

# Chapitre 14 : Langages formels

## Table des matières

<b>1</b>	<b>Langages réguliers</b>	<b>2</b>
1.1	Motivation . . . . .	2
1.1.1	introduction . . . . .	2
1.1.2	Exemple . . . . .	2
1.2	Langages . . . . .	3
1.2.1	Définition ( <i>alphabet</i> ) . . . . .	3
1.2.2	Définition ( <i>mot</i> ) . . . . .	3
1.2.3	Définition ( <i>concaténation</i> ) . . . . .	3
1.2.4	Définition ( <i>préfixe, suffixe, facteur, sous-mot</i> ) . . . . .	4
1.2.5	Langages . . . . .	4
1.3	Langages réguliers . . . . .	5
1.3.1	Opérations sur les langages . . . . .	5
1.3.2	Théorème ( <i>Lemme d'ARDEN</i> ) (H.P) . . . . .	5
1.3.3	Langage régulier . . . . .	7
1.3.4	Expressions régulières . . . . .	8

# 1 Langages réguliers

## 1.1 Motivation

### 1.1.1 introduction

On a souvent besoin de mettre en place une analyse de texte, même dans le cadre d'applications qui ne relèvent pas uniquement du traitement de texte.

Par exemple :

- La recherche d'un mot dans un texte (*cf* chap 11) ;
- analyser un document structuré afin de traiter de manière appropriée son contenu (exemple : compiler un programme, récupérer des données sérialisées (*cf* chap 11) dans un format particulier (ex : données brutes en CSV, fichiers de configuration en JSON ou en XML)) ;
- Reconnaître un encodage et le déchiffrer (exemple : QR-code).

Quelle que soit l'application, on a besoin d'un formalisme pour décrire la structure du texte et d'algorithmes efficaces capables d'analyser cette structure et d'extraire les données associées.

### 1.1.2 Exemple

Étant donné un fichier binaire, déterminer s'il contient la représentation binaire d'un entier non signé multiple de 3.

- Remarque : on n'utilise pas les types d'entiers natifs de C ou OCaml car ils ont une taille fixée qui peut être dépassée par le fichier.

Idée : on lit les bits un à un en effectuant les opérations associées modulo 3, en remarquant que

$$\begin{cases} \langle x0 \rangle_2 = 2 \langle x \rangle_2 \\ \langle x1 \rangle_2 = 2 \langle x \rangle_2 + 1 \end{cases}$$

On utilise cette table :

$\langle x \rangle_2 \bmod 3$	0	1	2
$\langle x0 \rangle_2 \bmod 3$	0	2	1
$\langle x1 \rangle_2 \bmod 3$	1	0	2

- Algorithme :



**Algorithm 1:**


---

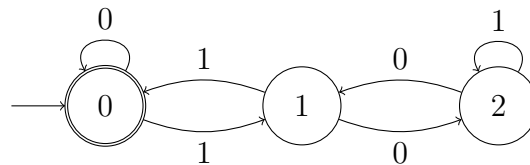
```

1  $x \leftarrow 0$ ;
2 for chaque bit  $b$  pris dans l'ordre do
3   if  $b = 0$  then
4      $x \leftarrow 2x \bmod 3$ ;
5   else
6      $x \leftarrow 2x + 1 \bmod 3$ 
7 return  $x = 0$ ;

```

---

• Représentation graphique :  $x$  ne peut prendre que trois valeurs différentes, appelées états, et on peut représenter les changements de valeur de  $x$  dans un graphe orienté dont les sommets sont les états, et les arcs sont étiquetés par le bit qui produit le changement de valeur de l'état source vers l'état cible.



On distingue de plus la valeur initiale par une flèche, et la valeur finale atteinte par un double cercle.

Cette représentation correspond au formalisme des *automates*, que nous verrons en ?? (page ??).

## 1.2 Langages

### 1.2.1 Définition (*alphabet*)

Un *alphabet* est un ensemble fini non vide, dont les éléments sont appelés lettres ou symboles.

Notation usuelle :  $\Sigma$ .

### 1.2.2 Définition (*mot*)

Soit  $\Sigma$  un alphabet.

Un *mot* sur  $\Sigma$  est une suite finie de symboles  $u = u_1 \cdots u_n$ , potentiellement vide.

Si  $n = 0$ , on note  $u = \varepsilon$ .

On note  $|u| = n$  la *longueur* du mot.

### 1.2.3 Définition (*concaténation*)

Soit  $\Sigma$  un alphabet, et  $u, v$  deux mots sur  $\Sigma$ .

On appelle *concaténation* de  $u$  et  $v$  le mot

$$uv = \begin{cases} v & \text{si } u = \varepsilon \\ u & \text{si } v = \varepsilon \\ u_1 \cdots u_n v_1 \cdots v_p & \text{si } \begin{cases} u = u_1 \cdots u_n \\ v = v_1 \cdots v_p \end{cases} \end{cases}$$

Exo

- $|uv| = |u| + |v|$
- La concaténation est une loi de composition interne associative et d'élément neutre  $\varepsilon$  sur l'ensemble des mots sur  $\Sigma$ .

### 1.2.4 Définition (*préfixe, suffixe, facteur, sous-mot*)

Soit  $\Sigma$  un alphabet, et  $u, v$  deux mots sur  $\Sigma$

- $v$  est un *préfixe* de  $u$  ssi  $\exists w$  mot tel que  $u = vw$
- $v$  est un *suffixe* de  $u$  ssi  $\exists w$  mot tel que  $u = wv$
- $v$  est un *facteur* de  $u$  ssi  $\exists x, y$  mots tels que  $u = xvy$
- $v$  est un *sous-mot* de  $u$  ssi  $\exists i_1 < i_2 < \cdots < i_k$  tels que si  $u = u_1 \cdots u_n$ , alors  $v = u_{i_1} \cdots u_{i_k}$

Exemple : si  $u = abc$ ,  $v = ac$  est un sous-mot de  $u$ , mais pas un facteur.

### 1.2.5 Langages

- Définition (*langage*) : un *langage* sur un alphabet  $\Sigma$  est un ensemble de mots sur  $\Sigma$ .
- Exemples :
  - l'ensemble de tous les mots, noté  $\Sigma^*$  (cf 1.3, page 5);
  - l'ensemble des écritures binaires de multiples des 3 ( $\Sigma = \{0, 1\}$ );
  - $\Sigma$  (si on voit les lettres comme des mots de longueur 1);
  - l'ensemble des code sources OCaml de programmes qui ne terminent pas.
- Problème : étant donné un langage  $L$  sur un alphabet  $\Sigma$ , on veut disposer d'une représentation formelle de  $L$  pour pouvoir étudier la question suivante : étant donné un mot  $u$ , a-t-on  $u \in L$ ?

C'est une question importante car souvent, comme dans les exemples en 1.1.1 (page 2), il faut pouvoir vérifier la structure d'un élément avant d'en extraire des données.

Malheureusement, on ne peut pas toujours répondre algorithmiquement à cette question. (cf chap 16 et la notion de décidabilité et le problème de l'arrêt), mais on peut y répondre pour une classe restreinte de langages.



## 1.3 Langages réguliers

### 1.3.1 Opérations sur les langages

Outre les opérations ensemblistes usuelles (intersection, union, complémentaire), on définit certaines opérations plus spécifiques aux langages.

- **Concaténation** : Soit  $\Sigma$  un alphabet, et  $L, L'$  deux langages sur  $\Sigma$ .

La concaténation de  $L$  et  $L'$  est le langage

$$LL' = \{uv \mid (u, v) \in L \times L'\}$$

Remarque :  $L\emptyset = \emptyset L = \emptyset$ .

- **Puissance** : Soit  $\Sigma$  un alphabet,  $L$  un langage sur  $\Sigma$ , et  $n \in \mathbb{N}$ .

La puissance  $n$ -ème de  $L$  est le langage

$$L^n = \begin{cases} \{\varepsilon\} & \text{si } n = 0 \\ LL^{n-1} & \text{si } n > 0 \end{cases}$$

- **Étoile de KLEINE** : Soit  $\Sigma$  un alphabet, et  $L$  un langage sur  $\Sigma$ .

L'étoile de KLEINE de  $L$  est le langage

$$L^* = \bigcup_{n \in \mathbb{N}} L^n$$

- Remarques :

- $\Sigma^*$  est bien l'ensemble de tous les mots : tout mot  $u = u_1 \cdots u_n$  est la caractérisation de ses lettres ( $\forall i \in \llbracket 1 ; n \rrbracket$ ,  $u_i \in \Sigma$ , donc  $u = u_1 \cdots u_n \in \Sigma^n \subseteq \Sigma^*$ ).
- On peut aussi définir la puissance  $n$ -ème d'un mot :

$$u^n = \begin{cases} \varepsilon & \text{si } n = 0 \\ uu^{n-1} & \text{si } n > 0 \end{cases}$$

Attention : ne pas confondre  $L^n$  et  $\{u^n \mid u \in L\}$ .

- On note aussi

$$L^+ = \bigcup_{n \in \mathbb{N}^*} L^n$$

$$\boxed{\text{Exo}} \quad L^+ = L^* \Leftrightarrow \varepsilon \in L.$$

### 1.3.2 Théorème (Lemme d'ARDEN) (H.P)

| Soit  $\Sigma$  un alphabet, et  $K, L$  deux langages sur  $\Sigma$ .

(1)  $K^*L$  est le minimum (pour l'ordre de l'inclusion) des solutions de l'équation

$$X = KX \cup L$$

d'inconnue un langage  $X$ .

(2) Si  $\varepsilon \notin K$ , alors  $K^*L$  est l'unique solution.

□

(1) On a :

$$\begin{aligned} K(K^*L) \cup L &= K \left( \bigcup_{n \in \mathbb{N}} K^n L \right) \cup L \\ &= K^+ L \cup K^0 L \\ &= K^* L \end{aligned}$$

Donc  $K^*L$  est une solution.

Puis soit  $X$  une solution. Montrons que  $K^*L \subset X$ .

Soit  $u \in K^*L$ . Par définition,

$$\exists n \in \mathbb{N}, \exists k_1, \dots, k_n \in K, \exists l \in L \mid u = k_1 \cdots k_n l$$

On montre par récurrence sur  $n$  que  $u \in X$  :

- $n = 0$  :  $u \in L \subset KX \cup L = X$ , donc  $u \in X$ .

- Hérédité : si  $\forall k_1 \cdots k_n \in K, \forall l \in L, k_1 \cdots k_n l \in X$ , considérons  $u = k_1 \cdots k_{n+1} l$ , avec  $k_1 \cdots k_{n+1} \in K$  et  $l \in L$ .

$$u = k_1 \underbrace{(k_2 \cdots k_{n+1} l)}_{\in X \text{ par H.R.}} \in KX \subseteq KX \cup L = X, \text{ donc } u \in X.$$

Finalement,  $K^*L \subset X$ , et  $K^*L$  est bien le minimum des solutions.

(2) On suppose  $\varepsilon \notin K$ . Soit  $X$  une solution.

On sait par (1) que  $K^*L \subseteq X$ .

Il suffit de montrer que  $X \subseteq K^*L$ .

Soit  $u \in X$  dont on note  $n$  la longueur.

On montre par récurrence que

$$\forall k \in \llbracket 0 ; n \rrbracket, X = \bigcup_{j=0}^k K^j L \cup K^{k+1} X$$

- $k = 0$  :

$$\bigcup_{j=0}^0 K^j L \cup K^{0+1} X = L \cup KX = X$$

- Hérédité : Soit  $k \in \llbracket 0 ; n-1 \rrbracket \mid X = \bigcup_{j=0}^k K^j L \cup K^{k+1} X$

$$\begin{aligned}
X &= KX \cup L \\
&= K \left( \bigcup_{j=0}^k K^j L \cup K^{k+1} L \right) \cup L \\
&= \bigcup_{j=0}^k K^{j+1} L \cup K^{k+1} L \cup L \\
&= \bigcup_{j=0}^{k+1} K^j L \cup K^{k+2} L
\end{aligned}$$

En particulier,

$$X = \bigcup_{j=0}^n K^j L \cup K^{n+1} X$$

Or  $\forall v \in K^{n+1} X$ ,  $|v| \geq n+1$

En effet,  $\exists k_1 \cdots k_{n+1} \in K$ ,  $\exists x \in X \mid v = k_1 \cdots k_{n+1} x$ .

Or  $\forall i \in \llbracket 1 ; n+1 \rrbracket$ ,  $k_i \neq \varepsilon$ , donc  $|k_i| \geq 1$ , donc

$$|v| = \sum_{i=1}^{n+1} |k_i| + |x| \geq n+1 + |x| \geq n+1$$

Donc, comme  $|u| = n < n+1$ ,

$$\bigcup_{j=0}^n K^j L \subseteq \bigcup_{j \in \mathbb{N}} K^j L = K^* L$$

Donc  $u \in K^* L$ , et  $X \subseteq K^* L$  ■

- Contre-exemple au point (2) si  $\varepsilon \in K$  : on prend  $K = \{\varepsilon\}$ , et  $L = \{a\}$ .  $K^* L = \{a\} = L$ , et tout  $X \mid a \in X$  est solution ( $KX \cup L = X \cup \{a\}$ ). Par exemple :  $\{a\} = K^* L$  ou  $\{a, aa\}, \{a, aa, aaa\}$ .

### 1.3.3 Langage régulier

La classe des *langages réguliers*, aussi appelés *langages rationnels*, est la famille des langages que l'on peut construire à partir de langages de base ( $\emptyset, \{\varepsilon\}, \{a\} \forall a \in \Sigma$ ) et des opérations dites *régulières* : union, concaténation et étoile de KLEINE. L'ensemble  $\text{Reg}(\Sigma)$  des langages réguliers sur l'alphabet  $\Sigma$  est défini inductivement par :

$$\begin{array}{c}
\frac{}{\emptyset \in \text{Reg}(\Sigma)} \quad \frac{a \in \Sigma}{\{a\} \in \text{Reg}(\Sigma)} \quad \frac{L \in \text{Reg}(\Sigma)}{L^* \in \text{Reg}(\Sigma)} \\
\\
\frac{l \in \text{Reg}(\Sigma) \quad l' \in \text{Reg}(\Sigma)}{L \cup L' \in \text{Reg}(\Sigma)} \quad \frac{l \in \text{Reg}(\Sigma) \quad l' \in \text{Reg}(\Sigma)}{LL' \in \text{Reg}(\Sigma)}
\end{array}$$

- Remarque :  $\{\varepsilon\} \in \text{Reg}(\Sigma)$  car  $\{\varepsilon\} = \emptyset^0 = \bigcup_{n \in \mathbb{N}} \emptyset^n = \emptyset^*$

$\Sigma \in \text{Reg}(\Sigma)$  car, en notant  $\Sigma = \{a_1 \cdots a_n\}$ , on a :

$$\Sigma = \bigcup_{k=1}^n \{a_k\}$$

Donc une récurrence sur  $n = |\Sigma|$  conclut.

De même, tout langage fini est régulier (Exo).

Il manque encore un formalisme pour décrire les langages réguliers.

### 1.3.4 Expressions régulières

Soit  $\Sigma$  un alphabet, et  $S = \{ (, ), |, *, \emptyset, \varepsilon \}$  un ensemble de symboles supposé disjoint de  $\Sigma$ .

L'ensemble  $\text{Regexp}(\Sigma)$  des *expressions régulières* sur  $\Sigma$ , aussi appelées *expressions rationnelles*, est l'ensemble des mots sur  $\Sigma \cup S$  défini inductivement par

$$\begin{array}{c} \frac{}{\emptyset \in \text{Regexp}(\Sigma)} \quad \frac{}{\Sigma \in \text{Regexp}(\Sigma)} \quad \frac{a \in \Sigma}{a \in \text{Regexp}(\Sigma)} \\[10pt] \frac{e \in \text{Regexp}(\Sigma) \quad f \in \text{Regexp}(\Sigma)}{(e|f) \in \text{Regexp}(\Sigma)} \quad \frac{e \in \text{Regexp}(\Sigma) \quad f \in \text{Regexp}(\Sigma)}{(ef) \in \text{Regexp}(\Sigma)} \\[10pt] \frac{e \in \text{Regexp}(\Sigma)}{(e^*) \in \text{Regexp}(\Sigma)} \end{array}$$

Remarques :

- $|$  est parfois noté  $+$ .
- On se passe de certaines parenthèses avec les règles de priorité :

$$* > \text{concaténation} > |$$

Par exemple,  $((a(b^*))|b)$  s'écrit  $ab^*|b$ .