

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

$L(a) = \{a\}$,

$L(r_1 r_2) = \{w_1 \cdot w_2 \mid w_1 \in \mathcal{R}(r_1) \text{ et } w_2 \in \mathcal{R}(r_2)\}$,
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 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

$w \in L(r).$

Des raccourcis

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

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

$$r? = \varepsilon|r$$

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

En Unix :

si $\Sigma = \{ a_1, \dots, a_n \}$ et $J = \{ a_{j_1}, a_{j_2}, \dots, a_{j_k} \} \subseteq \Sigma$:

$$[a_{i_1} a_{i_2} \dots a_{j_k}] = a_{j_1} | a_{j_2} | \dots | a_{j_k}$$

$$[a_i - a_j] = a_i | a_{i+1} | \dots | a_j.$$

Définitions régulières : un exemple

Les nombres non signés en Pascal :

```
chiffre → [0 - 9]
chiffres → chiffre+
fraction_opt → (.chiffres)?
exposant_opt → (E(+|-|ε)chiffres)?
nb → chiffres fraction_opt exposant_opt
```

On pose :

```
r(chiffre) = [0 - 9]
r(chiffres) = [0 - 9]+
r(fraction_opt) = ([0 - 9]+)?
r(exposant_opt) = (E(+|-|ε)[0 - 9]+)?
r(nb) = [0 - 9]+([0 - 9]+)?(E(+|-|ε)[0 - 9]+)?
```


Définitions régulières

De la forme

$$D_1 \rightarrow r_1$$

$$D_2 \rightarrow r_2$$

$$\vdots$$

$$D_n \rightarrow r_n$$

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

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

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

Les AFN

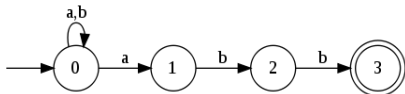
Un **A**utomate **F**ini **N**on-déterministe (AFN) est un tuple

$\mathcal{A} = \langle Q, \Sigma, \Delta, i, F \rangle$ 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.

Exemple

Représentation graphique :



C'est à dire :

- $Q = \{0, 1, 2, 3\}$,

- $\Sigma = \{a, b\}$,

- | Δ | a | b |
|----------|-------------|-------------|
| 0 | $\{0, 1\}$ | $\{0\}$ |
| 1 | \emptyset | $\{2\}$ |
| 2 | \emptyset | $\{3\}$ |
| 3 | \emptyset | \emptyset |

- $i = 0, F = \{3\}$.

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

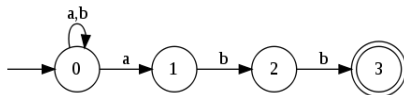
- $e_0 = i$ et $e_n \in F$,
- $e_i \in \Delta(e_{i-1}, a_i)$, pour $i = 1, \dots, n$.

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 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.

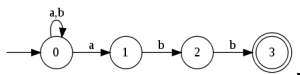
34, pages 3–41. Princeton University Press, Princeton, N.

J., 1956.

Par exemple,

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

où \mathcal{A} est l'automate



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.

```
1       $S_{courant} = \{i\}$ 
2      tant que  $w \neq \varepsilon$  faire
3           $a = tete(w)$  ( $a \in \Sigma$ ),  $S_{next} = \emptyset$ 

5          pout tout  $q \in S_{courant}$  faire
6               $S_{next} = S_{next} \cup \Delta(q, a)$ 

8           $w = queue(w)$ ,  $S_{courant} = S_{next}$ 

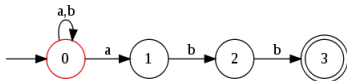
10     si  $S_{courant} \cap F \neq \emptyset$  accepter
11     sinon refuser
```

Complexité : $O(|w| \times |Q|)$.

Un calcul parallèle

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

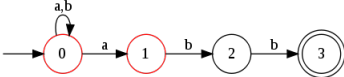
ababb



Un calcul parallèle

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

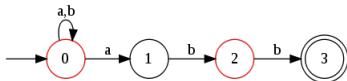
ababb



Un calcul parallèle

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

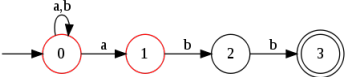
ababb



Un calcul parallèle

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

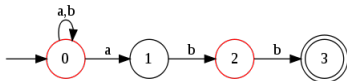
ababb



Un calcul parallèle

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

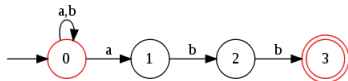
ababb



Un calcul parallèle

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

ababb



Les AFD

- Un **A**utomate **F**ini **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.
- $S_{courant}$ est toujours un singleton (ou l'ensemble vide)
- Le calcul de $w \in L(\mathcal{A})$ se fait en temps $O(|w|)$

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.

Idée de la preuve (AFN \rightarrow AFD)

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

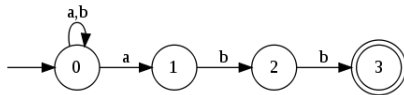
$$\mathcal{P}(\mathcal{A}) = \langle Q', \Sigma, \Delta', i', F' \rangle$$

où :

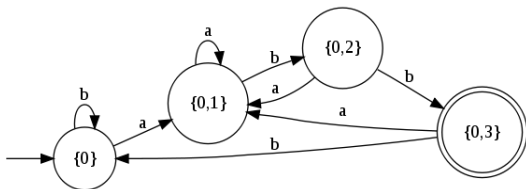
- ▶ $Q' = \mathcal{P}(Q)$,
- ▶ $\Delta'(S, a) = \bigcup_{s \in S} \Delta(s, a)$,
- ▶ $i' = \{i\}$ et $S \in F'$ ssi $S \cap F \neq \emptyset$.

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  $S_{courant} \in a\_explorer$ 

7       pour tout  $a \in \Sigma$  faire
8            $S_{next} = \emptyset$ 
9           pour tout  $s \in S_{courant}$ 
10               $S_{next} = S_{next} \cup \Delta(q, a)$ 
11              poser  $\Delta'(S_{courant}, a) = S_{next}$ 

13              si  $S_{next} \notin a\_explorer \cup visites$ 
14                  ajouter  $S_{next}$  à  $a\_explorer$ 

16              enlever  $S_{courant}$  de  $a\_explorer$ 
17              ajouter  $S_{courant}$  à  $visites$ 
```

Caveats

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

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

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.

Plan

Rappels de la théorie des langages

Expressions régulières

Automates finis

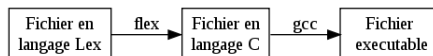
Les automates finis Déterministes

L'outil Lex

Principes de fonctionnement (des analyseurs lexicaux)

L'outil Lex

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 %%):

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

$$\begin{array}{c} \vdots \\ m_i \quad \{ a_i \} \\ \vdots \end{array}$$

où m_i est un motif et a_i est l'action associée à ce motif.

- ▶ Chaque m_i est une expression régulière sur l'alphabet ASCII + définitions,
 - ▶ Chaque a_i est un morceau de code C (à déclencher si le motif m_i est reconnu).
3. **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

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

Filtrer 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                {printf("Mot clé\n");}  
[a-zA-z]([a-zA-z0-9])* {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.

La loi de priorité des règles

Considérons des règles

$$\begin{array}{l} m_i \quad \{ a_i \} \\ \cdot \\ \cdot \\ \cdot \\ m_j \quad \{ a_j \} \end{array}$$

et supposons que

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

Faut-il déclencher l'action a_i ou a_j ?

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

Exemple

Considérons les deux règles

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

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 :

$\underbrace{\text{if } x := 3 \text{ then}}_{\text{erreur}} \underbrace{4 \text{ else } 4}_{\text{erreur}} \dots$

Plan

Rappels de la théorie des langages

Expressions régulières

Automates finis

Les automates finis Déterministes

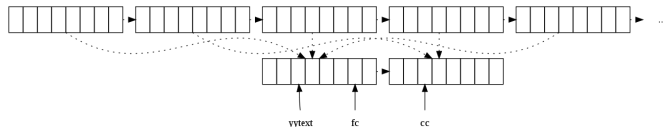
L'outil Lex

Principes de fonctionnement (des analyseurs lexicaux)

Le problème des I/O

- 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 caractères est beaucoup plus coûteux que transférer 1 fois 100 caractères.

Le double tampon circulaire



- Circularité : les déplacements se font modulo $2 * \text{TAILLETAMPON}$
- Trois pointeurs :
 - ▶ `yytext` : debut du lexème,
 - ▶ `fc` : fin courante du lexème,
 - ▶ `cc` : caractere 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 $cc - yytext > 0$ alors $cc - yytext < \text{TAILLETAMPON}$,
 - ▶ si $cc - yytext < 0$ alors $yytext - cc > \text{TAILLETAMPON}$.

Des expressions régulières aux automates

Soit

$$\begin{array}{c} m_1 \{ a_1 \} \\ \cdot \\ \cdot \\ \cdot \\ m_n \{ a_n \} \end{array}$$

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

```
1  repeter  
  
3      etat_courant = i  
4      ylex = cc = fc  
5      priorite = 1
```

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

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

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

24     cc++
```