

## 計算機による日本語文書校正のための基礎データの収集と解析

3B-6

その1 KWICによる特定表現の抽出とその結果

鈴木 恵美子 中山 妙子 吉岡 千草 山田圭以 福田恭子  
東京家政学院筑波短期大学

### 1. はじめに

コンピュータは数値計算をするためのものだと思われており、実際大量の数値データの処理に威力を発揮しているが、私たちのゼミでは、将来的には計算機（ワードプロセッサ）により、日本語文書の校正処理を行うことを目的として研究を進めている。それを行うためにはまず、大量の日本語文書の構成に関するデータと起こり得る誤りのデータを収集し、統計を取る必要があると考えている。

今回は、これらを行うための基本的なツールであるKWICのプログラムを開発した。また、これを用いて、特定の表現（論文や報告書によく使われる「～に関する」、「～に対する」）の出現状況やその表現の前後の単語関係を調べて、その結果を考察したのでそれについて述べる。

### 2. 大量の文書からの特定表現の抽出

図1 Volごとの総出現数(～に関する)

図2 Volごとの総出現数(～に対する)

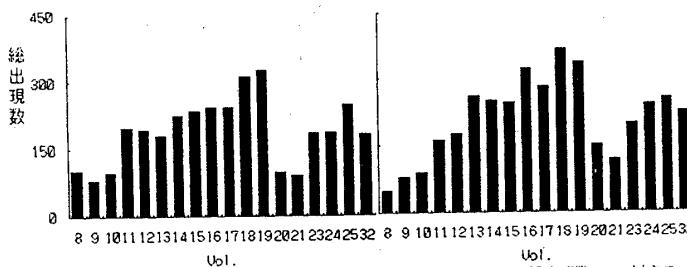


図3 No.に対しての総出現数(～に関する)

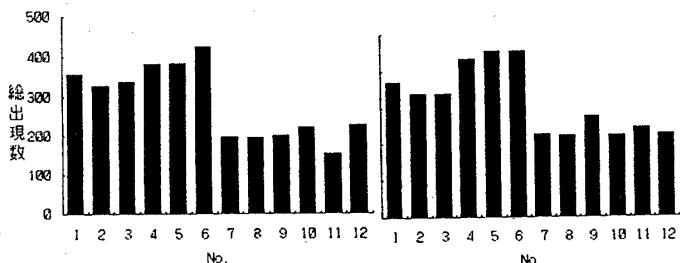


図1よりVolごとにみた「～に関する」の出現数は最高値がVol.19の326、最高値がVol.9の80で、範囲は246である。

図2よりVolごとにみた「～に対する」の出現数は最高値がVol.18の367、最高値がVol.8の49で、範囲は318である。

図3よりNo.ごとにみた「～に関する」の出現数は最高値がNo.6の425、最高値がNo.11の152で、範囲は273である。

図4よりNo.ごとにみた「～に関する」の出現数は最高値がNo.6の455、最高値がNo.8の221で、範囲は234である。

Collection of data and its analysis to be used for Japanese proofreading by using the computer  
Etsuko SUZUKI, taeko NAKAYAMA, chigusa YOSIOKA, kei YAMADA, yasuko FUKUDA  
Tokyo Kaseigakuin Tsukuba Junior College

本研究は、財團法人日本科学協会の毎年科学研究助成によって実施したもので

それぞれ200以上のひらきがあることがわかる。「～に関する…」が1つ出現するために必要な頁数は、Vol.8で2.55ページ、Vol.21で5.62ページと3ページ以上もあるのにに対し「～に対する…」が1つ出現するために必要な頁数は、Vol.8、Vol.10を除いては平均している。当初、この数字をもとに「～に関する…」「～に対する…」それぞれ5000例を収集することを目標として読み進めたが、物理的に入手できる論文誌の数も制限があったため、最終的に論文誌158冊、「～に関する…」が3398例、「～に対する…」が3800例で、収集作業は終了した。

全体的に見ると、古い年代に発行された論文誌よりも、新しい論文誌の方が「～に関する…」「～に対する…」の出現率が高いことがわかった。しかし、Vol.20 Vol.21では減少している。これは、雑誌の名前がVol.19までは「情報処理」だったが、Vol.20からは、「情報処理学会論文誌」と変わると共に内容が変化したためだと考えられる。

### 3. 入力ファイルの形態について

人手によって見つけられた、キーワードを含む一文をMS-DOSファイルにキーワードから入力してデータファイルを作成する。

- データを入力する際の注意とは、
1. 全て全角で入力する
  2. 公式は記号化して別に記録しておく  
 $a x^2 + b x + c \rightarrow (式A)$

大小区別無効INS  
akai5.dat [前画面] level-1 (X001:00001)  
1 スペース行列処理に関する通常のGauss消去法の消去順序を変更することにより、発生する非常要素数に関して、ほぼ消去の全域にわたり、かなり改良されたアルゴリズムを提案した。  
2 ここで、ADRは主記憶上の番地を、COUNTはアタイミングが真にならための回数に関する数値、いわゆる通過回数指定を表わしている。  
3. 6 ステップ数に関するデバッグ手法。  
4 分割の手続きStep4の実行回数に関する、次のような帰納法で証明される。  
5 またもちろん、 $d = m(x)$ について $Q(G) - Q(H)$ であっても、式数41050となることがあるので、 $Q(G) - Q(H)$ を $G, H$ に対するハッシュ値と見なして、 $G, H$ に関するデータレコードは順に並べておくことにする。  
6 18 領域解法として、従来より種々の方法が提案され、最近では領域分割の段階においても認識対象に関する簡単な事前情報（セマンティックス）を利用することができる。  
19 20 21 22 本論文において、事前情報や種々の要因に依存する経験的な知識等を用いない領域解法を提議し、本方法が領域分割（物体識別）に関する人間の視覚にはほぼ忠実であることを実験によって確かめた。  
23 D F A M に新たに追加された機能は、（1）シーク時間、回転待ち時間とデータ転送時間、応答時間および転送データ長に関する接針情報の収集とヒストグラムの作成、（2）ディスク・スケジューリング方式に基づく入出力要求のスケジュールである。  
24 R B (Request Block) は1個の入出力要求に関するすべての情報を含む制御表である。  
25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 にもかかわらず、小さなプログラムに関するこれらの実験は、価値あるもので、あったばかりでなく、当時としては不可欠でした。  
37 38 39 40 41 42 43 44 その一つの理由は、この結びつきがくも緊密である間は、たまたま存在している計算機の特性が、好みしいものであるかどうかにかかわりなく、プログラミングに関する思考を支配してしまう。  
一方、第二の原因とは、Natureや FloydやHoareの、プログラムの正確さを形式的方法で証明する可能性に関する論文によってもたらされた発見！とその応用に在るものがありました。

図5 入力ファイル

#### 4. 特定表現の抽出方法

ここで、K W I C のプログラムの流れについて述べる(図4参照)。

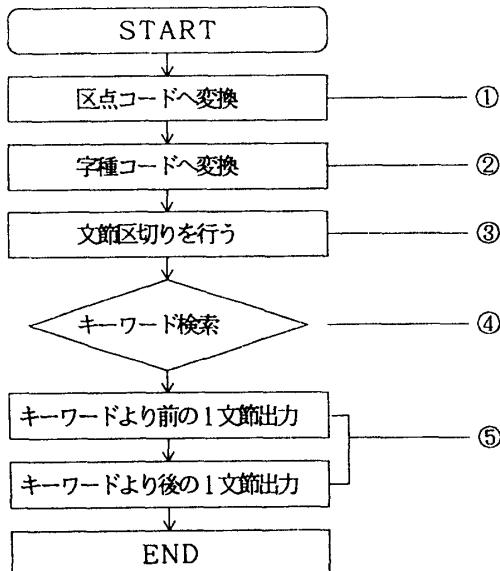


図6 KWICのフローチャート

2バイトコードのままで、文節区切りの判断が煩雑になるため、①の部分で全ての文字を区点コードに変換する。

次に②の部分で区点コードを見て、字種別にコード化する。

表1 区点コード別字種コード

	区点コード	字種コード
記号	0102~0390	6
ひらがな	0401~0483	0
カタカナ	0501~0586	1
漢字	1601~9404	2

ここで、我々が必要としていたのは、キーワードの前後一文節のみであるが、後の処理のことを考えて③の部分で予め、文全体を文節区切りしておく。

文節区切りの基準としては、日本語の文章は普通ひらがなから漢字に変わるので文節が区切られるので、それが出現したら文節として区切りとした。さらにその他にも文節が区切れると思われる場合(「私のコンピュータ」などは、ひらがなからカタカナに変わった時に文節が区切れる)についても文節区切りの条件とした。

私は犬と一緒に散歩へ行きました。  
 ↓ ②により  
 2 0 2 0 2 0 2 0 2 0 0 0 6  
 ↓ ③により  
 2 0 / 2 0 / 2 2 0 / 2 2 0 / 2 0 0 0 6  
 ↓ ③時点の出力  
 私は/犬と/一緒に/散歩へ/行きました。

図7 文節区切りの例

④の部分でキーワード(ここでは「～に関する」と「～に対する」)を本文中から探索する。

⑤の部分で、予め処理③で区切られた文章のうち前後一文節だけを出力する。

#### 5. 出力結果

下にK W I C のプログラムを実行した時の出力結果を掲げる

(キーワードは「～に関する…」)	
スパース行列処理	通常のGauss消去法の
回数	条件、ベッティング手法。
3, 6 ステップ数	～に関する～デバッギング
実行回数 i	～次のように
G, H 認識対象	～開する～データレコードは
(物体識別)	～開する～簡単な
伝送データ長	～開する～人間の
入出力要求	～開する～統計情報の
プログラム	～開する～すべて
プログラミング	～開する～ごれ考を
可能性正しさ	～開する～によってもたらされた
形式的証明	～開する～初期期の
構文法	～開する～形式理論の
証明進行	～開する～公理がありえられ
実仕事	～開する～のらら
品質審査	～開する～とをの
現製造業	～開する～私たちは
製造業	～開する～重要な
財形何学	～開する～ないし
信頼性	～開する～一つの
de que	～開する～所要の
既換行列	～開する～形式的手法を
マイクロプログラミング	～アルゴリズム(9)ムに
マイクロプログラム	～開する～情報書は
手法	～開する～成事をする
マイクロプロセッサ	～開する～仕
分散処理等	～開する～たる
使用経験等	～開する～から
生産技術	～開する～国際会議。
生産技術	～開する～一般論、
(5)製造技術	～開する～研究が
検査直接計測等	～開する～一般論は、
製造技術	～開する～各種の
主題分析	～開する～システムは、
パリック・ヘルス	～開する～システムのう
システム・スクリーニング	～開する～セッションのう
診断	～開する～もの
計算機アーキテクチャ	～開する～の討を
容量分析	～開する～研究が
自動解析	～開する～実際的手法について
情報処理	

図8 出力ファイル

#### 6. おわりに

今後の課題としては、今回のK W I C のプログラムでは入力形式に制限があったので、今後なくしていくことと、文節区切りの精度を上げることである。

また、次に発表する辞書を使用して辞書引きできるようにしたい。

#### 【参考文献】

- [1] 千早 耿一郎 「悪文の構造」, 木耳社
- [2] 情報処理学会第45回全国大会講演論文集
- [3] 田中 他 「日本文の自動分かち書き(ひら仮名文字列の分かち書き)」, 情報処理学会第20回全国大会講演論文集
- [4] Borland international, Inc. 「Turbo Pascal リファレンスガイド」, (株)マイクロソフトウェア アンシェイツ