

Traduction et Sémantique

Analyse lexicale

Luigi Santocanale
LIF, Université de Provence
Marseille, FRANCE

12 février 2010

Plan

Rappels de la théorie des langages
Expressions régulières
Automates finis
Les automates finis Déterministes

L'outil Lex

Principes de fonctionnement (des analyseurs lexicaux)

Plan

Rappels de la théorie des langages
Expressions régulières
Automates finis
Les automates finis Déterministes

L'outil Lex

Principes de fonctionnement (des analyseurs lexicaux)

Expressions régulières : syntaxe

Soit Σ un alphabet et Σ^* l'ensemble de mots sur Σ .

On définit $\mathcal{R}(\Sigma)$, l'ensemble des expressions régulières sur Σ , par induction :

- $\varepsilon \in \mathcal{R}(\Sigma)$,
- si $a \in \Sigma$, alors $a \in \mathcal{R}(\Sigma)$,
- si $r_1, r_2 \in \mathcal{R}(\Sigma)$ alors
 - ▶ $r_1 r_2 \in \mathcal{R}(\Sigma)$,
 - ▶ $r_1 | r_2 \in \mathcal{R}(\Sigma)$,
- si $r \in \mathcal{R}(\Sigma)$ alors $r^* \in \mathcal{R}(\Sigma)$.

Exemple : pour $\Sigma = \{a, b\}$,

$$(ab)^*abb, (ab)^*ba(ab)^*, (ab^*ba)^* \in \mathcal{R}(\Sigma).$$

Expressions régulières : sémantique

Un langage sur Σ est un sous-ensemble de Σ^* .
Pour $r \in \mathcal{R}(\Sigma)$ on définit $L(r)$ (le langage de r) :

$L(\varepsilon) = \{\varepsilon\}$, le singleton contenant le mot vide
 $L(a) = \{a\}$,

$L(r_1 r_2) = \{w_1 \cdot w_2 \mid w_1 \in \mathcal{L}(r_1) \text{ et } w_2 \in \mathcal{L}(r_2)\}$,
où \cdot denote la concaténation de mots

$L(r_1 | r_2) = L(r_1) \cup L(r_2)$,

$L(r^*) = \bigcup_{n \geq 0} \underbrace{L(r) \bullet \dots \bullet L(r)}_{n \text{ fois}}$,

c'est-à-dire :

$w \in L(r^*)$ ssi $\exists n \geq 0, \exists w_i \in L(r), i = 1, \dots, n$, tels que
 $w = w_1 \cdot \dots \cdot w_n$.

Des raccourcis

$r^+ = r(r)^*$

$L(r^+) = \{w \in \Sigma^* \mid \text{une ou plusieurs fois un mot filtré par } r\}$

$r? = \varepsilon | r$

$L(r?) = \{\varepsilon\} \cup L(r)$.

En Unix :

si $\Sigma = \{a_1, \dots, a_n\}$ et $J = \{a_{i_1}, a_{i_2}, \dots, a_{i_k}\} \subseteq \Sigma$:

$[a_{i_1} a_{i_2} \dots a_{i_k}] = a_{j_1} | a_{j_2} | \dots | a_{j_k}$
 $[a_i - a_j] = a_i | a_{i+1} | \dots | a_j$.

Des Exemples

$L(a^*) = \{w \in \Sigma^* \mid \text{seulement des } a, \text{ longueur arbitraire}\}$,
 $L((ab)^*) = \{\text{les suites de } ab\}$,
 $L(b?(ab)^*a?) = \{w \in \Sigma^* \mid \text{tout } a \text{ est suivi par } b \text{ et}$
tout b est suivi par $a\}$.

Soit $r \in \mathcal{R}(\Sigma)$ et $w \in \Sigma^*$. On dit que :

- le modèle r reconnaît le mot w ,
- le motif r filtre w ,

si

$w \in L(r)$.

Définitions régulières : un exemple

Les nombres non signés en Pascal :

chiffre $\rightarrow [0-9]$
chiffres $\rightarrow \text{chiffre}^+$
fraction_opt $\rightarrow (. \text{chiffres})?$
exposant_opt $\rightarrow (E(+|-|\varepsilon) \text{chiffres})?$
nb $\rightarrow \text{chiffres fraction_opt exposant_opt}$

On pose :

$r(\text{chiffre}) = [0-9]$
 $r(\text{chiffres}) = [0-9]^+$
 $r(\text{fraction_opt}) = (. [0-9]^+)?$
 $r(\text{exposant_opt}) = (E(+|-|\varepsilon) [0-9]^+)?$
 $r(\text{nb}) = [0-9]^+ (. [0-9]^+)? (E(+|-|\varepsilon) [0-9]^+)?$

De la forme

$$\begin{aligned} D_1 &\rightarrow r_1 \\ D_2 &\rightarrow r_2 \\ &\vdots \\ D_n &\rightarrow r_n \end{aligned}$$

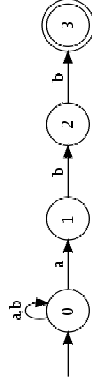
où $r_i \in \mathcal{R}(\Sigma \cup \{D_1, \dots, D_{i-1}\})$, pour $i = 1, \dots, n$.

On reconstruit des expressions régulières $r(D_i) \in \mathcal{R}(\Sigma)$:

$$r(D_i) = r_i[r(D_1)/D_1, \dots, r(D_{i-1})/D_{i-1}].$$

Exemple

Représentation graphique :



C'est à dire :

- $Q = \{0, 1, 2, 3\}$,
 - $\Sigma = \{a, b\}$,
 - Δ
- | Δ | a | b |
|----------|-------------|-------------|
| 0 | $\{0, 1\}$ | $\{0\}$ |
| 1 | \emptyset | $\{2\}$ |
| 2 | \emptyset | $\{3\}$ |
| 3 | \emptyset | \emptyset |
- $i = 0, F = \{3\}$.

Un Automate Fini Non-déterministe (AFN) est un tuple $\mathcal{A} = \langle Q, \Sigma, \Delta, i, F \rangle$ où

- Q est un ensemble fini d'états,
- Σ est son alphabet,
- $\Delta : Q \times \Sigma \rightarrow \mathcal{P}(Q)$ est la fonction de transition,
- $i \in Q$ est l'état initial,
- $F \subseteq Q$ sont les états finaux.

Acceptation par AFN

On dit qu'un mot

$$w = a_1 a_2 \dots a_n$$

est accepté par l'AFN \mathcal{A} s'il existe $n \geq 0$ et e_0, \dots, e_n tels que

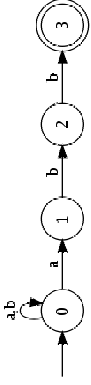
- $e_0 = i$ et $e_n \in F$,
- $e_j \in \Delta(e_{j-1}, a_j)$, pour $j = 1, \dots, n$.

On pose :

$$L(\mathcal{A}) = \{ w \in \Sigma^* \mid w \text{ est accepté par } \mathcal{A} \}.$$

Exemple

Pour \mathcal{A} l'automate



on a

$$L(\mathcal{A}) = \{w \in \{a, b\}^* \mid w \text{ se termine par } abb\}.$$

Le problème $w \in L(\mathcal{A})$ (I)

- le non déterminisme est un problème,
- ... qu'on peut se résoudre si on parcourt tous les chemins de i en parallèle.

- $S_{current} = \{i\}$
- tant que $w \neq \epsilon$ faire
- $a = \text{tete}(w)$ ($a \in \Sigma$), $S_{next} = \emptyset$
- pout tout $q \in S_{current}$ faire
- $S_{next} = S_{next} \cup \Delta(q, a)$
- $w = \text{queue}(w)$, $S_{current} = S_{next}$
- si $S_{current} \cap F \neq \emptyset$ accepter
- sinon refuser

Complexité : $O(|w| \times |Q|)$.

Le Théorème de Kleene

Théorème. Pour $L \subseteq \Sigma^*$, les faits suivants sont équivalents :

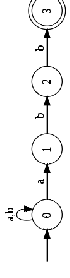
- il existe une expression régulière r tel que $L = L(r)$,
- il existe un AFN \mathcal{A} tel que $L = L(\mathcal{A})$.

S. C. Kleene.

Representation of events in nerve nets and finite automata. In *Automata studies*, Annals of mathematics studies, no. 34, pages 3–41. Princeton University Press, Princeton, N. J., 1956.

Par exemple,

$$L((a|b)^*abb) = L(\mathcal{A})$$

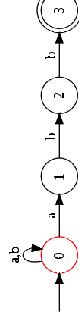


où \mathcal{A} est l'automate

Un calcul parallèle

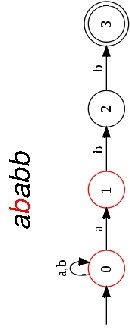
... avec l'automate \mathcal{A} et le mot $ababb$.

ababb



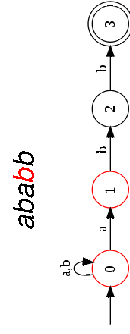
Un calcul parallèle

... avec l'automate \mathcal{A} et le mot $ababb$.



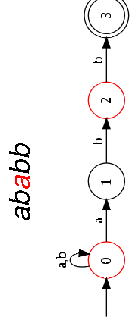
Un calcul parallèle

... avec l'automate \mathcal{A} et le mot $ababb$.



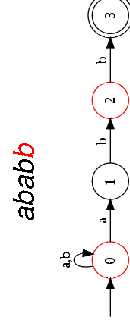
Un calcul parallèle

... avec l'automate \mathcal{A} et le mot $ababb$.

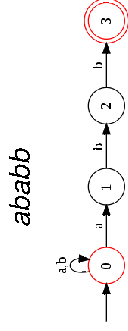


Un calcul parallèle

... avec l'automate \mathcal{A} et le mot $ababb$.



... avec l'automate \mathcal{A} et le mot $ababb$.



Le Théorème de Rabin-Scott

Théorème. Pour $L \subseteq \Sigma^*$, les faits suivants sont équivalents :

- il existe un AFN \mathcal{A} tel que $L = L(\mathcal{A})$.
- il existe un AFD \mathcal{D} tel que $L = L(\mathcal{D})$.



M. O. Rabin and D. Scott.
Finite automata and their decision problems.
IBM J. Res. Develop., 3 :114–125, 1959.

Les AFD

- Un Automate Fini **D**éterministe (AFD) est un AFN $\langle Q, \Sigma, \Delta, i, F \rangle$ tel que $\Delta(q, a)$ est vide ou un singleton.
- La table d'un AFD contient au plus un seul élément dans chaque case.
- *Scourant* est toujours un singleton (ou l'ensemble vide)
- La calcul de $w \in L(\mathcal{A})$ se fait en temps $O(|w|)$

Idée de la preuve (AFN \rightarrow AFD)

- Étant donné \mathcal{A} et pour tout mot w , on se pose la question si $w \in L(\mathcal{A})$
Comme avant, on fait ces calculs par en *parallèle*.
- Cela revient à construire/explore l'automate des sous-ensembles

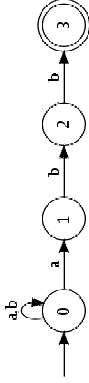
$$\mathcal{P}(\mathcal{A}) = \langle Q', \Sigma, \Delta', i', F' \rangle$$

où :

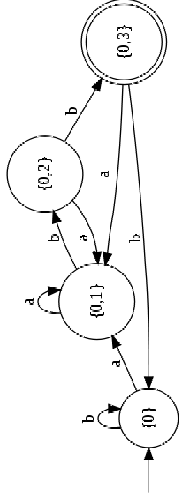
- ▶ $Q' = \mathcal{P}(Q)$,
- ▶ $\Delta'(S, a) = \bigcup_{s \in S} \Delta(s, a)$,
- ▶ $i' = \{i\}$ et $S \in F'$ ssi $S \cap F \neq \emptyset$.

Notre exemple

Pour l'automate \mathcal{A}



on obtient :



Un algorithme pour construire $\mathcal{P}(\mathcal{A})$

- 1 $a_explorer = \{ \{ i \} \}$
- 2 $visites = \emptyset$
- 4 tant que $a_explorer \neq \emptyset$
- 5 choisir $Scourant \in a_explorer$
- 7 pour tout $a \in \Sigma$ faire
- 8 $S_{next} = \emptyset$
- 9 pour tout $s \in Scourant$
- 10 $S_{next} = S_{next} \cup \Delta(q, a)$
- 11 poser $\Delta'(Scourant, a) = S_{next}$
- 13 si $S_{next} \not\subseteq a_explorer \cup visites$
- 14 ajouter S_{next} à $a_explorer$
- 16 enlever $Scourant$ de $a_explorer$
- 17 ajouter $Scourant$ à $visites$

Caveats

- L'automate $\mathcal{P}(\mathcal{A})$ peut avoir $2^{|\mathcal{Q}|}$ états.
- Pour le calculer on peut avoir besoin de $2^{|\mathcal{Q}|}$ étapes de calculs.
- On peut réduire (minimiser) le nombre d'états d'un AFD.
- Le langage

$$(a|b)^* \underbrace{a(a|b) \dots (a|b)}_{\text{rfois}}$$

est le langage d'un AFN de taille $O(n)$.
 Tout AFD qui reconnaît ce langage a taille au moins 2^n états.

Le problème $w \in L(\mathcal{A})$ (II)

- Calculer si $w \in L(\mathcal{A})$ en parallèle revient à calculer avec $\mathcal{P}(\mathcal{A})$.
- Cet algorithme pour décider si $w \in L(\mathcal{A})$ revient à construire $\mathcal{P}(\mathcal{A})$ à la volée.
- On peut l'améliorer en mettant les calculs dans une cache. On parle alors de construction paresseuse de $\mathcal{P}(\mathcal{A})$.
- Une solution radicalement différente est de construire d'abord tout $\mathcal{P}(\mathcal{A})$ et vérifier ensuite si $w \in L(\mathcal{P}(\mathcal{A}))$. Cela pose des problèmes d'espace, mais est très vite.

Rappels de la théorie des langages

- Expressions régulières
- Automates finis
- Les automates finis Déterministes

L'outil Lex

Principes de fonctionnement (des analyseurs lexicaux)

C'est un compilateur d'analyseurs lexicaux :



Écrit par Lesk en 1975.
Utilisé par nombreux compilateurs sous Unix.

Flex (fast lexical analyzer generator) :
écrit par Vern Paxson en 1987.
Ce n'est part du projet GNU.

Le langage Lex

Trois sections (séparées par %%) :

- Définitions** régulières (options, code C en préambule, ...)
- Suite de **règles** de la forme

$$\begin{array}{l} \vdots \\ m_j \{ a_j \} \\ \vdots \end{array}$$

où m_j est un motif et a_j est l'action associée à ce motif.

- ▶ Chaque m_j est une expression régulière sur l'alphabet ASCII + définitions,
 - ▶ Chaque a_j est un morceau de code C (à déclencher si le motif m_j est reconnu).
- Code C** : cette section définit des fonctions, à inclure dans le code C, dont on peut avoir besoin (`main`, `yywrap`, ...)

La loi du plus long lexème (leftmost longmost POSIX)

Soit w un préfixe de u :

$$u = w \cdot w'.$$

Supposons que

$$w \in L(m_i), \quad u \in L(m_j).$$

Filter w ou u ?

On filtre la chaîne de caractères la plus longue, on applique la règle lui associée :

$$m_j \quad \{ a_j \}$$

Un exemple

Considérons les deux règles

```
if || let  
[a-zA-z]([a-zA-z0-9])* {printf("Mot clé\n");}  
{printf("Identificateur\n");}
```

Si le flot d'entrée est de la forme

```
let ifin = 4 in ...
```

La chaîne `ifin` est reconnue comme un identificateur, (et non pas comme mot clé).

“Identificateur” est affiché à l'écran.

Exemple

Considérons les deux règles

```
if || let  
[a-zA-z]([a-zA-z0-9])* {printf("Mot clé\n");}  
{printf("Identificateur\n");}
```

Si le flot d'entrée est de la forme

```
if in = 4 then ...
```

La chaîne `if` est reconnue comme mot clé (et non pas comme identificateur). “Mot clé” est affiché à l'écran.

La loi de priorité des règles

Considérons des règles

```
mi { ai }  
.  
.  
.  
mj { aj }
```

et supposons que

$$w \in L(m_i) \cap L(m_j).$$

Faut-il déclencher l'action a_i où a_j ?

On déclenche l'action plus prioritaire, c'est-à-dire a_i .

Remarque

La loi du plus long lexème est plus forte que la loi des priorités entre règles.

Un erreur qui se produit si on oublie ce fait :

```
if | then | else {printf("Mot clé\n");}  
.  
.  
.  
.+ {printf("Un erreur trop long\n");}
```

La chaîne

```
if x := 3 then  
4 else 4
```

est divisé ainsi :

```
if x := 3 then 4 else 4 ...  
                  erreur
```

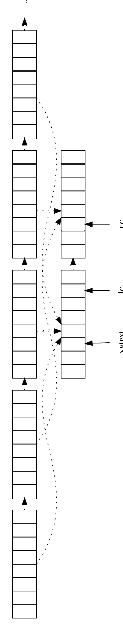
Rappels de la théorie des langages

- Expressions régulières
- Automates finis
- Les automates finis Déterministes

L'outil Lex

Principes de fonctionnement (des analyseurs lexicaux)

Le double tampon circulaire



- Circularité : les déplacements se font modulo $2 * \text{TAILLETAMPON}$
- Trois pointeurs :
 - ▶ yytext : début du lexème,
 - ▶ fc : fin courante du lexème,
 - ▶ cc : caractère courant.
- Quand cc s'approche à la fin d'un tampon, un nouveau tampon est transféré du flot vers la mémoire.
- Limite de la taille des lexèmes :
 - ▶ si $\text{cc} - \text{yytext} > 0$ alors $\text{cc} - \text{yytext} < \text{TAILLETAMPON}$,
 - ▶ si $\text{cc} - \text{yytext} < 0$ alors $\text{yytext} - \text{cc} > \text{TAILLETAMPON}$.

- Le flot de caractères en entrée est
 - le plus souvent – un fichier.
- Les transferts du disque à la mémoire vive sont coûteux.
- Transférer 100 fois un caractère est beaucoup plus coûteux que transférer 1 fois 100 caractères.

Des expressions régulières aux automates

Soit

$$m_1 \{ a_1 \}$$

$$\vdots$$

$$\vdots$$

$$m_n \{ a_n \}$$

un ensemble de règles.

1. On construit un AFD $\mathcal{D} = \langle Q, i, \Delta, F \rangle$ tel que

$$L(\mathcal{D}) = L(m_1 | \dots | m_n).$$

2. On construit aussi une fonction

$$\lambda : F \rightarrow \{1, \dots, n\}$$

telle que si $w = a_1 \dots a_k$ et

$$i \xrightarrow{a_1} q_1 \dots q_{k-1} \xrightarrow{a_k} q_k \in F$$

et $\lambda(q_k) = i$, alors $w \in L(m_i)$ et $w \notin L(m_j)$ pour $j < i$.

Fonctionnement de l'analyseur lexical

```
7  repeter
8  etat_courant =  $\Delta$ (etat_courant,*cc)

10 si etat_courant=erreur
11 si yylex < fc
12   excuter l'action  $\mathcal{A}_{priorite}$ 
13 sinon
14   /* yylex = fc, on a pas progressé */
15   excuter l'action par default
16   fc++
17   sortir de la boucle
```

Navigation icons: back, forward, search, etc.

36/38

Navigation icons: back, forward, search, etc.

37/38

```
19  sinon
20  si etat_courant  $\in F$ 
21  fc=cc
22  priorite =  $\lambda$ (etat_courant)

24  cc++
```

Navigation icons: back, forward, search, etc.

38/38

Navigation icons: back, forward, search, etc.

37/38