

Analyse Syntaxique

Chapitre extrait du livre : Compilation Cours et exercices

Auteur : H. DRIAS

Edition: OPU

Le rôle principal de l'analyse syntaxique est de vérifier si l'écriture du programme source est conforme avec la syntaxe du langage à compiler. Cette dernière est spécifiée à l'aide d'une grammaire.

Il existe plusieurs méthodes d'analyse appartenant à l'une des deux catégories d'analyse qui sont l'analyse descendante et l'analyse ascendante.

Dans l'analyse descendante nous essayons de dériver à partir de l'axiome de la grammaire, le programme source. il existe actuellement plusieurs méthodes qui parmi elles, nous distinguons les méthodes non déterministes ou avec retour arrière et les méthodes déterministes telles que l'analyse prédictive.

D'une façon opposée, l'analyse ascendante établit des réductions sur la chaîne à analyser pour aboutir à l'axiome de la grammaire. Comme pour l'analyse descendante, il existe des méthodes non déterministes peu efficaces et des méthodes déterministes. Dans toute méthode, un arbre syntaxique se construit soit explicitement soit implicitement.

Dans ce chapitre, nous effectuerons un tour d'horizon sur quelques méthodes d'analyse efficaces déterministes. Nous essayerons à la fin de fournir un aperçu sur le fonctionnement d'un générateur d'analyseur syntaxique nommé YACC.

3.1 RAPPELS.

Les notions de dérivation et ambiguïté seront constamment utilisés dans ce chapitre.

Soit la grammaire G donnée par:

$$\begin{array}{l} E \text{ ---> } E + E \\ \quad \quad \quad / \quad E * E \\ \quad \quad \quad / \quad a \end{array}$$

dérivation.

E dérive la chaîne $a * a + a$ et nous notons ceci par

$$E \xRightarrow{*} a * a + a$$

car il existe une suite de dérivations qui sont:

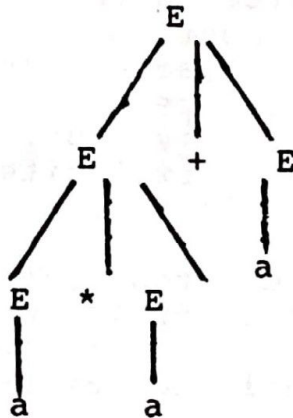
$$E \xRightarrow{*} E + E \xRightarrow{*} E * E + E \xRightarrow{*} a * E + E \xRightarrow{*} a * a + E$$

$$\xRightarrow{*} a * a + a$$

aboutissant à la chaîne $a * a + a$

ambiguïté.

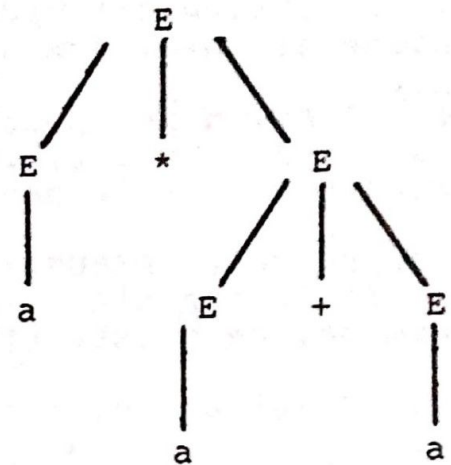
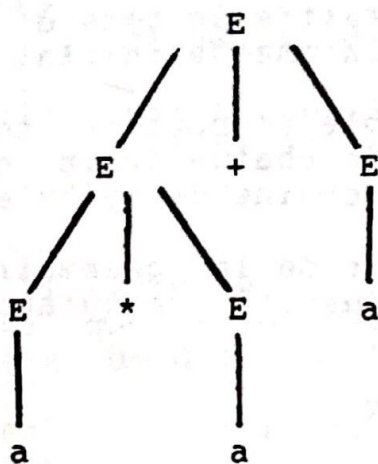
L'arbre



Schématise la dérivation précédente.

Une grammaire est dite ambiguë si pour une chaîne donnée engendrée par cette grammaire, il existe plusieurs arbres de dérivation distincts.

A la chaîne $a * a + a$ peuvent correspondre plusieurs arbres de dérivation parmi lesquels nous montrons les deux arbres suivants:



La grammaire G est donc ambiguë.

3.2 ANALYSE DESCENDANTE.

Les méthodes déterministes les plus efficaces connues à l'heure actuelle sont: l'analyse prédictive et la descente récursive. L'une et l'autre sont en quelque sorte équivalentes dans le sens où elles effectuent exactement les mêmes dérivations et par conséquent construisent le même arbre syntaxique pour une même chaîne acceptée par la grammaire. Lorsque le langage d'implémentation dispose de la récursivité, il est souhaitable et plus adéquat d'utiliser la descente récursive. Dans le cas contraire, l'analyse prédictive affirmera son existence.

3.2.1 Analyse Predictive.

Cette analyse utilise une table appelée table prédictive. Les lignes de cette table sont indicées par les symboles nonterminaux du langage, par contre, les colonnes de la table sont elles indicées par les symboles terminaux et le symbole '#'.

Les cases de la table contiennent éventuellement des règles de la grammaire. Pour la construction de cette table, deux notions qui sont DEBUT et SUIVANT sont nécessaires. Intuitivement le début d'une chaîne de symboles de la grammaire contient tous les terminaux (et éventuellement 'ε')

si la chaîne est vide) qui peuvent apparaître en tête d'une chaîne de terminaux dérivée à partir de la chaîne initiale.

Le SUIVANT d'un nonterminal est un ensemble comportant tous les terminaux (#, l'indicateur de fin de chaîne inclu) qui peuvent suivre ce nonterminal dans toute chaîne de symboles.

Pour calculer DEBUT(X) où X est un symbole de la grammaire, il faut exécuter les pas suivants jusqu'à ce qu'aucun terminal ne puisse être ajouté à DEBUT(X)

1-si X est un terminal alors DEBUT(X) = X

2-si X est un nonterminal et si $X \rightarrow a$ est une règle de la grammaire, a étant un terminal alors ajouter a à DEBUT(X).
Si $X \rightarrow \xi$ est une règle de la grammaire alors ajouter ξ à DEBUT(X).

3-si $X \rightarrow y_1 y_2 \dots y_m$ est une règle de grammaire et si pour $j=1, \dots, i-1$ DEBUT(y_j) contient ϵ alors ajouter tous les éléments différents de ϵ qui appartiennent à DEBUT(y_i) à DEBUT(X).
Si pour $j=1, \dots, m$ DEBUT(y_j) contient ϵ alors ajouter ϵ à DEBUT(X).

Pour déterminer SUIVANT(X) où X est un nonterminal, suivre les pas suivants jusqu'à ce qu'aucun élément ne puisse être ajouté à SUIVANT(X).

1-si X est l'axiome de la grammaire alors ajouter # (marqueur de fin de chaîne) à SUIVANT(X).

2-si $A \rightarrow XaB$ appartient à la grammaire et si a est un terminal alors ajouter a à SUIVANT(X)

3-si $A \rightarrow \alpha X$ est une règle de la grammaire alors ajouter tous les suivants de A à SUIVANT(X).

3.2.1.1 Algorithme De Construction De La Table Predictive.

Pour chaque règle $A \rightarrow \alpha$ de la grammaire où $\alpha \in \Sigma^+$ faire
pour chaque élément terminal a appartenant à DEBUT(α)

faire

placer dans la case (A, a) la règle $A \rightarrow a$.

Pour chaque règle de la grammaire de la forme $A \rightarrow \xi$ faire pour chaque terminal x dans $SUIVANT(A)$ faire placer dans la case (A, x) la règle $A \rightarrow \xi$.

Si la table, une fois construite est multidéfinie c'est à dire que si une case de la table est occupée par plus d'une règle, on dit que la grammaire n'est pas LL(1).

Dans le cas contraire c'est à dire lorsque dans toutes les cases de la table, il existe au plus une règle, on dit que la grammaire est LL(1).

exemple

Soit G la grammaire suivante:

$$E \rightarrow E + T$$

$$/ T$$

$$T \rightarrow T * F$$

$$/ F$$

$$F \rightarrow (E)$$

$$/ id$$

Figure 3.1: Grammaire des expressions arithmétiques.

Le DEBUT et le SUIVANT de chaque nonterminal sont montrés dans la table ci-dessous.

	DEBUT	SUIVANT
E	(id	# +)
T	(id	# +) *
F	(id	# +) *

La table prédictive de G est donc la suivante:

	+	*	()	id	#
E			E-->E+T E-->T		E-->T	
T			T-->T*F T-->F		T-->F	
F			F-->(E)		F-->id	

Cette table est multidéfinie. Par conséquent G n'est pas LL(1).

3.2.1.2 Définition Formelle D'une Grammaire LL(1). -

Une grammaire G est dite LL(1) si ses règles de production vérifient les conditions suivantes:

1-si $A \xrightarrow{*} \alpha$ et $A \xrightarrow{*} \beta$ sont deux règles de G alors on doit avoir ou bien

$\alpha \xrightarrow{*} \epsilon$ ou bien $\beta \xrightarrow{*} \epsilon$ mais pas les deux.

2-si $A \xrightarrow{*} \alpha$ et $A \xrightarrow{*} \beta$ sont deux règles de G alors
 $DEBUT(\alpha) \cap DEBUT(\beta) = \emptyset$

3-si $A \rightarrow \alpha$ et $A \rightarrow \beta$ sont deux règles de G et si $\beta \xrightarrow{*} \epsilon$
alors $\text{SUIVANT}(A) \cap \text{DEBUT}(\alpha) = \emptyset$.

Ces conditions entraînent l'unicité de l'existence d'une règle dans une case de la table prédictive.

3.2.1.3 Conditions Nécessaires Pour Qu'une Grammaire Soit LL(1). -

Pour qu'une grammaire soit LL(1), il faut que toutes ses règles soient non récursives à gauche et factorisées.

élimination de la récursivité à gauche.

La récursivité à gauche apparaît dans les règles de la forme
 $A \rightarrow A\alpha$ Pour l'éliminer, il faut substituer aux règles

$$\begin{aligned} A &\rightarrow A\alpha \\ A &\rightarrow \beta \end{aligned}$$

les règles:

$$\begin{aligned} A &\rightarrow \beta A' \\ A' &\rightarrow \alpha A' / \epsilon \end{aligned}$$

exemple:

La grammaire G de la figure 3.1 après élimination de la récursivité à gauche devient G' qui est:

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' / \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' / \\ F &\rightarrow (E) / id \end{aligned}$$

Figure 3.2: Grammaire G'.

Pour satisfaire la condition 2 d'une grammaire LL(1), nous devons factoriser toutes les règles de la grammaire. Cette opération revient à substituer aux règles de la forme :

$$\begin{aligned} A &\rightarrow \alpha \beta \\ A &\rightarrow \alpha \gamma \end{aligned}$$

les règles:

$$\begin{aligned} A &\rightarrow \alpha A' \\ A' &\rightarrow \beta / \gamma \end{aligned}$$

Et de façon générale, il faut remplacer les règles de la forme:

$$A \rightarrow \alpha\beta_1/\alpha\beta_2/\dots/\alpha\beta_n$$

par les règles:

$$\begin{aligned} A &\rightarrow \alpha A' \\ A' &\rightarrow \beta_1/\beta_2/\dots/\beta_n. \end{aligned}$$

exemple:

La grammaire G' de la figure 3.2 est déjà factorisée, de même est la grammaire de la figure 3.1.

3.2.1.4 Algorithme D'analyse. -

Nous utilisons une pile comportant les symboles de la grammaire et un vecteur contenant la chaîne à analyser terminée par le symbole #.

Initialement la configuration de l'analyse est:

(#s ----- contenu de la pile	, a1 a2...am#) ----- vecteur contenant la chaîne à analyser.
---------------------------------------	-----------------------------------------------------------------------

S étant l'axiome de la grammaire.

Avant d'analyser une chaîne, il faut convertir la grammaire décrivant le langage qui accepte la chaîne à analyser en une grammaire LL(1) et ensuite construire la table prédictive de cette dernière.

algorithme:

1-soient X l'élément au sommet de la pile et a_i le symbole courant de la chaîne à analyser.

a- si X est un terminal alors

-si X = ai alors dépiler X, avancer dans la chaîne et aller à 1.

-sinon si X = '#' et ai = '#' alors la chaîne est acceptée. stop.

-sinon si X <> ai alors la chaîne est erronée.

b- si X est un nonterminal alors soit (X--->α) la règle se trouvant dans la case (X, ai)

-dépiler X, empiler la partie droite de la règle et aller à 1.

exemple :

La table prédictive de G' est:

	+	*	()	id	#
E			E-> TE'		E->TE'	
T			T->FT'		T->FT'	
F			F->(E)		F->id	
E'	E->+TE'			E'-->ε		E'-->ε
T'	T'-->ε	T'-->*FT'		T'-->ε		T'-->ε

analysons la chaine id1 + id2*id3:

```

(#E          ,          id1+id2*id3#)
(#E'T       ,          id1+id2*id3#)
(#E'T'F     ,          id1+id2*id3#)
(#E'T'id    ,          id1+id2*id3#)
(#E'T'      ,          +id2*id3#)
(#E'        ,          + id2*id3#)
(#E'T+      ,          +id2 * id3# )
(#E'T       ,          id2*id3#)
(#E'T'F     ,          id2*id3#)
(#E'T'id    ,          id2*id3#)
(#E'T'      ,          * id3#)
(#E'T'F*    ,          *id3#)
(#E'T'F     ,          id3#)
(#E'T'id    ,          id#)
(#E'T'      ,          #)
(#E'        ,          #)
(#          ,          #)
o.k
    
```

L'arbre syntaxique produit par l'analyse de la chaîne $id_1+id_2*id_3$ est montré sur la figure 3.3.

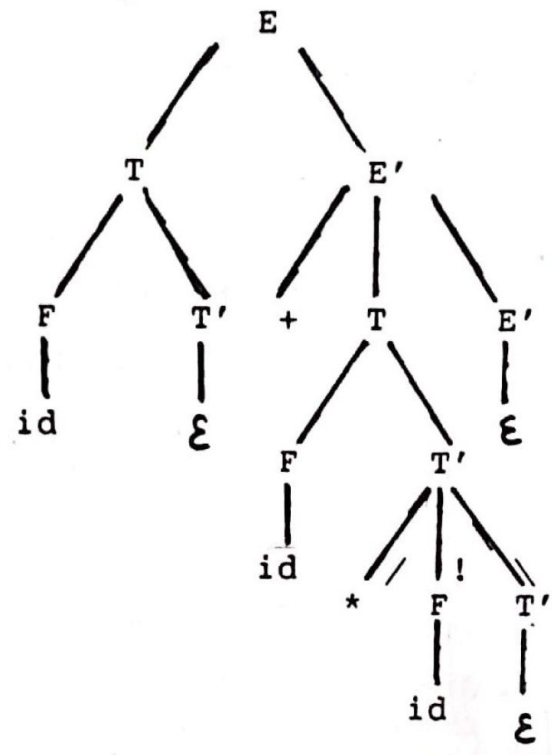


Figure 3.3.