

Partie II. Analyse descendante

Ce problème utilise le langage OCaml.

Dans ce problème, on étudie un algorithme d'analyse syntaxique, appelé analyse descendante, qui s'applique à certaines grammaires non contextuelles. Pour une grammaire donnée, on note N l'ensemble de ses non terminaux (notés avec des majuscules) et T l'ensemble de ses terminaux (notés avec des minuscules). Les lettres grecques (α, β, γ , etc.) désignent des mots de $(N \cup T)^*$. Le mot vide est noté ε . Une règle de production de la grammaire est notée $X \rightarrow \gamma$. Une dérivation immédiate est notée $\alpha \Rightarrow \beta$, ce qui signifie qu'il existe une règle de production $X \rightarrow \gamma$ avec $\alpha = \alpha_1 X \alpha_2$ et $\beta = \alpha_1 \gamma \alpha_2$. On note \Rightarrow^* la clôture réflexive transitive de la relation \Rightarrow , c'est-à-dire une dérivation en un nombre quelconque d'étapes, y compris zéro.

On prend en exemple une version simplifiée de la grammaire du langage de programmation LISP, avec un ensemble de quatre terminaux $T = \{\text{sym}, (,), \#\}$, un ensemble de trois non terminaux $N = \{S, L, E\}$, dont le symbole initial S , et les cinq règles de production suivantes :

$$\begin{array}{l} S \rightarrow L \# \\ L \rightarrow \varepsilon \\ \quad | \quad E L \\ E \rightarrow \text{sym} \\ \quad | \quad (L) \end{array}$$

Appelons G cette grammaire.

Question 11. Pour un entier $n \in \mathbb{N}$ arbitraire, donner un mot de longueur au moins n engendré par cette grammaire et la dérivation à gauche correspondante.

Correction : Il suffit de considérer le mot

$$\underbrace{\text{sym sym} \dots \text{sym}}_n \#$$

Sa dérivation à gauche est

$$\begin{array}{l} S \Rightarrow L \# \\ \Rightarrow E L \# \\ \Rightarrow \text{sym } L \# \\ \Rightarrow \text{sym } E L \# \\ \Rightarrow \text{sym sym } L \# \\ \vdots \\ \Rightarrow \text{sym} \dots \text{sym } L \# \\ \Rightarrow \text{sym} \dots \text{sym} \# \end{array}$$

Pour réaliser l'analyse syntaxique de ce petit langage avec OCaml, on se donne un type `token` pour les symboles terminaux `sym`, `(`, `)` et `#` :

```
type token = Sym | Lpar | Rpar | Eof
```

Ainsi `Sym` représente le terminal `sym`, `Lpar` le terminal `(`, `Rpar` le terminal `)` et `Eof` le terminal `#`.

Notre objectif est d'écrire une fonction `accepts: token list -> bool` qui détermine si un mot, donné comme une liste de terminaux, appartient ou non au langage de cette grammaire. Pour cela, on va commencer par construire trois fonctions, une pour chaque non terminal de la grammaire, avec les types suivants :

```
val parseS: token list -> token list
val parseL: token list -> token list
val parseE: token list -> token list
```

La fonction `parseX` reconnaît un préfixe maximal de la liste passée en argument comme étant un mot dérivé de `X` et renvoie le reste de la liste. S'il n'y a pas de tel préfixe, alors la fonction lève l'exception `SyntaxError`, que l'on suppose définie.

Pour définir ces trois fonctions, on se donne une table à deux entrées appelée *table LL*. Cette table indique, pour un non terminal que l'on cherche à reconnaître et pour un terminal au début de la liste, la règle de production à utiliser. Voici cette table pour notre grammaire :

	<code>sym</code>	<code>(</code>	<code>)</code>	<code>#</code>
<i>S</i>	<i>L#</i>	<i>L#</i>		<i>L#</i>
<i>L</i>	<i>EL</i>	<i>EL</i>	ε	ε
<i>E</i>	<code>sym</code>	<i>(L)</i>		

La fonction `parseS` est donc définie en suivant la première ligne de cette table et voici son code :

```
let rec parseS l = match l with
| (Sym | Lpar | Eof) :: _ ->
    (match parseL l with Eof :: q -> q | _ -> raise SyntaxError)
| [] | Rpar :: _ -> raise SyntaxError
```

En particulier, une case vide dans la table est interprétée comme un échec de l'analyse.

Question 12. Donner le code des fonctions `parseL` et `parseE`.

Correction :

```
and parseL l = match l with
| (Eof | Rpar) :: _ -> l
| (Sym | Lpar) :: _ -> parseL (parseE l)
| [] -> raise SyntaxError
and parseE l = match l with
| Sym :: q ->
    q
```

```

| Lpar :: q ->
  (match parseL q with Rpar :: r -> r | _ -> raise SyntaxError)
| [] | (Rpar | Eof) :: _ ->
  raise SyntaxError

```

Question 13. Donner le code de la fonction `accepts`.

Correction : Il faut rattraper l'exception `SyntaxError`, mais il faut également vérifier que tous les terminaux ont bien été consommés.

```

let accepts l =
  try parseS l = [] with SyntaxError -> false

```

Construction de la table. On va maintenant chercher à construire une telle table pour une grammaire arbitraire. On introduit pour cela une première notion : On dit qu'un non terminal X est *nul*, et on note $\text{NUL}(X)$, si le mot vide peut être dérivé de X , c'est-à-dire $X \Rightarrow^* \varepsilon$.

Question 14. Indiquer quels sont les symboles nuls de la grammaire G prise en exemple plus haut.

Correction : Le symbole L est trivialement nul, car on a la règle de production $L \rightarrow \varepsilon$.

Les symboles S et E ne sont pas nuls, car tout mot dérivé de S contient au moins le terminal `#`, et de même tout mot dérivé de E contient soit le terminal `sym`, soit le terminal `(`.

Calcul des symboles nuls. Pour déterminer $\text{NUL}(X)$, on propose l'algorithme suivant.

1. Initialement, on fixe $\text{NUL}(X)$ à `false` pour tout $X \in N$.
2. Pour chaque non terminal X , on affecte la valeur `true` à $\text{NUL}(X)$ s'il existe une production $X \rightarrow \varepsilon$ ou une production $X \rightarrow X_1 X_2 \dots X_p$ avec X_i des non terminaux et $\text{NUL}(X_i)$ pour tout i .
3. Si l'étape 2 a modifié au moins une valeur $\text{NUL}(X)$, alors on recommence l'étape 2.

Question 15. Montrer que cet algorithme termine.

Correction : La valeur de $\text{NUL}(X)$ n'évolue que dans le sens `false` vers `true`. Dès lors, le nombre de valeurs `false` ne fait que diminuer. S'il ne change pas, l'algorithme s'arrête. Sinon, il diminue strictement et constitue donc un variant de l'algorithme.

En particulier, le nombre d'étapes est borné par le nombre de non terminaux de la grammaire.

Question 16. Montrer que cet algorithme détermine bien la valeur de $\text{NUL}(X)$.

Correction : D'une part, on montre par récurrence sur le nombre d'étapes de l'algorithme que, si on obtient $\text{NUL}(X) = \text{true}$, alors effectivement $X \Rightarrow^* \varepsilon$. C'est clair initialement, car $\text{NUL}(X) = \text{false}$ pour tout X . Et si $\text{NUL}(X)$ devient true , alors soit $X \rightarrow \varepsilon$, soit $X \rightarrow X_1 X_2 \dots X_p$ avec $X_i \Rightarrow^* \varepsilon$ par hypothèse de récurrence.

D'autre part, on montre par récurrence sur le nombre d'étapes de la dérivation $X \Rightarrow^* \varepsilon$ qu'on obtient bien $\text{NUL}(X) = \text{true}$ par l'algorithme. Si $X \Rightarrow \varepsilon$, c'est qu'il existe une production $X \rightarrow \varepsilon$ et on aura bien $\text{NUL}(X) = \text{true}$ par l'algorithme. Sinon, $X \Rightarrow X_1 X_2 \dots X_p$ avec $X_i \Rightarrow^* \varepsilon$ et par hypothèse de récurrence, on aura $\text{NUL}(X_i) = \text{true}$ par l'algorithme, pour tout i , et donc $\text{NUL}(X) = \text{true}$.

Les premiers et les suivants. On introduit deux autres notions sur la grammaire. Pour un non terminal $X \in N$, on définit deux ensembles de terminaux :

- $\text{PREMIERS}(X)$ est l'ensemble des terminaux qui peuvent apparaître au début des mots dérivés depuis X , c'est-à-dire $\text{PREMIERS}(X) = \{t \in T \mid \exists \alpha \in (N \cup T)^*. X \Rightarrow^* t\alpha\}$;
- $\text{SUIVANTS}(X)$ est l'ensemble des terminaux qui peuvent apparaître après X dans une dérivation, c'est-à-dire $\text{SUIVANTS}(X) = \{t \in T \mid \exists \alpha, \beta \in (N \cup T)^*. S \Rightarrow^* \alpha X t \beta\}$.

Question 17. Donner les ensembles $\text{PREMIERS}(X)$ et $\text{SUIVANTS}(X)$ pour la grammaire prise en exemple plus haut (soit six ensembles au total).

Correction :

$$\begin{aligned}
 \text{PREMIERS}(S) &= \{\text{sym}, (, \#\} \\
 \text{PREMIERS}(L) &= \{\text{sym}, (\} \\
 \text{PREMIERS}(E) &= \{\text{sym}, (\} \\
 \text{SUIVANTS}(S) &= \{\} \\
 \text{SUIVANTS}(L) &= \{), \#\} \\
 \text{SUIVANTS}(E) &= \{\text{sym}, (,), \#\}
 \end{aligned}$$

Construction de la table. On admet que l'on peut calculer les ensembles $\text{PREMIERS}(X)$ et $\text{SUIVANTS}(X)$ pour toute grammaire. On étend NUL et PREMIERS sur tout mot de $(N \cup T)^*$ de

la manière suivante :

$$\begin{aligned}
\text{NUL}(\varepsilon) &= \text{true} \\
\text{NUL}(X_1 X_2 \dots X_p) &= \text{true} \quad \text{si } \text{NUL}(X_i) = \text{true} \text{ pour tout } i \\
\text{PREMIERS}(\varepsilon) &= \emptyset \\
\text{PREMIERS}(t\alpha) &= \{t\} \\
\text{PREMIERS}(X\alpha) &= \text{PREMIERS}(X) \quad \text{si } \text{NUL}(X) = \text{false} \\
\text{PREMIERS}(X\alpha) &= \text{PREMIERS}(X) \cup \text{PREMIERS}(\alpha) \quad \text{si } \text{NUL}(X) = \text{false}
\end{aligned}$$

où t est un terminal et X un non terminal.

On construit alors la table LL de la manière suivante : pour un non terminal $X \in N$ et un terminal $t \in T$, on indique la règle de production $X \rightarrow \gamma$ dans la case (X, t) de la table

- si $t \in \text{PREMIERS}(\gamma)$,
- ou si $\text{NUL}(\gamma)$ et $t \in \text{SUIVANTS}(X)$.

On peut alors se servir de cette table pour réaliser une analyse syntaxique dès lors que chaque case ne comporte qu'au plus une règle de production. *On pourra vérifier que l'on obtient bien la table donnée au début de cette partie avec les résultats obtenus aux questions 14 et 17.*

Un autre exemple. On considère une seconde grammaire G' , sur les mêmes symboles terminaux, avec des non terminaux $\{S', L', E'\}$ et un symbole initial S' :

$$\begin{array}{lcl}
S' & \rightarrow & L' \# \\
L' & \rightarrow & \varepsilon \\
& & | \quad L' E' \\
E' & \rightarrow & \text{sym} \\
& & | \quad (L')
\end{array}$$

Question 18. Montrer que G' reconnaît le même langage que G .

Correction : On comprend que G reconnaît les listes par la gauche et G' par la droite, le reste étant identique. On va montrer la double inclusion des langages.

On commence par un lemme : si $L \Rightarrow^* w \in T^* \setminus \{\varepsilon\}$ alors $w = w'u$ avec $L \Rightarrow^* w'$ et $E \Rightarrow^* u$. On le prouve par récurrence sur la longueur de la dérivation. La dérivation commence nécessairement par $L \Rightarrow EL$ car $w \neq \varepsilon$. Ensuite, elle se poursuit avec $E \Rightarrow^* w_1$ et $L \Rightarrow^* w_2$ et $w = w_1 w_2$. Par hypothèse de récurrence, $w_2 = w_3 u$ avec $L \Rightarrow^* w_3$ et $E \Rightarrow^* u$. Il vient donc $L \Rightarrow EL \Rightarrow^* w_1 w_3$ et on conclut en posant $w = w_1 w_3$.

Montrons maintenant que $X \Rightarrow^* w$ implique $X' \Rightarrow^* w$ pour chaque des trois non terminaux. On procède par récurrence sur la longueur de la dérivation $X \Rightarrow^* w$.

- pour $X = S$: la dérivation est $S \Rightarrow L\#$ avec $L \Rightarrow^* w'$ et $w = w'\#$. Par HR, $L' \Rightarrow^* w'$ et donc $S' \Rightarrow L'\# \Rightarrow^* w'\#$.

- pour $X = L$: Si $L \Rightarrow^* \varepsilon$ alors $L' \Rightarrow^* \varepsilon$. Sinon, on peut appliquer le lemme et $w = w'u$ avec $L \Rightarrow^* w'$ et $E \Rightarrow^* u$. Par HR, on a $L' \Rightarrow^* w'$ et $E' \Rightarrow^* u$, et donc $L' \Rightarrow L'E' \Rightarrow^* w'u = w$.
 - pour $X = E$, on a deux cas. Si $E \Rightarrow^* \text{sym}$ alors $E' \Rightarrow \text{sym}$ également. Si $E \Rightarrow (L) \Rightarrow^* (w')$ alors par HR $L' \Rightarrow^* w'$ et donc $E' \Rightarrow (L') \Rightarrow^* (w') = w$.
- L'autre inclusion se prouve de façon tout à fait similaire (avec un lemme similaire).
-

Question 19. Construire la table LL pour cette seconde grammaire. Permet-elle de coder un algorithme pour l'analyse syntaxique des mots générés par cette grammaire ?

Correction : Pour cette grammaire, on trouve

$$\begin{aligned}
 \text{PREMIERS}(S') &= \{\text{sym}, (, \#\} \\
 \text{PREMIERS}(L') &= \{\text{sym}, (\} \\
 \text{PREMIERS}(E') &= \{\text{sym}, (\} \\
 \text{SUIVANTS}(S') &= \{\} \\
 \text{SUIVANTS}(L') &= \{\text{sym}, (,), \#\} \\
 \text{SUIVANTS}(E') &= \{\text{sym}, (,), \#\}
 \end{aligned}$$

et on obtient du coup la table suivante :

	sym	()	#
S'	$L'\#$	$L'\#$		$L'\#$
L'	$\varepsilon/L'E'$	$\varepsilon/L'E'$	ε	ε
E'	sym	(L)		

Il y a deux cases dans la table où deux choix apparaissent, ce qui ne permet pas l'analyse.
