

Théorie des Langages

Matthew Coyle / Valentin Bayart

January 14, 2018

1 Introduction et concepts de base

- Introduction

Les **langages formels** ont été étudiés par :

- les **informaticiens** : langages de programmation (Avec une syntaxe (La définir, la vérifier) permettant de traduire ce même langage
- les **linguistes** : langues naturelles

Exemples de Langages :

- Entiers naturels (suite de chiffres parmi 0...9).
- Entiers naturels impairs (même représentation).
- Les mots français (du dictionnaire).
- Les identificateurs C++.
- Les phrases en français.
- Les programmes (syntaxiquement corrects) écrits en C++.

Tous ces exemples sont des ensembles ou des sous-ensembles d'un autre ensemble (lettres / alphabet / dictionnaire).

Points communs des langages formels :

- Chaque langage est un ensemble d'éléments appelés **mots** ou "**chaînes**".
- Chaque chaîne est une suite de **symboles** pris parmi un ensemble fini de symboles.
- Chaque chaîne est de longueur finie (même s'il n'y a pas de limite à cette longueur).

On étudie des modèles pour représenter de manière finie des langages :

- automates finis.
- expressions régulières.
- grammaires formelles.
- ...

Application pratiques:

- Recherche de "**motifs**" dans des fichiers.
- Traitement de texte.
- Modélisation de circuits.
- de machines à états.
- Compilation de langages de programmation.
- ...

- Concepts de Base

- Alphabets :

Un **alphabet** est un ensemble fini, non vide, de symboles.

On le note généralement Σ .

Exemples d'alphabets :

$$\Sigma_{entiers} = \{0,1,2,3,4,5,6,7,8,9\}$$

$$\Sigma_{mots} = \{a,b,c,\dots,z,',-\}$$

$$\Sigma_{ident} = \{a,\dots,z,A,\dots,Z,0,\dots,9,-\}$$

$$\Sigma_{prog} = \{int, float, bool, while, do, for, \dots, <, <=, >, >=, =, !=, +, -, /, *, :, \dots, 0, 1, \dots, 25, 26, 27, \dots, 12.56, \dots, a, b, toto, compteur, Tab, \dots\}$$

- Chaînes

Un **mot** ou une **chaîne** ω (omega) formé(e) sur un alphabet est une suite finie $s_1s_2\dots s_n$ de symboles de cet alphabet

La **chaîne vide**, noté ε (epsilon), est une chaîne ne contenant aucun symbole

La **longueur** d'une chaîne ω , notée $|\omega|$, est le nombre de symboles composant la chaîne ω

- **Opérations sur les chaînes** La concatenation de deux chaînes u et v , notée $u.v$ ou uv est la chaîne obtenue en écrivant les symboles de u suivis de ceux de v .

si $u = a_1a_2\dots a_n$ et $v = b_1b_2\dots b_p$

alors $uv = a_1a_2\dots a_nb_1b_2\dots b_p$

Propriétés :

$$- |u.v| = |u| + |v|$$

$$- \text{Associativité : } (u.v).w = u.(v.w)$$

$$- \varepsilon \text{ est l'élément neutre : } u.\varepsilon = \varepsilon.u = u$$

. **Puissance d'une chaîne** ω ω^k est la chaîne formée par la concaténation de k occurrences de ω

$$\omega^k = \underbrace{\omega\omega\omega\dots\omega}_k \text{ fois}$$

$$\omega^0 = \varepsilon$$

Un **préfixe** d'une chaîne ω ? est une suite, éventuellement vide, de symboles débutant ω .

Un **suffixe** de ω est une suite de symboles terminant ω .

$\forall x,y \mid \omega = x.y$, x est un préfixe de ω , y un suffixe.

Une **sous-chaîne** d'une chaîne ω est une suite de symboles apparaissant consécutivement dans ω .

Notation : $|\omega|_x$ est le nombre d'occurrences de la chaîne x dans la chaîne ω .

- Langage :

Un langage est un ensemble de chaînes.

Exemples de langages:

$$\{toto,titi,tata\}$$

$$\{1,11,101,1001\}$$

$\{1^n \mid n \geq 0\} = \{e, 1, 11, 111, 1111, 11111, \dots\}$ c'est un langage infini (nombre infini de chaînes) dont chaque chaîne est de longueur finie

Nombres binaires impaires: $\{1,11,101,111,1001,1011,\dots\}$

Nombres binaires premiers: $\{1,10,11,101,111,1011,\dots\}$

le **Langage vide**, noté \emptyset , ne contient aucune chaîne (ensemble vide).

Attention : $\emptyset \neq \{\varepsilon\}$

Le langage **plein**, noté Σ^* , contient toutes les chaînes que l'on peut former sur l'alphabet Σ
Remarque : $\Sigma^* = \Sigma^+ \cup \{\varepsilon\}$

- Opérations sur les langages :

l'**Union** de deux langages **A** et **B** est le langage, note $A \cup B$, composé de toutes les chaînes qui apparaissent dans l'un au moins de langages **A** ou **B**.

$$A \cup B = \{\omega \mid \omega \in A \text{ ou } \omega \in B\}$$

Propriétés :

- Commutativité : $A \cup B = B \cup A$
- Associativité : $(A \cup B) \cup C = A \cup (B \cup C)$
- \emptyset est élément neutre : $A \cup \emptyset = \emptyset \cup A = A$
- Idempotence : $A \cup A = A$

La **concaténation** de deux langages **A** et **B** est le langage, note $A.B$ ou AB , composé de toutes les chaînes formées par une chaîne de **A** concaténée à une chaîne de **B**.

$$A.B = \{u.v \mid u \in A \text{ ou } v \in B\}$$

Propriétés:

- Associativité: $(A.B).C = A.(B.C)$
- $\{\varepsilon\}$ est un élément neutre: $A.\{\varepsilon\} = \{\varepsilon\}.A = A$
- \emptyset est élément absorbant : $A.\emptyset = \emptyset.A = \emptyset$

Distributivité de la concaténation sur l'union:

- À gauche : $A.(B \cup C) = A.(B.C)$
- À droite : $(B \cup C).A = B.A \cup C.A$

. Puissance d'un langage A

A^k est le langage formé par la concaténation de **k** occurrences de **A**.

$$A^0 = \{\varepsilon\}$$

$$A^1 = A$$

$$A^n = \underbrace{AAA \dots AA}_{n \text{ fois}}$$

A^k : Mots formés par la concaténation de k mots de A

. Étoile de Kleene (fermeture ou cloture par .)

- la **fermeture de Kleene** d'un langage **A** est le langage, noté A^* , défini par : $A^* = A^0 \cup A^1 \cup A^2 \cup A^3 \dots$ soit

$$A^* = \bigcup_{i=0}^{\infty} A^i = A^0 \cup A^1 \cup A^2 \cup \dots$$

- la **fermeture positive** de **A** est le langage, noté A^+ , défini par : $A^+ = A^1 \cup A^2 \cup A^3 \cup \dots$ soit

$$A^+ = \bigcup_{i=1}^{\infty} A^i = A^1 \cup A^2 \cup A^3 \cup \dots$$

"mots formés par la concaténation de 1 ou plusieurs mots de A"

Remarque : $A^* = \{\varepsilon\} \cup A^+$ (

Propriétés : $A^+ = A.A^* = A^*.A$ (Possibilité : $A^+ = A^*$)

2 Modeles et Langages

- contextuels
- langages recursivement enumerables