

日本語文章からの意味フレーム自動生成システム SAGE(Semantic frame Automatic GEnerator) の開発研究

尾見 孝一郎 (15594035) 佐々木 毅 (15594057) 林田 和久 (15594094)

原 田 研 究 室

1. 初めに

過去に原田研究室において、日本語要求仕様を分析し、オブジェクト指向設計図を自動生成する研究を行っていた。このシステムはCAMEOと呼ばれる。設計要素の抽出や設計図への配置は既に完成しているが、日本語文を分析するのに必要なフレーム形式への変換は手作業であった。そこで、日本語文章を解析し、その内容を文中の語の意味と他の語との関係を、統一した表現形式である意味フレーム形式として自動生成するシステム「SAGE(Semantic frame Automatic GEnerator)」を開発する事にした。

なお、日本語文章を形態素に分解していく作業にはChasen(奈良先端大学松本研)、日本語文章を文節に分解し各文節と述部との従属関係を解析する作業にはBEST(奈良先端大学松本研)を、それぞれ使用する。

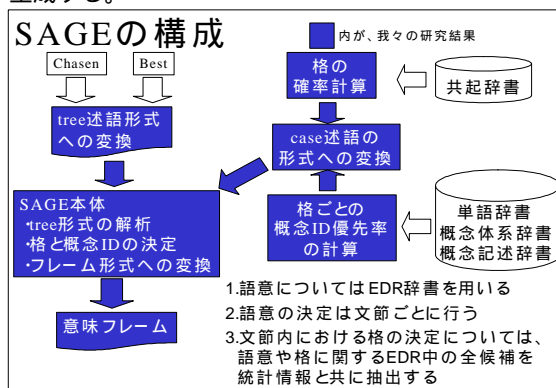
2. 意味フレーム自動生成システム -SAGE-

日本語文章を解析するためには、巨大な辞書が必要となる。我々はこれについて、良質の語意や格(語彙間の関係)に関する文例を多く持つEDR電子化辞書(以下EDR)を利用する事にした。

SAGEは

- ・ Tree 述語形式への変換ルーチン
- ・ EDR 辞書検索
- ・ SAGE 本体

より成り立ち、CAMEOへ受け渡す意味フレームを生成する。



2.1. SAGEデータ変換ルーチン

Chasen/BESTからの出力は、そのままではSAGEへの入力には使えない。そこで我々は、Prologの述語形式に対応するデータ構造としてTree述語を考案し、これに変換する事にした。

2.2. EDR検索ツール

これは二つのルーチンから成り立つ。

1) 格の確率計算ルーチン

文節内の助詞と述部との格関係をEDRから検索、統計を取る。これをある格を取る確率、即ち「格の確率」とする。

2) 格ごとの概念ID優先率計算ルーチン

・ 単語と品詞の情報から、EDRの単語辞書を検索して概念IDを得る。

・ 概念体系辞書を検索し、概念IDの上位概念を得る。

・ 概念記述辞書を検索し、概念間の格、及び概念間の距離を得る。

・ 概念間の距離から、同一格間における、その組み合わせが成りうる概念ID優先率を求める。

2.3. メイン推論ルーチン

一つの形態素は通常複数の概念IDを持つため、これを特定した文の構造表現(解釈木)は一つになるとは限らない。そこで、SAGEは最適な解釈木を得る指標として、確信度という概念を考えた。

確信度は、EDR検索ルーチンにより求められた概念ID確率(P_c)を基に、確信度 $CP = P_c$ として定義する。

各解釈木の確信度を比較し、確信度が最大となる解釈木を意味フレームの抽出に用いる。

なお、意味フレームは、frame(フレーム番号、文中単語、読み、EDR品詞情報、基本語、日本語品詞情報、活用形、概念ID、格情報、文章番号)で構成される。

3. 終わりに

我々は、SAGEにより日本語文章から意味フレームを自動生成することに成功した。この事により、SAGEはCAMEOを運用する際の有用な支援ツールになった。さらに、その他の自然言語処理問題に対しても適用が容易であることが示せた。

今後の課題としては、文脈解析(複文・重文解析)の精度を上げることにより、より多様な日本語文章に対処していくことが可能である。

4. 参考文献

長尾 真 編、「自然言語処理」、岩波書店、1996
『茶釜』 ホームページ奈良先端大学松本研究室

<http://cactus.aist-nara.ac.jp/lab/nlt/chasen.html>