

Synthèse sur les langages et automates

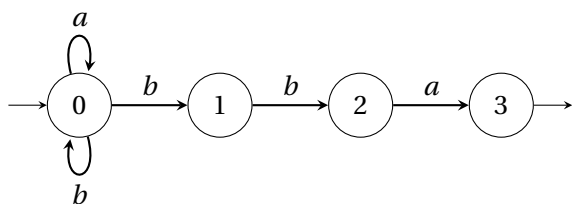
OPTION INFORMATIQUE - TP n° 4.3 - Olivier Reynet

À la fin de ce chapitre, je sais :

- Utiliser les définitions des mot et des langages
- Simplifier les expressions régulières
- Trouver le langage associé à une expression régulière
- Construire l'automate correspondant à un langage régulier
- Déterminer un AFND
- Transformer une expression régulière en AFND avec Thompson ou Berry-Sethi
- Transformer un automate en une expression régulière
- Montrer qu'un langage n'est pas régulier

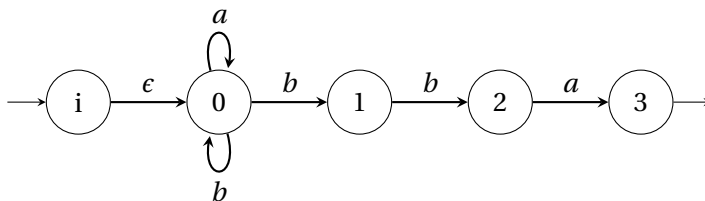
A Automate vers expression régulière

Pour les automates suivants, donner une expression régulière telle que le langage reconnaissable par l'automate est le langage dénoté par l'expression régulière. On éliminera les états après avoir normalisé l'automate en entrée et en sortie.

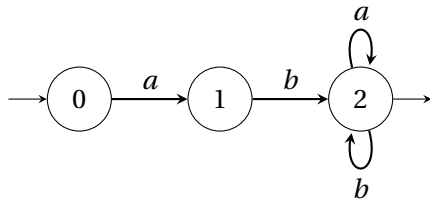


A1.

Solution : Normalisation :



Par élimination d'états : $(a|b)^* bba$



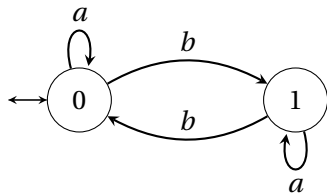
A2.

Solution : Normalisation :

```

    graph LR
        start(( )) --> 0((0))
        0 -- a --> 1((1))
        1 -- b --> 2((2))
        2 -- a --> 2
        2 -- b --> 2
        2 -- epsilon --> f((f))
        f --> end(( ))
    
```

Par élimination d'états : $ab(a|b)^*$



A3.

Solution : Normalisation :

```

    graph LR
        start(( )) --> i((i))
        i -- epsilon --> 0((0))
        0 -- a --> 0
        0 -- b --> 1((1))
        1 -- b --> 0
        1 -- a --> 1
        0 -- epsilon --> f((f))
        f --> end(( ))
    
```

Éliminations de 1 :

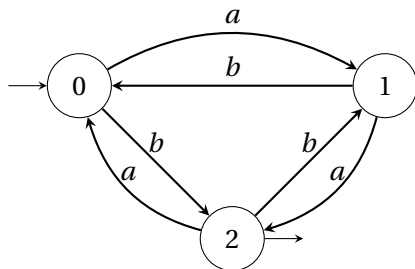
```

    graph LR
        start(( )) --> i((i))
        i -- epsilon --> 0((0))
        0 -- a --> 0
        0 -- ba*b --> 0
        0 -- epsilon --> f((f))
        f --> end(( ))
    
```

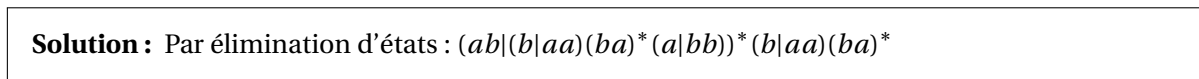
Fusion des expressions régulières à même destination depuis 0 :

```

    graph LR
        start(( )) --> i((i))
        i -- epsilon --> 0((0))
        0 -- "a|ba*b" --> 0
        0 -- epsilon --> f((f))
        f --> end(( ))
    
```



A4.



B Expression régulière vers automate

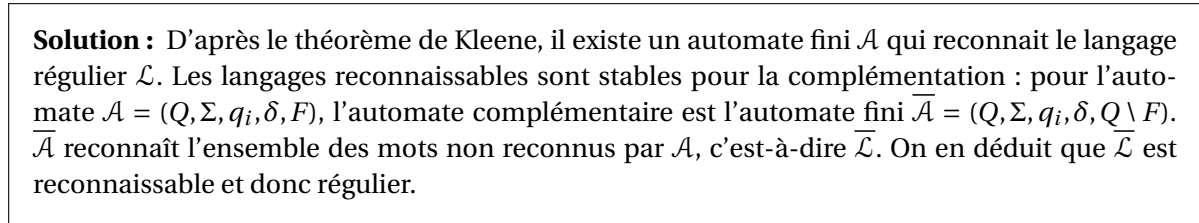
Pour chaque expression régulière, appliquer :

1. l'algorithme de Berry-Sethi, trouver l'automate de Glushkov et le déterminer si besoin.
2. l'algorithme de Thompson (méthode compositionnelle), trouver l'automate et éliminer les transitions spontanées.

- B1. $a|b$
- B2. a^*b
- B3. $(a|b)^*a$
- B4. $(a|b)a^*b^*$
- B5. $(a^*b)|(a(a|b)^*)$
- B6. $(ba|a)^*ab$
- B7. $(a|c)^*abb|(a|c)^*$
- B8. $(ba|a)^*(a|b)c$

C Stabilité des langages réguliers

- C1. Soit Σ un alphabet et \mathcal{L} un langage régulier sur Σ . Montrer que le complémentaire de \mathcal{L} , $C(\mathcal{L}) = \overline{\mathcal{L}} = \{w, w \in \Sigma^* \text{ et } w \notin \mathcal{L}\}$, est régulier.



- C2. Soit Σ un alphabet, \mathcal{L}_1 et \mathcal{L}_2 deux langages réguliers sur Σ . Montrer que $\mathcal{L}_1 \cap \mathcal{L}_2$ est un langage régulier.

Solution : On passe par la loi de Morgan : $\mathcal{L}_1 \cap \mathcal{L}_2 = \overline{\overline{\mathcal{L}_1} \cup \overline{\mathcal{L}_2}}$. Or, les langages réguliers sont stables pour l'union (par construction) et la complémentation (on vient le montrer). Donc, les langages réguliers sont stables pour l'intersection.

- C3. Soit Σ un alphabet, \mathcal{L}_1 et \mathcal{L}_2 deux langages réguliers sur Σ . Montrer que $\mathcal{L}_1 \setminus \mathcal{L}_2$ est un langage régulier.

Solution : On a $\mathcal{L}_1 \setminus \mathcal{L}_2 = \mathcal{L}_1 \cap \overline{\mathcal{L}_2}$. Comme les langages réguliers sont stables pour l'intersection et la complémentation, ils sont également stables pour la différence ensembliste.

- C4. Soit Σ un alphabet et \mathcal{L} un langage régulier sur Σ . On définit le mot miroir de $w = a_1 a_2 \dots a_n$ par $w^R = a_n a_{n-1} \dots a_1$ et le langage miroir $\mathcal{L}^R = \{u \in \Sigma^*, u^R \in \mathcal{L}\}$. Montrer que \mathcal{L}^R est régulier.

Solution : On considère un automate fini non déterministe $\mathcal{A} = (Q, \Sigma, Q_i, \Delta, F)$ associé à \mathcal{L} . Soit l'automate fini $\mathcal{A}^R = (Q, \Sigma, F, \Delta^{-1}, Q_i)$. On a noté Δ^{-1} l'ensemble des transitions obtenues en inversant le sens de chaque transition de Δ : par exemple, (q, a, q') devient (q', a, q) . \mathcal{A}^R reconnaît le langage miroir \mathcal{L}^R , car pour tout mot w de \mathcal{L} il existe un chemin acceptant dans \mathcal{A} qui correspond à un chemin acceptant w^R dans \mathcal{A}^R . \mathcal{L}^R est un langage reconnaissable donc régulier. Les langages réguliers sont stables par miroir.

- C5. Soit Σ un alphabet et \mathcal{L} un langage régulier sur Σ . Montrer que $\text{Pref}(\mathcal{L}) = \{u \in \Sigma^*, \exists v \in \Sigma^*, uv \in \mathcal{L}\}$, l'ensemble de mots préfixes de \mathcal{L} , est un langage régulier.

Solution : Soit $\mathcal{A} = (Q, \Sigma, q_i, \delta, F)$ un automate fini déterministe reconnaissant \mathcal{L} . Soit C l'ensemble des états co-accessibles de \mathcal{A} . Alors l'automate $\mathcal{A}_{pref} = (Q, \Sigma, q_i, \delta, C)$ est un automate fini qui reconnaît $\text{Pref}(\mathcal{L})$. Les langages réguliers sont donc stables pour l'opération préfixe.

- C6. Soit Σ un alphabet et \mathcal{L} un langage régulier sur Σ . Montrer que $\text{Suff}(\mathcal{L}) = \{v \in \Sigma^*, \exists u \in \Sigma^*, uv \in \mathcal{L}\}$, l'ensemble de mots suffixes de \mathcal{L} , est un langage régulier.

Solution : Soit $\mathcal{A} = (Q, \Sigma, q_i, \delta, F)$ un automate fini déterministe reconnaissant \mathcal{L} . Soit A l'ensemble des états accessibles de \mathcal{A} . Alors l'automate $\mathcal{A}_{suff} = (Q, \Sigma, A, \delta, F)$ est un automate fini non déterministe qui reconnaît $\text{Suff}(\mathcal{L})$. Les langages réguliers sont donc stables pour l'opération suffixe.

- C7. Soit Σ un alphabet et \mathcal{L} un langage régulier sur Σ . Montrer que $\text{Fact}(\mathcal{L}) = \{w \in \Sigma^*, \exists u, v \in \Sigma^*, u w v \in \mathcal{L}\}$, l'ensemble de mots facteurs de \mathcal{L} , est un langage régulier.

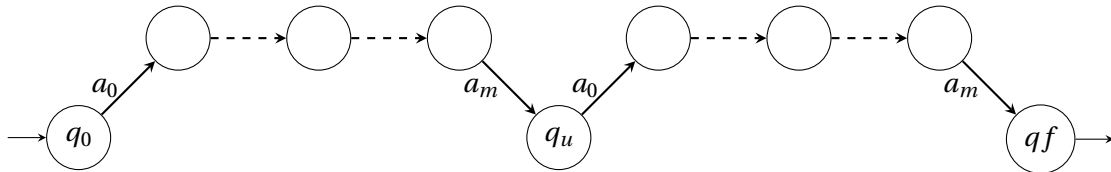
Solution : On remarque que $\text{Fact}(\mathcal{L}) = \text{Pref}(\text{Suff}(\mathcal{L}))$. Les langages réguliers sont donc stables pour l'opération facteur.

C8. Soit Σ un alphabet et \mathcal{L} un langage régulier sur Σ . On définit le langage racine $\sqrt{\mathcal{L}}$ par :

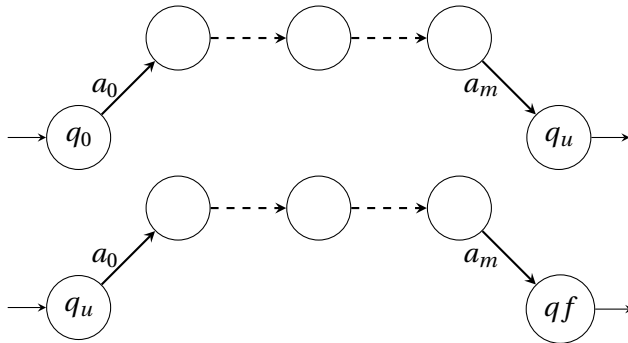
$$\sqrt{\mathcal{L}} = \{u \in \Sigma^*, uu \in \mathcal{L}\} \tag{1}$$

Montrer que $\sqrt{\mathcal{L}}$ est régulier en vous appuyant sur l'automate reconnaissant le langage \mathcal{L} .

Solution : Soit $\mathcal{A} = (Q, \Sigma, q_i, \delta, F)$ un automate fini déterministe complet reconnaissant le langage régulier \mathcal{L} . Soit $w = a_0 a_1 \dots a_m$ un mot de $\sqrt{\mathcal{L}}$. Nécessairement, le mot $w w$ est reconnu par \mathcal{A} . Cela signifie qu'il existe un état q et un chemin dans l'automate tel que :



Dire que w est un mot de $\sqrt{\mathcal{L}}$ est donc équivalent à dire que ce mot est reconnu à la fois par l'automate $\mathcal{A}_d = (Q, \Sigma, q_0, \delta, \{q_u\})$ et par l'automate $\mathcal{A}_f = (Q, \Sigma, q_u, \delta, \{q_f\})$.



Essayons de construire un automate qui reconnaît un mot u de $\sqrt{\mathcal{L}}$.

$$\mathcal{A}_q = (Q^2, \Sigma, (q_0, q), \delta_q, F_q) \tag{2}$$

où la fonction de transition δ_q est construite ainsi :

$$\delta_q((q_i, q_j), a) = (\delta(q_i, a), \delta(q_j, a)) \tag{3}$$

et l'ensemble des états accepteurs F_q est défini comme suit :

$$F_p = \{(q_i, q_f), q_i \in Q, q_f \in F\} \tag{4}$$

Créons alors l'automate $\mathcal{A}_r = \bigcup_{q \in Q} \mathcal{A}_q$ et montrons que $\mathcal{L}_{rec}(\mathcal{A}_r) = \sqrt{\mathcal{L}}$.

Démonstration. (\implies) Soit un mot u reconnu par \mathcal{A}_r . Cela signifie qu'il existe un automate \mathcal{A}_{q_u} qui reconnaît u . Pour cet automate, on a $\delta_q^*((q_0, q_u), u) = (q_u, q_f) \in F_{q_u}$ ce qui signifie que $\delta^*(q_0, u) = q_u$ et $\delta^*(q_u, u) = q_f \in F$. \mathcal{A} reconnaît donc le mot uu et u appartient au langage $\sqrt{\mathcal{L}}$.

(\impliedby) Soit un mot u du langage $\sqrt{\mathcal{L}}$. Dans l'automate \mathcal{A} , cela signifie qu'il existe un état q_u tel que $\delta^*(q_0, u) = q_u$ et $\delta^*(q_u, u) = q_f \in F$. Ceci peut s'écrire $(\delta^*(q_0, u), \delta^*(q_u, u)) = (q_u, q_f) \in F$. Or, par construction, il existe un automate \mathcal{A}_{q_u} de \mathcal{A}_r tel que $\delta_q^*((q_0, q_u), u) = (q_u, q_f) \in F_{q_u}$. Cet automate reconnaît u et donc \mathcal{A}_r reconnaît u . ■

$\sqrt{\mathcal{L}}$ est reconnu par un automate fini et est donc un langage régulier.

C9. Montrer que tout langage fini est un langage régulier.

Solution : Soit \mathcal{L} un langage fini de cardinal n . Chaque mot w_i de ce langage fini peut être représenté par une expression régulière e_i , celle qui utilise toutes les lettres nécessaires de l'alphabet pour le représenter. D'après la sémantique des expressions régulières, on a :

$$\mathcal{L} = \bigcup_{i \in \llbracket 0, n-1 \rrbracket} \mathcal{L}_{ER}(e_i) = \mathcal{L}_{ER}(e_0 | e_1 | \dots | e_{n-1})$$

Par définition, cela signifie que \mathcal{L} est un langage dénoté par une expression régulière. Par conséquent, un langage fini est toujours régulier.

D Lemme de l'étoile et non régularité

D1. Montrer que les langages ci-dessous ne sont pas réguliers.

(a) $\mathcal{L}_1 = \{a^p, p \text{ est premier}\}$

Solution : Supposons que \mathcal{L}_1 soit régulier. D'après le théorème de Kleene, il existe une automate fini à n états qui reconnaît le langage \mathcal{L}_1 .

Soit p un nombre premier strictement plus grand que n et $w = a^p$, un mot de le mot de \mathcal{L}_1 . D'après le lemme de l'étoile, il existe une décomposition de w en xyz tels que $|xy| \leq n$, $y \neq \epsilon$ et $xy^*z \subseteq \mathcal{L}_1$.

Soit k et j deux entiers tels que $k + j \leq n$ et $j > 0$. On pose $x = a^k$, $y = a^j$ et $z = a^{p-k-j}$. Une décomposition de w satisfaisant les conditions du lemme de l'étoile est nécessairement de la forme générale $w = xyz = a^k a^j a^{p-k-j}$. On peut donc itérer sur y . En particulier, on peut choisir d'élever y à la puissance $p + 1$ et obtenir u , un mot de \mathcal{L}_1 . $u = xy^{p+1}z = a^k (a^j)^{p+1} a^{p-k-j} = a^{(j+1)p}$. Or, $(j+1)p$ n'est pas premier et donc $u \notin \mathcal{L}_1$. On aboutit donc à une contradiction. \mathcal{L}_1 n'est pas un langage régulier.

(b) $\mathcal{L}_2 = \{w \in \Sigma^*, \Sigma = \{a, b\}, w \text{ possède autant de } a \text{ que de } b\}$

Solution : Supposons que \mathcal{L}_2 soit régulier. Soit $\mathcal{L} = \mathcal{L}_2 \cap \mathcal{L}_{ER}(a^* b^*)$. Comme $\mathcal{L}_{ER}(a^* b^*)$ est régulier car dénoté par une expression régulière et que l'intersection de deux langages réguliers est régulier, alors si \mathcal{L}_2 est régulier, le langage \mathcal{L} est régulier. On se propose donc de montrer que \mathcal{L} n'est pas régulier ce qui impliquera que \mathcal{L}_2 ne l'est pas non plus.

Or $\mathcal{L} = \{a^n b^n, n \in \mathbb{N}\}$ est le langage des puissances et on sait qu'il n'est pas régulier. Par conséquent, \mathcal{L}_2 non plus.

(c) $\mathcal{L}_3 = \{a^i b^j, i < j\}$

Solution : Supposons que \mathcal{L}_3 soit régulier. Soit \mathcal{A} un automate fini à n états qui reconnaît le langage \mathcal{L}_3 . Soit un mot w de \mathcal{L}_3 de longueur supérieure à n : $w = a^i b^j$ et $i + j > n$ et $n \leq i$. On peut lui appliquer le lemme de l'étoile.

Soit k et h deux entiers tels que $k + h \leq i$ et $h > 0$. On pose $x = a^k$, $y = a^h$ et $z = a^{i-k-h} b^j$. La décomposition $w = xyz$ est une forme générale des décompositions qui satisfont les

conditions du lemme de l'étoile puisque $n \leq i$. On peut donc itérer sur y . En particulier, on peut élever y à la puissance $j+1$ et obtenir un mot u de \mathcal{L}_3 : $u = xy^{j+1}z = a^k(a^h)^{j+1}a^{i-k-h}b^j = a^{i+hj}b^j$. Comme $h > 0$, $i + hj > j$. Donc $u \notin \mathcal{L}_3$. On aboutit à une contradiction. \mathcal{L}_3 n'est pas un langage régulier.

(d) $\mathcal{L}_4 = \{a^p, p \text{ n'est pas premier}\}$

Solution : Comme les langages réguliers sont stables par complémentation et que \mathcal{L}_1 n'est pas régulier, \mathcal{L}_4 n'est pas régulier.

D2. Montrer que l'ensemble des mots de Dyck \mathcal{D} n'est pas un langage régulier. On rappelle que \mathcal{D} est l'ensemble de mots bien parenthésés sur un alphabet fini de parenthèses ouvrantes et fermantes. Par exemple, sur la paire de parenthèses formée de (et), le mot $(())()$ est un mot bien parenthésé, alors que le mot $()()$ ne l'est pas.

Solution : En notant les parenthèses ouvrantes et fermantes a et b , dire que l'ensemble des mots de Dyck est un langage régulier implique que l'ensemble $\{a^n b^n, n \in \mathbb{N}\} \subset \mathcal{D}$ est un langage régulier. Or ce n'est pas le cas car il s'agit du langage des puissances qui n'est pas régulier (cf. cours).

Plus précisément, on peut écrire : $\{a^n b^n, n \in \mathbb{N}\} = \mathcal{D} \cap \mathcal{L}_{ER}(a^* b^*)$. Or, l'intersection de deux langages réguliers est un langage régulier. Donc si \mathcal{D} était régulier, le langage des puissances le serait aussi puisque $\mathcal{L}_{ER}(a^* b^*)$ est un langage régulier car dénoté par une expression régulière.

D3. Soit $\Sigma = \{a, b\}$ un alphabet à deux lettres et \mathcal{L}_{pal} l'ensemble des palindromes sur Σ . Montrer que \mathcal{L}_{pal} n'est pas un langage régulier.

Solution : Supposons que \mathcal{L}_{pal} soit régulier et n une constante d'itération de ce langage régulier.

Soit w un mot de \mathcal{L}_{pal} de longueur supérieure ou égale à n . Considérons le mot $w = a^n b a^n$ qui est un palindrome. Ce mot appartient au langage \mathcal{L}_i intersection de \mathcal{L}_{pal} et de $\mathcal{L}_{ER}(a^* b a^*)$. Comme $\mathcal{L}_{ER}(a^* b a^*)$ est un langage régulier et que l'intersection de deux langages réguliers est un langage régulier, si \mathcal{L}_{pal} est régulier alors \mathcal{L}_i l'est aussi.

On a $|w| = 2n + 1 > n$ et on peut appliquer le lemme de l'étoile à $w \in \mathcal{L}_i$. Il existe donc une décomposition w en xyz tels que $|xy| = k \leq n$ et $y \neq \epsilon$. Selon cette décomposition générale, il existe donc deux entiers i et $j > 0$ tels que $i + j = k$, $x = a^i$, $y = a^j \neq \epsilon$ et $z = a^{n-i-j} b a^n$. Le lemme de l'étoile nous dit que l'on peut itérer sur y et obtenir un mot de \mathcal{L}_i , $xy^* z \in \mathcal{L}_i$.

Soit le mot u obtenu en prenant la puissance 2 de y : $u = xy^2 z = a^i a^{2j} a^{n-i-j} b a^n = a^{n+j} b a^n$. Comme $j > 0$, ce mot n'est pas un palindrome et donc n'appartient pas au langage \mathcal{L}_i , ce qui est en contradiction avec le lemme de l'étoile : \mathcal{L}_i n'est pas régulier. Donc \mathcal{L}_{pal} n'est pas un langage régulier.