

9. Lemme de l'étoile (hors-contexte)

Enseignant: Arnaud Casteigts

Assistant: Alexandre-Quentin Berger

Nous avons déjà vu comment montrer qu'un langage n'est pas régulier en utilisant le lemme de l'étoile (*c.f.* Cours 5). Dans ce cours, nous allons utiliser une démarche analogue pour les langages hors-contextes, en utilisant une autre version du lemme de l'étoile.

9.1 Rappel du lemme de l'étoile pour les langages réguliers

Le lemme de l'étoile pour les langages *réguliers* établit une propriété que tout langage régulier doit satisfaire. Intuitivement, il établit que tout mot dépassant une certaine longueur doit contenir une partie (un *facteur*) qui peut être répété un nombre arbitraire de fois, tout en restant dans le langage. Nous le rappelons ici, en changeant le nom des variables pour mieux coller à la version d'aujourd'hui.

Lemme 9.1 (Lemme de l'étoile pour les langages réguliers). *Si un langage L est régulier, alors il existe une longueur k au delà de laquelle tout mot $z \in L$ peut s'écrire sous la forme uvw avec :*

1. $|v| > 0$,
2. $|uv| \leq k$,
3. $uv^i w \in L$ pour tout $i \in \mathbb{N}$.

Pour montrer qu'un langage L n'est pas régulier, il suffit donc de trouver un mot $z \in L$ tel que pour toute décomposition $z = uvw$, la répétition du facteur v peut nous faire sortir de L . Nous avons donné l'exemple du langage $L = \{a^n b^n \mid n \in \mathbb{N}\}$.

9.2 Lemme de l'étoile pour les langages hors-contextes

Un lemme similaire existe pour les langages hors-contextes, le voici :

Lemme 9.2 (Lemme de l'étoile pour les LHC). *Si L est un langage hors-contexte, alors il existe une longueur k au delà de laquelle tout mot $z \in L$ peut s'écrire $z = uvwxy$ avec :*

1. $|v| + |x| > 0$
2. $|vwx| \leq k$
3. $uv^i wx^i y \in L$ pour tout $i \in \mathbb{N}$.

En effet, si L est hors-contexte, on peut toujours faire cela. Prenons l'exemple du langage $L = \{a^n b^n \mid n \in \mathbb{N}\}$ (qui est hors-contexte) et un mot quelconque de L , par exemple $z = aaabbb$ (en l'occurrence, $k = 2$ pour ce langage). On peut décomposer $z = aa \cdot a \cdot \varepsilon \cdot b \cdot bb$, autrement dit $u = aa, v = a, w = \varepsilon, x = b$ et $y = bb$. On a bien $|vwx| \leq k$ et $|vx| > 0$. Enfin, il est facile de voir que $uv^i wx^i y \in L$ pour tout $i \in \mathbb{N}$.

Si L est un langage hors-contexte, alors on doit pouvoir faire cela pour *tous les mots* de L de longueur $\geq k$ (c'est le cas ici). Le lemme peut donc être utilisé dans l'autre sens : si l'on veut montrer qu'un langage L' n'est pas hors-contexte, il suffit de montrer qu'il existe un mot de L' pour lequel cela ne marche pas.

9.2.1 Exemple d'utilisation

Utilisons le lemme de l'étoile pour montrer que le langage $L = \{a^n b^n c^n \mid n \in \mathbb{N}\}$ sur l'alphabet $\Sigma = \{a, b, c\}$ n'est pas hors-contexte.

Par l'absurde, supposons d'abord que L est hors-contexte, le lemme nous dit qu'il existe un k tel que tout mot z de longueur $\geq k$ est décomposable en facteurs $uvwxy$ avec les trois propriétés vraies. Nous allons montrer que cela ne marche pas. Prenons par exemple le mot $z = a^k b^k c^k \in L$ et décomposons-le de sorte que v et x ne sont pas vides tous les deux (propriété 1 du lemme), il y a quatre possibilités pour la partie vwx :

- vwx ne contient que des a
- vwx ne contient que des b
- vwx ne contient que des c
- vwx contient deux types de symboles : a et b , ou b et c (il ne peut pas en contenir trois, car la décomposition satisfait $|vwx| \leq k$ (propriété 2), ce qui est trop petit pour couvrir trois symboles différents dans le mot choisi).

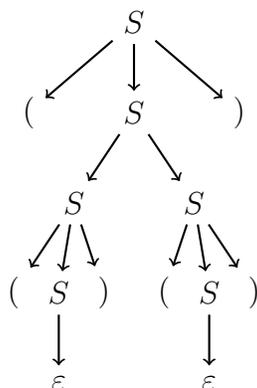
Examinons maintenant les conséquences de la répétition des facteurs v et x (propriété 3 du lemme). Si vwx ne contient que des a , on obtiendra un mot qui a trop de a (donc pas dans L , contradiction). Idem pour b et idem pour c . Reste le dernier cas, mais ici le symbole qui n'apparaît pas dans vwx se retrouvera en trop petite quantité. Dans tous les cas, on arrive donc à une contradiction, ce qui implique que L n'est pas hors-contexte.

9.2.2 Démonstration du lemme lui-même

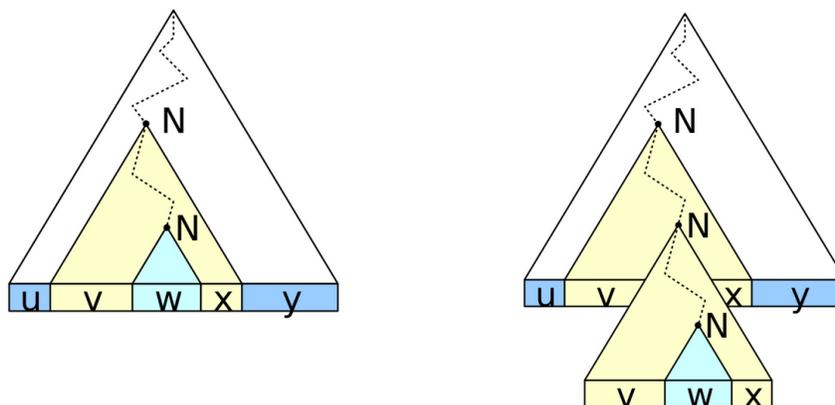
Nous allons maintenant donner les idées principales du lemme lui-même. Il n'est pas nécessaire de connaître la preuve en détail, mais une bonne compréhension de ces idées vous aidera certainement à bien l'utiliser (ce qui est l'objectif principal de ce cours).

Par définition, si un langage L est hors-contexte, alors il peut être engendré par une grammaire hors-contexte G . Si un mot w appartient à L , on peut donc l'obtenir par une dérivation de G et on peut également l'associer à un arbre de dérivation.

Par exemple, pour la grammaire $S \rightarrow (S) \mid SS \mid \varepsilon$, qui produit les mots bien parenthésés (c.f. Cours 6), le mot “ $((()))$ ” correspond à l’arbre suivant :



Dans cet exemple, il n’y a qu’une variable dans la grammaire (en l’occurrence, S), mais en général, une grammaire hors-contexte peut en avoir plusieurs. L’idée du lemme de l’étoile est la suivante : si un mot z est suffisamment long, alors pour n’importe quel arbre de dérivation qui lui correspond, au moins une variable doit être *répétée* dans un chemin de la racine vers les feuilles (et produire des symboles terminaux entre-temps). C’est clairement le cas dans l’exemple précédent, on peut aussi imaginer un cas plus général, par exemple l’arbre de dérivation ci-dessous (à gauche) pour une grammaire qui aurait plusieurs variables, dont la variable N :



Un mot z du langage peut alors être décomposé en facteurs $uvwxy$ tels que u correspond à la partie du mot produite à gauche de la première variable N (idem pour y à droite), et vw est la partie qui dérive de la première variable N (w étant la partie qui dérive de la deuxième occurrence de la variable N). Puisque la variable N est capable de se reproduire elle-même, elle pourrait également le faire un nombre arbitraire de fois avant de produire w , ce qui correspond bien aux mots de la forme uv^iwx^iz (dessin de droite). Ces mots doivent donc aussi appartenir au langage engendré par cette grammaire.