

## N-gramを用いた連鎖型共起表現の自動抽出法<sup>\*1</sup>

池原 悟 白井 諭 河岡 司

(NTTコミュニケーション科学研究所)

### 1. まえがき

最近、自然言語処理において、大量のコーパスや用例の重要性が指摘され、それを分析する技術の必要性が増大している。言語データの中から、使用頻度の高いフレーズや表現のパターンを抽出する方法としては、2単語の結びつきの強度に着目した方法<sup>(1)</sup>や単語間の距離に着目した方法<sup>(2)</sup>、結合単語数と出現回数を考慮した方法<sup>(3, 4)</sup>など<sup>(5)</sup>が提案されている。しかし、大量の言語データを対象とすると、計算量の問題から、任意の長さで、出現頻度の高い表現を自動的に発見して、抽出することは困難であった。

これに対して、最近、大量の言語データを対象に、任意の $n$ に対する $n$ -gram統計を高速に実行する方法が提案された<sup>(6)</sup>。この方法を用いれば、原文中に使用された文字列を、その長さ(文字数)の順かつ出現頻度の順に集計することが出来る。しかし、抽出する文字列間の相互関係が無視されているため、既に抽出された文字列の部分文字列が重複して抽出される。従って、抽出された文字列には、文法的、意味的にまとまりのない断片的な文字列が多数を占め、独立した表現の単位となるような文字列を選択する作業が大変であった。

本論文では、この問題を解決する方法として、言語データの中から、文字列を、それらが部分文字列として相互の重複のあるものは対象外とするという条件下で、文字数の多い順、かつ使用頻度の高い順に、漏れなく、自動的に抽出し、集計する方法を提案する。

### 2. 従来の方法とその問題点

#### (1) 文字列抽出の条件

自然言語表現に共起する表現としては、連語やフレーズのように連続した文字列を構成するもの(連鎖型共起表現)と、係り結び、呼応関係、特定の動詞と特定の名詞の組などのように、2種以上の文字列が、文中の離れた位置に現れるもの(離散型共起表現)がある。本報告では、前者の文字列を以下の条件で抽出す

ることとする。

第1の条件: 一致文字数の多い順に抽出する。

第2の条件: 出現頻度の高い順に抽出する。

第3の条件: ある文字列が抽出されたとき、その文字列内に含まれる部分文字列は、既に抽出済みと考え、それ以降は抽出の対象としない。

#### (2) 長尾・森の方法とその問題点

日本語文章データから、出現頻度の高い連続した文字列を抽出する長尾・森の方法<sup>(4)</sup>を要約すると以下の通りである。

##### [長尾・森の方法]

先頭文字から順に、0文字、1文字、2文字、...  
・ $N-1$ 文字を削除した $N$ 個の部分文字列(文字列単語と呼ぶ)に対して、各部分文字列の順に対応する $N$ 個のレコードを持った原文番地ファイルを作成し、文字列コード順にソートしたあと、前後の文字列の先頭からの一致文字数をカウントする方法で、以下の手順で構成される。

手順1: 原文番地ファイルの作成

手順2: 汎用ソートファイルの作成

手順3: 一致文字数のカウント

手順4: 文字列の抽出とカウント

この方法により、第1、第2の条件は満足されるが、目標とする第3の条件は満足されない。抽出された文字列の出現回数に対して、より長い文字列に含まれていた部分文字列の出現回数を差し引くなど、次数の異なる複数の $n$ -gram集計表を組み合わせる計算の方法が考えられるが、集計表の生成された時点では、抽出された文字列の原文中での相互関係の情報が失われているため、計算は不可能である<sup>\*2</sup>。

### 3. 連鎖型共起表現の抽出法

#### 3.1 重なる文字列の扱い

$n$ -gram文字列と $m$ -gram文字列の抽出を考えたとき、問題となるのは、図1に示すように、両者に何らかのオーバーラップがあるときである。

<sup>\*1</sup>Automatic Extraction of Uninterrupted Collocation by  $n$ -gram Statistics, Satoru Ikehara, Satoshi Shirai and Tsukasa Kawaoka, NTT Communication Science Laboratories.

<sup>\*2</sup>最近、単語列の場合、ある単語列 $\alpha$ の出現頻度を、その単語列の出現回数から、それを含むより大きな単語列 $\beta$ の出現回数を引くことにより評価する方法が使われていた例<sup>(4)</sup>もあるが、厳密には、この方法では正しい結果は得られない。

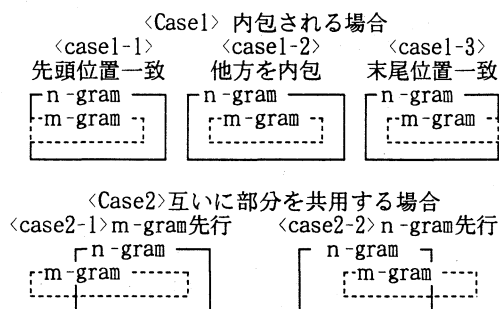


図1 抽出文字列の相互関係

### (1) 一方が他方を内包する場合

この場合は、さらに3つの場合に分けられるが、いずれの場合も、m-gram部分はn-gramとして抽出された文字列の部分に内包されるので、抽出の対象とならない。従って、n-gram文字列を抽出するとき、n-gram文字列として抽出した文字列の部分は、後の処理でm-gram文字列の対象としないよう「汎用ソートファイル」の該当レコードを無効化する必要がある。

### (2) 互いに部分を共有する場合

この場合は、さらに2つの場合に分けられるが、いずれの場合も、m-gram文字列とn-gram文字列は、原文上互いにその一部を共有するだけであるので、抽出対象文字列である。

## 3. 2 文字列抽出アルゴリズム

### (1) アルゴリズム

前章の議論をふまえ、原文データから、2回以上の出現回数を持つ固定的な(連続した)表現を文字列として、文字数の多い順に、かつ、重複なしに抽出するアルゴリズムを提案する。

#### <前準備>

まず、長尾等の方法で使用された「汎用ソートファイル」を拡張し、その各レコードに、「抽出文字数」の欄を設ける(図2参照)。

#### <文字列抽出アルゴリズム>

手順1：手順2：手順3： 長尾等の方法と同じ

手順4：「抽出文字数の記入」

上記で得た「拡張汎用ソートファイル」上の「抽出文字数」の欄に、自分のレコードの先頭から何文字までが抽出対象となっているかを記入する。具体的には、各レコードの抽出対象文字数の値を順に見て、自分のレコードの値が、直前のレコードと同一の数字か、それより大きい数字の場合は、一致文字数をコピーし、

逆に、自分の値が直前のレコードの一致文字数より小さいときは、直前のレコードの値をコピーすればよい。

### 手順5：「拡張原文番地ファイルの作成」

手順4で得られた拡張汎用ソートファイルを原文番号順に再ソートしたあと、各レコードに「採否表示」の欄を設け、それを「拡張原文番地ファイル」とする。

### 手順6：「有効無効判定処理」

拡張原文ファイル上で、以下の方法で、先頭レコードから順に、各レコードの有効判定を行い、その結果を、「採否表示」欄に記入する。

すなわち、先頭レコードから「抽出文字数」の値を調べ、iレコード目で、始めて0でない値nを発見したとする。そのとき、i番目のレコードの採否表示は「○(採用)」とし、それに続くi+1, i+2, ..., j+(n-1)のレコードの抽出文字数の値を調べる。それらの値が、それぞれ、n-1, n-2, ..., 1以下であるレコード採否表示を「×(不採用)」とし、それ以外のレコードの有効表示は「○」とする。次に、新たに「○」と判定されたレコードの番号をiとして、同様の処理を繰り返す。「○」と判定されたレコードのないときは、i+nのレコードまで進む。そのレコードの「抽出文字数」0の時は、はじめと同様、0出ない値を持つレコードまで進み、同様の手順を繰り返す。

### 手順7：「再拡張汎用ソートファイルの作成」

上記で得られた「拡張原文番地ファイル」を再度、「汎用ソートファイル」レコード順にソートし、「再拡張汎用ソートファイル」を作成する。

### 手順8：「抽出文字列集計処理」

「再拡張汎用ソートファイル」を参照して、以下の手順で、抽出する文字列を決定し、同時に、その出現回数を求める。

- ①採否表示が「○」となっているレコードを調べ、抽出文字数の最大のレコードを探す。そのレコードの一致文字数をnとする。
- ②抽出文字数の値がnの最初のレコードの番号を $i_1$ とすると、 $i_1$ から始まって、後方のレコードを見たとき、抽出文字数が連続してnとなるレコードの数が、そのn字の文字列の出現回数であるから、文字列と出現回数(但し、出現回数2以上)を集計表に記録する。次に現れた、抽出文字数がnのレコード $i_2$ に対して、前と同様にして文字列と出現回数を集計表に記録する。
- ③ $n = n - 1$ とし、②の手順を繰り返す。 $n = 0$ になったら処理は終了する。

〔原文データ〕： むかし むかしの おかしな おかし。 おかしのはなしは おかしな おはなし。  
 (抽出対象箇所) { ③ ⑥ ② ② ④ ① ④ } ○印：抽出される  
 語の文字列

原文番地ファイル

原文番地	文字列単語 〔ソートなし〕 (先頭部分)
1	むかしむかしの
2	かしむかしの
3	かしむかしの
4	むかしのおかし
5	むかしのおかし
6	かしのおかしな
7	かしのおかしな
8	かしのおかしな
9	かしのおかしな
10	かしのおかしな
11	かしのおかしな
12	かしのおかしな
13	かしのおかしな
14	かしのおかしな
15	かしのおかしな
16	かしのおかしな
17	かしのおかしな
18	かしのおかしな
19	かしのおかしな
20	かしのおかしな
21	かしのおかしな
22	かしのおかしな
23	かしのおかしな
24	かしのおかしな
25	かしのおかしな
26	かしのおかしな
27	かしのおかしな
28	かしのおかしな
29	かしのおかしな
30	かしのおかしな
31	かしのおかしな
32	かしのおかしな

〔手順2〕  
拡張汎用  
ソートファイル  
の作成

〔手順3〕→  
一致文字数  
のカウント

〔手順4〕→  
抽出文字数  
の記入

拡張汎用ソートファイル

抽出文字数	一致文字数	レコード番号	原文番地	文字列単語 〔ソートあり〕 (先頭部分)
5	5	1	8	おかしな
5	3	2	4	かしな
3	3	3	6	かしな
3	1	4	1	かしな
1	0	5	2	かしな
4	4	6	9	かしな
4	2	7	5	かしな
3	3	8	9	かしな
3	2	9	1	かしな
2	0	10	2	かしな
2	3	11	13	かしな
3	1	12	2	かしな
3	2	13	6	かしな
2	1	14	1	かしな
2	1	15	18	かしな
1	1	16	2	かしな
1	1	17	3	かしな
1	1	18	3	かしな
1	1	19	14	かしな
2	2	20	1	かしな
2	2	21	11	かしな
2	2	22	2	かしな
2	2	23	3	かしな
2	2	24	7	かしな
1	1	25	19	かしな
1	1	26	2	かしな
3	3	27	20	かしな
3	3	28	9	かしな
3	3	29	4	かしな
0	0	30	1	かしな
0	0	31	2	かしな
0	0	32	15	かしな

〔手順5〕  
拡張原文番  
地ファイルの  
作成

〔手順6〕→  
有効無効判  
定処理

拡張原文番地ファイル

採否表示	抽出文字数	一致文字数	レコード番号	原文番地	文字列単語
○	3	0	3	0	むかし
×	2	1	1	2	かし
×	1	0	7	3	かし
×	3	3	2	9	かし
×	3	1	8	4	かし
×	2	1	4	5	かし
×	2	2	1	6	かし
×	1	5	4	7	かし
×	4	3	2	8	かし
×	3	2	1	9	かし
×	2	3	2	10	かし
×	2	1	0	11	かし
×	0	0	0	12	かし
×	3	3	2	13	かし
×	2	1	0	14	かし
×	3	2	1	15	かし
×	1	3	2	16	かし
×	3	2	1	17	かし
×	2	1	0	18	かし
×	3	2	2	19	かし
×	1	3	2	20	かし
×	2	2	1	21	かし
×	1	5	3	22	かし
×	4	3	2	23	かし
×	3	2	1	24	かし
×	2	1	0	25	かし
×	1	3	2	26	かし
×	3	2	1	27	かし
×	2	1	0	28	かし
×	1	3	2	29	かし
×	3	2	1	30	かし
×	2	1	0	31	かし
×	1	0	0	32	かし

再拡張汎用ソートファイル

採否表示	抽出文字数	一致文字数	レコード番号	原文番地	文字列単語
○	5	5	1	8	かし
○	5	3	2	4	かし
○	3	3	3	6	かし
○	3	1	4	1	かし
×	4	4	5	2	かし
×	4	2	6	9	かし
×	3	3	7	5	かし
×	3	2	8	9	かし
×	2	2	9	1	かし
×	2	0	10	2	かし
×	3	1	11	13	かし
×	3	2	12	2	かし
×	3	1	13	6	かし
×	2	1	14	1	かし
×	2	1	15	18	かし
×	1	1	16	2	かし
×	1	1	17	3	かし
×	1	1	18	3	かし
×	1	1	19	14	かし
×	2	2	20	1	かし
×	2	2	21	11	かし
×	2	2	22	2	かし
×	2	2	23	3	かし
×	2	2	24	7	かし
×	1	1	25	19	かし
×	1	1	26	2	かし
×	3	3	27	20	かし
×	3	3	28	9	かし
×	3	3	29	4	かし
×	0	0	30	1	かし
×	0	0	31	2	かし
×	0	0	32	15	かし

〔手順7〕  
再拡張汎用  
ソートファイル  
の作成

〔手順8〕→  
抽出文字列  
集計処理

最終結果

〔原文番地〕 原文の単語から何番目の単語かを、レコード番号で抽出する。

〔レコード番号〕 拡張汎用ソートファイルのレコード番号。

〔一致文字数〕 直後のレコードの文字列と比べて、先頭から一致する文字数。

〔抽出文字数〕 自分のレコードで抽出する文字列の先頭から抽出する文字数。

〔採否表示〕 部分文字列として抽出するか否かの判定結果。

抽出された文字列

方法	本論文の方法		従来の方法(参考)	
gram数	文字列	出現回数	文字列	出現回数
5 gram	①おかしなお	2	おかしなお	2
4 gram	-----	-	おかしな かしなお	2 2
3 gram	②おかし	2	おかしな かしなの かしな はなし むかし	4
	③かしの	2		2
	④はなし	2		2
	⑤むかし	2		2
2 gram	--	-	おかしな かしなの かしな はなし はむか	4
	--	-		6
	--	-		2
	--	-		2
	--	-		2
	--	-		2
	--	-		2
	--	-		2
1 gram	-	-	おかしな のはむ	5
	-	-		6
	-	-		8
	-	-		4
	-	-		2
	-	-		3
	-	-		2
---	合 計	1 0	合 計	7 2

(2回以上出現した文字列で、部分的に相互重複のない文字列)

図2 連続型共起表現抽出アルゴリズム実施例

## (2) 例題検討

以上のアルゴリズムの適用例を図2に示す。この例では、従来のアルゴリズムで抽出される文字列の種類が24種類で、延べ抽出文字列数が72件であったのに対して、本論文の方法では、6種類、12件に絞られる。

## 4. 共起表現の抽出実験

日経産業新聞記事3月分(892万字)を対象に、連鎖型共起文字列の抽出実験を行った。その結果を、表1～4に示す。使用した計算機は、XEROX ARGOS 5270(SUN OS4.1.3)で、使用したメモリ量は、最大48MBである。

表1 抽出された文字列の種類と頻度(その1)

抽出対象	本論文の方法		長尾・森の方法	
	文字列の種類数	延べ出現回数	文字列の種類数	延べ出現回数
2文字以上	970,203	2,613,704	4,374,141	31,178,897
5文字以上	591,901	1,476,922	2,960,487	10,808,458
10文字以上	52,214	114,270	673,601	1,550,817
20文字以上	1,792	3,692	177,298	359,810

表2 抽出された文字列の種類と頻度(その2)

抽出対象	本論文の方法		長尾・森の方法	
	文字列の種類数	延べ出現回数	文字列の種類数	延べ出現回数
2回以上	970,203	2,613,704	4,377,087	39,588,291
5回以上	67,321	551,441	882,217	31,288,701
10回以上	12,351	217,934	372,291	28,050,199
20回以上	2,288	92,804	169,375	25,871,964
50回以上	285	37,850	62,991	22,209,875
100回以上	76	24,167	30,316	19,961,961
200回以上	20	16,771	14,363	17,759,432

表3 抽出された文字列の例(10gramの例)

本論文の方法	長尾・森の方法
することになっている(44) 第二次臨時行政調査会(35) することになりそうだ(19) 82ジャパニッシュ(17) しているのではない(16) ワシントン九日共同(14) サウジアラビア通商(14) として注目されている(14) ニューヨーク二十五日(13) を余儀なくされている(13) [合計 21,155種 延べ 47,336回]	したところによると、(273) らかにしたところによ(223) 明らかにしたところ(223) かにしたところによる(222) にしたところによると(222) 第二次臨時行政調査会(208) ることを明らかにした(191) 日明らかにしたところ(173) たことを明らかにした(120) ホテルニュージャバ(112) [合計 132,865種類 延べ 345,232回]

表4 出現頻度の高い文字列の例

高頻度文字列	200回以上
という(586),と述べた(512),としている(436),また(325),である(324),写真(315),しかし,(302),と語った(283),東京(281),価格(278),欧州共同体(277),しかし(274),ポイント(269),ひとこと(264),発売期間(259),また,(236),これは(220),このため(204),ただ,(201)	

これらの表から、以下のことが分かる。

①従来の方法に比べて、本論文の方法では、抽出される文字列の種類、出現回数共に大幅に抑制される。

例えば、2文字以上、2回以上の文字列では、抽出される種類が5分の1、延べ出現回数も約10分の1に抑制される。この効果は、文字数の大きい文字列ほど大きく、20文字以上の場合では、100分の1にもなる。

②本論文の方法で抽出された文字列は、文法的、意味的に表現の単位になるものが多く、断片的な文字列の抽出は大幅に抑制される。

処理時間では、最初の文字列単語のソート(汎用ソートファイルの作成)に、最も多くの時間(数十時間)がかかったが、その後の処理は、それに比べて短時間(約1～2時間)であった。

## 5. あとがき

言語コーパスなどの膨大な言語データから、使用頻度の高い表現を自動的に、長い順かつ頻度の高い順に、漏れなく発見し、集計する方法を提案した。

約890万字の新聞記事データに適用した例によれば、従来の方法に比べて、断片的な文字列の抽出が抑制され、抽出される文字列の種類は20～1%に、またそれらの述べ出現回数は14～1%程度に削減されることが分かった。

このように、文法的、意味的な表現の単位としての連鎖型共起表現が得られるようになったため、これを応用すれば、離散型共起表現も自動的に抽出することが可能となった。それについては、別稿で述べる。

なお、本論文では、日本語文字列への適用例を示したが、単語列や形態素解析結果への適用も同様である。

## 参考文献

- (1) Church, K.W. and Hanks, P. (1992): Word Association norms, Mutual Information and Lexicography, *Computational Linguistics*, Vol.16, No.1, pp.22-29
- (2) Smadja, F.A. and McKeown, K.R. (1992): Automatically Extracting and Representing Collocations for Language Generation, *Proceedings of the 28th Annual Meeting of the Association for Computational Linguistics*, pp.252-259
- (3) 北, 小倉, 森元, 矢野 (1993): 仕事量基準を用いたコーパスからの定型表現の自動抽出, 情報処理学会論文誌, Vol.34, No.9, pp.1937-1943
- (4) Kita, K., Kato, Y., Omoto, T. and Yano, Y. (1994): A Comparative Study of Automatic Extraction of Collocations from Corpora: Mutual Information vs. Cost Criteria, *Journal of Natural Language Processing*, Vol.1, No.1, pp.21-33
- (5) Smadja, F. (1993): Retrieving Collocations from Text: Xtract, *Computational Linguistics*, Vol.19, No.9, pp.143-177
- (6) Nagao, M. and Mori, S. (1994): A New Method of N-gram Statistics for Large Number of n and Automatic Extraction of Words and Phrases from Large Text Data of Japanese, *The Proceedings of the 15th International Conference on Computational Linguistics*, pp.611-615