

Conception et pratique de l'algorithmique

Module CPA

Philippe Trébuchet

Université Pierre et Marie Curie

Second semestre 2014-2015

Version du 10 janvier 2015

Première partie I

Les Expressions régulières

Plan

1

Définition

Les expressions régulières : généralités

- Langage de Description de chaînes de caractères.
- Utilisé pour reconnaître un motif dans un un corpus de texte
- plusieurs dialectes différents

Rappels : Les ERE Posix : Le tronc commun avec les autres dialectes

Un caractère non spécial se reconnaît lui-même.

- `.` signifie n'importe quel caractère.
- `[c1c2c3...]` signifie *un* caractère parmi c_1 , c_2 ou c_3
- `[^c1c2c3...]` signifie *un* caractère qui n'est pas c_1 , c_2 ou c_3
- `\ char_special` signifie le caractère spécial est considéré comme un caractère normal (`\` enlève la signification spéciale d'un caractère spécial).
- `^` signifie début de ligne
- `$` signifie fin de ligne.
- Dans un bloc crochet `[]` `$` et `^` (en dehors du premier caractère. perdent leur signification spéciale !

Les ERE Posix (Suite)

- les caractères spéciaux des ERE sont : `. + * { } () ^ $ [] | ?`
- `()` définition de bloc parenthésé. **pas de notion de rappel du i-ème bloc parenthésé**
- `(ERE1 | ERE2)` alternance : `ERE1` ou `ERE2`.

multiplicateur d'occurrences

- $ERE\{n\}$
- $ERE\{m, n\}$
- $ERE\{m, \}$
- $ERE\{0, n\}$
- ERE^+ signifie : au moins une fois l'ERE
- $ERE?$ zéro ou une fois l'ERE

Exemples d'ERE

- reconnaître une chaîne formée de `a` :
 - idem BRE `^a*$`
- reconnaître une chaîne contenant entre 3 et 5 `a` consécutifs.
 - `a{3,5}` (attention en ERE `{` et `}` sont des caractères spéciaux)
- reconnaître une chaîne contenant la chaîne `toto` répétée entre 3 et 5 fois d'affilée
 - `(toto){3,5}`
- reconnaître une chaîne commençant soit la chaîne `toto`, soit la chaîne `tutu`
 - `^(toto|tutu)` Les parenthèses sont obligatoires

Les outils GNU définissent en plus de RE précédentes un certain nombre d'extensions **pratiques** :

- `\b` match une chaîne vide à une frontière de mot.
- `\<` et `\>` matchent une chaîne vide en début et fin de mot.
- `\w` match un word constituant
- `\s` est un synonyme pour `[[:space:]]`

Des outils de filtrage puissants

- `grep` Get Regular Expression and Print
- **Synopsis** : `grep option BRE [noms de fichiers]`
- Si aucun nom de fichier n'est donné `stdin` est utilisé.
 - `-l` affiche le nom du fichier au lieu de la ligne.
 - `-E` utilise de ERE au lieu de BRE.
 - `-F` utilise des chaînes de caractères pour rechercher.
 - `-c` compte le nombre de matchs.
 - `-e` plusieurs motifs.
 - `-i` case insensitif.
 - `-q` quiet.
 - `-s` supprime les messages d'erreur.
 - `-v` inverse la sélection.

Exemples de `grep`

- Afficher les lignes contenant la chaîne `toto` :
 - `grep -f toto des_noms_de_fichiers`
- Afficher le nom des fichiers contenant des lignes commençant par `toto`
 - `grep -l '^toto' des_noms_de_fichiers`
- Tester si le fichier `fic` contient bien la chaîne de caractère `string` en majuscule ou minuscule et si oui afficher `ok`
 - `grep -iqs 'string' fic && echo ok`
- On peut aussi utiliser `grep` en conjonction avec l'option `exec` de `find` !
 - `find /etc/ -exec grep -il 'printer' {} \;`

Rappel Lexing/Parsing : Lexeur

- c'est un logiciel qui va reconnaître dans le texte une suite de caractères et va leur associer un jeton.
- est essentiellement un moteur d'inférence d'expressions régulières
- Les jetons (token) seront transmis au parser qui lui exécutera la grammaire.

(f)Lex

- Logiciel ancestral de génération de lexeurs

```
%{  
    Code Verbatim  
}%  
definition
```

```
%%  
ERE {ACTION}  
<context>ERE {ACTION}
```

```
%%  
Code
```

Lexeur à la main

```
#!/usr/bin/perl
sub lexer
{
    my $string=shift;
    return sub {
        return "LPAREN" if ($string=~/\G\( /cg);
        return "RPAREN" if ($string=~/\G\) /cg);
        return "PIPE" if ($string=~/\G\| /cg);
        return "LBRACKET" if ($string=~/\G\[ /cg);
        return "RBRACKET" if ($string=~/\G\] /cg);
        return "POINT" if ($string=~/\G\. /cg);
        return "STAR" if ($string=~/\G\* /cg);
        return "PLUS" if ($string=~/\G\+ /cg);
        return "QUESTIONMARK" if ($string=~/\G\? /cg);
        return "CHAR" if ($string=~/\G\w /cg);
        return "0";
    }
}
$lex=lexer($ARGV[0]);
print "$toto_" while ($toto=$lex->());
```

Rappel Lexing/Parsing : Parseur

- logiciel qui interprete le flux de token selon une grammaire.
- Parseur reccurssif descendant
- Parseur ascendant
- Production d'un *arbre de syntaxe* (AST).

- Logiciel ancestral de generation de parseur LALR (Ascendant)

```
%{  
    code verbatim  
}%  
definition de token  
definition de macros  
  
%%  
  
grammaire  
  
%%  
  
code verbatim
```

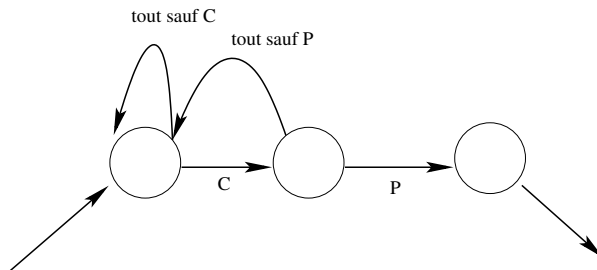

A la main

```
sub parser {  
  my $tmp;  
  my $lex=lexer(shift);  
  my $parse=sub {  
    my $res="";  
    my $token=$lex->();  
    if ($token=~ /LBRACKET/) {  
      $res.="[";  
      $res.=$token until (($token=$lex->())=~ /RBRACKET");  
      return $res."]";  
    }  
    if ($token eq /LPAREN/) {  
      $res.="(".$lex->() until ($lex->() eq "RPAREN");  
      return $res.")";  
    }  
    # if ($token=~ /  
    print "aaaaa";  
    return $res;  
  };  
}
```

Complexité

- Lexing : Linéaire !
- Parsing : Dépend beaucoup de la grammaire
 - de $\mathcal{O}(n^3)$.
 - à exponentielle.

Reconnaitre une chaîne de caractères



Deuxième partie II

Les automates

Définition

Un automate fini sur un alphabet \mathcal{A} est la donnée de :

- un ensemble d'états E
- une fonction de transition f
- un état initial i
- un ensemble d'états finaux O

Accepté ou pas ?

- Un mot m est accepté par l'automate si et seulement si en partant de l'état initial et en ne suivant que les transitions libellées par les lettres de m on peut aboutir à un état final.

Définition

L'ensemble des mots acceptés par un automate s'appelle le langage accepté.

Définition

Deux automates sont équivalents s'ils reconnaissent le même langage.

Automate déterministe vs non déterministe

definition

Si pour tout état e il n'y qu'au plus une transition étiquetée par une lettre a de l'alphabet, on dit que l'automate est déterministe. On dit que l'automate est non déterministe sinon.

On dit qu'un automate est complet si pour tout état e et toute lettre a , il existe au moins une transition qui part de e et qui est étiquetée par la lettre a .

- une transition est une ε -transition si elle n'est étiquetée par aucune lettre.

Un automate avec epsilon transition est un automate non déterministe.

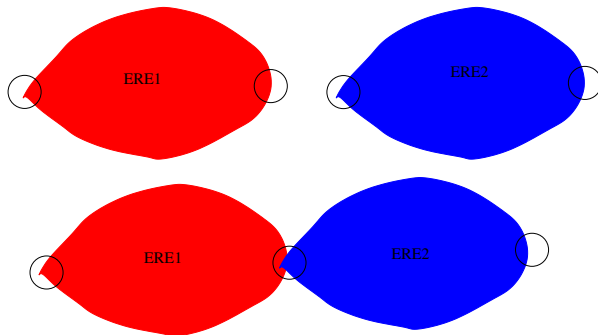
algorithme de reconnaissance

- Pour un automate fini déterministe il suffit de suivre les transitions. Le coût d'un match est linéaire en la taille de l'entrée.
- Pour un automate non déterministe, il faut explorer toutes les possibilités jusqu'à trouver un état acceptant ou épuiser toutes les possibilités. Le coût peut être exponentiel !

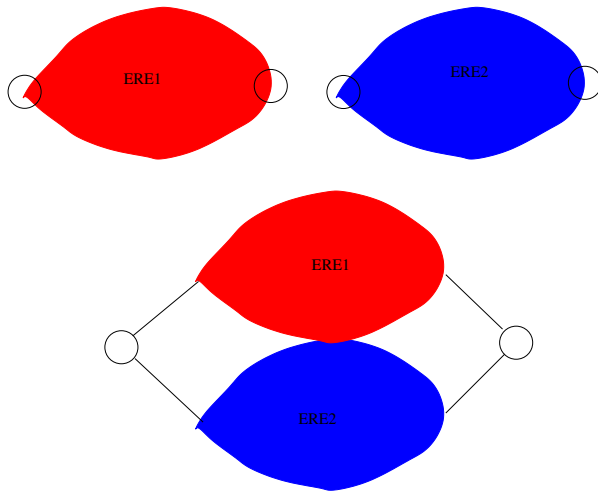
La construction de Thompson

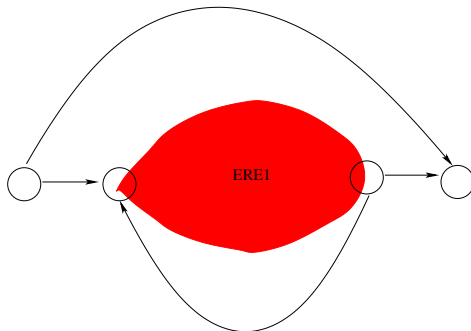
ERE1ERE2	On concatene
ERE*	une epsilon transition vers le debut de l'ERE
ERE+	ERE concaténé à ERE*
(ERE1—ERE2)	un état + 2 epsilon transitions vers ERE1 et ERE2
ERE ?	une epsilon transition vers la suite ou vers ERE

Concatenation



Union





Propriétés

Propriété

Si l'expression régulière est formée de n symbols (lettre + opérateurs), l'automate produit a $2n$ états.

Propriété

La traduction d'une expression régulière en automate est linéaire en la taille de l'expression.

Propriété

Chaque état de l'automate produit est la source (resp. la destination) d'une ou deux transitions (sauf les états initiaux et finaux).

Propriété

La construction précédente ne tient pas compte de l'associativité des opérateurs.

- les états peuvent être numérotés par des entiers
- la fonction de transition peut être représentée comme une matrice
 - d'entiers (déterministe)
 - de liste d'entiers (non déterministe)

Propriétés des automates

- Soit L un langage reconnaissable, il existe un automate déterministe le reconnaissant.
- Il existe des langages non reconnaissables.
- si L et L' sont des langages reconnaissables alors LL' l'est aussi

Lemme de l'étoile

Si L est un langage reconnaissable, alors il existe un entier N tel que tout mot m de taille supérieure à N peut se factoriser en trois sous mots u, v, w tels que, v est non vide, $m = uvw$ et $\forall n \in \mathbb{N}, uv^nw \in L$.

Opérations sur les automates

- Si on inverse les flèches d'un automate et qu'on échange les états initiaux et finaux on a un automate qui reconnaît le langage miroir.
- Il existe plusieurs automates pour le même langage.
- un automate complet est un automate pour lequel, pour chaque état e et pour chaque lettre $/$ il existe une transition partant de e et étiquetée par $/$.

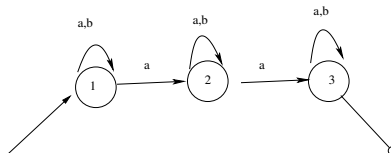
Proposition

Tout automate A est équivalent à un automate D complet et déterministe. Si A est fini avec n états, alors D peut être construit avec au plus 2^n états.

algorithme

- input : $A = (E, Q, I, F)$ un automate *sans epsilon transitions*
- output : B un automate fini deterministe equivalent à A .
- $Todo = I, E_B = \emptyset, Q_B = \emptyset, I_B = I$
- Tq $Todo$ est non vide faire
 - $S = pop(Todo)$
 - pour chaque lettre l
 - Soit $E = e_1, \dots, e_k$ l'ensemble des états accessibles depuis S par une transition étiquetée par l
 - si E n'est pas déjà dans E_B
ajouter E à E_B
 - ajouter une transition de S à E dans Q_B
- F_B est égal à l'ensemble des états de Q_B qui contiennent un état final de F .
- retourner $B = (E_B, Q_B, I_B, F_B)$

Exemple

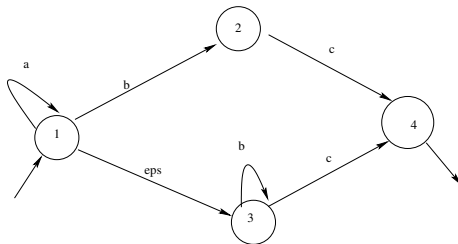


Proposition

Etant donné un automate fini avec *varepsilon*-transitions il existe un automate fini sans *varepsilon*-transitions reconnaissant le même langage.

- L'algorithme fonctionne en remplaçant les chemins de longueur 1 commençant par une *varepsilon*-transitions par une nouvelle transition.

exemple



Minimisation

- L'automate deterministe construit avant est tres tres gros.
- il existe des automates plus petit.
- il existe plusieurs algorithmes pour calculer le plus petit automate équivalent.

Définition

Soit A un automate, deux états x et y de A sont dits indistingables si et seulement si pour tout mot m tq xm est accepté alors ym l'est aussi et réciproquement.

Définition

Le transposé d'un automate est l'automate obtenu en inversant les transitions.

- L'automate transposé reconnaît le langage miroir
- l'automate transposé peut avoir des états inaccessibles

Proposition

Soit A un automate déterministe accessible. Le déterminisé du transposé de A est minimal.

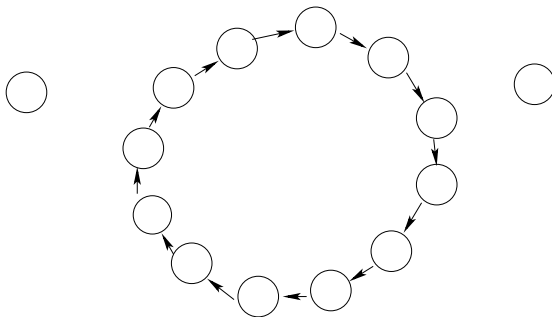
algo

Calculer $\text{det}(\text{Trans}(\text{det}(\text{trans}(A))))$

La construction de Moore

- On met les états finaux d'un côté et les autres de l'autre.
- On sature par la relation d'équivalence de distinguabilité.
- Complexité : $\mathcal{O}(n^2)$.

Contre exemple



Definition

- Soit E un ensemble d'états et a une lettre on note $a^{-1}E$ l'ensemble des q tels que une transition de q étiquetée par a amène dans E .
- Soit B un ensemble d'état on dit que B est scindé par la paire (E, a) Si :
 - $B_1 = B \cap a^{-1}E$
 - $B_2 = B \setminus B_1$

Sont tous les deux non vides. Sinon B est dit stable par (E, a) .

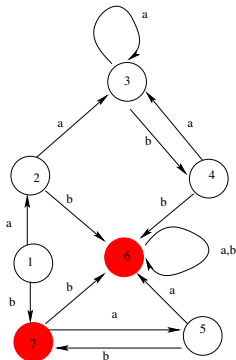
Principe de l'algorithme

- On part d'une partition en deux etat Finaux/Autres
- On maintient une liste de Paires scindantes
- Si une paire scinde un etat E en E_1, E_2 :
 - on remplace l'etat E par les etats E_1 et E_2
 - Si la paire (E, a) etait scindante on la remplace par les deux paires (E_1, a) et (E_2, a) .
 - Sinon on choisi *le plus petit* de B_1 et B_2 et on rajoute les paires (B_i, a) comme paire scindante. pour toute lettre.
 - On s'arrete quand il ne reste plus de paires scindantes.

La construction de Hopcroft

- $L = \emptyset$
- Si $|F| < |Q/F|$ alors
 - $C_0 = Q/F, C_1 = F$ Rajouter C_1 à L .
- sinon
 - $C_1 = Q/F, C_0 = F$ Rajouter C_1 à L .
- $P = \{C_0, C_1\}$
- Tant que $L \neq \emptyset$ faire
 - Extraire C de L
 - Scinder chaque ensemble C_i construit jusqu'à présent par C
 - Ajouter à L la plus petite coupure

exemple



Complexité

- complexité en $\mathcal{O}(n \log(n))$
- en pratique la constante est grosse
- pour des tailles usuelles on utilise plus volontier l'algorithme de Moore.

Des automates aux expressions

Definition

Un automate généralisé est un automate dont les transitions sont étiquetées par des langages rationnels.

Propriétés

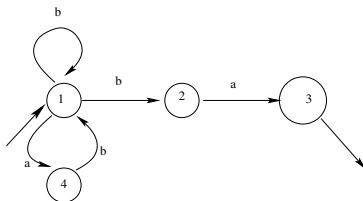
Le langage reconnu par un automate généralisé est un langage rationnel.

La construction de Brzozowski et McCluskey

Pour traduire un automate en expression :

- On commence par rajouter un unique etat initial et un unique etat final à l'automate
- tant qu'il reste un etat autre que ces deux là :
 - choisir un etat et eliminer les transisitions entrante et sortantes

Exemple



Troisième partie III

automates à pile

definition

Un automate à pile est un sextuplet :

- un alphabet de transition
- un ensemble d'etats
- un etat initial
- un ensemble d'etats finaux
- un alphabet de pile
- un ensemble de transitions

Les transitions sont données par des triplets $([i, P], a, [j, Q])$ signifiant ;

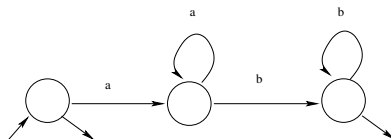
- on passe de l'etat i à l'etat j
- on depile P puis on empile Q

definition

On accepte un mot si et seulement si en partant de pile vide et en suivant les transitions données par le mot on abouti dans un etat acceptant ET avec une pile vide.

Exemple : $a^n b^n$

Un seul symbole de pile !



Théorème

Soit G un grammaire algébrique. Le langage L_G engendré par G peut être reconnu par un automate à pile A_p .

- c'est la construction faite par Yacc !
- il existe une construction inverse !

- états : un seul et unique s !
- alphabet de pile : les terminaux et les variables
- Pour chaque production de la forme $A \rightarrow B$ on ajoute une ε -transition $[s, A] \rightarrow [s, B]$
- Ajout des transitions pour depiler les symboles terminaux.

Définition

Un langage reconnu par un automate à pile est appelé langage algébrique.

Proposition

- l'intersection entre un langage algébrique et un langage régulier est algébrique.
- Il existe des langages algébriques dont l'intersection n'est pas algébrique.
- Les langages algébriques ne sont pas fermés par complémentation.

Lemme de l'étoile pour les langages algébriques

Si L est un langage engendré par une grammaire G alors il existe un entier N tel que tout mot de L de longueur supérieure à N peut se factoriser sous la forme $uxzyv$ tels que :

- x et y sont non vides
- $ux^nzy^n v \in L$ pour tout $n \in \mathbb{N}$.

Quatrième partie IV

Machine de Turing

Definition

Une machine de Turing sera définie par la donnée de :

- un alphabet fini.
- un ruban infini découpé en cases, chaque case pouvant contenir exactement un symbol de l'alphabet fini.
- d'une tête de lecture écriture pouvant se déplacer d'une case a la fois sur le ruban, lire la case où elle se trouve, écrire un symbole dans la case où elle se trouve et changer d'état. La tête ne peut prendre qu'un nombre fini d'états.
- un état initial pour la tête et un état final acceptant signifiant que le calcul s'est bien passé.

On dit qu'une machine de Turing est déterministe si étant donné un état de la tête et un symbole dans la case lue, on n'a qu'une seule action possible. Elle est non déterministe sinon.

Propriétés

- tout ce qu'on peut calculer avec une machine de turing à n rubans est calculable avec une machine à un seul ruban.
- tout ce qui peut se calculer avec un ruban infini dans les deux sens peut aussi être calculé avec un ruban infini dans une seule direction.

- un problème qui est accepté par une machine de Turing déterministe en un nombre polynomial d'opérations appartient à la classe de complexité P .
- Un Problème est dit NP s'il est accepté en temps polynomial par une machine de Turing non déterministe

Completude

- une reduction polynomiale d'un probleme A à un problème B est un calcul polynomial qui transforme la resolution du probleme A en celle du probleme B .
- Un probleme que l'on peut reduire à un autre est plus facile a resoudre.
- Les problemes les plus durs d'une classe sont appelés *complet*.

- Tous les problemes complets sont aussi durs les un que les autres.

Propriétés

On ne connait pas d'algorithme mieux qu'exponentiel pour resoudre un probleme NP-Complet.

Un probleme NP-Complet : SAT

Enoncé

Etant donné une formule logique sous forme normale conjonctive

$\bigvee_{i=1}^n v_i = v_1 \vee v_2 \vee \dots \vee v_n$ existe t'il une valeur des variable qui la satisfasse.

SAT est NP-Complet

Considérons maintenant un problème de NP de taille n , i.e. il existe une machine de Turing non-déterministe à q états sur un alphabet à s lettres résolvant le problème en temps polynomial, soit $P(n)$ le nombre d'opérations nécessaires. Posons

- pour $i \in [0, P(n)]$, $j \in [0, q - 1]$ on introduit $Q_{i,j}$ qui voudra dire qu'après la i -ième opération la machine se retrouve dans l'état j .
- pour $i \in [0, P(n)]$, $j \in [-P(n), P(n)]$ et $k \in [1, s]$ on introduit $S_{i,j,k}$ qui voudra dire qu'après la i -ième opération la case numéro j contient la lettre a_k .
- Pour $i \in [0, P(n)]$, $j \in [-P(n), P(n)]$ on introduit $T_{i,j}$ qui voudra dire qu'après la i -ième opération la machine lit la case j .

Démonstration

- clause stipulant qu'à l'étape i la machine de Turing est dans au moins un état
- a la i -ème opération la machine ne peut pas être dans l'état j et dans l'état k à la fois
- clauses qui sont vraies si et seulement si chaque case de la bande dont l'indice est compris entre $-P(n)$ et $P(n)$ contient exactement un symbol de l'alphabet.
- 'a l'étape i , la machine est en train de lire une et une seule case de la bande.
- a l'étape 0 la machine est dans l'état 1 et est en train de lire la case 0.
- implication qu'à l'étape i si la machine est dans l'état k , qu'elle lit la case j dans laquelle se trouve le symbole l , alors elle passe dans un état $Q_{i+1,t(k,l)}$, où $t(k,l)$ est la fonction exprimant le nouvel état de la machine.
- implication qu'à l'étape i si la machine est dans l'état k , qu'elle lit la case j dans laquelle se trouve le symbole l , alors elle écrit le symbole $u(k,l)$ dans
- sous les mêmes conditions la tête de la machine va à la case $j + d(k,l)$
- ...