

Langages et Automates : LA3
Partie 7 : Grammaires

Grammaire = mécanisme de génération de langages

Introduction à la notion de Grammaire

Grammaire = mécanisme de génération de langages

Plus général que les expressions régulières
ex : on pourra exprimer $\{a^n b^n, n \in \mathbb{N}\}$

Introduction à la notion de Grammaire

Grammaire = mécanisme de génération de langages

Plus général que les expressions régulières
ex : on pourra exprimer $\{a^n b^n, n \in \mathbb{N}\}$

Introduites notamment par Noam Chomsky pour formaliser les propriétés grammaticales des langues naturelles (traduction automatique)

Grammaire = mécanisme de génération de langages

Plus général que les expressions régulières
ex : on pourra exprimer $\{a^n b^n, n \in \mathbb{N}\}$

Introduites notamment par Noam Chomsky pour formaliser les propriétés grammaticales des langues naturelles (traduction automatique)

Elles ont un grand intérêt aussi en informatique, pour les notions de calculabilité, de complexité, ou pour la compilation de programmes.

En programmation, une étape cruciale est la phase de *compilation*,

Compilation = la traduction d'un programme d'un langage A (typiquement le langage de programmation, compréhensible par des humains) vers un langage B (typiquement un binaire compréhensible par l'ordinateur).

En programmation, une étape cruciale est la phase de *compilation*,

Compilation = la traduction d'un programme d'un langage A (typiquement le langage de programmation, compréhensible par des humains) vers un langage B (typiquement un binaire compréhensible par l'ordinateur).

Première étape de la compilation : vérifier par un algorithme si le lexique (les mots utilisés) et la syntaxe (la structure des phrases) est correcte.

En programmation, une étape cruciale est la phase de *compilation*,

Compilation = la traduction d'un programme d'un langage A (typiquement le langage de programmation, compréhensible par des humains) vers un langage B (typiquement un binaire compréhensible par l'ordinateur).

Première étape de la compilation : vérifier par un algorithme si le lexique (les mots utilisés) et la syntaxe (la structure des phrases) est correcte.

Pour vérifier la syntaxe, il faut donc spécifier et décrire la grammaire du langage A, c'est à dire la façon dont peuvent être construites les phrases autorisées.

Exemple

Avant de définir rigoureusement une grammaire, disons qu'il s'agit d'un ensemble de règles de substitution permettant de produire des phrases à partir d'un symbole de départ (souvent noté S).

Par exemple pour la langue française, une grammaire très sommaire pourrait être :

- $S := GN\ GV$
- $GN := Ar\ Ad\ N$
- $GN := GN\ PP\ GN$
- $GV := V$
- $GV := V\ GN$
- $Ar := le\ | \ la\ | \ un\ | \ une\ | \ \epsilon$
- $Ad := petit\ | \ grand\ | \ jeune\ | \ \epsilon$
- $V := regarde\ | \ mange$
- $PP := de\ | \ à$
- $N := Pierre\ | \ Toto\ | \ train\ | \ fille\ | \ viande\ | \ dessert$

(La barre verticale | signifie juste OU)

Avant de définir rigoureusement une grammaire, disons qu'il s'agit d'un ensemble de règles de substitution permettant de produire des phrases à partir d'un symbole de départ (souvent noté S).

Exemple

Pour générer la phrase "La jeune fille de pierre mange", il suffit d'appliquer une suite de substitutions suivante :

$$\begin{aligned}
 S &:= GN\ GV \\
 &:= GN\ PP\ GN\ GV \\
 &:= Ar\ Ad\ N\ PP\ GN\ GV \\
 &:= Ar\ N\ PP\ GN\ V \\
 &:= la\ N\ PP\ GN\ V \\
 &:= la\ fille\ PP\ GN\ V \\
 &:= la\ fille\ de\ GN\ V \\
 &:= la\ fille\ de\ Ar\ Ad\ N\ V \\
 &:= la\ fille\ de\ Pierre\ V \\
 &:= la\ fille\ de\ Pierre\ mange
 \end{aligned}$$

Exemple

Pour générer la phrase "Le petit viande mange un jeune train", il suffit d'appliquer une suite de substitutions suivante :

S := GN GV
:= Ar Ad N GV
:= Ar Ad N V GN
:= Ar Ad N V Ar Ad N
:= le Ad N V Ar Ad N
:= le petit N V Ar Ad N
:= le petit viande V Ar Ad N
:= le petit viande V Ar Ad train
:= le petit viande V un Ad train
:= le petit viande mange un Ad train
:= le petit viande mange un jeune train

Exemple

Comme ce qui nous interesse c'est la structure des phrases (la grammaire) et pas au lexique (quels mots sont autorisées ou pas), en fait les mots "le", "la", "Pierre", "toto" seront en fait représentées par des symboles : des lettres.

- $S := GN \text{ GV}$
- $GN := Ar \text{ Ad } N$
- $GN := GN \text{ PP } GN$
- $GV := V$
- $GV := V \text{ GN}$
- $Ar := a \mid b \mid c \mid d \mid \varepsilon$
- $Ad := e \mid f \mid g \mid \varepsilon$
- $V := i \mid j$
- $PP := k \mid l$
- $N := m \mid n \mid o \mid p \mid q \mid r$

Exemple

Comme ce qui nous interesse c'est la structure des phrases (la grammaire) et pas au lexique (quels mots sont autorisées ou pas), en fait les mots "le", "la", "Pierre", "toto" seront en fait représentées par des symboles : des lettres.

- $S := GN \text{ GV}$
- $GN := Ar \text{ Ad } N$
- $GN := GN \text{ PP } GN$
- $GV := V$
- $GV := V \text{ GN}$
- $Ar := le \mid la \mid un \mid une$
- $Ad := petit \mid grand \mid jeune \mid \varepsilon$
- $V := regarde \mid mange$
- $PP := de \mid \text{à}$
- $N := Pierre \mid Toto \mid train \mid fille \mid viande \mid dessert$

Exemple

Comme ce qui nous interesse c'est la structure des phrases (la grammaire) et pas au lexique (quels mots sont autorisées ou pas), en fait les mots "le", "la", "Pierre", "toto" seront en fait représentées par des symboles : des lettres.

- $S := V_1 \text{ V}_2$
- $V_1 := V_3 \text{ V}_4 \text{ V}_5$
- $V_1 := V_1 \text{ V}_6 \text{ V}_1$
- $V_2 := V_7$
- $V_2 := V_7 \text{ V}_1$
- $V_3 := a \mid b \mid c \mid d \mid \varepsilon$
- $V_4 := e \mid f \mid g \mid \varepsilon$
- $V_7 := i \mid j$
- $V_6 := k \mid l$
- $V_5 := m \mid n \mid o \mid p \mid q \mid r$

Comme ce qui nous interesse c'est la structure des phrases (la grammaire) et pas au lexique (quels mots sont autorisés ou pas), en fait les mots "le", "la", "Pierre", "toto" seront en fait représentées par des symboles : des lettres.

- $S := V_1 V_2$
- $V_1 := V_3 V_4 V_5$
- $V_1 := V_1 V_6 V_1$
- $V_2 := V_7$
- $V_2 := V_7 V_1$
- $V_3 := a \mid b \mid c \mid d \mid \varepsilon$
- $V_4 := e \mid f \mid g \mid \varepsilon$
- $V_7 := i \mid j$
- $V_6 := k \mid l$
- $V_5 := m \mid n \mid o \mid p \mid q \mid r$

Une phrase correcte ne contiendra plus de symboles rouges et sera produite à partir de ces règles de substitution.

Un autre exemple pourrait être la structure des instructions dans un langage

- **Identificateur** :=
- **Expression** := **Expression** + **Expression**
- ...
- **Instruction** := **Identificateur** = **Expression** ;
- **Instruction** := **InstructionConditionio**
- **InstructionConditionio** := **if**(**Condition**) **then** **Instruction** **else** **Instruction**
- **InstructionConditionio** := **while**(**Condition**) **then** **Instruction**

Définition

Une *grammaire algébrique (ou hors contexte)* est définie par un quadruplet $G = (\Sigma, V, S, P)$ où

- Σ est un ensemble de symboles dits *symboles terminaux*
- V est un ensemble (disjoint de Σ) de symboles dits *symboles non terminaux*
- $S \in V$ est un symbole non terminal particulier appelé *axiome*
- $P \subset (V \times (\Sigma \cup V)^*)$ est l'ensemble des *règles de production*

Définition

Une *grammaire algébrique (ou hors contexte)* est définie par un quadruplet $G = (\Sigma, V, S, P)$ où

- Σ est un ensemble de symboles dits *symboles terminaux*
- V est un ensemble (disjoint de Σ) de symboles dits *symboles non terminaux*
- $S \in V$ est un symbole non terminal particulier appelé *axiome*
- $P \subset (V \times (\Sigma \cup V)^*)$ est l'ensemble des *règles de production*

Une règle $(u, v) \in P$ s'écrira $u \rightarrow v$ et signifie "on peut substituer u par v ".

Définition

Une *grammaire algébrique (ou hors contexte)* est définie par un quadruplet $G = (\Sigma, V, S, P)$ où

- Σ est un ensemble de symboles dits *symboles terminaux*
- V est un ensemble (disjoint de Σ) de symboles dits *symboles non terminaux*
- $S \in V$ est un symbole non terminal particulier appelé *axiome*
- $P \subset (V \times (\Sigma \cup V)^*)$ est l'ensemble des *règles de production*

Une règle $(u, v) \in P$ s'écrira $u \rightarrow v$ et signifie "on peut substituer u par v ".

On utilisera par convention des lettres minuscules pour les symboles terminaux et majuscules pour les non terminaux.

Grammaire - Exemples

Considérons l'exemple $G_1 = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$
- $V = \{S\}$
- les règles de production P données par

$$S \rightarrow abS$$

$$S \rightarrow \varepsilon$$

On peut prouver que le langage généré est exactement ???

Définition

Une *grammaire algébrique (ou hors contexte)* est définie par un quadruplet $G = (\Sigma, V, S, P)$ où

- Σ est un ensemble de symboles dits *symboles terminaux*
- V est un ensemble (disjoint de Σ) de symboles dits *symboles non terminaux*
- $S \in V$ est un symbole non terminal particulier appelé *axiome*
- $P \subset (V \times (\Sigma \cup V)^*)$ est l'ensemble des *règles de production*

Une règle $(u, v) \in P$ s'écrira $u \rightarrow v$ et signifie "on peut substituer u par v ".

On utilisera par convention des lettres minuscules pour les symboles terminaux et majuscules pour les non terminaux.

Si plusieurs règles ont le même membre gauche, par exemple $u \rightarrow v$, $u \rightarrow w$, on utilisera la notation $u \rightarrow v|w$.

Grammaire - Exemples

Considérons l'exemple $G_1 = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$
- $V = \{S\}$
- les règles de production P données par

$$S \rightarrow abS$$

$$S \rightarrow \varepsilon$$

On peut prouver que le langage généré est exactement $\{(ab)^n, n \in \mathbb{N}\}$.

Considérons l'exemple $G_2 = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$
- $V = \{S\}$
- les règles de production P données par

$$S \rightarrow aSb \mid \varepsilon$$

On peut prouver que le langage généré est

Considérons l'exemple $G_3 = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$
- $V = \{S\}$
- les règles de production P données par

$$S \rightarrow aSa$$

$$S \rightarrow bSb$$

$$S \rightarrow a \mid b \mid \varepsilon$$

Le langage généré est :

Considérons l'exemple $G_2 = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$
- $V = \{S\}$
- les règles de production P données par

$$S \rightarrow aSb \mid \varepsilon$$

On peut prouver que le langage généré est $\{a^n b^n, n \in \mathbb{N}\}$.

Considérons l'exemple $G_3 = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$
- $V = \{S\}$
- les règles de production P données par

$$S \rightarrow aSa$$

$$S \rightarrow bSb$$

$$S \rightarrow a \mid b \mid \varepsilon$$

Le langage généré est : les palindromes.

Considérons l'exemple $G_4 = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$
- $V = \{S, X\}$
- les règles de production P données par
 - $S \rightarrow aX \mid SS \mid bSaa \mid \varepsilon$
 - $X \rightarrow SX \mid bSaS$

Le langage généré est

Considérons l'exemple $G_4 = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$
- $V = \{S, X\}$
- les règles de production P données par
 - $S \rightarrow aX \mid SS \mid bSaa \mid \varepsilon$
 - $X \rightarrow SX \mid bSaS$

Le langage généré est l'ensemble des mots w tels que $|w|_a = 2|w|_b$

Considérons l'exemple $G_4 = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$
- $V = \{S, X\}$
- les règles de production P données par
 - $S \rightarrow aX \mid SS \mid bSaa \mid \varepsilon$
 - $X \rightarrow SX \mid bSaS$

Le langage généré est l'ensemble des mots w tels que $|w|_a = 2|w|_b$

On peut aussi générer ce langage à l'aide des règles :

$$S \rightarrow SS \mid SbS \mid bSaa \mid \varepsilon$$

Considérons l'exemple $G_4 = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$
- $V = \{S, X\}$
- les règles de production P données par
 - $S \rightarrow aX \mid SS \mid bSaa \mid \varepsilon$
 - $X \rightarrow SX \mid bSaS$

Le langage généré est l'ensemble des mots w tels que $|w|_a = 2|w|_b$

On peut aussi générer ce langage à l'aide des règles :

$$S \rightarrow SS \mid SbS \mid bSaa \mid \varepsilon$$

Remarque

Deux grammaires distinctes peuvent générer le même langage. On dit alors que les grammaires sont **équivalentes**.

Pour certains langages on peut avoir besoin de grammaires plus générales que les grammaires hors contexte (ou algébriques)

ms

On peut pouvoir avoir des règles du type $aX \rightarrow aa$. (on peut remplacer X par un a que si X est déjà précédé d'un a).

Pour certains langages on peut avoir besoin de grammaires plus générales que les grammaires hors contexte (ou algébriques)

ms

On peut pouvoir avoir des règles du type $aX \rightarrow aa$. (on peut remplacer X par un a que si X est déjà précédé d'un a).

On dit que ces règles sont contextuelles, elles dépendent du contexte dans lequel on se trouve.

Pour certains langages on peut avoir besoin de grammaires plus générales que les grammaires hors contexte (ou algébriques)

ms

On peut pouvoir avoir des règles du type $aX \rightarrow aa$. (on peut remplacer X par un a que si X est déjà précédé d'un a).

On dit que ces règles sont contextuelles, elles dépendent du contexte dans lequel on se trouve.

Définition

Une *grammaire formelle* est définie par un quadruplet $G = (\Sigma, V, S, P)$ où

- Σ est un ensemble de symboles dits *symboles terminaux*
- V est un ensemble (disjoint de Σ) de symboles dits *symboles non terminaux*
- $S \in V$ est un symbole non terminal particulier appelé *axiome*
- $P \subset ((\Sigma \cup V)^+ \times (\Sigma \cup V)^*)$ est l'ensemble des *règles de production*

Considérons l'exemple $G_5 = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$
- $V = \{S, X, Y\}$
- les règles de production P données par

$$\begin{aligned} S &\rightarrow aSXa \\ S &\rightarrow aba \\ aX &\rightarrow Xa \\ bX &\rightarrow bb \end{aligned}$$

On peut prouver que le langage généré est exactement

Considérons l'exemple $G_5 = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$
- $V = \{S, X, Y\}$
- les règles de production P données par

$$\begin{aligned} S &\rightarrow aSXa \\ S &\rightarrow aba \\ aX &\rightarrow Xa \\ bX &\rightarrow bb \end{aligned}$$

On peut prouver que le langage généré est exactement $\{a^n b^n a^n, n \in \mathbb{N}\}$.

Langage généré

Définition

- Le langage généré par une grammaire G d'axiome S est l'ensemble de tous les mots $u \in \Sigma^*$ (que des symboles terminaux) tels que $S \Rightarrow u$. On le notera $L(G)$.
- Un langage $L \subset \Sigma^*$ tel qu'il existe une grammaire G avec $L = L(G)$ est dit **récursivement énumérable**.
- Un langage L est dit **algébrique** (ou **hors contexte**) si il existe une grammaire algébrique qui le génère.

On formalise mathématiquement ce dont nous avons parlé précédemment.

Définition

- Si $G = (\Sigma, V, S, P)$ est une grammaire et u et v deux mots de $(V \cup \Sigma)^*$, on dit que G **permet de dériver v à partir de u en une étape**, que l'on note $u \Rightarrow_G v$, si il existe x, y, u' , et v' dans $(V \cup \Sigma)^*$ tels que
 - $u = xu'y$
 - $v = xv'y$
 - $u' \rightarrow v'$ est dans P
- G permet de **dériver v à partir de u en plusieurs étapes**, noté $u \Rightarrow_G^* v$ si il existe u_0, \dots, u_k dans $(V \cup \Sigma)^*$ tels que
 - $u = u_0$
 - $v = u_k$
 - Pour tout $i < k$, G permet de dériver u_{i+1} à partir de u_i en une étape.

(Si il n'y a pas d'ambiguïté, et afin d'alléger les notations, on écrira souvent $u \rightarrow v$ et $u \Rightarrow v$ au lieu de $u \Rightarrow_G v$ et $u \Rightarrow_G^* v$)

Langage généré - Remarques

Lorsque l'on veut prouver qu'une grammaire G génère un langage L comme dans l'exemple ci dessus, il faut montrer deux choses :

- $L(G) \subset L$: tout mot généré par G appartient à L
- $L \subset L(G)$: tout mot de L peut être obtenu par les règles de production de G

Lorsque l'on veut prouver qu'une grammaire G génère un langage L comme dans l'exemple ci dessus, il faut montrer deux choses :

- $L(G) \subset L$: tout mot généré par G appartient à L
- $L \subset L(G)$: tout mot de L peut être obtenu par les règles de production de G

Deux grammaires distinctes pouvant générer le même langage, il se peut qu'un langage algébrique soit généré par une grammaire qui ne l'est pas. Pour montrer qu'un langage n'est pas algébrique, il faut bien montrer que toute grammaire qui le génère n'est pas algébrique.

Grammaire - Exemples

Reprenons l'exemple de la grammaire $G = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$, $V = \{S\}$
- les règles $S \rightarrow aSa \mid bSb \mid a \mid b \mid \varepsilon$

Soit L = l'ensemble des palindromes (cad tq $u = \text{miroir}(u)$).

1) Prouvons que $L \subset L(G)$ par recurrence sur la longueur des mots de L .

Reprenons l'exemple de la grammaire $G = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$, $V = \{S\}$
- les règles $S \rightarrow aSa \mid bSb \mid a \mid b \mid \varepsilon$

Soit L = l'ensemble des palindromes (cad tq $u = \text{miroir}(u)$).

Grammaire - Exemples

Reprenons l'exemple de la grammaire $G = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$, $V = \{S\}$
- les règles $S \rightarrow aSa \mid bSb \mid a \mid b \mid \varepsilon$

Soit L = l'ensemble des palindromes (cad tq $u = \text{miroir}(u)$).

1) Prouvons que $L \subset L(G)$ par recurrence sur la longueur des mots de L .

- Initialisation longueur 0 ou 1. Les mots ε, a, n sont tous dans L .

Reprenons l'exemple de la grammaire $G = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$, $V = \{S\}$
- les règles $S \rightarrow aSa \mid bSb \mid a \mid b \mid \varepsilon$

Soit L = l'ensemble des palindromes (cad tq $u = \text{miroir}(u)$).

1) Prouvons que $L \subset L(G)$ par récurrence sur la longueur des mots de L .

- Initialisation longueur 0 ou 1. Les mots ε, a, n sont tous dans L .
- Induction : Si u de longueur $n \geq 2$ palindrome, alors $u = ava$ ou $u = bvb$ où v est un palindrome plus court.

Par hypothèse de récurrence $v \in L(G)$ donc $S \xRightarrow{*} v$.

Mais alors $S \Rightarrow aSa \xRightarrow{*} ava$ et de même $S \xRightarrow{*} bvb$ et donc $S \xRightarrow{*} u$, c.a.d $u \in L(G)$.

Reprenons l'exemple de la grammaire $G = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$, $V = \{S\}$
- les règles $S \rightarrow aSa \mid bSb \mid a \mid b \mid \varepsilon$

Soit L = l'ensemble des palindromes (cad tq $u = \text{miroir}(u)$).

1) Prouvons que $L \subset L(G)$ par récurrence sur la longueur des mots de L .

- Initialisation longueur 0 ou 1. Les mots ε, a, n sont tous dans L .
- Induction : Si u de longueur $n \geq 2$ palindrome, alors $u = ava$ ou $u = bvb$ où v est un palindrome plus court.

Par hypothèse de récurrence $v \in L(G)$ donc $S \xRightarrow{*} v$.

Mais alors $S \Rightarrow aSa \xRightarrow{*} ava$ et de même $S \xRightarrow{*} bvb$ et donc $S \xRightarrow{*} u$, c.a.d $u \in L(G)$.

2) Prouvons que $L(G) \subset L$ par récurrence sur la longueur des mots de $L(G)$.

Reprenons l'exemple de la grammaire $G = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$, $V = \{S\}$
- les règles $S \rightarrow aSa \mid bSb \mid a \mid b \mid \varepsilon$

Soit L = l'ensemble des palindromes (cad tq $u = \text{miroir}(u)$).

1) Prouvons que $L \subset L(G)$ par récurrence sur la longueur des mots de L .

- Initialisation longueur 0 ou 1. Les mots ε, a, n sont tous dans L .
- Induction : Si u de longueur $n \geq 2$ palindrome, alors $u = ava$ ou $u = bvb$ où v est un palindrome plus court.

Par hypothèse de récurrence $v \in L(G)$ donc $S \xRightarrow{*} v$.

Mais alors $S \Rightarrow aSa \xRightarrow{*} ava$ et de même $S \xRightarrow{*} bvb$ et donc $S \xRightarrow{*} u$, c.a.d $u \in L(G)$.

2) Prouvons que $L(G) \subset L$ par récurrence sur la longueur des mots de $L(G)$.

- Initialisation longueur 0 ou 1. Les mots ε, a, n sont tous dans $L(G)$.

Reprenons l'exemple de la grammaire $G = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$, $V = \{S\}$
- les règles $S \rightarrow aSa \mid bSb \mid a \mid b \mid \varepsilon$

Soit L = l'ensemble des palindromes (cad tq $u = \text{miroir}(u)$).

1) Prouvons que $L \subset L(G)$ par récurrence sur la longueur des mots de L .

- Initialisation longueur 0 ou 1. Les mots ε, a, n sont tous dans L .
- Induction : Si u de longueur $n \geq 2$ palindrome, alors $u = ava$ ou $u = bvb$ où v est un palindrome plus court.

Par hypothèse de récurrence $v \in L(G)$ donc $S \xRightarrow{*} v$.

Mais alors $S \Rightarrow aSa \xRightarrow{*} ava$ et de même $S \xRightarrow{*} bvb$ et donc $S \xRightarrow{*} u$, c.a.d $u \in L(G)$.

2) Prouvons que $L(G) \subset L$ par récurrence sur la longueur des mots de $L(G)$.

- Initialisation longueur 0 ou 1. Les mots ε, a, n sont tous dans $L(G)$.
- Induction : Si u de longueur $n \geq 2$ est dans $L(G)$, alors la dérivation qui mène à u commence nécessairement par $S \Rightarrow aSa$ ou $S \Rightarrow bSb$. Ce qui signifie que $u = ava$ ou $u = bvb$ où v est tel que $S \xRightarrow{*} v$.
 $v \in L(G)$ et $|v| < |u|$, donc par hypothèse de récurrence, v palindrome, donc u palindrome.

Considérons l'exemple $G_6 = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b, +, -, /, \times\}$
- $V = \{S\}$
- les règles de production P données par

$S \rightarrow x y z$	$S \rightarrow S \times S$
$S \rightarrow S + S$	$S \rightarrow S/S$
$S \rightarrow S - S$	$S \rightarrow (S)$

On peut alors fabriquer toutes les expressions arithmétiques correctement parenthésées sur les variables x, y, z .

Considérons l'exemple $G_8 = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b, c\}$
- $V = \{S, X\}$
- les règles de production P données par
 - $S \rightarrow aXSc$
 - $S \rightarrow abc$
 - $Xa \rightarrow aX$
 - $Xb \rightarrow bb$

Le langage généré est

Considérons l'exemple $G_7 = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b, +, ., *\}$
- $V = \{S\}$
- les règles de production P données par

$S \rightarrow \varepsilon$	$S \rightarrow S.S$
$S \rightarrow (S)$	$S \rightarrow S^*$
$S \rightarrow S + S$	$S \rightarrow a b$

On peut alors fabriquer toutes les expressions rationnelles.

Considérons l'exemple $G_8 = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b, c\}$
- $V = \{S, X\}$
- les règles de production P données par
 - $S \rightarrow aXSc$
 - $S \rightarrow abc$
 - $Xa \rightarrow aX$
 - $Xb \rightarrow bb$

Le langage généré est $L = \{a^n b^n c^n \mid n \in \mathbb{N}\}$.

Considérons l'exemple $G_9 = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b, c, d, \dots, z, 0, 1, \dots, 9, \}$
- $V = \{< ident >, < lettre >, < chiffre >, < suiteCar >\}$
- $S = < ident >$
- les règles de production P données par :
 - $< ident > \rightarrow < lettre > < suiteCar >$
 - $< suiteCar > \rightarrow \varepsilon \mid < lettre > < suiteCar > \mid < chiffre > < suiteChar >$
 - $< lettre > \rightarrow a|b|c|\dots|z$
 - $< chiffre > \rightarrow 0|1|2|\dots|8|9$

Le langage généré est

Considérons l'exemple $G_9 = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b, c, d, \dots, z, 0, 1, \dots, 9, \}$
- $V = \{< ident >, < lettre >, < chiffre >, < suiteCar >\}$
- $S = < ident >$
- les règles de production P données par :
 - $< ident > \rightarrow < lettre > < suiteCar >$
 - $< suiteCar > \rightarrow \varepsilon \mid < lettre > < suiteCar > \mid < chiffre > < suiteChar >$
 - $< lettre > \rightarrow a|b|c|\dots|z$
 - $< chiffre > \rightarrow 0|1|2|\dots|8|9$

Le langage généré est l'ensemble des mots qui ne commencent pas par un chiffre (les noms de variables valides en Java par exemple)

Considérons l'exemple $G_{10} = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$
- $V = \{S, X, Y\}$
- les règles de production P données par
 - $S \rightarrow aSa \mid bSb \mid aTb \mid bTa$
 - $T \rightarrow aT \mid bT \mid \varepsilon$

Le langage généré est l'ensemble des mots non palindromes

Considérons l'exemple $G_{11} = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$
- $V = \{S, X, Y, T\}$
- les règles de production P données par
 - $S \rightarrow aXS \mid bYS \mid T$
 - $Xa \rightarrow aX$
 - $Xb \rightarrow bX$
 - $Ya \rightarrow aY$
 - $Yb \rightarrow bY$
 - $YT \rightarrow Tb$
 - $XT \rightarrow Ta$
 - $T \rightarrow \varepsilon$

Le langage généré est

Considérons l'exemple $G_{11} = (\Sigma, V, S, P)$ avec

- $\Sigma = \{a, b\}$
- $V = \{S, X, Y, T\}$
- les règles de production P données par
 - $S \rightarrow aXS \mid bYS \mid T$
 - $Xa \rightarrow aX$ • $YT \rightarrow Tb$
 - $Xb \rightarrow bX$ • $XT \rightarrow Ta$
 - $Ya \rightarrow aY$ • $T \rightarrow \varepsilon$
 - $Yb \rightarrow bY$

Le langage généré est $L = \{uu, u \in \Sigma^*\}$.

Grammaires régulières

Les grammaires régulières constituent une sous famille des grammaires algébriques, dont les règles sont simples.

Définition

Une grammaire est dite *régulière* (ou *linéaire à droite*) est une grammaire telle que les règles de production sont de la forme :

$$A \rightarrow aB \quad A, B \in V \quad (\text{nonterminaux})$$

ou $A \rightarrow \varepsilon$

Quels sont les langages correspondants ??

Les grammaires régulières constituent une sous famille des grammaires algébriques, dont les règles sont simples.

Définition

Une grammaire est dite *régulière* (ou *linéaire à droite*) est une grammaire telle que les règles de production sont de la forme :

$$A \rightarrow aB \quad A, B \in V \quad (\text{nonterminaux})$$

ou $A \rightarrow \varepsilon$

Grammaires régulières

Théoreme

Un langage est rationnel si et seulement si il est généré par une grammaire régulière.

Les symboles non terminaux correspondent exactement aux états de l'automate (S étant l'état initial) et les règles de production aux transitions. Pour toute règle $A \rightarrow \varepsilon$, l'état A est un état acceptant.

Théoreme

Un langage est rationnel si et seulement si il est généré par une grammaire régulière.

Les symboles non terminaux correspondent exactement aux états de l'automate (S étant l'état initial) et les règles de production aux transitions. Pour toute règle $A \rightarrow \varepsilon$, l'état A est un état acceptant.

Formellement : Pour la grammaire, $G = (\Sigma, V, S, P)$ on construit l'automate $\mathcal{A} = (\Sigma, V, S, F, \delta)$, où :

$$F = \{A \in V, P \text{ contient la règle de production } A \rightarrow \varepsilon\}$$

$$\delta(A, a) = B \text{ si et seulement si } P \text{ contient la règle de production } A \rightarrow aB$$

Hierarchie de Chomsky

Chomsky a défini une hiérarchie de complexité de langages par

- type 0 : non contraintes (langages récursivement énumérables)
- type 1 : contextuelles : $uXv \rightarrow uvw$ avec $w \in V^+$
- type 2 : hors-contexte ou algébriques
- type 3 : régulières

Dans le type 1, le fait important est que le membre droit d'une règle est plus long que le membre gauche. On dit parfois qu'elles sont monotones.

Théoreme

Un langage est rationnel si et seulement si il est généré par une grammaire régulière.

Les symboles non terminaux correspondent exactement aux états de l'automate (S étant l'état initial) et les règles de production aux transitions. Pour toute règle $A \rightarrow \varepsilon$, l'état A est un état acceptant.

Formellement : Pour la grammaire, $G = (\Sigma, V, S, P)$ on construit l'automate $\mathcal{A} = (\Sigma, V, S, F, \delta)$, où :

$$F = \{A \in V, P \text{ contient la règle de production } A \rightarrow \varepsilon\}$$

$$\delta(A, a) = B \text{ si et seulement si } P \text{ contient la règle de production } A \rightarrow aB$$

La réciproque consiste à faire la même dans le sens inverse (automate vers grammaire). Cela correspond à la construction du système d'équations linéaires gauche (voir chapitre automates).

Hierarchie de Chomsky

Chomsky a défini une hiérarchie de complexité de langages par

- type 0 : non contraintes (langages récursivement énumérables)
- type 1 : contextuelles : $uXv \rightarrow uvw$ avec $w \in V^+$
- type 2 : hors-contexte ou algébriques
- type 3 : régulières

Dans le type 1, le fait important est que le membre droit d'une règle est plus long que le membre gauche. On dit parfois qu'elles sont monotones.

En terme de machines, le type 3 correspond au AFD, le type 2 aux automates à pile (cf cours suivant) et les types 0 et 1 à des machines de Turing (avec une contrainte pour différencier les deux types).

Chomsky a défini une hiérarchie de complexité de langages par

- type 0 : non contraintes (langages récursivement énumérables)
- type 1 : contextuelles : $uXv \rightarrow uvw$ avec $w \in V^+$
- type 2 : hors-contexte ou algébriques
- type 3 : régulières

Dans le type 1, le fait important est que le membre droit d'une règle est plus long que le membre gauche. On dit parfois qu'elles sont monotones.

En terme de machines, le type 3 correspond au AFD, le type 2 aux automates à pile (cf cours suivant) et les types 0 et 1 a des machines de Turing (avec une contrainte pour différencier les deux types).

A partir de maintenant, on ne considèrera plus que les langages et grammaires algébriques.

Propriétés de Cloture

Théoreme

Les langages algébriques sont clos pour l'union, la concaténation et l'étoile de Kleene.

Si on a deux grammaires G_1 et G_2 de symboles initiaux S_1 et S_2 , il suffit de rajouter un symbole initial S et de rajouter une règle $S \rightarrow S_1 \mid S_2$ pour obtenir l'union des langages engendrés.

Théoreme

Les langages algébriques sont clos pour l'union, la concaténation et l'étoile de Kleene.

Propriétés de Cloture

Théoreme

Les langages algébriques sont clos pour l'union, la concaténation et l'étoile de Kleene.

Si on a deux grammaires G_1 et G_2 de symboles initiaux S_1 et S_2 , il suffit de rajouter un symbole initial S et de rajouter une règle $S \rightarrow S_1 \mid S_2$ pour obtenir l'union des langages engendrés.

On fera de même en rajoutante la règle $S \rightarrow S_1.S_2$ pour obtenir le produit.

Théoreme

Les langages algébriques sont clos pour l'union, la concaténation et l'étoile de Kleene.

Si on a deux grammaires G_1 et G_2 de symboles initiaux S_1 et S_2 , il suffit de rajouter un symbole initial S et de rajouter une règle $S \rightarrow S_1 \mid S_2$ pour obtenir l'union des langages engendrés.

On fera de même en rajoutante la règle $S \rightarrow S_1.S_2$ pour obtenir le produit.

Pour l'étoile, il suffit de rajouter la règle $S \rightarrow S_1S \mid \varepsilon$

Arbre de dérivation

Pour une grammaire algébrique $G = (\Sigma, V, S, P)$ un arbre de dérivation pour un mot w à partir du symbole S est un arbre enraciné tel que :

- la racine est étiquetée par S

Théoreme

Les langages algébriques sont clos pour l'union, la concaténation et l'étoile de Kleene.

Si on a deux grammaires G_1 et G_2 de symboles initiaux S_1 et S_2 , il suffit de rajouter un symbole initial S et de rajouter une règle $S \rightarrow S_1 \mid S_2$ pour obtenir l'union des langages engendrés.

On fera de même en rajoutante la règle $S \rightarrow S_1.S_2$ pour obtenir le produit.

Pour l'étoile, il suffit de rajouter la règle $S \rightarrow S_1S \mid \varepsilon$

Proposition

Les langages algébriques NE sont PAS clos pour l'intersection et le complémentaire

Exemple : $L_1 = \{a^n b^n c^m, (n, m) \in \mathbb{N}^2\}$ et $L_2 = \{a^n b^m c^m, (n, m) \in \mathbb{N}^2\}$ sont algébriques (pourquoi ?), mais leur intersection ne l'est pas (voir plus tard).

Arbre de dérivation

Pour une grammaire algébrique $G = (\Sigma, V, S, P)$ un arbre de dérivation pour un mot w à partir du symbole S est un arbre enraciné tel que :

- la racine est étiquetée par S
- Les feuilles sont étiquetées par des symboles terminaux.

Pour une grammaire algébrique $G = (\Sigma, V, S, P)$ un arbre de dérivation pour un mot w à partir du symbole S est un arbre enraciné tel que :

- la racine est étiquetée par S
- Les feuilles sont étiquetées par des symboles terminaux.
- Les noeuds internes sont étiquetés par des symboles non terminaux

Pour une grammaire algébrique $G = (\Sigma, V, S, P)$ un arbre de dérivation pour un mot w à partir du symbole S est un arbre enraciné tel que :

- la racine est étiquetée par S
- Les feuilles sont étiquetées par des symboles terminaux.
- Les noeuds internes sont étiquetés par des symboles non terminaux
- Pour chaque noeud interne, si X est son étiquette et Y_1, Y_2, \dots, Y_n sont les étiquettes de ses fils (**dans cet ordre**), $X \rightarrow Y_1 Y_2 \dots Y_n$ est une règle de production dans P .
- Le mot obtenu en lisant les feuilles de gauche à droite est le mot w .

Pour une grammaire algébrique $G = (\Sigma, V, S, P)$ un arbre de dérivation pour un mot w à partir du symbole S est un arbre enraciné tel que :

- la racine est étiquetée par S
- Les feuilles sont étiquetées par des symboles terminaux.
- Les noeuds internes sont étiquetés par des symboles non terminaux
- Pour chaque noeud interne, si X est son étiquette et Y_1, Y_2, \dots, Y_n sont les étiquettes de ses fils (**dans cet ordre**), $X \rightarrow Y_1 Y_2 \dots Y_n$ est une règle de production dans P .

Pour une grammaire algébrique $G = (\Sigma, V, S, P)$ un arbre de dérivation pour un mot w à partir du symbole S est un arbre enraciné tel que :

- la racine est étiquetée par S
- Les feuilles sont étiquetées par des symboles terminaux.
- Les noeuds internes sont étiquetés par des symboles non terminaux
- Pour chaque noeud interne, si X est son étiquette et Y_1, Y_2, \dots, Y_n sont les étiquettes de ses fils (**dans cet ordre**), $X \rightarrow Y_1 Y_2 \dots Y_n$ est une règle de production dans P .
- Le mot obtenu en lisant les feuilles de gauche à droite est le mot w .

A toute dérivation correspond un arbre de dérivation.

Exemple avec la grammaire sur $\Sigma = \{a, b\}$ donnée par les règles de production suivantes :

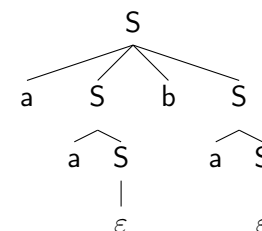
$$S \rightarrow aS \mid aSbS \mid \varepsilon$$

Exemple avec la grammaire sur $\Sigma = \{a, b\}$ donnée par les règles de production suivantes :

$$S \rightarrow aS \mid aSbS \mid \varepsilon$$

Voici un exemple de dérivation et d'arbre associé.

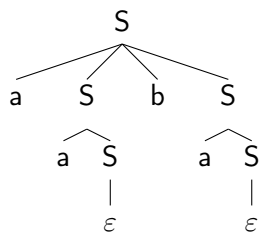
$$\underline{S} \rightarrow aSb\underline{S} \rightarrow aSba\underline{S} \rightarrow a\underline{S}ba \rightarrow aa\underline{S}ba \rightarrow aaba$$



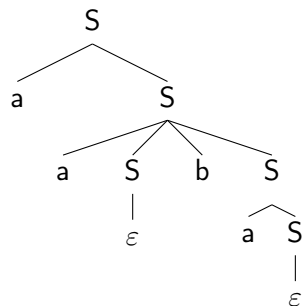
Dérivation la plus à gauche

Attention pour un même mot, il peut y avoir plusieurs dérivation avec des arbres différents

$$\underline{S} \rightarrow aSb\underline{S} \rightarrow aSba\underline{S} \rightarrow a\underline{S}ba \rightarrow aa\underline{S}ba \rightarrow aaba$$



$$\underline{S} \rightarrow a\underline{S} \rightarrow aa\underline{S}bS \rightarrow aab\underline{S} \rightarrow aaba\underline{S} \rightarrow aaba$$



Dérivation la plus à gauche

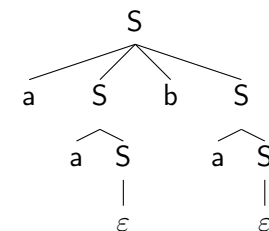
Attention aussi plusieurs dérivation d'un même mot peuvent donner le même arbre

$$\underline{S} \rightarrow aSb\underline{S} \rightarrow aSba\underline{S} \rightarrow a\underline{S}ba \rightarrow aa\underline{S}ba \rightarrow aaba$$

et

$$\underline{S} \rightarrow a\underline{S}bS \rightarrow aa\underline{S}bS \rightarrow aab\underline{S} \rightarrow aaba\underline{S} \rightarrow aaba$$

correspondent toutes les deux à l'arbre :



Pour un arbre donné, une dérivation correspond à un parcours de l'arbre. Pour en choisir une, on peut décider qu'à chaque étape de la dérivation, c'est le symbole non terminal le plus à gauche qui est dérivé. La dérivation obtenue est appelée **dérivation la plus à gauche**. Elle correspond donc à un parcours préfixe de l'arbre.

On est surtout intéressé par l'arbre, donc le fait qu'à un même arbre puisse correspondre plusieurs dérivations, c'est pas gênant, on peut toujours considérer la dérivation la plus à gauche.

On est surtout intéressé par l'arbre, donc le fait qu'à un même arbre puisse correspondre plusieurs dérivations, c'est pas gênant, on peut toujours considérer la dérivation la plus à gauche.

En revanche, le fait que pour un mot donné, il y ait plusieurs arbres et gênant, par exemple en compilation lorsque l'analyseur syntaxique va essayer de reconstruire l'arbre.

Définition

*Si tout mot de $L(G)$ possède une unique arbre de dérivation (ou de façon équivalente une unique dérivation la plus à gauche), on dit que la grammaire est **non ambiguë**.*

Un même langage peut avoir des grammaires qui l'engendrent qui sont ambiguës et d'autres qui ne le sont pas.

On est surtout intéressé par l'arbre, donc le fait qu'à un même arbre puisse correspondre plusieurs dérivations, c'est pas gênant, on peut toujours considérer la dérivation la plus à gauche.

En revanche, le fait que pour un mot donné, il y ait plusieurs arbres et gênant, par exemple en compilation lorsque l'analyseur syntaxique va essayer de reconstruire l'arbre.

La grammaire G donnée par l'axiome S et les règles

$$S \rightarrow aSb \mid aS \mid \varepsilon$$

génère les

La grammaire G donnée par l'axiome S et les règles

$$S \rightarrow aSb \mid aS \mid \varepsilon$$

génère les $a^m b^n$ avec $m \geq n$.

La grammaire G donnée par l'axiome S et les règles

$$S \rightarrow aSb \mid aS \mid \varepsilon$$

génère les $a^m b^n$ avec $m \geq n$.

C'est une grammaire ambiguë (pourquoi ?)

La grammaire G donnée par l'axiome S et les règles

$$S \rightarrow aSb \mid aS \mid \varepsilon$$

génère les $a^m b^n$ avec $m \geq n$.

C'est une grammaire ambiguë (pourquoi ?)

Cependant G' donnée par les règles

$$S \rightarrow aSb \mid aT \mid \varepsilon$$

$$T \rightarrow aT \mid \varepsilon$$

est non ambiguë et génère le même langage.

Mentionnons deux résultats sans en donner les preuves :

Lemme (Conséquence du lemme de Parikh)

Il existe des langages algébriques tels que toute grammaire les génère est ambiguë.

Un exemple est le langage $L = \{a^p b^q c^r, p = q \text{ ou } q = r\}$.

*Un tel langage est dit **inhéremment ambigu**.*

Théoreme

Il n'existe pas d'algorithme pour décider si une grammaire est ambiguë ou non.

La question cruciale une fois une grammaire définie, est de décider si un mot donné appartient ou non au langage décrit par cette grammaire.

La question cruciale une fois une grammaire définie, est de décider si un mot donné appartient ou non au langage décrit par cette grammaire.

Afin de résoudre cette question (et d'autres), il est souvent utile de transformer la grammaire afin de la mettre sous une forme un peu standardisée, on parle de forme normale.

Grammaires réduites

Définition

Un symbole non terminal X d'une grammaire algébrique, G est dit

- **Productif** si il existe un mot $u \in \Sigma^*$ tel que $X \xRightarrow{*} u$
- **Accessible** si $S \xRightarrow{*} uXv$ avec u et v dans V^* .
- **Utile** si il est productif et si $S \xRightarrow{*} uXv$ avec u et v ne contenant que des symboles terminaux ou des symboles non terminaux productifs.

Une grammaire est **réduite** si tous ses symboles non terminaux sont utiles.

On va voir au transparent suivant un algorithme pour supprimer les symboles non utiles d'une grammaire (et donc toutes les règles de production où ils apparaissent) sans changer le langage généré.

Grammaires réduites

Définition

Un symbole non terminal X d'une grammaire algébrique, G est dit

- **Productif** si il existe un mot $u \in \Sigma^*$ tel que $X \xRightarrow{*} u$
- **Accessible** si $S \xRightarrow{*} uXv$ avec u et v dans V^* .
- **Utile** si il est productif et si $S \xRightarrow{*} uXv$ avec u et v ne contenant que des symboles terminaux ou des symboles non terminaux productifs.

Une grammaire est **réduite** si tous ses symboles non terminaux sont utiles.

On va voir au transparent suivant un algorithme pour supprimer les symboles non utiles d'une grammaire (et donc toutes les règles de production où ils apparaissent) sans changer le langage généré.

Toute grammaire est donc équivalente à une grammaire réduite.

Définition

Un symbole non terminal X d'une grammaire algébrique, G est dit

- **Productif** si il existe un mot $u \in \Sigma^*$ tel que $X \xRightarrow{*} u$
- **Accessible** si $S \xRightarrow{*} uXv$ avec u et v dans V^* .
- **Utile** si il est productif et si $S \xRightarrow{*} uXv$ avec u et v ne contenant que des symboles terminaux ou des symboles non terminaux productifs.

Une grammaire est **réduite** si tous ses symboles non terminaux sont utiles.

On va voir au transparent suivant un algorithme pour supprimer les symboles non utiles d'une grammaire (et donc toutes les règles de production où ils apparaissent) sans changer le langage généré.

Toute grammaire est donc équivalente à une grammaire réduite.

Il s'agit de l'équivalent pour les grammaires des automates émondés (tous les états sont accessibles et co-accessibles)

Algorithme de Réduction d'une Grammaire

Pour réduire une grammaire on va appliquer l'algorithme suivant :

- 1 Déterminer l'ensemble des symboles non terminaux qui sont productifs
- 2 Supprimer les symboles non terminaux non productifs ainsi que toutes les règles où ils apparaissent

Pour réduire une grammaire on va appliquer l'algorithme suivant :

- 1 Déterminer l'ensemble des symboles non terminaux qui sont productifs

Algorithme de Réduction d'une Grammaire

Pour réduire une grammaire on va appliquer l'algorithme suivant :

- 1 Déterminer l'ensemble des symboles non terminaux qui sont productifs
- 2 Supprimer les symboles non terminaux non productifs ainsi que toutes les règles où ils apparaissent
- 3 Si S est improductif, la grammaire ne génère aucun mot on peut la remplacer par une grammaire triviale

Pour réduire une grammaire on va appliquer l'algorithme suivant :

- ➊ Déterminer l'ensemble des symboles non terminaux qui sont productifs
- ➋ Supprimer les symboles non terminaux non productifs ainsi que toutes les règles où ils apparaissent
- ➌ Si S est improductif, la grammaire ne génère aucun mot on peut la remplacer par une grammaire triviale
- ➍ Si S est productif, déterminer les symboles non terminaux accessibles

Pour réduire une grammaire on va appliquer l'algorithme suivant :

- ➊ Déterminer l'ensemble des symboles non terminaux qui sont productifs
- ➋ Supprimer les symboles non terminaux non productifs ainsi que toutes les règles où ils apparaissent
- ➌ Si S est improductif, la grammaire ne génère aucun mot on peut la remplacer par une grammaire triviale
- ➍ Si S est productif, déterminer les symboles non terminaux accessibles
- ➎ Supprimer les symboles non terminaux non accessibles ainsi que toutes les règles où ils apparaissent

Pour réduire une grammaire on va appliquer l'algorithme suivant :

- ➊ Déterminer l'ensemble des symboles non terminaux qui sont productifs
- ➋ Supprimer les symboles non terminaux non productifs ainsi que toutes les règles où ils apparaissent
- ➌ Si S est improductif, la grammaire ne génère aucun mot on peut la remplacer par une grammaire triviale
- ➍ Si S est productif, déterminer les symboles non terminaux accessibles
- ➎ Supprimer les symboles non terminaux non accessibles ainsi que toutes les règles où ils apparaissent

On construit l'ensemble de symboles de proche en proche, tout comme on déterminerait dans un automate les états qui sont co-accessibles.

- ➊ Poser $V_0 = \emptyset$

On construit l'ensemble de symboles de proche en proche, tout comme on déterminerait dans un automate les états qui sont co-accessibles.

- 1 Poser $V_0 = \emptyset$
- 2 Calculer successivement les ensembles

$$V_{i+1} = \{X \in V \text{ tels qu'il existe une règle } X \rightarrow u, \text{ où } u \in (\Sigma \cup V_i)^*\}$$

On construit l'ensemble de symboles de proche en proche, tout comme on déterminerait dans un automate les états qui sont co-accessibles.

- 1 Poser $V_0 = \emptyset$
- 2 Calculer successivement les ensembles

$$V_{i+1} = \{X \in V \text{ tels qu'il existe une règle } X \rightarrow u, \text{ où } u \in (\Sigma \cup V_i)^*\}$$

- 3 Arrêter lorsque $V_{i+1} = V_i$.

On construit l'ensemble de symboles de proche en proche, tout comme on déterminerait dans un automate les états qui sont co-accessibles.

- 1 Poser $V_0 = \emptyset$
- 2 Calculer successivement les ensembles

$$V_{i+1} = \{X \in V \text{ tels qu'il existe une règle } X \rightarrow u, \text{ où } u \in (\Sigma \cup V_i)^*\}$$

- 3 Arrêter lorsque $V_{i+1} = V_i$.
- 4 L'ensemble V_i est alors l'ensemble des symboles productifs.

Là encore on fait un parcours en largeur, ici à partir de S , pour déterminer les symboles accessibles.

- 1 Poser $W_0 = \{S\}$

Là encore on fait un parcours en largeur, ici à partir de S , pour déterminer les symboles accessibles.

- ❶ Poser $W_0 = \{S\}$
- ❷ Calculer successivement les ensembles

$$W_{i+1} = W_i \cup \{X \in V \text{ tels qu'il existe une règle } Y \rightarrow uXv, \text{ où } Y \in W_i\}.$$

Là encore on fait un parcours en largeur, ici à partir de S , pour déterminer les symboles accessibles.

- ❶ Poser $W_0 = \{S\}$
- ❷ Calculer successivement les ensembles

$$W_{i+1} = W_i \cup \{X \in V \text{ tels qu'il existe une règle } Y \rightarrow uXv, \text{ où } Y \in W_i\}.$$

- ❸ Arrêter lorsque $W_{i+1} = W_i$.

Là encore on fait un parcours en largeur, ici à partir de S , pour déterminer les symboles accessibles.

- ❶ Poser $W_0 = \{S\}$
- ❷ Calculer successivement les ensembles

$$W_{i+1} = W_i \cup \{X \in V \text{ tels qu'il existe une règle } Y \rightarrow uXv, \text{ où } Y \in W_i\}.$$

- ❸ Arrêter lorsque $W_{i+1} = W_i$.
- ❹ L'ensemble W_i est alors l'ensemble des symboles accessibles

Exemple avec la grammaire sur $\Sigma = \{a, b\}$ donnée par les règles de production suivantes :

$$S \rightarrow a \mid X$$

$$X \rightarrow XY$$

$$Y \rightarrow b$$

Exemple avec la grammaire sur $\Sigma = \{a, b\}$ donnée par les règles de production suivantes :

$$\begin{aligned} S &\rightarrow a \mid X \\ X &\rightarrow XY \\ Y &\rightarrow b \end{aligned}$$

- 1 Symboles productifs : $V_0 = \emptyset$,

Exemple avec la grammaire sur $\Sigma = \{a, b\}$ donnée par les règles de production suivantes :

$$\begin{aligned} S &\rightarrow a \mid X \\ X &\rightarrow XY \\ Y &\rightarrow b \end{aligned}$$

- 1 Symboles productifs : $V_0 = \emptyset$, $V_1 = \{S, Y\}$, $V_2 = V_1$

Exemple avec la grammaire sur $\Sigma = \{a, b\}$ donnée par les règles de production suivantes :

$$\begin{aligned} S &\rightarrow a \mid X \\ X &\rightarrow XY \\ Y &\rightarrow b \end{aligned}$$

- 1 Symboles productifs : $V_0 = \emptyset$, $V_1 = \{S, Y\}$,

Exemple avec la grammaire sur $\Sigma = \{a, b\}$ donnée par les règles de production suivantes :

$$\begin{aligned} S &\rightarrow a \mid X \\ X &\rightarrow XY \\ Y &\rightarrow b \end{aligned}$$

- 1 Symboles productifs : $V_0 = \emptyset$, $V_1 = \{S, Y\}$, $V_2 = V_1$
- 2 On supprime donc X pour obtenir la grammaire

$$\begin{aligned} S &\rightarrow a \\ Y &\rightarrow b \end{aligned}$$

Exemple avec la grammaire sur $\Sigma = \{a, b\}$ donnée par les règles de production suivantes :

$$\begin{aligned} S &\rightarrow a \mid X \\ X &\rightarrow XY \\ Y &\rightarrow b \end{aligned}$$

- ❶ Symboles productifs : $V_0 = \emptyset$, $V_1 = \{S, Y\}$, $V_2 = V_1$
- ❷ On supprime donc X pour obtenir la grammaire

$$\begin{aligned} S &\rightarrow a \\ Y &\rightarrow b \end{aligned}$$
- ❸ Symboles accessibles : $W_0 = \{S\}$,

Réduction d'une Grammaire - Exemple

Exemple avec la grammaire sur $\Sigma = \{a, b\}$ donnée par les règles de production suivantes :

$$\begin{aligned} S &\rightarrow a \mid X \\ X &\rightarrow XY \\ Y &\rightarrow b \end{aligned}$$

- ❶ Symboles productifs : $V_0 = \emptyset$, $V_1 = \{S, Y\}$, $V_2 = V_1$
- ❷ On supprime donc X pour obtenir la grammaire

$$\begin{aligned} S &\rightarrow a \\ Y &\rightarrow b \end{aligned}$$
- ❸ Symboles accessibles : $W_0 = \{S\}$, $W_1 = \{S\}$.
- ❹ On supprime donc Y pour obtenir la grammaire

$$S \rightarrow a$$

Exemple avec la grammaire sur $\Sigma = \{a, b\}$ donnée par les règles de production suivantes :

$$\begin{aligned} S &\rightarrow a \mid X \\ X &\rightarrow XY \\ Y &\rightarrow b \end{aligned}$$

- ❶ Symboles productifs : $V_0 = \emptyset$, $V_1 = \{S, Y\}$, $V_2 = V_1$
- ❷ On supprime donc X pour obtenir la grammaire

$$\begin{aligned} S &\rightarrow a \\ Y &\rightarrow b \end{aligned}$$
- ❸ Symboles accessibles : $W_0 = \{S\}$, $W_1 = \{S\}$.

Grammaires propres

Définition

Une grammaire algébrique est dite *propre* si elle ne contient

- pas de règle unitaire : $A \rightarrow B$ où A et B non terminaux.
- pas de règle ε : $A \rightarrow \varepsilon$

Définition

Une grammaire algébrique est dite *propre* si elle ne contient

- pas de règle unitaire : $A \rightarrow B$ où A et B non terminaux.
- pas de règle ε : $A \rightarrow \varepsilon$

Théoreme

Tout langage algébrique L ne contenant pas le mot vide peut être engendré par une grammaire propre

(si on s'intéresse à un langage algébrique qui contient le mot vide on peut dire une grammaire propre plus une règle $S \rightarrow \varepsilon$ et S n'apparaît pas dans une partie droite d'une règle de production).

Problème du mot

Corollaire (Décidabilité du Problème du mot)

Pour toute grammaire algébrique G et tout mot u , il existe un algorithme qui en temps fini répond oui si u peut être engendré par G , et non sinon

Une fois la grammaire propre, on remarque que toute règle est soit du type $X \rightarrow a$, soit son membre droit a longueur au moins 2. Cela implique qu'une suite de dérivation qui amène au mot u a longueur au plus $|u| - 1$. Il suffit donc de générer toutes ces dérivations possibles (il y en a bien un nombre fini) pour voir si u peut être engendré.

Corollaire (Décidabilité du Problème du mot)

Pour toute grammaire algébrique G et tout mot u , il existe un algorithme qui en temps fini répond oui si u peut être engendré par G , et non sinon

Problème du mot

Corollaire (Décidabilité du Problème du mot)

Pour toute grammaire algébrique G et tout mot u , il existe un algorithme qui en temps fini répond oui si u peut être engendré par G , et non sinon

Une fois la grammaire propre, on remarque que toute règle est soit du type $X \rightarrow a$, soit son membre droit a longueur au moins 2. Cela implique qu'une suite de dérivation qui amène au mot u a longueur au plus $|u| - 1$. Il suffit donc de générer toutes ces dérivations possibles (il y en a bien un nombre fini) pour voir si u peut être engendré.

Remarque 1 : ce nombre de dérivations peut être exponentiel.

Remarque 2 : ceci aurait aussi marché pour les grammaires de type 1

Corollaire (Décidabilité du Problème du mot)

Pour toute grammaire algébrique G et tout mot u , il existe un algorithme qui en temps fini répond oui si u peut être engendré par G , et non sinon

Une fois la grammaire propre, on remarque que toute règle est soit du type $X \rightarrow a$, soit son membre droit a longueur au moins 2. Cela implique qu'une suite de dérivation qui amène au mot u a longueur au plus $|u| - 1$. Il suffit donc de générer toutes ces dérivations possibles (il y en a bien un nombre fini) pour voir si u peut être engendré.

Remarque 1 : ce nombre de dérivations peut être exponentiel.

Remarque 2 : ceci aurait aussi marché pour les grammaires de type 1

Mentionnons sans preuve le résultat suivant

Théoreme

Il existe un algorithme polynomial pour résoudre ce problème pour les grammaires algébriques.

Nettoyage d'une Grammaire - Etape 1

- 1 On calcule les symboles annulables, cad les symboles X tels que $X \Rightarrow \varepsilon$.
Comme pour la réduction on pose $E_0 = \{\varepsilon\}$ et on calcule successivement les ensembles $E_{i+1} = \{X \in V \text{ tels que } (X \rightarrow u) \in P, \text{ où } u \in (E_i)^*\}$ jusqu'à ce que $E_{i+1} = E_i$.
- 2 Pour chaque règle $Y \rightarrow uXv$ où X annulable, on ajoute une règle $Y \rightarrow uv$.

- 1 On calcule les symboles annulables, cad les symboles X tels que $X \Rightarrow \varepsilon$.
Comme pour la réduction on pose $E_0 = \{\varepsilon\}$ et on calcule successivement les ensembles $E_{i+1} = \{X \in V \text{ tels que } (X \rightarrow u) \in P, \text{ où } u \in (E_i)^*\}$ jusqu'à ce que $E_{i+1} = E_i$.

Nettoyage d'une Grammaire - Etape 1

- 1 On calcule les symboles annulables, cad les symboles X tels que $X \Rightarrow \varepsilon$.
Comme pour la réduction on pose $E_0 = \{\varepsilon\}$ et on calcule successivement les ensembles $E_{i+1} = \{X \in V \text{ tels que } (X \rightarrow u) \in P, \text{ où } u \in (E_i)^*\}$ jusqu'à ce que $E_{i+1} = E_i$.
- 2 Pour chaque règle $Y \rightarrow uXv$ où X annulable, on ajoute une règle $Y \rightarrow uv$.
- 3 On supprime toutes les règles $X \rightarrow \varepsilon$

Ensuite on supprime toutes les règles unitaires.

- ① On identifie toutes les paires de symboles non terminaux X et Y tels que $X \xRightarrow{*} Y$ et $Y \xRightarrow{*} X$.

Ensuite on supprime toutes les règles unitaires.

- ① On identifie toutes les paires de symboles non terminaux X et Y tels que $X \xRightarrow{*} Y$ et $Y \xRightarrow{*} X$.
- ② Il s'agit d'une relation d'équivalence entre les symboles non terminaux. On choisit donc un représentant pour chaque classe et on le substitue à chaque symbole de sa classe dans toutes les règles de production.
- ③ Il en résulte peut-être des règles du type $X \rightarrow X$. On les supprime.

Ensuite on supprime toutes les règles unitaires.

- ① On identifie toutes les paires de symboles non terminaux X et Y tels que $X \xRightarrow{*} Y$ et $Y \xRightarrow{*} X$.
- ② Il s'agit d'une relation d'équivalence entre les symboles non terminaux. On choisit donc un représentant pour chaque classe et on le substitue à chaque symbole de sa classe dans toutes les règles de production.

Ensuite on supprime toutes les règles unitaires.

- ① On identifie toutes les paires de symboles non terminaux X et Y tels que $X \xRightarrow{*} Y$ et $Y \xRightarrow{*} X$.
- ② Il s'agit d'une relation d'équivalence entre les symboles non terminaux. On choisit donc un représentant pour chaque classe et on le substitue à chaque symbole de sa classe dans toutes les règles de production.
- ③ Il en résulte peut-être des règles du type $X \rightarrow X$. On les supprime.
- ④ Comme on a supprimé les symboles équivalents à l'étape précédente, la relation $X \xRightarrow{*} Y$ est bien un ordre partiel sur les symboles non terminaux (X est supérieur à Y).

Ensuite on supprime toutes les règles unitaires.

- ➊ On identifie toutes les paires de symboles non terminaux X et Y tels que $X \xRightarrow{*} Y$ et $Y \xRightarrow{*} X$.
- ➋ Il s'agit d'une relation d'équivalence entre les symboles non terminaux. On choisit donc un représentant pour chaque classe et on le substitue à chaque symbole de sa classe dans toutes les règles de production.
- ➌ Il en résulte peut-être des règles du type $X \rightarrow X$. On les supprime.
- ➍ Comme on a supprimé les symboles équivalents à l'étape précédente, la relation $X \xRightarrow{*} Y$ est bien un ordre partiel sur les symboles non terminaux (X est supérieur à Y).
- ➎ On prend le symbole le plus petit Y et on fait les opérations suivantes : pour chaque règle unitaire $X \rightarrow Y$ qu'on veut supprimer et pour chaque règle $Y \rightarrow u$, on ajoute une règle $X \rightarrow u$.

Nettoyage d'une Grammaire - Etape 2

Ensuite on supprime toutes les règles unitaires.

- ➊ On identifie toutes les paires de symboles non terminaux X et Y tels que $X \xRightarrow{*} Y$ et $Y \xRightarrow{*} X$.
- ➋ Il s'agit d'une relation d'équivalence entre les symboles non terminaux. On choisit donc un représentant pour chaque classe et on le substitue à chaque symbole de sa classe dans toutes les règles de production.
- ➌ Il en résulte peut-être des règles du type $X \rightarrow X$. On les supprime.
- ➍ Comme on a supprimé les symboles équivalents à l'étape précédente, la relation $X \xRightarrow{*} Y$ est bien un ordre partiel sur les symboles non terminaux (X est supérieur à Y).
- ➎ On prend le symbole le plus petit Y et on fait les opérations suivantes : pour chaque règle unitaire $X \rightarrow Y$ qu'on veut supprimer et pour chaque règle $Y \rightarrow u$, on ajoute une règle $X \rightarrow u$. Ensuite on peut supprimer la règle unitaire $X \rightarrow Y$.
- ➏ On recommence avec le nouveau symbole minimal jusqu'à ce qu'on ait tout nettoyé.

Ensuite on supprime toutes les règles unitaires.

- ➊ On identifie toutes les paires de symboles non terminaux X et Y tels que $X \xRightarrow{*} Y$ et $Y \xRightarrow{*} X$.
- ➋ Il s'agit d'une relation d'équivalence entre les symboles non terminaux. On choisit donc un représentant pour chaque classe et on le substitue à chaque symbole de sa classe dans toutes les règles de production.
- ➌ Il en résulte peut-être des règles du type $X \rightarrow X$. On les supprime.
- ➍ Comme on a supprimé les symboles équivalents à l'étape précédente, la relation $X \xRightarrow{*} Y$ est bien un ordre partiel sur les symboles non terminaux (X est supérieur à Y).
- ➎ On prend le symbole le plus petit Y et on fait les opérations suivantes : pour chaque règle unitaire $X \rightarrow Y$ qu'on veut supprimer et pour chaque règle $Y \rightarrow u$, on ajoute une règle $X \rightarrow u$. Ensuite on peut supprimer la règle unitaire $X \rightarrow Y$.

Nettoyage d'une Grammaire - Exemple

Prenons l'exemple de la grammaire $(\Sigma = \{a, b\}, V = \{S, A, B\}, S, P)$ ou les règles P sont

$$\begin{aligned} S &\rightarrow aAB \mid BA \mid b \\ A &\rightarrow BBB \mid a \\ B &\rightarrow AB \mid b \mid \varepsilon \end{aligned}$$

Prenons l'exemple de la grammaire ($\Sigma = \{a, b\}$, $V = \{S, A, B\}$, S, P) ou les règles P sont

$$S \rightarrow aAB \mid BA \mid b$$

$$A \rightarrow BBB \mid a$$

$$B \rightarrow AB \mid b \mid \varepsilon$$

Suite au tableau

Une grammaire algébrique est sous Forme Normale de Chomsky (FNC) si toute production est de la forme

$$A \rightarrow BC \quad A, B, C \in V$$

$$A \rightarrow a \quad a \in \Sigma$$

L'arbre de dérivation est donc un arbre binaire.

L'intérêt est théorique, cela permet de simplifier certaines preuves en évitant de devoir faire de nombreux cas.

Forme normale de Chomsky - Algorithme

On commence par transformer la grammaire sous forme réduite et propre.

Forme normale de Chomsky - Algorithme

On commence par transformer la grammaire sous forme réduite et propre.

Pour chaque lettre a dans Σ (i.e. chaque symbole terminal), on ajoute un symbole non productif X_a avec la règle $X_a \rightarrow a$

On commence par transformer la grammaire sous forme réduite et propre.

Pour chaque lettre a dans Σ (i.e. chaque symbole terminal), on ajoute un symbole non productif X_a avec la règle $X_a \rightarrow a$

Dans toutes les anciennes règles, on remplace a par X_a

On commence par transformer la grammaire sous forme réduite et propre.

Pour chaque lettre a dans Σ (i.e. chaque symbole terminal), on ajoute un symbole non productif X_a avec la règle $X_a \rightarrow a$

Dans toutes les anciennes règles, on remplace a par X_a

Toutes les règles sont donc maintenant de la forme $X \rightarrow a$ ou de la forme $X \rightarrow X_1 X_2 \dots X_n$.

Ensuite il suffit de remplacer toute règle $X \rightarrow X_1 X_2 \dots X_n$ avec $n > 2$ par

$$\begin{aligned} X &\rightarrow X_1 Y_1 \\ Y_1 &\rightarrow X_2 Y_2 \\ Y_2 &\rightarrow X_3 Y_3 \\ &\dots \\ Y_{n-2} &\rightarrow X_{n-1} X_n \end{aligned}$$

On commence par transformer la grammaire sous forme réduite et propre.

Pour chaque lettre a dans Σ (i.e. chaque symbole terminal), on ajoute un symbole non productif X_a avec la règle $X_a \rightarrow a$

Dans toutes les anciennes règles, on remplace a par X_a

Toutes les règles sont donc maintenant de la forme $X \rightarrow a$ ou de la forme $X \rightarrow X_1 X_2 \dots X_n$.

Commençons par les deux observations suivantes (faire un dessin!!!)

Lemme

Soit G est une grammaire algébrique ayant p symboles non terminaux. Si la dérivation d'un mot w a profondeur supérieure ou égale à $p + 1$ alors il existe un symbole non terminal X apparaissant deux fois sur une même branche dans la dérivation de w .

Commençons par les deux observations suivantes (faire un dessin !!!)

Lemme

*Soit G est une grammaire algébrique ayant p symboles non terminaux.
Si la dérivation d'un mot w a profondeur supérieure ou égale à $p + 1$ alors il existe un symbole non terminal X apparaissant deux fois sur une même branche dans la dérivation de w .*

Lemme

*Soit G est une grammaire sous forme normale de Chomsky, et soit w un mot généré par cette grammaire.
Si $|w| > 2^{p-1}$ alors tout arbre de dérivation pour ce mot dans cette grammaire a profondeur au moins $p + 1$.*

En combinant les lemmes précédent on peut démontrer :

Théorème (Lemme d'itération pour les langages algébriques)

Si L est algébrique, alors il existe un N tel que pour tout mot $w \in L$ de longueur au moins N , il existe une décomposition $w = xuyvz$ avec

- $uv \neq \varepsilon$
- $|uyv| < N$
- $\forall n \in \mathbb{N}, xu^n y v^n z \in L$

On prend $N = 2^{p-1} + 1$ ou p est le nombre de symboles non terminaux dans une grammaire sous forme de Chomsky qui engendre L .

Commençons par les deux observations suivantes (faire un dessin !!!)

Lemme

*Soit G est une grammaire algébrique ayant p symboles non terminaux.
Si la dérivation d'un mot w a profondeur supérieure ou égale à $p + 1$ alors il existe un symbole non terminal X apparaissant deux fois sur une même branche dans la dérivation de w .*

Lemme

*Soit G est une grammaire sous forme normale de Chomsky, et soit w un mot généré par cette grammaire.
Si $|w| > 2^{p-1}$ alors tout arbre de dérivation pour ce mot dans cette grammaire a profondeur au moins $p + 1$.*

En combinant les lemmes précédent on peut démontrer :

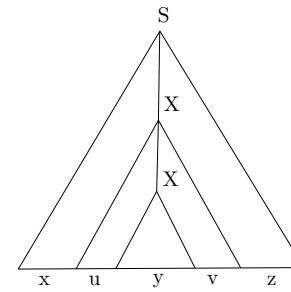
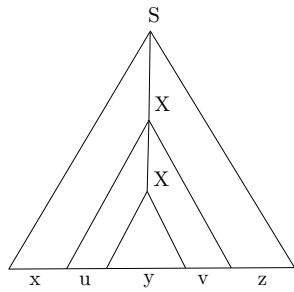
Théorème (Lemme d'itération pour les langages algébriques)

Si L est algébrique, alors il existe un N tel que pour tout mot $w \in L$ de longueur au moins N , il existe une décomposition $w = xuyvz$ avec

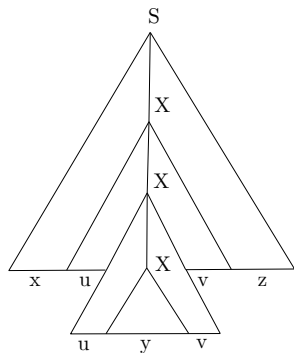
- $uv \neq \varepsilon$
- $|uyv| < N$
- $\forall n \in \mathbb{N}, xu^n y v^n z \in L$

On prend $N = 2^{p-1} + 1$ ou p est le nombre de symboles non terminaux dans une grammaire sous forme de Chomsky qui engendre L .

On trouve donc dans la dérivation de w deux symboles non productifs identiques dans une même branche partant de la racine.



Le symbole X apparaît au moins deux fois sur la même branche.
On considère le sous arbre enraciné en le X le plus bas sur la branche, et on lui substitue le sous arbre enraciné en X plus haut sur la branche. On obtient bien une dérivation du mot $xuuyvz$. On peut réitérer l'opération.



Le symbole X apparaît au moins deux fois sur la même branche.
On considère le sous arbre enraciné en le X le plus bas sur la branche, et on lui substitue le sous arbre enraciné en X plus haut sur la branche. On obtient bien une dérivation du mot $xuuyvz$. On peut réitérer l'opération.

Corollaire

le langage $\{a^n b^n c^n, n \in \mathbb{N}\}$ n'est pas algébrique.

Théoreme

Il existe un algorithme permettant de décider si une grammaire algébrique engendre au moins un mot.

Théoreme

Il existe un algorithme permettant de décider si une grammaire algébrique engendre au moins un mot.

La preuve est similaire à celle du lemme d'itération : si un mot u admet un arbre de dérivation, alors soit aucun symbole non productif n'apparaît plus d'une fois par branche, soit il existe un mot plus court qui admet aussi un arbre de dérivation.

Ainsi si le langage engendré est non vide, il existe un mot engendré tel qu'aucun symbole non productif n'apparaît plus d'une fois par branche.

Théoreme

Il existe un algorithme permettant de décider si une grammaire algébrique engendre au moins un mot.

La preuve est similaire à celle du lemme d'itération : si un mot u admet un arbre de dérivation, alors soit aucun symbole non productif n'apparaît plus d'une fois par branche, soit il existe un mot plus court qui admet aussi un arbre de dérivation.

Théoreme

Il existe un algorithme permettant de décider si une grammaire algébrique engendre au moins un mot.

La preuve est similaire à celle du lemme d'itération : si un mot u admet un arbre de dérivation, alors soit aucun symbole non productif n'apparaît plus d'une fois par branche, soit il existe un mot plus court qui admet aussi un arbre de dérivation.

Ainsi si le langage engendré est non vide, il existe un mot engendré tel qu'aucun symbole non productif n'apparaît plus d'une fois par branche.

Il y a un nombre fini de tels arbres, on peut tous les énumérer et voir si un produit un mot.