

## 非命題的意味解析のための日本語文末表現意味体系

本田 聖晃<sup>†</sup> 田辺 利文<sup>†</sup> 吉村 賢治<sup>†</sup> 首藤 公昭<sup>†</sup>  
<sup>†</sup>福岡大学工学部

## 1 はじめに

近年, Web の普及による文書情報の爆発的な増加に伴い, より精度の高い自然言語処理の必要性が認識されている. 特に, 意味を考慮した処理システムの需要が増加していることは Semantic Web 研究等の発展を見ても明らかである. 筆者らは, 意味を考慮した言語データ処理においては, 複単語表現 (Multiword Expression: MWE) を組み込んだシステム構築が不可欠であると考え. MWE に関しては, 2002 年の論文「Multiword Expressions: A Pain in the Neck for NLP」(Sag et al. 2002) を皮切りに 2007 年夏までに ACL において MWE のワークショップが計 4 回実施されるなど, 重要性が世界的に認識されつつある.

これまで筆者らは, 熟語性, 語彙の一体性, 確率的束縛性のうち少なくとも 1 つの性質を持つ単語列を MWE として収集してきた. 熟語性とは, 構成している単語の通常の意味から全体の意味を構成するのが難しいことを意味する. また, 語彙の一体性とは分離しにくさを, 確率的束縛性とは要素単語相互の確率的な縛りの強さを意味する. (Sag et al. 2002) では, WordNet1.7 での見出しの約 41% が MWE であること, また (田辺ら 2006) では, 日本語の述部における文末表現を構成する助動詞, 終助詞相当の MWE の出現比率が約 42% であることが報告されており, MWE の適切な処理が自然言語処理の質を向上する上で必要不可欠であることを示唆している. 発話者の主観を表す日本語の述部の文末表現には, 必要性を表す「なければならない」, 欲求を表す「たい」, 様態を表す「ようだ」のようにバラエティに富んでおり, かつ「なければならない/ようだった」のように複数個連続して用いられることも多く, 発話者の主観に基づく処理を行うためには機能語性 MWE を意識した文末表現の取り扱いが必要となる.

本論文では, 第 2 章で非命題的意味構造についての紹介を行い, 第 3 章で機能語性 MWE を組み込んだ体系としての日本語文末表現意味体系について述べる. 第 4 章で非命題的意味構造間の関係に基づく日本語文末表現意味体系の拡張について述べた後, 第 5 章でまとめ, 今後の課題について考察する.

## 2 非命題的意味構造

一般的に文の意味は, 命題的意味と非命題的意味からなっていると考えられ, 日本語においては, 非命題的意味を表す部分は述語に後接した, 時制, 判断, 否定, 話し手の態度など, 広義の様相情報を与える助動詞, 終助詞およびそれらに相当する機能語性 MWE により構成される場合が多い. 筆者らはそれ

らの表現を助述表現と呼び, これまで約 1,450 個の機能語性 MWE と約 50 個の助動詞, 終助詞を収録した辞書を作成している<sup>1</sup>. 機能語性 MWE としては, 例えば, 「て」と「いる」をまとめた「ている」, 「かも」と「しれ」と「ない」をまとめた「かもしれない」などを収録している. 助述表現は, 発話者の主観に基づく表現であるとみなすことができる. 筆者らは網羅性の高さを目標の一つと考えており, 例えば, (森田ら 1989) の収録表現のうち該当表現はすべてカバーしている. 機能語性 MWE を適切に設定することで, 日本語文の構造の大枠は次の生成規則で表すことができる.

$$(1) S_0 \rightarrow BP^* \cdot PRED$$

$$(2) S_i \rightarrow S_{i-1} \cdot e_i \quad (1 \leq i \leq n)$$

$S_0$ ,  $BP$ ,  $^*$ ,  $PRED$ ,  $e_i$  はそれぞれ骨格文, 文節, 閉包演算子, 述語, 助述表現を表す. (2) は,  $S_{i-1}$  と  $e_i$  とで文  $S_i$  が構成されること, 及び, 構造が左分岐であることを表している. 例えば, 日本語文「彼は動き始めていないかもしれない」では, 日本語文の構造の概形および対応する非命題的意味構造は図 1 のようになる<sup>2</sup>.

助述表現に対応する意味を意味関数として捉え, 並びを逆順にすることで, 非命題的意味構造を作ることができる. 一般的な非命題的意味構造は, 次のような入れ子型表現で表すことができる.

$$(3) M_m [M_{m-1} \cdots [M_2 [M_1 [S_0]]] \cdots]$$

但し,  $M_i (1 \leq i \leq m)$  は意味関数である. (1), (2) に示される構文構造と, (3) に示される意味構造とは一種の同型性があると言える. 助述表現が文の述部にいくつも並んだ複雑な文末表現の場合でも, 意味関数との対応をとることにより非命題的意味構造を求めることが可能である. このように (3) は, 構造のシンプルさと同時に対応可能な表現の多様さから工学的に重要な性質を示していると考えられ, 言語依存性も無いとされるため<sup>3</sup>, 言い換えや機械翻訳を行う際の間言語として有効であると考えられる.

<sup>1</sup> これらの数値は, 漢字・かななどの表現のゆれをまとめて 1 見出しとした値である.

<sup>2</sup> 図 1 の記号  $\cdot$  は通常の単語境界を表し,  $\prime$  は MWE による単語境界を表す. また, 動詞に後接する「はじめる」は動詞ととらえる考え方もあるが, アスペクト情報を含むと考えられるため本論文では助述表現としている.

<sup>3</sup> (Cinque 1999; Cinque 2006) は, 非命題的意味構造における意味関数の生起順序には, 言語に依存しない規則性が存在するのではと報告している.

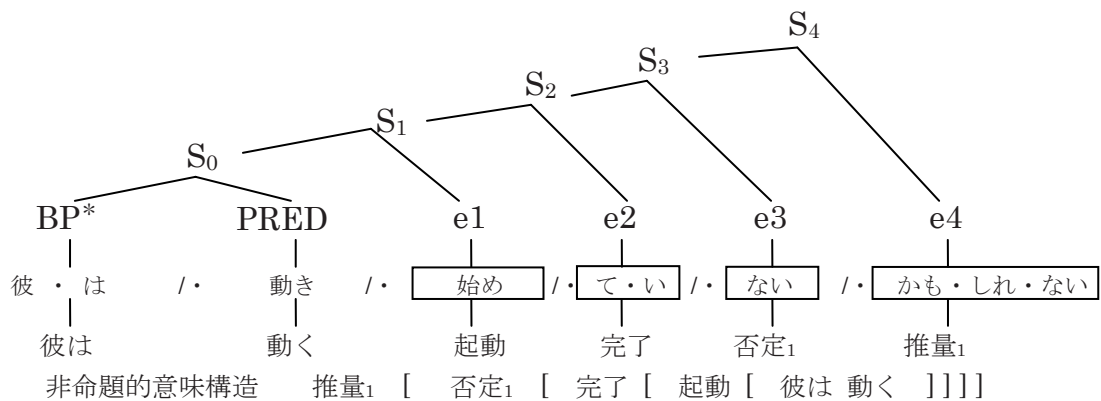


図1 日本語文の構造の概形及び非命題的意味構造

### 3 日本語文末表現意味体系

#### 3.1 意味関数の体系

助述表現に対応した意味関数は、ツリー状に体系化しており、意味関数間の類似度が定義できる。意味関数は

- ・最上位ノードを4種類設定
- ・より細分化された意味を下位ノードとした階層化<sup>4</sup>
- ・葉ノードを139種類設定<sup>5</sup>

して詳細に体系化した。詳細な体系化により、発話者の主観情報の詳細な取り扱いが可能になる。各ノード(意味)には助述表現の集合が対応づけられる。そのため、意味関数間で定義される類似度から、助述表現間の類似度を与えることができる。

#### 3.2 助述表現間の意味的類似度

助述表現、例えば、「なければならない」と「必要がある」の間の類似度の算出について考える。ここで、意味の集合を  $S_M$ 、助述表現の集合を  $S_E$  とおくと、関数  $f$  を、

$$(4) f: S_M \rightarrow 2^{S_E}$$

と形式的に表すことができる。(4)は意味の集合から助述表現の集合のべき集合へ写像する関数となる。ここで、 $S_M$  上の is-a 関係に基づいた意味

$m_1, m_2 \in S_M$  の類似度  $sim(m_1, m_2)$  を次のように定義する。

$$(5) sim(m_1, m_2) = \frac{2 \cdot d_c}{d_i + d_j}$$

ただし、 $d_c$  は、 $m_1, m_2$  の共通上位ノードのルートか

らの距離、 $d_i, d_j$  はそれぞれ、 $m_1, m_2$  のルートからの距離である。また、 $sim(m_1, m_2)$  が  $\alpha$  であるとき、

$$(6) f(m_1) = \{x_1, x_2, \dots, x_m\}$$

$$(7) f(m_2) = \{y_1, y_2, \dots, y_n\}$$

である  $x_i$  と  $y_j$  の間は、 $S_E$  上で類似度  $\alpha$  の関係にあると定義する。ここで、 $\alpha = 1$  の場合、 $x_i$  と  $y_j$  の間は同義関係にあると定義する。これらの定義は、任意の助述表現間の意味的類似度を算出することができること、および、同義関係にある助述表現間は基本的に言い換え可能であることを意味している。

#### 3.3 助述表現列間の意味的類似度

実際の日本語文末では、「なければならないようだった」などのように、助述表現が連続して助述表現列として出現する場合も多い。前節では助述表現間の関係として類似度を定義したが、助述表現列間でも類似度を定義できれば望ましい。そこで本節では、助述表現間から助述表現列間の類似度の拡張を考える。ここで、 $S_M$  の閉包を  $S_M^*$  (意味列の集合)、 $S_E$  の閉包を  $S_E^*$  (助述表現列の集合) とおく。ここで、関数  $g$  を、

$$(8) g: S_M^* \rightarrow 2^{S_E^*}$$

と形式的に表すことができる。(8)は意味列の集合から助述表現列の集合のべき集合へ写像する関数となる。(但し、 $S_M \subset S_M^*, S_E \subset S_E^*$  である。)(8)

は、 $m \in S_M, \alpha \in S_M^*$  において

$$(9) g(m\alpha) = f(m) \cdot g(\alpha)$$

$$(10) g(\lambda) = \lambda$$

と定義できる。但し、 $\lambda$  は空文字列、 $\cdot$  は結合演算を表す。ここで、意味列  $M_1, M_2 \in S_M^*$  の類似度  $sim(M_1, M_2)$  を定義する。意味列  $M_1, M_2$  間の類似度は、DP (Dynamic Programming) を用いた意味列間の最短距離計算 (DP マッチング) に基づき算出することができる。ここで、 $d(i, j)$  は意味列

<sup>4</sup> 例えば意味「必要性」を、「必要性<sub>1</sub>」「必要性<sub>2</sub>」・・・「必要性<sub>8</sub>」と細分化し、細分化されたそれぞれのノードを「必要性」の下位ノードとして設定している。

<sup>5</sup> 図1における意味関数は全て葉ノードを用いている。

```

begin
   $d(0,0) := 0;$ 
  for  $i := 1$  to  $m$  do
     $d(i,0) := d(i-1,0) + w(M_1,i);$ 
  for  $j := 1$  to  $n$  do
     $d(0,j) := d(0,j-1) + w(M_2,j);$ 
  for  $i := 1$  to  $m$  do
    for  $j := 1$  to  $n$  do
       $d(i,j) := \min \{d(i-1,j-1) + 2 * (1 - \text{sim}(a_i, b_j)), d(i-1,j) + w(M_1,i), d(i,j-1) + w(M_2,j)\};$ 
 $\text{sim}(M_1, M_2) := 1 - d(m, n) / \left( \sum_{k=1}^m w(M_1, k) + \sum_{l=1}^n w(M_2, l) \right);$ 
end

```

$w(M_1, i)$ :意味列  $M_1$  における意味  $a_i$  の重要度,  $0 \leq w(A, i) \leq 1$

$w(M_2, j)$ :意味列  $M_2$  における意味  $b_j$  の重要度,  $0 \leq w(B, j) \leq 1$

図2 意味列間類似度判定処理

$M_1 = a_1, a_2, \dots, a_i$  と意味列  $M_2 = b_1, b_2, \dots, b_j$  の意味列間距離を表すものとする. このプロセスは (i) 意味列  $a_1, a_2, \dots, a_{i-1}$  と  $b_1, b_2, \dots, b_{j-1}$  の最短距離  $d(i-1, j-1)$  に  $a_i$  と  $b_j$  の置き換えによる距離の増加分を加算したもの, (ii) 意味列  $a_1, a_2, \dots, a_{i-1}$  と  $b_1, b_2, \dots, b_j$  の最短距離  $d(i-1, j)$  に意味の重要度  $w(M_1, i)$  を加算したもの, および (iii) 意味列  $a_1, a_2, \dots, a_i$  と  $b_1, b_2, \dots, b_{j-1}$  の最短距離  $d(i, j-1)$  に意味の重要度  $w(M_2, j)$  を加算したものの, の3個の数値のうち最小値を意味列  $a_1, a_2, \dots, a_i$  と  $b_1, b_2, \dots, b_j$  の最短距離  $d(i, j)$  として求める部分である. 重要度  $w(M_1, i)$  は, 意味列間距離を算出する場合の寄与の割合であり, 例えば, 全ての意味に対して均等にする, 特定の意味を重視する, また, 最後尾の助述表現の意味のみを考慮するなど, アプリケーションに応じた柔軟な類似度算出が可能となる.  $i, j$  の値を, それぞれ1から  $m, 1$  から  $n$  まで1ずつ増やしながら処理を繰り返し実行することによって, 最終的に意味列  $M_1 = a_1, a_2, \dots, a_m$  と  $M_2 = b_1, b_2, \dots, b_n$  の最短距離として意味列間距離  $d(m, n)$  が求められる. 最短距離計算の詳細は(安武ら 1999)を参照されたい. 意味列間距離  $d(m, n)$  が求まると, 意味列間類似度  $\text{sim}(M_1, M_2)$  が求まる.

同様に,  $\text{sim}(M_1, M_2)$  が  $\alpha$  であるとき,

$$(8) \quad g(M_1) = \{X_1, X_2, \dots, X_m\}$$

$$(9) \quad g(M_2) = \{Y_1, Y_2, \dots, Y_n\}$$

である  $X_i$  と  $Y_j$  の間は,  $S_E^*$  上で類似度  $\alpha$  の関係に

あると定義する. 同様に  $\alpha = 1$  の場合,  $X_i$  と  $Y_j$  の間は同義関係にあると定義する. これらの定義は前節と同様, 任意の助述表現列間の意味的類似度を算出することができること, 及び, 同義関係にある助述表現列間は基本的に言い換えが可能であることを意味している.  $S_E^*$  による体系は意味的類似度を定義できる体系であり, この点でソーラスなどの意味体系と同一であるとみなすことができる.  $S_E^*$  の元は助述表現列であり, 助述表現列は日本語文末に存在しうることから, 本論文では  $S_E^*$  による体系を日本語文末表現意味体系と呼ぶことにする. 日本語文末表現意味体系は1500種の助述表現が複数個接続した表現による意味体系であり,  $n$  個の助述表現が接続した場合, 最大  $1500^n$  種の表現を網羅する大規模な意味体系であると考えられる<sup>6</sup>.

## 4 日本語文末表現意味体系の拡張

前章までで任意の助述表現列間に, 意味的類似度を定義できることを示した. しかし, 例えば「行かなければならないことはない」と「行かなくてもよい」のような助述表現列間の意味的類似度は文脈等に応じて高く見積もることが出来れば望ましい. そこで第4章では, 非命題的意味構造間に定義できる関係と, 日本語文末表現意味体系の拡張について考察する.

### 4.1 非命題的意味構造間の関係

#### 4.1.1 類似関係

非命題的意味構造間の関係の中でも特に類似関係として, 類似性規則が提案されている (田辺ら

<sup>6</sup> 助述表現間の接続には制約があり, 存在しうる助述表現列の種類は現在行っている.

2001).類似性規則は,論理的規則および語用論的規則に大別される.論理的規則には

(13) 否定<sub>1</sub>[否定<sub>1</sub>[S]] ⇐ [S]

「行かない/こと・は・ない」⇐「行く」

(14) 否定<sub>1</sub>[必要性<sub>1</sub>[S]] ⇐ 許容<sub>2</sub>[否定<sub>1</sub>[S]]

「食べ/な・けれ・ば・ならない/こと・は・ない」  
⇐「食べ/なく/て・も・よい」

など,語用論的規則には

(15) 願望[受動態<sub>3</sub>[S]] ⇐ 依頼<sub>1</sub>[S]

「見て/いただきたい」  
⇐「見て/ください」

(16) 疑問<sub>1</sub>[否定<sub>1</sub>[可能性<sub>2</sub>[S]]] ⇐ 依頼<sub>2</sub>[S]

「行く/こと・が・出来/ない/か」  
⇐「行って/ください」

などがある.

#### 4.1.2 類似以外の関係

非命題的意味構造間の関係として,類似以外の関係を挙げることもできる.次のような例を考える.

(17) 過去[否定[必要性<sub>3</sub>[XがYする]]]

→過去[XがYする]

「XがYする/べき・で・なかった」

→「XがYした」

(18) 願望[XがYする]

→否定[完了[XがYする]]

「XがYしたい」

→「XがYして/いない」

(17)は,「XがYするべきでなかった」ならば「XがYした」,(18)は,「XがYしたい」ならば「XがYしていない」はず,を表し,話し手の発話を受けた聞き手が推論しうる関係であると考えられる.

#### 4.2 日本語文末表現意味体系の拡張

非命題的意味構造間の関係を規則として適用することで,日本語文末表現意味体系の拡張を考える.非命題的意味構造間の関係は  $S_M^*$  から  $S_M^*$  への写像

(19)  $h: S_M^* \rightarrow S_M^*$

として形式的に表すことができる.ここで,(5)より,合成関数  $g \circ h$  は

(20)  $g \circ h: S_M^* \rightarrow 2^{S_E^*}$

となる.そのため,  $M_1 \in S_M^*$  に対して規則を適用した  $h(M_1)$  と,  $M_2 \in S_M^*$  の間での意味的類似度  $sim(h(M_1), M_2)$  を定義できる.ここで,  $sim(h(M_1), M_2) > sim(M_1, M_2)$  であり,  $sim(h(M_1), M_2)$  が  $\alpha$  であるとき,

(21)  $g(M_1) = \{X_1, X_2, \dots, X_m\}$

(22)  $g(M_2) = \{Y_1, Y_2, \dots, Y_n\}$

である  $X_i$  と  $Y_j$  の間は,  $S_E^*$  上で類似度  $\alpha$  の関係に

あると定義できる<sup>7</sup>.つまり,(19)を規則として適用することで助述表現列間の意味的類似度が高くなった場合,文脈などに応じて高いほうを優先することができる.このように,非命題的意味構造間の関係を規則として適用することで,日本語文末表現意味体系の拡張ができると考えている.

#### 5 おわりに

自然語文の非命題的意味は,対話理解,文脈モデルや話者の態度の推定などの自然言語処理で重要な役割を果たす.本論文では,日本語文において発話者の主観を表すと考えられる助述表現列について整理した日本語文末表現意味体系について報告した.日本語文末表現意味体系は大規模であり,任意の助述表現列間に対し意味的類似度を算出できる.また,日本語文末表現意味体系と,自立語シソーラスをうまく組み合わせることにより,例えば「教授が述べるに違いないだろう」と「先生が話すかもしれないらしい」の間の意味的類似度の算出も可能になり,文書分類や用例に基づく機械翻訳などの応用技術への貢献が期待できる.また,本論文では,助述表現列による非命題的意味について述べたが,「たぶん」「きっと」などのような副詞も非命題的意味を持ちうる.そのため,副詞をどのように非命題的意味構造に組み込むかも検討する必要がある.今後の課題として,非命題的意味構造間の関係の収集や,実際のテキストデータへの適用などが挙げられる.

#### 参考文献

- [1] Ivan A. Sag, Timothy Baldwin, Francis Bond, Ann Copestake and Dan Flickinger. “Multiword Expressions: A Pain in the Neck for NLP.” The Proc. of the 3rd CICLING, pp. 1-15. 2002.
- [2] 田辺利文, 本田聖晃, 高橋雅仁, 小山泰男, 吉村賢治, 首藤公昭. “日本語文末表現の取り扱いについて.” FIT2006, pp. 241-244. 2006.
- [3] 森田良行, 松木正恵. “日本語表現文型用例中心・複合辞の意味と用法.” アルク. 1989.
- [4] Guglielmo Cinque. “Adverbs and Functional Heads.” OXFORD UNIVERSITY PRESS. 1999.
- [5] Guglielmo Cinque. “Restructuring and Functional Heads.” OXFORD UNIVERSITY PRESS. 2006.
- [6] 安武満佐子, 小山泰男, 吉村賢治, 首藤公昭. “関係表現,助述表現の類似度を考慮した言語表現間の意味的類似度判定.” 福岡大学工学集報 第63号, pp.171-177. 9月. 1999.
- [7] 田辺利文, 吉村賢治, 首藤公昭. “日本語モダリティ表現とその言い換え.” 言語処理学会第7回年次大会ワークショップ論文集, pp. 47-50. 2001.

<sup>7</sup> ここでの  $h$  としては,類似関係のみ(類似性規則)を対象としており,類似以外の関係については今後の課題としている.