

179

統語情報を用いて文のわかち書き
を行なうプログラム。

田中穂積
(電子技術総合研究所)

§1. はじめに 西政語の場合には、単語ごとに区切って分かち書きすることが、長い間の慣習として固定し、それが正書法として確立されている。日本語の場合には、漢字と「かな」の使用により、語が連続していても区切りがつけ易く文の理解にさほど大きな影響を及ぼさないということもあり、分かち書きは、あまり行なわれなかった。計算機により日本語を処理する場合には、辞書の記載項目は単語単位であるから、連続した文字列の中から、いくつかの辞書記載項目を識別しなければならぬ。これが、分かち書きの問題である。分かち書きの最も一般的な手法は、最長一致により辞書記載項目を識別してゆくことである。この手法のみでは、分かち書きを自動化するプログラムとして不十分であることはよく知られている。(これについては、§3で考察する。) 漢字かなまじり文の表記を許すことにより、問題が解消することもあるが、それでもうまくゆかない例もある(たとえば文献(1)参照)。

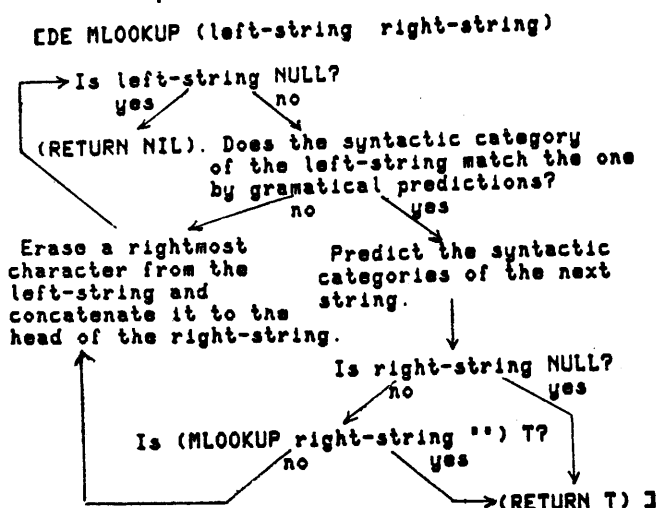
本稿では、最長一致の原則だけでなく、文法規則から得られる統語情報を用いながら、分かち書きと同時に解剖木をも作り出すプログラムを作成したので、その概要と動作例とを示す。統語情報の利用には、MITのPratt⁽²⁾⁽³⁾の開発したプログラムを使用した。なお、このプログラムは、Augmented Context Free Parser である。

§2. 分かち書きを自動化するプログラムの必要性 計算機に日本語を入力する場合、人に分かち書きさせて入力させるシステムが考えられる。しかし、その場合であっても、用言の活用変化に関し、用言を語幹と語尾に分離する操作が必要になる。これを避けるために、各用言毎に辞書記載項目の一部として活用変化表をもたせるシステムもあるが、メモリの使用量は増す。一方、わかち書きの形式自体、単語で切る(あまり実用的ではないといわれている)、文節で切る(詞辞の分離が必要になる)、両者の折衷法などあるが、単語や文節自体、学説によりまちまちで定説はない。我々が日本語文を書くときには、分かち書きしないことが普通であり、計算機に日本語を入力するとき、分かち書きの特定の形式を強制することは、必ずしも自然であるとはいえない。

§3. 統語情報の利用について

統語情報の利用により、従来不可能であった分かち書き

Fig. 1. Main flow of the program for the separation of words.

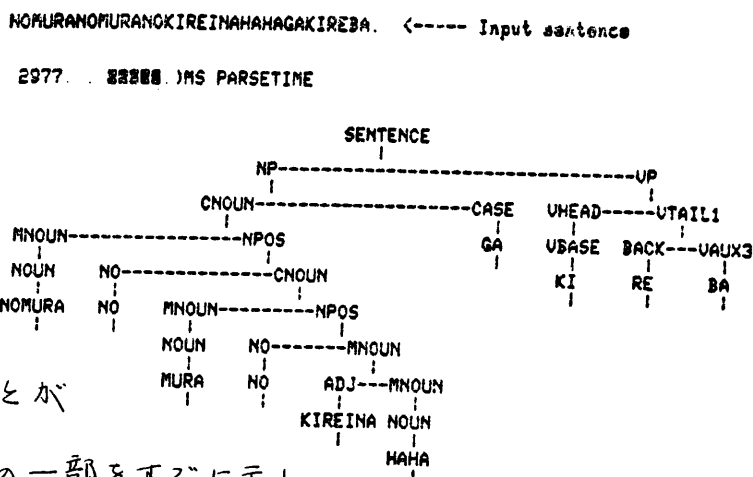


が計算機により可能になることを Fig. 2 の例につき説明する。プログラムは、左から右へと最長一致により辞書記載項目を探す。これにより最左端に NOMURA を認識する。ここで、同時に、文法規則から、次の単語の文法的カテゴリーを予測することができる。Fig. 2 の例の文法規則からは、それは Case であることが予測される。才 2 の辞書記載項目を最長一致で探すと、これも NOMURA であることがわかる。しかし、NOMURA の文法的カテゴリーは Noun であり予測と一致しないので、第 2 の NOMURA のサブパートのなかを最長一致で認識する。この Back up により、才 2 の NOMURA のなかの NO がそれに合致することを見出す。以下順に、このようにして正しい分岐を書きを自動的に得ると同時に、解剖木をも導出することができる。プログラムの概略を Fig. 1. に示す。これから明らかになるように、本プログラムは、統語情報を用いるが左から右へと文字列をスキャンしてゆくのが特長である。比較的文末に位置することの多い動詞を認識して、右から左へスキャンしてゆくシステム(4)もあるが、そうする必要はない。人間の行なう文の理解の仕組みは、むしろ、本稿で述べたような、左から右へ、統語情報を用いるが分解してゆく方式に近いものと想像される。心理学的な根拠を考察してみることも有意義であると思われる。

§4 プログラムの動作例

プログラムの動作例を Fig. 2 に示す。これは、計算機からの出力である。ここで特に注意してほしいことは、辞書には、「着る(又は「切る」)の語幹 KIRE」と、「切れる」の語幹 HIRE が記載されていることである。したがって、この例では、2ヶ所最長一致による誤りが生じ、統語情報により recover されていることがわかる。

Fig. 2. An example.



§5 おわりに

Fig. 2 に、その一部をすまに示すように、本稿で述べたプログラムを、用言の活用処理用として利用できる。また音声認識用として、音素系列から文を推測する場面に応用することも考えられる。

(謝辞) このプログラムは、Pratt の開発した LINGOL⁽²⁾⁽³⁾ をベースにしている。快く筆者に、プログラムのリストを与えられた Pratt 教授に深謝する。また、常日頃、有益な助言を与えて下さる淡路博推論研究室長と、電総研の LISP 作成者でプログラムのデバッグとスピード・アップに関し、協力して下さった東芝総研里川利明氏、LINGOL プログラムの解説に協力して下さった後藤典孝氏に深謝する。

参考文献 (1) 西村超彦:「日本語構文解析ヤマト」, 情報処理学会 CL2-2 (1975, 7.25)
 (2) V.R. Pratt: "A Linguistic Oriented Programming Language", IJCAI, 372-381 (1973)
 (3) V.R. Pratt: "LINGOL - A Progress Report", IJCAI, 422-428 (1975)
 (4) 長尾真他: "自然言語のためのプログラミング言語 PLATON", 情報処理, 15, 9, 654-661 (1974)

137

文法解析システムにおける
予測制御機構について.

田中穂積, 佐藤泰介, 元吉文男
(電子技術総合研究所)

§1. はじめに. より大きな言語理解システムを作成するためには, 柔軟で効率的なパーズング・システムが必要になる. 我々がベースにしたシステムは, Pratt の作成した LINGOL である.⁽¹⁾ その使用経験から, パーズング・システムとして, 特に機能拡張をはかる必要性が感ぜられたものは, つぎの2点である. (1) 形態素分析の強化, (2) 直接的な予測制御.

拡張にあたり, LINGOL の文法記述形式をそのまま保存し, しかも, Pratt のパーズ・アルゴリズムの高速性を生かす, ということを前提にした. (1) については前回, わかり書きを自動化するプログラムについて報告したが,⁽²⁾ このプログラムが形態素分析にそのまま応用できる. 今回報告する (2) の拡張とともに, CF (文脈自由) 文法規則的な記述により, かなり複雑な形態素分析が行なえる.⁽³⁾

§2. 予測制御の必要性. 予測制御の必要性を理解するために, つぎのような規則 (R1) を, CF 文法規則で記述してみる. (R1) --- "格助詞「は」, 「が」には, いかなる格助詞も後続しない." LINGOL のパーザは, 全ての CF 文法規則は無条件に適用可能であると考えている.⁽⁴⁾ したがって規則 (R1) はつぎのように記述される. 格助詞₁ = {が, は}, 格助詞₂ = {を, に, へ, ...} として,

$$\left. \begin{array}{l} \text{格助詞} \rightarrow \text{格助詞}_1, \text{格助詞} \rightarrow \text{格助詞}_2 \\ \text{格助詞}_2 \rightarrow \text{格助詞}_2 + \text{格助詞}_2, \text{格助詞} \rightarrow \text{格助詞}_2 + \text{格助詞}_1 \end{array} \right\} \text{---(R1)'}$$

これから, 格助詞に2つの文法カテゴリを設けたために CF 文法規則数が増大することが理解できる. 助詞の相互承接をさらに精密に規則化しようとするれば, 多数の(作為的な)文法カテゴリと, 多数の文法規則が必要になる.

これを避けるためには, 助詞の相互承接関係をチェックするプログラムを起動して, パーザが CF 文法規則の適用の可否を決めればよい. すなわち, CF 文法規則を条件付きにするのである. このような条件付きの CF 文法規則により, 規則 (R1) は, ただ1つの規則 (R1)" にすることができる. (最終的な形式は §4 の (R1)" に示す.) 格助詞 = {が, は, を, に, へ, ...} として,

$$\text{格助詞} \rightarrow \text{格助詞} + \text{格助詞} \text{ [ただし「が」, 「は」を除き適用可] ---(R1)"}$$

[] 内は, 左の CF 文法規則の適用の可否を決める条件である. このように, CF 文法規則に対し, 規則適用禁止条件を付加できるように拡張することにより,

- (1) 不必要な部分文法解析木の生成を禁止して, それに伴う不必要な文法カテゴリの予測を排除できる. これにより適切な予測が可能になる.
- (2) このような予測制御により, CF 文法規則数を制限し, パーズングの高速化をはかることができる. また適切なわかり書きの自動化が行なえる.

非公式な説明を終り, 以下に実際の予測制御機構について述べる.

§3. 予測制御機構. 条件付き CF 文法規則の記述形式はつぎのようである.

CF 文法規則_i [e_i]; ここで e_i は任意の S 式である. --- (1)

パーザの基本動作を LISP 的に記述するとつぎのようである.

$$\left. \begin{array}{l} (\text{COND} ((\text{NOT} (\text{EVAL } e_i)) <\text{CF 文法規則}_i \text{ を適用}>)) \\ (\text{T} <\text{CF 文法規則}_i \text{ の適用を禁止}>)) \end{array} \right\} \text{--- (2)}$$

e_i が NIL なら, 無条件に CF 文法規則 i が適用されることになる.

§ 4 予測制御用基本関数

(i) *pass-message* : (*pm* <message>)

適切な予測制御を行なうために, パージング制御のメッセージを送る. 関数値は常に NIL.

(例 1) $A \rightarrow B [(pm \ 'x)]$

ノード B が成立すれば, 無条件にノード A が成立する. この時メッセージ x が B より A に送られる.

(ii) *query-message* : (*qm* <message>)

関数 *pm* により送信されたメッセージが, そのノードにきているかを検査する. きていれば関数値は T, きていなければ NIL.

(例 2) $A \rightarrow B [(qm \ 'x)]$

ノード B が成立し, ノード B にメッセージがない時に限り, ノード A が成立する.

(例 3) $A \rightarrow B + C [(qm \ 'x)]$

ノード B が成立し, ノード B にメッセージがないときに限り, C を予測し, 予測文法規則を $A \rightarrow B + C$ に.

(iii) *memo-advice* : (*ma* <S式>), パーザに対する忠告を S 式としてメモしておく. メモされた S 式は, 予測カテゴリが満足された時点で再び評価される. 関数値は常に NIL.

(例 4) $A \rightarrow B + C [(ma \ (NOT \ (qm \ 'y)))]$, ノード B が成立すれば, 無条件に C を予測カテゴリとし, 予測文法規則を $A \rightarrow B + C$ にする. パージングが進み C が成立した時, ノード C にメッセージ y があるときに限り, A が成立する.

これらにより, (R1) はつぎのように書ける. 辞書に, 「を, 「に, 「へ... ← 格助詞, 「が, 「は ← 格助詞 [(~格)] ([] 内は, メッセージ ~ 格を, 格助詞というノードに送ることを意味.), として,

格助詞 → 格助詞 + 格助詞 [(qm \ '~格)] ----- (R1)''

§ 5. おわりに. 予測制御関数を組合せることにより複雑な予測制御を行なうことができる. (R1) に関していえば, (i) の e_i は任意の S 式が書けるから, 助詞の相互承接を記述した表の索引プログラムを書くことにより, CF 文法規則を増やさずに規則の精密化をはかることができる. パーザの制御用として, 意味に近い情報も利用できよう.

§ 6. 謝辞. 日頃有益な示唆をいただく渚一博推論機構研究室長に深謝する.

(1) Pratt, V. R.: "LINGOL--A Progress Report", Proc. IJCAI4, pp. 422-428 (1975).
 (2) 田中穂積: "構文情報を用いた自動わかり書きプログラム", 情報処理学会16回大会予稿179, (1975)
 (3) 田中穂積, 榎山昭一: "活用語尾を処理するプログラム", 本大会予稿, (1976)
 (4) 渚一博: "定理証明としてみた Earley/Pratt のパージングアルゴリズム", 本大会予稿, (1976)
 (5) Heidorn, G. E.: "Augmented Phrase Structure Grammar" in Schank and Nash-Webber ed.: "Theoretical Issues in Natural Language Processing", Cambridge, Mass., pp. 1-5, June, 1975.

