

後続要素を予告する表現の分析

木田敦子^{1,2} 山本 英子² 井佐原 均^{1,2}

1 通信・放送機構

2 通信総合研究所

{kida, eiko, isahara}@crl.go.jp

述語が文末に置かれる日本語の場合、文の終末まで行かないと文内容が確定しない。しかし、内容が複雑になり文が長くなると、早めに内容が否定なのか肯定なのか疑問なのかなどを知りたくなる。中世以前の日本語には、係助詞と文末の活用形とが形態的な呼応関係を持つ係り結びの用法があった。係り結びが消滅した現代語では、ある種の副詞などが古語の係助詞と似た役割を果たしており、後続要素を予告しているとの指摘がある。本稿では、補完類似度を用いて、大容量のテキストデータから半自動的に呼応関係を発見する手続きについて報告する。

Analysis of expression which projects the following elements beforehand

Atsuko Kida^{1,2} Eiko Yamamoto² Hitoshi Isahara^{1,2}

1 Telecommunication Advancement Organization

2 Communications Research Laboratory

In Japanese, a predicate appears in the sentence end. So, the content of the sentence is not fixed until the end of the sentence. However, when the contents become complicated and a sentence becomes long, listener wants to know the sentence type beforehand. Medieval times or before, there was usage of a KAKARI-MUSUBI in Japanese. KAKARI-MUSUBI projects the following elements beforehand. Now, there is no KAKARI-MUSUBI in Japanese. On the other hand, a certain kind of adverb is used to projects the following elements beforehand instead of KAKARI-JOSI. In this paper, we use CSM(Complementary Similarity Measure) to discover the pair of word which appears to KO-OU from large-scale text data half-automatically.

して形態的な拘束関係を持つ係り結びは衰退していった。言語類型上、文末に述語が現れる日本語の場合、文の終末まで行かないと文内容が確定しない。しかし、内容が複雑になりセンテンスが長くなると、早めに文内容が否定なのか肯定なのか疑問なのかなどを知りたくなる。[大野(1993)]は、「時間的に線状的に発展し連続していく言語表現の早い部分で、一文の行く手、肯定か否定か疑問かなどを予告しておこうとする」表現として、現代語には古語の係助詞に代わって「ある種の副詞」が存在することを示唆している。

3 主な論点と方法

本稿では、「時間的に線状的に発展し連続していく言語表現の早い部分で、一文の行く手、肯定か否定か疑問かなどを予告しておこうとする表現(=後続要素を予告する表現)」、つまり呼応関係が存在することを前提とする。その上で、呼応関係には具体的にどのようなものがあるのかを明らかにしていく。

そのための方法として、補完類似度[山本・梅村(2002)]を用い、大容量のテキストデータから半自動的に現代語における呼応関係を発見する。本稿では、この手続きについて報告する。

本稿の調査対象語は、『基礎日本語文法』で「提題助詞」「取り立て助詞」^{注1}「陳述の副詞」に分類されている以下の語とする。これらの語を呼応関係の「呼」要素と仮定して、調査を進める。調査データは、毎日新聞記事データ 10 年分(1991 年～2000 年)、日経新聞記事データ 11 年分(1990 年～2000 年)を使用する(表 1)。

[調査対象語]

「こそ」「しか」「さえ」「は」「も」「ばかり」「のみ」「すら」「なら」「くらい(ぐらい)」「だけ」「なんて」「決して」「おそらく(恐らく)」「たぶん(多分)」「ぜひ(是非)」「まるで」「もし」「きっと」

注1 「提題助詞」「取り立て助詞」を「係助詞」「副助詞」に分類する立場もある[山田(1936)]。また、一括して「取り立て助詞」と扱う場合[沼田(1986)]、「副助詞」として扱う場合(JUMAN、茶筌の辞書)もある。

4 共起ペアの選定

本節では、補完類似度を用いて調査対象語と共起しやすい語を選定し、共起ペアを得る手順について報告する。

- 1) 毎日新聞記事データ 10 年分(1991 年～2000 年)、日経新聞記事データ 11 年分(1990 年～2000 年)を一文ごとに区切る
- 2) 一文ごとに区切ったデータを JUMAN で形態素解析する
- 3) 特定の品詞の語を除外する
[除外する品詞 (JUMAN の品詞体系)]
 - 「未定義語」「特殊」
 - 「名詞」のうち「普通名詞」「固有名詞」「人名」「地名」「サ変名詞」
- 4) 新聞記事データ 1 年分ごとに補完類似度の計算を行う(表 2)
- 5) 各調査対象語ごとに 4 の計算結果中から当該語を含む行を抽出したデータを作成
- 6) 5 のデータを全年分結合し、類似度をキーにしてソートを行う(表 3)
- 7) 6 の上位 200 行から 15～20 語、多出している語を目視で選び共起語リストを作成する。原則として、二回出現した段階でリストに加える
- 8) 共起語を絞り込むために、7 のリストから以下のものを除外する

格助詞 (と, が, を, に, で)

格助詞は格関係を示す語。格関係を広義呼応関係と見ることできるが、その場合でも格助詞は「呼」要素にあたる。今回探しているのは「応」要素なので格助詞^{注2}は除外する。

副助詞 (は, も)

今回調査対象としている語は呼応関係の「呼」要素と仮定しているため、調査対象の語は除外する。

注2 除外対象としたのは必須格的な格助詞。周辺の格助詞は除外対象としていない。したがって、「まで(格助詞)」は除外していない。

[表 2 : 補完類似度の計算結果]

1528754.872677	の(接続助詞)	を(格助詞)
1064271.001864	に(格助詞)	を(格助詞)
730148.917629	する(動詞)	を(格助詞)
577524.855078	と(格助詞)	を(格助詞)
464473.537901	の(接続助詞)	は(副助詞)
385390.892772	と(格助詞)	いう(動詞)
375225.425650	して(動詞)	を(格助詞)
373562.222429	した(動詞)	を(格助詞)
372757.242754	と(格助詞)	は(副助詞)
371617.004656	で(格助詞)	を(格助詞)
360516.611664	が(格助詞)	に(格助詞)
335289.793261	など(副助詞)	を(格助詞)
323487.885893	と(格助詞)	に(格助詞)
318724.252488	の(接続助詞)	や(接続助詞)
318501.865737	に(格助詞)	よる(動詞)
315778.271689	の(接続助詞)	へ(格助詞)

[表 3 : 調査対象語ごとにまとめた計算結果の上位語]

33750.954537	ない(形容詞)	しか(副助詞)
31152.314548	ない(形容詞)	しか(副助詞)
31077.058861	ない(形容詞性述語接尾辞)	しか(副助詞)
30657.711741	ない(形容詞性述語接尾辞)	しか(副助詞)
29825.747268	ない(形容詞性述語接尾辞)	しか(副助詞)
29555.966954	ない(形容詞)	しか(副助詞)
28750.048788	は(副助詞)	しか(副助詞)
28680.546905	ない(形容詞)	しか(副助詞)
28564.763913	ない(形容詞性述語接尾辞)	しか(副助詞)
28495.449370	ない(形容詞性述語接尾辞)	しか(副助詞)
28488.144448	ない(形容詞)	しか(副助詞)
26863.918369	は(副助詞)	しか(副助詞)
25462.230780	ない(形容詞性述語接尾辞)	しか(副助詞)
25187.446116	は(副助詞)	しか(副助詞)
24770.013950	ない(形容詞)	しか(副助詞)
24713.236935	は(副助詞)	しか(副助詞)
24315.596584	は(副助詞)	しか(副助詞)
24089.164039	ない(形容詞性述語接尾辞)	しか(副助詞)
23084.440570	ない(形容詞)	しか(副助詞)
22380.548901	は(副助詞)	しか(副助詞)
22084.457745	と(格助詞)	しか(副助詞)
21970.863601	ない(形容詞性述語接尾辞)	しか(副助詞)
21621.360436	ない(形容詞)	しか(副助詞)

[表 4 : 共起ペア (出現率が低いもの)]

対象語	共起語	出現率(%)
まるで(副詞)	みたいでした(助動詞)	0.02
まるで(副詞)	みたいじゃ(助動詞)	0.08
たぶん(副詞)	必ず(副詞)	0.14
たぶん(副詞)	行か(動詞)	0.23
なら(副助詞)	まだしも(副詞)	0.23
さえ(副助詞)	おろか(副詞)	0.28
さえ(副助詞)	よければ(形容詞)	0.31
ばかり(副助詞)	こそ(時相名詞)	0.32
こそ(副助詞)	か(終助詞)	0.38
さえ(副助詞)	良ければ(形容詞)	0.43

[表 5 : 共起ペア (出現率が高いもの)]

対象語	共起語	出現率(%)
のみ(副助詞)	の(接続助詞)	70.51
決して(副詞)	ない(形容詞性述語接尾辞)	67.16
ぜひ(副詞)	の(接続助詞)	54.89
しか(副助詞)	ない(形容詞性述語接尾辞)	40.66
ぜひ(副詞)	たい(形容詞性述語接尾辞)	38.54
しか(副助詞)	ない(形容詞)	37.87
決して(副詞)	で(判定詞)	32.75
すら(副助詞)	ない(形容詞性述語接尾辞)	31.72
まるで(副詞)	の(判定詞)	31.69
ぜひ(副詞)	して(動詞)	28.68

助数詞

数詞

語の特質上、今回調査対象としている特定の語との共起関係が見られるが、あくまでも「共起」であることがはっきりしている。本稿で探そうとしているのは呼応関係なので、これらは除外する。

9) 調査対象語と共起語として挙げた語が同じ文中に出現する率を出す

$$\text{出現率} = \frac{\text{共起語と調査対象語を含む文の数}}{\text{調査対象語を含む文の数}}$$

以上の手続きを経て得た共起語と調査対象語を共起ペアとする(表 4, 表 5)。

5 共起ペアから呼応関係へ

第 4 節では、補完類似度を用いて共起ペアを選

定した。本節では、第 4 節で共起ペアとした語のうち呼応関係にあるペアを選定する手順を述べる。

1) 調査対象語と共起語の距離を出す

「呼」要素が何番目の形態素か? …(a)

「応」要素が何番目の形態素か? …(b)

「呼」要素と「応」要素の距離 = (b) - (a)

上記の結果で距離が 0 より大きいものは、調査対象語が前、共起語が後ろという位置関係になっている。この位置関係にあるものが、呼応関係の「応」要素になり得る。

以上の手続きを経て得られた語が、調査対象語に対する呼応関係の「応」要素の候補語である。

6 考察

呼応関係の「応」要素の候補語を調査対象語と同一文中に出現する率順に並べる(表 6)。このと

[表 6 : 呼応関係の「応」要素の候補語]

候補	呼要素	応要素	出現率	「呼」要素が前	「呼」要素が後
×	のみ(副助詞)	の(接続助詞)	70.51	1689	2774
	決して(副詞)	ない(形容詞性述語接尾辞)	67.16	696	0
×	ぜひ(副詞)	の(接続助詞)	54.89	639	1342
	しか(副助詞)	ない(形容詞性述語接尾辞)	40.66	17104	1611
	ぜひ(副詞)	たい(形容詞性述語接尾辞)	38.54	2441	60
	しか(副助詞)	ない(形容詞)	37.87	18115	874
	決して(副詞)	で(判定詞)	32.75	351	0
	すら(副助詞)	ない(形容詞性述語接尾辞)	31.72	1564	330
	まるで(副詞)	の(判定詞)	31.69	1944	39
	ぜひ(副詞)	して(動詞)	28.68	920	812
	なんて(副助詞)	ない(形容詞性述語接尾辞)	26.18	2882	807
	もし(副詞)	して(動詞)	24.71	1238	654
	おそらく(副詞)	だろう(助動詞)	23.62	494	328
	さえ(副助詞)	ない(形容詞性述語接尾辞)	23.33	2955	898
	もし(副詞)	ない(形容詞性述語接尾辞)	23.21	1595	164
	たぶん(副詞)	ない(形容詞性述語接尾辞)	21.81	437	96
	おそらく(副詞)	ない(形容詞性述語接尾辞)	21.08	581	87
	ぜひ(副詞)	ほしい(形容詞)	21.02	1392	11
×	だけ(副助詞)	して(動詞)	20.79	17732	23227
	もし(副詞)	する(動詞)	20.72	1339	229
	ばかり(副助詞)	ない(形容詞性述語接尾辞)	20.52	3503	1016
	もし(副詞)	こと(形式名詞)	20.48	1398	152
	まるで(副詞)	ように(助動詞)	20.35	1204	60
	すら(副助詞)	いる(動詞性接尾辞)	20.33	878	388
	だけ(副助詞)	いる(動詞性接尾辞)	20.16	24707	14723
×	のみ(副助詞)	して(動詞)	20.03	1343	1403
	のみ(副助詞)	する(動詞)	20.00	1563	1101
×	すら(副助詞)	して(動詞)	19.79	481	740
	さえ(副助詞)	いる(動詞性接尾辞)	19.67	2316	970
	まるで(副詞)	いる(動詞性接尾辞)	19.52	922	195
	きつと(副詞)	だろう(助動詞)	15.35	341	307

7 おわりに

以上、補完類似度を用いて、大容量のテキストデータから半自動的に呼応関係を発見する手続きについて報告した。最後に今後の課題について述べる。

(1) 本稿では、類似度計算の結果から共起語リストを作成する段階を目視にて手作業で行っている。このプロセスを自動化したい。客観性を保証すること、そして、手作業の負荷を減らして調査対象語の大幅拡張を可能にすることが狙いである。

(2) 今回は得られた結果を「応」要素の候補語として挙げるにとどめた。妥当性を検証するために、「応」要素の詳細な分析を行う必要がある。問題のある例の一つ挙げよう。[表 6]にある「まるで(副詞)」とその「応」要素の候補語の「の(判定詞)」は、同一文中に出現する率が 31.69%と高い(共起

きに上位にくる「決してーない」「たぶんーだろう」「おそらくーだろう」などは、従来から[益岡(1991)]などで呼応関係が指摘されており、また内省や直観である程度予測がつくものである。

これに対して、同様に上位に挙がっている「おそらくーない」の組み合わせは比較の見落とされやすい。また、「きつとーだろう」は予想よりも呼応関係の色が弱く、同一文中に出現する率が 15.35%、「きつと」が「だろう」より前に位置しているものが 341 件あるのに対して「だろう」が「きつと」の前に位置するものが 307 件だった。

このように、補完類似度を使うことにより、客観的データに基づいて呼応関係を調査することができ、従来指摘されていない呼応関係や直観では気づきにくい呼応関係の発見も可能になる。

ペア 294 組中 9 位)。さらに、調査対象語の「まるで」が「の」より前に位置しているものが 1944 件あるのに対して、「の」が「まるで」の前に位置するものは 39 件である。呼応関係の候補となる条件、調査対象語が前、共起語が後ろの位置関係は十分満たしている(比率は共起ペア 294 組中 42 位)。だが、実例を見ていくと、大部分が「まるで暴走族を賛美しているかのようだ」「まるでロシア皇帝のようだ」「まるで真珠の首飾りのようのようなものになっている。これは呼応関係ではなく、 α/β 形式で連体句を形成している例と見るべきだろう。このように現段階で指標としている数値だけでは判断しきれないものがある。候補語を詳細に検討していく必要がある。

(3) 本稿では調査対象語を一部の助詞と副詞に限ったが、今後は調査対象語を拡張していきたい。今回得られた「応」要素の候補語を調査対象語として「呼」要素を探す、副詞全般や助動詞などの他の品詞を調査対象語とするなどを予定している。この作業を繰り返すことによって、将来的には呼応関係テーブルの作成を目指したい。

参考文献

- [1] 大野晋：係り結びの研究，岩波書店（1993）。
- [2] 日本語形態素解析システム JUMAN：
[http://www-lab25.kuee.kyoto-u.ac.jp/nl-re
source/juman.html](http://www-lab25.kuee.kyoto-u.ac.jp/nl-re
source/juman.html)。
- [3] 沼田善子：第 2 章 とりたて詞，いわゆる日本語助詞の研究，pp. 107-225 凡人社（1986）。
- [4] 益岡隆志：モダリティの文法，くろしお出版（1991）。
- [5] 益岡隆志田，窪行則：基礎日本語文法—改定版一，くろしお出版（1992）。
- [6] 森重敏：日本文法通論，風間書房（1964）。
- [7] 山田孝雄：日本文法学概論，宝文館（1936）。
- [8] 山本英子，梅村恭司：コーパスの中の一対多関係を推定する問題における類似尺度，自然言語処理 Vol. 9 No. 2, pp. 45-75（2002）。

データ

- [9] 日経全文記事データベース 日本経済新聞 CD-ROM 版 1990 年～2000 年版
- [10] CD-毎日新聞データ集(1991 年～2000 年)，毎日新聞社。

[表 1 : 調査対象語 (抜粋)]

	文数	こそ	しか	さえ	は	も	ばかり	のみ	すら	なら	くらい	だけ	なんて
日経 1990	1630564	3285 (0.201%)	6597 (0.405%)	2897 (0.178%)	1077337 (66.071%)	369130 (22.638%)	3641 (0.223%)	2066 (0.127%)	775 (0.048%)	5070 (0.311%)	3157 (0.194%)	39565 (2.426%)	921 (0.056%)
日経 1991	1651399	3314 (0.201%)	6958 (0.421%)	2855 (0.173%)	1100327 (66.630%)	373395 (22.611%)	3738 (0.226%)	2059 (0.125%)	940 (0.057%)	4919 (0.298%)	2967 (0.180%)	40431 (2.448%)	980 (0.059%)
日経 1992	1519753	2629 (0.173%)	5695 (0.375%)	2135 (0.140%)	1015063 (66.791%)	336653 (22.152%)	2615 (0.172%)	1481 (0.097%)	720 (0.047%)	4039 (0.266%)	2358 (0.155%)	34978 (2.302%)	745 (0.049%)
日経 1993	1463967	2659 (0.182%)	5508 (0.376%)	1986 (0.136%)	976060 (66.672%)	325565 (22.239%)	2547 (0.174%)	1518 (0.104%)	684 (0.047%)	4201 (0.287%)	2160 (0.148%)	33614 (2.296%)	691 (0.047%)
日経 1994	1456242	2408 (0.165%)	5704 (0.392%)	1974 (0.136%)	973649 (66.860%)	324751 (22.301%)	2422 (0.166%)	1643 (0.113%)	716 (0.049%)	4298 (0.295%)	2101 (0.144%)	32687 (2.245%)	695 (0.048%)
日経 1995	1642744	2957 (0.180%)	6503 (0.396%)	2129 (0.130%)	1086618 (66.147%)	357844 (21.783%)	2790 (0.170%)	2080 (0.127%)	875 (0.053%)	4740 (0.289%)	2495 (0.152%)	36272 (2.208%)	925 (0.056%)
日経 1996	1590251	2863 (0.180%)	5930 (0.373%)	1875 (0.118%)	1026011 (64.519%)	327844 (20.616%)	2553 (0.161%)	1987 (0.125%)	720 (0.045%)	4797 (0.302%)	2234 (0.140%)	32085 (2.018%)	804 (0.051%)
日経 1997	1570514	2705 (0.172%)	5740 (0.365%)	1873 (0.119%)	1013445 (64.530%)	323881 (20.623%)	2413 (0.154%)	1952 (0.124%)	761 (0.048%)	4497 (0.286%)	2210 (0.141%)	31951 (2.034%)	833 (0.053%)
日経 1998	1573688	2779 (0.177%)	5862 (0.373%)	1812 (0.115%)	1009269 (64.134%)	323804 (20.576%)	2274 (0.145%)	2107 (0.134%)	792 (0.050%)	4816 (0.306%)	2018 (0.128%)	31426 (1.997%)	726 (0.046%)
日経 1999	1554893	2518 (0.162%)	5278 (0.339%)	1546 (0.099%)	997026 (64.122%)	317837 (20.441%)	2103 (0.135%)	2029 (0.130%)	707 (0.045%)	4551 (0.293%)	1973 (0.127%)	30353 (1.952%)	649 (0.042%)
日経 2000	1541892	2384 (0.155%)	4813 (0.312%)	1563 (0.101%)	986031 (63.949%)	307111 (19.918%)	1948 (0.126%)	2179 (0.141%)	615 (0.040%)	4322 (0.280%)	1942 (0.126%)	28545 (1.851%)	645 (0.042%)
毎日 1991	791126	1598 (0.202%)	3685 (0.466%)	1351 (0.171%)	490941 (62.056%)	147605 (18.658%)	1645 (0.208%)	1222 (0.154%)	544 (0.069%)	2588 (0.327%)	1441 (0.182%)	17382 (2.197%)	966 (0.122%)
毎日 1992	847910	1745 (0.206%)	4103 (0.484%)	1350 (0.159%)	510885 (60.252%)	157312 (18.553%)	1921 (0.227%)	1190 (0.140%)	553 (0.065%)	3019 (0.356%)	1784 (0.210%)	18987 (2.239%)	1175 (0.139%)
毎日 1993	800600	1796 (0.224%)	4099 (0.512%)	1357 (0.169%)	474666 (59.289%)	150250 (18.767%)	1880 (0.235%)	1132 (0.141%)	550 (0.069%)	2981 (0.372%)	1752 (0.219%)	18336 (2.290%)	1186 (0.148%)
毎日 1994	966405	2420 (0.250%)	4962 (0.513%)	1784 (0.185%)	560758 (58.025%)	179216 (18.545%)	2382 (0.246%)	1403 (0.145%)	679 (0.070%)	3753 (0.388%)	2440 (0.252%)	21417 (2.216%)	1747 (0.181%)
毎日 1995	972363	2588 (0.266%)	4904 (0.504%)	1805 (0.186%)	563107 (57.911%)	179112 (18.420%)	2430 (0.250%)	1598 (0.164%)	705 (0.073%)	3410 (0.351%)	2110 (0.217%)	20944 (2.154%)	1369 (0.141%)
毎日 1996	1114941	2846 (0.255%)	5711 (0.512%)	2307 (0.207%)	645473 (57.893%)	207730 (18.631%)	2735 (0.245%)	1808 (0.162%)	833 (0.075%)	4206 (0.377%)	2584 (0.232%)	24319 (2.181%)	1742 (0.156%)
毎日 1997	1207115	2922 (0.242%)	6025 (0.499%)	2339 (0.194%)	696312 (57.684%)	221328 (18.335%)	2814 (0.233%)	1891 (0.157%)	690 (0.057%)	4556 (0.402%)	2624 (0.217%)	26090 (2.161%)	1986 (0.165%)
毎日 1998	1243656	3057 (0.246%)	6466 (0.520%)	2106 (0.169%)	722369 (58.084%)	228458 (18.370%)	2914 (0.234%)	2016 (0.162%)	780 (0.063%)	4999 (0.402%)	2658 (0.214%)	26446 (2.126%)	2183 (0.176%)
毎日 1999	1168681	2923 (0.250%)	5919 (0.506%)	2060 (0.176%)	675683 (57.816%)	214598 (18.362%)	2672 (0.229%)	1805 (0.154%)	758 (0.065%)	4750 (0.406%)	2718 (0.233%)	24614 (2.106%)	1794 (0.154%)
毎日 2000	1160588	2944 (0.254%)	5786 (0.499%)	1939 (0.167%)	668833 (57.629%)	214185 (18.455%)	2595 (0.224%)	1844 (0.159%)	774 (0.067%)	4860 (0.419%)	2628 (0.226%)	24725 (2.130%)	1718 (0.148%)
合計	27469292	55340 (0.201%)	116248 (0.423%)	41043 (0.149%)	17269863 (62.870%)	5587609 (20.341%)	53032 (0.193%)	37010 (0.135%)	15171 (0.055%)	89372 (0.325%)	48354 (0.176%)	595167 (2.167%)	24480 (0.089%)