

Introduction à la théorie des langages

Claire.Lefevre@univ-angers.fr

1

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

- les informaticiens
⇒ langages de programmation
(définir syntaxe, vérifier la syntaxe d'un programme, le traduire en langage machine)
- les linguistes
⇒ langues naturelles
(les décrire et essayer de les traiter automatiquement)

2

- Exemples de langages

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

- Points communs

- Chaque langage est un ensemble d'éléments (« **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)

3

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

- Automates finis
- Expressions régulières
- Grammaires formelles
- ...

- Applications pratiques

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

4

Langages Concepts de base

5

Alphabets

- Un **alphabet** est un ensemble **fini**, non vide, de symboles
- On le note généralement Σ (sigma)

Exemples

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

$\Sigma_{\text{mots}} = \{a, b, \dots, z, \acute{e}, \grave{e}, \text{ç}, \grave{a}, \grave{u}, \dots, ', -\}$

$\Sigma_{\text{idents}} = \{a, \dots, z, A, \dots, Z, 0, \dots, 9, _ \}$

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

6

Chaînes

- Un **mot** ou une **chaîne** ω formé(e) sur un alphabet est une suite finie $s_1 s_2 \dots s_n$ de symboles de cet alphabet
- La **chaîne vide**, notée 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

7

Opérations sur les chaînes

- La **concaténation** de 2 chaînes u et v , notée $u.v$ ou uv , est la chaîne obtenue en écrivant les symboles de u suivis de ceux de v

$$\begin{array}{lcl} \text{si} & u = a_1 a_2 \dots a_n & \text{et} \quad v = b_1 b_2 \dots b_p \\ \text{alors} & uv = a_1 a_2 \dots a_n b_1 b_2 \dots b_p \end{array}$$

Propriétés :

- $|u.v| = |u| + |v|$
- Associativité : $(u.v).w = u.(v.w)$
- ϵ est élément neutre : $u.\epsilon = \epsilon.u = u$

- Puissances d'une chaîne ω**

- ω^k est la chaîne formée par la concaténation de k occurrences de ω

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

- $\omega^0 = \epsilon$

8

- 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

$$\forall x, y \quad \text{t.q.} \quad \omega = x.y \quad \begin{array}{l} x \text{ est un préfixe de } \omega \\ \text{et } y \text{ est un suffixe de } \omega \end{array}$$

- 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

9

Langages

- Un **langage** est un ensemble de chaînes

Exemples

- $\{\text{toto, titi, tata}\}$
- $\{1, 11, 101, 1001\}$
- $\{1^n \mid n \geq 0\} = \{\epsilon, 1, 11, 111, 1111, 11111, \dots\}$
c'est un langage infini (nombre infini de chaînes)
dont chaque chaîne est de longueur finie
- Nombres binaires impairs
 $\{1, 11, 101, 111, 1001, 1011, \dots\}$
- Nombres binaires qui sont premiers
 $\{1, 10, 11, 101, 111, 1011, \dots\}$

10

- Le **langage vide**, noté \emptyset , ne contient aucune chaîne
- Attention : $\emptyset \neq \{\epsilon\}$
- Le langage « plein », noté Σ^* , contient toutes les chaînes que l'on peut former sur l'alphabet Σ
- Σ^+ contient toutes les chaînes *non vides* sur Σ

$$\text{Rem : } \Sigma^* = \Sigma^+ \cup \{\epsilon\}$$

11

Opérations sur les langages

- L'**union** de 2 langages A et B est le langage, noté $A \cup B$, composé de toutes les chaînes qui apparaissent dans l'un au moins des 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 : $A \cup A = A$

12

- La **concaténation** de 2 langages **A** et **B** est le langage, noté **A.B** ou **AB**, composé de toutes les chaînes formées par une chaîne de **A** concaténée à une chaîne de **B** :

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

Propriétés :

- Associativité : $(A.B).C = A.(B.C)$
- $\{\epsilon\}$ est élément neutre : $A.\{\epsilon\} = \{\epsilon\}.A = A$
- \emptyset est élément absorbant : $A.\emptyset = \emptyset.A = \emptyset$
- Distributivité de la concaténation sur l'union :
 - À gauche : $A.(B \cup C) = A.B \cup A.C$
 - À droite : $(B \cup C).A = B.A \cup C.A$

13

- Puissances d'un langage A :**

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

- $A^0 = \{\epsilon\}$
- $A^1 = A$
- $A^n = \underbrace{A A A \dots A A}_{n \text{ fois}}$

A^k : « mots formés par la concaténation de k mots de A »

- Étoile de Kleene** (fermeture ou clôture par .)

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

« mots formés par la concaténation d'un nbre qq de mots de A »

La **fermeture positive** de A est le langage, noté A^+ , défini par :

$$A^+ = A^1 \cup A^2 \cup A^3 \cup \dots$$

« mots formés par la concaténation de 1 ou plusieurs mots de A »

Rem : $A^* = \{\epsilon\} \cup A^+$

14

Exercice : calculer A^+ et A^* pour les langages suivants

- $A = \{a, ab\}$
- $A = \{ab\}$
- $A = \{\epsilon, ab\}$
- $A = \{a^n \mid n \in \mathbb{N}\}$

Propriété : $A^+ = A.A^* = A^*.A$

Rem : il est possible que $A^+ = A^*$

15

- L'**intersection** de 2 langages A et B est le langage, noté $A \cap B$, composé des chaînes apparaissant à la fois dans A et dans B :

$$A \cap B = \{\omega \mid \omega \in A \text{ et } \omega \in B\}$$

- La **différence** de 2 langages A et B est le langage, noté $A \setminus B$ ou $A - B$, composé des chaînes de A n'apparaissant pas dans B :

$$A \setminus B = \{\omega \mid \omega \in A \text{ et } \omega \notin B\}$$

- Le **complémentaire** d'un langage A sur un alphabet Σ est le langage noté \overline{A} composé de toutes les chaînes de Σ^* n'apparaissant pas dans A :

$$\overline{A} = \Sigma^* \setminus A$$

16

Modèles et langages

- On va étudier 3 modèles pour représenter des langages
 - Automates finis** : « machines » qui permettent de déterminer si une chaîne donnée appartient oui ou non à un langage
=> **reconnaissance** d'un langage
 - Expressions régulières** : notation qui permet de définir exactement quelles chaînes constituent un langage
=> **spécification** d'un langage
 - Grammaires formelles** : notation récursive pour spécifier un lang.
- On classe les langages selon le type de modèles permettant de les représenter :
 - Les **langages réguliers** (modélisés par les automates finis et les expressions régulières)
 - Les **langages non contextuels** (modélisés par les grammaires non contextuelles)
 - ...

17