

3M-3

文脈自由文法の漸次的学習システム*

吉岡 智†

松本 雅史^ℓ

中村 克彦‡

東京電機大学大学院理工学研究科[§] 東京電機大学大学院理工学研究科[§] 東京電機大学理工学部[¶]

1 まえがき

われわれは、与えられた正負の文字列例から文脈自由文法を合成するような文法推論の方式の研究を続けており、この方式を実装した Synapse システムを開発している。現在、Synapse では生成規則の数が 10 個程度の文法を合成することができるが、生成規則数が多くなると探索空間が非常に大きくなるため、複雑な文法を合成することができないという大きな問題がある。

本報告では、Synapse システムの最新の文法合成の結果と漸次的学習方式によって探索時間の問題を解決する方法について述べる。

2 変形 Chomsky 標準形

任意の文脈自由文法 (CFG) $G = (N, T, P, S)$ をこれと等価な、Chomsky 標準形の文法に変換できることが知られている。規則集合探索プロセスの効率化のため、われわれは次の形式の規則のみからなる変形 Chomsky 標準形の文法を用いる。

$$A \rightarrow \beta\gamma \quad \beta, \gamma \in NUT.$$

この形式の文法は長さが 1 の文字列を含まないことを除いて一般の文脈自由言語を表すことができる。したがって、変形 Chomsky 標準形の文法は十分に一般的なものであるとみなすことができる。このように 1 つだけの形式の規則に制限することによって、終端記号の個数が少ない言語に対して一般に規則集合を少なくでき、また文法の合成を簡単化できる。さらに、記号列中の各記号 a を 2 つの記号 aa に置換えることによって、言語を $A \rightarrow a$ なる規則の代わりに変形 Chomsky 標準形の規則 $A \rightarrow aa$ を用いた文法で表すことができる。

Synapse: Automatic Synthesis of Context Free Grammars *

Satoru Yoshioka †

Masashi Matsumoto ^ℓ

Katsuhiko Nakamura ‡

Graduate School of Science and Engineering, Tokyo Denki University [§]Faculty of Science and Engineering, Tokyo Denki University [¶]

3 帰納的 CYK アルゴリズム

Synapse システムでは帰納的 CYK アルゴリズムによって正例の文字列と初期規則集合に対して、これを導出するための規則集合が合成される。与えられた文字列が規則集合から導出されるとき、帰納的 CYK アルゴリズムは通常の CYK アルゴリズムと同様に働く。文字列が規則集合から導出できないとき、このアルゴリズムは与えられた記号列を導出するために必要な規則を生成し、規則集合に追加する機能をもっている。以下に帰納的 CYK アルゴリズムの概要を示す。

1. 入力文字列 w に対し、CYK アルゴリズムを適用する。 w の構文解析に成功したならば終了。
2. 構文解析に失敗した場合、CYK アルゴリズムの実行において規則を適用した記号の対の集合 (テスト集合と呼ぶ) を変数 TS に保存する。
3. TS から記号の対を 1 つ選択し、これを右辺とした生成規則を合成して規則集合に加える (この規則の合成には対の選択に加えて、規則の左辺の非終端記号の選択があるが、これらの選択はこの後の処理からバックトラックされるたびにやり直される。) Step 1 に戻る。

Step 3 において、新たな非終端記号 A を用いて規則 $A \rightarrow BC$ を合成した場合、次に合成される規則は右辺に必ず A を含むという制限が可能となる。これによって規則集合の合成の 2 ~ 20 倍の効率化が可能となる。

4 Synapse システム

われわれは Prolog と C 言語を用いて Synapse システムを開発している。Prolog 版では、主にアルゴリズムや探索法の実験・改良を行い、C 言語版ではシステムの高速度および機能の充実を図っている。

Synapse は与えられた正負の例に対して、正例の文字列を導出し、負例の文字列を導出しない文脈自由文法を自動合成する。また、指定によりあいまい・非あいまいの両方の文法を合成することが可能である。システムの動作の概要は次の通りである。

