

910 - Langages algébriques. Exemples et applications

On suppose connus les résultats sur les automates finis.

I) Représentations

1) Grammaires algébriques

a) Définition et premières propriétés

Def 1: Une grammaire algébrique G est un triplet (Σ, Γ, P) où Σ et Γ sont des alphabets finis et P est une partie finie de $\Gamma \times (\Sigma \cup \Gamma)^*$. Les symboles de Σ sont appelés *terminaux* et ceux de Γ non terminaux ou variables. Les éléments de P sont appelés règles. On note $X \rightarrow u_1 + \dots + u_n$ pour $(X, u_1), \dots, (X, u_n) \in P$.

Ex 2: $G_1 = (\Sigma_1, \Gamma_1, P_1)$, $\Sigma_1 = \{a, b\}$, $\Gamma_1 = \{S\}$, $P_1 = \{S \rightarrow aSb, S \rightarrow \epsilon\}$

Rq 3: Par convention, on notera toujours les non terminaux en majuscules et les terminaux en minuscules. Ainsi, on pourra se contenter de donner P sans introduire d'ambiguïté.

Ex 4: $G_1 = \{S \rightarrow aSb + \epsilon\}$, $G_2 = \{P \rightarrow aI + \epsilon, I \rightarrow a\}$

$G_{3,n} = \{S \rightarrow ST + \epsilon, T \rightarrow a_1 S b_1 + \dots + a_n S b_n\}$

Def 5: On étend \rightarrow à $(\Sigma \cup \Gamma)^*$ par $X\beta \rightarrow \alpha u \beta$ si $X \rightarrow u$. On note \rightarrow^* la clôture réflexive transitive de \rightarrow . Si $u \rightarrow^* v$, on dit que u se dérive en v .

Ex 6: $S \rightarrow aSb \rightarrow aaSbb \rightarrow aaaSbbb$ donc G_1 .

Def 7: On note $L_G(u) = \{v \in (\Sigma \cup \Gamma)^* \mid u \rightarrow^* v\}$ et $L_G(u) = \bigcup_{v \in L_G(u)} \{v\}$. On appelle $L_G(u)$ le langage engendré par u dans G .

Ex 8: $L_{G_1}(S) = \{a^n b^n \mid n \in \mathbb{N}\}$, $L_{G_2}(P) = \{a^{2m} \mid m \in \mathbb{N}\}$, $L_{G_2}(I) = \{a^{2n+1} \mid n \in \mathbb{N}\}$

$D_n^* := L_{G_n}(S)$ est appelé langage de Dyck. C'est le langage des mots bien parenthésés (si on considère a_i comme une parenthèse ouvrante et b_i comme la parenthèse fermante correspondante).

Def 9: Un langage est algébrique ssi il est engendré par un non terminal dans une grammaire algébrique.

Lemme 10 (Fondamental) $(u_1, u_2 \rightarrow^* v) \Leftrightarrow (\exists k_1, k_2 \geq 0) \left(\begin{matrix} u_1 = v_1 + u_2 \\ u_2 = v_2 + u_1 \\ k = k_1 + k_2 \end{matrix} \right)$

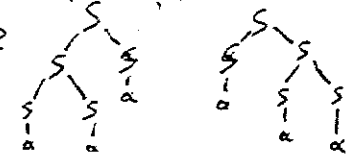
b) Arbres de dérivation

Def 11: Un arbre de dérivation partiel est un arbre fini dont les nœuds sont étiquetés par $\Sigma \cup \Gamma \cup \{\epsilon\}$ tel que si X étiquette un nœud et a_1, \dots, a_n étiquettent ses fils, alors $X \rightarrow a_1 \dots a_n$. Un arbre de dérivation est un arbre de dérivation partiel dont les feuilles sont étiquetées par $\Sigma \cup \{\epsilon\}$.

Def 12: La frontière d'un arbre de dérivation partiel est la concaténation de gauche à droite de ses feuilles.

Prop 13: $L_G(S)$ (resp. $L_G(S)$) est l'ensemble des frontières d'ordres de dérivation (resp. partiels) de racine étiquetée par S .

Exemple 14: $G_4 = \{S \rightarrow SS + a\}$



Def 15: Une grammaire est dite *ambiguë* s'il existe un mot qui est frontière de deux arbres de dérivation distincts dont la racine est étiquetée par un même non terminal.

Ex 16: G_4 est ambiguë.

Rq 17: $G_5 = \{S \rightarrow aS + a\}$ n'est pas ambiguë et $L_{G_5}(S) = L_{G_4}(S)$

c) Simplification des grammaires

Def 18: Une grammaire est dite *réduite* pour $S \in V$ ssi: $\forall T \in \Gamma, L_G(T) \neq \emptyset$ - $\forall T \in \Gamma, \exists u, v \in (\Sigma \cup \Gamma)^* / S \rightarrow^* uTv$

Def 19: Une grammaire est dite *propre* ssi elle ne contient aucune règle de la forme $S \rightarrow \epsilon$ ou $S \rightarrow T$.

Def 20: Une grammaire est en *forme normale quadratique* / de Chomsky ssi toutes ses règles sont de la forme $S \rightarrow UV$ ou $S \rightarrow a$.

Def 21: Une grammaire est en *forme normale de Greibach* (resp. normale quadratique) ssi toutes ses règles sont de la forme $S \rightarrow uV$ avec $u \in \Sigma \cup \{\epsilon\}$ (resp. $\Sigma \cup \Gamma \cup \{\epsilon\}$).

Prop 22: Pour toute grammaire algébrique de taille n engendrant L , on a

- L est engendré par une grammaire réduite de taille $O(n)$ calculable en temps $O(n)$.
- $L \setminus \{\epsilon\}$ est engendré par une grammaire normale de taille $O(n^2)$ calculable en temps $O(n^2)$.
- $L \setminus \{\epsilon\}$ est engendré par une grammaire en forme normale quadratique de taille $O(n^2)$ calculable en temps $O(n^2)$.
- $L \setminus \{\epsilon\}$ est engendré par une grammaire en forme normale de Greibach.

2) Automates à pile

Def 23: Un automate à pile est constitué: - d'un alphabet d'entrée Σ
 - d'un alphabet de pile Z avec un symbole initial z_0
 - d'un ensemble fini d'états Q dont un état initial q_0
 - de transitions de la forme $q, \gamma \xrightarrow{a} q', h$ avec $q, q' \in Q, \gamma \in Z \cup \{\epsilon\}, a \in \Sigma, h \in Z^*$.

Def 24: Une étape de calcul d'un automate à pile (AP) est une paire de configurations (C, C') notée $C \xrightarrow{a} C'$ telles que $C = (p, \gamma, w)$, $C' = (q, h, w)$ et $p, \gamma \xrightarrow{a} q, h$.
 Un calcul est une succession d'étapes de calcul: $C_0 \xrightarrow{a_1} C_1 \xrightarrow{a_2} \dots \xrightarrow{a_n} C_n$.
 Le mot $a_1 \dots a_n$ est l'étiquette du calcul.

Def 25: Un AP accepte un mot w par pile vide (resp. par état final) ssi $(q_0, z_0) \xrightarrow{w}^* (q, \epsilon)$ (resp. (q, w)) où $q \in F$, F étant un sous-ensemble distingué de Q , et w est quelconque.

Prop 26: L'ensemble des langages acceptés par pile vide est égal à l'ensemble des langages acceptés par état final.

Thm 27: Un langage L est algébrique ssi il existe un AP qui accepte L .

Def 28: Un AP est dit déterministe ssi pour tout C , toutes les étapes de calcul partant de C sont étiquetées par des symboles de Σ ou ϵ à deux près, ou il en existe une unique étiquette par ϵ .

Pr 29: Il n'y a pas équivalence des modes d'acceptation pour les AP déterministes.

II) Propriétés

1) Lemme d'itération

Lemme 30: (Pigeon) Pour toute grammaire et toute variable $X \in V$, il existe un entier $K \in \mathbb{N}$ tel que pour tout mot $f \in L_G^+(X)$ ayant au moins K lettres distinguées se factorise en $f = \alpha u \beta v \delta t \eta$:

- $S \rightarrow^* \alpha T \gamma$ et $T \rightarrow^* \beta T v + \beta$
- Soit α, u et β , soit β, v et δ contenant des lettres distinguées
- $\alpha \beta v$ contenant moins de K lettres distinguées

Corollaire 31: (Théorème de Bar-Hillel, Pumping and Shonit) Pour tout langage algébrique L , il existe $N \geq 0, k$ pour tout mot $f \in L$, si $|f| \geq N$ alors f se factorise en $f = \alpha u \beta v \delta$ avec $|u| > 0, |u \beta v| < N$ et $\alpha u^i \beta v^i \delta \in L$ pour tout $i \geq 0$.

Appli 32: $\{a^m b^n c^m \mid m, n \in \mathbb{N}\}$ n'est pas algébrique. Il peut représenter de la mise en forme en mode texte: \uparrow bits \downarrow

Appli 33: $\{a^m b^n c^m d^m \mid m, n \in \mathbb{N}\}$ n'est pas algébrique. Il peut représenter la séparation de deux procédures à suite m et n et leurs utilisations.

2) Propriétés de clôture

Def 34: Un morphisme de X^* dans Y^* est une fonction φ de X^* dans Y^* tel que $\varphi(\epsilon) = \epsilon$ et $\varphi(uv) = \varphi(u)\varphi(v)$. Un morphisme est entièrement déterminé par son action sur X . L'image de $L \subseteq X^*$ par φ est $\{\varphi(u) \mid u \in L\}$ et l'image inverse de $L \subseteq Y^*$ par φ est $\{u \mid \varphi(u) \in L\}$.

Def 35: Une substitution algébrique est une fonction $\sigma: X^* \rightarrow \mathcal{P}(Y^*)$ telle que $\sigma(\epsilon) = \{\epsilon\}$ et $\sigma(uv) = \sigma(u)\sigma(v)$ et pour tout $a \in X$, $\sigma(a)$ est algébrique.

Prop 36: L'ensemble des langages algébriques est stable par union, concaténation, passage à l'étoile, intersection avec un régulier, substitution algébrique, morphisme et morphisme inverse mais n'est pas fermé par complémentation ni par intersection.

c prouvé par séparation syntaxique et lexical

3) Théorème de Chomsky et Schützenberger

Thm 37: (Chomsky et Schützenberger) Un langage est algébrique si

$$L = \varphi(D_n^* \cap K) \text{ pour } n \in \mathbb{N}, \text{ un langage rationnel } K \text{ et un morphisme algébrique } \varphi (\forall u, \ell(u) \leq 1).$$

Lemme 38: Il existe un morphisme $\varphi: \Sigma_{3,2}^* \rightarrow \Sigma_{3,2}^* \text{ tq } D_n^* = \varphi^{-1}(D_2^*)$.

Corollaire 39: Un langage est algébrique ssi $L = \varphi(\varphi^{-1}(D_2^*) \cap K)$ pour φ morphisme et K langage rationnel.

Rq Les applications de la forme $X \mapsto \varphi(\varphi^{-1}(X) \cap K)$ sont les translations rationnelles et appartiennent naturellement aux de l'étude des langages rationnels.

IV) Sous-classes remarquables

L'ensemble \mathcal{A} des inclusions entre ces classes.

1) Langages rationnels

Prop 40: Tout langage rationnel est algébrique et l'inclusion est stricte.

2) Langages déterministes

Déf 41: Un langage algébrique est dit déterministe (resp. déterministe préfixe) ssi il est accepté par un AP déterministe par état final (resp. par pile vide).

Prop 42: Tout langage rationnel est déterministe et l'inclusion est stricte.

Prop 43: L'ensemble des langages algébriques déterministes est stable par complémentation et intersection avec un rationnel mais ni par union ni par intersection.

Prop 44: Un langage algébrique est déterministe-préfixe ssi il est déterministe et préfixe ($\forall n, m \in \mathbb{N}, m \neq n$ et pas préfixe strict).

Rq 45: Soit L un langage tq $L \cap \Sigma^* = \emptyset$. Alors L est préfixe.

pour bien comprendre il faut commencer par se demander de ce qu'il se passe avec D_1

3) Langages ombigus

Déf 46: Un langage est dit ombigu si toutes les grammaires algébriques qui l'engendrent le sont et non ombigu sinon.

Prop 47: Tout langage algébrique déterministe est non ombigu et l'inclusion est stricte.

Prop 48: L'inclusion des langages non ombigus dans les langages algébriques est stricte.

Prop 49: L'ensemble des langages algébriques non ombigus est stable par union disjointe, union avec un rationnel et intersection avec un rationnel mais pas par union.

IV) Problèmes de décision

1) Problèmes décidables

Th 50: La vacuité d'un langage algébrique (représenté par une grammaire algébrique) est décidable en temps $O(|G|)$.

Th 51: L'appartenance d'un mot w à un langage algébrique (représenté par une grammaire algébrique) est décidable en temps $O(|G||w|^3)$. en forme normale quadratique

2) Problèmes indécidables

Th 52: Les problèmes suivants sont indécidables: [DEV] (P. Padua)

- Pour deux grammaires, l'intersection des langages engendrés est-elle vide?
- Pour deux grammaires, engendrent-elles le même langage?
- Pour une grammaire engendrée-t-elle Σ^* ?
- Pour une grammaire, est-elle ombigu?

V) Applications

Déf 53: Soit G une grammaire. À chaque $A \rightarrow \alpha\beta \in P$, on associe $[A \rightarrow \alpha \cdot \beta]$ qu'on appelle item. On note It_G l'ensemble des items. Soit $S \in \Gamma$ une variable distinguée. On appelle automate des items l'automate dont l'alphabet d'entrée est Σ , dont l'alphabet de pile est It_G , dont les transitions sont

$$(E) [X \rightarrow \beta \cdot Y\gamma] \xrightarrow{\epsilon} [X \rightarrow \beta \cdot Y\bar{\gamma}] [Y \rightarrow \alpha] \text{ pour } Y \rightarrow \alpha \in P$$

$$(L) [X \rightarrow \beta \cdot a\gamma] \xrightarrow{a} [X \rightarrow \beta a \cdot \gamma]$$

$$(R) [X \rightarrow \beta \cdot Y\gamma] [Y \rightarrow \alpha] \xrightarrow{\epsilon} [X \rightarrow \beta Y \cdot \gamma]$$

dont l'état initial est $[S' \rightarrow \cdot S]$ et dont l'état final est $[S' \rightarrow S \cdot]$ où S' est une variable que l'on ajoute à Γ .

Thm 54: L'automate des items accepte $L_G(S)$. [DEV]

1) Analyse descendante (LL(k))

On part de S et on essaye d'appliquer des règles pour arriver au mot en faisant des dérivations gauches.

"Déf 55": Une grammaire est dite LL(k) si il est possible de savoir quelle règle utiliser pour dériver le non terminal le plus à gauche en connaissant les k prochains symboles du mot.

2) Analyse ascendante (LR(k))

On part du mot et on essaye d'appliquer les règles $X \rightarrow \alpha$ pour réduire le mot en S (en appliquant l'inverse de dérivations droites).

"Déf 56": LR(k) si on sait quand effectuer une lecture ou une réduction avec les k premières lettres après la partie du mot à potentiellement réduire.

Thm 57: Un langage algébrique est engendré par une grammaire LR(0) si il est déterministe-préfixe.

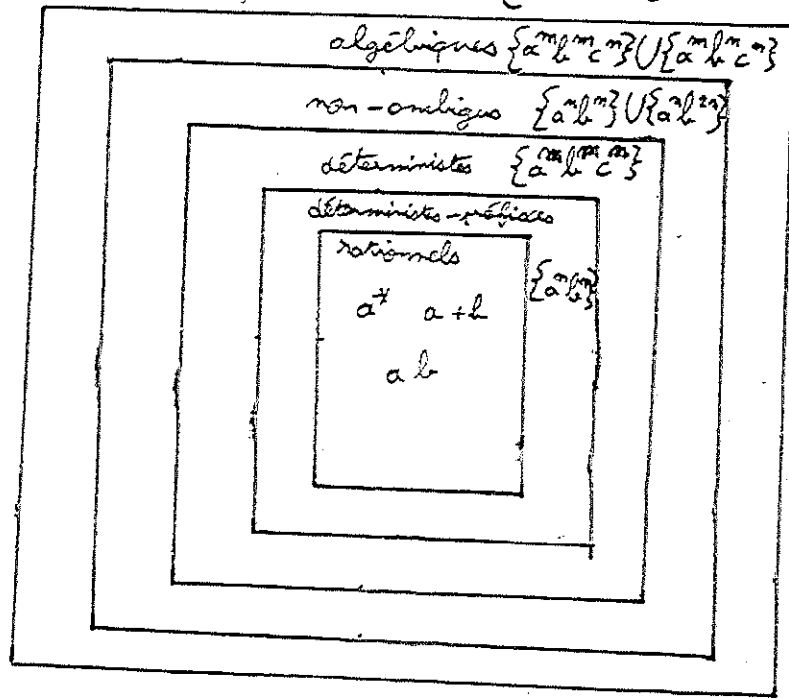
Thm 58: Un langage algébrique est engendré par une grammaire LR(k), $k \geq 1$ si il est engendré par une grammaire LR(1) si il est déterministe.

pas sûr énum.

récurivement énumérables

contexte sensible

$\{a^m b^m c^m\}$



Annexe A

Question : Décidables ?

- Algébrique? si on a une grammaire pas
- Régulier? si on a un algébrique?
- Pour L algébrique. Est-ce que L est déterministe? décidable

[Carton]. Les mots de la pile : langage régulier.
 Donc si on peut marquer le langage de pile est $< \infty$,
 alors le langage alg. est régulier. (nombre fini de configurations \Rightarrow automate fini)

Compléments : langages récurivement énumérables; thm de REISZ.

Rq: • Il manque au départ, une motivation, un contexte général.

• Il faut être enthousiaste dans la présentation du plan; articuler le discours ...

• Quand on étudie la complexité, il faut qu'il y ait qqch derrière, un but, un bilan...
Qu'il faut mettre en avant.

• Contre-exemple: lemme d'Ogden pas CS.

(Rq: Est-ce qu'un langage est algébrique // lemme d'Ogden: pas décidables?)

• Ptes de clôture: on peut s'en servir pour des raisonnements par l'absurde, ...

On ne sait pas dire vite $L_1 \in L_2$ pour les algébriques (Pb de décision III)

↳ du coup on peut commencer à distinguer des sous-classes. Si L_1 rationnel, L_2 déterministe, mieux ?

$L \subset R$? \rightarrow fautive: on décide $L \cap R^c = \emptyset$?

↳ Sous-classes remarquables: on veut améliorer les pts de clôture et/ou les problèmes de décision