

[B] p. 296

Ex 16: $*L = (a+b)^*aba$ est reconnue par: $\rightarrow \textcircled{1} \xrightarrow{a} \textcircled{2} \xrightarrow{b} \textcircled{3} \xrightarrow{a} \textcircled{4}$

$*L = (b+ab^*a)^*$ est reconnue par l'automate de l'exemple 11.

Lemme 17: Arden: Soient K et L deux langage, et l'équation

$X = KX + L$. Alors $\cdot \epsilon \notin K$: il existe une unique solution $X = K^*L$

$\cdot \epsilon \in K$: les solutions sont de la forme $X = K^*(L+Y)$ où $Y \in A^*$

Ex 18: L'automate $\textcircled{1} \xrightarrow{a} \textcircled{2} \xrightarrow{b} \textcircled{1}$ reconnaît $L = (b^*a)^*$

* L'automate $\textcircled{1} \xrightarrow{a} \textcircled{2} \xrightarrow{a} \textcircled{3} \xrightarrow{b} \textcircled{2}$ reconnaît $(a+ab^*a)^*b^*$

Théorème 19: Kleene: Un langage est rationnel si et seulement si il est reconnu par un automate fini.

II - Propriétés des langages rationnels

1) Quotients d'un langage

Prop et déf. 20: tout automate fini est équivalent (ie reconnaît le même langage) à un automate fini déterministe complet, ie $|I|=1$ et $\forall a \in A, \forall q \in Q, |\delta(q,a)|=1$

Déf 21: Soit L un langage; Soit u un mot. Le quotient à gauche de L par u est le langage $u^{-1}L = \{v \mid uv \in L\}$

Prop 22: Soit $a \in A$; Soient K et L deux langages:

$a^{-1}(K+L) = a^{-1}K + a^{-1}L$; $a^{-1}(KL) = (a^{-1}K)L + \epsilon(K)a^{-1}L$

($\epsilon(K) = \epsilon \in K, 0$ sinon); $a^{-1}(L^*) = (a^{-1}L)^*$

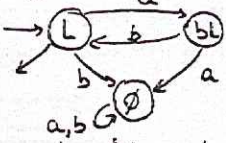
Ex 23: $a^{-1}(ab^*a+b)^* = b^*a(ab^*a+b)^*$

Déf 24: Soit L un langage rationnel. L'automate minimal reconnaissant L est: $\mathcal{A}_L = (Q, A, \delta, I, F)$ avec:

$Q = \{u^{-1}L \mid u \in A^*\}$, $I = \{L\}$, $F = \{u^{-1}L \mid u \in L\}$

$\delta(u^{-1}L, a) = (ua)^{-1}L$

Ex 25: $L = (ab)^*$



Prop 26: Un langage est rationnel si et seulement si il admet un nombre fini de quotients à gauche.

Prop 27: L'automate minimal de L est l'automate parmi les

automates déterministes complets ayant le moins d'états qui reconnaît L.

[C] p. 5

[B] p. 317

L'automate minimal peut être construit à partir d'algorithmes qui s'appuient sur:

Déf. 28: Congruence de NÉRODE: Soit $\mathcal{A} = (Q, A, \delta, I, F)$ un automate déterministe complet. $\forall q, q' \in Q$:

$q \sim q' \iff \forall w \in A^*, (q \cdot w \in F \iff q' \cdot w \in F)$

Prop 29: Soit \mathcal{A} un automate déterministe complet. L'automate minimal équivalent est isomorphe à \mathcal{A}/\sim où \sim est la congruence de NÉRODE de \mathcal{A} .

Ex 30: $\textcircled{1} \xrightarrow{a} \textcircled{2} \xrightarrow{b} \textcircled{3} \xrightarrow{a} \textcircled{4} \xrightarrow{b} \textcircled{3}$ équivaut à l'automate minimal suivant: $\textcircled{1} \xrightarrow{a} \textcircled{2} \xrightarrow{b} \textcircled{3}$

Rq: L'automate minimal est unique à isomorphisme près

2) lemme de l'étoile

Thm 31: Lemme de l'étoile: Soit L un langage rationnel.

1) $\exists N \in \mathbb{N}, \forall u \in L, |u| \geq N, \exists u = vt^k w, v \neq \epsilon, \forall n \in \mathbb{N}, vt^{n+k} w \in L$

2) $\exists N \in \mathbb{N}, \forall u \in L, \forall u = v_1 u_1 w, |u_1| \geq N: \exists u_2 = vt^k w, v \neq \epsilon, t^k \in A^*$

tg: $\forall n \in \mathbb{N}, v_1 vt^{n+k} w w_2 \in L$

3) $\exists N \in \mathbb{N}, \forall u \in L, u = v u_1 \dots u_k w, \text{ les } u_i \neq \epsilon, \exists 0 \leq j < k$

tg: $\forall u_1 \dots u_j (u_{j+1} \dots u_k)^* u_{j+1} \dots u_k w \in L$

Rq: Ce résultat sert à montrer que des langages ne sont pas rationnels.

Ex 33: $L = \{a^n b^n \mid n \in \mathbb{N}\}$ n'est pas rationnel

$\cdot L = \{a^n b^n \mid n \in \mathbb{N}\}$ satisfait 1) mais pas 2) donc n'est pas rationnel

$\cdot L = \{a^i b^i a^j b^j \mid i, j \in \mathbb{N}\}$ satisfait 1), 2) et 3) mais n'est pas rationnel!

La réciproque est fautive!

Rq 33: Une réciproque au lemme de l'étoile

Soit L un langage.

L et $A^* \setminus L$ vérifient le lemme de l'étoile par blocs (3)

si et seulement si L et $A^* \setminus L$ sont rationnels.

[C] p. 17

[C] p. 18

[S] p. 125

[S] p. 78

[S] p. 79

[S] p. 80

[S] p. 80

[S] p. 127

III - Applications.

[S] Les langages rationnels ont de nombreuses applications, notamment liés à l'analyse lexicale des langages de programmation et à la recherche de motifs.

1) Recherche de motif

[B] On considère un texte $t = t_1 \dots t_k$; on cherche une occurrence du motif $x = x_1 \dots x_m$ dans t .

[B] a. Algorithme naïf:

On compare le motif à chaque facteur de t de longueur m ; on s'arrête si on trouve une occurrence de x .

Complexité en pire cas: $O(mk)$ lorsque n petit devant k .

[B] b. Algorithme de Knuth-Morris et Pratt et automate de Simon.

On optimise le décalage dans la comparaison des motifs en considérant le plus grand suffixe u qui est aussi préfixe du motif x .

On peut également construire l'automate des occurrences et l'utiliser dans la recherche du motif

Complexité en pire cas: $O(m+k)$

Rq: Cette recherche de motif est utilisée pour rechercher les différents lexèmes dans le cadre de l'analyse lexicale.

[FB] 2) Séparation par automates.

Soient L et K deux langages finis. Soit $k \in \mathbb{N}$.

Existe-t-il un automate déterministe à k états tel que $L \subseteq L(A)$ et $T_n L(A) = \emptyset$?

Thm 34: Ce problème est NP-complet.

Références

[C]: Carton: Langages formels

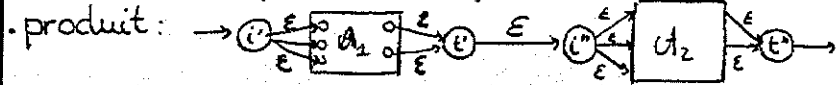
[S]: Sakharovitch: Éléments de théorie des automates

[B]: Beauquier-Berstel-Christienne: Éléments d'algorithmique

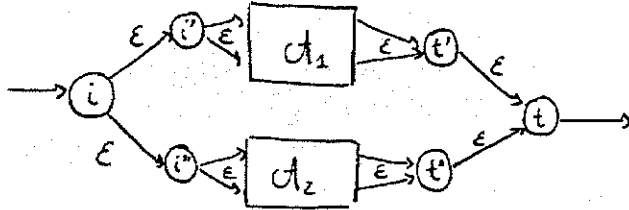
[FB]: Floyd-Beigel: Language of machines

Annexe 1 : Construction de Thompson :

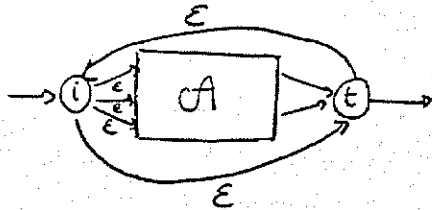
• $\rightarrow \text{O} \xrightarrow{a} \text{O} \rightarrow$ (cas de base)



• union :



• étoile :



Recherche de motif

Entrée : Un texte $t = t_1 \dots t_n$

Un motif $\pi = \pi_1 \dots \pi_m$

Sortie : des positions auxquelles apparaît le motif π dans t
(S'il n'apparaît pas, rien)

I Algorithme de Simon :

On construit l'automate du motif

lorsqu'on se trouve à l'état q , c'est qu'on a lu les q premières lettres du motif



On doit construire l'automate en pré-traitement :

def: Soit σ la fonction suffixe définie par :

$$\sigma : A^* \rightarrow [0, m]$$

$$u \mapsto \sigma(u) = \max \{ k : \pi_1 \dots \pi_k \sqsupseteq u \} \quad (\pi_1 \dots \pi_m \text{ est le motif})$$

(σ correspond à la longueur du plus long suffixe de u qui est préfixe de π .)

Rq: ϵ est suffixe de tout mot donc σ est bien défini

($\{k : \pi_1 \dots \pi_k \sqsupseteq u\}$ est non-vide avec la convention $\pi_1 \dots \pi_k = \epsilon$ lorsque $k=0$)

On définit l'automate du motif $\mathcal{A} = (Q, A, \delta, I, F)$ par

• $Q = [0, m]$, $I = \{0\}$, $F = \{m\}$

• $\forall q \in Q, \forall a \in A, \delta(q, a) = \sigma(\pi_1 \dots \pi_q a)$

(le plus grand préfixe de π qui soit suffixe

Rq: cet automate est déterministe complet)

Algorithme:

Simon (k, n):

$q \leftarrow 0$

pour $i = 1$ à n :

$q \leftarrow \delta(q, t_i)$

si $q = m$ alors imprimer "position $i - m$ ".

L'algorithme termine (facile) et est correct (ça se montre)

On peut montrer que sa complexité est en $O(n \cdot m^3 |A|)$

II. L'algorithme de Knuth-Morris et Pratt:

On veut améliorer l'algorithme précédent avec un prétraitement moins coûteux.

Pour cela, on va calculer une fonction préfixe π (qui dépendra du motif), qui exprimera les correspondances entre le motif et ses propres décalages.

Préfixe (n, n_m):

$\pi(1) \leftarrow 0$

$k \leftarrow 0$

pour $q = 2$ à m :

tant que $k > 0$ et $n_{k+1} \neq n_q$

$k \leftarrow \pi(k)$

si $n_{k+1} = n_q$

$k \leftarrow k + 1$

$\pi(q) \leftarrow k$

retourner π .

KMP (n, n_m, t_1, \dots, t_n):

$\pi \leftarrow \text{Préfixe}(n, \dots, n_m)$

$q \leftarrow 0$

pour $i = 1$ à n :

tant que $q > 0$ et $n_{q+1} \neq t_i$

$q \leftarrow \pi(q)$

si $n_{q+1} = t_i$

$q \leftarrow q + 1$

si $q = m$

imprimer "position $i - m$ ".

Cet algorithme termine (ok) et est valide (long).

Complexité:

On compte le nombre de comparaisons de lettres:

• Dans KMP:

On compte le nombre de tests positifs et négatifs.

• Si le test est positif: On incrémente i

Où $1 \leq i \leq n$ donc \rightarrow au plus n tests

• Si le test est négatif, i inchangé mais q diminue strict^T
donc $i - q$ augmente strictement à chaque test négatif
(+ reste inchangé à chaque test positif)

donc ($i - q = 0$ au début)

$\#(\text{tests} < 0) \leq i - q$ (à la fin)

$\leq i \leq n$

Donc la complexité est en $O(m)$.

• Dans Préfixe:

idem : au plus m tests positifs

• $\#(\text{test} < 0) \leq q - k$

$\leq m$

la complexité est en $O(m)$.

(\uparrow strict quand < 0 ,
 \uparrow quand > 0)

Au final, l'algorithme est en $O(m+m)$.

Automate quotient = automate minimal

Prop L'automate minimal de L est l'automate ayant le moins d'états parmi les automates déterministes complets.

donc 1) $\mathcal{A}_L = (Q, A, \delta, i, F)$ est déterministe et complet.

$$L = \{L\} \quad \text{et} \quad \forall u, v \in A^*, \forall a \in A \quad |\delta(u, a)| = |\delta(v, a)| = 1$$

2) il existe un nombre d'états minimal.

On commence par démontrer le lemme suivant:

Lemme: Soit L un langage rationnel et $\mathcal{A} = (Q, A, \delta, i, F)$ un automate déterministe complet accessible reconnaissant L .

Pour $q \in Q$, on pose $q(L) = \{w \in A^* / \delta'(q, w) \in F\}$.

$$\text{Alors } Q = \{q(L), q \in Q\}$$

Déf: un état $q \in Q$ est dit accessible si $\exists u \in A^*, \delta'(i, u) = q$.

\mathcal{A} est accessible si tout état de \mathcal{A} est accessible.

donc du lemme a) $\bigcap_{q \in Q} q(L) \subseteq \{q(L), q \in Q\}$.

Soit $u \in A^*$, comme \mathcal{A} est complet $\exists s \in Q, \delta'(i, u) = s$

$$\text{donc } w \in u^{-1}L \Leftrightarrow uw \in L \Leftrightarrow \delta'(i, uw) \in F \Leftrightarrow \delta'(s, w) \in F \Leftrightarrow w \in q(L)$$

$$\text{donc } \forall q, q(L) = u^{-1}L \in Q \quad \text{donc } Q = \{u^{-1}L, u \in A^*\} \subseteq \{q(L), q \in Q\}$$

b) $\bigcap_{q \in Q} q(L) \subseteq Q$

Soit $q \in Q$, comme \mathcal{A} est accessible $\exists u \in A^*, \delta'(i, u) = q$

on a comme précédemment $u^{-1}L = q(L)$ donc $\{q(L), q \in Q\} \subseteq Q$.

$$\text{Donc } Q = \{q(L), q \in Q\}$$

□

Retour à la preuve de la propriété:

d'après le lemme, tout automate déterministe complet accessible reconnaissant L

vérifie que $Q = \{q(L), q \in Q\}$ ie $|Q| \leq |Q|$.

de plus comme tout automate déterministe complet peut être rendu accessible en ajoutant un état initial, on a donc que \mathcal{A}

car l'automate déterministe complet reconnaissant L ayant le minimum d'états. (car les états non-accessibles ne font partie d'aucun chemin accepté par l'automate). \square .

Prop | Un langage est rationnel ssi il admet un nombre fini de quotients à gauche

\Rightarrow Soit L un langage rationnel, alors il existe un automate fini reconnaissant L noté \mathcal{A} , et est équivalent à un automate fini déterministe complet \mathcal{A}' .

~~Quelle que réduite le nombre d'états de \mathcal{A} , on peut supposer qu'il est minimal. Or à partir d'après le lemme utilisé dans la preuve~~

Comme \mathcal{A}' est l'automate minimal de L , il a le nombre minimal d'états parmi les automates déterministes complets. Donc $|\mathcal{Q}| = |\{u^{-1}L, u \in A^*\}| \leq |\text{états de } \mathcal{A}'| < \infty$

Donc L admet un nombre fini de quotients à gauche.

\Leftarrow Soit L un langage ayant un nombre fini de quotients à gauche alors l'automate \mathcal{A}' est un automate fini reconnaissant L donc L est rationnel. \square

Vérification que $L(\mathcal{A}') = L$.

On montre par récurrence sur la longueur du mot $u \in A^*$ que pour tout état $x \in \mathcal{Q}$,

$$\delta(x, u) = u^{-1}x.$$

• $n=1$: $u = a \in A$ Soit $x \in \mathcal{Q}$ alors $\exists v \in A^*, x = v^{-1}L$

$$\delta(x, a) = \delta(v^{-1}L, a) = (va)^{-1}L = \{w \in A^*, va \cdot w \in L\} = \{w \in A^*, a \cdot w \in v^{-1}L\} = \{w \in A^*, a \cdot w \in x\} = a^{-1}x$$

• on suppose le résultat vrai pour $k \leq n$ et on considère $u \in A^*, |u| = n+1$

on a donc $u = au'$ avec $a \in A$ et $u' \in A^*$ tel que $|u'| = n$

$$\text{Soit } x \in \mathcal{Q} \quad \delta(x, u) = \delta(x, au') = \delta(\delta(x, a), u') = \delta(a^{-1}x, u') \stackrel{\text{H.R. 1}}{=} u'^{-1}(a^{-1}x) \stackrel{\text{H.R. n}}{=} u^{-1}x$$

$$\text{or } u'^{-1}(a^{-1}x) = \{w \in A^*, u'w \in a^{-1}x\} = \{w \in A^*, au'w \in x\} = \{w \in A^*, u'w \in x\} = u^{-1}x$$

$$\forall u \in A^*, \forall x \in \mathcal{Q}, \delta(x, u) = u^{-1}x.$$

$$w \in L(\mathcal{A}') \stackrel{\text{I}=\{L\}}{\Leftrightarrow} \delta(L, w) \in F \stackrel{\text{Reconnaître}}{\Leftrightarrow} w^{-1}L \in F \stackrel{\text{dét. de } F}{\Leftrightarrow} w \in L.$$

$$\text{Ainsi } \underline{L(\mathcal{A}') = L.}$$