

Théorie des langages

Grammaires et langages algébriques

TELECOM Nancy (1A)

2019-2020

- Les langages algébriques (ou hors-contexte ou context free) sont les langages engendrés par les grammaires algébriques. Dans la hiérarchie (de Chomsky) des langages, ils correspondent à la classe immédiatement supérieure à celle des langages réguliers.
- Les langages algébriques ont des applications dans le domaine des langages de programmation (analyse syntaxique, compilation) et en linguistique.
- A TELECOM Nancy, on parle des langages algébriques :
 - dans le module de 1A "Mathématiques pour l'Informatique" : théorie des langages (MAI1), l'analyse syntaxique descendante (MAI2).
 - dans le module de Traduction 2A : analyse syntaxique ascendante, utilisation d'outils dédiés à l'analyse syntaxique, compilation des langages.
- En linguistique, les grammaires algébriques sont utilisées dans le traitement automatique des langues (TAL).

- Définitions et exemples (grammaire algébrique, dérivation, arbre syntaxique, mot généré, langage engendré)
- Ambiguïté d'une grammaire algébrique, d'un langage algébrique
- Expressions arithmétiques : notations fonctionnelle, préfixée, postfixée, infixée
- Langages réguliers et langages algébriques
 - Grammaires régulières à droite
 - Lemmes de l'étoile
- Propriétés des langages algébriques

Définition

On appelle grammaire algébrique un quadruplet $G = (N, T, \rightarrow, X)$ tel que :

- N est un vocabulaire (**fini**) auxiliaire de la grammaire appelé l'ensemble des **symboles non terminaux**. Intuitivement chaque élément de N doit se comprendre comme le nom d'une catégorie syntaxique de mots sur le vocabulaire T de la grammaire.
- T est un autre vocabulaire (**fini**) disjoint de N , les éléments de T s'appellent des **terminaux**. C'est sur ce vocabulaire que sont définis les mots engendrés par la grammaire.
- \rightarrow est une relation binaire **finie** de N vers $(N \cup T)^*$, chaque élément de la relation est une règle de la grammaire G et est notée $A \rightarrow \alpha$ où $A \in N$ et $\alpha \in (N \cup T)^*$.
- X est un élément distingué de N appelé l'axiome de la grammaire et représentant intuitivement la catégorie syntaxique la plus large du langage que l'on veut engendrer.

- Pour chaque règle $A \rightarrow \alpha$, A est un non terminal appelé, membre gauche de la règle, α est le membre droit de la règle.
- Par convention, les lettres majuscules dénotent généralement les non terminaux, et les lettres minuscules les terminaux. On adopte cette convention dans ce cours.
- Lorsque l'on a plusieurs règles avec le même membre gauche, il est courant de les écrire en factorisant le membre gauche commun. Par exemple si l'on a les trois règles suivantes :

$$A \rightarrow \alpha_1 \quad A \rightarrow \alpha_2 \quad A \rightarrow \alpha_3$$

on écrit

$$A \rightarrow \alpha_1 \mid \alpha_2 \mid \alpha_3$$

la barre verticale “|” dénote alors l'alternative.

Exemple 1 : grammaire (très simplifiée) de la langue française

- $N = \{ \text{< phrase >, < groupe nominal >, < groupe verbal >, < groupe-verbal-etre >, < groupe complement >, < nom >, < article >, < verbe >, < adjectif >, < preposition >, < attribut >, ...} \}$
- $T = \{ \text{le, la, chat, souris, table, est, mange, dans, sur, beau, belle, ...} \}$
Le vocabulaire T (l'ensemble des terminaux) de la grammaire est l'ensemble de tous les mots de la langue française.
- $X = \text{< phrase >}$, l'axiome est le concept syntaxique le plus général de la grammaire, ici c'est la phrase, correspondant au symbole < phrase > .

Exemple 2 : le langage de programmation Pascal

$N = \{ \text{< programme >, < partie declaration >,}$
 $\text{< partie instruction >, < declaration >, < instruction >,}$
 $\text{< affectation >, < identificateur >,}$
 $\text{< conditionnelle >, < iteration >, < procedure >, ...} \}$

$T = \{A, B, \dots, Z, a, \dots, z, (,), [,], +, *, /, ;, \dots\}$

$X = \text{< programme >}$

$\text{< programme >} \rightarrow \text{PROGRAM < identificateur >;}$
 $\text{< partie declaration >}$
 $\text{< partie instruction >}$

$\text{< identificateur >} \rightarrow \dots$

$\text{< partie declaration >} \rightarrow \dots$

$\text{< partie instruction >} \rightarrow \text{< declaration > |}$
 $\rightarrow \text{< declaration >; < partie instruction > |}$
 $\rightarrow \dots$

Définition

Soit $G = (N, T, \rightarrow, X)$ une grammaire algébrique. On définit sur $(N \cup T)^*$ les relations de réécriture \rightarrow , de dérivation $\xrightarrow{*}$ et de dérivation stricte $\xrightarrow{+}$ de la manière suivante :

- 1 $\alpha \rightarrow \beta$ si et seulement si, il existe $\alpha_1 \in (N \cup T)^*$, $\alpha_2 \in (N \cup T)^*$, $A \in N$, $\gamma \in (N \cup T)^*$ tels que $\alpha = \alpha_1 A \alpha_2$ et $\beta = \alpha_1 \gamma \alpha_2$ et $A \rightarrow \gamma$ est une règle de G .
- 2 $\alpha \xrightarrow{*} \beta$ si et seulement si, il existe un entier n et une suite finie $(\alpha_i)_{0 \leq i \leq n}$ telle que $\alpha = \alpha_0$ et $\beta = \alpha_n$ et pour tout i de $[0, n-1]$, $\alpha_i \rightarrow \alpha_{i+1}$.
- 3 $\alpha \xrightarrow{+} \beta$ si et seulement si, il existe un entier n non nul et une suite finie $(\alpha_i)_{0 \leq i \leq n}$ telle que $\alpha = \alpha_0$ et $\beta = \alpha_n$ et pour tout i de $[0, n-1]$, $\alpha_i \rightarrow \alpha_{i+1}$.

La suite $(\alpha_i)_{0 \leq i \leq n}$ s'appelle une dérivation de α à β (ou une dérivation au sens large si $n > 0$).

G définie par $N = \{X, A, B\}$, $T = \{a, b\}$, X l'axiome de la grammaire, les règles sont les suivantes :

$$X \rightarrow \varepsilon \mid aB \mid bA$$

$$A \rightarrow aX \mid bAA$$

$$B \rightarrow bX \mid aBB$$

Exemples de réécritures :

$$aBbA \rightarrow aaBBbA$$

$$aaBbA \rightarrow aaBbXbA$$

$$aaBbXbA \rightarrow aaBb\varepsilon bA = aaBbbA$$

Exemple de dérivation :

$$aBbA \rightarrow aaBBbA \rightarrow aaBbXbA \rightarrow aaBbbA \text{ d'où } aBbA \xrightarrow{*} aaBbbA$$

Mot et langage engendrés par une grammaire algébrique

Définition

Soit $G = (N, T, \rightarrow, X)$ une grammaire algébrique.

- On dit qu'un mot α de T^* est engendré par la grammaire G à partir de A ($A \in N$) si et seulement si $A \xrightarrow{*} \alpha$.
- On dit qu'un mot α de T^* est engendré par la grammaire G si et seulement si $X \xrightarrow{*} \alpha$ (α est engendré par G à partir de l'axiome X de la grammaire).
- Le langage engendré par la grammaire G à partir de A ($A \in N$) est l'ensemble des mots engendrés par G à partir de A , on le note $L(G, A)$.
- Le langage engendré par la grammaire G est le langage engendré par la grammaire à partir de l'axiome X de la grammaire, on le note $L(G)$. On a $L(G) = L(G, X) = \{\alpha ; \alpha \in T^* \text{ et } X \xrightarrow{*} \alpha\}$.
- Un langage engendré par une grammaire algébrique s'appelle un langage algébrique (ou hors-contexte).

Exemple de langage engendré par une grammaire algébrique

- Soit la grammaire suivante $G = (\{X\}, \{a, b\}, \rightarrow, X)$, où les règles sont définies par $X \rightarrow aXb \mid \varepsilon$
On a $L(G) = \{a^n b^n ; n \geq 0\}$.
- Soit la grammaire suivante $G = (\{X\}, \{a, b\}, \rightarrow, X)$ où les règles sont $X \rightarrow aXa \mid bXb \mid \varepsilon$.
Le langage engendré par G est $L(G) = \{w\tilde{w} ; w \in \{a, b\}^*\}$ (\tilde{w} est le mot obtenu en renversant le mot w).

Soit la grammaire algébrique $G = (\{X\}, \{a, b\}, \rightarrow, X)$, où \rightarrow est défini par

$$X \rightarrow XX \mid aXa \mid bXb \mid \varepsilon$$

Soit le mot $aabaab$, on veut montrer que $aabaab \in L(G)$

$$X \xrightarrow{1} XX \xrightarrow{2} aXaX \xrightarrow{4} aaX \xrightarrow{3} aabXb \xrightarrow{2} aabaXab \xrightarrow{4} aabaab$$

$$X \xrightarrow{1} XX \xrightarrow{3} XbXb \xrightarrow{2} XbaXab \xrightarrow{4} Xbaab \xrightarrow{2} aXabaab \xrightarrow{4} aabaab$$

il y a 8 autres dérivations possibles pour générer le mot $aabaab$ à partir de l'axiome X .

Remarque : le fait que l'ordre dans lequel les règles sont appliquées n'est pas important est une caractéristique des grammaires algébriques, aussi appelées grammaires hors-contexte (car un non-terminal peut être remplacé indépendamment de la chaîne de caractères qui l'entoure).

Définition

Soit $G = (N, T, \rightarrow, X)$, un arbre syntaxique pour la grammaire G est un arbre étiqueté par les éléments de $N \cup T \cup \{\varepsilon\}$ qui satisfait les conditions suivantes :

- la racine de l'arbre est étiquetée par X , l'axiome de G .
- chaque nœud interne est étiqueté par un élément de N . Chaque feuille est étiquetée par un élément de T ou par ε .
- pour tout nœud interne, si son étiquette est un non-terminal A et si ses descendants immédiats sont les nœuds n_1, \dots, n_k ayant respectivement pour étiquettes X_1, \dots, X_k alors $A \rightarrow X_1 \dots X_k$ doit être une règle de G .
- si un nœud est étiqueté par ε , alors ce nœud est le seul descendant immédiat de son prédécesseur (cette dernière règle évite l'introduction d'instances inutiles de ε dans l'arbre syntaxique).

Définition

Le mot généré par un arbre syntaxique est celui obtenu par la concaténation des feuilles de l'arbre prises de gauche à droite.

Exemple

Soit la grammaire algébrique $G = (\{X\}, \{a, b\}, \rightarrow, X)$, où \rightarrow est défini par

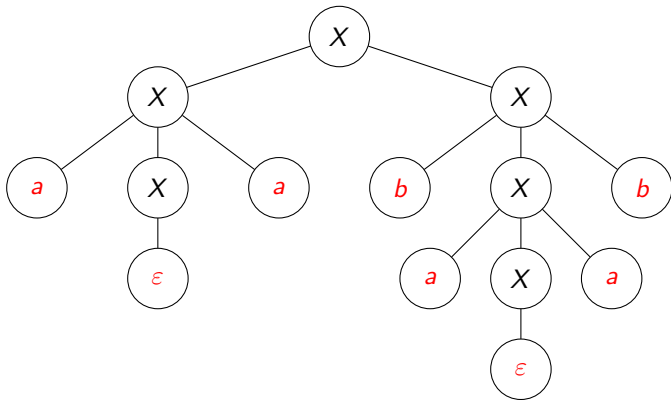
$$X \rightarrow XX \mid aXa \mid bXb \mid \varepsilon$$

Proposition

Soit G une grammaire algébrique, un mot w est généré par G si et seulement si il existe un arbre syntaxique de la grammaire G qui génère w .

Exemple d'arbre syntaxique

Un arbre syntaxique du mot *aabaab*



Ambiguïté d'une grammaire algébrique et d'un langage algébrique

Définition

Une grammaire G est ambiguë si et seulement si il existe un mot $w \in L(G)$ engendré par deux arbres syntaxiques différents.

Définition

Soit L un langage algébrique. On dit que L possède une ambiguïté inhérente si et seulement si toute grammaire algébrique engendrant L est ambiguë.

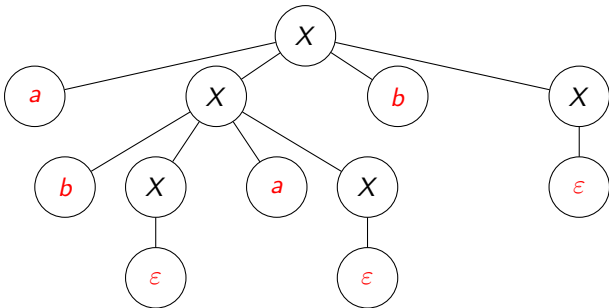
Exemple

Soit $G = (\{X\}, \{a, b\}, \rightarrow, X)$, la grammaire définie par les règles suivantes : $X \rightarrow aXbX \mid bXaX \mid \varepsilon$, est ambiguë.

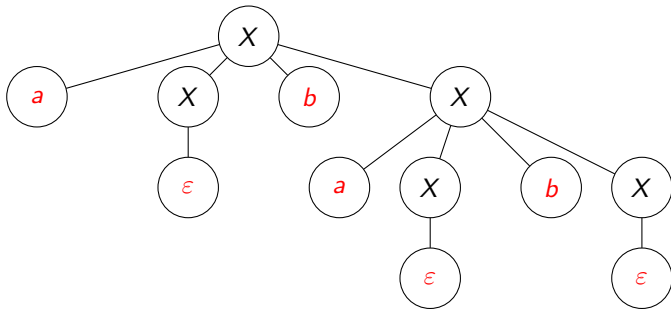
Le langage $L(G)$ ne possède pas d'ambiguïté inhérente (on peut montrer qu'il existe au moins une grammaire non ambiguë qui l'engendre).

Exemple : ambiguïté

Soit $G = (\{X\}, \{a, b\}, \rightarrow, X)$, la grammaire définie par les règles $X \rightarrow aXbX \mid bXaX \mid \varepsilon$ est ambiguë. Arbre syntaxique du mot **abab** :



Deuxième arbre syntaxique de *abab*



Le mot *abab* a deux arbres syntaxiques, il est ambigu, la grammaire G est donc ambiguë.

Les langages de programmation classiques (PASCAL, LISP, FORTRAN, C, ...) sont en première approximation des expressions arithmétiques compliquées...

On s'intéresse ici à la **syntaxe** des expressions arithmétiques. Pour construire les expressions arithmétiques on a besoin de noms de fonctions. Les symboles de fonctions sont classés selon leur **arité** (nombre d'arguments de la fonction) :

- $\{a, b, c\}$ les symboles de constantes (d'arité 0).
- $\{s, f, g\}$ les symboles d'arité 1.
- $\{+, *, /, -\}$ les opérateurs binaires.
- $\{\psi, \varphi\}$ les opérateurs d'arité 3.

On dispose d'un ensemble de symboles de variables

$$X = \{x, y, z, x', y', z'\}.$$

L'ensemble des symboles terminaux est constitué de l'ensemble des symboles de fonctions et des symboles de variables (plus éventuellement des symboles de ponctuation et des parenthèses selon les cas).

- $G = (\{X, A, B, C, D, V\}, T, \rightarrow, X)$ telle que
 - $T = \{x, y, z, x', y', z'\} \cup \{a, b, c\} \cup \{s, f, g\} \cup \{+, *, /, -\} \cup \{\psi, \varphi\} \cup \{(\textcolor{red}{(}, \textcolor{red}{)}, \textcolor{red}{.})\}$
 - \rightarrow est définie par :

$$X \rightarrow V \mid A \mid B(X) \mid C(X, X) \mid D(X, X, X)$$

$$V \rightarrow x \mid y \mid z \mid x' \mid y' \mid z'$$

$$A \rightarrow a \mid b \mid c$$

$$B \rightarrow s \mid f \mid g$$

$$C \rightarrow + \mid * \mid / \mid -$$

$$D \rightarrow \psi \mid \varphi$$

- $*(\phi(s(x), f(+ (a, y)), b), /(\psi(a, z, +(y, c)), *(x, y)))$ est un mot engendré par G

Remarque : les parenthèses et les virgules séparent les opérandes des fonctions, elles permettent une meilleure lisibilité des formules arithmétiques.

- $G = (\{X, A, B, C, D, V\}, T, \rightarrow, X)$ telle que
 - $T = \{x, y, z, x', y', z'\} \cup \{a, b, c\} \cup \{s, f, g\} \cup \{+, *, /, -\} \cup \{\psi, \varphi\}$
 - \rightarrow est définie par :

$$\begin{aligned}
 X &\rightarrow V \mid A \mid BX \mid CXX \mid DXXX \\
 V &\rightarrow x \mid y \mid z \mid x' \mid y' \mid z' \\
 A &\rightarrow a \mid b \mid c \\
 B &\rightarrow s \mid f \mid g \\
 C &\rightarrow + \mid * \mid / \mid - \\
 D &\rightarrow \psi \mid \varphi
 \end{aligned}$$

- $*\phi sxf + ayb/\psi az + yc * xy$ est un mot engendré par G

Remarque : notation minimale qui consiste à écrire les opérateurs avant les opérands. Cette notation est bien adaptée si l'on veut démontrer des résultats théoriques sur les expressions arithmétiques.

Notation polonaise inverse ou notation postfixée

- $G = (\{X, A, B, C, D, V\}, T, \rightarrow, X)$ telle que
 - $T = \{x, y, z, x', y', z'\} \cup \{a, b, c\} \cup \{s, f, g\} \cup \{+, *, /, -\} \cup \{\psi, \varphi\}$
 - \rightarrow est définie par :

$$X \rightarrow V \mid A \mid XB \mid XXC \mid XXXD$$

$$V \rightarrow x \mid y \mid z \mid x' \mid y' \mid z'$$

$$A \rightarrow a \mid b \mid c$$

$$B \rightarrow s \mid f \mid g$$

$$C \rightarrow + \mid * \mid / \mid -$$

$$D \rightarrow \psi \mid \varphi$$

- $xsay + fb\phi azyc + \psi xy * / *$ est un mot engendré par G

Remarque : les opérandes sont écrits avant les opérateurs. Cette notation est pratique si l'on veut évaluer les expressions, elle a été utilisée par certaine marque de calculatrice.

C'est la notation usuelle qui consiste à noter les opérateurs binaires sous forme infixée, elle nécessite d'utiliser des parenthèses et des virgules. Les opérateurs d'arité 1 ou supérieures à 2 sont sous forme préfixée.

L'écriture de ces grammaires se fera en exercice.

Exemple de mot : $\phi(s(x), f(a + y), b) * (\psi(a, z, y + c)/(x * y))$

Démonstration de $L(G) = E$: des exemples

- Soit la grammaire suivante $G = (\{X\}, \{a, b\}, \rightarrow, X)$, où les règles de G sont définies par $X \rightarrow aXb \mid \varepsilon$. On a $L(G) = \{a^n b^n ; n \geq 0\}$.
- Soit la grammaire suivante $G = (\{X\}, \{a, b\}, \rightarrow, X)$ où les règles sont $X \rightarrow aXa \mid bXb \mid \varepsilon$. Montrons que ce langage engendré par G est $L(G) = \{w\tilde{w} ; w \in \{a, b\}^*\}$ (\tilde{w} est le mot obtenu en renversant le mot w).

Preuve :

Soit un mot $u \in \{a, b\}^*$ de $L(G)$, on va montrer que u s'écrit $w\tilde{w}$.

$u \in L(G)$ donc $X \xrightarrow{*} u$, on effectue un raisonnement par récurrence sur la longueur n de la dérivation.

- si $n = 1$ alors $X \rightarrow \varepsilon$ d'où $u = \varepsilon$ et $u = \varepsilon\tilde{\varepsilon}$
- si $n > 1$ alors on a l'une des deux dérivations suivantes :

$$X \rightarrow aXa \rightarrow \dots \rightarrow u \text{ ou } X \rightarrow bXb \rightarrow \dots \rightarrow u$$

La première dérivation implique qu'il existe un mot v de $\{a, b\}^*$ tel que $u = ava$, on a donc $X \rightarrow aXa \rightarrow \dots \rightarrow ava$, et la dérivation $X \rightarrow \dots \rightarrow v$ a comme longueur $(n - 1)$, par hypothèse de récurrence il existe donc un mot s sur $\{a, b\}$ tel que $v = s\tilde{s}$ et par conséquent $u = as\tilde{s}a$, comme $\tilde{as} = \tilde{s}a$, en posant $w = as$, on a $u = w\tilde{w}$.

Même raisonnement avec la deuxième dérivation. . .

Réciproque : soit un mot u de la forme $w\tilde{w}$, montrons qu'il est engendré par G , on raisonne par récurrence sur la longueur du mot u .

- si $|u| = 0$ alors $u = \varepsilon$, ε est engendré par G car $X \rightarrow \varepsilon$
- si $|u| > 0$ alors $u = av\tilde{v}a$ ou $u = bv\tilde{v}b$, on suppose que $u = av\tilde{v}a$ (l'autre cas se traite de façon identique),
comme $|v\tilde{v}| < |u|$ par hypothèse de récurrence $v\tilde{v}$ est engendré par G ,
donc il existe une dérivation de la forme $X \rightarrow \dots \rightarrow v\tilde{v}$
et on a $X \rightarrow aXa \rightarrow \dots \rightarrow av\tilde{v}a = u$, d'où $u \in L(G)$

Remarque

Si G est une grammaire algébrique, démontrer $L(G) = E$ revient à montrer les **deux** inclusions $L(G) \subset E$ et $E \subset L(G)$.

- $L(G) \subset E$ se démontre par récurrence sur la longueur de la dérivation d'un mot α de $L(G)$
- $E \subset L(G)$ se démontre par récurrence sur la longueur d'un mot α de E

Définition

Soit $G = (N, T, \rightarrow, X)$ une grammaire algébrique, on dit que G est une grammaire régulière à droite si toutes ses règles sont de la forme :

$$A \rightarrow aB \quad \text{ou} \quad A \rightarrow a \quad \text{ou} \quad A \rightarrow \varepsilon$$

où $A, B \in N$ et $a \in T$.

Remarque

Les membres droits des règles de grammaires régulières à droite contiennent au plus un non terminal qui est situé à droite.

Proposition

- ① Soit $\mathcal{A} = (A, Q, q_0, \delta, T)$ un automate fini déterministe, le langage régulier $L(\mathcal{A})$ reconnu par \mathcal{A} est engendré par la grammaire algébrique suivante $G = (Q, A, \rightarrow, q_0)$ où la relation \rightarrow est définie par :
 - $q \rightarrow aq'$ si et seulement si $q' = \delta(q, a)$.
 - $q \rightarrow \varepsilon$ si et seulement si q est un état terminal.
- ② Soit $\mathcal{A} = (A, Q, I, E, T)$ un automate fini non déterministe, le langage régulier $L(\mathcal{A})$ reconnu par \mathcal{A} est engendré par la grammaire algébrique suivante $G = (Q \cup \{q_N\}, A, \rightarrow, q_N)$ où q_N est un nouvel état et la relation \rightarrow est définie par :
 - $q \rightarrow aq'$ si et seulement si $(q, a, q') \in E$.
 - $q \rightarrow q'$ si et seulement si $(q, \varepsilon, q') \in E$.
 - $q \rightarrow \varepsilon$ si et seulement si $q \in T$.
 - $q_N \rightarrow q_i$ si et seulement si $q_i \in I$.

Corollaire

Tout langage régulier est un langage algébrique.

Preuve de la proposition

On démontre la première partie de la proposition. Par récurrence sur la longueur du mot $\alpha \in A^*$ on montre l'équivalence suivante :

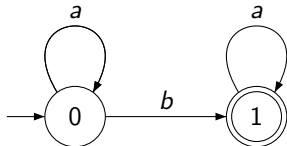
$$\delta(q, \alpha) = q' \text{ si et seulement si } q \xrightarrow{*} \alpha q'$$

pour tout mot α et tout état (ou non terminal) q, q' .

- ① Si le mot α est de longueur 0, $\alpha = \varepsilon$, le résultat est trivial, puisque $\delta(q, \varepsilon) = q$ et $q \xrightarrow{*} q$ sont vraies.
- ② Hypothèse de récurrence : $\forall (q, q') \in Q^2 \delta(q, \alpha) = q' \Leftrightarrow q \xrightarrow{*} \alpha q'$ pour $\alpha \in A^*$
 - ① Si $\delta(q, \alpha a) = \delta(\delta(q, \alpha), a) = q''$, si l'on pose $\delta(q, \alpha) = q'$ en appliquant l'hypothèse de récurrence on déduit (i) $q \xrightarrow{*} \alpha q'$, comme d'autre part $\delta(q', a) = q''$ en appliquant l'hypothèse de récurrence on a (ii) $q' \xrightarrow{*} a q''$. Des deux dérivations (i) et (ii), on déduit $q \xrightarrow{*} \alpha a q''$.
 - ② Si $q \xrightarrow{*} \alpha a q''$, comme la grammaire G est récursive à droite, la lettre a a été produite par une règle de la forme $q' \rightarrow a q''$, la dérivation peut alors s'écrire $q \xrightarrow{*} \alpha q' \xrightarrow{*} \alpha a q''$, par hypothèse de récurrence on a donc $\delta(q, \alpha) = q'$ et $\delta(q', a) = q''$, d'où l'on peut déduire $\delta(q, \alpha a) = q''$.

Il reste à montrer que les mots reconnus par l'automate \mathcal{A} sont exactement ceux générés par la grammaire G . α reconnu par \mathcal{A} ssi $\delta(q_0, \alpha) = q$ avec $q \in T$ ssi $q_0 \xrightarrow{*} \alpha q$ avec $q \in T$ ssi $q_0 \xrightarrow{*} \alpha$ car comme $q \in T$ il existe dans G une règle $q \rightarrow \varepsilon$. CQFD

Soit $\mathcal{A} = (\{a, b\}, \{0, 1\}, 0, \delta, \{1\})$ tel que $\delta(0, a) = 0$, $\delta(0, b) = 1$, $\delta(1, a) = 1$



La grammaire correspondant à cet automate est

$G = (\{X, A\}, \{a, b\}, \rightarrow, X)$ est

X correspond à l'état 0 et A à l'état 1.

Les règles sont :

$X \rightarrow aX \quad (\delta(0, a) = 0)$

$X \rightarrow bA \quad (\delta(0, b) = 1)$

$A \rightarrow aA \quad (\delta(1, a) = 1)$

$A \rightarrow \varepsilon \quad (1 \text{ est un état terminal})$

Proposition

Une grammaire $G = (N, A, \rightarrow, X)$ régulière à droite, engendre le langage régulier reconnu par l'automate fini suivant :

$\mathcal{A} = (A, N, \{X\}, E, T)$ tel que :

- 1 $(A, a, B) \in E$ ssi $A \rightarrow aB$.
- 2 $A \in T$ pour toute règle de la forme $A \rightarrow \varepsilon$.
- 3 S'il existe dans G des règles $A \rightarrow a$, ajouter à N et T un nouvel état q_N et dans E les transitions (A, a, q_N) (l'état q_N étant le même pour toutes les règles $A \rightarrow a$).

Cette proposition est la réciproque de la proposition précédente, on peut conclure qu'un langage est régulier si et seulement s'il peut être engendré par une grammaire régulière à droite.

Lemmes de l'étoile ou lemmes de pompage

Lemme de l'étoile (première version)

Soit L un langage régulier sur l'alphabet A , il existe un entier p tel que, pour tout mot α de L de longueur supérieure à p , il existe une factorisation de α en uvw ($u, v, w \in A^*$) avec $v \neq \varepsilon$, telle que pour tout n dans \mathbb{N} , $uv^n w$ est dans L .
Ce qui se formalise par :

$$(\forall \alpha \in L) |\alpha| \geq p \Rightarrow [(\exists u, v, w \in A^*) \alpha = uvw \text{ et } v \neq \varepsilon \text{ et } (\forall n \in \mathbb{N}) uv^n w \in L]$$

Démonstration

- 1 Si L est fini, on pose $p = \max_{\beta \in L} |\beta| + 1$, comme il n'existe pas de mot de longueur supérieure à p la prémisse de l'implication est fausse donc l'implication est trivialement vraie.
- 2 Si L est infini, soit p le nombre d'états d'un automate reconnaissant le langage L , soit α un mot de L de longueur supérieure (ou égale) à p , il existe un calcul partant de l'état initial de l'automate et aboutissant à un état final et dont l'étiquette est le mot α , comme $|\alpha| \geq p$ et que le nombre d'états de l'automate est p , il existe dans ce calcul un état q_i apparaissant deux fois et il existe donc un calcul d'origine q_i et d'extrémité q_i , soit v l'étiquette de ce calcul, α peut se décomposer en $\alpha = uvw$ avec $v \neq \varepsilon$ et $uv^n w \in L$ (car on peut "tourner" dans la boucle allant de q_i et menant à q_i un nombre quelconque de fois).

Lemme de l'étoile (deuxième version)

Lemme de l'étoile (deuxième version)

Soit L un langage régulier sur l'alphabet A , il existe un entier p tel que pour tout mot α de L et toute factorisation de la forme $\alpha = \gamma\beta\delta$ telle que $|\beta| \geq p$, il existe une factorisation $\beta = uvw$ avec $v \neq \varepsilon$ et telle que pour tout n de \mathbb{N} , le mot $\gamma uv^n w \delta$ appartient à L .

Démonstration

La démonstration est analogue à celle du premier lemme de l'étoile.

Remarque

Les lemmes de l'étoile énoncent des propriétés vérifiées par les langages réguliers, mais les propriétés décrites par ces lemmes ne sont pas des propriétés *caractéristiques* des langages réguliers, c'est-à-dire que certains langages peuvent vérifier ces propriétés sans être réguliers.

Par contre si un langage ne possède pas ces propriétés alors on peut déduire qu'il n'est pas régulier. C'est dans ce sens que l'on utilise souvent les lemmes de l'étoile.

Exemple d'utilisation du lemme de l'étoile

Soit le langage $L = \{a^i b^i, i \in \mathbb{N}\}$, montrons que L n'est pas un langage régulier en appliquant le lemme de l'étoile (version 2).

Soit p l'entier dont il est question dans le lemme de l'étoile, considérons le mot $a^p b^p$ et le facteur $\beta = a^p$, il existe une factorisation $\beta = uvw$ avec $v \neq \varepsilon$ cette factorisation est de la forme $a^q a^r a^s$ avec $q + r + s = p$, $r > 0$ et pour tout n de \mathbb{N} $a^q (a^r)^n a^s b^p \in L$, cette dernière condition ne peut être vérifiée puisque par exemple pour $n = 2$, $q + 2r + s \neq p$. L n'est donc pas un langage régulier.

Proposition

Si L et L' sont des langages algébriques alors $L \cup L'$, LL' , L^* sont des langages algébriques.

Remarques

- Pour montrer la proposition il suffit de construire des grammaires algébriques engendrant les langages $L \cup L'$, LL' et L^* à partir de grammaires engendrant L et L' (à voir en TD).
- L'intersection de deux langages algébriques n'est généralement pas un langage algébrique. De même le complémentaire d'un langage algébrique n'est généralement pas un langage algébrique.
- On peut montrer que l'intersection d'un langage régulier et d'un langage algébrique est un langage algébrique. (Mathématiques Discrètes. Pierre Marchand. Dunod).