

# 日本語全文情報の自動索引法

3G-5

岩淵 保 須田彰夫 中田安昭

丸善株式会社

## 1. はじめに

近年、ペーパーレスを合言葉に全文データベース・システムの構築が、盛んになってきた。構築時に問題になるのは、検索効率を高めるためのキーワード付与作業である。全文が既にコンピュータに採録されているので、その作業効率と標準化を考えると、自動索引システムを使用することになる。

実用化されている自動索引システムの問題点と、それを解決する為の方策と我々が実験したシステムの概要を報告する。

## 2. 問題点

実用化されている自動索引システムの基本的な問題点を2点上げる。

- (1) 大量のキーワードが抽出される。「1,000文字の文章から400文字(100ワード)程度のキーワードが抽出される。」<sup>1)</sup>との報告もある。従って、検索結果には、多くのゴミ情報が含まれてしまう。
- (2) 辞書類を多用するので、その質の良否がシステムの評価に影響を与える。良質な辞書の作成と、その維持をするために、メンテナンス作業が永続的に発生する。また、ある文章ではキーワードとして不要な語句も、他の文章では必要となる場合もあり得るので、上記の大量キーワード群を不要語辞書にて単純に消し込むことも問題である。

## 3. 解決策

我々が提案する解決策は、文章の主題を把握することで、辞書類の使用を極力排除しようとするものである。筆者らは、以前より自動抄録技術<sup>2)-4)</sup>を研究してきたが、その成果に基づき、自動索引法の実験システムを作成した。また、インデクサー(人間)は、文章中に存在しない語句をキーワードとして付与する。人間だからなせる技である。文章を抽象化することで似通った内容の文献を漏れなく検索することを可能にする。新聞記事では、15~20%、科学技術論文では、50~60%がこれらのキーワードである。我々は、これを想像キーワードと名付けて、その付与方法も実験した。

## 4. 自動索引システムの概要

今回、新聞記事を対象にして実験を試みた。図1は、その概要の流れ図である。なお、以降、KWは、キーワードの略号であり、単語とは、漢字・カタ仮名文字で構成される2文字以上の語句を言う。KWの抽出・生成ブロックで、より上位に位置するブロックでのKWの方が重要度が高い。以降、図1に準じて説明する。

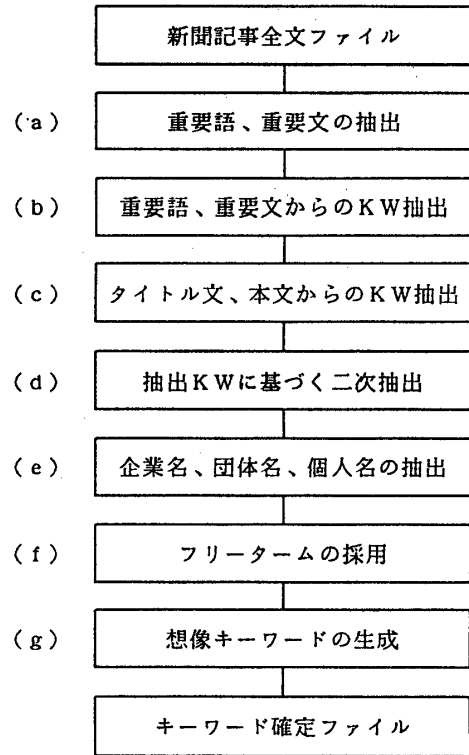


図1 自動索引システムの流れ図

- (a) 重要語、重要文の抽出  
新聞記事全文を読み込み、「は」と主題化を示す複合辞<sup>5)</sup>の前の単語(重要語、主題を構成する語句)を抽出する。  
次に、重要語の文章内での出現頻度をチェックして最重要語(最も頻度の多かった単語)を確定する。文章を文単位で切り出し、最重要語と重要語の含有をチェックして、重要文を(主題を述べている文)抽出する。
- (b) 重要語、重要文からのKW抽出  
重要語で、記事全文が1,300文字以下であるならば、その全てをKW候補単語とする。1,300文字を越えるときは、その出現頻度4回以上の単語をKW候補単語とする。  
重要文で、その先頭から2文目以内で、「が」、「で」、「は」の直前の単語をKW候補単語とする。但し、「では」は不採用とする。  
複合辞で仕手、仲介、根拠、原因、立場、視点、相関を示す語の直前の単語をKW候補単語とする。
- (c) タイトル文、本文からのKW抽出  
タイトル文(見出し文)に含まれる全ての単語をKW候補単語とする。また、漢字・カタ仮名文字種類が混在する単語は、文字種で切り分けて双方ともKW候補単語とする。

Automatic indexing method of Japanese language on full-text documents  
Tamotu IWABUCHI, Akio SUDA, Yasuaki NAKATA  
Maruzen Co., Ltd.

本文に関しては、括弧内の単語とその直前の単語をKW候補単語とする。但し、括弧内に「本社」、役職名がある時、括弧内のみ不採用とする。

「に」の直前の単語で2回以上出現の単語をKW候補単語とする。

「聞き手」があるときは「インタビュー」をKW候補単語とする。このような、誰でもが納得できる置き換え単語は、置き換えファイルに登録して置く。現在は、この一件だけである。

複合辞で確定、同格を示す直前の単語をKW候補単語とする。

「は」と「が」にはさまれた単語をKW候補単語とする。

同一文字種で3文字以上の「と」の前と後ろの単語をKW候補単語とする。

鍵括弧内の単語をKW候補語とする。文字種類が混在している単語であれば、文字種で切り分けた単語もKW候補語とする。

(a)～(c)で採用されたKW候補語をユニーク(一意)にする。

(d) 抽出KWに基づく二次抽出

インデクサーのKW付与結果を分析すると、主要KWをまず抽出し、次に文章中でそのKWの周辺から二次レベルのKWを抽出している。この癖を機械的に具現する。

全文中でKW候補語の直後の「が」に続く単語も、KW候補語とする。

全文中でキーワード候補語の直前に「と」、「や」、「、」、「・」、「の」があるとき、その直前の単語もKW候補語とする。

(a)～(d)で採用されたKW候補語をユニークにする。

(e) 企業名、団体名、個人名の抽出

(株)エレクトロニック・ライブラリで制作したKW集ファイルを利用する。このファイルは、ELINETを利用してユーザに提供されるKW集の機械可読版であり、検索頻度の高いKW(コントロールタム、統制語)と、企業名、団体名、個人名が収録されている。KW集は、1年単位で編集し直すとのことで、KW集ファイル(統制語辞書ファイル)は、副次的に生成される。KW集ファイルにて、文章中の企業名、団体名、個人名をKW候補語とする。

(f) フリータムの採用

(a)～(e)で採用されたKW候補語と、KW集ファイルを突合し、KW集ファイルにある語句は、コントロールタムKWとして確定する。存在しない語句は、ディスプレイ画面に表示して、インデクサーがフリータムKWとしての採用を選択する。

(g) 想像キーワードの生成

インデクサーが想像KWを付与する過程は、最も近い過去に生じた、内容類似の文章に付けたKWを想起し、その中から代表的なKWを付与していると考えられる。そこで、重要語から採用したKW(最も重要度の高いKW)を2個以上持つ文章を過去1年間の情報から検索し、その文章に付与されたKWのみを取り出す。次に、KWの出現頻度を取り、一定の頻度以上のKWを想像KWとして採用する。

## 5. 自動索引システムの評価

以降、コントロールタム=CT、フリータム=FT、想像KW=想像と略号で記す。

(例題) 日本経済新聞 90年3月12日朝刊 核心各論 =世界初の宇宙テーマパーク スペースワールド 社長 小池孜氏 地域活性化の先兵役に

(インデクサー付与)

リストラクチャリング(CT)・社長(CT)・観光開発(CT)・インタビュー(CT)・地域活性化(CT)・北九州市(CT)・新日本製鉄(CT)・日本興業銀行(CT)・福岡銀行(CT)・東京急行電鉄(CT)・JR九州(CT)・スペースワールド(CT)・テーマパーク(FT)・核心各論(FT)・小池孜(FT)・レジャーランド(想像)・福岡(想像)

(機械付与)

上記、CT、FTは全て付与される。想像KWでは、「福岡」は付与されるが、「レジャーランド」は、「レジャー施設」が採用される。また、上記以外に「宇宙」・「下関市」・「地域」が余分に切り出されてしまう。今後の改善点である。使用機器は、PC9801で行い、C言語で構築した。処理時間は、平均1分程度である。

## 6. おわりに

新聞記事を対象にして、全文情報自動索引技術の実験システムを報告した。実験で使用した全文情報は、10件である。今後、1日分の全件数で実証をし、改善していく予定である。また、科学技術文献への適用も実験してみたい。最後に、(株)エレクトロニック・ライブラリー佐藤取締役情報部長初めインデクサーの皆様にも多大なるご協力を賜りましたことを深謝いたします。

## 【参考文献】

- 1) 重要キーワード抽出方式とその活用方法 内山、中村 (東京電力システム研究所) 情報処理学会 データベース・システム研究会報告84-19 (1991. 7. 18)
- 2) 全文情報からの意味的情報の抽出と加工 岩淵、荒井 藍沢 (テレマティーク国際研究所) 情報処理学会第38回全国大会
- 3) 自動抄録法 岩淵、荒井、藍沢 (テレマティーク国際研究所) 電気学会通信研究会CMN-89-23 (1989. 7. 12)
- 4) 全文情報からの意味的情報の抽出と加工 自動抄録法「TELEMA-A」について 岩淵、荒井 藍沢 (テレマティーク国際研究所) 第26回情報科学技術研究集会発表論文
- 5) 日本語表現文型 用例中心・複合辞の意味と用法 森田良行、松木正恵 アルク
- 6) キーワード集ファイル エレクトロニック・ライブラリー ELINET
- 7) 新聞記事データベースにおけるキーワード自動抽出 神尾 (日本経済新聞社) 情報管理 Vol. 32 NO. 4 July 1989
- 8) 日本語の文法(上)、(下) 国立国語研究所 大蔵省印刷局
- 9) 自然言語の文法理論 郡司隆男 産業図書
- 10) 教師のための口語文法 渡辺正数 右文書院
- 11) 日本語概説 加藤彰彦、佐治圭三、森田良行 桜楓社
- 12) いわゆる日本語助詞の研究 奥津敬一郎、沼田善子、杉本武 凡人社