

Introduction: les langages algébriques, appelés aussi langages non-contexte, sont utiles pour l'analyse syntaxique des langages de programmation (ex: camlyacc).

I Grammaires algébriques.

1) Définitions

DEF: une grammaire algébrique est un triplet (A, V, P) où A, V alphabets finis et disjoints, et P partie finie de $V \times (AV)^*$.

Les symboles de A sont les terminaux, ceux de V les variables. P contient des règles $X \rightarrow w$ (notation).

EX: $A = \{a, b\}$, $V = \{S\}$, $P = \{S \rightarrow aS + E\}$, $G_{(1)} = (A, V, P)$.

REM: Par convention les lettres minuscules sont des terminaux, les majuscules des variables.

DEF: Soient $u, v \in (A+V)^*$. u se dérive en v , et on note $u \rightarrow v$, s'il existe $\alpha, \beta \in (A+V)^*$ et $X \in V$ tels que $v = \alpha X \beta$, $v = \alpha w \beta$ et $(X \rightarrow w) \in P$.

EX: $G_{(1)}$: $S \rightarrow aSb \rightarrow aaSbb \rightarrow aabb$.

REM: On note \rightarrow^* la clôture réflexive et transitive de \rightarrow .

DEF: Si $u \in (A+V)^*$, $\hat{L}_G(u) := \{v \in (A+V)^* \mid u \rightarrow^* v\}$
et $L_G(u) := \hat{L}_G(u) \cap A^*$.

$L_G(S)$ est appelé langage engendré par G .

DEF: un langage est dit algébrique s'il existe G grammaire algébrique et $S \in V$ tel que $L = L_G(S)$.

EX: $G_{(1)}$ est algébrique. $L_{G_{(1)}}(S) = \{a^m b^m \mid m \geq 0\}$.

EX: Langage de Dyck pour m parenthèses:

$$\begin{cases} A_m = \{a_1, \dots, a_m, \bar{a}_1, \dots, \bar{a}_m\} \\ S \rightarrow ST + E \\ T \rightarrow a_i S \bar{a}_i + \dots + a_m S \bar{a}_m \end{cases}$$

Lemme fondamental: Soit $G = (A, V, P)$ une grammaire algébrique, et u, v deux mots de $(A+V)^*$. On suppose $u = u_1 u_2$. Alors il existe $u \xrightarrow{k} v$ ssi $v = v_1 v_2$ et il existe $u_1 \xrightarrow{k_1} v_1$, $u_2 \xrightarrow{k_2} v_2$ avec $k = k_1 + k_2$.

2) Simplifications de grammaires

DEF: Une grammaire algébrique $G = (A, V, P)$ est réduite par $S_0 \in V$ si:

$$- \forall S \in V, L_G(S) \neq \emptyset$$

$$- \forall S \in V, \exists u, v \in (A+V)^*, S_0 \xrightarrow{*} u S v$$

PROP: Pour toute grammaire G , et $S_0 \in V$, il existe une grammaire G' réduite par $S'_0 \in V'$ telle que $L_G(S_0) = L_{G'}(S'_0)$.

Complexité: $O(m^2)$ où m est la taille de la grammaire.

DEF: G est propre si elle ne contient aucune règle $S \rightarrow \epsilon$ où $S \rightarrow S'$ pour $S, S' \in V$.

EX: $S \rightarrow aSb + ab$

PROP: Pour toute grammaire G et $S \in V$, il existe une grammaire G' et $S' \in V'$ propre telle que $L_{G'}(S') = L_G(S) \setminus \{\epsilon\}$.

DEF: (Forme normale de Chomsky) C'est une grammaire telle que toutes les règles sont de la

$$\text{forme } \begin{cases} S \rightarrow S_1 S_2 \text{ où } S_1, S_2 \in V \\ a \\ S \rightarrow a \text{ où } a \in A. \end{cases}$$

PROP: Pour toute grammaire G et SEV, il existe une grammaire sans forme normale de Chomsky G' , et $S'EV'$ telle que $L_{G'}(S') = L_G(S) \setminus \{\epsilon\}$.

Complexité: $O(m^2)$ où m est la taille de la grammaire.

Conséquence: algorithme pour savoir si un mot donné est engendré par une certaine grammaire G .

3) Systèmes d'équations

DEF: À $G = (A, V, P)$ grammaire où $V = \{X_1, \dots, X_m\}$ on associe le système $S(G)$: $L_i = \sum w(L)$ $1 \leq i \leq m$ où $w(L)$ morphisme de substitution $X_i \rightarrow w_i$ de $(A+V)^* \rightarrow P((A+V)^*)$ définit mutuellement par induction.

DEF: on note $L_G = (L_G(X_1), \dots, L_G(X_m))$.

PROP: $\forall w \in (A+V)^*$, $L_G(w) = w(L_G)$.

PROP: L_G est la solution minimale de $S(G)$ (par l'inclusion).

REM: On n'a pas unicité par contre: $X \rightarrow XX+E$ a pour solutions tous les K^* où $K \in A^*$. (on a unicité de minimal).

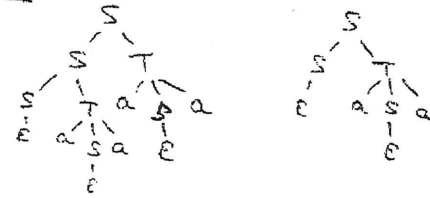
REM: systèmes d'équation: origine de la terminologie langage «algébrique».

4) Arbres de dérivation

DEF: Soit $G = (A, V, P)$ une grammaire algébrique. un arbre de dérivation est un arbre fini étiqueté par $A \cup V \cup \{\epsilon\}$ tel que si S est l'étiquette d'un nœud interne et si a_1, \dots, a_m sont les étiquettes de ses fils, alors $S \rightarrow a_1 \dots a_m$ est une règle de G . La frontière de l'arbre est le mot obtenu en

Concaténant les étiquettes des feuilles de gauche à droite.

EX: Pour $\{S \rightarrow ST+E, T \rightarrow aS\}$



DEF: Une grammaire est ambiguë s'il existe un mot ayant deux arbres de dérivation distincts avec même racine.

EX: $S \rightarrow Sa+aa$ donne S et S

TH: (lemme d'itération, Bar-Hillel, Perles et Shamir).

Pour tout langage algébrique L , il existe $N \geq 0$ tel que pour tout mot $f \in L$, si $|f| \geq N$, alors on peut trouver une factorisation $f = \alpha u \beta v \gamma$ tel que $|uv| > 0$, $|u\beta v| < N$ et $\alpha u^m \beta v^m \gamma \in L$ pour tout $m \geq 0$. (DEV)

Application: $L = \{a^m b^m c^m \mid m \geq 0\}$ n'est pas algébrique.

COR: La classe des langages algébriques n'est classe ni par intersection ni par complémentation.

II Propriétés des langages algébriques

1) Propriétés de clôture

PROP: Les langages algébriques sont clos par union, concaténation et étoile.

COR: Les langages rationnels sont algébriques.

DEF: une substitution $\sigma: A^* \rightarrow \mathcal{P}(B^*)$ où A, B sont deux alphabets, est algébrique si $\sigma(a)$ est algébrique pour tout $a \in A$.

PROP: Si $L \subseteq A^*$ est algébrique, et σ est une substitution algébrique $A^* \rightarrow \mathcal{P}(B^*)$, alors $\sigma(L) = \bigcup_{w \in L} \sigma(w)$ est algébrique.

PROP: Si L est algébrique, et k rationnel, alors kAL est algébrique.

2) Problèmes de décidabilité

PROP: Ces problèmes sont indécidables: (DEV)

a) pour deux grammaires G et G' , est-ce que

$$L_G(S) \cap L_{G'}(S') = \emptyset?$$

b) pour deux grammaires G et G' , est-ce que

$$L_G(S) = L_{G'}(S')?$$

c) pour une grammaire G , est-ce que

$$L_G(S) = A^*$$

d) pour une grammaire G , est-ce que G est ambiguë?

+ ajouter des décidables ($L = \emptyset$)

III. Automates à pile

1) Définitions

DEF: un automate à pile est constitué:

- d'un alphabet d'entrée A
- d'un alphabet de pile Z et un symbole initial $z_0 \in Z$.
- d'un ensemble fini d'états Q dont un état initial $q_0 \in Q$
- de transitions $q, z \xrightarrow{y} q', h$ avec $q, q' \in Q, y \in A \cup \epsilon, z \in Z$
 y est l'étiquette de la transition et $h \in Z$

DEF: une étape de calcul est une paire de configurations $(C, C') \in (Q \times Z^*)^2$, notée $C \xrightarrow{y} C'$ telle que $C = (p, z_0 w)$, $C' = (q, h w')$ et $p, z \xrightarrow{y} q, h$ est une transition de l'automate. un calcul est une suite d'étapes de calcul consécutives.

2) Lien avec les grammaires

TH: un langage $L \subseteq A^*$ est algébrique ssi il existe un automate à pile qui accepte L .

RS: Il y a plus modes d'acceptation équivalents (par état final, pile vide, etc...).

3) Automates à pile déterministes

DEF: un automate à pile (Q, A, Z, E, q_0, z_0) est déterministe

si pour toute paire $(p, z) \in Q \times Z$:

- soit il existe une unique $p, z \xrightarrow{y} q, h$ et pas de $p, z \xrightarrow{y'} q', h'$ pour $a \in A$

- soit il n'existe pas de $p, z \xrightarrow{y} q, h$, et pour chaque $a \in A$, il existe au plus un $p, z \xrightarrow{y} q, h$.

EX: $Q = \{q_0, q_1, q_2\}$ avec $q_0, z \xrightarrow{a} q_1, z_0$

$q_1, z \xrightarrow{a} q_1, z_0 z$ est déterministe

$q_1, z \xrightarrow{b} q_2, \epsilon$

$q_2, z \xrightarrow{b} q_2, \epsilon$

REM: Les différents modes d'acceptation ne sont plus équivalents.

PROP: Le complémentaire d'un langage algébrique déterministe en est un aussi.

PROP: Tout langage algébrique déterministe est non ambigu.
 → certaines prop. deviennent décidables ($L = A^*, L = L'$)

Références:

Carton Langages formels
 Dehornoy Automates de l'informatique
 Hopcroft & Ullman (pour les complexités)