

10/02
2015

WOL p 48
CAR p 75

WOL p 78

CAR p 76

WOL p 48

CAR p 73

CAR p 90

Langages algébriques. Exemples:

910

Insuffisance des langages rationnels: $\{a^n b^n\}$ n'est pas reconnaissable par automate!

Grammaires hors contextes et langages algébriques.

Def 1: Une grammaire hors contexte (ou algébrique) est un quadruplet (V, T, R, S) vérifiant:

- V , l'ensemble des variables, est fini (on notera ses éléments en majuscules)
- T , l'ensemble des terminaux, est fini (éléments notés en minuscule), avec $V \cap T = \emptyset$.
- $S \in V$ est le symbole de départ.
- R est un ensemble de règles: $R \subset V \times (VUT)^*$. On notera, pour $(X, u) \in R$, $X \rightarrow u$. X est le member de gauche et le mot u sur l'alphabet (VUT) est le member de droite.

Exemple 2: $T = \{a, b\}$, $V = \{S\}$, $R = \{S \rightarrow aSb, S \rightarrow \epsilon\}$.

Exemple 3: $T = \{a, b\}$, $V = \{S\}$, $R = \{S \rightarrow aSa, S \rightarrow bSb, S \rightarrow c\}$.

Def 4: Soit $G = (V, T, R, S)$ une grammaire algébrique, u et v deux mots sur $V+T$. On dit que u se dérive en v , noté $u \rightarrow v$, lorsque:

$\exists \alpha, \beta \in (A+V)^*, X \in V$ tels que:
 $u = \alpha X \beta$, $v = \alpha w \beta$, et $X \rightarrow w \in R$.

Remarque: l'appellation grammaire hors contexte est licite!

Def 5: Un langage est dit algébrique s'il est généré par une grammaire algébrique, id est: $L = \{v \in T^*; S \xrightarrow{*} v\}$, où $G = (V, T, R, S)$ est une grammaire algébrique.

Exemple 6: Le langage algébrique généré par la grammaire de l'exemple 2 est $\{a^n b^n, n \geq 0\}$, et celui de l'exemple 3 est $\{w \bar{w}^R, w \in \{a, b\}^*\}$.

Lemme 7 (fondamental): Soit $G = (V, T, R, S)$ une grammaire algébrique, u et v deux mots sur $(T+V)^*$. On sup. on que u se factorise en $u = u_1 u_2$. Alors il existe une dérivation de longueur h $u \xrightarrow{h} v$ ssi il existe une factorisation $v = v_1 v_2$ telle que:
 $u_1 \xrightarrow{h_1} v_1$, $u_2 \xrightarrow{h_2} v_2$ et $h_1 + h_2 = h$.

Exemple 8: G définie par $T = \{a, b\}$, S , $V = \{S, A\}$, $R = \{S \rightarrow SA + \epsilon, A \rightarrow aSb\}$ (langage de Dyck).

AA ne dérive $a b a a b b$ avec $A \rightarrow ab$ et $A \rightarrow a a b b$.

Def 9: Une grammaire algébrique. Un arbre de dérivation est un arbre fini étiqueté par VUT et vérifiant: si A est l'étiquette d'un noeud et a_1, \dots, a_n les étiquettes de tous ses fils, alors $A \rightarrow a_1 \dots a_n \in R$. La frontière de l'arbre est la concaténation de toutes les feuilles, de gauche à droite.

Exemple 10: cf annexe 1.

Def 10: Une grammaire G est dite ambiguë s'il existe au moins un mot ayant au moins deux arbres de dérivation.

Exemple 11: $V = \{S\}$, $T = \{a\}$, $R = \{S \rightarrow SS + a\}$. cf annexe 2.

Def 12: Une grammaire G est dite propre s'il n'existe aucune règle de la forme $A \rightarrow \epsilon$ ou $A \rightarrow B$ avec $A, B \in V$.

Exemple 13: la grammaire de l'exemple 2 n'est pas propre, mais la suivante l'est: $R = \{S \rightarrow aSb + ab\}$. Les langages reconnus sont les mêmes à ϵ près.

Def 14 (Th 16): Tout langage algébrique ne générant pas ϵ peut être engendré par une grammaire sous forme normale de Chomsky, c'est à dire que toutes les productions sont de la forme $A \rightarrow BC$ ou $A \rightarrow a$.

Exemple 15: $G = (\{S, A, B\}, \{a, b\}, A, S)$
avec les règles:
 $S \rightarrow bA + aB$
 $A \rightarrow bAA + aS + a$
 $B \rightarrow aBB + bS + b$

forme normale de Chomsky:	$S \rightarrow B'A + AB'$
	$A \rightarrow B'A'' + A'S + a$
	$B \rightarrow A'B'' + B'S + b$
	$B' \rightarrow b, A' \rightarrow a,$
	$B'' \rightarrow BB, A'' \rightarrow AA.$

Def 16 (Th 16): Tout langage algébrique ne générant pas ϵ peut être engendré par une grammaire sous forme normale de Equirbach, id est toutes les productions sont de la forme: $A \rightarrow \alpha \alpha$, où $\alpha \in V^*$.

Exemple 17 (langage de Dyck généralisé):
 $G = (\{S, A\}, \{a_1, \bar{a}_1, \dots, a_n, \bar{a}_n\}, R, S)$,
avec $R = \{S \rightarrow SA + \epsilon, A \rightarrow a_i S \bar{a}_i + \dots + a_n S \bar{a}_n\}$.

forme normale de Equirbach:	$S \rightarrow SA + \epsilon$
	$A \rightarrow a_i S X_1 + \dots + a_n S X_n$
	$X_i \rightarrow \bar{a}_i$

I

CAR p 90

CAR p 80

HU p 32

HU p 95

CAR p 77

AR
084

II Signifiés générales des langages algébriques.

Def 18: Système d'équation associé à une grammaire algébrique G:

Soit $n := |V|$, alors le système est constitué des n équations:

$$L_i = \sum_{X_i \rightarrow w} w(L_i), \text{ où } X_i \in V, \text{ en les variables } L_1, \dots, L_n. \quad (\text{noté } S_G)$$

Propo 19: Soit $G = (\{X_1, \dots, X_n\}, T, R, S)$. Alors le langage généré par G est la solution minimale du système défini précédemment. (*)

Exemple 20: $\begin{cases} X_1 \rightarrow X_1 X_2 + \varepsilon \\ X_2 \rightarrow a X_2 + b \end{cases}$ donne le système $\begin{cases} L_1 = L_1 L_2 + \varepsilon \\ L_2 = a L_2 + b \end{cases}$

Sa solution minimale est: $\{\varepsilon\} \cup \{(a+b)^* b\}$

Caractérisation des langages algébriques (conditions nécessaires):

WOL
088

Lemma 21 (Gogden): Soit $G = (V, T, R, S)$ et L le langage généré par V .

Il existe $K > 0$, telle que pour tout mot $f \in L$ avec $|f| > K$, il existe une factorisation $f = u v x y z$ avec v ou $y \neq \varepsilon$, $|u x y| \leq K$ et $u^i v^j x y^k z^l \in L$ pour tout $i, j, k, l \geq 0$.

HU

Exemple 22: $L = \{a^i b^j c^k \mid i \geq j, j \geq k \text{ et } i + k \geq n\}$ n'est pas algébrique.

P 130

Th 23 (Ben-Hellal, Sar, Shamir): Les langages algébriques sont clos par substitution, union, concaténation (morphisme).

P 135

Th 24: Si L est algébrique et K est rationnel, alors $L \cap K$ est algébrique.

Th 25: La classe des langages algébriques n'est pas close par intersection!

Exemple 26: $L_1 = \{a^i b^j c^i, i \geq 1, j \geq 1\}$ algébrique
 $L_2 = \{a^i b^j c^k, i \geq 1, j \geq 1, k \geq 1\}$ algébrique
 $L_3 = \{a^i b^j c^k, i \geq 1, j \geq 1, k \geq 1\}$ algébrique.

Le lemme d'Gogden montre que $L_1 \cap L_2$ n'est pas algébrique, mais

$$L_1 = L_2 \cap L_3!$$

AR
P 88

Def 27: Une grammaire étendue G est définie par (V, T, R, S) avec V et T deux alphabets finis disjoints, $R \subset V \times (T \cup V)^*$ telle que pour tout $X \in V$, $\{w \mid (X, w) \in R\}$ est rationnel.

Propo 28: L'ensemble des mots engendrés par une grammaire étendue est algébrique.

Def 29: On note \bar{L} l'image commutative d'un langage: $\bar{L} := \{\bar{w} \mid w \in L\}$,

où \bar{w} est la classe des anagrammes de w (on peut permuer les lettres).

Th 30 (Parikh): Tout langage algébrique L est commutativement équivalent à un langage rationnel R (i.e. $\bar{L} = \bar{R}$). DVT

CAR
P 86

III Automates à pile.

Def 31: Un automate à pile est défini par un septuplet $M = (Q, \Sigma, \Gamma, \Delta, z, s, f)$:

- Q est un ensemble fini d'états,
- Σ est l'alphabet d'entrée,
- $z \in \Gamma$ est le symbole initial de pile,
- Γ est l'alphabet de pile,
- $s \in Q$ est l'état initial,
- $F \subset Q$ est l'ensemble des états acceptés,

$\Delta \subset (Q \times \Sigma^* \times \Gamma^*) \times (Q \times \Gamma^*)$ est la relation de transition.

WOL
P 73

Def 32: La configuration (q', w', α') est dérivable en une étape de la configuration (q, w, α) (noté $(q, w, \alpha) \vdash (q', w', \alpha')$) si:

$$w = u w', \quad \alpha = \beta \delta \text{ et } \alpha' = \gamma \delta \text{ (sur la pile),}$$

$$((q, u, \beta), (q', \gamma)) \in \Delta$$

Def 33: Une exécution d'un automate à pile sur un mot w est une suite de configurations (finie ou infinie):

$$(s, w, z) \vdash (q_1, w_1, \alpha_1) \vdash \dots \vdash (q_n, w_n, \alpha_n) \vdash \dots$$

Def 34: Un automate à pile acceptant sur état final, accepte le mot w si: $(s, w, z) \vdash^* (p, \varepsilon, \gamma)$, avec $p \in F$.

Def 35: Un automate à pile acceptant sur pile vide, accepte le mot w si: $(s, w, z) \vdash^* (q, \varepsilon, z)$.

WOL
P 76

Def 36: Un langage L est reconnu par un automate à pile M si et seulement si l'ensemble des mots acceptés par M est L .

Propo 37: Un langage est reconnu par un automate acceptant sur état final si et si seulement il est reconnu par un automate acceptant sur pile vide.

(cf annonce 3)

Exemple 38: $Q = \{s, p, q\}$ $\Sigma = \{a, b\}$ $\Gamma = \{A\}$ $F = \{q\}$ et

$$A = \left\{ \begin{array}{l} (s, a, \varepsilon) \rightarrow (s, A) ; (s, \varepsilon, z) \rightarrow (q, \varepsilon) ; (s, b, A) \rightarrow (p, \varepsilon) ; \\ (p, b, A) \rightarrow (p, \varepsilon) ; (p, \varepsilon, z) \rightarrow (q, \varepsilon) \end{array} \right\} \text{ reconnaît } \{a^n b^n\}.$$

CAR
p 108

Th 39: Un langage L est algébrique si et seulement si il est reconnu par un automate à pile.

Def 40: Un automate à pile M est déterministe si pour toute paire $(q, \alpha) \in Q \times \Gamma$; soit il existe une unique transition $q, \alpha, \varepsilon \rightarrow q', \beta$ et il n'existe pas d'autre transition $q, \alpha, a \rightarrow q', \beta$, où $a \in \Sigma$; soit il n'existe pas de transition de la forme $q, \alpha, \varepsilon \rightarrow q', \beta$, et il existe au maximum un $a \in \Sigma$ tel que $q, \alpha, a \rightarrow q', \beta$.

Def 41: Un langage algébrique est déterministe si il est reconnu par un automate à pile déterministe acceptant son état final.

Prop 42: Le complémentaire d'un langage algébrique déterministe est un langage algébrique déterministe.

Prop 43: Si L est reconnu par un automate à pile déterministe acceptant, alors il est reconnu par un automate à pile déterministe un pile vide acceptant un état final.

⚠ La réciproque n'est pas vraie!

Prop 44: Tout langage algébrique déterministe est non ambigü.

Exemple 45: $L_1 = \{w c \bar{w}^R \mid w \in \{a, b\}^*\}$ est déterministe, mais (ADMES) $L_2 = \{w \bar{w}^R \mid w \in \{a, b\}^*\}$ ne l'est pas.

IV Analyse syntaxique.

Une partie du travail effectué par le compilateur: vérifier que la syntaxe est correcte que le mot appartient au langage.

Analyse descendante: on part d'un mot et de S , et on modifie S (ou la partie gauche dans l'exemple) avec les règles pour supprimer au fur et à mesure les lettres du mot.

Analyse ascendante: on essaie de réduire le mot de droite sur l'exemple, en utilisant les membres droits d'une dérivation.

AUT
p 127

Analyse descendante LL

- $S \rightarrow cT$ (1)
- $T \rightarrow aTbc$ (2)
- $T \rightarrow b$ (3)

- $(S, c a b b c)$
- $(Tc, c a b b c)$ (1)
- $(T, a b b c)$
- $(c b T a, a b b c)$ (2)
- $(c b T, b b c)$
- $(c b b, b b c)$ (3)
- $(c b, b c)$
- (c, c)
- (ϵ, ϵ)

Analyse ascendante LR

- $E \rightarrow E+T$ (1)
- $T \rightarrow (E)$ (2)
- $E \rightarrow T$ (3)
- $T \rightarrow a$ (4)
- $T \rightarrow b$ (5)

- $[1, b + (a + a)]$
- $[b, + (a + a)]$
- $[T, + (a + a)]$ (5)
- $[E, + (a + a)]$ (3)
- $[E+, (a + a)]$
- $[E+, (, a + a)]$
- $[E+, (a, + a)]$
- $[E+(T, + a)]$ (4)
- $[E+(E+T, ,)]$
- $[E+(E,)]$ (1)
- $[E+(E),)]$
- $[E+T,)]$ (2)
- $[E,)]$ (1)

Voici un algorithme qui teste l'appartenance d'un mot à un langage algébrique: l'algorithme Earle-Younger-Kasami DVI

Th 46: Il existe un algorithme qui teste l'appartenance à un langage algébrique en temps $O(|w|^3)$.

(*) Prop 19 bis: Soit G une grammaire propre; alors L_G est l'unique solution propre de S_G .

(**) Prop 23 bis: Soit G une grammaire; alors toute solution L propre du système S_G vérifie $L = L_G$.

H-U: Hopcroft - Ullman: Introduction to automata theory...

CAR: Barton, langages formels.

WOL: Wolper, introduction à la calculabilité.

AUT: Abutebort; théorie des langages et des automates.

F.B: Floyd - Szigel, Le langage des machines.

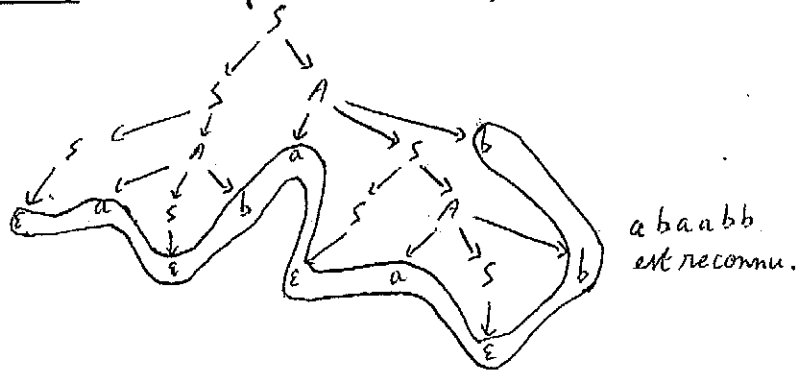
AUT
CAR p 125

WOL p 97

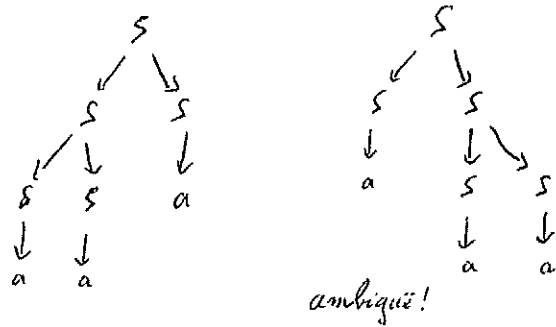
AUT p 133

FB
p 320

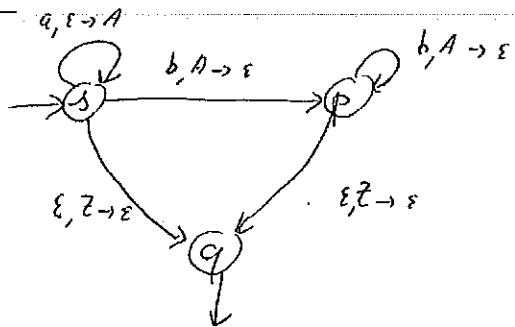
Annexe 1: avec la grammaire de l'exemple 8,



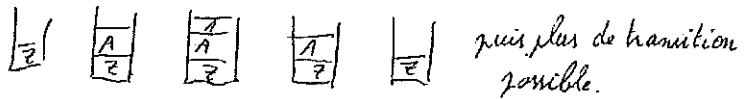
Annexe 2:



Annexe 3:



Pour reconnaître aabba, on a sur la pile:



Théorème de Parikh

Soit L un langage algébrique sur un alphabet A , alors il existe R rationnel sur A tel que $L = \bar{R}$ (où $\bar{\quad}$ désigne l'image commutative)

Preuve: par abus, on se permettra d'omettre le symbole de départ des grammaires, manipulées

I) Lemme: soit G une grammaire propre, alors toute solution L de

$$\bar{S}_G \text{ vérifie } L = \bar{L}_G$$

II) Proposition: soit G une grammaire généralisée, alors le système \bar{S}_G

admet une solution rationnelle, i.e. si $G = (A, V, P \{X_i \rightarrow S_i(X)\})$, où

$$V = \{X_0, \dots, X_m\} \text{ et } X = (X_0, \dots, X_m), \text{ il existe un } n\text{-uplet } R = (R_0, \dots, R_m)$$

de langages rationnels sur A tel que $\forall i \in \{0, \dots, m\} \quad \bar{R}_i = \bar{S}_i(R)$

III) Conclusion

I) Il s'agit de la proposition 29 bis du plan. La preuve est essentiellement la même que celle de la proposition 19 bis (voir Canton p.95)

II) On raisonne par récurrence sur le nombre de non-terminaux de G

Soit $\mathcal{P}(n)$: "Pour toute grammaire $G = (A, V, P)$ généralisée telle que $|V| = n$, \bar{S}_G admet une solution rationnelle".

Montrons $\mathcal{P}(1)$: soit $G = (A, \{X\}, \{X \rightarrow S(X)\})$ une grammaire généralisée.

$$\text{Soit } T = S(X) \cap A^* \text{ et } P(X) = S(X) \mid T = S(X) \cap ((A+X)^* X (A+X)^*)$$

On a donc $P(X) + T = S(X)$; $T \in \text{Rat}(A)$ et $P(X) \in \text{Rat}(A+X)$

Il existe $Q(X) \in \text{Rat}(A+X)$ tel que $\overline{P(X)} = \overline{Q(X)X}$ (par exemple, on peut choisir Q formé des mots de P privés de leur première occurrence de X)

Le système \bar{S}_G se réécrit alors, pour l'inconnue L :

$$\bar{L} = \overline{Q(L)L + T}$$

Soit donc $R = Q(T)^* T$, $R \in \text{Rat}(A)$ (car $Q(T)$ est l'image de $Q(X)$ par un morphisme rationnel)

$$\text{et } \overline{Q(R)R + T} = \overline{Q(Q(T)^* T)Q^* T + T} = \overline{Q(T)Q(T)^* T} = \overline{Q(T)^* T} = \bar{R}$$

donc R est solution de \bar{S}_G , d'où $\mathcal{P}(1)$

Soit $n \geq 1$, on suppose $\mathcal{P}(A)$ pour $k \leq n$

Soit $G = (A, \{X_0, \dots, X_m\} = V, S = \bigcup_{i=0}^m \{X_i \rightarrow S_i(X)\})$ une grammaire généralisée à $n+1$ non-terminaux, où X désigne (X_0, \dots, X_m) .

On notera $V' = \{X_1, \dots, X_m\}$ et $X' = (X_1, \dots, X_m)$.

On définit $G_0 = (A+V', \{X_0\}, \{X_0 \rightarrow S_0(X_0, X')\})$.

Par hypothèse de récurrence, $\overline{S_0}$ admet une solution $R_0(X') \in \text{Rat}(A+V')$, i.e. $\overline{S_0(R_0(X'), X')} = R_0(X') \{*\}$

Soit maintenant $G' = (A, V', \bigcup_{i=1}^m \{X_i \rightarrow S_i(R_0(X'), X')\})$

Par hypothèse de récurrence, $\overline{S_i}$ admet une solution rationnelle $R' = (R'_1, \dots, R'_m)$

i.e. $\forall i \in \{1, \dots, m\} R_i \in \text{Rat}(A)$ et $\overline{R'_i} = \overline{S_i(R_0(R'), R')}$

De plus, $\{*\}$ entraîne en particulier $\overline{S_i(R_0(R'), R')} = \overline{R_0(R')}$

et $R_0(R') \in \text{Rat}(A)$ en tant qu'image de $R_0(X') \in \text{Rat}(A+V')$ par un morphisme de $A+V'$ dans $\text{Rat}(A)$

donc $R = (R_0(R'), R')$ est solution rationnelle de $\overline{S_0}$, d'où $\mathcal{P}(n+1)$.

$\mathcal{P}(n)$ est vrai pour tout n , ce qui conclut la preuve de II)

III) ; Soit L un langage algébrique sur un alphabet A .

Soit G une grammaire propre reconnaissant $L \in \mathcal{E}$, $G = (A, \{X_0, \dots, X_m\}, P, X_0)$

Alors si $L_0 = (L_0, \dots, L_m)$, on a $L \in \mathcal{E} = L_0$

Soit $R = (R_0, \dots, R_m)$ une solution rationnelle de $\overline{S_0}$. On a $\overline{R} = \overline{L_0}$ donc en particulier $\overline{R_0} = \overline{L_0} = \overline{L \in \mathcal{E}}$

d'où $\overline{L} = \overline{R_0 \in \mathcal{E}(L)}$, avec $R_0 \in \mathcal{E}(L)$ rationnel sur A .