

6. Grammaires formelles

Enseignant: Arnaud Casteigts

Assistants: A.-Q. Berger & M. Marseloo

Moniteurs: N. Beghdadi & E. Bussod

Nous entamons une partie du cours à cheval entre la linguistique et l'informatique, à savoir l'étude des grammaires formelles. Ces dernières permettent de décrire des langages de manière générative, par l'utilisation de *règles de production*. Les grammaires sont plus puissantes que les expressions régulières et les automates finis, elles permettent notamment de représenter des langages non réguliers (et bien plus).

6.1 Grammaires formelles en général

Intuitivement, une grammaire formelle est un ensemble de règles qui permettent de créer des motifs textuels par remplacement successifs. Chaque règle est de la forme $\alpha \rightarrow \beta$ où α est le motif à remplacer (appelé **partie gauche** de la règle) et β est le motif qui va remplacer α (appelé **partie droite**). Voici un exemple de grammaire comportant deux règles de production :

Grammaire G_1 :

$$S \rightarrow \mathbf{a}S\mathbf{b}$$

$$S \rightarrow \varepsilon$$

Il y a deux types de symboles dans une grammaire : des **symboles terminaux** et des **variables** (aussi appelés symboles non-terminaux). Les symboles terminaux sont généralement représentés en minuscule, comme ici **a** et **b** et sont issus d'un alphabet Σ comme ceux que l'on connaît déjà. Les variables sont représentées en majuscule, comme ici S , et font partie d'un autre alphabet V . Ces variables ont vocation à être remplacées ultérieurement. La grammaire spécifie aussi quel est le symbole initial à partir duquel on va générer tous les mots, ici S également.

La génération consiste donc à partir du symbole initial et d'appliquer des règles jusqu'à ce qu'il n'y ait que des symboles terminaux. On a alors **généré** (ou **engendré**) un mot du langage. Lorsque plusieurs règles sont applicables, on peut utiliser l'une ou l'autre.

Dans cet exemple, nous partons de S et nous avons le choix d'appliquer la première ou la seconde règle. Si nous utilisons la seconde règle directement, cela produit le mot vide ε et nous avons terminé car il n'y a plus de variable à remplacer. Si nous utilisons plutôt la première règle, cela produit le mot intermédiaire **aSb**. Nous pouvons alors recommencer en

remplaçant le S de ce mot. Si on utilise ici la seconde règle, qui remplace le S du milieu par ε , cela donnera le mot \mathbf{ab} et nous aurons terminé. Si on continue à utiliser la première règle, cela donnera le mot \mathbf{aaSb} , et ainsi de suite. Quel est l'ensemble de tous les mots que l'on peut obtenir? On se convaincra assez facilement qu'il s'agit des mots $\varepsilon, ab, aabb, aaabbb, \dots$, autrement dit, le langage $L(G_1) = \{a^n b^n \mid n \in \mathbb{N}\}$, qui nous le savons, n'est pas régulier! (c.f. cours précédent).

Dans leur version la plus générale, les grammaires formelles peuvent utiliser des règles de productions qui transforment n'importe quel motif (combinant variables et symboles terminaux) en n'importe quel autre motif. Ces grammaires étant très expressives, on se restreint souvent à des versions particulières qui sont moins puissantes. Le type de grammaire le plus répandu est celui des grammaires *hors-contexte*, qui sont très utilisées en théorie des langages de programmation (notamment pour le parsing et la compilation), ainsi que dans les protocoles réseaux.

6.2 Cas particulier : Grammaires hors-contexte

Les grammaires **hors-contexte** (aussi appelées grammaires non contextuelles) sont des cas particuliers de grammaires formelles, qui imposent que la partie gauche de chaque règle consiste en une variable seulement (et rien d'autre), c'est le cas de l'exemple donné plus haut. On les appelle *hors-contexte* parce qu'elles ne tiennent pas compte de ce qu'il y a autour de la variable à remplacer.

Formellement, une grammaire hors-contexte est définie par :

- Un alphabet V de variables (symboles non terminaux),
- Un alphabet Σ de symboles terminaux,
- Un symbole de départ dans V ,
- Un ensemble fini de règles de production \mathcal{P} de la forme $\alpha \rightarrow \beta$, où $\alpha \in V$ et $\beta \in (V \cup \Sigma)^*$.

Par exemple, la grammaire G_1 donnée plus haut correspond à $V = \{S\}, \Sigma = \{\mathbf{a}, \mathbf{b}\}$, le symbole de départ $S \in V$ et les règles de productions $\mathcal{P} = \{S \rightarrow \mathbf{aSb}, S \rightarrow \varepsilon\}$.

Lorsque plusieurs règles ont la même partie gauche, on peut les regrouper en utilisant une barre verticale, comme $S \rightarrow \mathbf{aSb} \mid \varepsilon$, qui se lit " S donne \mathbf{aSb} ou ε ". (Il s'agit quand même de deux règles différentes.) Et pour être encore plus compact, on écrit parfois directement : $G_1 = (\{S\}, \{\mathbf{a}, \mathbf{b}\}, S, \{S \rightarrow \mathbf{aSb} \mid \varepsilon\})$.

L'application d'une ou plusieurs règles d'une grammaire est appelé une **dérivation**. Voici une dérivation possible de la grammaire G_1 :

$$S \Rightarrow \mathbf{aSb} \Rightarrow \mathbf{aaSb} \Rightarrow \mathbf{aaSbbb} \Rightarrow \mathbf{aaabbb}.$$

Notez la différence entre les symboles \rightarrow (spécification de la règle, dans la définition de la grammaire) et \Rightarrow (son application sur un mot). Nous appellerons **mot intermédiaire**

les mots qui contiennent au moins une variable, et **mot terminal** ceux qui n'ont que des symboles terminaux. S'il existe une dérivation d'un mot w_1 vers un mot w_2 (terminal ou non), on note $w_1 \Rightarrow^* w_2$ et on dit qu'on peut dériver w_2 depuis w_1 . L'ensemble des mots terminaux que l'on peut dériver depuis le symbole de départ définit le **langage engendré** par la grammaire. Autrement dit, le langage $L(G)$ engendré par une grammaire G avec symbole de départ S est :

$$L(G) = \{w \in \Sigma^* \mid S \Rightarrow^* w\}$$

Un langage est **hors-contexte** s'il existe une grammaire hors-contexte (GHC) qui l'engendre.

6.3 Exemples et autres notions

Voici quelques exemples de grammaires hors contexte. Nous en profitons pour illustrer quelques notions supplémentaires, comme les *arbres de dérivation*, les *dérivations gauches* (ou *droite*), et les *grammaires ambiguës*.

6.3.1 Mots binaires se terminant par 0 (nombres pairs)

On peut représenter le langage des mots binaires se terminant par 0 par une grammaire hors-contexte $G_2 = (\{S\}, \{0, 1\}, S, \{S \rightarrow 0S \mid 1S \mid 0\})$

Exemple de dérivation : $S \Rightarrow 1S \Rightarrow 10S \Rightarrow 101S \Rightarrow 1010$

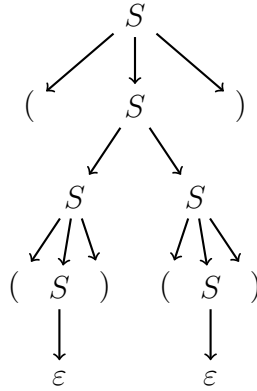
On peut se convaincre assez facilement que tout mot binaire se terminant par 0 peut être généré par cette grammaire. Autrement dit, $L(G_2)$ est le langage régulier correspondant à l'expression $(0 \cup 1)^*0$. Comme déjà évoqué, les grammaires hors-contexte peuvent générer des langages non réguliers, mais elles peuvent bien sûr aussi générer des langages réguliers ! (Qui peut le plus, peut le moins...)

6.3.2 Mots bien parenthésés

On peut représenter l'ensemble des mots bien parenthésés par une grammaire $G_3 = (V = \{S\}, \Sigma = \{(\, ,)\}, S, \mathcal{P} = \{S \rightarrow (S) \mid SS \mid \varepsilon\})$. Cette grammaire produit tous les mots constitués de parenthèses correctement entrelacées, p.ex. $()$, $((\,))$, $(\,)(\,)$, $((\,)(\,))$, $((\,))(\,)$, ... Elle ne produit pas, en revanche, les mots tels que $(\,(\,)$, ou $(\,(\,(\,)$, ou encore $(\,)(\,)(\,)$.

Exemple de dérivation : $S \Rightarrow (S) \Rightarrow (SS) \Rightarrow ((S)S) \Rightarrow ((\,)S) \Rightarrow ((\,)(S)) \Rightarrow ((\,)(\,))$

On peut représenter la dérivation d'un mot par un **arbre de dérivation**, comme suit :



En fait, cet arbre représente *plusieurs* dérivations possibles, mais qui sont toutes équivalentes. Par exemple, dans la dérivation précédente, nous avons systématiquement utilisé la variable la plus à gauche pour dériver. Une telle dérivation est appelée **dérivation gauche**. Elle correspond à un parcours en profondeur de la gauche vers la droite dans l'**arbre de dérivation**. Mais nous aurions aussi pu effectuer une **dérivation droite**, correspondant à un parcours en profondeur de la droite vers la gauche comme suit :

$$S \Rightarrow (S) \Rightarrow (SS) \Rightarrow (S(S)) \Rightarrow (S()) \Rightarrow ((S)()) \Rightarrow (()())$$

Ou encore prendre les variables dans un ordre arbitraire... La notion d'arbre de dérivation est importante et mérite d'être bien comprise. Nous y reviendrons plus bas.

Notez que $L(G_3)$ n'est pas un langage régulier. Vous le prouverez peut-être en exercices en utilisant le lemme de l'étoile. En revanche, c'est bien évidemment un langage hors-contexte, puisque la grammaire qui le génère est hors-contexte.

6.3.3 Expressions arithmétiques

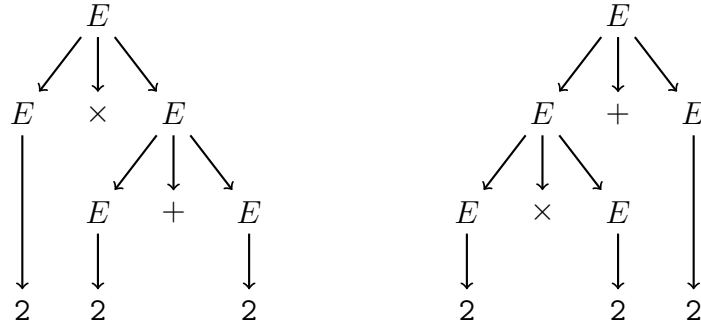
Voici un exemple de règles qui produisent des expressions arithmétiques avec les opérations $+$ et \times à partir du chiffre 2, l'alphabet terminal étant $\Sigma = \{\times, +, 2, (,)\}$:

$$E \rightarrow E + E \mid E \times E \mid (E) \mid 2$$

Par exemple, voici une dérivation gauche qui produit le mot $2 \times (2 + 2)$:

$$E \Rightarrow E \times E \Rightarrow 2 \times E \Rightarrow 2 \times (E) \Rightarrow 2 \times (E + E) \Rightarrow 2 \times (2 + E) \Rightarrow 2 \times (2 + 2).$$

Comme précédemment, nous aurions pu aussi utiliser une dérivation droite, ces dérivations revenant in fine au même arbre de dérivation. Cependant, il est cette fois possible qu'un même mot admette *deux arbres* de dérivations différents ! Prenons l'exemple de $2 \times 2 + 2$ (sans les parenthèses). On peut l'obtenir via les arbres :



Si cela peut se produire pour au moins un mot, alors on dit que ce mot est ambigu et que la grammaire elle-même est une **grammaire ambiguë**. En général, on essaie d’éviter cela, par exemple, dans les langages de programmation, il est préférable qu’il n’y ait qu’une façon possible d’analyser syntaxiquement (= de “parser”) le code.

Le langage humain est naturellement ambigu, p.ex : “j’ai vu sa main avec un miroir” peut signifier “j’ai vu [sa main avec un miroir]” ou “[j’ai vu sa main] avec un miroir”.

Dans certains cas, une grammaire peut être modifiée pour être rendue non-ambigüe, mais ce n’est pas toujours possible. Ici, c’est possible, la grammaire suivante engendre le même langage, sans ambiguïté (au prix d’utiliser trois variables au lieu d’une) :

$$\begin{aligned}
 E &\rightarrow E + T \mid T \\
 T &\rightarrow T \times F \mid F \\
 F &\rightarrow (E) \mid 2
 \end{aligned}$$

6.3.4 Langage humain ?

On pourrait utiliser des variables du genre $V = \{\text{PHRASE, GROUPE NOMINAL, GROUPE VERBAL, COMPLEMENT, NOM, VERBE}\}$ et des terminaux $\Sigma = \{\text{a...z, } _ \}$ avec des règles du type :

$$\begin{aligned}
 \text{PHRASE} &\rightarrow \text{GROUPE NOMINAL GROUPE VERBAL} \\
 \text{GROUPE VERBAL} &\rightarrow \text{VERBE COMPLEMENT} \\
 \text{VERBE} &\rightarrow \text{manger, dormir, ...} \\
 &\dots
 \end{aligned}$$

Ce type de grammaires est utilisée en linguistique pour modéliser les langues naturelles. Elles ont longtemps servi de base au traitement automatique des langues (pour la traduction, la rédaction, etc.). Elles ont récemment été supplantées par les méthodes statistiques (réseaux de neurones, notamment), mais peuvent continuer à être utilisées en conjonction avec ces dernières.

6.4 Cas particulier : grammaires régulières

(deuxième semaine)

Les grammaires **régulières** sont un cas particulier de grammaires hors-contexte, qui imposent une restriction supplémentaire sur la forme des règles. La partie gauche n'est pas affectée (on a toujours $\alpha \in V$), mais la partie droite β est plus restreinte. En fait, on distingue deux types de grammaires régulières : les grammaires régulières à gauche et les grammaires régulières à droite (les deux ayant la même expressivité). Les **grammaires régulières à gauche** imposent que pour toute règle $\alpha \rightarrow \beta$, la partie droite β doit être de la forme Xs , ou bien s , ou bien ε , avec $X \in V$ et $s \in \Sigma$. Pour les **grammaires régulières à droite**, c'est l'inverse : β est de la forme sX , ou s , ou ε . Autrement dit, chaque règle ne peut générer qu'au plus une variable et au plus un symbole terminal, qui sera toujours *du même côté*.

Par exemple, la grammaire G_2 (déjà vue plus haut) qui génère les mots binaires se terminant par 0 est une grammaire régulière. Ses trois règles sont $S \rightarrow 0S \mid 1S \mid 0$, ce qui respecte bien les contraintes d'une grammaire régulière (en l'occurrence, à droite). On aurait aussi pu générer ce langage avec une grammaire régulière à gauche, en utilisant deux variables et quatre règles comme suit : $G'_2 = (\{S, T\}, \{0, 1\}, S, \{S \rightarrow T0, T \rightarrow T0 \mid T1 \mid \varepsilon\})$.

6.4.1 Équivalence entre grammaires régulières et langages réguliers

Les grammaires régulières correspondent exactement aux langages réguliers. Cela peut se démontrer (comme souvent) par des transformations. En l'occurrence, n'importe quel AFD peut être converti en grammaire régulière et n'importe quelle grammaire régulière peut être convertie en AFN (et donc en AFD, par détermination).

Transformation d'un AFD vers une grammaire régulière

Soit un AFD $A = (Q, \Sigma, \delta, q_0, F)$. On peut construire une grammaire régulière à droite $G = (V, \Sigma, S, \mathcal{P})$ qui génère exactement $L(A)$, en la définissant comme suit :

- Les variables correspondent aux états de l'automate ($V = Q$),
- Les symboles terminaux correspondent à l'alphabet de l'automate,
- Le symbole de départ S correspond à l'état initial q_0 ,
- Pour chaque transition (q_i, s, q_j) , on crée une règle $q_i \rightarrow sq_j$,
- Pour chaque état final $q \in F$, on crée une règle $q \rightarrow \varepsilon$.

Il existe une construction similaire pour transformer l'automate A en une grammaire régulière à gauche.

Transformation d'une grammaire régulière vers un AFN

Soit une grammaire régulière à droite $G = (V, \Sigma, S, \mathcal{P})$, on peut construire un AFN $A = (Q, \Sigma, \delta, q_0, F)$ qui reconnaît exactement $L(G)$, en le définissant comme suit :

- Les états de l'automate correspondent aux variables V + un état final q_f ,
- L'état initial q_0 correspond à la variable de départ S ,
- L'alphabet de l'automate correspond aux symboles terminaux de la grammaire,
- Pour chaque règle de la forme $A \rightarrow sB$, on ajoute une transition depuis l'état correspondant à A vers l'état correspondant à B avec le symbole s ,
- Pour chaque règle de la forme $A \rightarrow s$, on ajoute une transition depuis l'état correspondant à A vers q_f avec le symbole s ,
- Pour chaque règle de la forme $A \rightarrow \varepsilon$, on ajoute une ε -transition depuis l'état correspondant à A vers q_f .

Il existe une construction similaire pour transformer une grammaire régulière à gauche vers un AFN.

6.5 Hiérarchie de Chomsky (aperçu)

La hiérarchie de Chomsky met en relation quatre types de grammaires différentes avec quatre types de machines capable de connaître les langages correspondants. Traduction : context-free = hors-contexte ; push-down automaton = automate à pile.

