

階層的機能表現辞書に基づく新聞記事中の機能表現の調査・分析

長坂泰治，坂本明子，宇津呂武仁，森下洋平

(筑波大学大学院システム情報工学研究科)

松吉俊（奈良先端科学技術大学院大学情報科学研究科）

土屋雅稔（豊橋技術科学大学情報メディア基盤センター）

Analysis of Japanese Functional Expressions in News Corpus based on a Hierarchical Lexicon

Taiji Nagasaka, Akiko Sakamoto, Takehito Utsuro, Yohei Morishita

(University of Tsukuba)

Suguru Matsuyoshi (Nara Institute of Science and Technology)

Masatoshi Tsuchiya (Toyohashi University of Technology)

概要

本研究では、松吉による大規模日本語機能表現辞書に収録されている、約 17,000 種類の機能表現を対象として、日本語文中の機能表現を検出することを目的としている。本研究では、大規模日本語機能表現辞書を利用した集約的検出方式を提案している。この方式では、新聞記事に 50 回以上出現する機能表現のうち、代表的な表現 300 を介して、それらに対応する 5,000 の派生的な表現の検出が可能となる。また、これまでの数百表現規模の調査の結果、およそ 1/3 の機能表現は、機能的用法と内容的用法が適度な割合で存在し、残りの 2/3 については、機能的用法に偏っていることが分かっている。本研究では、新聞記事 1 年分（約 100 万文）に出現する機能表現に対して、人手で用法判定することで、代表的な表現 300 の内訳 100 と 200 の区分を明らかにする。

1 はじめに

機能表現¹とは、「にあたって」や「をめぐって」のように、2つ以上の語から構成され、全体として1つの機能的な意味をもつ表現である。一方、この機能表現に対して、それと同一表記をとり、内容的な意味をもつ表現が存在することがある。例えば、文(1)と文(2)には「にあたって」という表記の表現が共通して現れている。

- (1) 出発する にあたって、荷物をチェックした。
- (2) ボールは壁 にあたって、跳ね返った。

文(1)では、下線部はひとかたまりとなって「機会が来たのに当面して」という機能的な意味で用いられている。それに対して、文(2)では、下線部に含まれている動詞「あたる」は、動詞「あたる」本来の内容的な意味で用いられている。このような表現においては、機能的な意味で用いられている場合と、内容的な意味で用いられている場合とを識別する必要がある。

我々はこれまでに、現代語複合辞用例集 [国研 01](以下、用例集) 中の代表的複合辞一覧に基づいて、それらの派生形である 337 種類の機能表現を規定し、その用例データベース(日本語複合辞用例データベース [土屋 06, 土屋 07a], 以下、用例データベース)を作成した。また、それらの用例データベースを訓練事例として、機械学習により機能表現の検出・係り受け解析を行う方式を提案した [土

¹機能表現は、複数形態素からなる複合辞と一つの形態素からなる機能語から構成されるが、本稿では、複合辞と同等の意味で機能表現という用語を用いる。

表 1: 機能表現辞書の 9 つの階層

階層		分類数	表現数		
			合計 (L^9 表現数)	助動詞 型以外	助動詞型
L^1	見出し語	—	341 (488)	281	207
L^2	意味	88	435 (488)	281	207
L^3	派生 (格助詞型, 接続助詞型, 連体助詞型, 接続詞型, 助動詞型, 形式名詞型,とりたて詞型, 提題助詞型)	8	555	348	207
L^4	機能語の交替	—	774	492	282
L^5	音韻的変化	38	1,187	633	554
L^6	とりたて詞の挿入	18	1,810	659	1151
L^7	活用	—	6,870	659	6211
L^8	「です/ます」の有無	2	9,722	895	8827
L^9	表記のゆれ	—	16,801	1360	15411

屋 07b, 注連 07] . また , 機能表現の異形の語構成パターンを網羅することにより , (日本語機能表現一覧 [松吉 07] , 以下 , 機能表現一覧²⁾)を作成した .

ここで , [土屋 07b, 注連 07] の機械学習による機能表現検出においては , 一つの表現あたり 50 例程度の訓練用例に対して , 人手で機能的・自立的等の用法判定を行う必要がある . しかし , 機能表現一覧の全機能表現 16,801 種類に対して , それだけの規模の作業を行うことは容易ではない . そこで , [長坂 08] では , 機能表現一覧の階層性を利用し , 階層において下位に位置する機能表現 (以下 , 派生的表現) について , 用法が類似するより上位の表現 (以下 , 代表的表現) に言い換えた後 , 用法判定を行う方式を提案した .

一方 , [長坂 09] では , [長坂 08] の提案をふまえて , 機能表現一覧中の情報のうち , 特に文体の情報に注目し , 代表的表現および派生的表現の区別を整理した . さらに , 每日新聞 1995 年分のテキストデータ中において , 機能表現一覧の機能表現の出現頻度調査を行い , [土屋 07b, 注連 07] の機械学習による機能表現検出において必要となる訓練事例 (出現頻度 50 以上) が存在する機能表現の規模を推定した .

本稿では , [長坂 09] をふまえて , 每日新聞 1995 年分のテキストデータ中に 50 回以上出現する代表的表現の表記を対象として , 人手で機能的用法・自立的用法の判定作業を行った . さらに , 各機能表現表記に対して , 機能的用法・自立的用法の両方が適度な割合で混合して出現し , 機械学習によつて機能表現検出を行う必要のある機能表現表記の割合を求めた結果について報告する .

2 階層的機能表現辞書

機能表現一覧 [松吉 07] は , 9 つの階層構造をなしており , 各階層は , 表 1 に示されるような観点によって分類されている . 同表に , 各階層における機能表現数が示されており , 図 1 に階層構造の一部をそれぞれ示す .

また , 機能表現の文体に着目し , 文体ごとの機能表現の振る舞いについて述べる . 文体とは , 機能表現一覧中の表現に付与されている情報であり , 常体 , 堅い文体 , 口語体 , 敬体の 4 種類がある . 表 2 にそれぞれの文体における表現例を示す .

²<http://kotoba.nuee.nagoya-u.ac.jp/tsutsuji/>

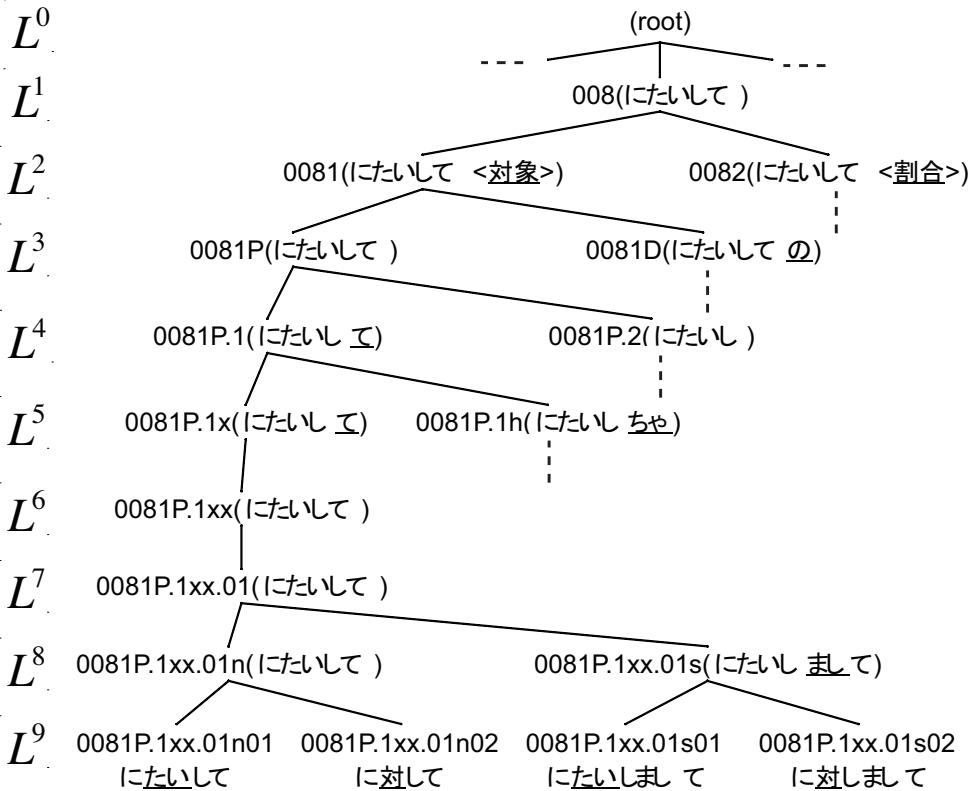


図 1: 機能表現辞書階層構造の一部

表 2: 文体の種類

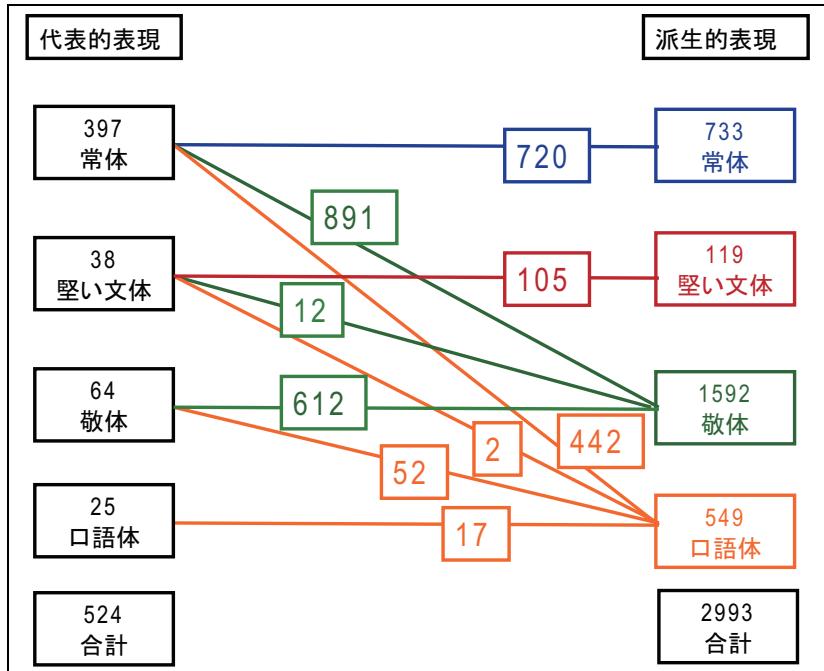
文体	表現例
常体	について
堅い文体	につき
口語体	についちゃ
敬語体	につきまして

3 代表的表現への集約

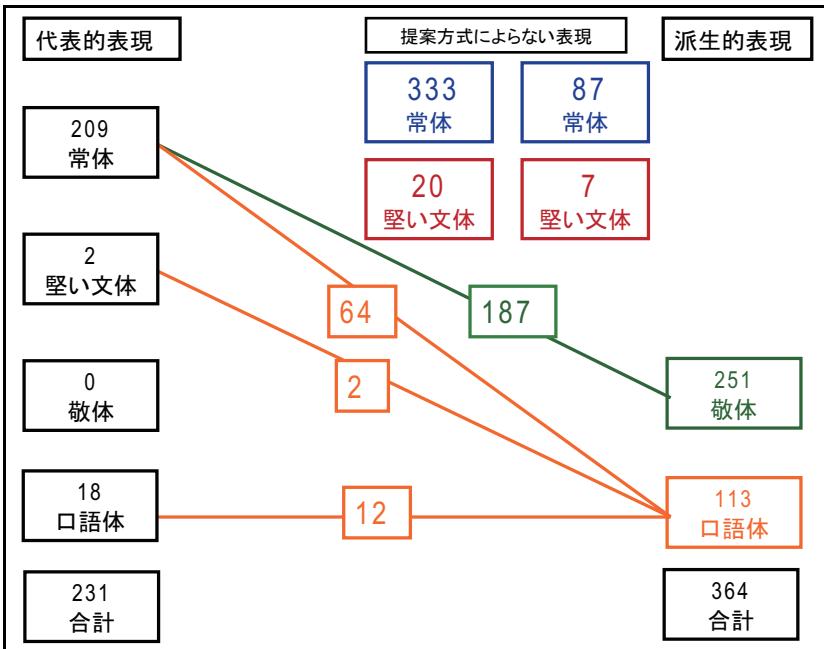
3.1 基本的な考え方

[長坂 08] で提案した代表的表現への集約方式においては、階層の上位に位置する代表的表現は、 L^4 階層相当の 1,000 表現程度の規模とする。そして、機能表現一覧において、代表的表現を除く表現はすべて、言い換えの対象の表現となる。本研究では、これらの表現を派生的表現と定義する。派生的表現を代表的表現に言い換える際には、以下の制約を課す。

- 機能表現の語頭の無声・有声の制約により前接する活用語の活用型が制限される場合は、この制限を保持する。
- 機能表現の仮名表記・漢字表記の違いを保持する。
- 助動詞型の機能表現の場合には、言い換え前後で活用形を保持する。



(a) (助動詞型基本形)



(b) (助動詞型以外)

図 2: 代表的表現への集約における文体ごとの機能表現数

3.2 文体ごとの機能表現数

[長坂 08] では，助動詞型の場合，派生的表現に 3,000 表現を，代表的表現 500 表現に集約した。一方，助動詞型以外の場合，文体を常体および堅い文体で見てみると，代表的表現が 353 表現であるのに対して，派生的表現が 94 表現と少ないため，集約方式の効果はほとんど得られない。そこで，本論文では，助動詞型以外の機能表現のうち，常体および堅い文体の表現は，敬体や口語体の派生的表現に対する代表的表現を除いて，集約方式の対象表現から外す。

表 3: 毎日新聞 1995 年において、50 回以上出現する機能表現数の分布

	助動詞型 (基本形)		助動詞型以外			合計
	代表的表現	派生的表現	代表的表現	派生的表現	その他	
常体	164	38	87	0	178	467
堅い文体	8	3	0	0	9	20
敬体	14	42	0	1	0	57
口語体	7	37	1	13	0	58
合計	193	120	88	14	187	602

表 4: 複数の ID を持つ機能表現表記数 (機能表現一覧全体 / 每日新聞 1995 年分で頻度 50 以上)

	助動詞型 (基本形)		助動詞型以外			合計
	代表的表現	派生的表現	代表的表現	派生的表現	その他	
常体	57/44	13/5	27/19	0/0	136/73	233/141
堅い文体	6/3	14/0	0/0	0/0	12/5	32/8
敬体	2/0	72/8	0/0	64/1	0/0	136/9
口語体	7/4	33/13	3/1	35/9	0/0	78/27
合計	72/51	132/26	30/20	99/10	148/78	479/185

以上をふまえて、各体ごとに、代表的表現、派生的表現の数を示したものを、図 2³ に示す。両図の節点の数字は、機能表現の表記数を表している。また、枝の数字は、端点における体の間で言い換え組を形成する表現の数である。

4 新聞記事における機能表現の分布

4.1 機能表現表記数の分布

毎日新聞 1995 年の一年分の中で、50 回以上出現する機能表現表記数の調査を行った結果を表 3 に示す。機能表現一覧の中には、同一の表記に対して複数の ID が存在する表現があることをふまえて、表 3 には、機能表現表記単位で集計した結果を示す。また、表 4 に、複数の ID を持つ機能表現の表記数の一覧を示す。

表 5: 每日新聞 1995 年において、50 回以上出現する代表的表現と派生的表現の組数の分布および組の例 (助動詞型基本形)

文体の組み合わせ	組数	組の例
常体の代表的表現-常体の代表的表現	25	てよい-てもよい
常体の代表的表現-敬体の代表的表現	30	なければならない-なければならないです
常体の代表的表現-口語体の代表的表現	20	ても仕方がない-ても仕方ない
合計	75	

³助動詞型では、活用形を基本形に限定した数値を算出している。

表 6: 機能的用法・自立的用法の分布

	助動詞型（基本形）		助動詞型以外			合計
	代表的表現	派生的表現	代表的表現	派生的表現	その他	
常体	66.2/28.1/5.7	0/100/0	57.1/32.5/10.4	—	53.0/39.1/7.9	58.5/33.9/7.6
堅い文体	83.3/16.7/0	—	—	—	14.3/57.1/28.6	46.2/38.5/15.3
敬体	100/0/0	85.7/0/14.3	—	—	—	94.4/0/5.6
口語体	20.0/60.0/20.0	33.4/33.3/33.3	0/0/100	0/100/0	—	23.5/47.1/29.4
合計	67.7/26.7/5.6	50.0/27.8/22.2	56.4/32.1/11.5	0/100/0	51.3/39.9/8.8	58.3/33.1/8.6

x : 機能的用法の割合
 $90\% \leq x \leq 100\%$ となる表現の割合 (%) / $10\% < x < 90\%$ となる表現の割合 (%)
/ $0\% \leq x \leq 10\%$ となる表現の割合 (%)

表 3 より、全 602 表現のうち、常体の機能表現が 467 表現と大きい割合を占めていることが分かる。表 3 の中で、代表的表現と派生的表現の関係にある組の数および例を表 5 に示す。

4.2 機能的用法・自立的用法の分布

これまでに、[土屋 06, 土屋 07b, 注連 07] で述べたように、全機能表現表記のうち、特に、機能的用法・自立的用法の両方が適度な割合で混合して出現する機能表現表記に対してのみ、機械学習によって機能表現検出を行う必要がある。また、[土屋 06] で報告したように、[土屋 06] の用例データベースの範囲においては、毎日新聞 1995 年の一年分の中で、50 回以上出現する機能表現表記 187 表記のうち、機能的用法・自立的用法の両方が適度な割合で混合して出現する機能表現表記の割合は約 3 分の 1 程度であった。

一方、本節では、表 3 中の代表的表現、派生的表現、および「その他」の表現（助動詞型以外の場合のみ）を対象として、機能的用法・自立的用法の用法判定作業を行った結果を表 6 に示す⁴⁵。この結果から、機能的用法・自立的用法の両方が適度な割合で混合して出現する（機能的用法の割合 x が、 $10\% < x < 90\%$ となる）機能表現表記の割合は、代表的表現、派生的表現および「その他」の全体では約 3 分の 1 程度であることがわかる。しかし、「その他」のみでは、約 40% と多くなっている。

一方、表 7 には、毎日新聞 1995 年分を用いて代表的表現の用法判定のための検出器の教師あり学習を行った場合に、集約方式によって検出対象となる派生的表現のうち、毎日新聞 1995 年分に 50 回以上出現しない表現の数を示す。これらの合計 1,213 表現は、毎日新聞 1995 年分に 50 回以上出現しないため、毎日新聞 1995 年分の範囲では、用法判定のための検出器の教師あり学習ができない。しかし、代表的表現への集約方式を用いることにより、用法判定が可能となる。

5 関連研究

[松吉 08]においては、機能表現一覧 [松吉 07] 中の機能表現を対象として、意味を保存する言い換えが可能な機能表現の分類を規定している。一方、本論文では、機能表現の用法判定の性能を保ったまま、代表的表現への言い換えを行うという、より緩い制約のもとでの機能表現の言い換えが目的で

⁴ 本稿の執筆時点では、表 6 には、表 3 中の 602 表現のうち、417 表現のみを対象として用法判定作業を行った結果を示す。残りの 185 表現は、417 表現のいずれかの表現との間で文中に出現する文字位置が重複することが多く、用法判定作業箇所の集計の際に集計漏れとなっている。これらの表現に対する集計の復元は容易に行うことができる。

⁵ [土屋 06]において、用法判定の際に用いられた 6 種類の用法判定ラベルのうち、他の機能表現や慣用表現の表記と交差する位置に、別の機能表現表記の文字列が重複する場合の判定ラベルとして「判定ラベル B」を用いている。今回の用法判定作業箇所のうち、「判定ラベル B」となった箇所は全 56033 篇所中、合計 3505 篇所であったが、そのうち 1309 篇所については、前後の形態素の表記もしくは品詞を参照すれば一意に「判定ラベル B」と判定可能であった。そこで、これらの箇所については、今後、機械的に用例文収集結果から除外することとし、その結果不足した用例文については、再収集および用法の再判定を行うこととする。なお、本稿の執筆時点では、上述の 1309 篇所を除外したところまで作業が進んでいるため、表 6 にはその段階における集計結果を載せている。

表 7: 每日新聞 1995 年に 50 回以上出現する代表的表現に対する派生的表現のうち，毎日新聞 1995 年に 50 回以上出現するものを除いた表現の数

文体	助動詞型（基本形）	助動詞型以外
常体	347	8
堅い文体	52	2
敬体	537	78
口語体	164	25
合計	1100	113

ある。また、代表的表現への言い換えを介した機械翻訳の研究としては、内容語と口語的な機能表現を扱った [山本 01, 山本 02]、機能表現一覧 [松吉 07] の機能表現を対象とした [坂本 09] がある。

6 おわりに

本稿では、[長坂 09] をふまえて、毎日新聞 1995 年分のテキストデータ中に 50 回以上出現する代表的表現の表記を対象として、人手で機能的用法・自立的用法の判定作業を行った。さらに、各機能表現表記に対して、機能的用法・自立的用法の両方が適度な割合で混合して出現し、機械学習によって機能表現検出を行う必要のある機能表現表記の割合を求めた結果について報告した。今後は、機械学習により用法判定を行うための訓練・評価データを整備し、提案方式の実装および評価を行う。

参考文献

- [国研 01] 国立国語研究所：現代語複合辞用例集（2001）。
- [松吉 07] 松吉俊、佐藤理史、宇津呂武仁：日本語機能表現辞書の編纂、自然言語処理、Vol. 14, No. 5, pp. 123–146 (2007)。
- [松吉 08] 松吉俊、佐藤理史：文体と難易度を制御可能な日本語機能表現の言い換え、自然言語処理、Vol. 15, No. 2, pp. 75–99 (2008)。
- [長坂 08] 長坂泰治、宇津呂武仁、土屋雅穂：大規模日本語機能表現辞書の階層性を利用した機能表現検出、言語処理学会第 14 回年次大会論文集, pp. 837–840 (2008)。
- [長坂 09] 長坂泰治、宇津呂武仁、松吉俊、土屋雅穂：大規模階層辞書を利用した日本語機能表現の集約と解析、言語処理学会第 15 回年次大会論文集, pp. 328–331 (2009)。
- [坂本 09] 坂本明子、宇津呂武仁、松吉俊：日本語機能表現の集約的英訳、言語処理学会第 15 回年次大会論文集, pp. 654–657 (2009)。
- [注連 07] 注連隆夫、土屋雅穂、松吉俊、宇津呂武仁、佐藤理史：日本語機能表現の自動検出と統計的係り受け解析への応用、自然言語処理、Vol. 14, No. 5, pp. 167–197 (2007)。
- [土屋 06] 土屋雅穂、宇津呂武仁、松吉俊、佐藤理史、中川聖一：日本語複合辞用例データベースの作成と分析、情報処理学会論文誌、Vol. 47, No. 6, pp. 1728–1741 (2006)。

- [土屋 07a] 土屋雅穂, 注連隆夫, 松吉俊, 宇津呂武仁, 佐藤理史, 中川聖一：機能表現を考慮した日本語
係り受け解析器学習のためのコーパス作成, 言語処理学会第 13 回年次大会論文集, pp. 510–513
(2007).
- [土屋 07b] 土屋雅穂, 注連隆夫, 高木俊宏, 内元清貴, 松吉俊, 宇津呂武仁, 佐藤理史, 中川聖一：機械
学習を用いた日本語機能表現のチャンキング, 自然言語処理, Vol. 14, No. 1, pp. 111–138 (2007).
- [山本 01] 山本和英, 白井諭, 坂本仁, 張玉潔：SANDGLASS: 両言語換言機構を基軸とする音声翻訳,
言語処理学会第 7 回年次大会発表論文集, pp. 221–224 (2001).
- [山本 02] 山本和英：換言と言語変換の協調による機械翻訳モデル, 言語処理学会第 8 回年次大会発
表論文集, pp. 307–310 (2002).