

Heur (1^{er}) Alg (1^{er})

Interprétation: Les langages algébriques sont 3 fois les classes de la hiérarchie de Chomsky et sont utiles dans l'analyse syntaxique des langages de programmation, notamment dans la reconnaissance des lignes sources.

I Grammaires algébriques. Définitions et propriétés

1) Définitions de Grammaires

Def: Une grammaire algébrique est un quadruplet $G = (\Sigma, V, P, S)$ où:

- Σ est un alphabet de terminaux fini.
- V est un ensemble fini de non-terminaux.
- $P \subset V \times (\Sigma + V)^*$ un ensemble fini de règles de dérivation.
- $S \in V$ un axiome.

* Un mot w est dérivé au n -ième pas si $w \in X_n$ et $x = x'xy$ où $(x, y) \in P$, mot $x \rightarrow w$.

* Un mot w est dérivé si $w \in X_n$ pour un n quelconque. L'ensemble des mots dérivés est noté $L(G)$.

* Le langage engendré par G est l'ensemble des mots dérivés par G , noté $L(G)$.

* Un langage est dit algébrique s'il est engendré par une grammaire algébrique.

Exemple: Soit $G = (\Sigma, V, P, S)$ la grammaire telle que:

$$V = \{S, A, B\}, \Sigma = \{a, b\}, P = \{S \rightarrow A, S \rightarrow B, A \rightarrow aA, B \rightarrow aB, A \rightarrow \epsilon, B \rightarrow \epsilon\}$$

$$L(G) = \{a^n, n \geq 0\} \cup \{b^n, n \geq 0\}$$

1^{er} exemple: $a^n \in L(G) \rightarrow a^n \in a^* \rightarrow a^n \in L(G)$

Théorème fondamental: Tout $G = (\Sigma, V, P, S)$ est une grammaire algébrique, où $\epsilon \in (\Sigma + V)^*$, $v \in (\Sigma + V)^*$.

Un langage qui ne satisfait pas $L_1, L_2 = L_1$.

gammes gram alg

Alors: $A \rightarrow v \Leftrightarrow A_1 \xrightarrow{1} v, A_2 \xrightarrow{2} v, \dots, A_n \xrightarrow{n} v$

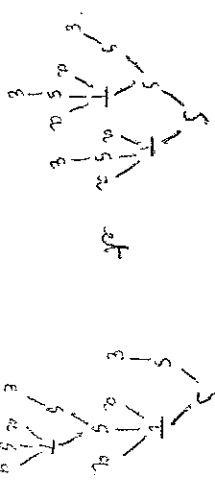
2) Notion de dérivation ambiguë

Def: Soit $G = (\Sigma, V, P, S)$ une grammaire algébrique

* Un mot w est dit ambiguë si son arbre de dérivation n'est pas unique. Ex: $\Sigma = \{a, b\}$, $V = \{A, B\}$, $P = \{A \rightarrow aA, A \rightarrow aB, B \rightarrow aB, B \rightarrow aA\}$.

* Les arbres de dérivation d'un mot w sont notés T_w . On dit qu'un mot w est ambiguë si il existe un mot w qui a deux arbres de dérivation distincts.

Exemple: Les grammaires de départ pour $\{S \rightarrow ST+T, T \rightarrow a, S \rightarrow \epsilon\}$ est ambiguë.



noté avec arbre de dérivation de a

3) Normes de dérivation

* Forme normale: Une grammaire (Σ, V, P, S) est dite normale si

$$\forall A \in V, L(A) \neq \emptyset \text{ et } \forall A \in V, \exists u, v \in (\Sigma + V)^*, S \xrightarrow{*} uAv$$

Propriété: Toute grammaire normale est équivalente à une grammaire normale engendrant le même langage.

* Forme normale: Une grammaire $G = (\Sigma, V, P, S)$ est dite normale si P ne contient pas de règles $A \rightarrow \epsilon$ ou $A \rightarrow B, A, B \in V$.

preuves constructives (algèbre)

Forme normale de Greibach

ENC de dérivée $S \rightarrow W$ et dérivement de degré m ($m=1$)

Propriété : Pour toute grammaire G , il existe une grammaire Type G'

telle que $L(G) = L(G') \setminus \{\epsilon\}$

• Une grammaire de Greibach : Une grammaire est nommée grammaire normale de Greibach si toutes les règles sont de la forme :

$$A \rightarrow A_1 A_2 \dots A_n \epsilon V \text{ où } A \rightarrow \epsilon, a \in \Sigma$$

Propriété : Pour toute grammaire G , il existe une grammaire normée de Greibach G' telle que $L(G) = L(G') \setminus \{\epsilon\}$

Exemple : La grammaire G définie par $S \rightarrow aSb + \epsilon$ est réduite,

Tout grammaire associée est $S \rightarrow aSb + \epsilon$ et sa forme de Greibach est $\left\{ \begin{array}{l} S \rightarrow AS_1 + A0 \quad A \rightarrow a \\ S_1 \rightarrow SB \quad 0 \rightarrow \epsilon \end{array} \right\}$

II Propriétés des langages algébriques

1) Forme d'Uspen

Théorème : Pour tout langage algébrique L , il existe $K \in N$ tel que pour tout $f \in L$ il y a un mot K lettre distinguée qui se factorise en $f = xuxy$ avec :

- $\forall u \in N$ et $\forall x, y \in L$
- un mot x, u, y, A où A est uniquement des lettres distinguées
- u et v contiennent des lettres de K lettre distinguée

Application : $\{a^n b^m c^n, m \geq 0\}$ n'est pas algébrique

2) Propriétés de clôture

Propriété : La classe des langages algébriques est close par union, concaténation et étoile

Propriété : L'intersection de deux langages algébriques et le complémentaire d'un langage algébrique ne sont pas algébriques

Exemple pour l'intersection :

$$\{a^n b^m, m \geq 0\} \cap \{a^m b^n, m, n \geq 0\} = \{a^n b^n, n \geq 0\}$$

3) Problèmes décidables et indécidables

Propriété : Il existe un algorithme permettant de vérifier si un langage algébrique est vide

Il existe un algorithme permettant de vérifier si un mot appartient à un langage algébrique

C'est l'algorithme de Earle Younger, Karasani **DVPT**

Propriété : Les problèmes suivants sont indécidables :

- deux langages algébriques sont-ils égaux ?
- l'intersection de deux langages algébriques est-elle vide ?
- un langage algébrique est-il universel ?

Revenir au problème de l'intersection de Post

4) Théorème de Parikh

Théorème : Tout langage algébrique est commutativement équivalent à un langage rationnel **DVPT**

III Automates à pile

1) Définition : Un automate à pile est une machine à états finis avec un objet de la forme $\{Q, \Sigma, \Gamma, A, q_0, \delta, F\}$ où :

- Q ensemble fini d'états
- Σ un alphabet d'entrée
- Γ un alphabet de pile
- A un ensemble de mots initiaux
- $q_0 \in Q$ un état initial
- $F \subseteq Q$ un ensemble de mots finaux
- δ une fonction de transition de la forme $(q, Z) \xrightarrow{a} (q', \gamma)$ où $a \in \Sigma$ et $\gamma \in \Gamma^*$

Un mot est accepté si il existe une suite de transitions ayant un terme de F et un mot dans un état de F

* Le langage reconnu par l'automate est énumérable des mots

[C] p 94

[C] p 153

[C] p 86

[C] p 104

[C] p 93

[C] p 92

Developpement : Théorème de Parikh.

Def: Deux mots w et w' sont anagrammes si

$w = a_1 \dots a_n$ et il existe σ une permutation de $\{1, \dots, n\}$ telle que $w' = a_{\sigma(1)} \dots a_{\sigma(n)}$

L'ensemble des anagrammes de w est noté \bar{w} .

On appelle image commutative d'un langage L , notée \bar{L} , l'ensemble $\{\bar{w}, w \in L\}$.

Pour deux langages L et M , on dit qu'ils sont commutativement équivalents si $\bar{L} = \bar{M}$.

Th: (Parikh) Tout langage algébrique L est commutativement équivalent à un langage rationnel R .

Pour prouver cela on va s'appuyer sur beaucoup d'outils et de lemmes intermédiaires.

Pour d'abord on définit les opérations de substitution.

pour une grammaire $G = (A, V, P)$, $V = \{x_1, \dots, x_n\}$ et un n -uplet de langages $L = (L_1, \dots, L_n)$ on définit par induction:

$$\bullet \quad \varepsilon(L) = \{\varepsilon\}; \quad a(L) = \{a\} \text{ pour } a \in A; \quad X_i(L) = L_i$$

$$xy(L) = x(L)y(L) \text{ pour } x, y \in (A+V)^*; \quad K(L) = \bigcup_{w \in K} w(L) \text{ pour } K \subseteq (A+V)^*$$

On introduit ensuite le système d'équations associé à une grammaire.

def: $G = (A, V, P)$, $V = (x_1, \dots, x_n)$ on a alors le système $S(L)$

formé des n équations

$$L_{x_i} = \sum_{x_i \rightarrow w} w(L) \quad \text{pour } i \leq n$$

On note alors L_G le n -uplet $(L(x_1), \dots, L(x_n))$

et on a une proposition "évidente"

prop: $\forall w \in (A+V)^*$ on a $L_G(w) = w - (L_G)$

preuve par induction sur les substitutions.

Ensuite on se donne le lemme:

lemme: $G = (A, V, P)$ et L solution de $S(G)$

Si $w, w' \in (A+V)^*$ vérifiant $w \xrightarrow{*} w'$ alors $w(L) \supset w'(L)$

preuve: Il suffit de démontrer le résultat pour une dérivation et le

lemme s'en déduit.

si $w = u x_i v \rightarrow u y v = w'$, $x_i \rightarrow y$ règle de la grammaire

L est solution de $S(G)$ donc $L_i = \sum_{x_i \rightarrow \delta} \delta(L)$ et donc $L_i \supset y(L)$

d'où $w(L) \supset w'(L)$

■

de ce lemme on déduit:

prop: $G = (A, V, P)$, $V = \{x_1, \dots, x_n\}$ alors L_G est la relation minimale de $S(G)$ pour l'inclusion.

preuve: - on vérifie que c'est bien une solution:

$$L_G(x_i) = \sum_{x_i \rightarrow \alpha} L_G(\alpha) = \sum_{x_i \rightarrow \alpha} \alpha(L_G) \text{ d'après la propriété précédente}$$

- par minimalité découle elle du lemme précédent.

■

On se penche ensuite sur le cas des grammaires propres!

i.e. une grammaire dont tous les L_i ne contiennent pas ϵ .

prop: G propre alors L_G est l'unique relation propre de $S(G)$

preuve: G propre donc L_G propre et d'après la prop. précédente c'est une solution de $S(G)$.

On introduit la relation sur les langages \mathcal{L} avec $\mathcal{L} \in \mathcal{P}(W)$

telle que $K \preceq K'$ si $\{w \in K, |w| \leq l\} = \{w \in K', |w| \leq l\}$.

On étend cette relation aux n -uplets, comparés par composants.

Soient alors L et L' deux solutions propres de G .

On va montrer que $L \stackrel{L'}{\sim} L'$ pour tout L et donc $L = L'$.

par récurrence, pour $L = \emptyset$, L et L' sont donc en accord.

maintenant supposons $L \stackrel{L'}{\sim} L'$ pour $L \geq 0$ et montrons $L \stackrel{L+1}{\sim} L'$ c'est à dire

$$L_{i+1} \stackrel{L+1}{\sim} L'_{i+1} \text{ pour } 1 \leq i \leq n.$$

On introduit $S_i = \{w \mid x_i \rightarrow w\}$, alors par définition $L_{i+1} = S_i(L)$

On veut donc montrer $S_i(L) \stackrel{L+1}{\sim} S_i(L')$.

Soit donc $1 \leq i \leq n$ fixé et $w \in S_i$, $w = w_1 \dots w_p$ $w_j \in (A+V)$

$p \neq 0$ car la grammaire est propre.

Soit $u \in w(L)$ de longueur inférieure à $L+1$ alors $u = u_1 \dots u_p$

avec chaque u_j étant soit une lettre de A soit, si $u_j = x_k$ est un mot dans L_k .

de plus $|u_j| \leq L$ sinon $p=1$ ce qui voudrait dire que $x_i \rightarrow x_k \in P$

impossible car la grammaire est propre. Donc par hypothèse $u_j \in L'_k$

et donc $u \in w(L')$ d'où $L_{i+1} \subseteq L'_{i+1}$, l'inclusion inverse se montre

par symétrie des rôles. et on obtient le résultat. \square

On définit alors le système d'équations en commutatif.

cela part que si $\bar{u} = \bar{v}$ en commutatif $\bar{u}(L) = \bar{v}(L)$ et par

extension aux langages si $\bar{T} = \bar{K}$, $\bar{T}(L) = \bar{K}(L)$ on peut donc

définir de manière cohérente $\bar{S}(L)$ comme

$$\bar{L}_i = \sum_{x_i \rightarrow w} \bar{w}(L) \text{ pour } 1 \leq i \leq n$$

prop: \bar{L}_0 est l'unique solution propre de $\bar{S}(L)$ si la grammaire est propre.

La grammaire est la même que la précédente.

on construit grammaire à $n+1$ variables X_0, X_1, \dots, X_n
 et que l'on a le résultat pour les nombres inférieurs.

les règles de G s'écrivent

$$X_i \rightarrow S_i(X_0, \dots, X_n) \quad \text{pour } 0 \leq i \leq n.$$

ou chaque S_i est rationnel sur $A+X$ avec $X = \{X_0, \dots, X_n\}$.

soit $X' = X \setminus \{X_0\}$, on introduit la grammaire G_0 en considérant
 X_1, \dots, X_n comme des lettres terminales et on se propose que les
 règles associées à X_0 dans G se

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

que d'après le résultat précédent on a $R_0(X')$ solution rationnel du système $\overline{S(G_0)}$

donc
$$R_0(X') = S_0(R_0(X'), X')$$

soit alors la grammaire G' obtenue en remplaçant X_0 dans
 les autres règles par $R_0(X')$ donc

$$G' = (A, X', \{X_i \rightarrow S_i(R_0(X'), X') \mid 1 \leq i \leq n\})$$

La grammaire G' vérifie les hypothèses de récurrence donc il existe


un langage rationnel R_1, \dots, R_n sur A tel que $R' = (R_1, \dots, R_n)$

est solution du système $\overline{S(G')}$ donc $R_i = S_i(R_0(R'), R')$ $1 \leq i \leq n$

Il reste à vérifier $R = (R_0(R'), R_1, \dots, R_n)$ est solution de $\overline{S(G)}$.

En substituant X' par R' dans l'équation $R_0(X') = S_0(R_0(X'), X')$

on obtient une équation complétant les n équations $R_i = S_i(R_0(R'), R')$

et ceci termine la preuve. 

En posant enfin le théorème de Zariski :

soit L algébrique, $\mathfrak{a} = (A, V, P)$, $L = L_{\mathfrak{a}}(X, a)$

$L_{\mathfrak{a}}$ est l'unique solution de SLW quitte à se restreindre à une sous-algèbre propre.

On a aussi d'après le th précédent une solution rationnelle. d'où le résultat par unicité.

Si au départ $L_{\mathfrak{a}}$ est gas propre on applique le résultat à $L \setminus \mathbb{Q}$.

C QFD

ref: Cartan p. 23 à 30

Dir: algorithme CYK.

Locke, Kasami et Younger.

Prend en entrée un mot et une grammaire en forme normale quadratique G et décide en temps cubique si w est engendré par G .

Il utilise la programmation dynamique.

$$G = (A, V, P) \quad V = \{S_0, \dots, S_{m-1}\}$$

Soit $w = a_1 \dots a_m$ alors pour tous entiers $1 \leq i \leq j \leq m$

on note $w[i, j]$ le facteur $a_i \dots a_j$ de w . On note

alors $E_{i, j}$ l'ensemble des variables $S \in V$ telles que

$w[i, j] \in L_G(S)$. L'algorithme calcule tous les ensembles $E_{i, j}$.

Le mot w appartient alors à $L_G(S_0)$ si $S_0 \in E_{1, m}$.

L'initialisation se fait pour $i = j$. alors $w[i, j]$ est réduit

à la lettre a_i . Comme la grammaire est sous forme normale

quadratique, $S \in E_{i, i}$ si et seulement si $S \rightarrow a_i$ est une règle de G .

Ensuite, pour $i < j$, $S \in E_{i, j}$ si et seulement si il existe une

règle $S \rightarrow S_1 S_2$ de G et un entier $i \leq k < j$ tel que $w[i, k] \in L_G(S_1)$

et $w[k+1, j] \in L_G(S_2)$. (découle du lemme fondamental).

Donc $E_{i, j}$ peut être calculé à partir des ensembles $E_{i, k}$ et $E_{k+1, j}$.

On calcule donc les $E_{i, j}$ par récurrence sur la différence $j - i$.

$i \backslash j$	1	i		n
1	$E_{1,1}$			
		$E_{2,2}$		
		$E_{3,3}$		
			$E_{i,i}$	
j			$E_{i,j}$	$E_{j,j}$
n	$E_{1,n}$			$E_{n,n}$

les lignes sup j et les col i m

$$E_{i,j} = U \{S\}$$

$S \rightarrow s_1 s_2$
 tel que
 $s_1 \in E_{i,k}$ et
 $s_2 \in E_{k,j}$

On a donc l'algorithme:

C-Y-K pour $(w = a_1 \dots a_n) =$

pour $1 \leq i \leq j \leq n$ faire

$E_{i,j} \leftarrow \emptyset$ (*initialisation*)

fin

pour i de 1 à n faire

 pour toute règle $S \rightarrow a$ faire

 si $a_i = a$ alors $E_{i,i} \leftarrow E_{i,i} \cup \{S\}$ (*cas $i=j$ *)

 fin

 pour $1 \leq d \leq n-i$ faire (*cas $i < j$; $d = j - i$ *)

 pour $1 \leq i \leq n-d$ faire

 pour $i \leq k \leq i+d$ faire

 pour toute règle $S \rightarrow s_1 s_2$ faire

 si $s_1 \in E_{i,k}$ et $s_2 \in E_{k,i+d}$ alors

$E_{i,i+d} \leftarrow E_{i,i+d} \cup \{S\}$

 fin

 fin

 fin

 fin

Count \rightarrow ajouter ex