

# 解析事例を用いた係り受け曖昧性の解消

3P-5

NTT情報通信網研究所

武石 英二

## 1. はじめに

柔軟な自然言語処理システムを実現するためには、対象とする文章から言語知識、解析知識を取得し、文章の固有性や特徴をシステムに反映するための枠組みが必要だと考えられる。

そのため、機械翻訳における変換規則の獲得や意味解析における意味表現の獲得の枠組みとして、解析・変換結果を蓄積し、それらを直接利用する手法が提案されている<sup>1)・2)</sup>。一方、単語間の共起関係や係り受け関係から解析ルールを取得する手法も検討されている<sup>3)</sup>。しかし、前者の枠組みを統語構造の解析に適用したものは少なく、後者は係りと受けの2項間関係や動詞を中心とした構造を対象とするものが多いため、文全体の構造的な特徴を取得することは難しいと考えられる。

このような背景のもとで、解析事例を用いた係り受け解析の検討を進めている。これは、係りと受けの2項間の関係を解析対象とし、その前後の文構造を解析対象をとりまく環境(単語の共起や文構造の固有性)として捉えようとするものである。

## 2. 解析事例を用いた係り受け解析の構成

係り受け解析方式の全体構成を図1に示す。これは、事例に基づく手法の一般モデル<sup>4)</sup>をもとにしたものである。点線内が実際に解析を行う部分であり、点線外が解析の成否をフィードバックする部分である。事例に基づく手法の精度は、事例データの固有性に依存する部分が大きいので、処理結果をフィードバックする機構が必須だと考えられる。本稿では、実解析部分のうちの<sup>5)</sup>部分について述べる。

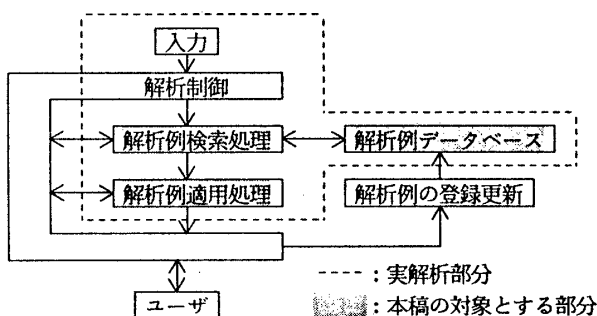


図1. 解析事例を用いた係り受け解析の構成

## 3. 解析事例の表現方法

入力となる文節列とその解析結果を対して保存しておくだけでは、適用することのできる文が大幅に限定される。したがって、解析事例そのものを一般化したり、解析事例間を組織化することによって、適用することのできる範囲の拡大を図る必要がある。

本検討では、係り受けネットワークによる解析事例の表現を採用した。解析事例を、文節中の自立語をノード、主として付属語によって表現される係り受け関係をリンクとしたネットワークとしてとらえ、すべての解析事例に対する係り受けネットワークを組み合わせたものを解析事例全体と考える。また、正しい係り受け関係だけでなく、誤った係り受け関係もリンクとして保存する(実際には、既存の係り受け解析<sup>6)</sup>から出力される曖昧性を含んだ係り受け構造の正誤を人手によって判断して用いた)。以下では、正しい係り受け関係を示すリンクをリンクタイプA、誤った係り受け関係を示すリンクをリンクタイプBと呼ぶ(図2)。

リンク $l$ が始点ノード $n$ 、終点ノード $n'$ 、係り受け関係 $k$ 、リンクタイプ $t$ 、リンクの出現する解析事例の全体 $\{s\}$ から構成され

ると考えると、解析事例の全体 $X$ は、リンクの全体 $X = \{l\}$ 、 $l = (n, n', k, t, \{s\})$ とみることができる。

以下では、リンク $l$ の始点ノード、終点ノード、係り受け関係、リンクタイプ、リンクの出現文の全体を、それぞれ $n(l)$ 、 $n'(l)$ 、 $k(l)$ 、 $t(l)$  [ $=A$ または $B$ ]、 $\{s\}$ と表す。

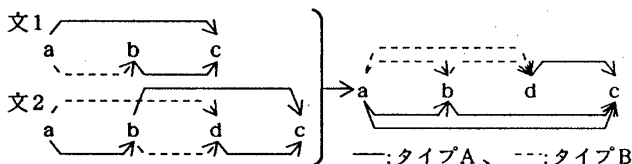


図2. 解析事例の表現方法

解析事例全体をネットワークとして表現することによって、

- ・着目する係り受けをもつ解析事例を容易に取得することができる、
- ・共通の部分構造をもつ解析事例を容易に取得することができる、
- ・それによって、類似した解析事例の中から適切な解析事例を選択することができる、

という効果を期待することができる。また、正しい係り受け関係だけでなく、誤った係り受け関係も蓄積することによって、

- ・不適切な解析事例の適用を抑止することができる、
- ・誤った係り受け関係を積極的に適用することによって、係り先の限定、特定を行うことができる、

という効果を期待することができる。

## 4. 解析事例による係り受け曖昧性の解消効果

3章で述べた解析事例の表現方法を用いて、係り受け曖昧性の解消法について検討し、評価実験によって有効性の検証を行った。概要は次の通りである。

既存の係り受け解析の結果として、文節とその係り先文節(一般に複数個)のリストが出力される。係り先文節が複数個出力された文節に対して、その文節および係り先候補の文節をキーとして蓄積した解析事例データからもっとも適切だと考えられる解析事例を選択し、解析事例中の係り受け関係を用いて係り先文節の特定を行う。また、曖昧性の解消は、文頭文節から順に1文節ずつ行う。

### 4.1. ノード間、リンク間、文間の類似度

#### (1) ノード間類似度

既存の解析処理に使用している日本語辞書<sup>7)</sup>には、単語の意味的な属性(意味カテゴリ; 約2800分類)が記述されている。この意味カテゴリを40に粗く再分類したカテゴリを用いてノード間の類似度を求める。一般に1つのノード $n$ には、複数のカテゴリが付与されており、これを $C(n)$ で表す。2つのノード $n_1$ と $n_2$ の類似度 $F$ を、ノードカテゴリの一致度によって次のように求める。

$$F(n_1, n_2) = \frac{\# \{C(n_1) \cap C(n_2)\} + \alpha f(n_1, n_2)}{\max \{\#C(n_1), \#C(n_2)\} + \alpha}$$

ここで、 $\alpha$ は定数、 $\#A$ は $A$ の個数を示しており、 $f(n_1, n_2)$ は $n_1$ と $n_2$ の表記が一致する場合に1、一致しない場合に0をとる。

#### (2) リンク間類似度

2つのリンク $l_1$ 、 $l_2$ の類似度 $G$ を、次のように求める。

$$G(l_1, l_2) = \begin{cases} \frac{F(n(l_1), n(l_2)) + F(n'(l_1), n'(l_2))}{2} & [l_1, l_2 \text{ のリンクタイプ、係り受け関係が一致する場合}] \\ 0 & [\text{その他の場合}] \end{cases}$$

#### (3) 文間類似度

処理対象文のノードの並びを $S = (n_1, \dots, n_p)$ 、解析事例を構成する文節に対応するノードの並びを $S' = (n'_1, \dots, n'_q)$ とすると、 $S$ と $S'$ の文間類似度 $H$ を次のように求める。

$$H(S, S') = \frac{\sum_{i=1}^p \max_{j=1}^q \{F(n_{i,}, n'_{j,})\}}{p}$$

4.2. 解析例の選択、適用

処理対象文のノードの並びを  $S = (n_{1,}, \dots, n_{p,})$ 、曖昧性解消の対象ノードを  $n_m$ 、その係り先候補のノードを  $n_{s,1}, \dots, n_{s,r}$ 、 $n_m$  から  $n_{s,1}$  への係り受けを  $l_{m,1}, \dots$ 、その係り受け関係を  $k$  で表す。

(1) 解析例の選択

$[F(n_m, n(l)) > \beta_1, F(n'_{s,i}, n'(l)) > \beta_2 \text{ なる } n'_{s,i} [1 \leq i \leq s_r]]$  が存在する、 $k(l) = k$  …(※)

を満たすリンク  $l$  をもつ解析例を適用候補として選択する。 $\beta_1, \beta_2$  は定数(しきい値)である。

(2) 解析例の適用

選択された解析例を次の手順で適用する。

① 選択された解析例のうち、処理対象文との文間類似度をもっとも大きい解析例を選択する。

② 解析例内に(※)を満たすリンクが複数個 ( $l_1, \dots, l_t$ ) 存在する場合には、

$$G(l_1, \dots, l_t) = \max_{j=1}^t \max_{i=1}^r G(l_{m,}, l_{s,i,})$$

を与えるリンク  $l_v$  をすべて選択する。

③  $l_v$  のリンクタイプによって、 $n_m$  から  $n_{s,i}$  への係り受けの成否を決定する。

$t(l_v) = A$  のとき、 $n_m \rightarrow n_{s,i}$  : 成立

$t(l_v) = B$  のとき、 $n_m \rightarrow n_{s,i}$  : 不成立

4.3. 評価実験

(1) 実験対象データの概要

2種類の計算機マニュアル(1, 2)を用いて曖昧性解消の評価実験を行った。マニュアル1の1章分(計算機の構成説明の部分)を解析対象文とした。解析対象文の概要を表1に示す。表1において多義数とは、曖昧性のある箇所(総数を示しており、候補2, 3, 4とは、それぞれ係り先の候補数が2, 3, 4である箇所数を示している。また、マニュアル2の対応する章を解析例データX、マニュアル2の他の1章分を解析例データYとし、解析例データX、X+Yの2種類の解析例データを用いて実験を行った。解析例データの概要を表2に示す。表2において、各数値は個数を表しており、A、BとはリンクにおけるタイプA、Bの内数である。

表1. 解析対象文の概要

総文数	多義数	候補2	候補3	候補4
67	95	78	12	5

表2. 解析例データの概要

データ	総文数	ノード	リンク	A	B
X	67	322	635	516	119
X+Y	133	560	1194	940	254

(2) 実験内容

表3に示す3種類の実験を行った。×は解析例の間の順位付けを行わないことを示している。また、前節、前々節で導入した定数の値として、 $\alpha = 3, \beta_1 = \beta_2 = 0$  を用いた。

表3. 評価実験内容

	解析例選択方法	解析例適用方法
実験1	表記一致	×
実験2	ノード間類似度F	×
実験3	ノード間類似度F	文間類似度H

(3) 実験結果

実験結果を図3に示す。太線がデータX+Yに対する結果を示しており、○は正しい係り先を認定した(適用成功)件数、丸数字はタイプBのリンクが適用された内数、破線は同一類似度をもつ複数の矛盾するリンクが選択されたため、係り先が一意に決定しなかった件数を示している。また、折れ線は適用成功率を示している。

$$(\text{適用成功率}) = \frac{(\text{適用成功件数})}{(\text{適用成功件数} + \text{適用失敗件数})} \times 100$$

解析例間の順位付けを導入した場合、適用成功率をあまり低下させず(Xの場合で81.8%から80.4%、X+Yの場合で87.7%から79.4%)に、

選択成功数がXで20から61、X+Yで24から66へ改善されたことから、ノード間類似度、文間類似度が有効に機能していることがわかる。

解析例データYの追加による効果については、選択成功件数、適用成功件数に増加が見られることから、ある程度の有効性を確認することができる。しかし、着目する係り受け曖昧性とあまり関連がない解析例が雑音となって適用が失敗する場合があります、また解析例数の増加に比べ、選択成功数の増加はそれほど大きくない。

リンクタイプBが適用された件数も比較的多い(32%から50%)ことから、誤った係り受け関係を保存することが有効であることがわかる。

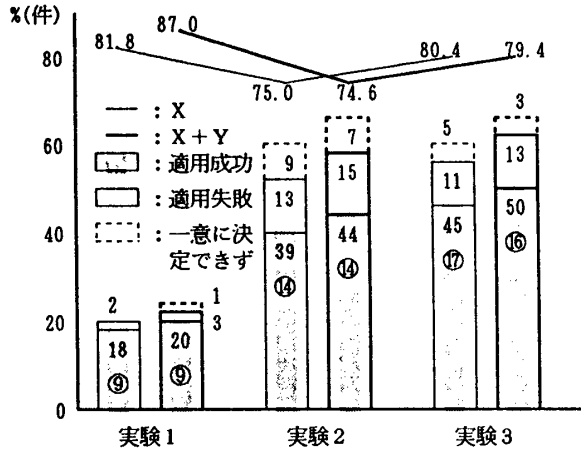


図3. 評価実験結果

(4) 考察

・解析例の選択、適用方法の改良

現在、粗分類して用いている意味カテゴリーの利用方法の改良、リンクの出現頻度の導入によって、ノード間類似度の算出方法を精密化すること、解析例として保存された構造を利用することによって、文間類似度の算出方法を改善することが考えられる。

・解析の制御

どこから解析(曖昧性の解消)を行うかを決定する、解析例中の構造やすでに解析済みの部分構造を利用することによって、解析例の部分構造を適用したり複数の解析例を組み合わせて適用するなど、解析を制御する機構について検討する必要がある。

・解析知識の利用

現状では、たとえば「接続した機器」と「機器を接続する」の対応関係を認定できない、「AとBを利用する」に対して「Aを利用する」を認定できないなどの問題点がある。態やアスペクト、並列関係などの利用方法についても検討する必要がある。

5. おわりに

解析例を用いた係り受け解析の構成、解析例の表現方法および解析例を用いた係り受け曖昧性の解消効果について述べた。今後は、4章で考察した課題について検討を進めるとともに、2章で示した係り受け解析全体の構成についても具体化を進めていく予定である。

謝辞

本検討に対して数多くの貴重な助言を頂いたNTT情報通信網研究所、林良彦主任研究員に深く感謝致します。

参考文献

- 渡辺: 文法的情報を含んだ翻訳例を用いるトランスファーステム, 情処第43回全国大会論文集, 2H-4 (1991).
- 島津: コーパス解析に基づく事例ベースパーザ, 情処人工知能研究会資料, 75-10 (1991).
- 小松他: 係り受け事例からの結合値フレーム獲得の試み, 情処第43回全国大会論文集, 3H-2 (1991).
- 小林: 事例ベース推論の研究課題, 情処人工知能研究会資料, 75-4 (1991).
- 白井: 日英翻訳システムALT-J/Eにおけるテーブル駆動型係り受け解析法, 情処第34回全国大会論文集, 5W-6 (1987).
- 池原他: 言語における話者の認識と多段翻訳方式, 情処論文誌, Vol. 28, No. 12 (1987).