

言語運用のモデルに基づく実時間の構文解析

橋田浩一, 山田尚勇(東京大学理学部)

1. はじめに

ここで紹介する研究の目標は、人間の言語処理の過程を明らかにすることと、その挙動をまねるようなシステムを計算機上に実現することにある。人間の言語処理の外的な特徴のうち、われわれが注目したのは、それが普通は実時間で行なわれる、という点である。

自然言語を構文解析する手法として、今までさまざまなもののが提案されてきたが、それらの大部分は入力文の長さに関して線型より大きな時間計算量をもっている。たとえば、ATN[1]やDCG[2]などの仮説駆動(hypothesis driven)型のやりかたは、一般に入力文の長さの指數関数にあたる時間計算量をもち、LINGOL[3]などで用いられているデータ駆動(data driven)型のアルゴリズムは三次多項式の時間計算量をもっている。すなわち、これらの手法は、時間計算量に関して人間の振舞いに十分近いものとは言えない。

従来の方法で実時間の処理が行なえないのは、それらが言語能力(linguistic competence)の反映である文法の記述に偏り、言語運用(linguistic performance)に関する十分な洞察に基づいていないためであるとわれわれは考える。特に欠けていると思われるるのは、人間の短期記憶の有限性[4,5]に関する配慮である。

われわれはまず、言語運用に見られるある側面を短期記憶の有限性と関係づける代数的なモデルを作成した。これは発話のモデルで、人間が発することのできる全ての文を生成する能力を持つ有限オートマトンである。これを変型することによって実時間の構文解析システムを実現できるわけである。ただし、残念ながらこのシステムは人間が言語を認識する過程のモデルとしては十分でない。

ここでの議論は、対象とする自然言語の文法が文脈自由な形に書けることを前提しているが、最近の拡張文脈自由文法(generalized phrase structure grammar)[6,7]についての議論にも見るよう、これは文法が変型文法の枠内で記述できるという前提と同程度の一般性を持っているとわれわれは考える。

さらに、われわれの手法は、ほぼ純粹に言語運用論的なものであり、文法に対しては、それが文脈自由な形に書けること以外にいかなる条件をも課さない。したがつ

て、文法の記述は文脈自由である限りここでの議論とは独立に論じられる。特に、構文解析の際に複数個の可能な構文木のうちどれを選ぶかというような問題は、たとえば確率的文脈自由文法[8]の枠組の中で処理するべきもので、われわれの手法が関知するところではないと考えられる。また、われわれの議論は全く構文論的なもので、意味的な処理はここで述べる手法とは独立に扱われる。

なお、以下の議論のうち定理の証明は省略した。詳細については橋田[9]を参照して頂きたい。

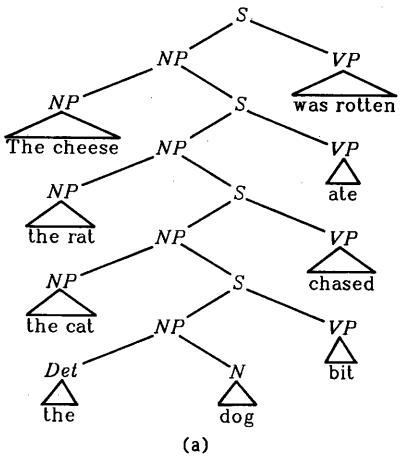
2. 準備

ここでは以後の議論に必要な用語を定義する。

まず、対象とする自然言語を L と決める。 L を記述する文脈自由文法を G とする。ただし、 G が生成する言語 $L(G)$ が L に一致するわけではなく、 $L(G)$ から、実際の運用には用いられない表現を除いたものが L である。

$G = \langle VN, VT, P, S \rangle$ としよう。 VN は非終端記号の集合、 VT は終端記号の集合、 P は生成規則の集合、 $S \in VN$ は開始記号である。 G はチョムスキーモルヒニク形であるとすると、つまり、生成規則は $A \rightarrow BC$ または $A \rightarrow a$ という形に限られる。ここで大文字は非終端記号、小文字は終端記号とする。 $A \in VN$, $a \in VT$, $\alpha \in VT^* \cup VN^*$ に対し、 $(A \rightarrow \alpha) \in P$ を単に $A \rightarrow \alpha$ と書く。

句とは、構文木ないしはその葉の並びの終端記号列のことである。句 Q が句 R の中に入れ子になっているとは、終端記号列 α, β が存在し、 $R = \alpha Q \beta$ かつ α, β が空列でないことである。このときさらに、 α も β も空列でなければ、 Q は R の中に真の入れ子になっていると言う。句 Q の（真の）入れ子の深さとは、 $0 < i < d$ について $P_{i+1} \in P_i$ の中に（真の）入れ子になっている $P_0 = Q, P_1, \dots, P_d$ が存在するような d の最大値のことである。構文木の中の路(path) p において、節点 A が曲がり角であるとは、 A と A の親 B が p 上にあり、 A は B の左(右)の子であって、 A の右(左)の子も p 上にあることを言う。 G はチョムスキーモルヒニク形であるから、句 Q が句 R の中に真の入れ子になっている条件は、 R の根から Q の根までの路が曲がり角を持つことである。



(a)

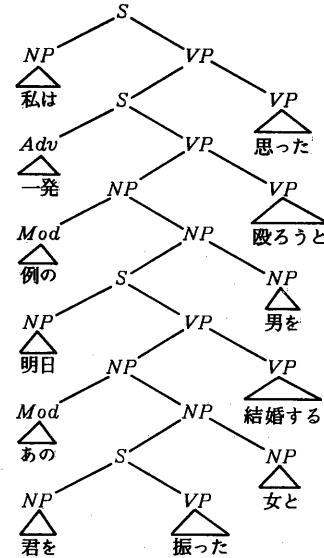


図1

(b)

3. 真の入れ子に関する制約

われわれは理想化された言語能力の反映として文脈自由文法を考えているわけだが、人間の情報処理能力が制限されているために、そのような文法によって生成される「正しい」文でありながら、実際の言語運用には現われないものがある。たとえば、次のような文がそれである。

(a) "The cheese the rat the cat the dog bit
chased ate was rotten."

(b) 「私は一発例の明日あの君を振った女と結
婚する男を殴ろうと思った。」

図1にこれらの文の構造を示す。これを見てすぐにわかるように、どちらの文もかなり大きな真の入れ子の深さを持っている。(a)における真の入れ子の深さは4、(b)においては5である。

一般に、真の入れ子の深さが増すほど、文の生成および認識が人間にとて難しくなる、ということは古くからよく知られており、実際には、真の入れ子の深さがある定数を越えるような表現は使われないと考えられる。たとえば上の例文のようなものは、通常の会話には現われない。

自然な言語運用においては真の入れ子の深さがある定数で抑えられる、という定立を、われわれは真の入れ子に関する制約と呼ぶことにする。以下では、この定立が

正しいと仮定した上で話を進めることにする。

真の入れ子によって引き起こされる処理の難しさは、いわゆる garden path 、つまり局所的な構文上の曖昧さによるそれとは異なる。実際、上の例文はgarden pathを含まない。さらに、garden path を含む文が発話されることがしばしばあるのに対し、真の入れ子が深い文は発話されることすらないように見える。

4. 短期記憶の有限性

人間の記憶は、いくつかの階層を持っている。それらの階層のうち、短期記憶 (short term memory) あるいは一時記憶 (temporary memory) と呼ばれるものは、パターン認識を経て意味づけのなされた情報が初めてに入るところであり、かつ言語の発話、認識に伴うような即時的な演算が行なわれる場であると考えられている[5]。

ところが、この短期記憶には一度に 7±2個 の情報の塊 (chunk) しか入らない。ただし1個の塊の情報量は可変で、たとえば二進数字の並びを3個づつの組に分け、それぞれの組を0~7の数字だと思って憶えると、普通にそのまま憶えたときの約3倍の長さの並びを正確に憶えることができる[4]。とは言え、塊は無限に大きくなるわけではなく、情報の種類と個人ごとに1個の塊が持てる情報量は決まっているようである。

5. 言語運用のモデル

ここでは、短期記憶の有限性を満たすことによって真の入れ子に関する制約を説明するような、言語運用の代

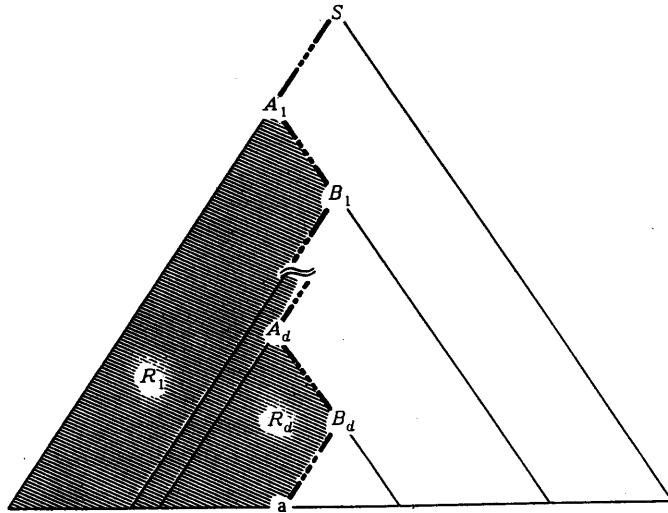


図2

数的なモデルについて述べる。われわれはこのモデルを発話に関するモデルとして提示するわけであるが、そのようなモデルは、人間が発話できるGに関して正しい文であれば生成できる、という条件を満足すべきであり、人間が発話できない文はなるべく生成できない方がよい。発話できるかできないかという規準としてわれわれがここで考えているのは、真の入れ子の深さである。

さて、人間が図2のような構文木を持つ文を発話しており、その葉aの直前までの語列を発声し終えたところであるとしよう。すると、図の斜線を引いた部分の構造はすでに決まっていると考えてよいだろう。これに対して白い部分の構造は話者にとっても未定であり得る。話者は斜線部分と白い部分の文法的な整合を図りつつ発話を続けるものと考えらるが、そのために短期記憶内にはいきなる構文的な情報を貯えておけばよいだろうか？

図2のA₁, B₁, …, A_d, B_dという節点は、斜線部分と白い部分を過不足なく連結している。したがって、これらを短期記憶に入れておけば文全体の文法的な正しさを管理しつつ文を生成できそうである。さらに、根Sから葉aに致る路において、B₁, A₂, …, B_{d-1}, A_dが曲がり角であることに注意して頂きたい（A₁は根Sに一致するかも知れず、またB_dはaの親かも知れないでの、これらは曲がり角とは限らない。）。したがって、短期記憶の有限性、つまりdの最大値の有限性は、真の入れ子に関する制約に符号すると予想できる。

以上の議論をもとに、発話のモデルΦ(G, D)を定

式化する。Φ(G, D)は有限オートマトンと見なすことができる。その内部状態をsという変数で表わすことにしてよう。

$$s \in \bigcup_{i=0}^{2D+1} VN^i$$

である。Φ(G, D)の動きは次の手続きによって定められる。

- (1) s := < > (空リスト) とし、(2)へ行く。
- (2) s = < A₁, B₁, …, A_d, B_d > であるとする。下の(2a), (2b) のうち実行可能なものをひとつだけ実行して(3)へ行く。
 - (2a) A → a なる A ∈ VN, a ∈ VT が存在するとき、
s := < A₁, B₁, …, A_d, B_d, A >
とし、a を出力する。
 - (2b) d ≥ 1 であり、B_d → a なる a ∈ VT が存在するとき、
s := < A₁, B₁, …, A_d > とし、a を出力する。
- (3) s = < A₁, B₁, …, A_d, B_d, C > であるとする。下の(3a), (3b), (3c) のうち実行可能なものをひとつだけ実行する。
 - (3a) d < D であり、A → C B なる A, B ∈ VN が存在するとき、
s := < A₁, B₁, …, A_d, B_d, A, B >
として(2)へ行く。

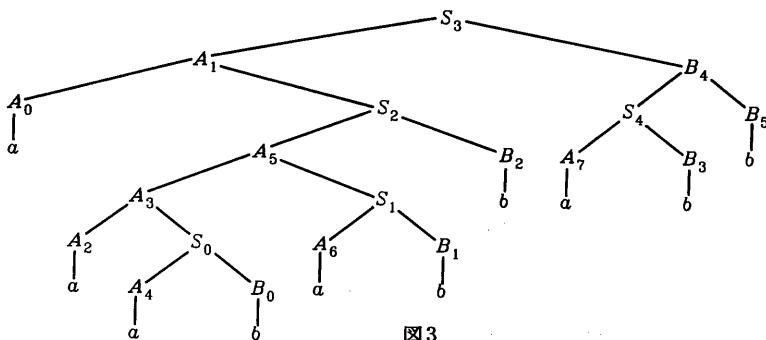


図3

- (3b) $d \geq 1$ であり、 $Bd \rightarrow C$ B なる $B \in VN$ が存在するとき、
 $s := < A_1, B_1, \dots, A_d, B >$
 として(2)へ行く。

(3c) $s = < S >$ なら終了する。

図3に示される構文木をもつ文を $\Phi(G, D)$ が生成する様子を下に示す。

s の 値	出力	規則
$< >$		
$< A_0 >$	a	(2a)
$< A_1, S_2 >$		(3a)
$< A_1, S_2, A_2 >$	a	(2a)
$< A_1, S_2, A_3, S_0 >$		(3a)
$< A_1, S_2, A_3, S_0, A_4 >$	a	(2a)
$< A_1, S_2, A_3, B_0 >$		(3b)
$< A_1, S_2, A_3 >$	b	(2b)
$< A_1, S_2, A_5, S_1 >$		(3a)
$< A_1, S_2, A_5, S_1, A_6 >$	a	(2a)
$< A_1, S_2, A_5, B_1 >$		(3b)
$< A_1, S_2, A_5 >$	b	(2b)
$< A_1, B_2 >$		(3b)
$< A_1 >$	b	(2b)
$< S_3, B_4 >$		(3a)
$< S_3, B_4, A_7 >$	a	(2a)
$< S_3, B_4, S_4, B_3 >$		(3a)
$< S_3, B_4, S_4 >$	b	(2a)
$< S_3, B_5 >$		(3b)
$< S_3 >$	b	(2b)
		(3c)

実際の運用に現われる真の入れ子の深さの最大値に1を加えた値が D だと考えれば、 $\Phi(G, D)$ は人間の発話できる正しい文を全て生成することができる、ということが、上の定理から言える。

しかしながら、話はむしろ逆で、 $\Phi(G, D)$ として抽象化されるような言語処理過程の結果として、真の入れ子に関する制約があるのだ、と考えるべきである。われわれは $\Phi(G, D)$ の s を人間の短期記憶ないしはその一部の抽象化と考えているので、上の手続きの各ステップ中の d を、その瞬間の記憶負荷と呼ぼう。ある文の記憶負荷と言った場合には、その文を $\Phi(G, D)$ が生成するときの d の最大値を指すものとする。もちろんその場合、その文のひとつの構文木が想定されているものとする。 $\Phi(G, D)$ を発話のモデルと考えるからには、短期記憶の有限性によって制約を受けるのは真の入れ子の深さではなく、記憶負荷であると考えなければならない。

このことは次のふたつの例文が何とか発話できるものであることによって例証される。

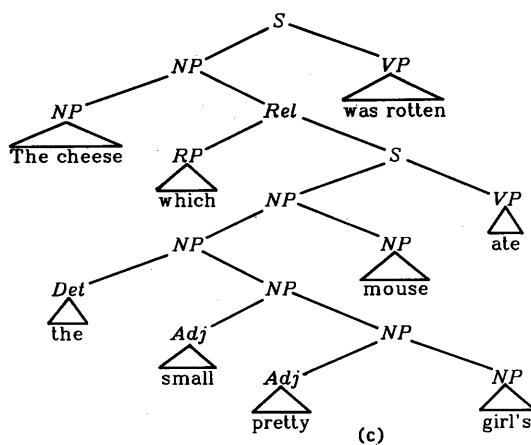
(c) "The cheese which the small pretty girl's mouse ate was rotten."

(d) 「私は是非昨日君が私が君にしぶしぶ貸した本を捨ててしまったわけを知りたい。」

図4にこれらの構文木を示す。

(c), (d)の真の入れ子の深さはそれぞれ4、5であり、(a), (b)と一致している。ところが、(a), (b)の記憶負荷がそれぞれ5、6であるのに対し、(c), (d)ではそれぞれ2、3である。つまり、発話の難しさを引き起こしているのは、真の入れ子の深さというよりはむしろ記憶

定理1. $\Phi(G, D)$ が生成する文字列は $L(G)$ の文であり、真の入れ子の深さが $2D$ を越えない構文木を持つ。また、真の入れ子の深さが $D-1$ を越えない構文木を持つ文は $\Phi(G, D)$ によって生成される。



負荷なのである。

そういう文脈の中で $\Phi(G, D)$ を特徴づけるのが次の定理である。実際には、定理1はこれから導かれる。

定理2. $\Phi(G, D)$ がある文字列を生成できる条件は、
その文字列が G で生成される文であり、かつその記憶負荷がある構文木に関して D を越えないことである。

6. 実時間の構文解析

この節では、上の $\Phi(G, D)$ をもとにして、自然言語の構文解析を実時間で行うアルゴリズムを設計する。

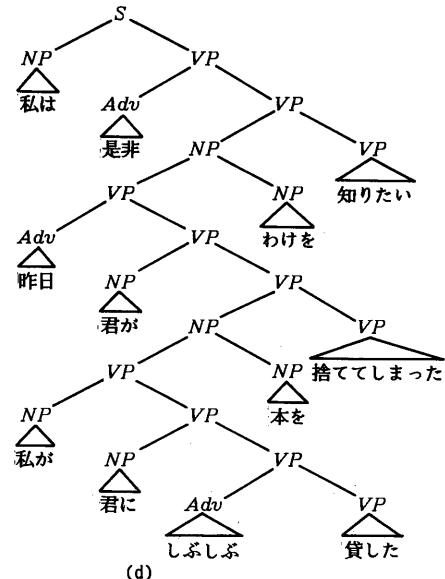
前に示した例文を見ると、認識の難しさについても真の入れ子の深さよりはむしろ記憶負荷が大きな原因になっているように思われる。そこでわれわれは、ここで設計するアルゴリズムが満足すべき要件として、次のふたつを考える。

- (ア) そのアルゴリズムがある文字列 α を受理するならば、 α は $L(G)$ の文であり、このとき記憶負荷がDを越えないような α の構文木がつくれられる。

(イ) $L(G)$ の文 α が記憶負荷がDを越えない構文木を持てば、そのアルゴリズムは α を受理する

$\Phi(G, D)$ の出入力関係を入れかえ、さらに構文木をつくる手続きを加えれば、定理2によって、上の(ア)、(イ)を満足するアルゴリズムができるだろう。

以下で述べるアルゴリズムにおいて構木を作る際に、未完成の木の部分構造は、図2のR₁, ..., R_d のような構造のリストで表わすことにする。一般に、構木の部分木であって、右端の葉が非終端記号で、他には欠けたところのないものを左木と呼ぶ。左木の右端の葉をそ



の左木の尾と呼ぶ。構文木の集合を TREE, 左木の集合を LTREE と書くことにする。図2の R_i は A_i が根であり B_i が尾であるような左木である。また、

$$LTLIST = \bigcup_{d=0}^{\infty} LTREE^d, TLIST = LTLIST \times TREE$$
とする。

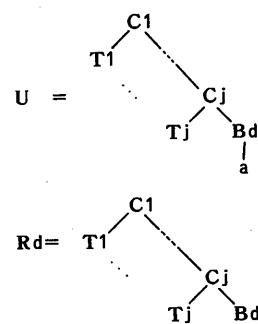
アルゴリズム1-

2つの変数 $|T| \in LTLIST$, $T \in TLIST$ を用いる。

- (1) LTL : = {< >} として(2)へ。
 (2) まず TL : = ϕ とし、 $M \in LTL$ なる全ての M に対して、
 次の(2a), (2b)を実行して(3)へ。ただし、
 $M = < R1, \dots, Rd >$ であるとし、次の入力は
 $a \in VT$ であるとする。

- (2a) $A \rightarrow a$ なる全ての $A \in VN$ に対し、
 $TL := TL \cup \{ < R1, \dots, Rd, \frac{A}{a} > \}$ とする。

(2b) $d \geq 1$ かつ $Bd \rightarrow a$ ならば
 $TL := TL \cup \{ N \}$ とする。ただし、
 $N = < R1, \dots, Rd-1, U >$

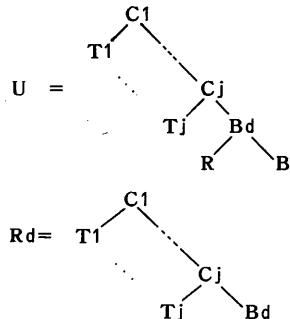


である。

(3) 入力が残っているとき：

まず $LTL := \phi$ とし、 $M \in TL$ なる全ての M に
対して、次の(3a), (3b)を実行して(2)へ。ただし、 $M = \langle R_1, \dots, R_d, R \rangle$ とし、木 R
の根を $C \in VN$ とする。

- (3a) $d < D$ ならば、 $A \rightarrow C B$ なる全ての
 $A, B \in VN$ に対して、
 $LTL := LTL \cup \{ \langle R_1, \dots, R_d, \frac{A}{R} B \rangle \}$ と
する。
- (3b) $d \geq 1$ ならば、 $B_d \rightarrow C B$ なる全ての $B \in VN$
に対し、 $LTL := LTL \cup \{ N \}$ とする。ただし、
 $N = \langle R_1, \dots, R_{d-1}, U \rangle$



である。

入力が残っていないとき：

$\langle T \rangle \in TL$ かつ根が S であるような T を全て
出力して終わる。

定理3. アルゴリズム1は(ア)、(イ)を満足する。

ところが、上のアルゴリズムにおいて、 LTL と TL の
濃度は入力の長さに関して指數関数的に増大するので、
その時間計算量は指數関数になる。そこで、(ア)、(イ)
に違反しない範囲で LTL の要素を削除することを考える。

左木 R の根を $RM(R)$ 、尾を $TM(R)$ と書き、

$\sigma(\langle R_1, \dots, R_d \rangle)$

$= \langle RM(R_1), TM(R_1), \dots, RM(R_d), TM(R_d) \rangle$
としよう。

アルゴリズム2。

このアルゴリズムは上のアルゴリズム1の(3)の後
に(4)を実行するようにしたものである。

(4) $M, N \in LTL, M \neq N$ なる全ての M, N に対し、

$\sigma(M) = \sigma(N)$ ならば、 $LTL := LTL - \{M\}$ として
(2)へ行く。

定理4. アルゴリズム2は(ア)、(イ)を満たし、入力
の長さに関して線型の時間計算量を持つ。

アルゴリズム2は(4)において、 $\Phi(G, D)$ の同じ状
態に対応する LTL の異なる要素のひとつを捨てるこ
とにによって、 LTL の濃度を $\sum_{d=0}^D |VN|^{2^d}$ で抑えているわけであ
る。しかし、 $|VN|$ は少くとも10程度、 D は3ぐらい
だろうから、この上限はかなり大きなものになってしま
う。実時間の処理を実現するためには、アルゴリズム
が入力を順方向に処理しその時間計算量が入力の長さに
関して線型であれば数学的には十分であるが、現実問題
としては、時間計算量の比例定数が十分小さくなければ
ならない。

そこで、上のアルゴリズムの(4)を改良して LTL の濃度
をさらに小さく保つことを考える。

$\Psi = (Q, \Sigma, q_I, F, \delta)$ を非決定性有限オートマト
ンとしよう。 Q は状態の集合、 Σ はアルファベットの集
合、 $q_I \in Q$ は初期状態、 $F \subseteq Q$ は最終状態の集合、
 $\delta : Q \times \Sigma^* \rightarrow 2^Q$ は遷移関数である。このような Ψ に対
し、

関数 $SUF_\Psi : Q \rightarrow 2^{\Sigma^*}$ を

$$SUF_\Psi(q) = \{ \alpha \in \Sigma^* \mid \delta(q, \alpha) \cap F \neq \emptyset \}$$

で定めよう。さらに、 Q における半順序 Ψ を

$$q \stackrel{\Psi}{\ll} r \iff SUF_\Psi(q) \subseteq SUF_\Psi(r)$$

によって定める。

定理5. $q \stackrel{\Psi}{\ll} r$ のとき、しかもそのときに限り、
 $|\alpha| < 2^{|q|}$ なる全ての $\alpha \in \Sigma^*$ に対して、

$\delta(q, \alpha) \cap F \neq \emptyset$ ならば $\delta(r, \alpha) \cap F \neq \emptyset$ である。

この定理によって、 \triangleleft の構造を陽に求めるアルゴリズ
ムが存在することがわかる。そこで、 \triangleleft を単に \triangleleft と
書くことにすると、次のアルゴリズム3は確かにアルゴ
リズムになっている。

アルゴリズム3。

このアルゴリズムは(4)を除いてアルゴリズム2と同
じである。

(4) $M, N \in LTL, M \neq N$ なる全ての M, N に対し、

$\sigma(M) \triangleleft \sigma(N)$ ならば、 $LTL := LTL - \{M\}$ として
(2)へ。

定理6. アルゴリズム3も(ア)、(イ)を満たし、入力
の長さに関して線型の時間計算量を持つ。

\triangleleft の構造はアルゴリズム3の実行に先立って求めて

おけばよいか、アルゴリズム3の効率は \triangleleft の構造を求める計算の複雑さには影響されない。

ところが、定理5を見ると、 \triangleleft の構造を求めるには少くとも $O(2^{|Q|})$ 程度の手間がかかりそうである。この場合は $|Q| > |VN|^2$ であるから、 $|VN|=10, D=3$ としても、この計算はとても手に負えない。

そういうわけで、 \triangleleft を正確に計算することは諦めざるを得ない。そのかわり、アルゴリズム3の(4)において、何らかの方法で $\sigma(M) \triangleleft \sigma(N)$ がわかったものについてだけ $LTL := LTL - \{M\}$ とすることにしよう。

一般に、次の定理に示されるように、 \triangleleft を隣接する非終端記号の組ごとの関係に分解することができる。

定理7. 次の条件(A), (B), (C)が成り立つとき、

$\langle A_1, B_1, \dots, A_d, B_d \rangle \triangleleft$

$\langle C_1, D_1, \dots, C_d, D_d \rangle$ である。

ただし、wider-than(B, A, D, C) という述語は、

$E_0 = B, E_n = A$ で、 $0 \leq j < n$ に対して

$E_j \rightarrow E_{j+1} F_j$ であるような $n \geq 0$ と

$E_0, \dots, E_n, F_0, \dots, F_{n-1} \in VN$ が存在するならば、

$G_0 = D, G_n = C$ で、 $0 \leq j < n$ に対して

$G_j = G_{j+1} F_j$ であるような $G, \dots, G_n \in VN$ が存在する。

という意味であるとする。

(A) wider-than(S, A_1, S, C_1)。

(B) $1 \leq i \leq d$ に対して、

wider-than($B_i, A_{i+1}, D_i, C_{i+1}$)。

(C) 任意の $a \in VT, A \rightarrow a$ なる任意の $A \in VN$ に対

し、wider-than(B_d, A, D_d, C)かつ $C \rightarrow a$ なる $C \in VN$ が存在する。

wider-than(B, A, D, C)の値を全ての B, A, D, C の組み合わせに対してあらかじめ求めておいて、上の定理を使って $\sigma(M) \triangleleft \sigma(N)$ がわかるかどうかを調べるというやりかたは、 \triangleleft を正確に求めるよりは、はるかに現実的である。しかし、実際の自然言語に見られる構文的な特徴を考慮することによって、さらに効率のよい方法が得られる。

たとえば、文法 G が次の性質を満たすと仮定しよう。

(LR) 全ての A, B, C, E, F に対して、もし $A \rightarrow F B$ ならば、 $C \rightarrow A E$ のとき、しかもそのときに限り

$C \rightarrow B E$ である。

これを直観的な言い方に直すと：

句全体の文法的な性質は、その句の右端の部分によって決定される。

日本語、朝鮮語、トルコ語など、修飾語が被修飾語に先立つような諸言語は、およそその性質を満たしている。こうした言語は世界中の言語の半数近くを占めている[10]ので、(LR)を前提としたやり方は広い適用可能性をもっていると考えられる。

定理8. G が(LR)を満たすとき、次の(A), (B) が成り立つ。

$\langle A_1, B_1, \dots, A_j, B_j \rangle \triangleleft$

$\langle C_1, D_1, \dots, C_k, D_k \rangle$ である。

ただし、 $A, B \in VN$ に対して、

$\{\alpha \in VT \cup VN^2 \mid A \rightarrow \alpha\} \subseteq \{\alpha \in VT \cup VN^2 \mid B \rightarrow \alpha\}$

を $A \triangleleft B$ と書く。

(A) $A_1 = S$ ならば $C_1 = S$ である。

(B) 自然数から自然数への関数 f で、

$1 \triangleleft f(1) \triangleleft f(2) \triangleleft \dots \triangleleft f(k) \triangleleft j$ かつ

$1 \leq i \leq k$ に対して $Bf(i) \leq Df(i)$ となるものが存在する。

日本語や朝鮮語とは逆に、修飾語が被修飾語より後に置かれる言語に関する限り、左右を逆転した議論によって、やはり上の定理と同程度に便利な結果を得ることができる。修飾の向きに左右両方がある英語などの言語については、これらの結果を混用することによって対処できる。

7. 実験

前節で述べたアルゴリズムが本当に実時間で動くことは、実験による検証を必要とする。そこでわれわれは、日本語を対象に、このアルゴリズムに基づいた構文解析システムを作成し、簡単な実験を行なった。日本語は最も馴染みの深い言語であるのみならず、性質(LR)をほぼ完全に満足するという、比較的すっきりした構造をもっているので、最初の実験の対象として適切であると考えたわけである。

われわれが用いた文法は各文節を終端記号とする文脈自由文法である。この文法は非常に簡単なもので、陳述副詞と述語の対応、並列句の扱いなどの細かい記述は含んでいないが、助詞と用言の係り受けなどの基本的なものはほぼ盛り込んである。

実際に時間がかかるのは、構文木を作る手続きよりも文節を認識する手続きであるが、それを含めてもわれわ

れの実験システムは十分速いことがわかった。解析された文は文節の数が15個程度までのものであるが、LTL の濃度は文節の数が5,6 個を過ぎたあたりから安定しており、この傾向はいかに長大な文であろうとも変わらないと考えられる。したがって、このシステムは実時間の構文解析をしていると見なすことができる。

8. 認識のモデルに向けて

(LR)が満たされているとき、定理8からも想像されるように、図2の斜線部分と白い部分の整合性を管理するために、R₁, ..., R_dの尾だけを短期記憶に入れておけばよい。つまり、統語構造が左から右への修飾関係のみによって構成されている言語では、構文木を生成しているときに短期記憶に入っている非終端記号の数は、記憶負荷に等しい。

統語構造が右から左への修飾関係だけで構成されるような言語に関しても、同様の議論が成立する。すなわち、記憶負荷と同じ個数の非終端記号で生成の際の構文的な整合性が管理できる。

したがって、人間の短期記憶における1個の情報の塊がちょうど1個の非終端記号に対応すると仮定すれば、少くとも左から右への修飾だけで文ができる言語、および右から左への修飾だけで文ができる言語に関しては、統語的な制御に用いられる情報の塊の数は記憶負荷程度だということになる。

記憶負荷の値は文法によって変わるが、たとえば日本語の場合、ごく普通に作られた文法に関して、生成に際しても認識に際しても、記憶負荷の最大値は高々4程度と思われる。したがって、短期記憶中に3±2個の情報の塊を別に確保することができ、これらは意味的な処理や構文的な曖昧さを解決するための処理に使われているのであろう。

ところが、6節で述べたわれわれのアルゴリズムは、以上の議論にあてはまらない。つまり、LTL の元の中に含まれる異なる左木の個数が短期記憶中の情報の塊の数に一致すべきであるにもかかわらず、その数が陽に7±2で抑えられていかない。実験においても、その数は最大10以上に達していることがわかった。

これはわれわれのアルゴリズムが条件(ア)、(イ)を満足するために構文的な曖昧性の解決を完全に行なっているからである。ところが人間が文を認識する際には、少くとも条件(イ)は満たされていない。Marcus[11]はここに着目し、garden path を排除することによって、人間

が普通に行なう認識手順を決定性の手続きで模倣することを試みている。われわれの方法は garden path を排除しておらず、Marcusの方法は真の入れ子の深さの有限性に合致していないので、両者は相補う関係にあり、これらを結びつけることによって、言語の認識過程のモデルができるのではないかとわれわれは考えている。

【参考文献】

- [1] Woods, W.A., "Transition Network Grammars for Natural Language Analysis", *CACM*, Vol. 13, (October 1970), pp. 591-606.
- [2] Pereira, F.C.N., and D.H.D. Warren, "Definite Clause Grammar for Language Analysis - A Survey and a Comparison with Augmented Transition Networks", *Artificial Intelligence*, Vol. 13, (1980), pp. 231-278.
- [3] Pratt, V.R. "LINGOL - A Progress Report" *Advance Papers of the Fourth International Joint Conference on Artificial Intelligence*, Vol. 1, (1975), pp. 422-428.
- [4] Miller, G.A., "The Magical Number Seven, Plus or Minus Two: Some Limits on Our Capacity for Processing Information", *The Psychological Review*, Vol. 63, No. 2 (March 1956), pp. 81-97.
- [5] Klatzky, R.L., *Human Memory: Structures and Processes* (2nd edition), (1975), W. H. Freeman and Company.
邦訳：「記憶のしくみ I, II」
箱田；中譯
心理学叢書、サイエンス社
- [6] Gazdar, G., "Phrase Structure Grammar", in P. Jacobson and G. K. Pullum (eds.) *On the Nature of Syntactic Representation*, (1981), Reidel: Dordrecht.
- [7] Gunji, T., *A Phrase structural Analysis of the Japanese Language*, M.A. thesis, Ohio State University, (1981).
- [8] Wetherell, C.S., "Probabilistic Languages; A Review and Some Open Questions", *Computing Surveys*, Vol. 12, No. 4 (December 1980), pp. 361-379.
- [9] Hasida, K., *A Real-Time Syntactic Analysis of Natural Languages*, M.S. Thesis, University of Tōkyō, (January 1983).
- [10] Greenberg, J.H., "Some Universals of Grammar with Particular Reference to the Order of Meaningful Elements", in Greenberg (ed.) *Universals of Grammar*, (1963), pp. 73-113.
- [11] Marcus, M.P., *A Theory of Syntactic Recognition for Natural Language*, (1980), MIT Press: Cambridge, Mass.