

On suppose connues les définitions de :
 - automate fini déterministe (complet)
 - automate fini non déterministe
 - automate fini non déterministe avec ϵ -transitions
 et l'équivalence entre tous ces automates.
 Par la suite, Σ désignera un alphabet fini.

1) Langages rationnels et leurs caractérisations

La classe des langages rationnels est le premier niveau de la hiérarchie de Chomsky et peut être définie de beaucoup de manières différentes. En voici quelques unes :

1) Une première définition [Car] p.53

Définition [Langages rationnels]

La classe $\text{Rat}(\Sigma^*)$ des langages rationnels (sur Σ) est la plus petite famille de langages telle que :

- * $\emptyset \in \text{Rat}(\Sigma^*)$
- * $\forall a \in \Sigma \quad \{a\} \in \text{Rat}(\Sigma^*)$ (finie)
- * $\text{Rat}(\Sigma^*)$ est close par union, par produit (i.e. concaténation) et par l'étoile.

Exemples : $\Sigma \in \text{Rat}(\Sigma^*)$ car il s'écrit $\Sigma = \bigcup_{a \in \Sigma} \{a\}$

- Les mots de longueur paire forment un langage rationnel car égal à $(\Sigma \Sigma)^*$.
- Le langage des mots contenant le facteur $abba$ est rationnel car égal à $\Sigma^* \{a-b-b-a\} \Sigma^*$.

2) Expressions rationnelles [Sak] p.133

Définition [Expressions rationnelles] :

Une expression rationnelle sur Σ est une

formule obtenue inductivement de la manière suivante (on suppose que Σ ne contient pas les lettres $\emptyset, \epsilon, +, \cdot, ^*$ ($\epsilon, \emptyset, \cdot, ^*$)).
 • \emptyset, ϵ et a , pour tout a dans Σ sont des expressions rationnelles.

• Si E et F sont deux expressions rationnelles alors $(E+F), (E \cdot F)$ et $(E)^*$ sont des expressions rationnelles.

(Pour une meilleure lisibilité on enlèvera les parenthèses qui n'impliquent pas d'ambiguïté au niveau des expressions en imposant des priorités entre les opérateurs)

Ex : $(a+b)^* a b b a (a+b)^*$

Définition À chaque expression rationnelle E on fait correspondre un langage noté $L(E)$ tel que :

- * $L(\emptyset) = \emptyset$ $L(\epsilon) = \{\epsilon\}$ (not vide) $L(a) = \{a\}$ pour $a \in \Sigma$
- * $L(E+F) = L(E) \cup L(F)$ $L(E \cdot F) = L(E) \cdot L(F)$
- et $L(E^*) = (L(E))^*$.

Proposition Un langage est rationnel ssi il est dénoté par une expression rationnelle. Il faut distinguer expression = objet syntaxique et le langage qu'elle dénote = sémantique.

Ex : $(a+b)^*$ et $(a^*b)^+ a^*$ sont des expressions différentes mais les langages qu'elles dénotent sont les mêmes.

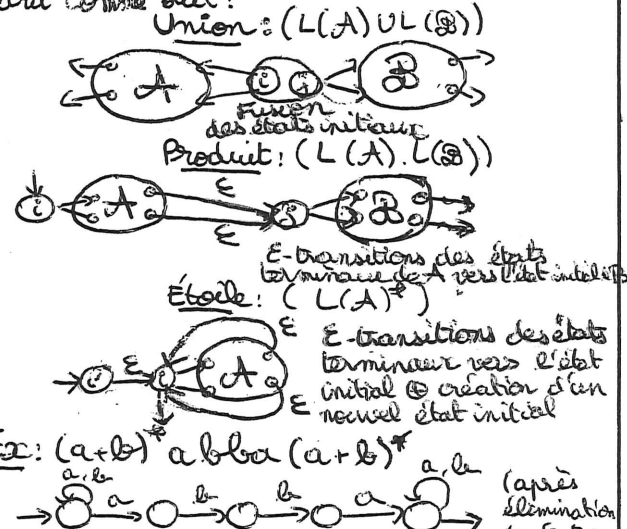
Définition Deux expressions sont dites équivalentes si elles dénotent le même langage.

Problème : Comment déterminer si deux expressions données sont équivalentes ?

3) Théorème de Kleene

Théorème de Kleene : Un langage L est rationnel ssi il existe un automate (fini) A tel que $L = L(A)$.

Des expressions aux automates [Sak] p.155
 La transition peut se faire de manière standard comme suit :



À noter : On peut directement obtenir un automate déterministe grâce à la méthode de Brzozowski (automate des dérivées)

Des automates aux expressions [Car] p.35

Algorithme de McNaughton-Yamada :

On étiquette Q , ensemble des états d'un automate A , par $1, 2, \dots, n$. Pour $(q, q') \in Q^2$ on définit alors $L_{q,q'}^{(k)} = \{a_1 \dots a_n \mid q \xrightarrow{a_1} \dots \xrightarrow{a_n} q' \text{ et } 1, \dots, n-1, n\}$

On a les formules :
 $L_{q,q'}^0 = \{a \mid q \xrightarrow{a} q'\} \cup \{\epsilon \mid q = q'\}$
 $L_{q,q'}^{(k+1)} = L_{q,q'}^{(k)} + L_{q,r+1}^{(k)} (L_{r+1,r+1}^{(k)})^* L_{r+1,q'}$

Et d'autre part : $L(A) = \bigcup_{i,j} L_{i,j}^{(n)}$
 Élimination de Gours : Soit l'équation $X = AX + B$ où l'inconnue X est un langage.
 1. Si $\epsilon \notin A, X = A^*B$.

2, $\exists \alpha \in \mathbb{A}$ X est de la forme $X = A^*(B)^*$
 où $P \subset \Sigma^*$.

Principe de la méthode: On appelle pour $p \in \mathbb{A}$ $L_p = \{w \mid \exists \beta \in \mathbb{A} \ p \xrightarrow{w} \beta\}$
 On a donc $L(A) = L_i$, où i est l'état initial.
 On résout alors le système

$$\begin{cases} L_p = \sum_{q \rightarrow p} a L_q + \epsilon \end{cases}$$

si ϵ est un état final

on résoudra par substitutions successives les équations grâce au lemme d'Arden.

4) Autres caractérisations

a) Automate boustrophédon [Sak, p 155]
 Un automate boustrophédon est une machine de Turing dont la tête qui se déplace sur la bande est seulement une tête de lecture (et pas d'écriture).

Théorème: Un langage est rationnel ssi il est reconnu par un automate boustrophédon.

(La preuve requiert la finitude des quotients à gauche pour un langage régulier.)

b) Grammaires linéaires [Vidp, p 222]

Une grammaire est dite linéaire à droite (resp à gauche) ssi ces productions sont de la forme

$$A \rightarrow wB \quad \text{ou} \quad A \rightarrow w \quad \text{avec } A, B \text{ non terminaux, } w \text{ une chaîne de terminaux}$$

Proposition Un langage est rationnel ssi il est engendré par une grammaire linéaire (à gauche ou à droite)

Ex: $(a+b)^* abba (a+b)^*$ est engendré par

$$S \rightarrow aS \mid bS \mid abbaT$$

$$T \rightarrow aT \mid bT \mid \epsilon$$

(l'axiome est S)

II Conséquences du théorème de Kleene

1) Propriétés de clôture [Ry, sel p 33]

Propriété: la classe des langages rationnels est stable par:

- union (finie)
- intersection (finie)
- complémentation
- concaténation (finie)
- étoile
- image miroir

Théorème: L'image d'un langage rationnel par un transducteur fini reste un langage rationnel.

Corollaire: Pour L un langage rationnel, les ensembles suivants sont rationnels:

- préfixes(L) = $\{x \mid \exists y \in \Sigma^* \ xy \in L\}$
- suffixes(L) = $\{y \mid \exists x \in \Sigma^* \ xy \in L\}$
- infixes(L) = $\{y \mid \exists x, z \in \Sigma^* \ xyz \in L\}$

Corollaire 2: L'image d'un langage rationnel par un morphisme (i.e. une application h telle que $\forall x, y \in \Sigma^* \ h(xy) = h(x).h(y)$) est un langage rationnel.

Propriété: Pour L langage rationnel et K un langage quelconque, le quotient à gauche de L par K i.e. $K^{-1}L = \{y \mid \exists x \in K \ xy \in L\}$ et le quotient à droite de L par K i.e. $LK^{-1} = \{x \mid \exists y \in K \ xy \in L\}$ sont des langages rationnels.

2) lemme(s) de l'étoile et variantes [Sak, p 155]

Définition [Facteur itérant] Un mot v est dit facteur itérant d'un langage L lorsque il existe un v dans Σ^* tels que $u v^* w \subset L$

lemme de l'étoile 1: Dans un langage rationnel, tout mot suffisamment long contient un facteur itérant non vide.

Application: $\{a^n b^n \mid n \in \mathbb{N}\}$ n'est pas rationnel.

lemme de l'étoile 2: Dans un langage rationnel L, tout facteur suffisamment long contient un facteur itérant non vide. i.e. il existe un entier N tel que pour tout β dans L et toute factorisation $\beta = g_1 h g_2$ telle que $|h| \geq N$, il existe une factorisation de $h = uv^*w$, avec $v \neq \epsilon$ telle que $g_1 u v^* w g_2 \subset L$.

Application: $\{a^n b^n \mid n \in \mathbb{N}\} \cup \Sigma^* b a \Sigma^*$ n'est pas rationnel.

lemme de l'étoile par bloc: Si L est rationnel, il existe $N \geq 0$ tel que pour tout $\beta \in L$ et toute factorisation de la forme $\beta = u v_1 v_2 \dots v_r w$, où $v_i \neq \epsilon$

il existe (j, k) , $0 \leq j < k \leq N$, tel que: $u v_1 v_2 \dots v_j (v_{j+1} \dots v_k)^* v_{k+1} \dots v_r w \subset L$

Application: $\{(aabb)^n (abbb)^m \mid n, m \in \mathbb{N}\} \cup \Sigma^* \{aaa, bbb\} \Sigma^*$ n'est pas rationnel.

Théorème: [Ry, sel p 36] Un langage L est régulier ssi il existe $N \geq 0$ tel que pour tout $w \in L$ avec $|w| \geq N$ il existe une factorisation $w = xy^iz$, $y \neq \epsilon$ tel que:

$$\forall i \in \mathbb{N} \forall v \in \Sigma^* \ (wv \in L \Leftrightarrow xy^i z \in L)$$

Théorème: L'ensemble des facteurs itérants d'un langage rationnel forme un langage rationnel. **DVP**

3) Quotients [Bq, p 312]

Définition [quotient gauche, quotient droit] Les quotients gauche et droit d'un langage L par un mot w sont respectivement:

$$w^{-1}L = \{v \mid wv \in L\} \quad \text{et} \quad Lw^{-1} = \{u \mid uw \in L\}$$

Notons que $\epsilon^{-1}L = L$ et $(wv)^{-1}L = v^{-1}(w^{-1}L)$

Proposition Soit L un langage reconnu par un automate A, on note pour $q \in Q(A)$ $L_q = \{w \in \Sigma^* \mid q \xrightarrow{w} \text{état final de } A\}$
 Alors les L_q décrivent l'ensemble des quotients gauches de L.

*: notion supposée connue

Théorème: Un langage est rationnel si il admet un nombre fini de quotients gauches.

Ex: $L = (a+b)^* ab (a+b)^*$; $a^{-1}L = L + b(a+b)^*$; $b^{-1}L = L$; $(ab)^{-1}L = (a+b)^*$; $(aa)^{-1}L = L$; $(aba)^{-1}L = (abb)^{-1}L = (a+b)^*$

III Minimisation des automates

Motivation de cette partie!

* Obtenir un procédé efficace (principalement au niveau de l'occupation mémoire) pour déterminer si un mot donné appartient à un langage rationnel fixé ou pas. (aspect pratique)

* Pouvoir tester l'équivalence de deux expressions rationnelles (aspect théorique)

1) Automate minimal et équivalence de Nerode [Bq p.315]

Définition [Équivalence de Nerode]

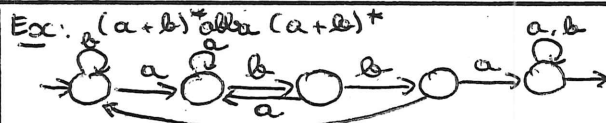
Pour un automate $A = (Q, \Sigma, \delta, i, F)$ on définit une relation d'équivalence sur Q appelée équivalence de Nerode définie comme

$p \sim q$ ssi $L_p = L_q$
ssi $\forall u \in \Sigma^* \delta^*(p, u) \in F \Leftrightarrow \delta^*(q, u) \in F$

On notera $A/\sim = (Q/\sim, \Sigma, \bar{\delta}, i, F)$

l'automate A quotienté par \sim (Q/\sim est l'ensemble des classes de la relation \sim et $\bar{\delta}: (\bar{a}, a) \mapsto \bar{\delta}(a)$ est bien définie car \sim est régulière à droite)

Théorème: L'automate minimal est l'automate ayant le moins d'états parmi les automates déterministes complets qui reconnaissent L . Il est unique à isomorphisme près. Il est égal à A/\sim , mais aussi à l'automate des quotients.



Application: On peut décider si oui ou non deux expressions données sont équivalentes

2) Calcul de l'automate minimal [Lar p.65]

Méthode de Moore: (On fixe $d = (Q, \Sigma, \delta, i, F)$ det.)
Proposition: L'équivalence de Nerode est la relation d'équivalence sur Q régulière à droite (i.e. $p \sim q \Leftrightarrow \delta(p, a) \sim \delta(q, a)$) la plus grossière qui est plus fine que d, F, F^c .

On définit alors une suite (v_i) telle que $q \sim v_0 q' \Leftrightarrow (q \in F \Leftrightarrow q' \in F)$
 $q \sim v_{i+1} q' \Leftrightarrow (q v_i, q' \text{ et } \forall a \in \Sigma \delta(q, a) \sim v_i \delta(q', a))$
Si $v_i \sim v_{i+1}$, alors v_i est l'équivalence de Nerode (l'algorithme est en $O(mn^2)$)

Algorithme de Hopcroft:

Définition [Stabilité] Soient $P \subseteq Q, R \subseteq Q$ et $a \in \Sigma$

P est dit stable pour (R, a) ssi $Pca^{-1}R$ ou $Pca^{-1}\bar{R}$.
Sinon P est dit loisné pour (R, a) , i.e. $P \cap a^{-1}R \neq \emptyset$ et $P \cap a^{-1}\bar{R} \neq \emptyset$
Une partition \mathcal{P} de Q est stable pour (R, a) ssi $\forall P \in \mathcal{P} P$ est stable pour (R, a)
 \mathcal{P} est dit stable (tout court) ssi $\forall R \in \mathcal{P} R$ est stable pour (R, a) .

Proposition Une relation d'équivalence sur Q est régulière à droite ssi la partition associée est stable.

Théorème: L'algorithme dit de Hopcroft minimise l'automate d en une complexité $O(mn \log n)$ où $m = \text{card}(\Sigma)$ et $n = \text{card}(Q)$

IV Applications

1) En biologie, en mathématiques
• À l'origine introduit pour modéliser les neurones

• Peut être utilisé pour résoudre des problèmes de combinatoire... Exemples à l'oral

2) Reconnaissance de motifs [Sak p.63]

On veut chercher un mot u dans un texte. Concrètement on veut reconnaître le langage $\Sigma^* u \Sigma^*$ (ou $\Sigma^* u$) (exemple repris durant toute cette leçon!) avec $u = abba$

On note $\text{Pref}(u) =$ ensemble des préfixes de u

$B_u = (\text{Pref}(u), \Sigma, \delta, \{ \epsilon \}, \{ u \})$ avec $\delta(p, a) =$ le plus long g préfixe de pa dans $\text{Pref}(u)$

Théorème: B_u reconnaît $\Sigma^* u$. De plus, il s'agit de l'automate minimal.

3) Analyse lexicale [WH p.209]

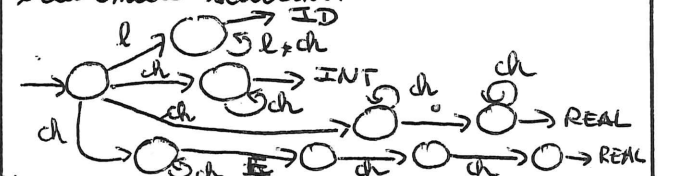
C'est la première étape d'un compilateur. Elle précède l'analyse syntaxique.

Elle est chargée de lire un programme source (sous forme d'une suite de caractères) afin de le décomposer en une succession d'unités lexicales, appelées lexèmes.

Ex de lexèmes: Entier, Identificateur, mots réservés au langage (comme begin, if, int)

Pour ex, on veut reconnaître le langage $\underline{l(ch+l)^*} + \underline{ch} \underline{ch^*} + \underline{ch} \underline{ch^*} \underline{ch^*} + \underline{ch} \underline{ch^*} \underline{E} \underline{ch} \underline{ch}$

ID INT REAL REAL
où $ch = \{1, 2, \dots, 9\}$ et $l = \{a, b, \dots, z\}$
On peut alors déterminer & minimiser l'automate suivant.



ANNEXES

À propos des références

[Car] Carton : Un des plus durs sur le sujet à mon avis

[Sak] Sakaravitch : Complet, bien expliqué, j'aime beaucoup, il a juste un humour assez léger

[Rog-Sal] Rogyberg - Salama (Hanoïloise of Formal Languages) Livre sur théorie, en anglais, mais pas mal expliqué.

[Bq] Beauquier (Éléments d'Algèbre théorique) Si, si, on parle des automates dans le Beauquier, j'aime bcp la partie sur la minimisation, c'est un des ouvrages de plus durs sur le sujet. Mais aussi.

[Wol] Wolper

Apparemment intéressant (d'après Leclerc) mais j'en ai pas eu le temps.

[Wol] Wilhelm (Les compléments) No comment.

J'aurais emprunté le Autbert mais je n'aurais pas du tout aimé.

À propos du plan : la partie apparemment incontournable serait la minimisation des automates - Moi je suis convaincu que ça a sa place ici, à voir de près votre devoir (et votre avis).

Je déconseillerais de parler de dérivation des expressions régulières / Reuteurs d'états / langages sous étude. Sur bon théorique avec peu d'applications pratique, pas sûr que ça passe au jury - (quoi que la dérivation soit à discuter...)

Bien sûr tout n'est pas indispensable dans ce que j'ai mis - J'ai été même parfois (un peu) redondant - Il aurait également fallu que je précise un peu mieux les notations utilisées

À propos des développements les développements que je propose, à savoir la régularité des patrons itérants et l'algorithme de Hopcroft sont plutôt longs et pas du genre faciles je pense - le second a de gros overwriting de preuve se passe dans 3 lignes (dont celle sur Divers pour régner) Un investissement rentable? (Car il y a du profit). J'aurais peut-être aussi à :

- TR de Kleene

- Localisation / Automates des péages (dep d'Availia que)

- Équivalence avec les grammaires linéaires

- Méthode de Brzozowski - complétude autour d'espace de l'équivalence de 2 expr régulières?