

Grammaires hors-contexte. Applications à l'analyse syntaxique.

I. Grammaires hors-contexte [CAR]

Définition 1 $G = (S_0, T, V, P)$ est une grammaire hors-contexte où :

- ▶ T (terminaux) et V (variables ou non-terminaux) sont des alphabets, finis et disjoints.
- ▶ $S_0 \in V$ est l'axiome, le non-terminal duquel on commence à dériver
- ▶ P est un sous-ensemble fini de $V \times (T \cup V)^*$ appelées règles de production.

On les note $X \rightsquigarrow w_1 \dots w_n$.

Notation 2 Par convention, les mots dont les symboles appartiennent à :

- ▶ T sont notés par des minuscules (a, b, \dots)
- ▶ V sont notés par des majuscules (A, B, \dots)
- ▶ $T \cup V$ sont notés par des lettres grecques (α, β, \dots)

Définition 3 Relation de dérivation On peut dériver $s = \alpha X \beta$ en $s' = \alpha \omega \beta$ si $X \rightsquigarrow \omega$. On note cette dérivation $s \rightarrow s'$.

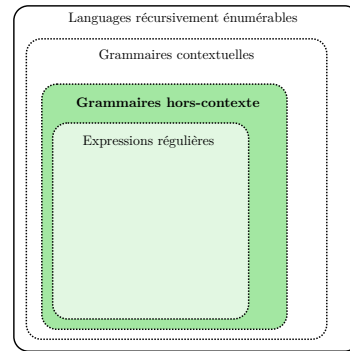
Notation 4 Clôture De plus, on note \rightarrow^* la clôture réflexive et transitive de \rightarrow . Si $s_1 \rightarrow \dots \rightarrow s_2$, alors $s_1 \rightarrow^* s_2$.

Définition 5 Mot engendré Un mot $w \in T^*$ est engendré par la grammaire G si $S \rightarrow^* w$.

Définition 6 Le langage engendré par une grammaire G à partir du non-terminal X est noté $L_G(X) = \{w \in T^* \mid X \rightarrow^* w\}$. Par convention, $L_G = L_G(S_0)$ dénote le langage engendré par G à partir de son axiome.

Exemple 7 Langage des mots bien parenthésés La grammaire

$$S \rightsquigarrow ST \quad S \rightsquigarrow \varepsilon \quad T \rightsquigarrow (S)$$



reconnaît le langage des mots bien parenthésés, aussi appelé langage de Dyck.

Définition 8 Dérivations Soit un mot $s = wX\alpha$ (resp αXw) avec $w \in T^*$. Appliquer une règle $X \rightsquigarrow x$ tels que $s \rightarrow wx\alpha$ (resp $s \rightarrow \alpha xw$) est une dérivation à gauche (resp droite).

Lemme 9 Fondamental Pour $\alpha = \alpha_1 \alpha_2$ et β quelconque,

$$\alpha \xrightarrow{k_1+k_2} \beta \Leftrightarrow \exists \beta_1, \beta_2, k_1, k_2 \text{ tq } \beta = \beta_1 \beta_2 \text{ et } \alpha_1 \xrightarrow{k_1} \beta_1, \alpha_2 \xrightarrow{k_2} \beta_2$$

Cela signifie que l'on peut dériver un mot à gauche ou à droite de manière indépendante et équivalente.

Remarque 10 Les grammaires hors-contexte sont aussi appelées « algébriques », car on peut montrer que les langages engendrés par les grammaires hors-contextes sont des solutions de systèmes d'équations polynomiales.

A. Transformations

Définition 11 Forme normale de Chomsky Une grammaire $G = (S_0, T, V, P)$ est en forme normale de Chomsky si toutes ses règles sont d'une des formes suivantes :

$$A \rightsquigarrow BC \quad \text{où } B, C \in V$$

$$A \rightsquigarrow a \quad \text{où } a \in T$$

Remarque 12 La forme normale de Chomsky est également appelé forme normale quadratique

B. Arbre de dérivation [TIGER]

Définition 13 Arbre de dérivation L'application d'une règle de dérivation $X \rightarrow \alpha_1 \alpha_2 \dots \alpha_n$ peut-être vu comme un arbre ayant pour racine X et pour sous-arbres des arbres de dérivation des α_i produits.

On peut ainsi représenter les applications successives de règles de dérivation par un arbre de dérivation ayant pour racine S .

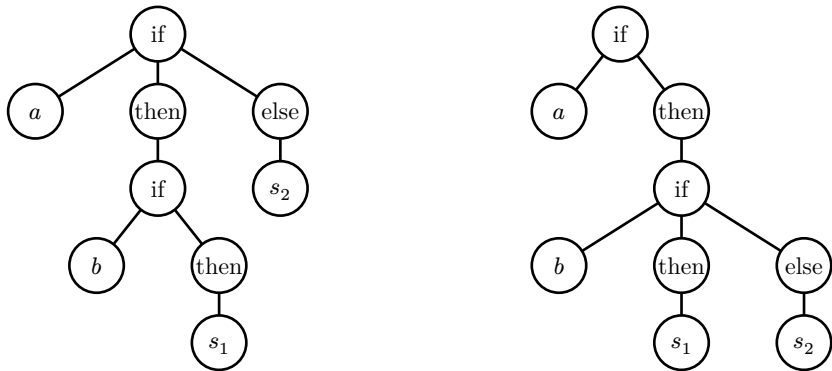
Définition 14 Ambiguïté d'une grammaire [TIGER] Une grammaire G est ambiguë lorsqu'il existe au moins un mot appartenant

au langage engendré L_G pour lequel il existe deux arbres de dérivation distincts.

Exemple 15 Sinon en suspens Avec la grammaire

$$\begin{aligned} S &\rightsquigarrow \text{if } E \text{ then } S \text{ else } S & S &\rightsquigarrow \text{if } E \text{ then } S \\ S &\rightsquigarrow s_1 \mid s_2 & E &\rightsquigarrow a \mid b \end{aligned}$$

le mot « if a then if b then s_1 else s_2 » possède deux arbres de dérivation :



Définition 16 Langage inhéremment ambigu On parle bien de grammaire ambiguë et non de langage. Il est possible de construire deux grammaire distinctes ayant le même langage engendré, avec l’une ambiguë et l’autre non.

Cependant, il existe de langages **inhéremment ambiguës**, pour lesquels toute grammaire engendrant ce langage est ambiguë.

Exemple 17 La grammaire

$$\begin{aligned} S &\rightsquigarrow M & S &\rightsquigarrow U \\ M &\rightsquigarrow \text{if } E \text{ then } M \text{ else } M & M &\rightsquigarrow s_1 \mid s_2 \\ U &\rightsquigarrow \text{if } E \text{ then } S & E &\rightsquigarrow a \mid b \end{aligned}$$

reconnaît le même langage que précédemment, mais n’est pas ambiguë.

Exemple 18 Le langage $L = \{a^m b^m c^n \mid m, n \geq 0\} \cup \{a^m b^n c^n \mid m, n \geq 0\}$ est inhéremment ambigu.

Théorème 19 Bar-Hillel, Perles et Shamir Pour tout langage algébrique L , il existe $N > 0$ tel que pour tout mot $f \in L$, si $|f| \geq$

N alors on peut trouver une factorisation $f = \alpha u \beta v \gamma$ telle que $|uv| > 0, |u \beta v| < N$ et $\alpha u^n \beta v^n \gamma \in L$ pour tout $n \geq 0$.

Remarque 20 Ce théorème est analogue au lemme de l’étoile dans le contexte des grammaires hors-contexte.

C. Automates à pile

Ouverture 21 Automate à pile On peut étendre la définition d’un automate pour y ajouter une pile : les transitions dépendent alors de l’état, de la lettre lue, et du sommet de la pile. On peut également choisir d’empiler ou de dépiler des éléments.

Propriété 22 Un langage $L \subseteq A^*$ est algébrique si et seulement s’il existe un automate à pile qui accepte L .

II. Analyse syntaxique [TIGER]

DEV **Algorithme 23** Cocke-Younger-Kasami (CYK) L’algorithme CYK est un algorithme d’analyse syntaxique sur les grammaires sous forme normale de Chomsky utilisant la programmation dynamique. Sa complexité en $O(|w|^3 + |P|)$ et sa flexibilité le rendent utile pour certaines tâches.

Remarque 24 En pratique, l’algorithme CYK est trop lent pour être utilisable dans un compilateur. On préférera des algorithmes moins généraux, mais traitant un sous-ensemble choisi de langages en un temps de traitement linéaire.

A. Analyse syntaxique LL(1)

Définition 25 Ensembles premiers, suivants et annulable Soit une grammaire $G = (S_0, T, V, P)$, on définit :

- ▶ Annulable(X), pour $X \in V$ est vrai si $X \rightarrow^* \varepsilon$
- ▶ Premiers(γ), pour $\gamma \in (T \times V)^*$ est l’ensemble des terminaux pouvant commencer les mots dérivants de γ
- ▶ Suivants(X), pour $X \in V$ est l’ensemble des terminaux pouvant immédiatement suivre X , c’est à dire $t \in \text{Suivants}(X)$ s’il existe une dérivation de la forme $Y \xrightarrow{*} \omega X t$, avec $Y \in V$ et $\omega \in (T \times V)^*$.

Algorithme 26 Calcul des ensembles premiers, suivants et annulables. On peut calculer ces ensembles à l'aide d'un algorithme de recherche de point fixe. On initialise premiers et suivants aux ensembles vides, et annulables à tout faux.

```

pour tout terminal Z
    premiers[Z] ← {Z}
répéter
    pour chaque règle  $X \rightarrow Y_1 Y_2 \dots Y_k$ 
        si  $Y_1 \dots Y_k$  are tous annulable (ou si  $k = 0$ )
            alors annulable[X] ← vrai
        pour i de 1 à k, j_i de i + 1 à k
            si  $Y_1 \dots Y_{i-1}$  sont tous annulable (ou si  $i = 1$ )
                alors premiers[X] ← premiers[X] ∪ premiers[Y_i]
            si  $Y_{i+1} \dots Y_k$  sont tous annulable (ou si  $i = k$ )
                alors suivants[Y_i] ← suivants[Y_i] ∪ suivants[X]
            si  $Y_{i+1} \dots Y_{j_i-1}$  sont tous annulable (ou si  $i + 1 = j_i$ )
                alors suivants[Y_i] ← suivants[Y_i] ∪ premiers[Y_{j_i}]
tant que premiers, suivants, et annulables n'ont pas changé

```

Définition 27 Méthode LL(1) La méthode d'analyse syntaxique LL(1) (pour Left to right, Leftmost derivation, 1, de gauche à droite, dérivation gauche, en regardant 1 token) consiste à déterminer quelle règle de la grammaire appliquer en lisant 1 token, de gauche à droite.

On dit alors qu'un langage est LL(1) si il est reconnu par un analyseur LL(1).

Propriété 28 Grammaire LL(1) [COMPI] Une grammaire G est dite LL(1) si et seulement si, pour tout non-terminal N , avec les règles de productions $N \rightsquigarrow \alpha_1, \dots, N \rightsquigarrow \alpha_n$

- Les ensembles $\text{Premiers}(\alpha_1), \dots, \text{Premiers}(\alpha_n)$ sont disjoints deux à deux
- Si de plus on a $\text{Annulable}(N)$, alors $\text{Suivants}(N)$ est disjoint de chaque $\text{Premiers}(\alpha_i)$

Exemple 29 La grammaire

$S \rightsquigarrow +SS \quad S \rightsquigarrow T \quad T \rightsquigarrow \text{id} \quad T \rightsquigarrow \text{num}$

est LL1. On peut construire la table suivante :

	+	id	num
S	$S \rightsquigarrow +SS$	$S \rightsquigarrow T$	$S \rightsquigarrow T$
T		$T \rightsquigarrow \text{id}$	$T \rightsquigarrow \text{num}$

qui nous donne, pour chaque état, pour chaque token lu, quelle règle doit être appliquée.

Définition 30 Récursivité à gauche Une grammaire est récursive à gauche si elle contient une règle de production de la forme : $A \rightsquigarrow A\omega$.

Proposition 31 Certains langages ne sont pas analysables par un algorithme LL(1). C'est notamment le cas des grammaires récursives à gauche.

Transformation 32 Élimination de la récursion à gauche [TIGER 3.2] Certains langages sont LL(k), mais leur grammaire peut-être donnée récursive à gauche. On peut transformer ces grammaires de sorte à pourvoir les donner en entrée d'un analyseur LL(k).

Proposition 33 Complexité de l'analyse LL(1) L'algorithme est de complexité linéaire en la taille du mot d'entrée.

Exemple 34 LR(k) Il existe d'autres méthodes d'analyse syntaxique, par exemple en faisant des dérivations droite au lieu de dérivations gauche, donnant les algorithmes LR(k).

Pratique 35 Générateur d'analyseur syntaxique En pratique, on utilise des programmes comme Yacc, Bison, Menhir ou antlr pour générer automatiquement un analyseur syntaxique à partir d'une grammaire.

Pratique 36 Résolution d'ambiguïté Certaines ambiguïtés inhérentes peuvent se résoudre avec des règles ad-hoc.

Exemple 37 Problème du sinon en suspens Les grammaires acceptant à la fois les formes `if () then ()` et `if () then () else ()` sont souvent ambiguës, mais on peut résoudre ce problème en choisissant de toujours lier un `else` au `if` le plus proche.

<u>Grammaires hors-contexte. Applications à l'analyse syntaxique.</u>	
<u>I. Grammaires hors-contexte [CAR]</u>	
1 Def	8 Def Dérivations
2 Not	9 Lemme Fondamental
3 Def Relation de dérivation	10 Rem
4 Not Clôture	A. Transformations
5 Def Mot engendré	11 Def Forme normale de Chomsky
6 Def Le langage engendré	12 Rem
7 Ex Langage des mots bien parenthésés	B. Arbre de dérivation [TIGER]
	13 Def Arbre de dérivation
	14 Def Ambiguïté d'une grammaire [TIGER]
15 Ex Sinon en suspens	20 Rem
	C. Automates à pile
	21 Ouverture Automate à pile
	22 Prop
	II. Analyse syntaxique [TIGER]
16 Def Langage inhéremment ambigu	23 Algo Cocke-Younger-Kasami (CYK)
	24 Rem
17 Ex	A. Analyse syntaxique LL(1)
18 Ex	25 Def Ensembles premiers, suivants et annulable
19 Thm Bar-Hillel, Perles et Shamir	
26 Algo Calcul des ensembles premiers, suivants et annulables.	30 Def Récursivité à gauche
	31 Prop
27 Def Méthode LL(1)	32 Transformation Élimination de la récursion à gauche [TIGER 3.2]
	33 Prop Complexité de l'analyse LL(1)
	34 Ex LR(k)
28 Prop Grammaire LL(1) [COMPI]	35 Prat Générateur d'analyseur syntaxique
	36 Prat Résolution d'ambiguïté
29 Ex	37 Ex Problème du sinon en suspens

Remarque

Bibliographie

[CAR] O. Carton & G. G. Gagnees, *Langages formels: Calculabilité et complexité.*

[TIGER] A. Appel, *Modern compiler implementation.*

[COMPI] F. Schwarzentruher & R. Legendre, *Compilation : analyse lexicale et syntaxique.*