

Langages rationnels : Exemples et applications

Les instances d'un problème donné peuvent être représentées par une chaîne de caractères :

- nombres (séquences de bits)
- image (suite des intensités de chaque point)
- son (suite de nombres représentant les tonalités)

Ainsi, beaucoup de problèmes se réduisent à décider l'appartenance d'un élément à un ensemble. Si l'on peut décrire l'ensemble à l'aide d'une expression rationnelle, la théorie des langages rationnels intervient pour résoudre ce problème. Celle-ci est très liée à la théorie des automates finis.

I Généralités [Roz, Wol, HOP]

1) Langage sur un alphabet

Définition 1: Un alphabet Σ est un ensemble non vide fini de symboles, appelés lettres.

exemple : $A = \{0, 1\}$; $B = \{a, b, c, e\}$ sont des alphabets.

Définition 2: Un mot sur Σ est une séquence fini de lettres de Σ .

exemple :

- ϵ = mot vide et un mot sur tout alphabet
- $w_1 = 001$; $w_2 = 101$ sont des mots sur A .
- $w_3 = abce$; $w_4 = aec$ sont des mots sur B .

Notations :

- $|w|$ est la longueur du mot, i.e. le nombre de symboles qui le compose. ($|\epsilon| = 0$)
- Σ^* est l'ensemble des mots sur Σ .
- $w_1 w_2 = w_1 \cdot w_2$ est la concaténation de w_1 et w_2 : les symboles de w_1 suivis de ceux de w_2 .

Définition 3: Un langage L sur Σ est un ensemble de mots sur Σ . C'est une partie de Σ^* .

2) Langages et expressions rationnelles

Définition 4: On définit les expressions rationnelles de manière inductive :

- Cas de base : ϕ , ϵ , $a \in \Sigma$ sont des expressions rationnelles

- Induction : Soient E_1, E_2 deux expressions rationnelles.

- $E_1 + E_2 = E_1 \cup E_2$ est une expression rationnelle : clôture par union
- $E_1 \cdot E_2 = \{w_1 w_2 \mid (w_1, w_2) \in E_1 \times E_2\}$ est une expression rationnelle : clôture par concaténation
- E_1^* est une expression rationnelle : clôture par passage à l'étoile.

Définition 5: On définit par le même procédé d'induction le langage associé à une expression rationnelle E noté $L(E)$:

- Cas de base : $L(\phi) = \phi$; $L(\epsilon) = \{\epsilon\}$; $L(a) = \{a\}$ $a \in \Sigma$

- Induction : si E_1 et E_2 sont deux expressions rationnelles, alors $L(E_1 + E_2) = L(E_1) \cup L(E_2)$, $L(E_1 \cdot E_2) = L(E_1) \cdot L(E_2)$, $L(E_1^*) = L(E_1)^*$.

exemples : $L((0+1) \cdot 1 \cdot (0+1) \cdot (0+1)^*)$ est l'ensemble des mots de longueur impaire sur A .
 $L(((abc)^* + (cab)^*) + (b \cdot (c \cdot b)^*) + (c \cdot (b \cdot c)^*))$ est l'ensemble des mots où les lettres sont alternées.

Définition 6: Un langage est dit rationnel s'il existe une expression rationnelle E telle que $L = L(E)$. On notera $\text{Rat}(\Sigma)$ l'ensemble des langages rationnels sur Σ .

Remarques :

- $\text{Rat}(\Sigma)$ est la plus petite classe de langages stable par \cup , \cdot , * et contenant $\{\phi, \epsilon, a\}$ lorsque $a \in \Sigma$.
- Plusieurs expressions peuvent décrire un même langage :
 $L((ab)^* + (ba)^* + a(ba)^* + b(ab)^*) = L((\epsilon+ab)(ab)^*(\epsilon+ba))$
ce qui invalide la définition suivante.

Définition 7: E_1 et E_2 sont dites équivalentes si elles sont associées au même langage : $L(E_1) = L(E_2)$.

Proposition 2: L'ensemble des sous-mots d'un langage L quelconque est rationnel

Proposition 3: Il existe des langages non rationnels

Problèmes :

- Peut-on déterminer si deux expressions rationnelles sont équivalentes ?
- Pour un mot sur Σ et un langage sur Σ , peut-on déterminer si $w \in L$?

I Liens avec les automates finis [Car; Sab; Hop]

Les automates finis donnent une autre définition des langages rationnels et sont un outil parfois plus facile à manier pour certaines preuves.

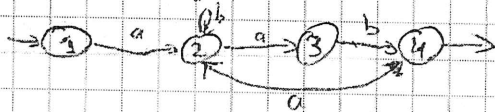
1) Préliminaires

Définition 10: Un automate fini \mathcal{A} sur Σ est la donnée d'un quintuplet $(Q, \Sigma, \delta, I, F)$, avec Q ensemble fini, $I, F \subseteq Q$, $\delta \subseteq Q \times \Sigma \times Q$. Les éléments de Q sont appelés les états, ceux de I les états initiaux, ceux de F les états finaux. Les éléments (q, a, q') de δ sont les transitions de l'automate.

Remarques:

- Généralement, on prendra $Q = \{1, \dots, n\}$, $n \in \mathbb{N}^*$ car le nom des états n'a pas sur l'étude de l'automate.
- On représente un automate \mathcal{A} par un graphe avec arêtes étiquetées et des étiquettes sur les états initiaux et finaux.

Exemple: $\mathcal{A} = (\{1, 2, 3, 4\}; \{a, b\}; \{(1, a, 2), (2, a, 3), (3, b, 4), (4, a, 2)\}; \{1\}, \{4\})$ est représenté par:



Définition 11: Un chemin dans \mathcal{A} est une suite de transitions consécutives $(q_0, a_1, q_1, a_2, q_2, \dots, a_n, q_n)$ où q_0 et q_n sont les états de départ et d'arrivée.

• Un chemin est dit **acceptant** si $q_0 \in I$ et $q_n \in F$.

• Si $q_0 \stackrel{a_1 \dots a_n}{\rightarrow} q_n$ est un chemin, on dit que $a_1 \dots a_n$ est l'étiquette du chemin.

Définition 12: Un mot est dit **accepté** par \mathcal{A} si il est l'étiquette d'un chemin acceptant dans \mathcal{A} . On note $L(\mathcal{A})$ l'ensemble des mots acceptés par \mathcal{A} .

Exemple: $\mathcal{A} = \rightarrow 1 \xrightarrow{a} 2 \xrightarrow{b} 2 \xrightarrow{a} 2$ $L(\mathcal{A}) = a(bb^*a)^*(a+b)$.
Le chemin $c = (1, a, 2) (2, b, 2) (2, a, 2) (2, b, 2)$ est acceptant.

Définition 13: Un automate est dit **émondé** si par tout état passe un chemin acceptant.

Définition 14: Un automate est dit **normalisé** si $|I| = |F| = 1$ et que'il n'y a pas de transitions de la forme (q, a, i) ou (f, a, q) si $I = \{i\}$ et $F = \{f\}$.

Proposition 15: Pour tout automate \mathcal{A} , il existe un automate normalisé \mathcal{A}' tel que $L(\mathcal{A}) = L(\mathcal{A}') \setminus \{\epsilon\}$.

2) Langages rationnels et reconnaissables
Lemme 16: (Ardon) L'équation $X = AX + B$ d'inconnue X avec A, B deux langages sur Σ fixé admet pour solution:
• Si $\epsilon \notin A$, $X = A^*B$
• Si $\epsilon \in A$, \exists langage $P \subseteq \Sigma^*$, $X = A^*(B+P)$

Théorème 17: (Kleene) [DEV]

Un langage L est rationnel si et seulement si il est reconnaissable par automate fini, c'est-à-dire qu'il existe \mathcal{A} tel que $L = L(\mathcal{A})$.

Lemme de l'étoile 18: Soit $L \in \text{Rat}(\Sigma)$, alors $\exists N \in \mathbb{N}, \forall p \in \Sigma^*$

$$(|p| \geq N) \Rightarrow (\exists (u, v, w) \in \Sigma^3, \forall \epsilon, f = uv^\epsilon w)$$

$$|uv| \leq N \text{ et } uv^*u \in L$$

Lemme de l'étoile par bloc 19: Soit $L \in \text{Rat}(A)$, alors $\exists N \in \mathbb{N}, \forall f \in L$, pour toute

$$f \text{ décomposition } f = uv_1 \dots v_k w, u, w \neq \epsilon, \exists 0 \leq j < k \leq N,$$

$$uv_1 \dots v_j (v_{j+1} \dots v_k)^* v_{j+1} \dots v_k w \in L.$$

Applications: $\{a^*b^* \mid a \neq b\}$ n'est pas rationnel.

• L'ensemble des palindromes n'est pas rationnel.
Cela donne une preuve à la proposition 3 qui est difficile à montrer sans les automates.

• $\text{Rat}(\Sigma)$ est stable par \cup, \cap, \setminus .

Proposition 20: L'appartenance d'un mot à un langage rationnel est décidable.

Corollaire 21: L'équivalence de deux expressions rationnelles est décidable.

3) Quotients et automates minimaux [Car]

Définition 22: Soit $L \subseteq \Sigma^*$ un langage sur Σ . Le quotient à gauche de L par un mot $u \in \Sigma^*$ est le langage $u^{-1}L = \{v \in \Sigma^* \mid uv \in L\}$.
Le quotient à gauche de L par $K \subseteq \Sigma^*$ est $K^{-1}L = \bigcup_{k \in K} k^{-1}L$.

Proposition 23: ~~Soit $L \subseteq \Sigma^*$~~ pour tout mot u, v et w et tous langages K et L , on a les relations suivantes:

$$(\varepsilon L) = \varepsilon(L)$$

- $w^{-1}(KL) = w^{-1}K \cup w^{-1}L$
- $a^{-1}(KL) = (a^{-1}K)L \cup \varepsilon(K)a^{-1}L$
- $w^{-1}(KL) = (w^{-1}K)L \cup \sum_{uv=w} \varepsilon(u^{-1}K/v^{-1}L)$
- $a^{-1}(L^*) = (a^{-1}L)^*$
- $w^{-1}L^* = \sum_{uv=w} \varepsilon(u^{-1}L^*)(v^{-1}L)^*$
- $(uv)^{-1}L = v^{-1}(u^{-1}L)$

Définition 24: Un automate $A = (Q, \Sigma, \delta, I, F)$ est déterministe si:

- $|I| = 1$
- $(p, a, q), (p, a, q') \in \delta \Rightarrow q = q'$

Lemme 25: Soit $A = (Q, \Sigma, \delta, \{i\}, F)$ un automate déterministe tel que $L(A) = L$.

Pour tout mot $u \in \Sigma^*$, on a $u^{-1}L = \{v \mid q \xrightarrow{v} f \text{ avec } f \in F\}$

Définition 26: Soit $L \in \text{Rat}(\Sigma)$. L'automate minimal de L est un automate $A_L = (Q, \Sigma, \delta, I, F)$ où $Q = \{a^{-1}L \mid a \in \Sigma^*\}$, $I = \{\varepsilon\}$, $F = \{a^{-1}L \mid a \in L\}$, $\delta = \{a^{-1}L \xrightarrow{a} (a^{-1}L) \mid a \in \Sigma\}$.

Lemme 27: $L(A_L) = L$

Corollaire 28: $(L \in \text{Rat}(\Sigma)) \iff \exists L \text{ admet un nombre fini de quotient à gauche}$

III) Applications des langages rationnels [BBC], [Car]

1) Recherche de motifs

On veut rechercher un mot $u \in \Sigma^*$ dans un texte. Concrètement, on veut reconnaître $\Sigma^* u \Sigma^*$

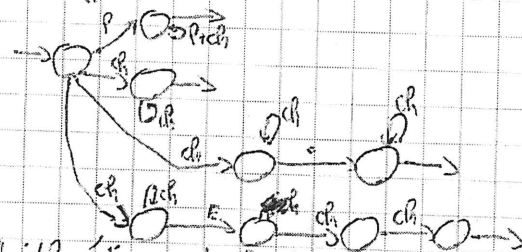
→ biologie : recherches de motifs dans une chaîne de nucléotides.
→ informatique : traitement de texte, recherche ou remplacement d'une chaîne de caractères par une autre.
Résolution du problème par l'algorithme de Knuth-Morris-Pratt.

2) Analyse lexicale

Il s'agit de la première étape d'un compilateur avant l'analyse syntaxique. On décompose le programme source du compilateur en unités lexicales, c'est-à-dire que l'on regroupe les chaînes qui codent le même élément du programme (entier, mot du langage : begin, else, ...).

exemple : On veut reconnaître $\{ (ch+P)^* + ch \cdot ch^* + ch \cdot ch^* \cdot ch^* + ch \cdot ch^* \cdot ch \}$ où $ch \in \{a, b, \dots, z\}$ et $\{ \{ \}, \{ \}$.

Il suffit de minimiser :



3) Arithmétique de Presburger. [DEV] [Car]

On se place dans le cadre de la logique du premier ordre sur \mathbb{N} . On considère les opérations $+$ et $=$ uniquement. On appelle ceci l'arithmétique de Presburger.

Définition 29: On dit que une théorie logique est décidable si le problème de savoir si une formule close est vraie est décidable.

Théorème 30 (Presburger): La théorie du premier ordre des entiers munis de l'addition est décidable.

Ref: [BBC]: Beauquier, Harel, Orléanne [Car] Caron [Wol] Wolper [Sak] Sakurai
[Hop] Hopcroft - Ullman [Roz] Reyzens - Salomaa