

# LRパーサを用いた文字列置換アルゴリズム

2C-2

津田 和彦<sup>†</sup>

中村 雅巳<sup>†</sup>

青江 順一<sup>††</sup>

<sup>†</sup>住友金属工業(株)

<sup>††</sup>徳島大学

## 1. はじめに

著者は、これまで日本語文書の短縮法についての研究を行ってきた(1)(2)。著者が提案している文書短縮法は、より少ない文字数で表現できる形態素列を短縮規則として定義し、文書中より短縮規則が適用できる部分を検索・置換することで実行される。しかし、短縮規則が増加すると文書中の特定部分で複数の短縮規則が適用でき一般的なストリング・パターン・マッチング手法では、適用可能な全ての短縮規則を検出できないという問題点が生じてきた。

そこで、文脈依存文法の構文規則とこの短縮規則が類似し、また、文脈依存言語に対するLRパーサの還元動作と文書短縮のための置換動作が類似している点に着目して、バックトラックのない文字列置換アルゴリズムを提案する。本手法で得られた文字列置換マシンはバックトラックが不要であるという特徴に加え、従来のコード最適化マシンにおいて不可能であったより良い規則の選択が可能となる。

## 2. 文脈依存文法と置換規則

### 2.1 文脈依存文法

文脈依存文法(context-sensitive grammar)を  $G=(V_N, V_T, P, S)$  で表す。ここで、 $V_T, V_N$  はそれぞれ非終端語(terminals), 終端語(nonterminals)の有限集合、 $P$  は  $a \rightarrow \beta$  なる生成規則(productions)の有限集合である。 $S$  は  $V_N$  の要素で始記号(start symbol)を表す。 $V = V_N \cup V_T$  とし、 $V$  上のすべての有限長の記号列の集合を  $V^*$  とする。 $V^+$  は  $V^*$  から空集合  $\epsilon$  を除いた集合を表し、本章では特に断らない限り、 $V, V^+, V^*$  の要素をそれぞれ  $X, Y, Z, x; a, \beta, \gamma, \delta; \dots$  で表す。

### 2.2 置換規則

文脈依存文法の構文規則  $A_{p1} A_{p2} \dots A_{pnp} \rightarrow B_{p1} B_{p2} \dots B_{pmp}$  ( $0 \leq mp < np$ ) の左右を交換したものが置換規則と見なすことができる。このようにして得られた第  $p$  番目の置換規則を、

$$B_{p1} B_{p2} \dots B_{pmp} \rightarrow A_{p1} A_{p2} \dots A_{pnp} \quad (0 \leq mp < np)$$

と定義する。

以上で定義された置換規則の集合を  $P_{OPT}$  で表し、 $P_{OPT}$  の最大値を  $v$  とする。そして、LRパーサの構成法をこの置換規則  $P_{OPT}$  に応用して得られたマシンを最適化マシンOPTと呼ぶ。例として以下の規則を用いる。

$$P_{OPT} = \{ NN \rightarrow \epsilon, MN \rightarrow NM, DN \rightarrow ND, NA \rightarrow S, NS \rightarrow A, LIT \rightarrow V \}$$

また、これらの規則より作成した解析表を図1に示す。

### 2.3 マシンOPTの解析表の構成法

マシンOPTの解析表の構成法を説明するために、次のようにOPT項を  $[a \rightarrow \beta \cdot \gamma, \delta]$  と定義する。

$a \rightarrow \beta \cdot \gamma$  はマシンOPTの解析が  $\beta$  まで終了し、 $\gamma$  を次に解析しようとしている状態を表す。また、 $\delta$  は  $W^*$  の要素である。

(定義 1) OPT項の集合  $J$  に対する閉包  $CLOSURE(J)$  は次のようにして得られる。

- 1)  $CLOSURE(J)$  を  $J$  に初期設定する。
  - 2)  $[a \rightarrow \beta \cdot \gamma, \delta] \in CLOSURE(J)$  に対して、 $a'$  が  $\gamma \delta$  の接頭辞である置換規則  $a' \rightarrow \beta'$  が  $P_{OPT}$  に含まれるなら  $[a' \rightarrow \beta', \delta']$  を  $CLOSURE(J)$  に加える。ただし、 $\delta'$  は  $\gamma \delta = a' \delta'$  を満足する。
  - 3)  $[a \rightarrow \beta \cdot \gamma, \delta] \in CLOSURE(J)$  に対して、 $\gamma \delta$  が  $a'$  の接頭辞である置換規則  $a' \rightarrow \beta'$  が  $P_{OPT}$  に含まれるなら  $[a' \rightarrow \beta', \epsilon]$  を  $CLOSURE(J)$  に加える。
  - 4) 上記2), 3) を新たに加えるものがなくなるまで繰り返す。
- 最適化マシンOPTにおいてバックトラックをなくすためには、一般的LRパーサのCLOSUREの定義に加えて、さらに定義1の3)の手順を追加しなくてはならない。これは、次のコード列の変化において、バックトラックを生じさせないためである。

$$\begin{aligned} \beta & \beta' \dots \dots (1) \\ \beta & a' \dots \dots (2) \\ \beta & \gamma \delta \dots \dots (3) \\ a & \delta \dots \dots (4) \end{aligned}$$

$\beta'$  が  $a'$  に置換され、 $a' = \gamma \delta$  であり、 $\beta \gamma$  が  $a$  に置換されることを表す。状態から状態への遷移を考えるために  $[a \rightarrow \beta \cdot A \gamma, \delta]$  を含んだ閉包から  $[a \rightarrow \beta \cdot A \cdot \gamma, \delta]$  を含んだ閉包への  $A$  による遷移を次のように定義する。

(定義 2)  $J$  をOPT項の集合、 $A \in W$  としたとき、次のGOTO関数を定義する。  
 $GOTO(J, A) = CLOSURE(\{ [a \rightarrow \beta \cdot A \gamma, \delta] \mid [a \rightarrow \beta \cdot \gamma, \delta] \in J \})$   
置換規則  $P_{OPT}$  に対するすべての状態集合  $K_{OPT}$  は次の手順により求まる。

- 1)  $K_{OPT}$  を  $\{ CLOSURE(\{ [a_1 \rightarrow \beta_1, \epsilon], [a_2 \rightarrow \beta_2, \epsilon], \dots, [a_v \rightarrow \beta_v, \epsilon] \}) \}$  と初期設定する。
- 2) 以下を新たな集合が  $K_{OPT}$  に加えられなくなるまで繰り返す。  
任意の  $J \in K_{OPT}$  と  $A \in W$  に対し  $GOTO(J, A)$  が空集合でなく  $K_{OPT}$  に入っていないければ、それを  $K_{OPT}$  に加える。  
定義2の1) で初期化される状態は初期状態と呼ばれ  $J_0$  で表す。置換規則  $P_{OPT}$  には定義2の1) のようにすべての規則に対する状態を加える必要がある。

置換規則  $P_{OPT}$  に対して状態集合  $K_{OPT}$  が求めれば、マシンOPTの動作を決定できるが、 $K_{OPT}$  においてもLRパーサと同様にR・R矛盾をもつ状態が存在する。LRパーサではこの矛盾を先読みにより解決しているが、マシンOPTでは次の関数EFF(efficientの略)を定義し、この矛盾の解決法を与える。

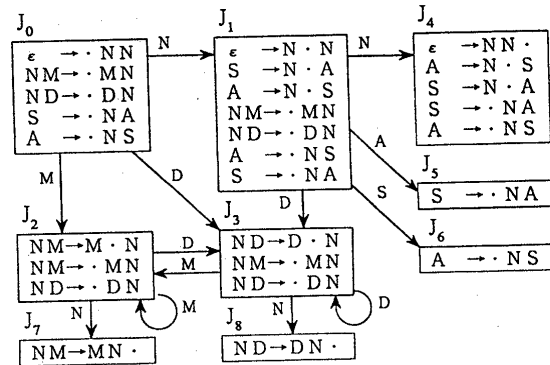


図1 マシンOPTの解析表

A String Replacement Algorithm by Using LR Parser.

Kazuhiko TSUDA<sup>†</sup>, Masami NAKAMURA<sup>†</sup>, Jun-ichi AOE<sup>††</sup>

<sup>†</sup> Sumitomo Metal Industries, LTD. <sup>††</sup> University of Tokushima

$a \rightarrow \beta \in P_{OPT}$ ならば、 $i\beta$ は $\beta$ に含まれる $W$ の要素数である。OPT項の第2項を省略して説明する。一般的にR・R矛盾をもつ状態 $J$ は、 $[\gamma \rightarrow \beta \cdot]$ なる1つのOPT項とこの $\beta$ を右辺に接尾辞をもつ規則 $\gamma_i \rightarrow \gamma'_i \beta$ ,  $1 \leq i \leq n$ によるOPT項 $[\gamma_i \rightarrow \gamma'_i \beta \cdot]$ から成り立っていると考えられる。またR・S矛盾をもつ状態 $J$ は、シフト動作に関係した $[a_j \rightarrow a'_j \cdot \delta_j]$ ,  $1 \leq j \leq m$ と還元動作に関係した $[\gamma_i \rightarrow \gamma'_i \beta \cdot]$ ,  $1 \leq i \leq n$ から成ると考えらる。従って、以後の説明では図2の(a)と(b)のR・R矛盾とR・S矛盾の一般形を用いる。

L Rパーサでは、これらの矛盾に対して先読みにより一意的に正しい動作を決定されなければならないが最適化マシンにおいてはより効率的なコードを得るようにこれらの矛盾に対して動作を決定しなければならない。すなわち、この関数EFFは局所的にコードが最短になる動作を決定するために使用される。

以下に $K_{OPT}$ からマシンOPTの動作を与えるための定義を与える。関数ACTIONと同様に使用されるが、動作の矛盾をもつ状態 $J$ に対してのみ、先読みコード列の要素の集合LOOK(J)を定義し、LOOK(J)に先読みが含まれる場合の動作をFIRST(J)、そうでない場合の動作をSECOND(J)で表す。

(定義 3)

- 1) 状態 $J$ が動作矛盾をもたない場合
  - a)  $[a \rightarrow \beta \cdot \gamma, \delta] \in J$ かつ $\gamma \neq \epsilon$ ならば  
ACTION(J)=shiftを定義する。
  - b)  $[a \rightarrow \beta \cdot, \delta]$ かつ $a \rightarrow \beta$ が第 $p$ 番目の規則ならば、  
ACTION(J)=reduce  $p$ を定義する。
- 2) 状態 $J$ がR・R矛盾をもつ場合 (図2の(a)参照)
 

$\gamma_1 \rightarrow \gamma'_1 \beta, \gamma_2 \rightarrow \gamma'_2 \beta, \dots, \gamma_n \rightarrow \gamma'_n \beta$ の中で関数EFFの値が最大のものを $\gamma_i \rightarrow \gamma'_i \beta$ ,  $1 \leq i \leq n$ と仮定し、 $\gamma \rightarrow \beta$ を第 $p$ 番目、 $\gamma_i \rightarrow \gamma'_i \beta$ を第 $q$ 番目の規則と仮定する。

  - a)  $EFF(\delta_1 \rightarrow \delta_2 \gamma \delta) > EFF(\gamma_i \rightarrow \gamma'_i \beta)$ の場合  
ACTION(J)=conflict, LOOK(J)={ $\delta$ },  
FIRST(J)=reduce  $p$ ,  
SECOND(J)=reduce  $q$ と定義する。
  - b)  $EFF(\delta_1 \rightarrow \delta_2 \gamma \delta) \leq EFF(\gamma_i \rightarrow \gamma'_i \beta)$ の場合  
ACTION(J)=reduce  $q$ と定義する。
- 3) 状態 $J$ がR・S矛盾をもつ場合 (図2の(b)参照)
 

$\gamma_1 \rightarrow \delta_1, \gamma_2 \rightarrow \delta_2, \dots, \gamma_n \rightarrow \delta_n$ の中で関数EFFの値が最大のものを $\gamma_i \rightarrow \delta_i$ ,  $1 \leq i \leq n$ と仮定し、この規則番号を第 $p$ 番目とする。また、 $a_1 \rightarrow \beta_1 \delta_1, a_2 \rightarrow \beta_2 \delta_2, \dots, a_m \rightarrow \beta_m \delta_m$ は関数EFFの値の大きい順にソーティングされているものと仮定する。

  - a)  $EFF(a_i \rightarrow \beta_i \delta_i) > EFF(\gamma_i \rightarrow \delta_i)$   
 $\geq EFF(a_{i+1} \rightarrow \beta_{i+1} \delta_{i+1})$ ,  $1 \leq i \leq m$ の場合  
ACTION(J)=conflict, LOOK(J)={ $\delta_1, \delta_2, \dots, \delta_c$ },  
FIRST(J)=shift, SECOND(J)=reduce  $p$ と定義する。
  - b)  $EFF(\gamma_i \rightarrow \delta_i) \geq EFF(a_i \rightarrow \beta_i \delta_i)$ の場合  
ACTION(J)=reduce  $p$ と定義する。

図2(a)の状態 $J$ に対して、 $n+1$ 種類の還元動作が可能であるが、本論文で新しく提案する矛盾解決法は、局所的に短いコードを生成する還元動作を選び出すことである。最も局所的な矛盾解決法は、 $\gamma_1 \rightarrow \gamma'_1 \beta, \gamma_2 \rightarrow \gamma'_2 \beta, \dots, \gamma_n \rightarrow \gamma'_n \beta$ の中で関数EFFの値が最大の規則 $(\gamma_i \rightarrow \gamma'_i \beta)$ による還元動作を選び出すことである。すなわち、この選択はR・R矛盾に関係した規則の情報だけに依存している。

しかし、定義6の2)では、規則 $\gamma \rightarrow \beta$ に対する局所的な先読み情報を利用して矛盾を解決している。状態 $J$ のOPT項 $[\gamma \rightarrow \beta \cdot]$ は、状態 $J_1$ のOPT項 $[\delta_1 \rightarrow \delta_2 \cdot \gamma \delta]$ の閉包から得られたものであるので $\gamma \rightarrow \beta$ による還元動作は、先読みコード列 $\delta$ が存在すれば、 $\delta_1 \rightarrow \delta_2 \gamma \delta$ による還元動作が起こることが保証される。従って $EFF(\delta_1 \rightarrow \delta_2 \gamma \delta) > EFF(\gamma_i \rightarrow \gamma'_i \beta)$ であって、状態 $J$ の先読みコード列が $\delta$ であれば、 $\gamma \rightarrow \beta$ による還元動作を選択すべきである。また、 $EFF(\delta_1 \rightarrow \delta_2 \gamma \delta) < EFF(\gamma_i \rightarrow \gamma'_i \beta)$ ならば $\gamma_i \rightarrow \gamma'_i \beta$ による還元動作が選択されるべきである。

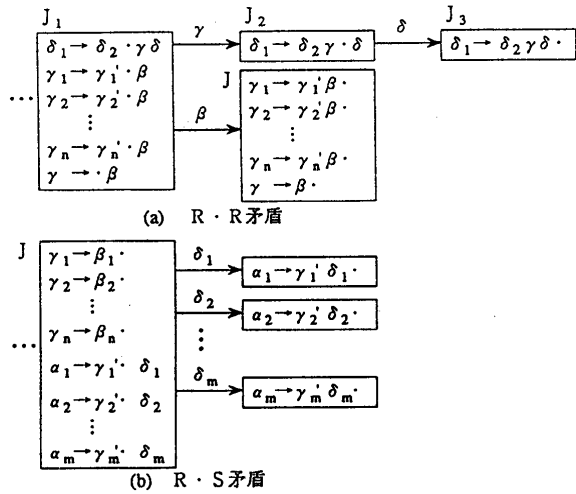


図2 R・R矛盾とR・S矛盾

図2(b)の状態 $J$ において、コード最適化では、 $[a_j \rightarrow a'_j \cdot \delta_j]$ ,  $1 \leq j \leq m$ なる項に対するシフト動作は、後に $a_j \rightarrow a'_j \cdot \delta_j$ による還元動作が起こることが保証されない。従って、状態 $J$ に対する最も局所的な矛盾解決法は、 $\gamma_1 \rightarrow \beta_1, \gamma_2 \rightarrow \beta_2, \dots, \gamma_n \rightarrow \beta_n$ の中で関数EFFの値が最も大きい規則 $(\gamma_i \rightarrow \beta_i)$ による還元動作を選び出すことである。

しかし、定義3の3)では、シフト動作に対する局所的な先読み情報を利用して矛盾を解決している。EFF $(\gamma_i \rightarrow \beta_i) \geq EFF(a_i \rightarrow \beta_i \delta_i)$ ならば、先読みなしに規則 $\gamma_i \rightarrow \beta_i$ による還元動作が選択されるべきである。これに対しEFF $(a_i \rightarrow \beta_i \delta_i) > EFF(\gamma_i \rightarrow \beta_i) \geq EFF(a_{i+1} \rightarrow \beta_{i+1} \delta_{i+1})$ ,  $1 \leq i \leq m$ であって $\delta_k$ ,  $1 \leq k \leq i$ なる先読みコード列が存在すれば、 $a_k \rightarrow \beta_k \delta_k$ による還元動作が保証されるので、シフト動作が選択されるべきである。ただし、この $\delta_k$ が存在しない場合は、 $\gamma_i \rightarrow \beta_i$ による還元動作が選択される。これを実行するマシンOPTの解析は、2つのスタックを用いて行われるが、詳細はD.A.Waltersの文献(4)を参照されたい。

3. まとめ

本稿では、バックトラックが不要かつより良い規則の選択が可能となる文字列置換アルゴリズムを提案した。本アルゴリズムは、局所的コード最適化にも応用可能である。今後は、本アルゴリズムを文書短縮に適用し評価を早急に進めて行きたい。また、局所的コード最適化にも応用し、A.S.Tanenbaum(3)による最適化マシンとの比較・評価も行う予定である。

参考文献

- (1) 津田和彦, 中村雅巳, 青江順一: “形態素置換による文書短縮法”, 信学論Vol.J75-D-II, No.3, pp.619-627, (1992-03)
- (2) 津田和彦, 中村雅巳, 青江順一: “文書短縮のための文字列置換アルゴリズム”, 情処学NL究報, 85-5, (1991-09)
- (3) A.S.Tanenbaum, H.Staveren and J.W.Stevenson: “Using peephole optimization on intermediate code”, ACM Trans. Prog. Lang. Syst., 4, 1, pp.21-36, (1982)
- (4) D.A.Walters: “Deterministic context-sensitive language: part 1 part 2” Inf. Control, 17, pp.14-40, pp.41-61, (1970)
- (5) J.W.Davidson and C.WcFraser: “The design and application of a retargetable peephole optimizer” ACM Trans. Prog. Lang. Syst., 2, 2, pp.191-202, (1980)