

A-2-1 意味解析に基づく並列名詞句の構造解析

Coordinate Structure Analysis based on Semantic Analysis

田村直良, 田中穂積 (東京工業大学 工学部 情報工学科)

本論文では、自然言語の意味解析に重要である並列構造解析について述べる。本方式では、並列助詞「と」と連体助詞「の」、名詞のみから成る名詞句を対象とする。評価原理としては、連体助詞「の」の修飾関係が成り立つ導出木について類似度等を考慮した評価関数を用いて、その値の最大なものを意味構造とするものである。名詞の数4以下の58例について本方式を適用したところ、95%の正解率で正しい構造決定が行えた。

1. はじめに

自然言語において、並列構造は、文の曖昧性を増やす主たる要因であり、文の意味解析、意味抽出には、文が持つ並列構造を正確に把握することが非常に重要な問題となる。並列構造決定の手法として、[6]では、構文的な10個の「手がかり」を用いた名詞句の解析手法を述べている。構文的な範囲に納めた理由として、機会翻訳のような応用では、個々の名詞ごとに詳細な意味記述が困難であることをあげている。しかし、おそらく人間は、主として意味解析から並列構造を決定していることであろう。よって、このためにどの程度の解析、意味記述が必要か、その結果どの程度の判断ができるかについて調べることも興味ある題材である。

自然言語において、並列される文の要素としては、文、用言、体言その他ほとんどの要素が考えられるが、我々の研究では、簡単のために、並列助詞「と」と、連体助詞「の」、名詞からなる名詞句のみを対象とし、さらに、「と」の出現の個数は、高々一個と制限する。名詞句のみからその並列構造を決定しようとする場合、その名詞句の前後の修飾関係、文や文脈の意味等を考慮する必要があり、また、根本的に曖昧性のある句もあって、一般にはすべての名詞句についての構造決定は不可能である。本研究での目的は、上記の形式の名詞句のみから、連体助詞「の」が結ぶ名詞の意味関係と、並列される名詞の類似度から、人間が判断できる程度に解析できるような解析手法の確立を目指すことにある。

2. 並列構造決定の評価原理

この章では、本方式の解析手法の原理について説明する。

2.1 扱う名詞句の形式と意味構造

まず、処理対象とする名詞句の形式と、この形式から得られる意味構造の一般形式を定義する。我々は、名詞が連体助詞「の」のみにより連結されている名詞句の意味構造を、前の名詞が直後の名詞に係るという一次元的な構造に仮定している。こうすると、(2-1)式に対する意味構造は、図1のような構文木で代表される。(2-2)式は、図1に対応している。

[定義] 並列名詞句の形式と意味構造

a) 評価対象とする並列名詞句は、次の形式である。

$$A_1 \text{の} \dots \text{の} A_m \text{と} B_1 \text{の} \dots \text{の} B_n \quad (2-1)$$

ここで、各 $A_i, B_j (1 \leq k \leq m, 1 \leq l \leq n)$ は名詞、「の」は助詞の「の」を表す。

b) (2-1)式で表される並列名詞句の意味構造を

$$A_1 \text{の} \dots \text{の} (A_k \text{の} \dots \text{の} A_m \text{と} B_1 \text{の} \dots \text{の} B_l) \text{の} \dots \text{の} B_n \quad (2-2)$$

のように定義する ($1 \leq k \leq m, 1 \leq l \leq n$)。ここで、「 $A_k \text{の} \dots \text{の} A_m$ 」と「 $B_l \text{の} \dots \text{の} B_l$ 」が並列する名詞句であり (A_m と B_l が並列する名詞である)、「 $A_1 \text{の} \dots \text{の} A_{k-1}$ 」は、並列名詞句「 $A_k \text{の} \dots \text{の} A_m$ と $B_l \text{の} \dots \text{の} B_l$ 」に係り、さらにこの並列名詞句は、「 $B_{l+1} \text{の} \dots \text{の} B_n$ 」に係る。

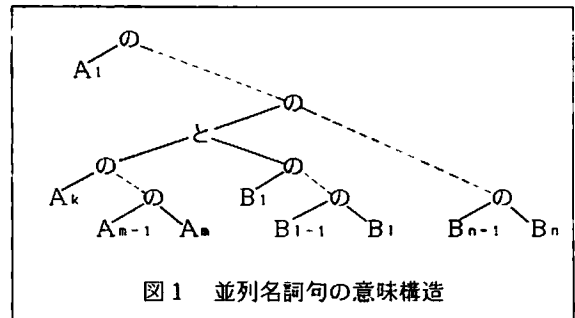


図1 並列名詞句の意味構造

名詞句の並列構造を決定することは、定義1のもとで、 $k (1 \leq k \leq m)$ と $l (1 \leq l \leq n)$ を求めることに帰着される。

2.2 連体助詞「の」の意味解析

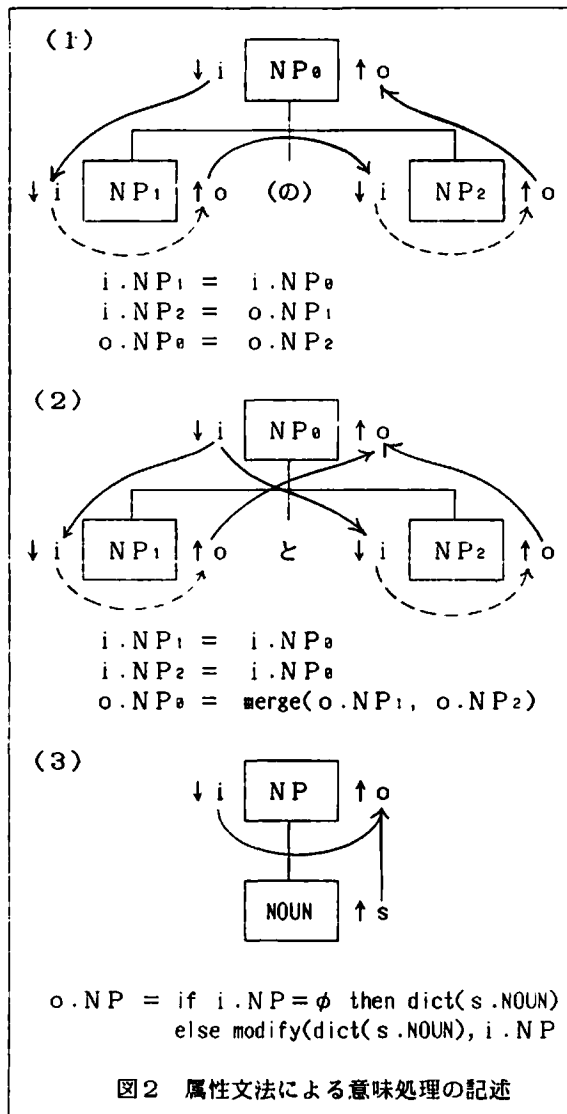
連体助詞「の」による接続を、前の名詞が後の名詞を修飾するもの仮定しているため、名詞句全体では、修飾関係による「意味の流れ」に対応する情報の列ができる。この「流れ」は、導出木の根より出発し、導出木上を左から右へ inorder に traverse しているが、「と」に対応するノードでは、左右の部分木へ分岐し、このノードで両部分木で得られた情報がマージされて上へ返される。また、導出木の葉の部分、つまり各名詞においては、「の」による修飾関係の解析がなされて、その結果の情報が列に加えられる。この様な処理を形式的に記述するために、属性文法を用いる(図2)。

図2において、NP は名詞句を表す非終端記号、NOUN、「の」、「と」は、名詞、および助詞「の」、「と」に対応する終端記号である。また、非終端記号 NP は、相続属性(導出木の根から葉の方向へ流れる属性) i と、合成属性(導出木の葉から根の方向へ流れる属性) o を

持つ。また、終短記号NOUNは、名詞を構成する文字列に対応する合成属性 s を持つ。終短記号が表す文字列の領域を S で、意味の領域を M で表すと、関数 $dict$ 、 $modify$ 、 $merge$ は、それぞれ、以下のような領域で定義される。

$dict : S \rightarrow M$
 $modify : M \times M \rightarrow M$
 $merge : M \times M \rightarrow M$

ここで、 $dict$ は、終短記号が表す文字列について辞書を引き、対応する意味を得るための関数である。 $modify$ は、第1引数が表す名詞が、第2引数が表す名詞句に連体助詞「の」により修飾されたときの意味を返す関数である。 $merge$ は、第1引数、第2引数ふたつの意味構造を統合した並列構造を返す関数である。このうち、関数 $modify$ は、両引数が意味がある修飾関係のときのみ定義されているものとする。この制限は、不合理な修飾関係を持つ導出を除外するためにも利用される。



2.3 並列構造評価関数

前節の手法により得られる導出木は、関数 $modify$ によ

り一応意味のあるもののみ選択される。一般に、一つの名詞句からは複数の導出木が得られるが、これらの中からどれが最も適したものであるかを決定するために、導出木にある値を対応づける並列構造評価関数 $eval$ を導入する。この関数 $eval$ は、図1の形式の導出木に対して、次のような原理に基づいて導出木の意味的な「正しさ」の目安を計算する。

図1において、

- (a) 並列する名詞 A_n と B_1 の類似度が大きいほどの $eval$ の値は大きい。
 - (b) 並列する名詞 A_n と B_1 の類似度が大きいときは k の値が大きいほど $eval$ の値は大きく、類似度が小さいときは k の値が小さいほど $eval$ の値は大きい。
- (a) は、並列する名詞を選択するための手がかりであり (l の決定)、(b) は、並列する名詞の類似度が大きいときは前方にある名詞は両方の並列名詞に係る傾向が強く、類似度が低いときは前方の名詞は「と」の直前の名詞のみに係る傾向が強いとの観察による (k の決定)。以上から評価関数 $eval$ を次のように定義する。

[定義] 評価関数 $eval$

$$eval(t) = k \cdot a_1(A_n, B_1) + a_2(A_n, B_1) \quad (2-3)$$

ここで、 t を図1の形式の導出木とし、 A_n と B_1 は並列する名詞句である。 k は並列する名詞句に係る名詞句の長さ、 a_1 、 a_2 は類似度を表す関数である。ただし、両関数とも類似度が大きくなるにつれて値が大きくなる単調増加関数であるが、 a_2 が常に正の値をとるのに対して、 a_1 では、類似度が小さいときには負の値を取る。

この関数 $eval$ では、「の」による修飾関係については考慮されていない。これについては、第5章で触れることにする。

3. 並列構造解析システム

この章では、解析システムの実現について述べる。

3.1 概要

システムはLangLAB上に構築されている。評価手順は、1) 構文解析、2) 連体助詞「の」の修飾関係評価、3) 類似度および関数 $eval$ の評価、と進むが、これらの処理は、入力された単語の列に対して構文解析が非決定的であることから、組合せ的に可能なすべての意味構造を作り出すように行われる。「の」の修飾関係評価では、2.2で述べたように、意味情報が導出木を $inorder$ に $traverse$ する形で流れる。このため、これまでの導出木が下から生成されるボトム・アップ形のパーザでは構文解析と同時に意味情報の評価を行えなかった。我々は、BUP-XGを応用した手法[1]を用いて、構文解析用文法に「の」の修飾関係評価の部分の補強項として記述している。この記述は、同時に不合理な「の」の修飾関係を除外するための制約条件にもなっている。1) ~ 3) の結果得られた意味構造について $eval$ を評価し、評価値が最大の構造を求める結果として出力する。なお、この方式で必要な名詞の意味記述は、「の」の修飾関係を決定するためと類似度の計算のために必要な、シソーラス上での分

類位置と、「束ね」の名詞であるかどうかの指定だけである。

3.2 「AのB」の解析

「AのB」の解析は、前節で述べたように、解析により得られた意味構造のうち、「AのB」の意味として妥当でないものを除外するために用いられる。ここでの方式は、[3,4]をDCKR[5]で実現したものである。概念の検索は、シソーラス[9]をDCKRで記述したものをを用いている。

「AのB」の解析の方法の概略を以下に示す。

case 1.

Aが、「人-有意志体」、「社会的集団」、「物」、「行為活動」、「抽象概念」、「時」、「場所」で、
Bが、「行為活動」のとき、

Aが格要素、Bが述語の名詞形なる構造とする。

case 2.

Bが「時」、「場所」、「目的」などのとき、

Aに対しBが相対的な場所、時等を示す構造とする。

case 3.

Aが「物」、「人-有意志体」、「社会的集団」、「抽象概念」で、

Bが「属性」、「精神的活動源や心情」のとき、

BがAの数量、性質等を表す構造とする。

case 4.

Aが、「行為活動」、

Bが、「人-有意志体」、「社会的集団」、「物」、「行為活動」、「抽象概念」、「時」、「場所」のとき、

Aが格要素、Bが述語の名詞形なる構造とする。

case 5.

その他の場合（所有関係、連辞関係、方法関係などの関係をDCKRにより単純化したもの。詳細については省略する）、

AとBとがある述語の格要素を表す構造とする。

3.3 評価関数 evalの実現

まず、二つの名詞AとBのシソーラス上の距離を定義する。

【定義】 距離

名詞AとBのシソーラス上の距離を次のように定義する。

1. 同一の名前のとき、0。

2. シソーラス中の高さ*i*の部分木にAとBが含まれるとき、*i*。

我々は、この距離を用いて関数eval中の a_1 、 a_2 を表1のように決定した。

数値の決定は数回の実験から適当なものを選んだものであるが、 a_2 の値を大きくとることにより並列される名詞の選択を優先させ、次に「の」の修飾構造を決めようとの意図による。

評価関数 evalは、

$$\text{eval}(t) = k \cdot a_1(A_n, B_1) + a_2(A_n, B_1) \quad (2-3)$$

のように定義されているが、並列構造の直後の名詞がい

わゆる「束ね」の名詞[7]であるときは、評価値が高くなるようにして（具体的には、500を加えている）、このような構造がより尤らしいと判断できるようにしている。

距離	a_1	a_2
0	200	3000
1	110	2000
2	40	1500
3	10	500
4	-150	100
5	-200	0

表1 関数 a_1 と a_2

4. 実験

[8]中の65例の名詞句を用いて実験を行なったところ、evalの評価値について次の5とおりの結果が得られた。

1. 正しい意味構造に対して一番大きい値が与えられ、かつ誤った意味構造に対してはそれより小さい値が与えられているもの ……40個。

例) 「育てるためのくふうと努力」

2. 正しい意味構造が2個あり（人間にとっても曖昧な句）両者に同じ評価値が与えられているもの ……4個。

例) 「小学生と高校生の三人の子供たち」

3. 正しい意味構造に対して一番大きい値が与えられているが、同時に誤った意味構造にも同じ値が与えられているもの ……4個。

例) 「経済無視と国民生活の悪化」

4. 「束ね」の構造として正しい意味構造にのみ一番大きな値が与えられたもの ……9個。

例) 「目と目の間」

5. 正しい意味構造に一番大きな値が割り当てられていないもの ……8個。

例) 「世界の世論と米国民」

正しい意味構造に一番大きな値を与えているものを正しい構造決定とするなら、上記実験における正解率は、88%（1～4）である。また、誤った意味構造に最大の値を割り当てたものを不正解とするならば、正解率は、82%となる（1, 2, 4）。さらに、名詞の数が4個以下の名詞句58例のみに制限するならば、正解率は、それぞれ95%、87%であった。上記2については、人間にも曖昧性があると感じられる名詞句の意味構造に同じ評価値が与えられており、興味深いところである。

図3に解析結果の一部を示す。

5. 問題点および限界

本方式では、かなりの単純化を行っているために、問題点もある。この章では、前章で述べた3、5について考察する。

1. 関数 a_1 と a_2 の選び方の違いにより生じる不正解

例: 「若さと政務次官の肩書」

解: 「(若さと政務次官)の肩書」

これは、関数 a_1 と a_2 の選び方により生じる不正解である。 a_1 と a_2 を帰れば一つの不正解は訂正できるが、その他の不正解が新たに生ずる恐れがある。統計的手法

などで最適の a_1 と a_2 を決定できるだろうが、それだけではこの種の不正解を完全に取り去ることはできない。

2. 矛盾する修飾関係を受け入れたことにより生じる不正解

例：「反政府の仏教系と政府支持のカトリック系労組の分裂メーデー」

解：「反政府の（仏教系と政府支持のカトリック系労組）の分裂メーデー」

本方式では、修飾関係の分類までは行っているが、属性スロットのチェックなどは行っていない。これに対処するとすると、名詞の意味記述がソーラス上の分類関係だけでは収まらなくなってしまう。

3. 長い名詞句

例：「企業の中の技術革新の進展と職種の変化」

解：「（企業の中の技術革新の進展と職種の変化）」

本方式では、名詞の数に応じて意味構造の組合せ数は、 $O(n^2)$ で増える（式2-2において k と l の選び方）。これに対して、評価関数は単純過ぎる。

4. eval では、「(AとB)のC」と「(AとBのC)」の区別がつかない

例：「横滑りと視界の良くないこと」

解：「（横滑りと視界の良くないこと）」

「（横滑りと視界の良くないこと）」

本方式では、「(AとB)のC」と「(AとBのC)」については、「AのC」、「BのC」が同一の修飾関係であるかどうかでしか両名詞句の区別がつかない。「の」の修飾関係を eval に導入することが必要となろう。

6. まとめ

本報告では、意味解析を用いた並列名詞句の構造決定の手法について述べた。本方式は、「の」の修飾関係が成り立っている意味構造に対して、並列する名詞の類似度等を基にした評価関数 eval を用いて意味構造に正しさの順位をつけるものである。

本方式により65例の名詞句について構造決定を行ったところ、正解に対して一番大きな評価値が与えられなかったものが8例（12%）であった。また、名詞句中の名詞の数を4以下の58例に限れば、3例（5%）であった。

〔参考文献〕

1. 奥村, 田村, 徳永, 田中: BUP系自然言語解析システム上でのトップ・ダウンな情報の制御について, 本大会予稿集, 1986.
2. 首藤, 吉村, 津田: 日本語技術文における並列構造, 情報処理学会論文誌, Vol.27, No.2, 1986.
3. 島津, 内藤, 野村: 日本語文意味構造の分類—名詞句構造を中心に, 自然言語処理研究会, 47-4, 1985.
4. 島津, 内藤, 野村: 助詞「の」が結ぶ名詞の意味関係のsubcategorization, 自然言語処理, 53-1, 1986.
5. 田中, 小山, 奥村: 知識表現言語DCKRとその応用, コンピュータ・ソフトウェア, Vol.3, 4号, 1986.
6. 長尾, 辻井, 田中, 石川: 科学技術論文における並列句とその解析, 自然言語処理研究会, 36-4, 1983.
7. 水谷, 田中: 語の並列結合子, 計量国語学, No.63, 1972.
8. 電総研: 新編日本語品詞列集成, 1979.
9. 内田裕士, 荻野孝野: ICOT Technical Memo, to appear.

<<意味構造>>

No. 1

並列構造

1540, 0+, {くふう#3, individual_of(くふう)}, {努力#4, individual_of(努力)}

意味

{くふう#3, (case5ed: {育てるため#2, individual_of(育てるため)}, individual_of(くふう))},

{努力#4, (case5ed: {育てるため#2, individual_of(育てるため)}, individual_of(努力))}

解析木

1-名詞句

1-名詞句

1-1-名詞

1-1-1-名詞 — 育てるため

1-助詞 — の

1-名詞句

1-1-名詞句

1-1-1-名詞

1-1-1-1-名詞 — くふう

1-助詞 — と

1-名詞句

1-1-名詞

1-1-1-名詞 — 努力

図3 出力例