t par v .

La relation de dérivation $\xrightarrow{\Gamma}^*$ est la clôture réflexive et transitive de $\xrightarrow{\Gamma}$. Le langage produit par Γ est l'ensemble $\mathcal{L}_{\Gamma} = \{t \in \mathcal{T}(\mathcal{F}) : S \xrightarrow{\Gamma}^* t\}$.

Définition 3

Une grammaire Γ est réversible si et seulement si :

1. il n'existe pas deux productions $X \xrightarrow{\Gamma} c[Y]$ et $X \xrightarrow{\Gamma} c[Z]$ avec la même tête X et telles que $Y \neq Z$,
2. il n'existe pas deux productions $X \xrightarrow{\Gamma} t$ et $Y \xrightarrow{\Gamma} t$ avec la même partie droite et telles que $X \neq Y$.

Définition 4

Une GRA Γ est normale si chaque production est de la forme $X \xrightarrow{\Gamma} a$, ou de la forme $X \xrightarrow{\Gamma} f(X_1, \dots, X_n)$.

Théorème 5

Un langage \mathcal{L} est réversible si et seulement si il est produit par une GRA réversible et normale.

Exemple 6

Soit la GRA $G_1 = (\{Xor, Not, true, false\}, \{q_0, q_1\}, q_1, \xrightarrow{\mathcal{A}_1}, \xrightarrow{\mathcal{A}_1} = \{q_1 \xrightarrow{G_1} true, q_0 \xrightarrow{G_1} false, q_1 \xrightarrow{G_1} Not(q_0), q_1 \xrightarrow{G_1} Xor(q_0, q_1), q_1 \xrightarrow{G_1} Xor(q_1, q_0), q_0 \xrightarrow{G_1} Not(q_1), q_0 \xrightarrow{G_1} Xor(q_0, q_0), q_0 \xrightarrow{G_1} Xor(q_1, q_1)\}$. Le langage $\mathcal{L}_{\mathcal{A}_1}$ est produit par G_1 qui est réversible et normale.

2.5 Apprentissage des langages réguliers d'arbres réversibles

Une *présentation positive* d'un langage \mathcal{L} est une suite t_1, \dots, t_n, \dots qui énumère les éléments de \mathcal{L} . Soit Ω une classe d'automates (ou de grammaires) donnée. Un algorithme d'apprentissage A prend pour entrée, une partie finie t_1, \dots, t_n d'une présentation positive d'un langage \mathcal{L} et propose un automate $A(t_1, \dots, t_n)$ en sortie. Un algorithme d'inférence A converge vers \mathcal{L} si il existe une borne N à partir de laquelle le langage produit par $A(t_1, \dots, t_n)$, $n > N$ est toujours égal à \mathcal{L} . Une classe de langages \mathbb{L} est *identifiable* si et seulement si il existe un algorithme d'apprentissage A tel que pour tout langage \mathcal{L} de la classe et toute présentation positive de \mathcal{L} , A converge vers \mathcal{L} .

Théorème 7

La classe des langages réguliers d'arbres réversibles est identifiable.

Nous allons démontrer ce théorème dans le prochain chapitre.

3 APPRENTISSAGE DES LANGAGES RÉGULIERS D'ARBRES RÉVERSIBLES

3.1 Une caractérisation algébrique des langages réguliers d'arbres réversibles

Une relation d'équivalence \equiv est *close par contexte* si pour tous termes t et u dans $\mathcal{T}(\mathcal{F})$, $t \equiv u$ si pour tout contexte $c[\diamond]$, $c[t] \equiv c[u]$. Une congruence \equiv sur $\mathcal{T}(\mathcal{F})$ est une relation d'équivalence close par contexte.

Pour un AEFD donné \mathcal{A} , la relation d'équivalence $\equiv_{\mathcal{A}}$ est définie par : $t \equiv_{\mathcal{A}} u$ si $t \xrightarrow{\mathcal{A}}^* q$ et $u \xrightarrow{\mathcal{A}}^* q$. On voit facilement que $\equiv_{\mathcal{A}}$ est close par contexte et donc est une congruence.

Lemme 8

Soit \mathcal{A} un AEFD réversible. Pour tout contexte $c[\diamond]$ et pour tous termes t et u tels que $c[t]$ et $c[u]$ soient dans $\mathcal{L}_{\mathcal{A}}$, on a $t \equiv_{\mathcal{A}} u$.

Démonstration. Par induction sur la taille du contexte.

Cas de base : Supposons que $c[\diamond] = \diamond$. Par hypothèse, t et u sont dans $\mathcal{L}_{\mathcal{A}}$. Puisque \mathcal{A} est réversible, il existe un unique état final q_f et donc, on a $t \xrightarrow{\mathcal{A}}^* q_f$ et $u \xrightarrow{\mathcal{A}}^* q_f$.

D'où $t \equiv_{\mathcal{A}} u$.

Étape d'induction : Supposons que $c[\diamond] = c'[\mathfrak{f}(t_1, \dots, t_n, \diamond, t_{n+1}, \dots, t_m)]$. Par hypothèse, il existe un état q_i avec $t_i \xrightarrow{\mathcal{A}}^* q_i$. Puisque \mathcal{A} est réversible, les états q_1, \dots, q_m déterminent une unique transition dont la partie gauche est $\mathfrak{f}(q_1, \dots, q_n, q, q_{n+1}, \dots, q_m)$. Ceci implique que $t \xrightarrow{\mathcal{A}}^* q$ et $u \xrightarrow{\mathcal{A}}^* q$, donc $t \equiv_{\mathcal{A}} u$.

□

Pour tout langage d'arbres régulier \mathcal{L} , la congruence $\equiv_{\mathcal{L}}$ est définie par : $t \equiv_{\mathcal{L}} u$ si pour tout contexte $c[\diamond], c[t] \in \mathcal{L}$ si et seulement si $c[u] \in \mathcal{L}$. D'après le théorème de Myhill-Nerode, pour tout automate \mathcal{A} qui reconnaît \mathcal{L} , le nombre de classes de $\equiv_{\mathcal{L}}$ est inférieur ou égal au nombre de classes de $\equiv_{\mathcal{A}}$.

En conséquence, l'AEFD minimal (relativement au nombre d'états) $\mathcal{A} = \langle \mathcal{F}, \mathcal{Q}, \mathcal{Q}_F, \xrightarrow{\mathcal{L}} \rangle$ qui reconnaît \mathcal{L} est défini comme suit. Soit $[t]$ la classe d'équivalence de t pour $\equiv_{\mathcal{L}}$. Posons alors $\mathcal{Q} = \{[t] : t \in \text{sub}(\mathcal{L})\}$ et $\mathcal{Q}_F = \{[t] : t \in \mathcal{L}\}$. Pour chaque état $[t_1], \dots, [t_n]$ et pour chaque symbole n-aire $f \in \mathcal{F}$, il existe une transition $f([t_1], \dots, [t_n]) \xrightarrow{\mathcal{L}} [f(t_1, \dots, t_n)]$.

Théorème 9

Un langage régulier d'arbres \mathcal{L} est réversible si et seulement si pour tout contexte $c[\diamond]$ et pour tous termes t et u , si $c[t]$ et $c[u]$ sont dans \mathcal{L} alors on a $t \equiv_{\mathcal{L}} u$.

Démonstration. Il existe un AEFD réversible \mathcal{A} qui reconnaît \mathcal{L} . D'après le théorème de Myhill-Nerode, si $t \equiv_{\mathcal{A}} u$ alors $t \equiv_{\mathcal{L}} u$. Le lemme 8 permet de conclure.

Inversement, considérons le AEFD minimal $\mathcal{A}_{\mathcal{L}}$ qui reconnaît \mathcal{L} . $\mathcal{A}_{\mathcal{L}}$ possède un unique état final car par hypothèse, tous les arbres de \mathcal{L} appartiennent à la même classe d'équivalence.

Supposons maintenant l'existence de deux transitions de la forme $f(p_1, \dots, p_n, q, p_{n+1}, \dots, p_m) \xrightarrow{\mathcal{L}} q''$ et $f(p_1, \dots, p_n, q', p_{n+1}, \dots, p_m) \xrightarrow{\mathcal{L}} q''$. Puisque $\mathcal{A}_{\mathcal{L}}$ est minimal, il existe des termes t_1, \dots, t_m tels que $t_i \xrightarrow{\mathcal{L}}^* p_i$. Ceci permet de considérer le contexte $c[\diamond] = f(t_1, \dots, t_n, \diamond, t_{n+1}, \dots, t_m)$ tel que $c[q] \xrightarrow{\mathcal{L}}^* q_f$ et $c[q'] \xrightarrow{\mathcal{L}}^* q_f$, où q est l'état final de \mathcal{A} . Alors, par hypothèse, on a $q = q'$ et donc $\mathcal{A}_{\mathcal{L}}$ est réversible. On conclut que \mathcal{L} est réversible.

□

La démonstration ci-dessus induit les conséquences suivantes.

Corollaire 10

Un langage régulier d'arbres est réversible si et seulement si l'AEFD minimal \mathcal{A} est réversible.

Exemple 11

Soit $\mathcal{L} = \{f(g^n(a)) : n \geq 0\} \cup \{g(f^n(a)) : n \geq 0\}$. Le langage régulier d'arbres \mathcal{L} est reconnu par l'AEFD \mathcal{A} dont les transitions sont :

$$a, \xrightarrow{\mathcal{A}} A, g(A), \xrightarrow{\mathcal{A}} B, f(A), \xrightarrow{\mathcal{A}} C, g(B), \xrightarrow{\mathcal{A}} G, f(C), \xrightarrow{\mathcal{A}} S, g(G), \xrightarrow{\mathcal{A}} G, f(F), \xrightarrow{\mathcal{A}} F, g(F), \xrightarrow{\mathcal{A}} S, f(G), \xrightarrow{\mathcal{A}} S. \text{ Les états finaux sont } S, B \text{ et } C.$$

L'AEFD \mathcal{A} est minimal et donc isomorphe à $\mathcal{A}_{\mathcal{L}}$. En conséquence du corollaire 10, \mathcal{L} n'est pas un langage régulier d'arbres réversible (bien qu'union de deux langages réguliers d'arbres réversibles).

3.2 Ensembles caractéristiques

Soit \mathcal{L} un langage régulier d'arbres réversible reconnu par l'automate minimal $\mathcal{A}_{\mathcal{L}}$. On associe chaque état q à un terme $\text{rt}(q)$ de $\mathcal{T}(\mathcal{F})$ tel que $\text{rt}(q) \xrightarrow[\mathcal{L}]^* q$ et à un contexte minimale $C_q[\diamond]$ (relativement à la taille), tel que $C_q[q] \xrightarrow[\mathcal{L}]^* q_f$.

L'ensemble d'exemples caractéristiques $CS(\mathcal{L})$ est le plus petit ensemble qui contient

1. le terme $C_q[\text{rt}(q)]$ pour tout état q ,
2. le terme $C_q[\mathfrak{f}(\text{rt}(q_1), \dots, \text{rt}(q_n))]$ pour toute transition $\mathfrak{f}(q_1, \dots, q_n) \xrightarrow{\mathcal{L}} q$.

On voit facilement que $CS(\mathcal{L}) \subseteq \mathcal{L}$. On remarquera également que $C_{q_f}[\diamond] = \diamond$, si q_f est un état final.

Théorème 12

Soient \mathcal{L}_1 et \mathcal{L}_2 deux langages réversibles sur $\mathcal{T}(\mathcal{F})$ et supposons que $CS(\mathcal{L}_1) \subseteq \mathcal{L}_2$. Alors $\mathcal{L}_1 \subseteq \mathcal{L}_2$.

Démonstration. Supposons que $\mathcal{A}' = \langle \mathcal{F}, \mathcal{Q}_1, \{q_{f_1}\}, \xrightarrow{1} \rangle$ et $\mathcal{A}_2 = \langle \mathcal{F}, \mathcal{Q}_2, \{q_{f_2}\}, \xrightarrow{2} \rangle$ soient les AEFD minimaux qui reconnaissent respectivement \mathcal{L}_1 et \mathcal{L}_2 .

Pour tout état $q \in \mathcal{Q}$, $C_q[\text{rt}(q)] \in CS(\mathcal{L}_1) \subseteq \mathcal{L}_2$. Alors, $\text{rt}(q) \in \text{sub}(\mathcal{L}_2)$ et il existe une unique fonction $\theta : \mathcal{Q}_1 \rightarrow \mathcal{Q}_2$ telle que $\text{rt}(q) \xrightarrow{2}^* \theta(q)$, car \mathcal{A}_2 est déterministe.

Le théorème 9, nous permet de faire une observation importante. Pour tout état $t \in \mathcal{T}(\mathcal{F})$, si $C_q[t]$ est dans \mathcal{L}_2 , alors $t \xrightarrow{2}^* \theta(q)$.

Nous montrons par induction sur la taille du terme $t \in \mathcal{T}(\mathcal{F})$, que pour chaque état $q \in \mathcal{Q}$, si $C_q[t] \in \mathcal{L}_1$ alors $C_q[t] \in \mathcal{L}_2$. En conséquence, puisque $C_{q_1}[t] = t$, si $t \in \mathcal{L}_1$ alors $t \in \mathcal{L}_2$.

Cas de base : Supposons que t est un symbole d'arité nulle. Alors il existe une transition $t \xrightarrow{1} q$ et $C_q[t] \in CS(\mathcal{L}_1) \subseteq \mathcal{L}_2$.

Étape d'induction : Supposons que $t = \mathfrak{f}(t_1, \dots, t_n)$. Par hypothèse, il existe un état q_i dans \mathcal{Q}_1 , tel que $t_i \xrightarrow{1}^* q_i$. Puisque $C_{q_i}[t_i] \in \mathcal{L}_1$, nous pouvons appliquer l'hypothèse d'induction pour obtenir $C_{q_i}[t_i] \in \mathcal{L}_2$. De plus, $t_i \xrightarrow{2}^* \theta(q_i)$. Nous avons,

$$C_q[\mathfrak{f}(t_1, \dots, t_n)] \xrightarrow{2}^* C_q[\mathfrak{f}(\theta(q_1), \dots, \theta(q_n))] \quad (1)$$

et d'autre part, puisque $C_q[\mathfrak{f}(\text{rt}(q_1), \dots, \text{rt}(q_n))] \in \mathcal{L}_2$, l'on a

$$C_q[\mathfrak{f}(\text{rt}(q_1), \dots, \text{rt}(q_n))] \xrightarrow{2}^* C_q[\mathfrak{f}(\theta(q_1), \dots, \theta(q_n))] \quad (2)$$

$$\xrightarrow{2}^* q_{f_2} \quad (3)$$

-
- R1** SI $f(p_1, \dots, p_n, q, p_{n+1}, \dots, p_m) \rightarrow p$ et $f(p_1, \dots, p_n, q', p_{n+1}, \dots, p_m) \rightarrow p$ alors $q = q'$.
- R2** SI $f(q_1, \dots, q_n) \rightarrow q$ et $f(q_1, \dots, q_n) \rightarrow q'$ alors $q = q'$.
- R3** si q_{f_1} et q_{f_2} sont deux états finaux alors $q_{f_1} = q_{f_2}$.

FIG. 1 – Les règles de réduction

En combinant (1) et (3), on obtient

$$C_q[f(t_1, \dots, t_n)] \xrightarrow{*} q_{f_2} \quad (4)$$

et donc $C_q[t] \in \mathcal{L}_2$.

□

D'après Angluin (Angluin, 1982), les exemples caractéristiques sont des ensembles *telltale*; cette observation nous permet de conclure d'ores et présent, que la classe des langages réguliers d'arbres réversibles est identifiable.

3.3 Un algorithme d'apprentissage effi cace

L'algorithme d'apprentissage fonctionne comme suit. L'entrée est un ensemble fini S , d'exemples positifs, c'est à dire d'éléments d'un langage inconnu.

Commençons par définir l'automate d'arbres préfixe $PTA(S) = \langle \mathcal{F}_0, \mathcal{Q}_{F_0}, \rightarrow \rangle$. Pour chaque sous-terme t de $sub(S)$, il existe un état noté $[t]$ dans \mathcal{Q} . L'ensemble des états finaux \mathcal{Q}_0 contient chaque état $[t]$, où $t \in S$. Pour chaque sous-terme $f(t_1, \dots, t_n)$ de $sub(S)$, il existe une transition $f([t_1], \dots, [t_n]) \xrightarrow{0} [f(t_1, \dots, t_n)]$. On voit facilement que $\mathcal{L}_{PTA(S)} = S$.

Puis, une suite $\mathcal{A}_0 = PTA(S), \mathcal{A}_1, \dots, \mathcal{A}_n$ d'AEFN est calculée par applications successives d'une des règles de la figure 1, jusqu'à ce qu'aucune règle ne puisse plus s'appliquer. On remarque qu'un automate \mathcal{A}_{i+1} est obtenu à partir d'un automate \mathcal{A}_i par assimilation de deux états distincts de \mathcal{A}_i . Le processus termine après n étapes, puisque à chaque étape, le nombre d'états diminue strictement. On note $nf(S) = \mathcal{A}_n$.

Lemme 13

L'AEFD $nf(S)$ est réversible.

Démonstration. Le processus s'arrête si aucune règle n'est applicable, ce qui implique que $nf(S)$ est réversible.

□

Un homomorphisme d'automates entre les AEFN $\mathcal{A}' = \langle \mathcal{F}, \mathcal{Q}', \mathcal{Q}_{F'}, \rightarrow \rangle$ et $\mathcal{A}_2 = \langle \mathcal{F}, \mathcal{Q}_2, \mathcal{Q}_{F_2}, \xrightarrow{2} \rangle$ est une fonction θ qui applique \mathcal{Q}_1 à \mathcal{Q} et qui a la propriété suivante: si

$f(q_1, \dots, q_n) \xrightarrow{1} q$ alors $f(\theta(q_1), \dots, \theta(q_n)) \xrightarrow{2} \theta(q)$, pour chaque transition de \mathcal{A}' et $\theta(q_f) \in \mathcal{Q}_{F_2}$ si $q_f \in \mathcal{Q}_{F_1}$.

On dit que $\mathcal{A}' \sqsubseteq \mathcal{A}_2$ si il existe un homomorphisme d'automates de \mathcal{A}' vers \mathcal{A}_2 . On remarque que si $\mathcal{A}' \sqsubseteq \mathcal{A}_2$ alors $\mathcal{L}_{\mathcal{A}'} \subseteq \mathcal{L}_{\mathcal{A}_2}$.

Lemme 14

Si $S \subseteq \mathcal{L}$ alors $PTA(S) = \mathcal{A}_0 \sqsubseteq \mathcal{A}_{\mathcal{L}}$.

Démonstration. Définissons θ comme suit. $\theta([t])$ est la classe d'équivalence de t .

□

Lemme 15

Soit \mathcal{A} un AEFD réversible. Supposons que \mathcal{A}_{i+1} par application d'une des règles

R1, **R2** ou **R3** à \mathcal{A} si $\mathcal{A}_i \sqsubseteq \mathcal{A}$ alors $\mathcal{A}_{i+1} \sqsubseteq \mathcal{A}$.

Démonstration. Soit θ l'homomorphisme d'automates de \mathcal{A}_i vers \mathcal{A} . Il y a trois cas possibles.

La règle **R1** a produit \mathcal{A}_{i+1} . Puisque \mathcal{A} est réversible, on a $\theta(q) = \theta(\hat{q})$.

La règle **R2** a produit \mathcal{A}_{i+1} . Puisque \mathcal{A} est déterministe, on a $\theta(q) = \theta(\hat{q})$.

La règle **R3** a produit \mathcal{A}_{i+1} . Puisque \mathcal{A} possède un unique état final, on a $\theta(q_{f_1}) = \theta(q_{f_2})$.

Nous concluons que θ est un homomorphisme d'automates de \mathcal{A}_{i+1} vers \mathcal{A} .

□

Lemme 16

Pour tout langage réversible \mathcal{L} , si $S \subseteq \mathcal{L}$ alors $\mathcal{L}_{PTA(S)} \subseteq \mathcal{L}$.

Démonstration. Conséquence des lemmes 15 et 14.

□

Exemple 17

Considérons les données $\{f(a,b), f(g(a),b), f(g(g(a)),b)\}$. L'algorithme commence par calculer l'automate $PTA(S)$ comme :

$PTA(S) = \langle \mathcal{F}, \mathcal{Q}_0, \mathcal{Q}_{F_0}, \xrightarrow{\circ} \rangle$, où :

- \mathcal{F} est l'ensemble $\{f, g, a, b\}$,
- \mathcal{Q}_0 est l'ensemble $\{[a], [b], [f(a,b)], [g(a)], [f(g(a),b)], [g(g(a))], [f(g(g(a)),b)]\}$,
- \mathcal{Q}_{F_0} est l'ensemble $\{[f(a,b)], [f(g(a),b)], [f(g(g(a)),b)]\}$,
- $\xrightarrow{\circ}$ est l'ensemble $\{a \xrightarrow{\circ} [a], b \xrightarrow{\circ} [b], f([a],[b]) \xrightarrow{\circ} [f(a,b)], g([a]) \xrightarrow{\circ} [g(a)], f([g(a)],[b]) \xrightarrow{\circ} [f(g(a),b)], g([g(a)]) \xrightarrow{\circ} [g(g(a))], f([g(g(a))],[b]) \xrightarrow{\circ} [f(g(g(a)),b)]\}$.

En appliquant trois fois la règle **R3**, on obtient : $[f(a,b)] = [f(g(a),b)] = [f(g(g(a)),b)]$ et le nouvel ensemble $\{a \rightarrow [a], b \rightarrow [b], f([a],[b]) \rightarrow [f(a,b)], g([a]) \rightarrow [g(a)], f([g(a)],[b]) \rightarrow [f(a,b)], g([g(a)]) \rightarrow [g(g(a))], f([g(g(a))],[b])$

$\rightarrow [f(a,b)]$. Nous appliquons maintenant deux fois la règle **R1** pour conclure $[a] = [g(a)] = [g(g(a))]$ et produire l'automate final : $\langle \mathcal{F}, \mathcal{Q}, \mathcal{Q}_F, \rightarrow \rangle$, où :

- \mathcal{F} est l'ensemble $\{f, g, a, b\}$,
- \mathcal{Q} est l'ensemble $\{[a], [b], [f(a,b)]\}$,
- \mathcal{Q}_F est l'ensemble $\{[f(a,b)]\}$,
- \rightarrow est l'ensemble $\{a \rightarrow [a], b \rightarrow [b], f([a],[b]) \rightarrow [f(a,b)], g([a]) \rightarrow [a]\}$.

Théorème 18 (Démonstration du théorème 7)

Supposons que \mathcal{L} est un langage régulier d'arbres réversible. Pour chaque présentation d'exemples positifs t_1, \dots, t_n, \dots de \mathcal{L} , il existe une étape N telle que pour tout $n > N$, $\mathcal{L}_{nf}(t_1, \dots, t_n)$ est l'AEFD réversible minimal qui reconnaît \mathcal{L} . En d'autres termes, la classe des langages réguliers d'arbres réversible est identifiable.

Démonstration. Il existe une étape N telle que pour tout $n > N$, t_1, \dots, t_n contient $CS(\mathcal{L})$. D'après le lemme 16, $\mathcal{L}_{nf}(t_1, \dots, t_n) \subseteq \mathcal{L}$. Or $\mathcal{L}_{nf}(t_1, \dots, t_n)$ est un langage réversible. En appliquant le théorème 12, on conclut que $\mathcal{L} \subseteq \mathcal{L}_{nf}(t_1, \dots, t_n)$. Nous avons donc $\mathcal{L} = \mathcal{L}_{nf}(t_1, \dots, t_n)$.

□

L'algorithme d'apprentissage est incrémentale et s'exécute en temps quadratique en fonction de la taille des données d'entrée.

4 APPRENDRE À PARTIR D'EXEMPLES STRUCTURÉS

Nous donnons dans ce chapitre une illustration de notre méthode d'apprentissage en considérant l'apprentissage à la limite des langages de mots.

4.1 Langages algébriques

Rappelons brièvement la définition des grammaires algébriques (CFG). Une grammaire algébrique G est un quadruplet $\langle \Sigma, \mathcal{N}, \xrightarrow{G}, S \rangle$ où Σ est l'alphabet, \mathcal{N} est l'ensemble des symboles non-terminaux, et $S \in \mathcal{N}$ est le symbole de départ. Σ^* est l'ensemble des mots générés sur l'alphabet Σ . \xrightarrow{G} est l'ensemble des règles de production de la forme : $X \xrightarrow{G} w$ où $X \in \mathcal{N}$ et $w \in (\Sigma \cup \mathcal{N}^*)$. Le langage \mathcal{L}_G défini par G est l'ensemble $\mathcal{L} = \{w \in \Sigma^* : S \xrightarrow{G}^* w\}$. $\mathcal{D}(G)$ est l'ensemble des arbres de dérivation de G .

Il existe une étroite relation entre les langages algébriques et les grammaires régulières d'arbres. Pour commencer, le langage de mots obtenu à partir des arbres d'un langage régulier d'arbres, en lisant pour chaque arbre, le mot formé par les étiquettes des feuilles, est un langage algébrique. En effet, soit Σ l'ensemble des symboles d'arité nulle de \mathcal{F} . On définit la fonction **yield** de $\mathcal{T}(\mathcal{F})$ dans Σ^*

comme suit : $\mathbf{yield}(a) = a$ et $\mathbf{yield}(f(t_1, \dots, t_n)) = \mathbf{yield}(t_1) \dots \mathbf{yield}(t_n)$. Si $\mathcal{L} \subseteq \mathcal{T}(\mathcal{F})$ est un langage d'arbres, alors $\mathbf{yield}(\mathcal{L}) = \{\mathbf{yield}(t) : t \in \mathcal{L}\}$.

Théorème 19

Si \mathcal{L} est un langage régulier d'arbres alors $\mathbf{yield}(\mathcal{L})$ est un langage algébrique.

Ensuite, l'ensemble des arbres de dérivation d'une grammaire algébrique est un langage régulier d'arbres.

Théorème 20

Soit G une grammaire algébrique. L'ensemble $\mathcal{D}(G)$ des arbres de dérivation de G est un langage régulier d'arbres.

Construisons une GRA $\Gamma = \langle \mathcal{F}, \mathcal{X}, \xrightarrow{\Gamma}, S \rangle$, à partir d'une grammaire $G = \langle \Sigma, \mathcal{N}, \xrightarrow{G}, S \rangle$. À chaque lettre de Σ correspond un symbole d'arité nulle de \mathcal{F} . À chaque règle $X \xrightarrow{G} w$, où w est un mot de taille n , correspond un symbole X_n de \mathcal{F} d'arité n . Tout non-terminal de \mathcal{N} est une variable de \mathcal{X} . À chaque production $X \xrightarrow{G} w_1, \dots, w_n$, correspond une production $X \xrightarrow{\Gamma} X_n(Y_1, \dots, Y_n)$, où $Y_i = w_i$ si w_i est dans \mathcal{N} et $Y_i = [w_i]$ sinon. Enfin, pour chaque lettre $a \in \Sigma$, on a $[a] \xrightarrow{\Gamma} a$.

4.2 Exemples structurés

Dans le cas des langages algébriques, le problème d'inférence d'une grammaire compatible avec un langage, à partir de la simple présentation d'une séquence d'éléments de ce langage est difficile. Ainsi, plusieurs travaux ont introduit l'idée d'apprentissage, non plus à partir d'éléments d'un langage, mais à partir de structures intégrant de l'information sur la manière dont ont été produits ces éléments. On parle alors d'*exemples structurés* et l'on passe d'un problème d'apprentissage de langages à un problème d'apprentissage de grammaires. Citons par exemple, les travaux de Sakakibara (Sakakibara, 1992) pour les grammaires algébriques et de Kanazawa (Kanazawa, 1998) pour les grammaires catégorielles. Considérons deux cas :

1. une présentation complète d'une CFG G est une suite t_1, \dots, t_n, \dots d'arbres de dérivation de G ,
2. une fonction de désétiquetage sk est une fonction définie par : $sk(a) = a$ et $sk(f(t_1, \dots, t_n)) = \sigma_n(sk(t_1), \dots, sk(t_n))$ Une présentation squelette d'une grammaire G est une suite $sk(t_1), \dots, sk(t_n), \dots$ de l'ensemble des arbres de dérivation désétiquetés de G .

Une donnée d'apprentissage peut donc être un langage régulier d'arbres, image par homomorphisme de l'ensemble des arbres de dérivation d'une grammaire de mots. Nous pouvons alors considérer une h -présentation d'une CFG G définie comme suit. Soit h un homomorphisme d'arbres, une h -présentation d'une CFG

G est une suite $h(t_1), \dots, h(t_n), \dots$, où t_1, \dots, t_n, \dots est une énumération de l'ensemble des arbres de dérivation de G .

4.3 Identification de grammaires

Soient $h(\mathcal{D}(G)) = \{h(t) : t \in \mathcal{D}(G)\}$, Ω une classe de grammaires et h un homomorphisme d'arbres. Pour toute h -présentation $h(t_1), \dots, h(t_n), \dots$, l'algorithme d'apprentissage A converge vers $G \in \Omega$ si il existe une borne N telle que pour tout $n > N$, le langage produit $A(h(t_1), \dots, h(t_n))$ soit égal à $h(\mathcal{D}(G))$.

Une classe Ω de grammaires est *identifiable* à partir de h -présentations si et seulement si il existe un algorithme d'apprentissage A tel que pour toute grammaire G et toute h -présentation de celle-ci, A converge vers G .

4.4 Les travaux de Sakakibara

Sakakibara a démontré que la classe des grammaires algébriques réversibles est identifiable à partir de présentations squelettes.

Définition 21

Une CFG est réversible si et seulement si

(Reset-free) il n'y a pas deux productions $X \xrightarrow{G} wYv$ et $X \xrightarrow{G} wZv$ telles que Y, Z soient des non-terminaux et $Y \neq Z$,

(Déterministe) il n'y a pas deux productions distinctes $X \xrightarrow{G} t$ et $Y \xrightarrow{G} t$ possédant la même partie droite.

Un langage algébrique réversible est un langage produit par une CFG réversible.

Théorème 22 (Sakakibara (Sakakibara, 1992))

La classe des grammaires algébriques réversibles est identifiable à partir de présentations squelettes.

Démonstration. Les arbres d'une présentation squelette est un langage régulier d'arbres réversible. Le théorème 7 implique donc que cette classe est identifiable. \square

4.5 L'apprentissage des langages algébriques de mots réversibles

Définition 23

La classe \mathcal{R} des langages algébriques de mots est le plus petit ensemble tel que pour tout langage \mathcal{L} de \mathcal{R} , il existe un langage régulier d'arbres réversible \mathcal{L} avec $\mathcal{L} = \text{yield}(\mathcal{L}_1)$. On dit alors que \mathcal{L}_1 est une présentation d'arbres réversibles de \mathcal{L} .

Le théorème 7 nous permet d'établir le résultat suivant.

Théorème 24

La classe \mathcal{R} est identifiable à partir de présentations d'arbres réversibles.

Afin d'illustrer ce résultat et de le comparer avec les travaux de Sakakibara nous donnons des exemples.

Exemple 25

La grammaire Γ dont les productions suivent est réversible. $S \xrightarrow{\Gamma} \varepsilon(X, Y), X \xrightarrow{\Gamma} g(A, X, B), X \xrightarrow{\Gamma} c, Y \xrightarrow{\Gamma} g(A, Y, B), Y \xrightarrow{\Gamma} d, A \xrightarrow{\Gamma} a, B \xrightarrow{\Gamma} b$. D'après le théorème 7, \mathcal{L} est identifiable.

De plus, le théorème précédent implique que le langage algébrique $yield(\mathcal{L})$ est identifiable à partir de toute énumération de \mathcal{L} qui constitue un langage régulier d'arbres réversible. Néanmoins, \mathcal{L} n'est pas l'ensemble des arbres de dérivation d'un langage algébrique. Nous avons $g \xrightarrow{*} acb$ et $g \xrightarrow{*} adb$ et, d'autre part nous devons avoir la règle $\varepsilon \rightarrow gg$. On devrait donc avoir $acbacb$ dans $yield(\mathcal{L}_\Gamma)$, ce qui n'est pas le cas.

Exemple 26

Ce deuxième exemple montre que l'apprentissage peut être facilité par la présentation des arbres de dérivation complets. La grammaire G suivante n'est pas réversible car G n'est pas déterministe. $S \xrightarrow{G} ab, S \xrightarrow{G} aXb, X \xrightarrow{G} ab$. L'ensemble $\mathcal{D}(G)$ des arbres de dérivation est un langage réversible. En effet, il peut être généré par la grammaire réversible minimale suivante. $S \xrightarrow{G} S_2([a],[b]), S \xrightarrow{G} S_3([a],X,[b]), X \xrightarrow{G} X_2([a],[b]), [a] \xrightarrow{G} a, [b] \xrightarrow{G} b$.

5 LES LANGAGES D'ARBRES DE DÉPENDANCES

Suivant Dikovsky et Modina (Dikovsky & Modina, 2000), nous présentons une classe des langages d'arbres de dépendances réversibles introduite par Hays (Hays, 1961) et Gaifman (Gaifman, 1965). Une *grammaire de dépendances lexicalisée* (LDG) Γ est un quadruplet $\langle \Sigma, N, P, S \rangle$, où Σ est un alphabet, N est l'ensemble des non-terminaux, $S \in N$ est le symbole de départ et chaque règle de production

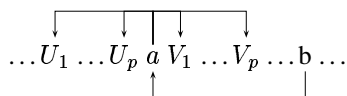
de P est de la forme $X \rightarrow U_1 \dots U_p a V_1 \dots V_q$, où $X \in N$ et chaque U_i et V_j est dans $\Sigma \cup N$. La partie droite peut être vue à la fois comme un arbre étiqueté ordonné de hauteur 1, dont la racine est a , ou comme un mot $U_1 \dots U_p a V_1 \dots V_q$. Ainsi, les nœuds des arbres sont totalement ordonnés.

Les *arbres de dépendances partiels* sont définis comme suit.

1. S est un arbre de dépendances partiel généré par Γ .

2. Si $\dots X \dots b \dots$ est un arbre de dépendances partiel généré par Γ , et si

$X \rightarrow U_1 \dots U_p a V_1 \dots V_q$ est une production de Γ , alors



est un arbre de dépendances partiel g en er e par Γ .

Un *arbre de d ependances* g en er e par une LDG Γ est un arbre de d ependances partiel g en er e par Γ dont tous les n euds sont  etiquet es par des symboles terminaux. Le langage $\mathcal{D}(\Gamma)$ est l'ensemble des arbres de d ependances g en er es par Γ .

On note $\alpha a \beta$ pour $\alpha_1 \dots \alpha_p a \beta_1 \dots \beta_q$.

D efinition 27

Une *grammaire LDG* Γ est r eversible si et seulement si les conditions suivantes sont v erifi ees.

1. Si $X \rightarrow U a V$ et si $Y \rightarrow U a V$, alors $X = Y$.

2. Si $X \rightarrow \alpha Y \beta a \gamma$ et si $X \rightarrow \alpha Z \beta a \gamma$,
alors $Y = Z$, o u $Y, Z \in N$.

3. Si $X \rightarrow \alpha a \beta Y \gamma$ et si $X \rightarrow \alpha a \beta Z \gamma$,
alors $Y = Z$, o u $Y, Z \in N$.

La classe des langages d'arbres de d ependances r eversibles est la classe g en er ee par les LDG r eversibles.

Th eor eme 28

La classe des langages d'arbres de d ependances r eversibles est identifiable.

D emonstration. Nous adaptons l'algorithme d'ecrit en 3.3 de mani ere  a prendre en compte l'ordre des n euds et  a consid erer les symboles d'arit e multiple.

□

6 LES GRAMMAIRES CATÉGORIELLES CLASSIQUES

Bar-Hillel a défini les *grammaires catégorielles*, classe largement utilisée dans le domaine linguistique. Les grammaires catégorielles sont une logique non-commutative, fragment de calcul de Lambek, lui-même fragment de la logique linéaire. Les *types* sont définis à partir d'un ensemble de types atomiques par : un type atomique de *Atoms* est un type et si A et B sont des types, alors A/B et $A \setminus B$ sont des types. Une grammaire catégorielle est un quadruplet $\langle \Sigma, Atoms, Lex, S \rangle$, où Σ est un alphabet et Lex est une fonction qui associe chaque lettre de Σ à un ensemble de types. S est le type final.

Un mot u_1, \dots, u_n est de type A si il existe un type $T_i \in Lex(u_i)$ pour chaque lettre u_i tel que $T_1 \dots T_n \rightarrow A$ par applications des règles suivantes : $A(A \setminus B) \rightarrow B$ et $(A/B)B \rightarrow A$. Un mot est généré par une grammaire si et seulement si il est du type S .

Remarque 29

Nous avons vu que les règles de productions produites par \setminus et $/$ sont réversibles. Plusieurs auteurs ont étendu les grammaires catégorielles classiques par ajout de nouvelles règles. Ces nouvelles règles comme la composition produisent elles aussi des règles réversibles et donc, ces nouvelles grammaires sont elles aussi identifiées.

Les arbres de dérivation sont définis comme suit.

- Une lettre u de Σ est un arbre de dérivation étiqueté par un type $T \in Lex(u)$.
- Si t_1 et t_2 sont deux arbres de dérivation étiquetés respectivement par A et $A \setminus B$, alors $\setminus(t_1, t_2)$ est un arbre de dérivation étiqueté par B .
- Si t_1 et t_2 sont deux arbres de dérivation étiquetés respectivement par A/B et B , alors $\setminus(t_1, t_2)$ est un arbre de dérivation étiqueté par A .

Un arbre de dérivation désétagué est un arbre de dérivation auquel les étiquettes ont été retirées. Dans le livre (Kanazawa, 1998), Kanazawa a étudié les grammaires rigides. Ces grammaires sont des grammaires catégorielles classiques pour lesquelles, à chaque lettre correspond un et un seul type. C'est à dire que Lex est une fonction de Σ dans l'ensemble des types.

Théorème 30 (Kanazawa)

La classe des grammaires catégorielles rigides est identifiable à partir des arbres de dérivation désétagués.

Démonstration. L'ensemble des arbres de dérivation désétagués est un langage régulier d'arbres. Pour s'en rendre compte, construisons une GRA normale Γ . À chaque sous-type A d'un type de lettre, correspond un état $[A]$. À chaque lettre u de Σ , correspond une règle $[Lex(u)] \xrightarrow{\Gamma} u$. À chaque sous-type A d'un type de lettre, correspond une règle : $[B] \xrightarrow{\Gamma} \setminus([A], [A \setminus B])$ ou $[A] \xrightarrow{\Gamma} /([A/B], [B])$. On

remarque que ces règles sont réversibles et puisque la grammaire est rigide, Γ est réversible.

Nous concluons par le théorème 7.

□

RÉFÉRENCES

- ANGLUIN D. (1982). Inference of reversible languages. *Journal of the ACM*, **29**, 741–765.
- BIERMANN A. & FELDMAN J. (1972). On the synthesis of finite state machines from samples of their behavior. *IEEE Trans. on Computers*, **21**, 592–597.
- COMON H., DAUCHET M., GILLERON R., JACQUEMARD F., LUGIEZ D., TISON S. & TOMMASI M. (1997). Tree automata techniques and applications. Available on: <http://www.grappa.univ-lille3.fr/tata>.
- DIKOVSKY A. & MODINA L. (2000). Dependencies on the other side of the curtain. *Traitement automatique des langues*, **41**(1), 67–96.
- EDWARDS J. J., GONZALEZ R. C. & THOMASON M. G. (1976). An algorithm for inference of tree grammars. *International Journal of Computer and Information Sciences*, **5**(2), 145–164.
- FERNAU H. (2001). Learning xml grammars. *Machine Learning and Data Mining in Pattern Recognition*, p. 73–87.
- FUKUDA H. & KAMATA K. (1984). Inference of tree automata from sample set of trees. *International Journal of Computer and Information Sciences*, **13**(3), 177–196.
- GAIFMAN H. (1965). Dependency systems and phrase structure systems. *Information and Control*, **8**(3), 304–337.
- HAYS D. (1961). Grouping and dependency theories. In *National symp. on machine translation*.
- KAMATA K. (1988). Inference methods for tree automata from sample set of trees. *IEEE International Conference on Systems, Man and Cybernetics*, p. 490–493.
- KANAZAWA M. (1998). *Learnable classes of Categorical Grammars*. CSLI.
- KNUUTILA T. & STEINBY M. (1994). The inference of tree languages from finite samples: an algebraic approach. *Theoretical Computer Science*, **129**, 337–367.
- LEE L. (1996). *Learning of Context-Free Languages: A survey of the literature*. Rapport interne, Harvard University.
- LEVINE B. (1981). Derivatives of tree sets with applications to grammatical inference. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **3**(3), 285–293.
- MÄKINEN E. (1997). *Inferring regular languages by merging nonterminals*. Rapport interne, University of Tampere.
- PINKER S. (1994). *The language instinct*. Harper.
- SAKAKIBARA Y. (1992). Efficient learning of context free grammars from positive structural examples. *Information and Computation*, **97**, 23–60.
- SAKAKIBARA Y. (1997). Recent advances of grammatical inference. *Information and Computation*, **185**, 15–45.