

## 擬似キーワード相関法による重要キーワードと重要文の抽出

亀田 雅之

(株)リコー・研究開発本部・情報通信研究所

kameda@ic.rdc.ricoh.co.jp

### 概要

キーワード候補(擬似キーワード)間の部分文字列の重複に基づき、日本語文書中のキーワードやキーセンテンスを重要度付きで抽出する。

キーワードの重要度は、部分文字列の重複を利用して修正単語頻度及び単語長と字種から概算予測した構成単語数により評価する。また、文ごとに含まれるキーワード候補群同士の部分文字列の重複の総和により得た文間の関連度に基づく文マトリクスを利用して、文の重要度評価を行う。新聞記事や論文等の文書に有効で、重要度に基づいた適切なサイズの重要キーワードリストや抄録の作成や、関連度に基づいた関連文を提示する機能を実現できる。

### 1 はじめに

我々は、日本語文書の「読み」の支援として、軽量・高速な日本語解析系QJP[1]を利用した日本語文書読解支援系QJR[2]の機能を探っている<sup>1</sup>。QJRの4つの支援機能<sup>2</sup>のうち、Screening支援では、キーワードが文書の対象や分野を端的に類推する手掛かりになることに着目し、簡易判定したキーワード候補群を提示して文書のふるい分け支援として利用し、Skip Reading支援では、文書中の重要文を識別表示することで、飛ばし読みを支援することを狙った。

我々は、このキーワードと重要文を、文書検索システムでの検索式に対する検索文書の適合度や、検索文書をキーワード群や抄録等で内容を縮約してユーザに提示する支援機能にも有用であると注目している。特に、キーワードや文が重み付け(重要度付与)されているとこれらの扱いの可能性が拡がる。

キーワードは、従来、文書を検索するためのキーとして付与されてきた。自動的なキーワード抽出としては、構成単語や形態素のキーワード性を分析する方法

<sup>1</sup>我々のアプローチは、深い意味理解や内容理解には至らず、表面から形態素/構文レベルでの比較的幅広い処理での実現を目指す

<sup>2</sup>1.Screening(文書のふるい分け)支援、2.SkipReading(飛ばし読み)支援、3.Skimming(走り読み)支援、4.Analytic Reading(分析的な読み)支援

[3]があるが、一般には、文書内の単語の単語頻度が重要な指標とされている[4]。しかし、単語頻度では、同一の単語の出現がないと計数されない。一方、QJRでは、複合語は、それだけ特殊性/専門性が高く、文書を特徴付けると考え、単語長に着目し、処理コストの小さい方法としてQJPによる形態素解析結果の品詞と単語長に基づいて文書中のキーワード候補(擬似キーワードと呼ぶ; 図1)を判定した<sup>3</sup>。

重要文抽出に関連する抄録の手法としては、接続詞等の手掛かり語により論理的な構成を分析し、重要な文を判定したり[5]、キーワードを多く含む文を抽出することが考えられる。前者は、文脈解析的側面が強く、難しい。この他<sup>4</sup>に、キーワードや抄録生成については、単語間の共起[6]や分類語彙表に基づく関連性[7]等を利用した研究がある。

本稿では、日本語文書中のキーワードや重要文(キーセンテンス)を重要度付きで抽出するために、擬似キーワード間の部分文字列の照合に基づいた擬似キーワード相関法と呼ぶ手法を主体にした方法を提案する。

#### 通常兵器製造の工業製品 輸出額割が始動 4ヶ国対象

通常兵器の部品や加工機械に転用できる工業製品の輸出規制が二十日、日本でも始まった。英米などの主要先進七カ国(G7)の合意に基づいた調査であり、イラクなど四カ国を対象にして、対共産圏輸出禁制委員会(ココム)のリストを準用する。G7は既に対象となる品目、国を広げるための話し合いを始めており、冷戦終結で変わった輸出規制に発展しそうだ。

調査対象となる国は、イラン、イラク、リビア、朝鮮民主主義共和国(北朝鮮)の四カ国である。北朝鮮は、既に共産圏として特定地域に指定されているため、新たに追加されるのは三カ国である。また、イラクは脱資源で技術移譲が取られている。今回の措置で輸出に大きな変化が出るのはイランとリビアの二国になりそうだ。

輸出貿易管理令などに基づいて定められたコンピュータや工作機械などの機械品目を輸出する時には、運送會に許可申請を行う。その際、特定地域に指定されている国に対しては、明らかに裏取引とわかる場合でなければ許可が下りず、事実上、機械品は輸出できない。

[文書出典:朝日新聞 1993年1月21日]

図1. QJRでの擬似キーワードの識別表示

<sup>3</sup>尚、単語頻度の考慮は必要であり、Screening支援では、擬似キーワードを原文と同じ位置で抽出表示する他に、単語頻度ごとに単語長による評価で簡単に順序付けした一覧表で示す方法を示した

<sup>4</sup>[2]では、重要文の簡便な判定方法としてSkip Readingと呼ばれる手法で取られている文書や段落の主意を表しやすい位置にある段落や文(先頭/末尾段落、段落内の先頭/末尾の文)に着目する経験的な方法を紹介した

## 2 擬似キーワード相関法

### 2.1 擬似キーワード

文書におけるキーワードを文書内容の対象や分野を示す語<sup>5</sup>と考える。QJRでは、数名詞や機能性の形式名詞や副詞名詞を除いた2文字以上の名詞をキーワード候補(擬似キーワード)とした。2文字以上としたのは、一般に1文字からなる和語名詞は、キーワードとはなりにくいという推測に基づく。本稿での擬似キーワードも同じ基準で選ぶものとする(図1)。

### 2.2 擬似キーワード相関法

キーワードの評価において単語頻度を指標とするのは、出現数の多い単語がキーワード性を高めているという観点である<sup>6</sup>。しかし、たとえば、図1の文書中のキーワード候補(擬似キーワード)の「輸出規制」、「対共産圏輸出統制委員会」、「規制」、「規制対象」、「輸出」、「輸出貿易管理令」といった出現単語については、「輸出」や「規制」という共通する構成単語がお互いに関連し、キーワード性を高め、文書を特徴付けていると見ることができるにもかかわらず、単語頻度だけではこうした寄与が見逃されてしまう。

擬似キーワード相関法<sup>7</sup>の基本的なアイデアは、こうした寄与を考慮することにある。擬似キーワード相関法では、上記の共通の構成単語の検出を単語間の部分文字列の照合で代替し<sup>8</sup>、その程度を単語間の関連度とみなし、また、以降で示すように単語頻度の修正として加味するものである。

擬似キーワード相関法の定式化として、擬似キーワード  $W_i$  の  $W_j$ に対する関連度  $RW_{W_i}$  を次のように与える。

$$RW_{W_i}(W_j) = CW(W_i, W_j) / LW(W_i)$$

$CW(W_i, W_j)$ :  $W_i$  と  $W_j$  の重複文字列長  
 $LW(W_i)$ :  $W_i$  の文字列長

上記の例では、「輸出規制」と「対共産圏輸出統制委員会」との組では、 $CW$ (輸出規制, 対共産圏輸出統制委員会) が 3(「輸出」と「制」の重複) であるので、 $RW_{\text{輸出規制}}(\text{対共産圏輸出統制委員会})$  と  $RW_{\text{対共産圏輸出統制委員会}}(\text{輸出規制})$  は、各々の文字列長で除して  $3/4$  と  $3/11$  となる。

<sup>5</sup>キーワードは「何」について述べられているかを示す語であつて、「どうした」という内容を示すものではない

<sup>6</sup>対象文書群全体での頻度との比をとった方がよい指標となる

<sup>7</sup>擬似キーワードに限定するものではないが、擬似キーワードに適用することで効果を上げていることからこの名称を与えた

<sup>8</sup>簡便な部分文字列の重複の計数によるのは、形態素解析系とした用いたQJPが複合語の分割を行なわないということによる

さらに、この関連度を、文や段落、文書全体といった言語単位に含まれる擬似キーワード群に拡張する。即ち、擬似キーワード群  $G_i$  の  $G_j$ に対する関連度  $RG_{G_i}(G_j)$  を次のように与える。

$$RG_{G_i}(G_j) = CG(G_i, G_j) / LG(G_i)$$

$$CG(G_i, G_j) = \sum_{w_n \in G_i, w_m \in G_j} CW(W_n, W_m) \quad ^9$$

$$LG(G_i) = \sum_{w_n \in G_i} LW(W_n) \quad ^{10}$$

## 3 キーワードの重み付け

キーワードの重み付けでは、従来の拡張として、前節の擬似キーワード相関法に基づく修正単語頻度に、予測構成単語数という指標を加え、擬似キーワードの重みを評価する。

### 3.1 修正単語頻度

修正単語頻度は、上記の擬似キーワード相関法に基づいて複合単語の構成単語レベルでの単語重複を単語頻度に加味して補正したもので、従来、キーワード抽出で利用されていた単語頻度に代わるものである。

修正単語頻度は、次のような類推に基づく。たとえば、4文字単語「輸出規制」は、「対共産圏輸出統制委員会」との3文字の重複によって、単語頻度を  $3/4$  だけ増やすと考える。一方、11文字単語「対共産圏輸出統制委員会」は、3回出現する「輸出規制」と3文字の重複があることから、 $3/11 \cdot 3 = 9/11$  だけ単語頻度を増やすとする。同様に他の擬似キーワードのすべての寄与を合わせて単語頻度に反映する。

これを上記の擬似キーワード相関法の定式化の拡張として扱う。即ち、擬似キーワード相関法で対象とする擬似キーワード  $W_i$  の文書全体の擬似キーワード群  $D$ に対する関連度  $RG_{W_i}(D)$  をもって、修正単語頻度  $FW(W_i)$  とする。即ち、

$$FW(W_i) = \sum_{W_j \in D} CW(W_i, W_j) / LW(W_i).$$

これによれば、「輸出規制」の単語頻度は 3 だが、文書中の他の擬似キーワードとの重複の寄与により、 $FW$ (輸出規制) は  $26/4=6.5$  となる。また「対共産圏輸出統制委員会」の単語頻度は 1 だが、文書中の他の擬似キーワードとの重複の寄与により、 $FW$ (対共産圏輸出統制委員会) は  $24/11=2.2$  となる<sup>11</sup>。

<sup>9</sup>  $G_i$  と  $G_j$  に含まれる擬似キーワード同士の重複文字列長の総和

<sup>10</sup>  $G_i$  に含まれる擬似キーワードの文字列長の総和

<sup>11</sup>  $FW_{W_i}(S)$  は、 $W_i$  と  $W_j$  が一致しない場合は、 $CW(W_i, W_j)$  に 0 を与えるとすると、通常の単語頻度となる

### 3.2 予測構成単語数

QJRでは、単語が複合しているということは、それだけ単語の意味の限定が強まるとして、単語長をキーワード性の判定に用いた。この考え方を一般化し、単語長に代え、構成単語数という指標を考える。

擬似キーワード  $W_i$  の構成単語数として、単語長とその単語の表記文字種に応じた 1 構成単語の平均長から概算予測した予測構成単語数  $NW(W_i)$  を用いる<sup>12</sup>。

$$NW(W_i) = LW(W_i)/lw(T(W_i))$$

$lw(T)$  : 文字種 T からなる構成単語の平均長

$T(W_i)$  :  $W_i$  の表記文字種

たとえば、 $lw(\text{漢字}) = 2$  なら、前記の例では、 $NW(\text{輸出規制})$  は  $4/2=2$ 、 $NW(\text{対共産圏輸出統制委員会})$  は  $11/2=5.5$  となる<sup>13</sup>。

### 3.3 擬似キーワードの重み付け

上記の修正単語頻度と予測構成単語数は、いずれも値が大きい程、キーワード性が高くなると考えられる。このことから、擬似キーワード  $W_i$  の重み付けの評価式  $WW(W_i)$  の 1 例として、 $c_f, c_n$  を重み係数として、次のように 2 要素の線形和を考える。

$$WW(W_i) = c_f \cdot FW(W_i) + c_n \cdot NW(W_i)$$

### 3.4 実験

上記評価式の 2 つの重み係数の何組かの組合せ試行のうち  $c_f = c_n (= 10)$  の場合に、下記表 1 に示す再現率/正解率が比較的良好な値を示した。図 2 に、その係数での図 1 文書の擬似キーワードに対する重み付けの例を示す。尚、ここでは、 $lw$  は統計値ではなく、経験的な概数 [ $lw(\text{漢字}) = 2, lw(\text{カタカナ}) = 5$ ] を用いた。

表 1 は、177 新聞記事ごとに人手で付与した上位 10 キーワードに対し、単語頻度のみ、修正単語頻度のみ、予測構成単語数のみ、及びこれらの線形和によって各々重み付けされた上位の擬似キーワードに基づく部分一致照合での再現率/正解率を示したものである。単語頻度に対しての修正単語頻度及び予測構成単語数の単独の重み付けの再現率/正解率が高く、さらに、これらの線形和による重み付けの再現率/正解率がさらに高いことが確認できる。

尚、他の単語に完全に含まれる単語は、上位キーワード抽出の際にリストから除いた (cf. 最長語併合 [6])。

<sup>12</sup> ここでは、形態素解析を行う QJP が複合名詞を分割しないことからこうした方法をとる

<sup>13</sup> 複数の表記文字種からなる複合語（例「擬似キーワード相関法」等）は、その文字種ごとの部分の予測構成単語数の和とする

## 4 文の重み付け

文の重み付けでは、各文を文中の擬似キーワード群で代表し、擬似キーワード相関法により文間の関連度を求め、関連度に基づき文の重要度を評価する。

### 4.1 文相関法

文内の擬似キーワード群に擬似キーワード相関法を適用する（特に、文相関法と呼ぶ）。文  $S_i$  の  $S_j$  に対する関連度  $RS_{S_i}(S_j)$  を次のように与える。

$$RS_{S_i}(S_j) = CS(S_i, S_j)/LS(S_i)$$

$$CS(S_i, S_j) = \sum_{W_m \in S_i, W_n \in S_j} CW(W_m, W_n) \quad ^{14}$$

$$LS(S_i) = \sum_{W_m \in S_i} LW(W_m) \quad ^{15}$$

### 4.2 文の重要度の指標

$RS_{S_i}(S_j)$  自体は、文間の関連度を示すものであるが、これを用いて、文の重要度の指標として、次のような指標を導入する。

平均関連度 :  $ARS_{S_i} = \sum_{S_j \in D, i \neq j} RS_{S_i}(S_j)/n'$

カバレージ :  $CRS_{S_i} = \sum_{S_j \in D, i \neq j} \delta(RS_{S_i}(S_j))/n'$

$n'$  : 文書内文数-1     $\delta(x) = 0(x=0), 1(x \neq 0)$

平均関連度は他の文との関連度の平均値であり、カバレージは他の文とどの程度広く関連しているかを示し、各観点での文の重要度に関わる。

( 1 )	85 (3)	輸出規制	( 10 )	40 (3)	イラク
( 2 )	77 (1)	対共産圏輸出統制委員会	( 13 )	35 (1)	加工機械
( 3 )	75 (1)	規制対象	( 13 )	35 (1)	禁輸措置
( 4 )	68 (1)	規制品目	( 13 )	35 (1)	許可申請
( 5 )	59 (1)	朝鮮民主主義共和国	( 13 )	35 (1)	工作機械
( 5 )	59 (1)	輸出貿易管理令	( 17 )	30 (2)	G 7
( 7 )	47 (1)	通常兵器関連	( 17 )	30 (1)	話し合い
( 8 )	45 (1)	主要先進七ヶ国	( 17 )	30 (1)	経済制裁
( 9 )	42 (2)	北朝鮮	( 17 )	30 (2)	リビア
( 10 )	40 (2)	特定地域	( 17 )	30 (2)	イラン
( 10 )	40 (2)	工業製品	( 17 )	30 (1)	冷戦終結

\* ( 10 ) 40 (2) 工業製品 : (順位) 重み (単語頻度) 擬似キーワード

図 2. キーワードの重み付け抽出 (上位20位)

表 1. 修正単語頻度、予測構成単語数による  
キーワード抽出の評価(再現率/正解率)

抽出数(平均数)	上位5(5.00)	上位10(9.99)	上位20(19.02)
単語頻度 $F_w$	27.9%/46.7%	48.8%/42.2%	73.5%/35.9%
修正単語頻度 $FW$	31.5%/55.7%	53.1%/48.6%	76.8%/38.6%
予測構成単語数 $NW$	33.7%/54.8%	54.7%/46.8%	78.5%/38.5%
$F_w + NW$	34.8%/57.1%	56.9%/48.6%	79.2%/38.7%
$FW + NW$	35.6%/60.3%	57.5%/52.0%	79.5%/39.6%

※対象文書：新聞記事177    ※正解評価方法：部分一致照合

※正解キーワード：記事ごとに人手でつけたキーワード上位10個

※再現率：正解キーワード中、抽出キーワードで照合できた割合

※正解率：抽出キーワード中、正解キーワードで照合できた割合

<sup>14</sup>  $S_i$  と  $S_j$  に含まれる擬似キーワード同士の重複文字列長の総和

<sup>15</sup>  $S_i$  に含まれる擬似キーワードの文字列長の総和

### 4.3 実験

図3に、文間の関連度  $RS_{S_i}(S_j)$  を要素とする文マトリクス及び上記2指標値とそれらの積/和値を適宜、正規化して示した。第n行には第n文の他の文に対する関連度が示されている。見出し文は局所的に強い関連などに対し、第5,7文は全体に広く関連している。

図4は、2指標値の積  $ARS_{S_i} \cdot CRS_{S_i}$  を重要度とした場合の上位文の抽出例である。経験的に重要と見られる見出しや先頭/末尾の段落や文の選択 [Skip Reading手法] に近い。図5では、関連度により第1文の見出しの関連文を抽出表示した。短い見出し文の内容を敷衍する本文中の説明文が抽出されている。

## 5 まとめと今後の展開

擬似キーワード相関法は、文書の部分/全体の内容をその中に含まれる擬似キーワード群で代表させた上で、部分文字列照合手段により相互の関連度を得る手法である。本稿では、擬似キーワード相関法を提案し、これをベースに、擬似キーワードと文書全体の擬似キーワード群により修正単語頻度、文内の擬似キーワード群同士で文間の関連度を導き、キーワードと文の重み付けへの応用を示した。充分な評価には至っていないが、実験によりその効果を確認した<sup>16</sup>。また、文のカバレージや関連文といった新しい観点を示した。

本手法は、表層レベルの簡易なアプローチながら、文マトリクスに見られるような文脈レベルの手掛かりを与える。さらに、段落や章・節といった単位への拡張も可能である。

適切なキーワードや重要文の抽出には、従来の手法 ([3][5][6][7] 等)との組合せも検討すべきであるが、本手法は、単独でも小コストで高い効果が見込める。特に、論文や新聞記事等の漢語や外来語が多い文書に有効であると考えられる。

今後は、これらを利用して、QJRの支援機能を改善・拡張する一方、文書検索での検索式に対する検索文書の評価や文書の縮約提示(図6<sup>17</sup>)等に利用する予定である。また、重み付けや抄録の適切な評価法の検討や大規模な評価実験の実施も今後の課題となる。

謝辞 キーワード評価のためのキーワードデータ及びツールを提供してくれた小川泰嗣研究員に感謝する。

<sup>16</sup>キーワードの重み付けでは、予測構成単語数という単純な指標の有効性も確認した

<sup>17</sup>擬似キーワード及び文の各重み付けを用いた縮約例である。重み値の閾値を変更することで縮約サイズを変更できる

[ i ] : .. RS1(j).. < ARS1 , CRS1   Ai*C1, Ai+C1 >	文の番頭
[ 1 ] : *008000000000 < 7 [1], 9 [12]   0 [12], 8 [12] >	通常兵器関連の
[ 2 ] : 0*0777300373 < 33 [ # ], 63 [ 5 ]   21 [ # ], 48 [ 3 ] >	輸出規制が始動
[ 3 ] : 00*0AA000000 < 27 [ 2 ], 27 [ 9 ]   7 [ 7 ], 27 [ 9 ] >	4ヶ国対象
[ 4 ] : 420*#22100131 < 14 [ 4 ], 72 [ 3 ]   10 [ 3 ], 43 [ 5 ] >	通常兵器の部品
[ 5 ] : 0111*#22101111 < 9 [10], 90 [ # ]   8 [ 5 ], 50 [ 2 ] >	英米などの主要
[ 6 ] : 02124*200131 < 16 [ 3 ], 72 [ 3 ]   11 [ 2 ], 44 [ 4 ] >	G 7 は既に対象
[ 7 ] : 011132*21111 < 12 [ 6 ], 90 [ 1 ]   10 [ 3 ], 51 [ # ] >	規制対象となる
[ 8 ] : 0000305*00004 < 10 [ 8 ], 27 [ 9 ]   2 [ 10 ], 19 [ 10 ] >	北朝鮮は、既に
[ 9 ] : 00002020*200 < 4 [ 12 ], 27 [ 9 ]   1 [ 11 ], 15 [ 11 ] >	また、イラクは
[ 10 ] : 020222302*20 < 12 [ 6 ], 63 [ 5 ]   7 [ 7 ], 37 [ 7 ] >	今回の措置で輸
[ 11 ] : 0202221001*2 < 10 [ 8 ], 63 [ 5 ]   6 [ 9 ], 37 [ 7 ] >	輸出貿易管理令
[ 12 ] : 0202223004* < 14 [ 4 ], 63 [ 5 ]   8 [ 5 ], 38 [ 6 ] >	その際、特定地

j : 123456789ABC ※ A, B, C : 10, 11, 12 ※ [ ] 内: 順位

図3. 文マトリクス

### 輸出規制が始動

通常兵器の部品や加工機械に転用できる工業製品の輸出規制が二十日、日本でも始まった。

G 7 は既に対象となる品目、国を広げるための話し合いを始めており、冷戦終結で変わった新たな輸出規制に発展しそうだ。

規制対象となる国は、イラン、イラク、リビア、朝鮮民主主義共和国（北朝鮮）の四ヵ国である。

図4. 重要文の抽出（上位4文）

### [1] 通常兵器関連の工業製品

→ [4] 通常兵器の部品や加工機械に転用できる工業製品の輸出規制が二十日、日本でも始まった。

図5. 第1文の関連文の抽出

### 【擬似キーワードベースト8】

輸出規制／対共産圏輸出統制委員会／規制対象／規制品目／朝鮮民主主義共和国／輸出貿易管理令／通常兵器関連／主要先進七ヵ国

### 【重要文ベースト2】

輸出規制が始動／G 7 は既に対象となる品目、国を広げるための話し合いを始めており、冷戦終結で変わった新たな輸出規制に発展しそうだ。

図6. 文書の縮約表示

## 参考文献

- [1] 亀田雅之：「大量・高速な日本語解析ツール『簡易日本語解析系 QJP』」、言語処理学会 第1回年次大会、1995.
- [2] 亀田雅之：「日本語文書読解支援系 QJP の検討」、情処学会 自然言語処理研究会報告 110-9、1995.
- [3] 小川泰嗣、望主、別所：「複合語キーワードの自動抽出法」、情処学会 自然言語処理研究会報告 97-15、1993.
- [4] 久保田：「キーワード抽出装置」、特開昭 63-244259、1987.
- [5] 住田一男、小野、三池：「対話的文書検索のための文書構造解析」、情処学会 自然言語処理研究会報告 97-11、1993.
- [6] 原正己：「単語共起と語の部分一致を利用したキーワード抽出法の検討」、情処学会 自然言語処理研究会報告 106-1、1995.
- [7] 鈴木斎、増山、内藤：「語彙の結束性に基づいた文章抄録法」、情処学会 自然言語処理研究会報告 98-9、1993.