

Théorie des langages et automates

Expressions régulières

Clément Moreau^{1,2}

¹BRED – Banque Populaire — clement.moreau@bred.fr

²Université de Tours

- 1 Algèbre de Kleene
- 2 Conversion AFN \rightarrow ER
- 3 Conversion ER \rightarrow AFN
- 4 Pattern matching
- 5 Lemme de l'étoile

Table of Contents

- 1 Algèbre de Kleene
- 2 Conversion AFN \rightarrow ER
- 3 Conversion ER \rightarrow AFN
- 4 Pattern matching
- 5 Lemme de l'étoile

- On rappelle qu'un langage L sur un alphabet Σ est un ensemble de mots $x \in \Sigma^*$.
- Un langage est dit *régulier* si il est reconnu par un automate fini.
- Une expression régulière E est une description algébrique d'un langage. On note celui-ci $L(E)$.

Définition

Une algèbre de Kleene est une structure algébrique : $\mathcal{K} = (\Sigma, +, \cdot, *, \emptyset, \varepsilon)$ dotée des axiomes suivants :

$$\text{A1} \quad a + \emptyset = a$$

$$\text{A2} \quad a + b = b + a$$

$$\text{A3} \quad (a + b) + c = a + (b + c)$$

$$\text{A4} \quad a + a = a$$

$$\text{A5} \quad a\emptyset = \emptyset a = \emptyset$$

$$\text{A6} \quad a\varepsilon = \varepsilon a = a$$

$$\text{A7} \quad a(b + c) = ab + ac$$

$$\text{A8} \quad (a + b)c = ac + bc$$

$$\text{A9} \quad a(bc) = (ab)c$$

$$\text{A10} \quad (a + b)^* = (a^*b)^*a^*$$

$$\text{A11} \quad (ab)^* = \varepsilon + a(ba)^*b$$

$$\text{A12} \quad (a^*)^* = a^*$$

où $a, b, c \in \Sigma$.

Expression régulière

On appelle *expression régulière* (ER) toute expression E sur \mathcal{K} dont la formule s'écrit au moyen d'opérateurs de l'algèbre de Kleene et obéissant aux règles inductives suivantes :

- Tout symbole $a \in \Sigma$ est une expression régulière,
- Si $E \in \mathcal{K}$, alors $(E)^* \in \mathcal{K}$
- Si $E, F \in \mathcal{K}$, alors $(E + F) \in \mathcal{K}$
- Si $E, F \in \mathcal{K}$, alors $(EF) \in \mathcal{K}$

Exemples

- $E_1 = (0 + 1)^* 01(0 + 1)^* \Leftrightarrow L(E_1) = \{x01y \mid x, y \in \{0, 1\}^*\}$
- $E_2 = ('-' + \varepsilon)([1, 9][0, 9]^* + 0)(\varepsilon + ('' [0, 9]^* [1, 9])) \Leftrightarrow L(E_2) = \mathbb{D}$

Deux expressions régulières E, F sont équivalentes si et seulement si :

$$L(E) \subseteq L(F) \wedge L(F) \subseteq L(E)$$

On raisonne par double inclusion, ou par réduction à l'aide des axiomes vus précédemment.

Exemples

Montrer que les expressions régulières suivantes sont équivalentes :

① $a^* = (aa)^* + a(aa)^*$

② $a(ba)^* = (ab)^*a$

Théorème

Pour toute expression régulière $E \in \mathcal{K}$, $L(E) \in \mathbf{Reg}$.

Corollaire

- Toute expression régulière peut-être convertie en automate fini.
 - Par l'algorithme de Thomson
 - Par dérivée partielle d'Antimirov
- Tout automate fini peut-être converti en ER.
 - Par méthode d'élimination d'états de Brzozowski et McCluskey
 - Par les équations d'Arden

Table of Contents

- 1 Algèbre de Kleene
- 2 Conversion AFN \rightarrow ER**
- 3 Conversion ER \rightarrow AFN
- 4 Pattern matching
- 5 Lemme de l'étoile

Équation d'Arden

Lemme d'Arden

Soient A, B deux langages sur un alphabet Σ . Soit l'équation :

$$X = AX + B$$

Le langage $L = A^*B$ est le plus petit langage qui est solution.

Équation d'Arden

Lemme d'Arden

Soient A, B deux langages sur un alphabet Σ . Soit l'équation :

$$X = AX + B$$

Le langage $L = A^*B$ est le plus petit langage qui est solution.

Démonstration.

- ① L est solution. En effet, on a :

$$\begin{aligned} AL + B &= A(A^*B) + B \\ &= (AA^* + \varepsilon)B \\ &= A^*B = L. \end{aligned}$$

- ② Le langage A^*B est la plus petite solution.

En développant l'équation et par récurrence, on obtient $X = A^{N+1}X + \sum_{n=0}^N A^n B$.

Dès lors, $X = \lim_{N \rightarrow \infty} A^{N+1} + \sum_{n=0}^N A^n B = A^*B$.



Conversion AFN \rightarrow ER

Soit $N = (Q, \Sigma, \Delta, S, F)$ un automate fini.

Pour $q \in Q$, X_q est le langage des mots acceptés avec q comme état de départ.

$$X_q = \begin{cases} \sum_{(q,a,p) \in \delta} aX_p + \varepsilon & \text{si } q \in F \\ \sum_{(q,a,p) \in \delta} aX_p & \text{sinon} \end{cases}$$

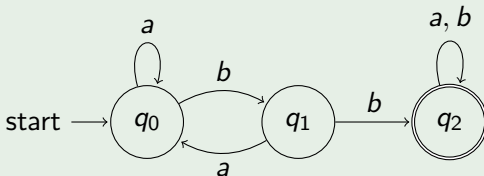
Où $(q, a, p) \in \delta \Leftrightarrow \delta(q, a) = p$.

On a alors :

$$L(N) = \sum_{s \in S} X_s$$

Exemples

Soit l'automate N suivant :



On cherche à résoudre le système d'équations d'Arden (S) suivant :

$$(S) \Leftrightarrow \begin{cases} L_0 = aL_0 + bL_1 & (1) \\ L_1 = aL_0 + bL_2 & (2) \\ L_2 = aL_2 + bL_2 + \varepsilon & (3) \end{cases}$$

Exemples

D'après le lemme d'Arden, l'équation (1) $\Leftrightarrow L_0 = a^*L_1$. Il vient que :

$$(S) \Leftrightarrow \begin{cases} L_0 = a^*bL_1 & (1) \\ L_1 = a^+bL_1 + bL_2 & (2) \\ L_2 = (a + b)^* & (3) \end{cases}$$

Exemples

D'après le lemme d'Arden, l'équation (1) $\Leftrightarrow L_0 = a^*L_1$. Il vient que :

$$(S) \Leftrightarrow \begin{cases} L_0 = a^*bL_1 & (1) \\ L_1 = a^+bL_1 + bL_2 & (2) \\ L_2 = (a+b)^* & (3) \end{cases}$$

$$\Leftrightarrow \begin{cases} L_0 = a^*bL_1 & (1) \\ L_1 = (a^+b)^*bL_2 & (2) \\ L_2 = (a+b)^* & (3) \end{cases}$$

Exemples

D'après le lemme d'Arden, l'équation (1) $\Leftrightarrow L_0 = a^*L_1$. Il vient que :

$$(S) \Leftrightarrow \begin{cases} L_0 = a^*bL_1 & (1) \\ L_1 = a^+bL_1 + bL_2 & (2) \\ L_2 = (a+b)^* & (3) \end{cases}$$

$$\Leftrightarrow \begin{cases} L_0 = a^*bL_1 & (1) \\ L_1 = (a^+b)^*bL_2 & (2) \\ L_2 = (a+b)^* & (3) \end{cases}$$

$$\Leftrightarrow \begin{cases} L_0 = a^*b(a^+b)^*b(a+b)^* & (1) \\ L_1 = (a^+b)^*b(a+b)^* & (2) \\ L_2 = (a+b)^* & (3) \end{cases}$$

On a finalement

$$L(N) = a^*b(a^+b)^*b(a+b)^*$$

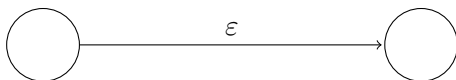
Table of Contents

- 1 Algèbre de Kleene
- 2 Conversion AFN \rightarrow ER
- 3 Conversion ER \rightarrow AFN**
- 4 Pattern matching
- 5 Lemme de l'étoile

Algorithme de Thomson

L'algorithme de Thomson construit un automate équivalent à une expression régulière E par induction :

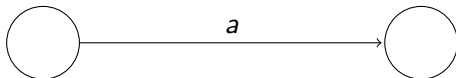
- Si $E = \varepsilon$



- Si $E = \emptyset$

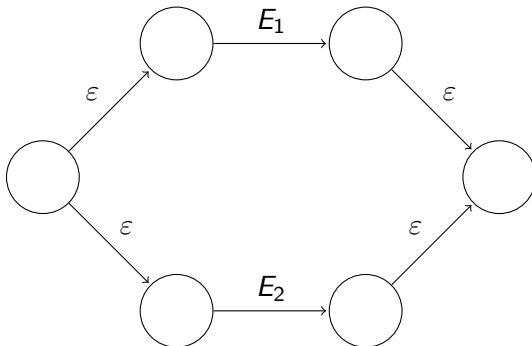


- Si $E = a \in \Sigma$

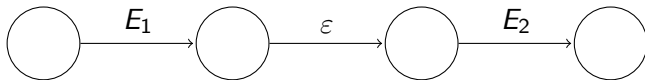


Algorithme de Thomson

- Si $E = E_1 + E_2$

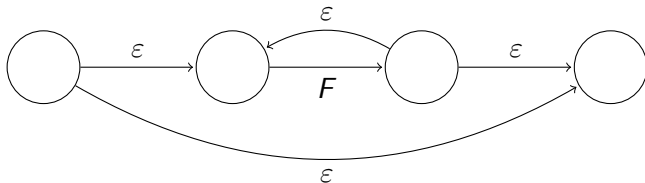


- Si $E = E_1 E_2$



Algorithme de Thomson

- Si $E = (F^*)$



Exemples

- 1 Appliquer l'algorithme de Thomson sur l'expression régulière :

$$E = (0 + 1)^*1(0 + 1)$$

- 2 Supprimer les ϵ -transitions inutiles.

Table of Contents

- 1 Algèbre de Kleene
- 2 Conversion AFN \rightarrow ER
- 3 Conversion ER \rightarrow AFN
- 4 Pattern matching**
- 5 Lemme de l'étoile

- Les expressions régulières ont un intérêt pratique pour le contrôle de format de documents et la *recherche de motif* dans un texte.
- Elles sont couramment utilisées au sein des différents langages de programmation (Java, C, Python, etc.) et des systèmes de type UNIX.

Regex syntaxe

Symbole	Description
.	N'importe quel symbole
[...]	Ensemble de caractères (Ex : [0-9a-z], [0-Z])
[^...]	Ensemble complémentaire
^	Début du mot
\$	Fin de mot
	Ou
(...)	Groupe de capture
*	0, 1 ou plusieurs occurrences
+	1 ou plusieurs occurrences
?	0 ou 1 occurrence
{...}	Comptage (Ex. $a\{3\} = aaa$, $a\{1,5\} = \{a, aa, \dots, aaaaa\}$)

Attention

Cette symbologie peut varier selon les langages. Elle est valable ici pour Java.

Regex syntaxe

Symbole	Description
<code>\d</code>	Ensemble des chiffres
<code>\s</code>	Ensemble des espaces
<code>\t</code>	Tabulation
<code>\w</code>	Ensemble alphanumérique (chiffres, lettres et <code>_</code>)
<code>\n</code>	Saut de ligne
<code>\p{L}</code>	Ensemble des lettres
<code>[[:alpha :]]</code>	Ensemble des lettres
<code>[[:digit :]]</code>	Ensembles des chiffres
<code>[[:alnum :]]</code>	Ensemble des lettres et chiffres et <code>_</code>
<code>[[:space :]]</code>	Ensemble des espaces
<code>[[:punct :]]</code>	Ensembles des symboles de ponctuation
<code>[[:lower :]]</code>	Ensemble des lettres minuscules
<code>[[:upper :]]</code>	Ensemble des lettres majuscules

Exemples

Donner le sens des regex ou l'expression correspondante :

- `\d{1,3}\.\d{1,3}\.\d{1,3}\.\d{1,3}`

Exemples

Donner le sens des regex ou l'expression correspondante :

- `\d{1,3}\.\d{1,3}\.\d{1,3}\.\d{1,3}`
- `<([a-z]+)[^>]>`

Exemples

Donner le sens des regex ou l'expression correspondante :

- `\d{1,3}\.\d{1,3}\.\d{1,3}\.\d{1,3}`
- `<([a-z]+)[^>]>`
- Un numéro de téléphone français

Exemples

Donner le sens des regex ou l'expression correspondante :

- `\d{1,3}\.\d{1,3}\.\d{1,3}\.\d{1,3}`
- `<([a-z]+)[^>]>`
- Un numéro de téléphone français
- Une adresse mail de l'université de Tours.

Pattern matching et expressions régulières

Il existe plusieurs méthodes en Java selon l'utilité qu'on a des regex :

- Vérifier le bon format d'un texte.
- Chercher un fragment obéissant à un pattern donné.
- Remplacer ou extraire de l'information.

```
import java.util.regex.Pattern;

public class Regex {
    public static void main(String[] args) {
        String s = ...; // Some string
        String regex = ...; // Some regex
        Pattern.matches(regex, s); // s match the regex
        Pattern.compile(regex).matcher(s).find() // s contains a substring
            corresponding regex
        s.replaceAll(regex, _s) // replace all substrings maching regex
            by _s in s
    }
}
```

Pattern matching et expressions régulières

*Sous le pont Mirabeau coule la Seine
Et nos amours
Faut-il qu'il m'en souviennne
La joie venait toujours après la peine
[...]*

*L'amour s'en va comme cette eau
courante
L'Amour s'en va
Comme la vie est lente
Et comme l'Espérance est violente*

Exemples

Déterminer les fragments de code permettant de :

- 1 Trouver les mots qui commencent par une majuscule mais qui ne sont pas des débuts de phrase.

Pattern matching et expressions régulières

*Sous le pont Mirabeau coule la Seine
Et nos amours
Faut-il qu'il m'en souviennne
La joie venait toujours après la peine
[...]*

*L'amour s'en va comme cette eau
courante
L'Amour s'en va
Comme la vie est lente
Et comme l'Espérance est violente*

Exemples

Déterminer les fragments de code permettant de :

- 1 Trouver les mots qui commencent par une majuscule mais qui ne sont pas des débuts de phrase.
- 2 Remplacer le mot "lente" par "rapide".

Pattern matching et expressions régulières

*Sous le pont Mirabeau coule la Seine
Et nos amours
Faut-il qu'il m'en souvienne
La joie venait toujours après la peine
[...]*

*L'amour s'en va comme cette eau
courante
L'Amour s'en va
Comme la vie est lente
Et comme l'Espérance est violente*

Exemples

Déterminer les fragments de code permettant de :

- ➊ Trouver les mots qui commencent par une majuscule mais qui ne sont pas des débuts de phrase.
- ➋ Remplacer le mot "lente" par "rapide".
- ➌ Vérifier que le texte ne contient pas de ponctuation (hormis [...]).

Pattern matching et expressions régulières

*Sous le pont Mirabeau coule la Seine
Et nos amours
Faut-il qu'il m'en souvienn
La joie venait toujours après la peine
[...]*

*L'amour s'en va comme cette eau
courante
L'Amour s'en va
Comme la vie est lente
Et comme l'Espérance est violente*

Exemples

Déterminer les fragments de code permettant de :

- 1 Trouver les mots qui commencent par une majuscule mais qui ne sont pas des débuts de phrase.
- 2 Remplacer le mot "lente" par "rapide".
- 3 Vérifier que le texte ne contient pas de ponctuation (hormis [...]).
- 4 Vérifier que la première et la dernière phrase du premier couplet sont des rimes en "eine".

Pattern matching et expressions régulières

*Sous le pont Mirabeau coule la Seine
Et nos amours
Faut-il qu'il m'en souvienn
La joie venait toujours après la peine
[...]*

*L'amour s'en va comme cette eau
courante
L'Amour s'en va
Comme la vie est lente
Et comme l'Espérance est violente*

Exemples

Déterminer les fragments de code permettant de :

- 1 Trouver les mots qui commencent par une majuscule mais qui ne sont pas des débuts de phrase.
- 2 Remplacer le mot "lente" par "rapide".
- 3 Vérifier que le texte ne contient pas de ponctuation (hormis [...]).
- 4 Vérifier que la première et la dernière phrase du premier couplet sont des rimes en "eine".

Pattern matching et expressions régulières

```
// s equals the poem

/** Question 1 */
Pattern p = Pattern.compile("[^\\n[A-z]] [A-Z] [a-z]+");
Matcher m = p.matcher(s);
while (m.find()) {
    System.out.println(s.substring(m.start(0)+1, m.end(0)));
}

/** Question 2 */
s.replaceAll(" lente", " rapide");

/** Question 3 */
Pattern.compile("[?!.,;]").matcher(s).find() &&
    !Pattern.compile("\\[\\].\\.\\.\\.\\.\\.\\.\\.").matcher(s).find();

/** Question 4 */
Pattern.matches("(\\n|.)+eine\\n(\\n|.)+eine\\n\\[\\].\\.\\.\\.\\.\\.\\.\\.\\. (\\n|.)*", s);
```

Table of Contents

- 1 Algèbre de Kleene
- 2 Conversion AFN \rightarrow ER
- 3 Conversion ER \rightarrow AFN
- 4 Pattern matching
- 5 Lemme de l'étoile

Certains langages ne sont pas reconnaissables. On ne peut pas les exprimer à l'aide d'un automate.

Exemples

- $L_1 = \{a^n b^n \mid n \in \mathbb{N}\}$
- $L_2 = \{w \in \{a, b\}^* \mid w \text{ possède autant de } a \text{ que de } b\}$
- $L_3 = \{w \in \{a, b\}^* \mid w \text{ possède un nombre différent de } a \text{ et de } b\}$

Comment montrer qu'un langage n'est pas reconnaissable ?

Certains langages ne sont pas reconnaissables. On ne peut pas les exprimer à l'aide d'un automate.

Exemples

- $L_1 = \{a^n b^n \mid n \in \mathbb{N}\}$
- $L_2 = \{w \in \{a, b\}^* \mid w \text{ possède autant de } a \text{ que de } b\}$
- $L_3 = \{w \in \{a, b\}^* \mid w \text{ possède un nombre différent de } a \text{ et de } b\}$

Comment montrer qu'un langage n'est pas reconnaissable ?

Réponse : Par l'absurde, à l'aide du lemme de l'étoile ou grâce aux propriétés de clôture.

Lemme de l'étoile

Soit $L \in \mathbf{Reg}(\Sigma^*)$. Soit $N \geq 0$ et $w \in L$, $\exists x, y, z$ tel que $w = xyz$ avec $0 < |y| \leq N$, alors $\forall n \geq 0, xy^n z \in L$

Précisons que L est supposé de cardinal infini.

Le lemme de l'étoile énonce que tout mot assez long et provenant de L comporte une partie non vide qui peut être itérée à l'infini.

Démonstration.

Soit $L \in \mathbf{Reg}$, d'après le théorème de Kleene, il existe un automate fini M tel que $L(M) = L$. Soit N le nombre d'états de M .

Lemme de l'étoile

Démonstration.

Soit $L \in \mathbf{Reg}$, d'après le théorème de Kleene, il existe un automate fini M tel que $L(M) = L$. Soit N le nombre d'états de M .

Soit $w \in L$, tel que $|w| \geq N$. On sait également que $w \in L(M)$, c'est-à-dire qu'il existe un chemin $q_0 \xrightarrow{w} q_t$ avec q_0 comme état initial et q_t appartenant aux états finaux.

Lemme de l'étoile

Démonstration.

Soit $L \in \mathbf{Reg}$, d'après le théorème de Kleene, il existe un automate fini M tel que $L(M) = L$. Soit N le nombre d'états de M .

Soit $w \in L$, tel que $|w| \geq N$. On sait également que $w \in L(M)$, c'est-à-dire qu'il existe un chemin $q_0 \xrightarrow{w} q_t$ avec q_0 comme état initial et q_t appartenant aux états finaux.

On pose $w = a_1 a_2 \dots a_N w'$ avec $a_{i \in [1, N]} \in \Sigma$ et $w' \in \Sigma^*$. Soient q_1, q_2, \dots, q_N les états successifs atteints au cours de la lecture des N premiers symboles. On a donc :

$$q_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} q_2 \xrightarrow{a_3} \dots \xrightarrow{a_N} q_N \xrightarrow{w'} q_t$$

Lemme de l'étoile

Démonstration.

Soit $L \in \mathbf{Reg}$, d'après le théorème de Kleene, il existe un automate fini M tel que $L(M) = L$. Soit N le nombre d'états de M .

Soit $w \in L$, tel que $|w| \geq N$. On sait également que $w \in L(M)$, c'est-à-dire qu'il existe un chemin $q_0 \xrightarrow{w} q_t$ avec q_0 comme état initial et q_t appartenant aux états finaux.

On pose $w = a_1 a_2 \dots a_N w'$ avec $a_i \in [1, N] \in \Sigma$ et $w' \in \Sigma^*$. Soient q_1, q_2, \dots, q_N les états successifs atteints au cours de la lecture des N premiers symboles. On a donc :

$$q_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} q_2 \xrightarrow{a_3} \dots \xrightarrow{a_N} q_N \xrightarrow{w'} q_t$$

Comme M possède N états et comme de q_0 à q_N on parcourt $N + 1$ états, alors d'après le principe des tiroirs, il existe deux entiers i, j tels que $0 \leq i < j \leq N$ avec $q_i = q_j$. On pose q cet état et $x = a_1 \dots a_i$, $y = a_{i+1} \dots a_j$, $z = a_{j+1} \dots a_N w'$. On a alors le chemin de w peut être factorisé tel que :

$$q_0 \xrightarrow{x} q \xrightarrow{y} q \xrightarrow{z} q_t$$

On a exhibé un cycle au niveau du mot y , il vient que $\forall n \geq 0, xy^n z \in L$. De plus, on a $|y| = j - i$, donc $0 < |y| \leq N$. □

Exemples

On rappelle que $L_1 = \{a^n b^n \mid n \geq 0\}$. Soit :

$$w = \underbrace{aa \dots a}_N \underbrace{bb \dots b}_N = a_1 \dots a_N b_1 \dots b_N$$

On suppose que L_1 est régulier.

Posons : $x = a_1 \dots a_{N-k}$, $y = a_{N-k+1} \dots a_N$ et $z = b_1 \dots b_N$ tel que $w = xyz$ et $|y| = k$.

Alors, d'après le lemme de l'étoile : $\forall n \geq 0, xy^n z \in L$.

Il vient que $N = N + k(n - 1) \Leftrightarrow 0 = k(n - 1)$

Or $k > 0$, la propriété n'est pas valable pour tout $n \geq 0$. L'hypothèse de départ est absurde, L_1 n'est pas régulier.

Théorèmes de clôture

On montre que les autres langages ne sont pas reconnaissables grâce aux propriétés de clôture.

Théorème : Clôture par \cap et \neg

- Soient L_1, L_2, L_3 tels que $L_1 \cap L_2 = L_3 \notin \mathbf{Reg}$, alors $L_1 \notin$ ou $L_2 \notin \mathbf{Reg}$.
- Soit $L \notin \mathbf{Reg}$, alors $\neg L \notin \mathbf{Reg}$.

Théorèmes de clôture

On montre que les autres langages ne sont pas reconnaissables grâce aux propriétés de clôture.

Théorème : Clôture par \cap et \neg

- Soient L_1, L_2, L_3 tels que $L_1 \cap L_2 = L_3 \notin \mathbf{Reg}$, alors $L_1 \notin$ ou $L_2 \notin \mathbf{Reg}$.
- Soit $L \notin \mathbf{Reg}$, alors $\neg L \notin \mathbf{Reg}$.

Exemples

On sait que $L_1 \notin \mathbf{Reg}$.

- $L_2 \cap a^*b^* = L_1 \longrightarrow$ Donc L_2 n'est pas reconnaissable.
- $\neg L_3 = L_2 \longrightarrow$ Donc L_3 n'est pas reconnaissable.

Morphisme

Soient deux monoïdes (E, \star, e) et $(F, *, f)$. Un *morphisme* $\varphi : E \rightarrow F$ est une application qui vérifie :

- $\forall x, y \in E, \varphi(x \star y) = \varphi(x) * \varphi(y)$
- $\varphi(e) = f$

Exemples

Soit $\Sigma = \{0, 1\}$ et le morphisme $\varphi : \Sigma^* \rightarrow \Sigma^*$ tel que

$$\begin{cases} \varphi(0) = 01 \\ \varphi(1) = 10 \end{cases}$$

Alors, $\varphi(011) = \varphi(0)\varphi(1)\varphi(1) = 011010$

Théorème : Clôture par morphisme

Soit $L \in \mathbf{Reg}$ et φ un morphisme, alors $\varphi^{-1}(L) \in \mathbf{Reg}$.

Ce théorème peut être utilisé pour montrer qu'un langage n'est pas régulier.

Exemples

Soit $L_4 = \{w \in \{a, b\}^* \mid w \text{ a 2 fois plus de } a \text{ que de } b\}$

Théorème : Clôture par morphisme

Soit $L \in \mathbf{Reg}$ et φ un morphisme, alors $\varphi^{-1}(L) \in \mathbf{Reg}$.

Ce théorème peut être utilisé pour montrer qu'un langage n'est pas régulier.

Exemples

Soit $L_4 = \{w \in \{a, b\}^* \mid w \text{ a 2 fois plus de } a \text{ que de } b\}$

Si L_4 est régulier, alors $L_4 \cap a^*b^* = \{a^{2n}b^n \mid n \geq 0\}$ l'est aussi.

Or, soit $\varphi(0) = aa$ et $\varphi(1) = b$, alors $\varphi^{-1}(L_4 \cap a^*b^*) = \{0^n1^n \mid n \geq 0\}$ que l'on sait irrégulier.

Dès lors L_4 est irrégulier.