

Feuille d'exercices n°12 - Grammaires non contextuelles

Notions abordées

- exemple de grammaires non contextuelles
- suite de dérivations, arbre de dérivation
- langage d'une grammaire non contextuelle, preuve par induction
- exemple de langages non contextuels : langage de Dyck, mots de Łukasiewicz
- grammaire décrivant des listes, décrivant des formules sous FNC
- régularité des langages non contextuels
- grammaire propre

Exercice 1 : Exemples de dérivations

Soit $\Sigma = \{a, b\}$, $\mathcal{V} = \{S, D, T, X\}$, et \mathcal{G} la grammaire non contextuelle de symbole initial S , et contenant les règles de production ci-contre ♣.

$$\begin{array}{l} S \rightarrow XSX \mid D \\ D \rightarrow aTb \mid bTa \\ T \rightarrow XTX \mid X \mid \varepsilon \\ X \rightarrow a \mid b \end{array}$$

Q. 1 Parmi les dérivations suivantes, lesquelles sont vraies ?

$$(a) T \Rightarrow aba \quad (b) T \xrightarrow{*} aba \quad (c) T \Rightarrow T \quad (d) T \xrightarrow{*} T$$

Q. 2 Parmi les dérivations suivantes, lesquelles sont vraies ?

$$(a) XXX \xrightarrow{*} aba \quad (c) T \xrightarrow{*} XX \quad (e) D \xrightarrow{*} \varepsilon \\ (b) X \xrightarrow{*} aba \quad (d) T \xrightarrow{*} XXX$$

Q. 3 Donner 3 mots de $\mathcal{L}(\mathcal{G})$. Justifier leur appartenance à $\mathcal{L}(\mathcal{G})$ par une dérivation.

Q. 4 Donner 3 mots de Σ^* qui ne sont pas dans $\mathcal{L}(\mathcal{G})$.

Q. 5 Donner une description en français de $\mathcal{L}(\mathcal{G})$. En première lecture de cet énoncé, on ne demande pas de preuve du résultat avancé.

Exercice 2 : Arbres de dérivations

Soit $\Sigma = \{a, b, c, \dots, z, +, (,)\}$, $\mathcal{V} = \{S, P, F, T\}$, et la grammaire \mathcal{G} de symbole initial S et de règles de production ci-contre.

$$\begin{array}{l} S \rightarrow S + P \mid P \\ P \rightarrow PF \mid F \\ F \rightarrow (S) \mid T \\ T \rightarrow a \mid b \mid c \mid \dots \mid z \end{array}$$

Q. 1 Donner un arbre de dérivation dont la production est $ab + bc$.

Q. 2 Donner un arbre de dérivation de $b + (a + b)c + a(bc)$.

Q. 3 Que représente cette grammaire, i.e. que représentent les mots engendrés par \mathcal{G} ?

♣. Lorsque plusieurs règles de production ont le même membre gauche, on peut les représenter factorisées en séparant les différents membres droits par $|$. Par exemple, on note $A \rightarrow w|w'$ pour les deux règles $A \rightarrow w$ et $A \rightarrow w'$.

Exercice 3 : Construction de grammaires

- Q. 1** Soit $\Sigma = \{a, b\}$. Donner une grammaire non contextuelle \mathcal{G}_i de langage L_i ci-dessous.
- $L_0 = \{w \in \Sigma^* \mid |w| \geq 3\}$;
 - $L_1 = \{w \in \Sigma^* \mid |w|_a \geq 3\}$;
 - $L_2 = \{w \in \Sigma^+ \mid w \text{ commence et finit par la même lettre}\}$;
 - $L_3 = \emptyset$;
- Q. 2** Soit $\Sigma = \{a, b\}$. Donner une grammaire non contextuelle \mathcal{G}_i de langage L_i ci-dessous♣.
- $L_4 = \{w \in \Sigma^* \mid |w| \equiv 1[2]\}$;
 - $L_5 = \{w \in \Sigma^* \mid |w| \equiv 1[2] \text{ et sa lettre du milieu est } a\}$;
 - $L_6 = \{w \in \Sigma^* \mid w \text{ est un palindrome}\}$;
 - $L_7 = L_6^C$;
- Q. 3** Donner une grammaire non contextuelle \mathcal{G}_i de langage L_i ci-dessous.
- $L_8 = \{a^n b^n \mid n \in \mathbb{N}\}$;
 - $L_9 = \{a^n b^m \mid (n, m) \in \mathbb{N}^2, m \geq n\}$;
 - $L_{10} = \{a^n b^m \mid (n, m) \in \mathbb{N}^2, m \geq 2n\}$;
- Q. 4** Soit $\Sigma = \{a, b, \#\}$. Donner une grammaire non contextuelle \mathcal{G}_{13} de langage L_{13} ci-dessous.
- $L_{10} = \{w\#x \mid (w, x) \in (\{a, b\}^*)^2, \bar{w}^R \text{ est un sous-mot de } x \text{ où } \bar{w}^R \text{ désigne le miroir de } w\}$.

Exercice 4 : Langage de Dyck

Le langage des mots de Dyck est le langage des mots sur $\Sigma = \{(\,)\}$ bien parenthésés, c'est-à-dire le plus petit langage $L \subseteq \Sigma^*$ tel que :

- $\varepsilon \in L$;
- $\forall u \in L, (u) \in L$;
- $\forall (u, v) \in L \times L, u \cdot v \in L$.

- Q. 1** Donner une grammaire non contextuelle \mathcal{G} reconnaissant le langage de Dyck.
- Q. 2** Montrer que les mots u de $\mathcal{L}(\mathcal{G})$ sont tels que :
- $|u|_{(} = |u|_{)}$
 - pour tout préfixe v de u , $|v|_{(} \geq |v|_{)}$
- Q. 3** Montrer que les mots $u \in \Sigma^*$ vérifiant ces deux conditions sont bien parenthésés.
- Q. 4** En déduire une fonction OCAML testant si un mot est dans $\mathcal{L}(\mathcal{G})$.
- Q. 5** Démontrer que ce langage n'est pas régulier.

Exercice 5 : Reasonner par induction sur une grammaire

Soit $\Sigma = \{a, b\}$, $\mathcal{V} = \{S\}$ et la grammaire \mathcal{G} contenant les règles de production :

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

- Q. 1** Montrer que ba n'est sous-mot d'aucun mot de $\mathcal{L}(\mathcal{G})$.
- Q. 2** Donner $\mathcal{L}(\mathcal{G})$.
- Q. 3** Donner une grammaire reconnaissant $\mathcal{L}(\mathcal{G})$ et assurant que chaque mot est obtenu par une unique suite de dérivations immédiates depuis le symbole initial.

♣. Lorsque la définition d'un langage fait intervenir le complémentaire, celui-ci s'entend par rapport à Σ^* .

Exercice 6 : Ambiguïté

Soit $\Sigma = \{a, b\}$, $\mathcal{V} = \{S\}$, et la grammaire \mathcal{G} de symbole initial S et de règles de production :

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

Q. 1 Montrer que \mathcal{G} est ambiguë.

Q. 2 Montrer que $\mathcal{L}(G) = \{v \in \Sigma^* \mid \text{pour tout } u \text{ préfixe de } v, |u|_a \geq |u|_b\}$.

Q. 3 Donner une grammaire non ambiguë reconnaissant $\mathcal{L}(G)$.

Exercice 7 : Listes OCAML

On considère l'alphabet $\Sigma = \{[,], ;, \text{true}, \text{false}\}$. Donner une grammaire non contextuelle dont le langage est l'ensemble des listes OCAML de booléens♣.

Exercice 8 : Forme normale conjonctive

Soit un ensemble fini de variables propositionnelles $\mathcal{Q} = \{p, q, \dots\}$. Soit \mathcal{O} l'ensemble des symboles $\{\&, |, -\}$ ♥. Soit l'alphabet $\Sigma = \mathcal{Q} \cup \mathcal{O}$. Soit finalement le langage FNC sur l'alphabet Σ des mots qui décrivent une formule de la logique propositionnelle sous forme normale conjonctive.

Par exemple FNC contient :

- ε qui représente la conjonction vide ;
- p qui représente la conjonction réduite à une disjonction, elle-même réduite au seul littéral p ;
- $p|-q$ qui représente la conjonction réduite à une disjonction : $(p \vee \neg q)$;
- $p|q\&r|-s\&-p$ qui représente la forme normale conjonctive : $(p \vee q) \wedge (r \vee \neg s) \wedge \neg p$

Q. 1 Donner une grammaire non contextuelle non ambiguë \mathcal{G} de langage FNC.

On se munit des types OCAML suivants, permettant la représentation des formes normales conjonctives en OCAML.

```
1 | type var = char           (* une variable *)
2 | type lit = var * bool     (* un littéral *)
3 | type cls = lit list      (* une clause *)
4 | type fnc = cls list      (* une FNC *)
```

♣. On ne décrit pas l'ensemble des expressions OCAML de type `bool list` car on limite les expressions booléennes autorisées à `true` et `false` seulement

♥. Jouant respectivement les rôles de $\{\wedge, \vee, \neg\}$

Q. 2 Grâce aux types précédemment définis, donner les expressions OCAML représentant les éléments suivants.

- la variable p
- le littéral p
- le littéral $\neg q$
- la clause $p \vee q$
- la formule $p \vee q$
- $p \wedge q$
- $p \wedge (q \vee \neg p)$

Q. 3 Donner 5 fonctions OCAML mutuellement récursives :

- parse_f (s: string) (start: int) (len: int) : fnc
- parse_c (s: string) (start: int) (len: int) : fnc
- parse_d (s: string) (start: int) (len: int) : cls
- parse_l (s: string) (start: int) (len: int) : lit
- parse_v (s: string) (start: int) (len: int) : char

prenant en arguments une chaîne de caractères s , un indice de début $start$ et une longueur len et retournant la forme normale conjonctive (resp. la FNC non vide, la clause disjonctive, le littéral, la variable) représentée par la chaîne de caractères `String.sub s i len`. On déclenchera une erreur dans l'éventualité où la chaîne de caractères n'est pas de la forme attendue.

Exemples : `|| (parse_v "p|q&r|-s&-p" 2 1) s'évalue en 'q'`
`|| (parse_l "p|q&r|-s&-p" 9 2) s'évalue en ('p', false)`
`|| (parse_l "p|q&r|-s&-p" 4 4) s'évalue en [('r', true); ('s', false)]`

Exercice 9 : Un lemme d'itération

Dans cet exercice on s'intéresse à un équivalent du lemme de l'étoile pour les langages non contextuels. Ce résultat donne des contraintes que doivent vérifier les langages engendrés par des grammaires non contextuelles, ainsi on pourra établir que certains langages ne peuvent être engendrés par des grammaires non contextuelles en montrant qu'ils ne vérifient pas lesdites propriétés.

Soit une grammaire $\mathcal{G} = (\mathcal{V}, \Sigma, P, S)$. On suppose dans la suite que $P \neq \emptyset$. Soit donc $p = \max(\{1\} \cup \{|w| \mid V \rightarrow w \in P\})$ la longueur du plus long mot apparaissant à droite d'une règle, ou 1 si tous les mots à droite des règles sont le mot vide.

Q. 1 Donner et prouver une majoration sur la longueur des mots de $(\Sigma \cup \mathcal{V})^*$ admettant un arbre de dérivation de hauteur h .

Q. 2 En déduire qu'il existe un entier N tel que si w est un mot de $\mathcal{L}(\mathcal{G})$ de longueur $\geq N$, alors tout arbre de dérivation de w admet une branche contenant au moins deux fois le même symbole non terminal.

Q. 3 En déduire qu'il existe un entier N^\clubsuit tel que pour tout mot $w \in \mathcal{G}$ tel que $|w| \geq N$, il existe des mots u, x, c, y et v de Σ^* tels que :

- i) $w = uxcyv$
- ii) $|xy| \geq 1$
- iii) $|xcy| \leq N$
- iv) $\forall n \in \mathbb{N}, ux^n cy^n v \in \mathcal{L}(\mathcal{G})$

Q. 4 En déduire que le langage $L = \{a^n b^n c^n \mid n \in \mathbb{N}\}$ n'est pas engendré par une grammaire non contextuelle.

Q. 5 En déduire que l'ensemble des langages non contextuels n'est stable ni par intersection, ni par complémentaire.

\clubsuit . évidemment cet entier dépend de la grammaire fixée en début d'exercice

Exercice 10 : Grammaires propres

Grammaires propres. $\mathcal{G} = (\mathcal{V}, \Sigma, P, S)$ est dite **propre** dès lors que P ne contient aucune règle de la forme :

1. $V \rightarrow \varepsilon$ avec $V \in \mathcal{V}$
2. $V \rightarrow V'$ avec $(V, V') \in \mathcal{V}^2$

- Q. 1** Proposer un algorithme prenant en argument une grammaire \mathcal{G} et calculant une grammaire reconnaissant le langage $\mathcal{L}(\mathcal{G}) \setminus \{\varepsilon\}$ et ne contenant aucune règle de la forme 1. ci-dessus.
- Q. 2** Proposer un algorithme prenant en argument une grammaire \mathcal{G} ne contenant pas de règles de production de la forme 1. ci-dessus et la transformant en une grammaire ne contenant aucune règle de la forme 2. ci-dessus (on s'efforcera de ne pas produire une grammaire contenant des règles de la forme 1.). Conclure.
- Q. 3** Démontrer que dans une grammaire contextuelle propre $\mathcal{G} = (\mathcal{V}, \Sigma, P, S)$, si $S \Rightarrow w_1 \Rightarrow w_2 \Rightarrow w_3 \dots \Rightarrow w_n$ alors $|w_n|_{\mathcal{V}} + 2|w_n|_{\Sigma} \geq 1 + n^{\clubsuit}$.
- Q. 4** En déduire que l'ensemble des langages des grammaires non contextuelles est décidable.

Exercice 11 : Les langages réguliers sont non contextuels

Dans cet exercice, on se propose de donner deux preuves alternatives du résultat de cours : tout langage régulier est le langage d'une grammaire non contextuelle. Ces preuves reposent sur deux définitions qu'on peut donner d'un langage régulier, à savoir que c'est le langage d'un automate, ou que c'est le langage d'une expression régulière.

1. Avec des automates

Soit L un langage régulier et $\mathcal{A} = (\Sigma, Q, I, F, \delta)$ un automate[♡] le reconnaissant. On considère alors la famille d'automates $(\mathcal{A}_q)_{q \in Q} = (\Sigma, Q, \{q\}, F, \delta)$. On se munit par ailleurs de la famille de symboles non terminaux $(X_q)_{q \in Q}$, qu'on note être finie.

- Q. 1** Proposer un ensemble de règles de production P tel qu'en notant $\overset{*}{\Rightarrow}$ la relation de dérivation induite par P , on a $\forall q \in Q, \mathcal{L}(\mathcal{A}_q) = \{u \in \Sigma^*, X_q \overset{*}{\Rightarrow} u\}$. Démontrer cette propriété.
- Q. 2** Conclure.

2. Avec des expressions régulières

Soit L un langage régulier sur un alphabet fini Σ et e une expression régulière de langage L . On définit inductivement les **sous-expressions régulières** de e comme suit.

$$\text{ssexp}(e) = \{e\} \cup \begin{cases} \text{ssexp}(e_1) \cup \text{ssexp}(e_2) & \text{si } e = e_1|e_2 \text{ ou } e = e_1 \cdot e_2 \\ \text{ssexp}(e_1) & \text{si } e = (e_1)^* \\ \emptyset & \text{sinon, i.e. si } e = \varepsilon \text{ ou } e \in \Sigma \text{ ou } e = \emptyset \end{cases}$$

On se munit alors de la famille de symboles non terminaux $(X_r)_{r \in \text{ssexp}(e)}$ (qu'on remarque être finie).

- Q. 3** Proposer un ensemble de règles de production P tel qu'en notant $\overset{*}{\Rightarrow}$ la relation de dérivation induite par P , on a $\forall r \in \text{ssexp}(e), \mathcal{L}(r) = \{u \in \Sigma^*, X_r \overset{*}{\Rightarrow} u\}$. Démontrer cette propriété.
- Q. 4** Conclure.

[♣]. Lorsque w est un mot sur un alphabet et A un sous-ensemble de cet alphabet, $|w|_A$ désigne le nombre de lettres de w appartenant à A .

[♡]. sans ε -transitions

Exercice 12 : Mots de Łukasiewicz [CENTRALE 2008]

On considère l'ensemble de deux symboles $\Sigma = \{\square, \circ\}$, Σ contient donc uniquement le symbole "carré" : \square et le symbole "rond" : \circ . On s'intéresse dans cet exercice à un sous-ensemble de Σ^* nommé **mots de Łukasiewicz**. On rappelle qu'étant donné un mot $w = w_1 w_2 \dots w_n \in \Sigma^*$, la notation $|w|$ désigne la longueur de w , n ici.

On considère la fonction $va : \Sigma \rightarrow \{-1, 1\}$ telle que $va(\square) = -1$ et $va(\circ) = 1$. À un mot $w = w_1 w_2 \dots w_n \in \Sigma^*$ et à un entier $i \in \llbracket 0, n \rrbracket$ on associe la valeur : $\bar{w}^i = \sum_{k=1}^i va(w_k)$, en prenant comme convention que $\sum_{k=1}^0 va(w_k) = 0$.

Exemples : $\left\| \begin{array}{l} \overline{\circ\circ\square\square}^4 = va(\circ) + va(\circ) + va(\square) + va(\square) = 1 + 1 - 1 - 1 = 0 \\ \bar{\varepsilon}^0 = 0 \end{array} \right.$

On définit alors l'ensemble des **mots de Łukasiewicz** comme étant l'ensemble des mots $w \in \Sigma^*$ tels que $\forall i \in \llbracket 0, |w| - 1 \rrbracket, \bar{w}^i \geq 0$ et $\bar{w}^{|w|} = -1$

Par exemple, $w = \circ\circ\square\square \in \mathcal{L}$. En effet les valeurs des \bar{w}^i sont les suivantes.

i	0	1	2	3	4	5
\bar{w}^i	0	1	2	1	0	-1

Q. 1 Donner l'ensemble des éléments de \mathcal{L} de longueur 1, 2, 3.

Q. 2 Donner (en justifiant le résultat) l'ensemble des éléments de \mathcal{L} de longueur paire.

On suppose définis en OCaml les types et fonctions ci-contre. Le type `sym` est utilisé pour représenter l'ensemble Σ : le symbole \circ est représenté par le constructeur `R` (pour "rond") et \square est représenté par le constructeur `C` (pour "carré"). Le type `word` est utilisé pour représenter les mots de Σ^* . La fonction `va` implémente la fonction `va` définie ci-dessus.

```

1 | type sym = | R | C
2 | type word = sym list
3 |
4 | let va (s: sym) : int =
5 |   match s with
6 |   | R -> 1
7 |   | C -> -1

```

Q. 3 L'ensemble des mots de Łukasiewicz est-il régulier ?

Q. 4 Définir une fonction `est_luka : word -> bool` permettant de tester si un mot est de Łukasiewicz, *i.e.* permettant de tester, pour $w \in \Sigma^*$, si $w \in \mathcal{L}$. On proposera une implémentation ayant une complexité pire cas en $\mathcal{O}(n)$.

Q. 5 On appelle préfixe strict d'un mot $w \in \Sigma^*$, un mot $u \in \Sigma^*$ tel qu'il existe un mot $v \in \Sigma^* \setminus \{\varepsilon\}$ tel que $w = u \cdot v$. Montrer qu'un préfixe strict d'un mot de \mathcal{L} n'est pas un mot de \mathcal{L} .

Q. 6 Montrer que si $u \in \mathcal{L}$ et $v \in \mathcal{L}$ alors $\circ \cdot u \cdot v \in \mathcal{L}$.

Q. 7 Montrer que si $w \in \mathcal{L}$ et $|w| \neq 1$ alors il existe un unique couple $(u, v) \in \mathcal{L} \times \mathcal{L}$ tel que $w = \circ \cdot u \cdot v$.

Q. 8 Montrer qu'il existe une grammaire contextuelle non ambiguë reconnaissant l'ensemble des mots de Łukasiewicz.

Q. 9 Donner une fonction `decompose : word -> word * word` prenant en argument un mot de \mathcal{L} de taille différente de 1 (et de 0 car $\varepsilon \notin \mathcal{L}$) et calculant la décomposition de la Q. 7.

Exemples : $\left\| \begin{array}{l} \text{decompose } [R; R; C; C; C] = ([R; C; C], [C]) \\ \text{decompose } [R; R; C; C; R; C; C] = ([R; C; C], [R; C; C]) \end{array} \right.$

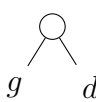
On définit maintenant l'ensemble des arbres binaires \mathbb{B} **sans étiquette** par les deux règles de construction :

```

1 | type btree =
2 |   | N of btree * btree      (* Nœud interne *)
3 |   | E                      (* Arbre vide *)

```

Autrement énoncé :

- $\square \in \mathbb{B}$
- Si $g \in \mathbb{B}$ et $d \in \mathbb{B}$ alors  $\in \mathbb{B}$

Q. 10 Donner une bijection naturelle de \mathbb{B} dans \mathcal{L} , induite par la **Q. 8**

Q. 11 Définir en OCAML deux fonctions `to_word : btree -> word` et `from_word: word -> btree` telles que $\forall w \in \mathcal{L}, (\text{to_word } (\text{from_word } w)) = w$ et $\forall b \in \mathbb{B}, (\text{from_word } (\text{to_word } b)) = b$.

Exemples : $\left\| \begin{array}{l} \text{from_word } [R; R; C; C; R; C; C] = N (N (E, E), N (E, E)) \\ \text{to_word } (N (N (E, E), N (E, E))) = [R; R; C; C; R; C; C] \end{array} \right.$

Q. 12 Démontrer que $\forall w \in \mathcal{L}, |w| \geq 3 \Rightarrow \exists (u, v) \in \Sigma^*, w = u \cdot \square \square \square \cdot v$.