

4. Expressions régulières

Enseignant: Arnaud Casteigts

Assistants: A.-Q. Berger & M. Marseloo

Moniteurs: N. Beghdadi & E. Bussod

4.1 Expressions régulières

Les expressions régulières (ER) permettent de décrire les langages réguliers. Pour rappel la famille des langages réguliers est définie comme suit :

- \emptyset et $\{\varepsilon\}$ sont des langages réguliers,
- Pour tout $s \in \Sigma$, $\{s\}$ est un langage régulier,
- Si L est un langage régulier, alors L^* l'est aussi,
- Si L_1 et L_2 sont des langages réguliers, alors $L_1 \cup L_2$ et $L_1 \circ L_2$ le sont aussi.

Les expressions régulières consistent à représenter ces opérations de manière simplifiée. Notamment, les accolades sont omises et la concaténation \circ est généralement omise aussi. Par ailleurs, des parenthèses peuvent être utilisées pour gérer les priorités entre opérations, comme pour les expressions arithmétiques classiques.

Par exemple, sur l'alphabet $\Sigma = \{a, b, c\}$, l'expression $(a \cup b)^*ca^+$ désigne le langage $(\{a\} \cup \{b\})^* \circ \{c\} \circ \{a\}^+$, à savoir tous les mots qui contiennent exactement un c , précédé de n'importe quel préfixe sans c et suivi d'un nombre arbitraire (mais strictement positif) de a . On a donc $(a \cup b)^*ca^+ = \{ca, aca, bca, caa, aaca, abca, baca, bbca, caaa, aacaa, abcaa, \dots\}$

Outre leur intérêt dans l'étude des langages formels, les expressions régulières ont de nombreuses applications, notamment pour la recherche de motif textuels dans les fichiers. Plusieurs langages de programmation et commandes UNIX les supportent nativement.

4.1.1 Quelques exemples

Voici quelques exemples d'ERs sur l'alphabet $\Sigma = \{0, 1\}$:

- 0^*10^* = mots contenant un seul 1.
- $\Sigma^*1\Sigma^*$ = mots contenant au moins un 1.
- $\Sigma^*001\Sigma^*$ = mots contenant le facteur 001.
- $1^*(01^+)^*$ = mots dans lesquels chaque 0 est suivi d'au moins un 1.
- $(\Sigma\Sigma)^*$ = mots de longueur paire.

- $(\Sigma\Sigma\Sigma)^* =$ mots dont la longueur est multiple de trois.
- $01 \cup 10 = \{01, 10\}$.
- $0\Sigma^*0 \cup 1\Sigma^*1 \cup 0 \cup 1 =$ mots qui commencent et finissent par le même symbole.
Les priorités sont implicites ici, il faut savoir que la concaténation est prioritaire sur l'union. Cette expression se lit donc comme $(0\Sigma^*0) \cup (1\Sigma^*1) \cup 0 \cup 1$.
- $(0 \cup \varepsilon)(1 \cup \varepsilon) = \{\varepsilon, 0, 1, 01\}$.

Remarque : on autorise l'usage de Σ dans ces expressions, désignant n'importe quel facteur de longueur 1 sur Σ ". Ici, Σ^* est donc synonyme de $\{0 \cup 1\}^*$, autrement dit n'importe quel mot de longueur arbitraire fait de symboles de Σ .

4.2 Expressions régulières \rightarrow AFN

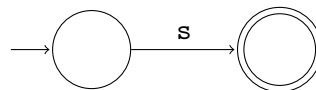
Dans cette section, nous allons montrer que n'importe quelle expression régulière peut être transformée en un AFN qui reconnaît exactement le même langage. Les AFNs peuvent donc reconnaître n'importe quel langage régulier. Et comme nous avons déjà vu (cours n°3) que tout AFN peut être déterminisé en un AFD, il en va de même des AFDs.

L'idée est de construire un AFNs graduellement, en utilisant les mêmes opérations que celles qui définissent les langages/expressions régulières, à savoir :

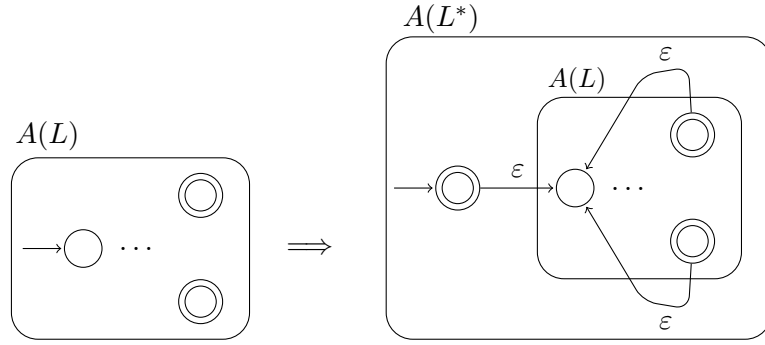
- Les langages \emptyset et $\{\varepsilon\}$ peuvent être reconnus par les automates suivants (à gauche et à droite, respectivement) :



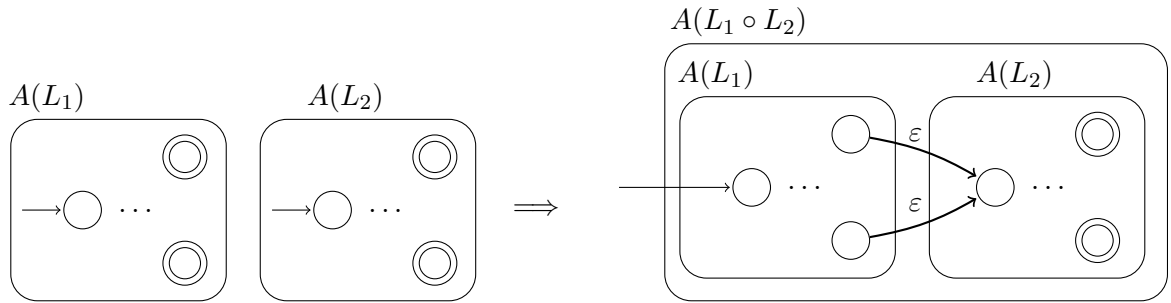
- Pour chaque symbole $s \in \Sigma$, le langage $\{s\}$ peut être reconnu par l'automate suivant :



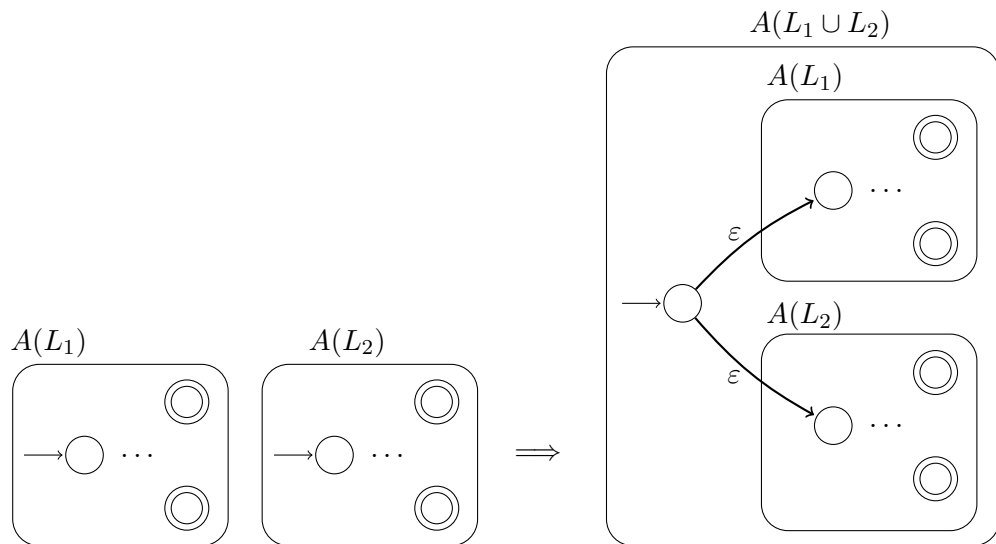
- Si L est reconnu par un AFN $A(L)$, alors L^* est reconnu par l'automate $A(L^*)$ suivant :



- Si L_1 est reconnu par un AFN $A(L_1)$ et L_2 est reconnu par un AFN $A(L_2)$, alors $L_1 \circ L_2$ est reconnu par l'automate $A(L_1 \circ L_2)$ suivant :



- Si L_1 est reconnu par un AFN $A(L_1)$ et L_2 est reconnu par un AFN $A(L_2)$, alors $L_1 \cup L_2$ est reconnu par l'automate $A(L_1 \cup L_2)$ suivant :



4.2.1 Commentaires

Pour l'opération L^* , la construction consiste à :

- Créer un nouvel état initial qui est également final et ajouter une ε -transition depuis cet état vers l'état initial de $A(L)$.
- Créer une ε -transition depuis chaque état final de $A(L)$ vers l'état initial de $A(L)$.

Le fait que l'ancien état initial soit aussi final permet d'entrer "zéro fois" dans l'automate et donc de reconnaître $L^0 = \{\varepsilon\}$. Le fait que les états finaux de $A(L)$ renvoient vers l'état initial de $A(L)$ permet de répéter le passage dans l'automate autant de fois qu'on veut, et donc de reconnaître L^k pour tout $k \geq 1$. On reconnaît donc bien L^* .

Pour l'opération $L_1 \circ L_2$, la construction consiste à :

- Brancher tous les états sortants de l'automate reconnaissant L_1 à l'état initial de l'automate reconnaissant L_2 (via des ε -transitions), ce qui permet de passer dans le second automate à chaque fois que le préfixe lu jusqu'à présent est dans L_1 .
- Rendre non-initial l'ancien état initial de l'automate reconnaissant L_2 , pour forcer à un lire un mot de L_1 avant de commencer à lire un mot de L_2 .
- Rendre les états finaux de L_1 non-finaux, pour forcer à lire ensuite un mot de L_2 .

Pour l'opération $L_1 \cup L_2$, la construction consiste à :

- Remplacer les deux états initiaux par un seul nouvel état initial, qui y sera relié par des ε -transitions. Cela permet d'explorer en parallèle (indépendamment) les deux automates, donc d'accepter un mot si au moins l'un des deux automates l'aurait accepté.

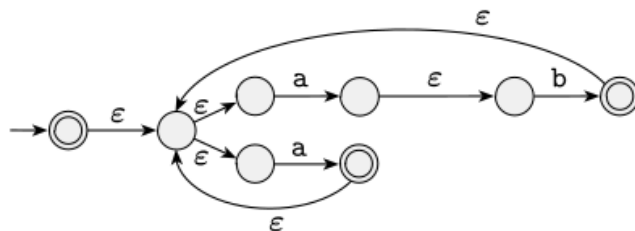
Ainsi, pour tout langage régulier (et donc pour toute expression régulière le représentant), il existe un AFN qui reconnaît exactement ce langage. Et comme nous avons vu que tout AFN peut être déterminisé en AFD, cela donne la chaîne de transformation suivante :

$$\text{ER} \longrightarrow \text{AFN} \longrightarrow \text{AFD}.$$

Il ne reste plus qu'à refermer le cycle en montrant que tout AFD peut être transformé en une expression régulière, et nous aurons alors montré que tous ces formalismes sont équivalents.

4.3 Exemple

$$(ab \cup a)^* = ((\{a\} \circ \{b\}) \cup \{a\})^*$$



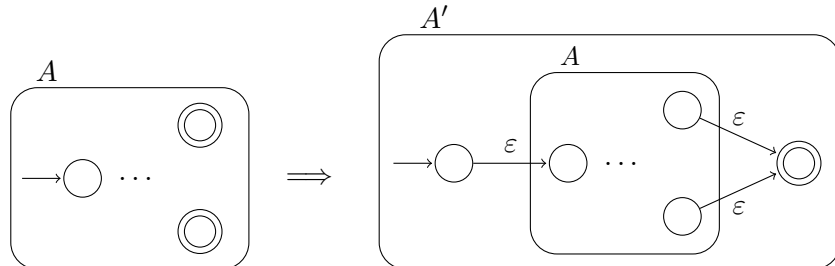
4.4 AFD \rightarrow Expressions régulières

(début du cours suivant)

Nous allons voir ici une autre transformation qui convertit un AFD quelconque en une expression régulière qui représente exactement le même langage. La première étape consiste à transformer notre AFD A en un AFN A' qui a les propriétés suivantes :

- L'état initial de A' n'a aucune transition *entrante*.
- Il n'y a qu'un état final de A' , et il n'a aucune transition *sortante*.

Cette transformation est simple, il suffit d'ajouter un nouvel état initial qui remplace l'ancien et qui le relie par une ε -transition. De même, il suffit d'ajouter un nouvel état final qui remplace tous les anciens, avec une ε -transition provenant de chacun d'entre eux. Clairement, ce nouvel automate reconnaît le même langage que A .



La deuxième étape consiste à supprimer, les uns après les autres, tous les états de A dans l'automate A' , en le modifiant au fur et à mesure jusqu'à ce qu'il ne reste plus que l'état initial et l'état final, avec une seule transition entre les deux. Cette transition indiquera (magie) l'expression régulière que nous cherchons. À ce stade, l'AFN n'est plus vraiment un AFN de base, mais un modèle plus général qui permet d'associer chaque transition à des morceaux d'expressions régulières.

L'intuition est la suivante. Supposons que nous avons, quelque part dans l'automate A' , trois états q , q_1 et q_2 avec, par exemple, les transitions suivantes (dessin de gauche) :



En regardant le dessin de gauche, on comprend que l'état q et ses transitions permettent, plus globalement, d'aller de q_1 à q_2 en lisant n'importe quel facteur correspondant à l'expression régulière ac^*b (s'il n'y avait pas de boucle sur q , ce serait juste ab). Sur cet exemple, q ne sert à rien d'autre. On peut donc le supprimer ainsi que ses transitions, en le remplaçant par une nouvelle transition de q_1 à q_2 correspondant à l'expression $E = ac^*b$ (dessin de droite).

Si vous avez compris l'opération ci-dessus, vous avez l'intuition principale (sinon, relisez le paragraphe précédent). Le cas général est juste un peu plus complexe, car :

1. Il pourrait déjà exister une transition de q_1 à q_2 .
→ Si une transition entre q_1 et q_2 existe déjà avec une autre expression E' , il suffit d'y ajouter la nouvelle expression en faisant l'union $E' \cup E$ (signifiant qu'on peut aller de q_1 à q_2 en lisant le motif E' ou le motif E).
2. L'état q pourrait être utilisé comme intermédiaire entre plusieurs couples d'états, pas seulement q_1 et q_2 .
→ Quand on veut supprimer un état q , il faut considérer *tous les couples d'états possibles* qui utilisent q comme intermédiaire, et ajouter toutes les transitions correspondantes, avant de pouvoir supprimer q .
3. Un état q_1 pourrait avoir des transitions vers q et depuis q (flèches dans les deux sens). L'état q serait donc un intermédiaire pour aller de q_1 à q_1 lui-même.
→ C'est la même chose qu'entre q_1 et q_2 , sauf que la transition ajoutée (ou sur laquelle on utilise l'union) est une boucle sur q_1 .

Ces éléments en tête, voici l'algorithme général : On notera ici S l'ensemble des états à supprimer, qui sont tous les états de Q'_A sauf son nouvel état initial et nouvel état final :

Tant que $S \neq \emptyset$:

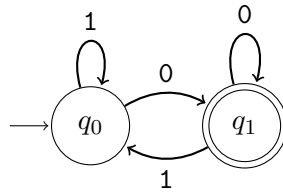
- | Choisir un état q de S (arbitrairement)
- | Pour chaque état $q_1 \neq q$ de Q'_A :
 - | | Pour chaque état $q_2 \neq q$ de Q'_A :
 - | | | Créer l'expression E correspondant aux transitions $q_1 \rightarrow \hat{q} \rightarrow q_2$ (si applicable)
 - | | | Ajouter (ou créer) E comme transition entre q_1 et q_2 (comme discuté ci-dessus)
- | Enlever q de S et de $Q_{A'}$

À l'issue de ces opérations, il ne reste plus que l'état initial et l'état final dans $Q_{A'}$ et la transition entre ces deux états correspond bien à une expression qui décrit le langage reconnu par A . Nous ne le démontrerons pas ici, mais cela pourrait se faire par induction sur le nombre d'états de $Q_{A'}$, en montrant qu'après chaque suppression d'état, le langage

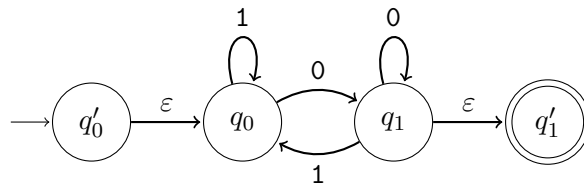
reconnu par l'automate a été préservé.

4.4.1 Exemple

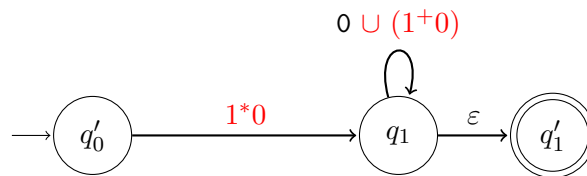
Voyons maintenant un exemple concret. Supposons que nous avons l'AFD suivant (en l'occurrence, qui reconnaît les nombres pairs sur l'alphabet $\Sigma = \{0, 1\}$). Nous voulons le transformer en expression régulière qui décrit le même langage.



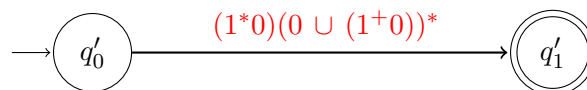
Ajout des nouveaux états (initial et final) :



Nous allons commencer par supprimer q_0 (choix arbitraire). Cet état sert d'intermédiaire de deux manières différentes : (1) de q'_0 vers q_1 , avec pour expression correspondante $\varepsilon 1^* 0 = 1^* 0$, on ajoute donc une transition $1^* 0$ de q'_0 à q_1 ; (2) de q_1 vers lui-même, avec pour expression $1 1^* 0 = 1^+ 0$, on ajoute donc cette expression (via une union) sur la boucle existante de q_1 . Cela donne l'automate :



Supprimons maintenant q_1 . Cet état sert d'intermédiaire de q'_0 à q'_1 avec l'expression $(1^* 0)(0 \cup (1^+ 0))^* \varepsilon$. On ajoute donc l'expression $(1^* 0)(0 \cup (1^+ 0))^*$ de q'_0 à q'_1 . Cela donne l'automate :



À ce stade, nous avons terminé, et l'expression sur la dernière transition doit correspondre au langage qui était reconnu par l'automate de départ. Autrement dit, $(1^*0)(0 \cup (1^+0))^*$ décrit le langage des nombres pairs. Ce n'est pas évident à première vue, car cette expression est différente de celle que l'on connaît : $(0 \cup 1)^*0$. Mais elles décrivent bien le même langage. Tout d'abord, observons que $(0 \cup (1^+0)) = (1^*0)$, donc $(1^*0)(0 \cup (1^+0))^* = (1^*0)(1^*0)^* = (1^*0)^+$. Autrement dit, le mot vide est exclu, et tout facteur comportant des 1 doit être suivi d'au moins un 0, ce qui correspond bien aux nombres pairs.

4.4.2 Conclusion

Au cours de ces dernières séances, nous avons montré que tout AFN peut être transformé en un AFD équivalent (les AFDs étant, par ailleurs, des cas particuliers d'AFNs) ; que toute expression régulière (ER) peut être transformée en AFN ; et enfin que tout AFD peut être transformé en ER. Nous avons donc montré, transitivement, que tous ces outils ont la même expressivité : ils décrivent chacun à sa façon des langages réguliers, ni plus ni moins.

