

## 前後の段落との共起を利用した文章の結束性の測定

山崎 誠

yamazaki@ninja.ac.jp

国立国語研究所

## 1. はじめに

国立国語研究所が中心となって 2006 年度から構築してきた「現代日本語書き言葉均衡コーパス」(Balanced Corpus of Contemporary Written Japanese、以下 BCCWJ と略す) が 2011 年に完成し、それを利用した日本語研究の展開が期待されている。BCCWJ の特徴として、多様な日本語を収録していることやアノテーションの充実が挙げられる。それらを生かした研究が今後発多く発表されると推測される。本稿では BCCWJ のアノテーション情報を利用したテキストの特徴を捉える試みを紹介する。

## 2. 文章の結束性

結束性 (cohesion) とは、文章をひとつの統一体としてまとめあげるために必要な性質のひとつである。結束性について最初に詳細に研究を行ったのは Halliday & Hasan(1976)である。それによると、結束性とは次のような性質を有しているとされる。

「結束性が生じるのは、談話のある要素の解釈 (INTERPRITATION) が別の要素の解釈に依存する場合である。一方を効果的に解釈するためには他方に頼らなければならないという意味で、一方は他方を前提 (PRESUPPOSE) とする。こういうことが生じるとき、結束関係が成立する。その結果、前提語と被前提語という 2 つの要素が、少なくとも潜在的には、統合されて 1 つのテキストになるのである。」(邦訳 p.5)

また、結束性には文法的結束性と語彙的結束性があり、前者の手段として「指示」「代用」「省略」が、後者には「再叙 (reiteration)」と「コロケーション」がある<sup>1</sup>。再叙には以下の 4 つのタイプがある。

- (a) 同一語 (繰り返し)
- (b) 同義語 (または近似同義語)
- (c) 上位語
- (d) 一般語

<sup>1</sup> 文法的結束性と語彙的結束性の中間の性質を持つものとして「接続」が挙げられている。

結束性の考え方は、外国語学習における作文の評価などにも適用されている。

本発表では、語彙的結束性を表す手段のうちもっとも単純な同一語の繰り返しを観察することによって、文章全体の結束性の様子を探るものである。

## 3. データ

本発表では、2011 年 12 月にリリースされた『現代日本語書き言葉均衡コーパス』の DVD 版を使用した。Disk1 の M-XML フォルダに含まれる xml ファイルを対象とした。この xml ファイルは可変長サンプルと固定長サンプルを統合したもので、短単位、長単位の形態論情報のタグのほか可変長部分には文章構造のタグを含んでいる<sup>2</sup>。

本発表ではこの xml ファイルで<paragraph>というタグが付けられた部分を対象にそこに含まれる短単位の形態論情報をもとに分析を行う。結束性を見るには文も妥当な単位であるが、BCCWJ に付与された文を表すタグ<sentence>は見出しや図表のキャプションにも付与されており、通常の本文との区別をしなければならないため、今回の調査では確実に本文部分を表している<paragraph>タグを対象とした。<paragraph>タグを含むサンプル数は表 1 のとおりである。

表 1 対象サンプル数

媒体	全サンプル数	P サンプル数
出版書籍	10,117	9,742
雑誌	1,996	1,767
新聞	1,473	1,457
図書館書籍	10,551	10,369
白書	1,500	1,496
教科書	412	0
広報紙	354	354
ベストセラー	1,390	1,374
Yahoo!知恵袋	91,445	0
Yahoo!ブログ	52,680	0

<sup>2</sup> タグの詳細については小木曾ほか(2011)を参照。

韻文	252	0
法律	346	56
国会会議録	159	159
合計	172,675	26,774

教科書、Yahoo!知恵袋、Yahoo!ブログ、韻文<paragraph>タグを用いていないため、対象サンプル数はゼロである。なお、<paragraph>タグの問題点については西部ほか(2011:232)を参照されたい。

表2は、対象となったサンプルの延べ語数、段落数、1段落あたりの延べ語数、1段落あたりの異なり語数のそれぞれの平均値である。1段落あたりの延べ語数を見てみると国会会議録の値が大きい。これは国会会議録における段落の認定(1発言が1段落)が影響しているものである。なお、語数には補助記号、空白、助詞、助動詞は含まれていない。

表2 各媒体の延べ語数等

	SN	SP	PN	PV
出版書籍	1384.61	43.76	50.51	37.06
雑誌	891.17	29.81	40.05	33.27
新聞	334.33	9.28	38.78	33.33
図書館書籍	1450.16	54.53	45.76	34.70
白書	1793.10	29.32	64.74	44.33
広報紙	2903.53	103.14	28.14	23.39
ベストセラー	1404.46	69.30	29.52	24.28
法律	219.50	6.93	24.04	15.03
国会会議録	17885.87	144.06	151.30	76.21

SN:サンプルの延べ語数、SP:サンプルの段落数

PN:段落の延べ語数、PV:段落の異なり語数

#### 4. 結束性の算出方法

本発表ではある段落とそれに隣接する段落との間で共通して現れる語の多寡に着目した。語の単純な繰り返しの扱うことのメリットは、他の結束性を表す現象と比べて正確な把握がしやすいこと、また、頻繁に起きる現象であるため、観察がしやすいことである。一方、デメリットとしては観察結果が「語」の単位認定基準に依拠してしまうこと及び同じ語か異なる語かだけの把握にとどまり、意味的な関係が把握できないことである。共通する語だけでなく、類義語等まで含めた計測方法としてHoey(1991)やKároly(2002)があるが、扱っているデータ量は多くない。大量のデータを使って自動的に計

測するには語の繰り返しがもっとも適している。

本稿では、以下の式により結束性の度合いを計り、共起語率と名付けた。

$$C(a, b) = \frac{F(a, b)}{N_a}$$

C(a,b): 段落aの段落bに対する共起語率。

F(a,b): 段落aと段落bとで共通して現れる語の延べ語数を段落a内で数えた数。

N<sub>a</sub>: 段落aの延べ語数。

共起語率は、水谷(1980)の非対称類似度をもとにした指標である。そのため連続する2つの段落の間の共起語率に2つの値が存在する。後続の段落に対する共起語率と前接の段落に対する共起語率である。上述の式では、b=a+1のとき、後続段落に対する共起語率となり、b=a-1のとき、前節段落に対する共起語率となる。ただし、文章の冒頭の段落の前接段落及び最後の段落の後続段落は存在しないため、便宜的にその場合の共起語率は0とする。

この方法で共起語率を測るにはひとつ制約がある。それは、文章が2つ以上の段落から構成されていることである。そのため表1で対象としたサンプルから計340サンプルが除外される。

なお、計測対象からは言語表現とは見なさない補助記号、空白、及び文章の結束性には影響を及ぼさない助詞、助動詞を除外した。

#### 5. 結果

表3は、段落あたりの共起語の数と共起語率の平均値である。後続段落との共起語率と前接段落との共起語率とはほぼ等しい値を示している。このことは、どの媒体もそれぞれ同程度の依存関係でつなが

表3 共起語の数と共起語率

	後続段落との共起語数	後続段落との共起語率	前接段落との共起語数	前接段落との共起語率
出版書籍	12.98	0.22	12.74	0.22
雑誌	6.89	0.16	6.82	0.16
新聞	5.99	0.15	5.84	0.16
図書館書籍	10.49	0.19	10.36	0.19
白書	20.00	0.31	19.84	0.31
広報紙	5.19	0.18	5.13	0.17
ベストセラー	5.49	0.15	5.47	0.15
法律	12.16	0.48	12.31	0.47
国会会議録	40.45	0.30	39.01	0.30

っていると解釈できる。個々に眺めてみると、法律、白書、国会会議録の共起語率が高く、新聞、ベストセラー、雑誌の共起語率が低いことが分かる。

表4 NDC 別の共起語の数と共起語率

	後続段落との共起語数	後続段落との共起語率	前節段落との共起語数	前節段落との共起語率
0 総記	12.97	0.22	12.95	0.22
1 哲学	17.55	0.25	17.73	0.24
2 歴史	14.80	0.21	14.60	0.21
3 社会科学	15.02	0.24	14.84	0.24
4 自然科学	14.32	0.24	13.96	0.24
5 技術・工学	10.72	0.22	10.56	0.21
6 産業	11.03	0.21	10.82	0.21
7 芸術・美術	12.02	0.20	11.98	0.20
8 言語	10.40	0.21	10.17	0.20
9 文学	5.07	0.12	4.97	0.12
分類なし	3.46	0.13	3.45	0.13

表4は、図書館書籍のデータについて、NDC（日本十進分類法）別の共起語数と共起語率を算出したものである。図書館書籍全体では共起語率は0.19であったが、NDC 別に見ると「9 文学」と「分類なし」の値が他と比べて低いことが分かる。「分類なし」についてはデータを見ていないので理由は分からないが、「9 文学」は会話文のような短い段落が多いため、共起語率が低くなったと推測される（表3のベストセラーの値の低さもそれに起因しているであろう）。それを確かめるために、1 段落あたりの延べ語数の平均と共起語率の平均との相関を見てみよう。図1にその結果を示す。正の相関が認められ、決定係数は0.799 と高い値を示した。

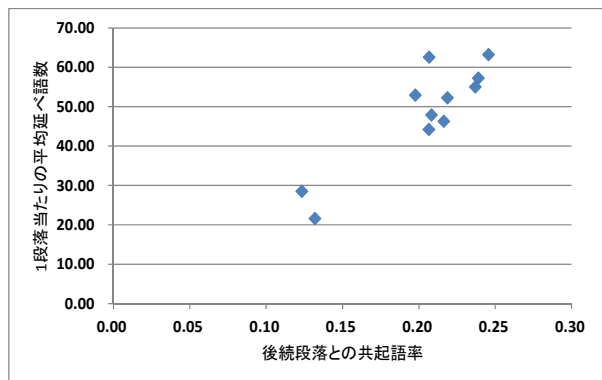


図1 段落の延べ語数と共起語率との相関

## 6. 文章中の共起語率の推移

共起語率の値はひとつの文章中でどのような変化を示すのだろうか。白書の例を見てみよう。表5はOW1X\_00000（昭和54年版経済白書）というサンプルである。

表5 文章中の共起語率の推移

段落番号	延べ語数	後続段落との共起語数	後続段落との共起語率	前接段落との共起語数	前接段落との共起語率
1	45	17	0.378	0	0
2	68	19	0.279	29	0.426
3	55	22	0.4	15	0.273
4	126	46	0.365	41	0.325
5	59	10	0.169	34	0.576
6	74	27	0.365	20	0.27
7	77	32	0.416	14	0.182
8	40	26	0.65	25	0.625
9	64	15	0.234	36	0.563
10	22	10	0.455	8	0.364
11	71	29	0.408	14	0.197
12	57	9	0.158	27	0.474
13	25	8	0.32	6	0.24
14	37	13	0.351	7	0.189
15	82	40	0.488	28	0.341
16	74	20	0.27	45	0.608
17	18	6	0.333	7	0.389
18	36	3	0.083	7	0.194
19	20	6	0.3	2	0.1
20	118	71	0.602	14	0.119
21	128	29	0.227	56	0.438
22	52	11	0.212	15	0.288
23	30	14	0.467	11	0.367
24	73	33	0.452	16	0.219
25	72	21	0.292	27	0.375
26	47	12	0.255	12	0.255
27	26	9	0.346	9	0.346
28	77	21	0.273	13	0.169
29	33	3	0.091	14	0.424
30	9	4	0.444	3	0.333
31	126	32	0.254	9	0.071
32	28	1	0.036	17	0.607
33	13	2	0.154	1	0.077
34	45	11	0.244	2	0.044

35	73	10	0.137	17	0.233
36	59	14	0.237	6	0.102
○37	52	18	0.346	8	0.154
38	64	9	0.141	24	0.375
39	22	12	0.545	6	0.273
40	81	46	0.568	20	0.247
41	88	37	0.42	50	0.568
42	46	18	0.391	21	0.457
43	60	21	0.35	19	0.317
○44	68	27	0.397	22	0.324
45	103	11	0.107	27	0.262
46	19	5	0.263	7	0.368
47	48	30	0.625	9	0.188
48	53	32	0.604	22	0.415
49	68	15	0.221	38	0.559
○50	34	12	0.353	9	0.265
51	58	21	0.362	16	0.276
52	27	8	0.296	14	0.519
53	50	20	0.4	9	0.18
54	58	14	0.241	20	0.345
55	37	11	0.297	9	0.243
56	62	0	0	13	0.21

段落番号の前に○を付けた 9 箇所は、文章中で小見出しが立っていて、内容が切り替わっていると思われる箇所である。その部分における後続段落との共起語率と前接段落との共起語率とを比べてみると、9 箇所のうち 8 箇所が後続段落との共起語率が前接段落との共起語率を上回っている（残りの 1 箇所は同じ値）。このことは、新規の内容になった最初の段落は、新しい話題を展開させるため、その次の段落との結束性が高くついていると言えるのではないだろうか。

逆に○の直前の段落は、あるまとまりの最後の段落を意味する。この部分の後続段落と前接段落の共起語率はどうなっているかというと、9 箇所中 6 箇所が前接段落との共起語率の値のほうが高い。これは一つの例にすぎないが、このような文章中での共起語率の推移を利用して段落のまとまりを自動的に推測することに応用出来る可能性がある。

## 7. まとめと今後の課題

本発表では非常に単純な指標である共起語率を用いて文章の結束性の度合いを観察した。その結果、法律、白書、国会会議録のように結束性の高い文章

と新聞、ベストセラー、雑誌のように結束性の低い文章があることが分かった。NDC 別に観察したデータでは、文学の結束性が低いという結果になった。これは文学に会話文が多いことの現れと考えることができる。

また、文章中の共起語率の推移をみることにより文章のセグメンテーションへの応用が考えられることを示した。

今後の課題として以下の 3 点を挙げる。これらを通じて文章における結束性について客観的な記述を目指したい。

(1)西部ほか(2011:232)によると、サンプルを構成する文がすべて段落に分割される訳でないと指摘されている。また、<paragraph>の認定は行頭の空白をもとに自動的に認定しているとのことなので段落の実態を確認して分析に問題がないかどうか確認する必要がある。

(2)段落と文の両方を利用した結束性の測定の方法を探る。

(3)指示詞や接続詞など文法的結束性の手段との相関を調べること。

## 〔謝辞〕

本研究は国立国語研究所の共同研究プロジェクト「テキストにおける語彙の分布と文章構造」による研究成果の一部である。

## 〔参考文献〕

- Halliday, M.A.K. and Hasan, R.(1976) *Cohesion in English*. Longman (邦訳『テキストはどのように構成されるか』、大修館書店、1997 刊)
- Hoey,Michael.(1991) *Patterns of Lexis in Text*. Oxford University Press.
- Károly,Krisztina.(2002) *Lexical Repetition in Text*. Peter Lang.
- 小木曾智信、間淵洋子、前川喜久雄(2011)『現代日本語書き言葉均衡コーパス』における形態論情報付き XML フォーマット」、言語処理学会第 17 回年次大会予稿集、pp.352-355.
- 西部みちる、大島一、間淵洋子、小林正行、田島孝治、高田智和、山口昌也(2011)『現代日本語書き言葉均衡コーパス』における電子化テキストの構築」、国立国語研究所内部報告書(LR-CCG-10-03)
- 水谷静夫(1980)「用語類似度による歌謡曲仕分『湯の町エレジー』『上海帰りのリル』及びその周辺』『計量国語学』12(4)、pp.145-161.