

多重文脈自由文法のある部分クラスに対する 効率の良い構文解析法について

2 E - 6

中西 隆一* 関 浩之* 藤井 譲** 嵩 忠雄*

*大阪大学基礎工学部情報工学科 **大阪大学教養部

1. まえがき

自然言語の構文記述のための形式文法として、多重文脈自由文法（以下mcfgと略記）が提案されている⁽³⁾。mcfgでは、respectively文や倒置文のような互いに割り込んだ構文を自然に記述することができる。その生成能力は、文脈自由文法（以下cfgと略記）や、同じく自然言語の構文記述のためにPollardによって導入されたhead grammar⁽²⁾よりも真に大きく、文脈規定文法よりも真に小さい。また、mcfg Gに対して、入力系列 α がGの生成する言語に属するかどうかを判定する時間計算量 $O(n^2)$ (n は入力系列長) のアルゴリズムが提案されている⁽⁴⁾。ここで、 e はGのみに依存して定まる定数であるが、その上界は知られていない。

本稿では、mcfgのある部分クラスに対する時間計算量 $O(n^2)$ の構文解析アルゴリズムを提案する。この部分クラスに属するmcfgの生成する言語のクラスは、あいまいでないcfgの生成する言語をすべて含み、言語 $\{a_1^p b_1^q a_2^r b_2^s \dots a_m^p b_m^q | p, q \geq 1\}$ ($m \geq 1$)、 $\{w^p w^q | p, q \in \{0, 1\}^*\}$ 等も含む。

2. 多重文脈自由文法

正整数 m に対して、 m -多重文脈自由文法 (m -mcfg) と略記。 m を明示しないときはmcfg) は次の[M1]～[M5]を満たす5字組 $G = (N, O, F, P, S)$ として定義される⁽³⁾。

[M1] N は非終端記号の有限集合。各 $A \in N$ について、 m 以下の正整数 $d(A)$ が定まっているとする。

[M2] $O = \bigcup_{i=1}^m (T^*)^i$ 。 T は N と互いに素である終端記号の有限集合。

[M3] $F = \bigcup_{0 < r \leq m; 0 \leq q; 0 < d_1, \dots, d_q \leq m} F\{d_1, d_2, \dots, d_q; r\}$ 。

$$0 < r \leq m; 0 \leq q; 0 < d_1, \dots, d_q \leq m$$

ここで、 $F\{d_1, d_2, \dots, d_q; r\}$ は、次の条件 [f1], [f2] をともに満たす $(T^*)^{d_1} \times (T^*)^{d_2} \times \dots \times (T^*)^{d_q}$ から $(T^*)^r$ への写像からなるある集合。 $x_i = (x_{i1}, x_{i2}, \dots, x_{id_i})$ を $(T^*)^{d_i}$ の上の変数とし、 $X \triangleq \{x_{ij} | 1 \leq i \leq q, 1 \leq j \leq d_i\}$ とおく。

[f1] f の第 h 成分 f^h ($1 \leq h \leq r$) の値は T^* の定系列と X に属するいくつかの変数の連接で表される。すなわち、

$$f^h[x_1, x_2, \dots, x_q] \triangleq \alpha_{h0} z_{h1} \alpha_{h1} z_{h2} \dots z_{hv_h} \alpha_{hv_h} \quad (*1)$$

ここで、 $\alpha_{hk} \in T^*$ ($0 \leq k \leq v_h$)、 $z_{hk} \in X$ ($1 \leq k \leq v_h$)。

[f2] X の各変数 x_{ij} ($1 \leq i \leq q, 1 \leq j \leq d_i$) に対して、 $h = 1$ から q までにわたって、(*1)の右辺に x_{ij} が現れる総回数は 1 以下である。

[M4] P は $A_\theta \rightarrow f[A_1, A_2, \dots, A_q]$ の形の規則の有限集合。ここで、

$A_i \in N$ ($0 \leq i \leq q$)、 $f \in F\{d(A_1), d(A_2), \dots, d(A_q); d(A_\theta)\}$ 。特に、 $q=0$ のときの規則を終端規則と呼び、それ以外の規則を非終端規則と呼ぶ。

[M5] $S \in N$ 。 S は始記号と呼ばれる。

$A \in N$ について、 $L_G(A)$ を次の条件 [L1] と [L2] を満たす最小集合として定義する。

[L1] P に終端規則 $A \rightarrow \alpha$ があれば、 $\alpha \in L_G(A)$ 。

[L2] $\alpha_i \in L_G(A_i)$ ($1 \leq i \leq q$) かつ $A \rightarrow f[A_1, A_2, \dots, A_q] \in P$ であるとき、 $f[\alpha_1, \alpha_2, \dots, \alpha_q] \in L_G(A)$ 。

[M4] から、 $L_G(A) \subseteq (T^*)^{d(A)}$ が成り立つ。 $L_G(A)$ の元を A から生成される連句という。 $L(G) \triangleq L_G(S)$ と定義し、 $L(G)$ を mcfg G の生成する言語と呼ぶ。

[例1] $G = (\{S, A, B\}, T \cup T^2, F, P, S)$ 、但し、 $T = \{a, b, c, d\}$ 、 $P = \{S \rightarrow f_1[A, B], A \rightarrow f_2[A] | (a, c), B \rightarrow f_3[B] | (b, d)\}$ ($A \rightarrow \alpha$ | β は $A \rightarrow \alpha$ 、 $A \rightarrow \beta$ を表す)、 $f_1[(x_1, x_2), (y_1, y_2)] \triangleq x_1 y_1 x_2 y_2$ 、 $f_2[(x, y)] \triangleq (ax, cy)$ 、 $f_3[(x, y)] \triangleq (bx, dy)$ とする。 $L_G(A) = \{(a^p, c^p) | p \geq 1\}$ 、 $L_G(B) = \{(b^q, d^q) | q \geq 1\}$ 、 $L(G) = L_G(S) = \{a^p b^q c^p d^q | p, q \geq 1\}$ である。□

G における導出木を次の[T1]～[T3]のように定義する。

[T1] 終端規則 $A \rightarrow \alpha$ に対して、根の節点（ラベル A ）がただ一つの子節点（ラベル α ）をもつ木は α の導出木である。

[T2] α の導出木を τ_i ($1 \leq i \leq q$)、その根のラベルを A_i とし、 $A \rightarrow f[A_1, A_2, \dots, A_q] \in P$ とする。根のラベルが A (必要なら、 A の代りに適用規則名)、根から出る枝数が q 、 i 番目の枝の端点（ラベル A_i ）以下に τ_i を部分木としてもつ木は $f[\alpha_1, \alpha_2, \dots, \alpha_q]$ の導出木である。

[T3] それ以外に導出木はない。

次の補題が成り立つ。

[補題1⁽³⁾] m -mcfg G が与えられたとき、 $L(G)$ を生成し、次の情報無損失条件 [f3]⁽³⁾、および条件 [N1] をともに満たす m -mcfg を構成できる。

[f3] X のどの変数も、ちょうど一つの h について、式 (*1) の右辺にちょうど一回現れる。

[N1] 各規則の右辺には、同じ非終端記号は 2 回以上現れない。

以下では、上の条件 [f3], [N1] をともに満たすような mcfgのみを考える。

3. mcfgから派生するcfg

任意の mcfg $G = (N, O, F, P, S)$ に対して、次の(1), (2)を満たす

An Efficient Parsing Algorithm for a Subclass of Multiple Context-free Grammars

Ryuichi Nakanishi, Hiroyuki Seki, Mamoru Fujii and Tadao Kasami

Osaka University

すcfg $G' = (N', T, P', S^{[1]})$ を, G から派生するcfgという.

(1) $N' = \{A^{[h]} \mid A \in N, 1 \leq h \leq d(A)\}$.

(2) 規則 $\rho : A \rightarrow f[A_1, \dots, A_n]$ 右辺の $f \in P$ の定義式(*1)において, $Z_{hk} = x_p(hk)_{q(hk)} \quad (1 \leq h \leq d(A), 1 \leq k \leq v_h)$ とする.

$\rho : A \rightarrow f[A_1, \dots, A_n] \in P \quad (=)$

各 $h \quad (1 \leq h \leq d(A))$ について, $R_h : A^{[h]} \rightarrow$

$\alpha_{h0} A_{p(h1)}^{[q(h1)]} \dots A_{p(hv_h)}^{[q(hv_h)]} \alpha_{hv_h} \in P'$.

上で, G' の規則の組 $(R_1, \dots, R_{d(A)})$ は G の規則 ρ から派生するという.

[例2] 例1のmcfg G から派生するcfgは, $G' = (S^{[1]}, A^{[1]}, A^{[2]}, B^{[1]}, B^{[2]}, T, P', S^{[1]})$ である. ここで, $P' = \{S^{[1]} \rightarrow A^{[1]}B^{[1]}A^{[2]}B^{[2]}, A^{[1]} \rightarrow aA^{[1]}|a, A^{[2]} \rightarrow cA^{[2]}|c, B^{[1]} \rightarrow bB^{[1]}|b, B^{[2]} \rightarrow dB^{[2]}|d\}$ である. \square

cfg G_1 に対して, G_1 の生成する言語を $L(G_1)$ と書く. G を任意のmcfgとし, G' を G から派生するcfgとする. 定義より, $L(G) \subseteq L(G')$ であることが容易に示せる.

t を根のラベルが始記号であるような G' における任意の導出木とする. t における連句頂点ベクトルおよび, t に対応する G における導出木プレフィクスを以下の[C1]～[C3]のように定義する. 高さ k の導出木プレフィクスは, G' の終端記号の組または非終端記号をラベルとしてもつ頂点からなる高さ(根から葉頂点までの道の最大長) k の木であり, 非終端記号をラベルとしてもつ深さ i ($1 \leq i < k$) の頂点には, t における深さ i の連句頂点ベクトルが対応づけられる.

[C1] 唯一の頂点 v_0 からなる木は, t に対応する G における高さ0の導出木プレフィクスである. 但し, t の根を r とするとき, v_0 には連句頂点ベクトル(r)が対応づけられる.

[C2] t を t に対応する高さ $k-1$ の導出木プレフィクスであるとする. v を t におけるラベルが終端記号でない深さ $k-1$ の任意の頂点とし, v に対応づけられている連句頂点ベクトルを (v_1, \dots, v_n) とする. v_i ($1 \leq i \leq n$) における適用規則を R_i とし, G' の規則の組 (R_1, \dots, R_n) を派生すると仮定する(条件[N1]より, このような規則 ρ は高々一つしか存在しない). ρ を非終端規則とする. 各 h ($1 \leq h \leq s$) について, v_1, \dots, v_n の子頂点で, $A_h^{[j]} \quad (1 \leq j \leq d(A_h))$ をラベルにもつ頂点を u_{hj} とするとき(条件[f3]より, このような頂点は必ず存在し唯一つに定まる), $U_h = (u_{h1}, \dots, u_{hd(A_h)})$ は深さ k の連句頂点ベクトルである.

また, 上の各頂点 v に対して, 次の(a), (b)を行って得られる木は, t に対応する高さ k の導出木プレフィクスである. (a) ρ が非終端規則のとき, 各 h ($1 \leq h \leq s$) について, v の第 h 子として, 新しい頂点 u_h を追加する. 但し, u_h には上の連句頂点ベクトル U_h を対応づける.

(b) ρ が終端規則のとき, 終端記号列の組 $f[]$ をラベルとしてもつ頂点を, v の唯一の子頂点として加える.

[C3] 連句頂点ベクトルおよび導出木プレフィクスは上の[C1], [C2]で定義されるもの以外に存在しない. \square

t を G' における α の導出木とし, k_t を t の高さとするとき, 上の定義から, t に対応する高さ k_t の G における導出木プレ

フィクスは, G における α の導出木である.

4. 構文解析アルゴリズム

G を任意のmcfgとし, G から派生するcfgを G' とする. G' があいまいでない⁽¹⁾と仮定する. 以下で述べるアルゴリズム PARSE(α)では, 与えられた系列 α に対し, まず, G' に対して(あいまいでないcfgに対する構文解析法を用いて) α を構文解析し, その結果得られる導出木 t に対して, t に対応する G における導出木プレフィクスを根頂点からトップダウンに求めていく. 以下にアルゴリズムを示す.

[手続きPARSE(α)]

(入力) α : 終端記号列.

(出力) $\alpha \in L(G)$ のとき G における α の導出木,
 $\alpha \notin L(G)$ のとき "NO".

(手続き) 次の(1)～(3)を順に行う.

(1) α をcfg G' に対して構文解析し, その結果が $\alpha \notin L(G')$ であれば, "NO"を出力して停止. $\alpha \in L(G')$ のとき, 得られた導出木を t とし, t の高さを k_t とする.

(2) $k=0$ から k_t まで順に, 以下の(2a)を行う.

(2a) 導出木プレフィクスの定義に従って, ($k \geq 1$ のときは高さ $k-1$ の導出木プレフィクスを用いて,) t に対応する G における高さ k の導出木プレフィクスを求める. もしそれが存在しなければ, "NO"を出力して停止する.

(3) G における α の導出木が得られたのでそれを出力する. \square

PARSE(α)の時間計算量を評価する. 次の性質が成り立つ.
[性質2] G' をあいまいでないcfgとし, $\alpha \in L(G')$ とする. G' における α の導出木の頂点の個数は $O(n)$ ($n = |\alpha|$) である. (証明は容易) \square

G' があいまいでないという仮定から, PARSEのステップ

(1) は $O(n^2)$ ($n = |\alpha|$) で実行できる⁽¹⁾. ステップ(2a)において, 高さ $k-1$ の導出木プレフィクスから, 深さ k の連句頂点ベクトルおよび高さ k の導出木プレフィクスを求める操作は, t における深さ $k-1$ と k の頂点の個数の和のオーダで求められる. 従って, ステップ(2)全体は, t の頂点の個数のオーダで実行できる. これは性質2より $O(n)$ である. 以上より, 上のPARSEの時間計算量は $O(n^2)$ である.

文 献

- (1) J. Hopcroft and J. D. Ullman : "Introduction to Automata Theory, Languages, and Computation", Addison-Wesley, Reading, Mass., U.S.A. (1979).
- (2) C. J. Pollard: "Generalized phrase structure grammars, head grammars, and natural language", Ph.D. dissertation, Stanford University (1984).
- (3) 嵩, 関, 藤井: "一般化文脈自由文法と多重文脈自由文法", 信学論(D), J71-D, 5, pp. 758-765 (昭63-05).
- (4) 嵩, 関, 藤井: "Head Languageおよび多重文脈自由言語の所属問題", 信学論(D), J71-D, 6, pp. 935-941 (昭63-06).