

## Analyse de Données Structurées - Cours 6

Ralf Treinen



Université Paris Diderot  
UFR Informatique  
Laboratoire Preuves, Programmes et Systèmes

treinen@pps.univ-paris-diderot.fr

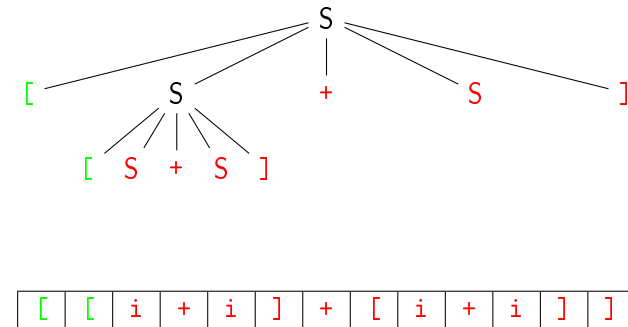
11 mars 2015

© Ralf Treinen 2015

## Grammaires LL(1)

- Intuition derrière les grammaires LL(1) : dans la construction d'une dérivation gauche, le symbole suivant de l'entrée nous indique quelle règle de la grammaire appliquer au non-terminal le plus à gauche de l'arbre de dérivation.
- Conséquence : toute grammaire LL(1) (ou même LL(k)) est non-ambiguë (oublié de préciser au dernier cours).
- Un critère très simple : Si tous les côtés droites de la grammaire pour le *même* non-terminal commencent sur des terminaux différents, alors la grammaire est LL(1).

## Construction d'un arbre de dérivation



Choisir règle (1) : c'est la seule qui peut produire à partir de S un mot qui commence sur i.

## Une caractérisation des grammaires LL(k)

## Théorème

Une grammaire  $G = (V_T, V_N, S, R)$  est LL(k) ssi

- si  $A \rightarrow \beta$  et  $A \rightarrow \gamma$  sont deux règles différentes
- et  $S \rightarrow^* wA\alpha$  une dérivation gauche,  $w \in V_T^*$
- alors  $\text{FIRST}_k(\beta\alpha) \cap \text{FIRST}_k(\gamma\alpha) = \emptyset$
- Preuve omise (conséquence immédiate de la définition)
- Attention, ce critère général prend en compte le "contexte" dans lequel on peut obtenir A.

## Un meilleur critère pour être LL(1)

- Définition de  $FIRST_1(\alpha)$  : l'ensemble des symboles avec lesquelles un mot terminal dérivé à partir de  $\alpha$  peut commencer (plus  $\epsilon$  dans le cas où  $\alpha \rightarrow^* \epsilon$ ).
- Vu au dernier cours : calcul de  $FIRST_1$  dans le cas où aucun côté droite est  $\epsilon$ .
- Meilleur critère : Si tous les côtés droites de la grammaire pour le *même* non-terminal ont des ensembles  $FIRST_1$  disjointes, alors la grammaire est LL(1).

## Exemple vu au dernier cours

- Grammaire  $G = (\{a, (, ), +\}, \{F, S\}, S, R)$  où  $R$  est

$$F \rightarrow a$$

$$S \rightarrow (F+S)$$

$$S \rightarrow F$$

- Le première critère simple ne s'applique pas.
- On obtient pour les côtés droites des règles :

$$FIRST_1(a) = \{a\}$$

$$FIRST_1((F+S)) = \{($$

$$FIRST_1(F) = \{a\}$$

## Reconnaître la fin de l'entrée

- Le programme vu au dernier cours a un défaut : il accepte aussi des expressions correctes, avec un texte quelconque ajouté à la fin :  $(a+(a+a))\%\$ \# @$
- Solution :
  - l'analyse lexicale envoie un jeton qui signale la fin de l'entrée (par exemple, EOF)
  - remplacer l'axiome par  $S'$ , avec une règle

$$S' \rightarrow S \text{ EOF}$$

où  $S$  est l'ancien axiome de la grammaire.

## Le même exemple avec reconnaissance de la fin

- Grammaire  $G = (\{a, (, ), +, \text{EOF}\}, \{F, S, S'\}, S', R)$  où  $R$  est

$$F \rightarrow a$$

$$S \rightarrow (F+S)$$

$$S \rightarrow F$$

$$S' \rightarrow S \text{ EOF}$$

- On obtient pour les côtés droites des règles :

$$FIRST_1(a) = \{a\}$$

$$FIRST_1((F+S)) = \{($$

$$FIRST_1(F) = \{a\}$$

$$FIRST_1(S \text{ EOF}) = \{a, ($$

## Un exemple un peu plus avancé

- ▶ Une grammaire pour les expressions arithmétiques partiellement parenthésées.
- ▶ Terminaux :  $\{i, +, *, (, ), EOF\}$
- ▶ Règles :

$$\begin{aligned} S &\rightarrow E EOF \\ E &\rightarrow E+T \mid T \\ T &\rightarrow T*F \mid F \\ F &\rightarrow (E) \mid i \end{aligned}$$

- ▶ Axiome : S

## Un critère pour ne pas *pas* être LL(k)

### Définition

Une grammaire  $G = (V_T, V_N, S, R)$  est *récursive à gauche* s'il y a un non-terminal  $N \in V_N$  tel que  $N \rightarrow^+ N\alpha$  pour un  $\alpha \in (V_T \cup V_N)^*$ .

- ▶  $\rightarrow^+$  : dérivation en au moins une étape.
- ▶ Exemple : notre grammaire pour les expressions partiellement parenthésées, car  $E \rightarrow E+T$

### Lemme

Si la grammaire  $G$  est récursive à gauche, alors  $G$  n'est pas *pas* LL(k), pour aucun  $k \in \mathbb{N}$ .

## L'exemple des expressions partiellement parenthésées

- ▶ Cette grammaire est non-ambiguë ☺
- ▶ Intuition : Les "+" peuvent être produites seulement à partir du E. Tout mot engendré par E est protégé par des parenthèses ( et ).
- ▶ (Il y a aussi une preuve formelle mais je vous en fait grâce.)
- ▶ Cette grammaire, est-elle aussi LL(1) ? Ou au moins LL(k) pour quelque  $k \in \mathbb{N}$  ?
- ▶ Elle n'est pas LL(k), pour aucun  $k$  ! ☹

## Preuve

- ▶ Supposons pour l'absurde que  $G$  est LL(k) et récursive à gauche.
- ▶ Il y a donc une dérivation gauche  $S \rightarrow^* wX\gamma$  (sinon  $X$  n'est pas atteignable).
- ▶ Il y a aussi une règle  $X \rightarrow X\alpha$  (hypothèse simplificatrice), et une règle différente  $X \rightarrow \beta$  (sinon  $X$  est non-productif).
- ▶ Il existe donc une dérivation gauche

$$S \rightarrow^* wX\gamma \rightarrow^* wX\alpha^k\gamma$$

- ▶ Par le théorème du début du cours :

$$\text{FIRST}_k(X\alpha^{k+1}\gamma) \cap \text{FIRST}_k(\beta\alpha^k\gamma) = \emptyset$$

## Preuve (2)

- ▶ On a :

$$\text{FIRST}_k(X\alpha^{k+1}\gamma) \cap \text{FIRST}_k(\beta\alpha^k\gamma) = \emptyset$$

- ▶ Donc, grâce à la règle  $X \rightarrow \beta$  :

$$\text{FIRST}_k(\beta\alpha^{k+1}\gamma) \cap \text{FIRST}_k(\beta\alpha^k\gamma) = \emptyset$$

- ▶ Contradiction (deux cas :  $\alpha \rightarrow^* \epsilon$  ou pas).

## La grammaire transformée

- ▶ Une grammaire pour les expressions arithmétiques partiellement parenthésées.
- ▶ Terminaux :  $\{i, +, *, (, ), \text{EOF}\}$
- ▶ Règles :

$$\begin{aligned} S &\rightarrow E \text{ EOF} \\ E &\rightarrow T E' \\ E' &\rightarrow \epsilon \mid +E \\ T &\rightarrow F T' \\ T' &\rightarrow \epsilon \mid *T \\ F &\rightarrow (E) \mid i \end{aligned}$$

- ▶ Axiome : S

## Quoi faire ?

- ▶ On peut transformer la grammaire en une grammaire équivalente (qui définit le même langage), et qui est LL(1).
- ▶ C'est toujours possible, mais il y a deux inconvénients :
  - ▶ la grammaire résultante peut être plus grande ;
  - ▶ la structure de l'arbre de dérivation peut changer.
- ▶ Il y a un troisième inconvénient : la transformation peut introduire des règles  $N \rightarrow \epsilon$ , il faut donc adapter la technique à ce cas.

## Explication de la transformation

- ▶ Les deux règles originales pour le non-terminal E :

$$\begin{aligned} E &\rightarrow T \\ E &\rightarrow E+T \end{aligned}$$

- ▶ Dans la grammaire d'origine, le non-terminal E engendre une séquence non-vide de non-terminaux T, séparés par des +.
- ▶ Dans la nouvelle grammaire, le non-terminal E' engendre la suite de cette séquence après un T :

$$\begin{aligned} E &\rightarrow T E' \\ E' &\rightarrow \epsilon \mid +E \end{aligned}$$

- ▶ Pareil pour le non-terminal T.

## But d'une analyse d'une grammaire

- ▶ Détecter des propriétés des non-terminaux dans une grammaire, et des grammaires.
- ▶ Exemple : non-terminaux atteignables ou pas, productifs ou pas, pouvant engendrer le mot vide ou pas ? Grammaires  $LL(1)$ ,  $LL(k)$  ou pas ?
- ▶ Défis : récurrence entre des les non-terminaux d'une grammaire.
- ▶ Dans certains cas on maîtrise une descente récursive *sur un mot donné* pour la construction d'une dérivation - cas des grammaires  $LL(k)$ .
- ▶ Dans une analyse de la grammaire, on a pas de mot d'entrée fixe donné mais on s'intéresse à une propriété générale de la grammaire.

## Calcul d'un point fixe

- ▶ Notre définition de productivité est correcte, mais il faut organiser le calcul différemment :
  - ▶ Au début, on pose que tout non-terminal  $N$  pour lequel existe une règle  $N \rightarrow w$ , où  $w \in V_T^*$ , est productif.
  - ▶ Puis, s'il y a une règle  $N \rightarrow \alpha$ , où on a déjà reconnu tous les non-terminaux dans  $\alpha$  comme étant productifs, alors on pose que  $N$  est aussi productif.
  - ▶ On s'arrête si on ne peut plus ajouter d'information de cette façon.
- ▶ Il s'agit d'un point fixe !

## Illustration du problème d'une descente récursive

- ▶ Un non-terminal  $N$  est *productif* s'il y a une règle  $N \rightarrow \alpha$  telle que tous les non-terminaux dans  $\alpha$  sont productifs. (C'est en particulier vrai quand  $\alpha$  ne contient que des terminaux.)
- ▶ Imaginez les règles suivantes de la grammaire :

$$\begin{aligned} A &\rightarrow BC \mid a \\ B &\rightarrow Cb \\ C &\rightarrow Ac \end{aligned}$$

- ▶ Il faut éviter une descente récursive qui mène dans une boucle :  
 $A \text{ productif ?} \rightarrow B \text{ productif ?} \rightarrow C \text{ productif ?}$   
 $\rightarrow A \text{ productif ?} \rightarrow \dots$

## Points fixes

- ▶ En maths : un *point fixe* d'une fonction  $f: D \rightarrow D$  est une valeur  $x \in D$  telle que  $f(x) = x$ .
- ▶ Application à l'algorithmique :
  - ▶ le domaine  $D$  est l'ensemble de toutes les affectations possibles à des variables d'intérêt.  
Dans l'exemple : toutes les affectations possibles de  $P$  du type  $V_N \rightarrow \text{boolean}$ .
  - ▶ la fonction  $f$  est une mise à jour des variables.  
Dans l'exemple : propagation de l'information de productivité d'un non-terminal à un autre.
  - ▶ on ne cherche pas un point fixe quelconque, mais on commence avec une valeur initiale, puis on applique  $f$  jusqu'à un point fixe.

## Calcul d'un point fixe

Nous utiliserons dans nos algorithmes la construction “do ... until  $X_1, \dots, X_n$  fix”. Exemple productivité :

```
for all  $N \in V_N$  :  
  if exists  $(N \rightarrow \alpha) \in R$  tel que  $\alpha \in V_T^*$   
  then  $P(N) = \text{true}$   
  else  $P(N) = \text{false}$   
do  
  for all  $(N \rightarrow \alpha) \in R$  :  
    if  $P(M)$  for all non-terminals  $M$  in  $\alpha$   
    then  $P(N) = \text{true}$   
until  $P$  fix
```

## Autre exemple symboles atteignables

- ▶ Un non-terminal  $N$  est *atteignable* (à partir de l'axiome  $S$ ) s'il existe une dérivation  $S \rightarrow^* \alpha N \beta$ .
- ▶ Calcul des non-terminaux atteignables : transparent suivant.
- ▶ Une grammaire est *réduite* si
  - ▶ tous ses non-terminaux sont productifs et
  - ▶ tous ses non-terminaux sont atteignables.
- ▶ Dans la suite nous supposons que toutes les grammaires sont réduites.

## La construction do ... until ... fix

### Le code

```
// code initialisation de X  
do  
  // code mise a jour de X  
until X fix
```

### peut être traduit vers :

```
// code initialisation de X  
do  
   $X_{old} = X$   
  // code mise a jour de X  
while ( $X \neq X_{old}$ )
```

## Calcul des non-terminaux atteignables

### Algorithme

Donnée une grammaire  $G = (V_T, V_N, S, R)$ .

```
for all  $N \in V_N$  :  
   $A(N) = \text{false}$   
 $A(S) = \text{true}$   
do  
  for all  $(N \rightarrow \alpha) \in R$  :  
    if  $A(N)$   
    then for all non-terminals  $M$  in  $\alpha$  :  
       $A(M) = \text{true}$   
until  $A$  fix
```

### Lemma

Pour tout  $N \in V_N$  :  $A(N) == \text{true}$  ssi  $N$  est atteignable.

## Non-terminaux qui peuvent produire $\epsilon$

### Algorithme

Donnée une grammaire  $G = (V_T, V_N, S, R)$ .

```
EPS = {N ∈ V_N | (N → ε) ∈ R}
do
  EPS = EPS ∪ {N ∈ V_N | (N → N_1 ... N_n) ∈ R,
                        N_1, ..., N_n ∈ EPS}
until EPS fix
```

### Lemme

Pour tout  $N \in V_N$  :  $N \in \text{EPS}$  ssi  $N \rightarrow^* \epsilon$ .

## Comment calculer $FIRST_1$ dans le cas général?

- Imaginez une règle  $A \rightarrow B C d E$
- Tous les non-terminaux peuvent a priori engendrer  $\epsilon$ .
- Si  $B \notin \text{EPS}$  : dans cette règle, seulement  $FIRST_1(B)$  peut contribuer à  $FIRST_1(A)$ .
- Si  $B \in \text{EPS}$  :  $FIRST_1(C)$  peut aussi contribuer à  $FIRST_1(A)$ .
- Si  $B \in \text{EPS}$  et  $C \in \text{EPS}$  :  $d$  doit être dans  $FIRST_1(A)$ .
- Dans aucun des cas,  $FIRST_1(E)$  ne peut contribuer car il se trouve derrière le terminal  $d$ .

## Calcul de EPS sur l'exemple

$E \rightarrow$	$T E'$	$E' \rightarrow$	$\epsilon \mid +E$
$T \rightarrow$	$F T'$	$T' \rightarrow$	$\epsilon \mid *T$
$F \rightarrow$	$(E) \mid i$	$S \rightarrow$	$E E O F$

- Initialisation :  $\text{EPS} = \{E', T'\}$
- C'est déjà un point fixe!

## Le calcul de $FIRST_1$ dans le cas général

- Pour des raisons techniques, il est plus facile de calculer d'abord une variante de  $FIRST_1$  sans  $\epsilon$  :

$$Fi(N) = \{a \in V_T \mid N \rightarrow^* aw, w \in V_T^*\}$$

- Puis on en obtient  $FIRST_1(N)$  en utilisant  $\text{EPS}$  que nous avons déjà calculé.

## Calcul de $FIRST_1$ dans le cas général

### Algorithme

Donnée une grammaire  $G = (V_T, V_N, S, R)$ .

```

pour tout  $N \in V_N$ :
   $Fi(N) = \{a \in V_T \mid (N \rightarrow N_1 \dots N_n a \alpha) \in R, \\ N_1, \dots, N_n \in EPS\}$ 
do
  pour tout  $(N \rightarrow N_1 \dots N_n M \alpha) \in R$ 
  tel que  $M \in V_N, N_1, \dots, N_n \in EPS$ :
     $Fi(N) = Fi(N) \cup Fi(M)$ 
until F fix
  
```

### Lemme

Pour tout  $N \in V_N$  :  $Fi(N) = FIRST_1(N) - \{\epsilon\}$

## Calcul de $FIRST_1$ dans le cas général

Pour tout  $N \in V_N$  :

$$FIRST_i(N) = \begin{cases} Fi(N) \cup \{\epsilon\} & \text{si } N \in EPS \\ Fi(N) & \text{sinon} \end{cases}$$

## Calcul de Fi sur l'exemple

$$\begin{aligned}
 E &\rightarrow T E' & E' &\rightarrow \epsilon \mid +E \\
 T &\rightarrow F T' & T' &\rightarrow \epsilon \mid *T \\
 F &\rightarrow (E) \mid i & S &\rightarrow E EOF \\
 EPS &= \{E', T'\}
 \end{aligned}$$

### Calcul

	<i>Initial</i>	<i>Iter1</i>	<i>Iter2</i>	<i>Iter3</i>
S	$\emptyset$	$\emptyset$	$\emptyset$	$\{i, ()\}$
E	$\emptyset$	$\emptyset$	$\{i, ()\}$	$\{i, ()\}$
E'	$\{+\}$	$\{+\}$	$\{+\}$	$\{+\}$
T	$\emptyset$	$\{i, ()\}$	$\{i, ()\}$	$\{i, ()\}$
T'	$\{*\}$	$\{*\}$	$\{*\}$	$\{*\}$
F	$\{i, ()\}$	$\{i, ()\}$	$\{i, ()\}$	$\{i, ()\}$

## Calcul de $FIRST_1$ sur l'exemple

$$\begin{aligned}
 E &\rightarrow T E' & E' &\rightarrow \epsilon \mid +E \\
 T &\rightarrow F T' & T' &\rightarrow \epsilon \mid *T \\
 F &\rightarrow (E) \mid i & S &\rightarrow E EOF \\
 EPS &= \{E', T'\}
 \end{aligned}$$

### Calcul

	<i>Fi</i>	<i>FIRST<sub>1</sub></i>
S	$\{i, ()\}$	$\{i, ()\}$
E	$\{i, ()\}$	$\{i, ()\}$
E'	$\{+\}$	$\{+, \epsilon\}$
T	$\{i, ()\}$	$\{i, ()\}$
T'	$\{*\}$	$\{*, \epsilon\}$
F	$\{i, ()\}$	$\{i, ()\}$



## Calcul de $FIRST_1$ dans le cas général

On étend maintenant  $FIRST_1$  à des mots de terminaux et non-terminaux :

- ▶  $FIRST_1(\epsilon) = \{\epsilon\}$
- ▶ pour tout  $a \in V_T$  :  $FIRST_1(a\alpha) = \{a\}$
- ▶ pour tout  $N \in V_N$  :

$$FIRST_1(N\alpha) = \begin{cases} FIRST_1(N) & \text{si } N \notin EPS \\ (FIRST_1(N) \setminus \{\epsilon\}) \cup FIRST_1(\alpha) & \text{si } N \in EPS \end{cases}$$

## La fonction $FOLLOW_k$

### Définition

Soit  $G = (V_T, V_N, S, R)$  une grammaire,  $k \in \mathbb{N}$ . La fonction  $FOLLOW_k : V_N \rightarrow 2^{V_T^*}$  est définie par

$$FOLLOW_k(N) = \{w \mid S \rightarrow^* \beta N \gamma, \beta, \gamma \in (V_T \cup V_N)^*, w \in FIRST_k(\gamma)\}$$

### Remarques

- ▶  $FOLLOW_k(N)$  est l'ensemble de tous les mots de terminaux de longueur  $k$  qui peuvent, dans des mots de  $\mathcal{L}(G)$ , suivre à un mot dérivé de  $N$ .

## Nous avons besoin de plus d'information !

- ▶ Le calcul de  $FIRST_1$  n'est plus suffisant pour savoir quelle production appliquer !
- ▶ Exemple :  $E' \rightarrow \epsilon \mid +E$   
Si nous voyons  $+$  alors il faut utiliser la deuxième alternative pour réécrire  $E'$ . Mais quand faut-il appliquer la première ?
- ▶ Il nous manque une information : quel sont les symboles terminaux qui peuvent *suivre* à un mot produit par un non-terminal ?

## Calcul de $FOLLOW_1$

### Algorithme

Donnée une grammaire  $G = (V_T, V_N, S, R)$ .

```
for all  $N \in V_N$ :  
  Fo(N) = { $a \in V_T \mid (M \rightarrow \dots NN_1 \dots N_k a \dots) \in R, N_1, \dots, N_n \in EPS$ }  
  for all  $(M \rightarrow \dots NN_1 \dots N_k N' \dots) \in R$   
  with  $N_1, \dots, N_k \in EPS$  :  
    Fo(N) = Fo(N)  $\cup$  Fi(N')  
  do  
    for all  $(M \rightarrow \dots NN_1 \dots N_k) \in R$   
    with  $N_1, \dots, N_n \in EPS$  :  
      Fo(N) = Fo(N)  $\cup$  Fo(M)  
until Fo fix
```

### Lemme

Si  $(S \rightarrow N \text{ EOF}) \in R$ , alors  $\forall N \in V_N : Fo(N) = FOLLOW_1(N)$

## Calcul de $FOLLOW_1$ sur l'exemple

$$\begin{aligned} E &\rightarrow T E' & E' &\rightarrow \epsilon \mid +E \\ T &\rightarrow F T' & T' &\rightarrow \epsilon \mid *T \\ F &\rightarrow (E) \mid i & S &\rightarrow E EOF \\ EPS &= \{E', T'\} & Fi(E') &= \{+\} & Fi(T') &= \{*\} \end{aligned}$$

### Calcul

	Initial	Iter1	Iter2
S	$\emptyset$	$\emptyset$	$\emptyset$
E	$\{EOF, )\}$	$\{EOF, )\}$	$\{EOF, )\}$
E'	$\emptyset$	$\{EOF, )\}$	$\{EOF, )\}$
T	$\{+\}$	$\{+, EOF, )\}$	$\{+, EOF, )\}$
T'	$\emptyset$	$\{+\}$	$\{EOF, ), +\}$
F	$\{*\}$	$\{*, EOF, ), +\}$	$\{*, EOF, ), +\}$

## Comment choisir la règle dans l'analyse lexicale

Soit  $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$  une alternative. Il y a deux cas :

1. Soit, aucun des  $FIRST_1(\alpha_i)$  ne contient  $\epsilon$  :  
comme avant :
  - ▶ On choisit la règle  $A \rightarrow \alpha_i$  quand le symbole suivant est dans  $FIRST_1(\alpha_i)$  (ils sont tous disjoints).
  - ▶ Erreur si aucun tel  $i$  existe
2. Soit  $\alpha_i$  avec  $\epsilon \in FIRST_1(\alpha_i)$  :
  - ▶ si le symbole suivant est dans  $FIRST_1(\alpha_j)$  : choisir  $A \rightarrow \alpha_j$ , pour  $1 \leq j \leq n$ .
  - ▶ si le symbole suivant est dans  $FOLLOW_1(A)$  : choisir  $A \rightarrow \alpha_i$ .
  - ▶ sinon Erreur.

## Finalement (!) : le critère pour être LL(1)

### Théorème

La grammaire  $G = (V_T, V_N, S, R)$  est LL(1) ssi pour toutes les alternatives  $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$  :

- ▶  $FIRST_1(\alpha_1), \dots, FIRST_1(\alpha_n)$  sont disjoints entr'eux ;
- ▶ Si  $\epsilon \in FIRST_1(\alpha_i)$ , alors pour tous  $j \neq i$  :

$$FIRST_1(\alpha_j) \cap FOLLOW_1(A) = \emptyset$$

### Remarque

Condition (1) implique qu'au plus un des ensembles  $FIRST_1(\alpha_i)$  contient  $\epsilon$ .