

Les Langages

Chapitre 4

1	Introduction	1
1.1	Généralités	1
1.2	Algorithme naïf	1
2	Langages	2
2.1	Alphabet et mots	2
2.2	Langages	4
3	Langages rationnels et expressions régulières	6
3.1	Langages rationnels	6
3.2	Expressions régulières	7
3.3	Implémentation des expressions régulières	9

1 Introduction

1.1 Généralités

Dans ce chapitre nous allons étudier la recherche de motifs. Cela a beaucoup d'applications dans de nombreux domaines. Le principe est de rechercher un motif dans des données de dimension 1 (une suite de caractères) ou de dimension 2 (une image). On peut penser aux exemples suivants :

- Un visage dans une photo
- Un QR code
- Une suite de lettres (ou un motif¹) dans une chaîne de caractères

Dans la suite nous travaillerons essentiellement sur la recherche d'un motif dans une chaîne de caractères. C'est en particulier cela qui sert à un compilateur pour identifier si un programme ne comporte pas d'erreurs de syntaxes.

1.2 Algorithme naïf

L'algorithme principal est le suivant (en prenant pour motif un mot).

- On parcourt le texte.
- Pour chaque position possible du motif dans le texte on teste si la position est une occurrence du motif.
- Si, pour une position donnée on trouve le motif (en testant de gauche à droite) on retourne la position.
- Dans le cas contraire on teste la position suivante.
- Si on est arrivé au bout des positions possibles le motif n'est pas détecté.

Implémentation en Caml :

- En Caml il existe un type `char` pour les caractères. Un caractère est défini entre apostrophes.
- Il existe aussi un type `string` pour les chaînes de caractères. Une chaîne de caractères est définie entre guillemets.
- On accède à la i -ème lettre d'une chaîne `w` par la commande `w.[i]`
- Les chaînes de caractères ne sont pas mutables.

Exercices :

1. Ecrire une fonction `trouve : string -> string -> bool` telle que `trouve texte mot` renvoie un booléen selon que le mot `mot` se trouve ou non dans le texte `texte`.

1. La notion de motif sera précisée plus loin.

Définition 2.1.2 (Concaténation)

Soit w_1 et w_2 deux mots sur \mathcal{A} , la concaténation de w_1 et w_2 notée $w_1.w_2$ ou w_1w_2 est le mot obtenu en mettant bout à bout les mots w_1 et w_2 . On a donc

$$\forall k \in \llbracket 0; |w_1| + |w_2| - 1 \rrbracket, (w_1w_2)[k] = \begin{cases} w_1[k] & \text{si } k < |w_1| \\ w_2[k - |w_1|] & \text{si } k \geq |w_1| \end{cases}$$

Proposition 2.1.3

- L'opération \cdot est une opération interne sur \mathcal{A}^* .
- L'opération \cdot est associative.
- L'opération \cdot admet le mot vide ε comme élément neutre.

Remarque : On dit que (\mathcal{A}^*, \cdot) est un monoïde (ou magma associatif unifié).

Exercice : (\mathcal{A}^*, \cdot) est-il un groupe?

Proposition 2.1.4

Les fonctions longueur $|\cdot|$ et longueur en un caractère $a : |\cdot|_a$ sont additives :

$$\forall (w_1, w_2) \in (\mathcal{A}^*)^2, |w_1w_2| = |w_1| + |w_2| \text{ et } |w_1w_2|_a = |w_1|_a + |w_2|_a.$$

On dit que ce sont des morphismes de monoïdes de (\mathcal{A}^*, \cdot) dans $(\mathbf{N}, +)$

Définition 2.1.5 (Puissances)

Soit w un mot, on définit ses puissances par

$$w^0 = \varepsilon ; w^1 = w, \text{ et } \forall n \in \mathbf{N}, w^{n+1} = w.w^n.$$

Définition 2.1.6 (Facteurs, préfixes et suffixes)

Soit w un mot sur \mathcal{A}^* .

- S'il existe des α, β et u tels que $w = \alpha u \beta$, on dit que u est un facteur de w .
- Si de plus, $\alpha = \varepsilon$ on dit que u est un préfixe de w .
- A l'inverse si $\beta = \varepsilon$ on dit que u est un suffixe de w .

Exemple : Si $w = abababac$ alors ab est un préfixe de w , ba est un facteur de w et c un suffixe de w .

Proposition 2.1.7

Le monoïde \mathcal{A}^* est régulier à droite et à gauche. Cela signifie que si u, v et w appartiennent à \mathcal{A}^* alors

$$uv = uw \Rightarrow v = w \text{ et } vu = wu \Rightarrow v = w$$

2.2 Langages**Définition 2.2.8**

Soit \mathcal{A} un alphabet fini. On appelle langage sur \mathcal{A} une partie de \mathcal{A}^* . L'ensemble des langages est noté $\mathcal{L}(\mathcal{A})$.

Remarque : On a donc $\mathcal{L}(A) = \mathcal{P}(\mathcal{A}^*)$.

Exemples :

1. L'ensemble vide \emptyset est un langage
2. Le singleton $\{\varepsilon\}$ est un langage. Il ne faut pas le confondre avec le précédent.
3. Chaque singleton $\{a\}$ où a est un élément de \mathcal{A} ou plus généralement $\{w\}$ où w est un élément de \mathcal{A}^* est un langage.
4. On prend $\mathcal{A} = \{a, b\}$. On peut considérer le langage des mots qui commencent par a et finissent par b .
5. Sur un alphabet \mathcal{A} on peut considérer l'ensemble de mots de longueur impaire.

Remarque : La notion de langage va être ce qui va correspondre à nos motifs. On peut se poser la question de chercher dans un texte un mot appartenant à un langage donné ou inversement en se donnant un mot savoir s'il appartient ou non à un langage. On peut par exemple considérer les langages suivants

- Pour \mathcal{A} qui désigne l'ensemble des signes ASCII. On peut considérer le langage des verbes conjugués à l'imparfait du subjonctif.
- Pour $\mathcal{A} = \{0, 1\}$ on peut considérer le langage des mots qui correspondent à un nombre (en binaire) divisible par 3.

Définition 2.2.9 (Opérations sur les langages)

On définit sur $\mathcal{L}(A)$ deux opérations :

- L'union noté \cup : L'union des langages L_1 et L_2 est le langage $L_1 \cup L_2$. On le note aussi $L_1 + L_2$.
- La concaténation : La concaténation L_1 et L_2 notée L_1L_2 (ou $L_1.L_2$) est l'ensemble des mots obtenus comme concaténation d'un mot de L_1 et d'un mot de L_2 :

$$L_1L_2 = \{w_1w_2 \in \mathcal{A}^* \mid w_1 \in L_1, w_2 \in L_2\}$$

Remarques :

1. La concaténation s'appelle aussi produit.
2. On peut aussi définir l'intersection et la différence de deux langages comme étant les opérations ensemblistes classiques mais c'est moins utile.
3. Il arrive que l'on note juste a le langage $\{a\}$. Par exemple on notera $a\mathcal{A}^*$ pour la concaténation du langage $\{a\}$ et du langage \mathcal{A}^* qui est l'ensemble des mots qui commencent par a .

Proposition 2.2.10 (Règles de calculs)

Soit L, A et B trois langages,

1. $L.(A + B) = L.A + L.B$ et $(A + B).L = A.L + B.L$
2. $L + \emptyset =$
3. $L.\emptyset =$
4. $L.\{\varepsilon\} = \{\varepsilon\}.L =$

Démonstration :

Exemple : Si on considère $L_1 = \{a\}^*$ les mots qui ne s'écrivent qu'avec a et $L_2 = \{b\}^*$ les mots qui ne s'écrivent qu'avec b . Le langage $L_1.L_2$ est l'ensemble des mots de la forme : $a \cdots ab \cdots b$.

Définition 2.2.11 (Puissance d'un langage)

Soit L un langage, on définit ses puissances par

$$L^0 = \{\varepsilon\} ; L^1 = L, \text{ et } \forall n \in \mathbf{N}, L^{n+1} = L.L^n.$$

Exercice : Soit L un langage. Comparer L^2 et $\{w^2 \mid w \in L\}$.

Définition 2.2.12 (Opération étoile)

Soit L un langage on note L^* le langage des mots obtenus en concaténant un nombre fini (éventuellement nul) de mots de L . On a

$$L^* = \bigcup_{k \in \mathbf{N}} L^k.$$

En particulier, on a toujours $\varepsilon \in L^*$.

Remarque : Cette opération s'appelle aussi étoile de Kleene.

Exemple : Pour $\mathcal{A} = \{a, b, c\}$. Si on pose $L_1 = \{a\}$, $L_2 = \{ab, ac\}$ et $L = L_1.(L_2)^*$. Les mots suivants appartiennent-ils à L ?

— ab

- a
- $aababac$
- $aabaac$

3 Langages rationnels et expressions régulières

3.1 Langages rationnels

Définition 3.1.13 (Langages rationnels)

Soit \mathcal{A} un alphabet fini.

On note $R(\mathcal{A})$ le plus petit sous-ensemble de $\mathcal{L}(\mathcal{A})$ stable par les opérations suivantes :

- union de deux langages (et donc d'un nombre fini)
- concaténation de deux langages (et donc d'un nombre fini)
- étoile

et contenant

- le langage vide,
- le langage $\{\varepsilon\}$,
- les langages $\{a\}$ pour $a \in \mathcal{A}$.

Les langages de $R(\mathcal{A})$ s'appellent les langages rationnels.

Remarque : Cela signifie que les langages rationnels sont les langages que l'on obtient à partir des « atomes », \emptyset , $\{\varepsilon\}$ et $\{a\}$ pour $a \in \mathcal{A}$ en appliquant les trois opérations.

Exemples :

ATTENTION

1. Si L est un langage rationnel. Un sous-langage de L n'est pas nécessairement rationnel.
2. Un langage rationnel n'est pas nécessairement stable par $.$ et \star .

3.2 Expressions régulières

Si on se donne un langage L , il n'est pas aisé de trouver un mot du langage dans un texte ou savoir si un mot donné appartient (ou pas) au langage.

On va donc représenter les langages par des formules.

Commençons par la définition formelle de ces formules. Elle sont définies par induction comme les formules logiques.

Définition 3.2.14 (Expressions régulières)

Soit \mathcal{A} un alphabet fini. On définit récursivement l'ensemble $\mathcal{E}(\mathcal{A})$ des expressions régulières par :

- \emptyset est une expression régulière
- ε et toute lettre de l'alphabet sont des expressions régulières que l'on dira atomiques
- si e_1 et e_2 sont des expressions régulières la somme $(e_1|e_2)$ est une expression régulière
- si e_1 et e_2 sont des expressions régulières le produit $(e_1.e_2)$ est une expression régulière
- si e est une expression régulière $e\star$ est une expression régulière.

Remarques :

1. On les appellera aussi expressions rationnelles ou motifs.
2. On peut voir $\mathcal{E}(\mathcal{A})$ comme la clôture inductive de l'ensemble $B = \mathcal{A} \cup \{\varepsilon, \emptyset\}$ par les constructeurs $|$ et $.$ d'arité 2 et le constructeur \star d'arité 1.
3. Pour minimiser les parenthèses, on considère que \star est prioritaire sur le produit, lui même prioritaire sur la somme. Par exemple l'expression :

$$((a|(b\star)).b)|a \text{ se note } (a|b\star).b|a$$

4. Il arrive que l'on note $+$ pour $|$.

On va associer à toute expression régulière un langage.

Définition 3.2.15 (Sémantique des expressions régulières)

On définit par induction structurelle une application L de $\mathcal{E}(\mathcal{A})$ dans $\mathcal{L}(\mathcal{A})$ qui associe à chaque expression régulière un langage.

- $L(\emptyset)$
- $L(\varepsilon) = \{\varepsilon\}$ et pour tout a de \mathcal{A} , $L(a) = \{a\}$
- $\forall (e_1, e_2) \in \mathcal{E}(\mathcal{A})^2$, $L(e_1|e_2) = L(e_1) \cup L(e_2)$
- $\forall (e_1, e_2) \in \mathcal{E}(\mathcal{A})^2$, $L(e_1.e_2) = L(e_1) \cdot L(e_2)$
- $\forall e \in \mathcal{E}(\mathcal{A})$, $L(e\star) = L(e)^*$

Exemples :

1. Si $\mathcal{A} = \{a, b\}$, $L(a.b\star)$ est le langage $\{a\}.\{b\}^*$ des mots de la forme ab^k pour $k \in \mathbf{N}$ (éventuellement $k = 0$)
2. Si $\mathcal{A} = \{a, b\}$, $L((a|b)\star)$ est \mathcal{A}^* .

Exercice :

- Décrire les mots du langage $L(ab\star a\star)$.

— Définir une expression régulière décrivant le langage $\{a^n b^p \mid n \in \mathbf{N}, p \in \mathbf{N}^*\}$.

Notation : Il existe de nombreuses versions d'expressions régulières avec des syntaxes plus enrichies. Malheureusement il n'existe pas de normalisation universelle. Voici quelques exemples les plus courants :

- Le symbole $.$ désigne toute lettre de l'alphabet. C'est l'expression qui est la somme de toutes les constantes atomiques sauf le mot vide. Par exemple pour $\mathcal{A} = \{a, b, c\}$, $.$ désigne $(a|b|c)$.
- Si e est une expression, le symbole e^+ désigne la répétition au moins une fois d'un motif de e . C'est donc ee^* .
- On note pour finir $e^?$ pour $(e|\varepsilon)$ à savoir la répétition au plus une fois d'un motif.

Définition 3.2.16

Soit $L \subset \mathcal{P}(\mathcal{A}^*)$ un langage. On dit que L est un langage régulier s'il existe une expression régulière e telle que $L = L(e)$.

Proposition 3.2.17

Soit $L \subset \mathcal{P}(\mathcal{A}^*)$ un langage.
Il est régulier si et seulement s'il est rationnel.

Démonstration :

□

Exemple : On a $\mathcal{A} = \{a, b, c\}$. Déterminer une expression régulière e telle $L(e)$ soit le langage des mots dans lesquels c n'apparaît jamais juste à gauche d'un b .

ATTENTION

Comme dans le cas des formules logiques, il est important de faire la différence entre la syntaxe d'une expression régulière (son écriture) et la sémantique (le langage qui lui est associé).

Deux expressions régulières différentes peuvent définir le même langage. Par exemple $(a|b)^*$ et $\varepsilon|(a|b).(a|b)^*$

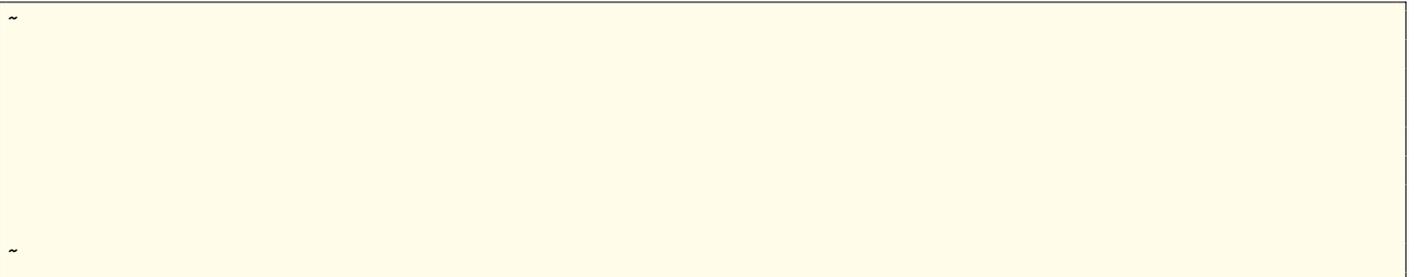
Définition 3.2.18

Deux expressions régulières définissant le même langage sont dites sémantiquement équivalentes.

Notation : Si e_1 et e_2 sont des expressions régulières équivalentes (c'est-à-dire $L(e_1) = L(e_2)$) on note $e_1 \equiv e_2$.

3.3 Implémentation des expressions régulières

On peut utiliser l'implémentation



Dans l'implémentation ci-dessus :

- On utilise `Mot of string` plutôt que `Lettre of char` ce qui permet d'utiliser `Mot("ab")` plutôt que `Produit(Lettre('a'),Lettre('b'))`.
- On utilise des listes pour les produits et les sommes en utilisant l'associativité. La encore on peut utiliser `Somme([Mot("a"); Mot("b") ; Mot("c")])` plutôt que `Somme(Somme(Mot("a"),Mot("b")),Mot("c"))`.