

## 専門用語の自動収集システムについて

吉村賢治・山下明男・日高 達・吉田 将  
(九州大学 工学部)

### 1. まえがき

科学技術文献の機械翻訳や機械支援翻訳に関する研究の進展に伴い、専門用語辞書の開発が望まれている。しかし、日本語における専門用語の大部分は複合語と外来語で、その数が膨大であること、新しい用語の出現数が多いことなどから、専門用語を手だけで収集することは困難であり、専門用語収集の機械化が必要である。筆者らは、現在、科学技術文献の表題と抄録文から専門用語を収集するシステムを開発している。この専門用語の自動収集システムは、次のようなサブ・システムで構成されるが、本稿では実験システムが完成している専門用語の自動抽出について報告する。

- (1) 専門用語の自動抽出
- (2) 専門用語辞書の作成・管理
- (3) シソーラスの作成・管理

専門用語を表題、抄録文等から自動収集するための手法としては、不要語(stop word)を除去する方法が一般的であり、実用性も高いと考えられる。不要語除去の具体的方法として次の二つの方法がある。

- (1) 不要語テーブルを用いて文字列のパターン・マッチングを行なう方法<sup>①</sup>。
- (2) 単語辞書と文法規則のテーブルを用いて文法解析を行なう方法。

方法(1)は文法解析を行わないため処理は簡単であるが、不要語テーブルを作成するために大量のデータを調査しなければならない。また、不要語テーブルとのマッチングだけで不要語を除去した場合、平仮名一文字で表わされる助詞などが原因で誤った切断が生じる。これを避けて専門用語抽出の精度を上げるためには、不要語テーブルとのマッチングで生じる切断の誤りを抑制するために切断禁止規則のテーブル等が必要となり、テーブル間の干渉の管理問題などが発生する。

一方、方法(2)には上記のような問題はないが、文法解析を行なうために処理が複雑になり、大規模な自立語辞書を使用した場合、処理時間も長くなる。また、一般に大部分の専門用語はシステムの単語辞書に登録されていないため、未登録語を含んだ日本語文を処理できる文法解析システムが必要である。

本稿で報告する専門用語の自動収集システムは日本語文の形態素解析システム<sup>②</sup>を利用している。この形態素解析システムは未登録語を含む日本語文の解析が可能であり、この特徴を利用して上に述べた大規模な自立語辞書の検索に要する時間の削減を図っている。

専門用語の自動抽出システムは図-1に示すモジュール

ルから構成されている。以下、図-1の各モジュールについて説明し、日本科学技術情報センターの科学技術文献速報・電気工学編VOL.25,NO.6から抽出した表題及び抄録文に対して行なった実験とその結果について述べる。なお、ここで抽出の対象としている専門用語の品詞は名詞だけである。

### 2. 形態素解析

#### 2.1 形態素解析の手順

このシステムでは、参考文献(2)で報告した形態素解析アルゴリズムを用いている。本章ではアルゴリズムの概略、及び一般の形態素解析システムとの相違点について述べる。

この形態素解析アルゴリズムでは、まず入力文を文頭から文末に走査し、入力文を構成しているすべての解析単位(入力文が分解される解析の単位。一般には単語であるが、以下では解析単位と呼ぶ。)について、5項

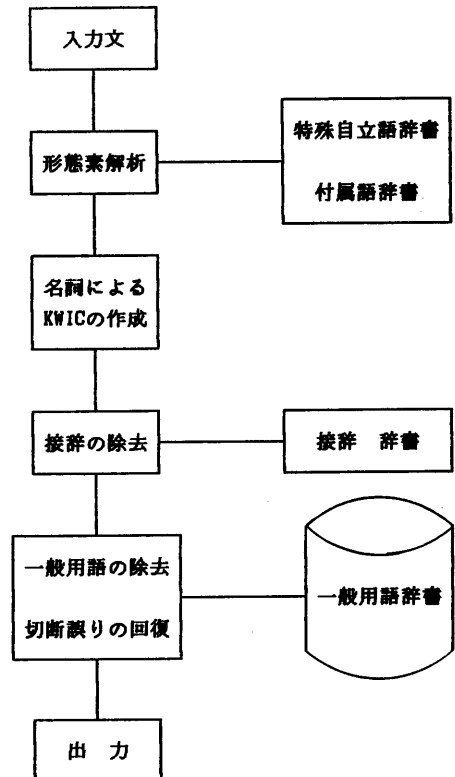


図-1 システムの構成

組  $(i, j, \alpha, k, a)$  を作成する。この5項組をアイテムと呼び、一つの入力文に対して作成されるアイテムの集合をパーズリストと呼ぶ。ここで  $i, j$  は入力文における解析単位の始点と終点、 $\alpha$  は解析単位の文法情報である。専門用語の自動収集システムにおける解析単位はその素性によって表-1に示す10種類に分類される。 $a$  はその種類を示し、表-1の右の記号を値として持つ。また、解析単位にはその種類に応じてコストが設定されており、 $k$  は文頭からその解析単位までのコストの総和である。このコストの総和は連続する解析単位が文法的な接続規則をみたしている列についてのみ求める。

表-1において、(vi) ~ (viii) が未登録語に対処するための解析単位であり、漢字仮名混り語の未登録語は (vi) と (vii) で、外来語の未登録語は (vii) で処理される。また、連結記号とは文字列を連結して専門用語を構成するハイフン等の記号である。

ここで、解析単位の種類を表わす記号を添字にして各解析単位のコストを  $C_J, C_N, C_F, C_G, C_H, C_K, C_R, C_M, C_C, C_S$  で表わす。筆者らは参考文献(3)でもっともらしい解析結果から抽出するためのヒューリスティックスとして文節数最小法を提案した。解析結果のコストの総和が小さいものから出力するとき、このヒューリスティックスに従うために

$$C_J, C_N, C_H > C_F, C_G$$

とする。また、ほとんどの表記は正書法に従っており、本来の未登録語は稀にしか出現しないという仮定から、

$$C_H > C_J, C_N$$

$$C_K, C_R, C_M > C_J, C_N$$

とする。以上のことと字種の違いを考慮して、各解析単位のコストは次の関係を満たすように設定する。

$$C_M > C_K, C_R > C_H > C_J, C_N > C_F, C_G$$

なお、本稿で報告する実験は、

$$C_C, C_S = 0$$

$$C_F = 1$$

$$C_G = 2$$

	解析単位	記号
(i)	正書表記の自立語	J
(ii)	数詞	N
(iii)	付属語	F
(iv)	活用語尾	G
(v)	非正書表記の自立語	H
(vi)	漢字	K
(vii)	アルファベット・片仮名列	R
(viii)	平仮名	M
(ix)	連結記号	C
(x)	記号	S

表-1 解析単位の種類

$$C_J, C_N = 4$$

$$C_H, C_K, C_R = 6$$

$$C_M = 14$$

で行なったものである。

形態素解析システムでは、作成されたパーズリストを文末側から文頭側に走査して、入力文の解析結果を抽出するが、専門用語の自動収集システムでは3章で述べる手順に従ってパーズリストから名詞(列)をキーとしたKWICレコードを作成する。

## 2.2 特殊自立語辞書

参考文献(2)の形態素解析システムでは単語数約83,000、見出し語数約170,000の自立語辞書<sup>4)</sup>を使用しているが、1章で述べたように文法解析を行なう専門用語の自動収集システムでは二次記憶上に置かれた自立語辞書の検索時間が問題になる。そこで、この専門用語の自動収集システムにおける形態素解析では、上記の辞書の代わりに主記憶上に常駐することができる小規模の自立語辞書(これを特殊自立語辞書と呼ぶ)を使用している。

表-2に示すように、特殊自立語辞書には、連体詞、接続詞、副詞、形容詞など専門用語やその構成要素でない814単語が登録されている。なお、特殊自立語辞書の作成にあたっては首藤のデータ<sup>5)</sup>を参考にしており、多くの慣用句的な表現を単語として登録している。

この特殊自立語辞書中の単語は、切断禁止文字列を除いて不要語であり、形態素解析の結果、専門用語を含む一般の自立語は未登録語として出力される。ここで切断禁止文字列とは、副詞“絶対”によって切断されては不都合である“絶対値”のような文字列で、形態素解析では漢字(K)と同様に扱う。

## 3. 名詞によるKWICの作成

パーズリストからはコストの総和が最小となる解析結果のみを抽出する。この解析結果抽出の過程で名詞の可能性のある未登録語をキーとするKWICレコードを作成する。いま、解析単位の種類を示す記号でそれぞれの解析単位を表わすと、KWICレコードのキーになる解析単位の列Xは次のように記述される。

品詞	単語数
副詞	335
連体詞	209
接続詞	160
名詞	67
形容詞	35
切断禁止文字列	8

表-2 特殊自立語辞書の内容

$X := \{K, M, R\}$   
 $X := X \cdot X$   
 $X := X \cdot C \cdot X$   
 $X := N \cdot C \cdot X$   
 $X := N \cdot X$   
 $X := X \cdot G^M$

ここで、括弧 “(”、)”” は、その中の一つを選択することを表わし、中点 “.” は連結を表わす。また、 $G^M$  は5段動詞の連用形の活用語尾である。

KWICレコードの作成は文献単位に行なう。一つの文献のKWICレコードは、表題から作成されたものと抄録文から作成されたものに分け、4章で述べる接辞の除去、5章で述べる一般用語の除去と切断誤りの回復を行なった後、それぞれをキーの部分でソートして表示する。KWICレコードを作成した段階の出力例を図-2に示す。

この段階で抽出されたものは特殊自立語辞書に登録している名詞以外の名詞(列)である。ここでは、日本科学技術情報センター発行の科学技術文献速報・電気工学編 Vol.25, No.6の先頭から取り出した100文献に対して行なった抽出実験の結果を示す。

対象：日本科学技術情報センター・科学技術文献速報  
 電気工学編 Vol.25, No.6  
 文献番号 E82060001~E82060100

総文数 . . . . . 428文  
 名詞の総数 a . . . . . 2339個  
 抽出できなかった名詞の個数 b . . . . . 66個  
 誤って抽出した文字列の個数 c . . . . . 124個

いま、抽出されるべきものの個数をa、抽出されるべきものの中で抽出されなかったものの個数をb、抽出されるべきもの以外で抽出されたものの個数をcとして、

$$\text{抽出率} = (a - b) / a$$

$$\text{抽出の精度} = (a - b) / (a - b + c)$$

と定義すると、KWICレコードの作成段階における名詞の抽出率は97.2%、抽出の精度は94.8%である。

抽出に失敗した名詞は殆どのものが、図-3に示すような漢字と平仮名で混ざ書きされたものである。また、誤って抽出されたcに相当するものは、殆どが形容動詞の語幹であった。

2つのはん関数の順序積の平均値。

Kuboの順序指数 (J. Phys. Soc. Japan, 1962, 17) の自然な一般化である順序付けられたはん関数を導入。交換する確率過程の2つのはん関数の積の平均に関する周知の定理 (特にFurutsu - Novikovの定理) を非交換確率過程の場合に一般化。閉じた量子系のじょう乱解析での応用を示す。

応答の高次平均が順序キュミュラントで表現され、物理系のGauss特性を最大限利用できる。

FROM	TITLE
2つの 2つのはん関数の 2つのはん関数の順序積の	はん関数 順序積 平均値
FROM	ABSTRACT
Kuboの Kuboの順序指数 ( c. Japan, 1962, 17) の自然な 17) の自然な一般化である順序付けられた 然な一般化である順序付けられたはん関数を 交換する 交換する確率過程の2つの 交換する確率過程の2つのはん関数の 率過程の2つのはん関数の積の平均に関する の2つのはん関数の積の平均に関する周知の はん関数の積の平均に関する周知の定理 (特に 理 (特にFurutsu - Novikovの Furutsu - Novikovの定理) を非交換確率過程の場合に 閉じた 閉じた量子系のじょう 閉じた量子系のじょう乱解析での	Kubo 順序指数 J. Phys. Soc. Japan 一般化 はん関数 導入 確率過程 はん関数 積 平均 周知 定理 Furutsu - Novikov 定理 非交換確率過程 一般化 量子系 乱解析 応用 応答 高次平均 順序キュミュラント 物理系 Gauss特性 利用
応答の 応答の高次平均が 高次平均が順序キュミュラントで表現され、 が順序キュミュラントで表現され、物理系の 表現され、物理系のGauss特性を最大限	

の順序積の平均値。  
の平均値。  
。

の順序指数 (J. Phys. Soc. Jap  
(J. Phys. Soc. Japan, 19  
, 1962, 17) の自然な一般化である順  
である順序付けられたはん関数を導入。  
を導入。

の2つのはん関数の積の平均に関する周知の  
積の平均に関する周知の定理 (特にFur  
の平均に関する周知の定理 (特にFur  
に関する周知の定理 (特にFurutsu -  
の定理 (特にFurutsu - Noviko  
(特にFurutsu - Novikovの定  
の定理) を非交換確率過程の場合に一般化。  
) を非交換確率過程の場合に一般化。  
の場合に一般化。

のじょう乱解析での応用を示す。  
での応用を示す。  
を示す。  
の高次平均が順序キュミュラントで表現され  
が順序キュミュラントで表現され、物理系  
で表現され、物理系のGauss特性を最大  
のGauss特性を最大限利用できる。  
を最大限利用できる。  
できる。

図-2 名詞によるKWICの作成結果

#### 4. 接辞の除去

2章で述べた形態素解析では接辞処理を行っていないため、不要な接辞が連続した名詞は、そのままの形でKWICレコードのキーになっている。そこで、表-3に示す接頭語辞書、接尾語辞書を用いて、キーの先頭にある不要な接頭語を前の文脈文字列へ、キーの末尾にある不要な接尾語を後ろの文脈文字列へ各々移動する。このとき、名詞の一部が接辞と一致してキーから取り除かれることもあるが、この誤りの殆どは次のステップで回復される。図-2に示したKWICレコードのキーから接辞を除去した後の出力例を図-4に示す。

#### 5. 一般用語の除去と切断誤りの回復

このステップまでに作成されたKWICレコードの多くは一般用語がキーになっている。このステップでは一般用語辞書を用いて、このようなKWICレコードの削除を行なう。また、形態素解析の段階で生じた専門用語の切断誤りの回復も同時に行なう。

現在、一般用語辞書としては、参考文献(2)の形態素解析システムで使用している自立語辞書から正書表記の名詞だけを抽出して構成したものを使用している。この一般用語辞書の見出し語数は79,102個で、処理の都合により、文頭側からの検索に用いる正引き辞書と文末側からの検索に用いる逆引き辞書を用意している。辞書のデ

磁気しゃ▲へ▲い  
二乗平均平方根の広▲が▲り  
じょう▲乱解析  
不変はめ▲込み理論

図-3 抽出を誤った名詞の例

	接頭語	接尾語
1	各	やすい
2	諸	下
3	当該	以外
4	本	間
5	両	時
6		上
7		側
8		全体
9		中
10		等
11		内
12		後半
13		分
14		用

表-3 接辞辞書

-----	FROM	TITLE	-----
	2つの	はん関数	
	2つのはん関数の	順序積	
	2つのはん関数の順序積の	平均値	
-----	FROM	ABSTRACT	-----
	Kuboの	Kubo	
	Kuboの順序指数 (	順序指数	
	c. Japan, 1962, 17) の自然な	J. Phys. Soc. Japan	
	17) の自然な一般化である順序付けられた	一般化	
	然な一般化である順序付けられたはん関数を	はん関数	
	交換する	導入	
	交換する確率過程の2つの	確率過程	
	交換する確率過程の2つのはん関数の	はん関数	
	積	積	
	交換する確率過程の2つのはん関数の積の	平均	
	率過程の2つのはん関数の積の平均に関する	周知	
	の2つのはん関数の積の平均に関する周知の	定理	
	はん関数の積の平均に関する周知の定理 (特に	Furutsu - Novikov	
	理 (特にFurutsu - Novikovの	定理	
	Furutsu - Novikovの定理) を	非交換確率過程	
	ikovの定理) を非交換確率過程の場合に	一般化	
	閉じた	量子系	
	閉じた量子系のじょう	乱解析	
	閉じた量子系のじょう乱解析での	応	
	応答	応答	
	高次平均	高次平均	
	順序キュミュラント	順序キュミュラント	
	物理系	物理系	
	Gauss特性	Gauss特性	
	利	利	

の順序積の平均値。  
の平均値。

の順序指数 (J. Phys. Soc. Jap (J. Phys. Soc. Japan, 19, 1962, 17) の自然な一般化である順序付けられたはん関数を導入。

の2つのはん関数の積の平均に関する周知の積の平均に関する周知の定理 (特にFurutの平均に関する周知の定理 (特にFurutに関する周知の定理 (特にFurutsuの定理 (特にFurutsu - Noviko (特にFurutsu - Novikovの定の定理) を非交換確率過程の場合に一般化。 ) を非交換確率過程の場合に一般化。の場合に一般化。

のじょう乱解析での応用を示す。での応用を示す。用を示す。の高次平均が順序キュミュラントで表現されが順序キュミュラントで表現され、物理系ので表現され、物理系のGauss特性を最大のGauss特性を最大限利用できる。を最大限利用できる。用でできる。

図-4 接辞除去の結果

ータ構造としては拡張B-treeとTRIEを併用しており<sup>4)</sup>、容量は正引き辞書が1,121Kbyte、逆引き辞書が1,178Kbyteである。

### 5.1 アルゴリズム

一般用語の除去と切断誤り回復のアルゴリズムを簡潔に記述するために、ここで幾つか記述上の約束と関数の定義を行なう。

- (i) 文字列 $w$ の長さ(文字数)を $|w|$ で表わす。
- (ii) 二つの文字列 $v, w$ の連結を $v \cdot w$ で表わす。
- (iii) 文字列 $v$ が文字列 $w$ の最左部分列(prefix)であることを $v \sqsubseteq w$ で表わす。
- (iv) 正引き一般用語辞書の見出し文字列を要素とする集合を $D_L$ で表わし、逆引き一般用語辞書の見出し文字列を要素とする集合を $D_R$ で表わす。
- (v) 文字列 $s$ と文字列の集合 $D$ に対して、 $w \sqsubseteq s$ と $w \in D$ を満たす総ての $|w|$ における最大値を $\text{Max}(s, D)$ で表わす。
- (vi) 文字列 $s$ における右側から長さ $i$ の部分文字列を $\text{Right}(s, i)$ で表わし、左側から長さ $i$ の部分文字列を $\text{Left}(s, i)$ で表わす。
- (vii) 文字列 $s$ の並びを左右逆にしてできる文字列を $\text{Reverse}(s)$ で表わす。

ここで、空文字列 $\epsilon$  ( $|\epsilon| = 0$ )は任意の文字列 $s$ と文字列の集合 $D$ について、 $\epsilon \sqsubseteq s$ 、 $\epsilon \in D$ である。以上の記号と関数を用いて、一つのKWICレコードに対して一般用語の除去と切断誤りの回復を行なうアルゴリズムを右に示す。

右に示すアルゴリズムを総てのKWICレコードに対し適用して、一般用語の除去と切断誤りの回復を行なう。なお、この処理に先立ち、キーの末尾が”的”であるKWICレコードは総て除去しておく。また、アルゴリズムの記述中、KWICレコードにおける前の文脈文字列を $C_L$ で、後の文脈文字列を $C_R$ で、キーの文字列を $Key$ で表わす。

図-4に示したKWICレコードに対して、一般用語の除去と切断誤りの回復を行ない、キーで昇順にソートした結果を図-5に示す。図-5において”【”と”】”で囲まれたキーは、一般用語として除去されたキーであり、キーの文字列中にスペースを含むキーは、その位置で切断誤りの回復が行なわれたキーである。

一般用語の除去と切断誤りの回復を行なった時点で、専門用語の抽出を完了する。3章で述べた名詞の抽出実験と同一の入力データに対して行なった、専門用語の抽出実験の結果を次に示す。

専門用語の総数	a	.....	1403個
抽出できなかった専門用語の個数	b	....	6個
誤って抽出した文字列の個数	c	.....	23個
部分的に抽出した専門用語の個数	d	....	9個
余分な文字列を伴って抽出した専門用語の個数	e	....	9個

[一般用語の除去と切断誤り回復のアルゴリズム]

#### Step. 1 (右向き走査の初期設定)

$s \leftarrow Key \cdot C_R$ .  
 $p \leftarrow 0$ .

$p \geq |Key|$  になるまで Step. 2 を繰り返す。

#### Step. 2 (右向き走査)

$\text{Max}(s, D_L) = 0$  ならば  $p \leftarrow p + 1$ .  
 $\text{Max}(s, D_L) \neq 0$  ならば  
 $p \leftarrow p + \text{Max}(s, D_L)$ .  
 $s \leftarrow \text{Right}(s, |Key \cdot C_R| - p)$ .

#### Step. 3 (切断誤りの回復)

$p > |Key|$  ならば  
 $Key \leftarrow \text{Left}(Key \cdot C_R, p)$   
 $C_R \leftarrow s$ .

#### Step. 4 (左向き走査の初期設定)

$s \leftarrow \text{Reverse}(C_L \cdot Key)$ .  
 $p \leftarrow 0$ .  
 $Count \leftarrow 0$ .

$p \geq |Key|$  になるまで Step. 5 を繰り返す。

#### Step. 5 (左向き走査1)

$\text{Max}(s, D_R) = 0$  ならば  $p \leftarrow p + 1$ .  
 $\text{Max}(s, D_R) \neq 0$  ならば  
 $p \leftarrow p + \text{Max}(s, D_R)$ .  
 $s \leftarrow \text{Right}(s, |C_L \cdot Key| - p)$ .  
 $Count \leftarrow Count + 1$ .

$\text{Max}(s, D_R) \neq 0$  のあいだ Step. 6 を繰り返す。

#### Step. 6 (左向き走査2)

$p \leftarrow p + \text{Max}(s, D_R)$ .  
 $s \leftarrow \text{Right}(s, |C_L \cdot Key| - p)$ .  
 $Count \leftarrow Count + 1$ .

#### Step. 7 (切断誤りの回復)

$p > |Key|$  ならば  
 $Key \leftarrow \text{Right}(C_L \cdot Key, p)$   
 $C_R \leftarrow \text{Reverse}(s)$ .

#### Step. 8 (一般用語の除去)

$Count = 1$  ならば、KWICレコードを除去する。

[アルゴリズム終わり]

ここで、

$$\text{抽出率} = (a - b - d - e) / a$$

$$\text{抽出の精度} = (a - b - d - e) / (a - b + c)$$

とすると、専門用語の抽出率は 98.3%、専門用語抽出の精度は 97.1%である。専門用語抽出の精度が100%でないため、この段階で人間による用語の選択作業が必要になる。このとき、KWICレコードのキーとして部分的にでも抽出されている専門用語は、人間が介入することにより正しく抽出することができる。このようなKWICレコードも正しく抽出されたものとして扱おうと、

$$\text{抽出率} = (a - b) / a$$

$$\text{抽出の精度} = (a - b) / (a - b + c)$$

となり、抽出率は99.6%、抽出の精度は 98.4%になる。

全く抽出できなかった専門用語は次のものである。

- (1) 開放共振器近くの
- (2) 斜め入射の
- (3) 飛しよう
- (4) 部分的線形近似がえられるが
- (5) 飽和蒸気圧にて
- (6) 巨大電子なだれを

これらのうち、(5) を抽出できなかった原因は付属語辞書の不備である。

本稿で報告した実験システムは、九州大学大型計算機センターの FACOM M382 上にPL/Iを用いて実現した。この環境で、一文献の表題と抄録文を処理するのに要する時間の平均は1439.3msecであった。一文献の表題と抄録文の平均文字数は 174.4文字である。参考に図-5に示した結果を求めるのに要した時間は1585msecである。この内、右向きの走査に要した時間は 696msecで、左向きの走査に要した時間は 779msecであった。

## 6. あとがき

本稿では形態素解析システムを利用した専門用語の自動抽出システムについて報告した。専門用語抽出の精度については、特殊自立語辞書、一般用語辞書のデータを整備することにより、さらに高めることが可能である。

大規模な自立語辞書を用いた形態素解析を利用して専門用語を抽出する場合に対する、抽出の精度及び処理時間の比較はまだ行っていないが、処理時間は殆ど改善されていないようである。この原因は切断誤りの回復処理において、二重の辞書引きを行なうことにある。5章のアルゴリズムを、

$$s \leftarrow \text{Key} \cdot \text{CR}.$$

$$\text{Max}(s, \text{DL}) \geq |\text{Key}| \text{ならば}$$

KWICレコードを除去する。

<p>FROM TITLE</p> <p>2つのはん関数の順序積</p> <p>2つのはん関数の順序積の【平均値】</p>	<p>はん関数</p> <p>順序積</p> <p>【平均値】</p>	<p>の順序積の平均値。</p> <p>の平均値。</p> <p>。</p>
<p>FROM ABSTRACT</p> <p>はん関数の積の平均に関する周知の定理 (特にが順序キュミュラントで表現され、物理系の Kubo の順序指数 (閉じた量子系の交換する確率過程の2つの17) の自然な一般化である順序付けられた ikov の定理) を非交換確率過程の場合に c. Japan, 1962, 17) の自然な閉じた量子系のじょう乱解析での交換する確率過程の2つの応答の高次平均の順序積の平均に関する Kubo の交換する確率過程の2つのはん関数の積の平均に関する周知の理 (特に Furutsu - Novikov の然な一般化である順序付けられたはん関数を Furutsu - Novikov の定理) を高次平均が順序キュミュラントで表現され、交換する確率過程の2つのはん関数の積の平均に関する周知の理 (特に Furutsu - Novikov の定) を非交換確率過程の場合に一般化。の場合に一般化。の Gauss 特性を最大限利用できる。に関する周知の定理 (特に Furutsu - のじょう乱解析での応用を示す。</p>	<p>Furutsu - Novikov Gauss 特性 J. Phys. Soc. Japan Kubo じょう乱解析はん関数はん関数【一般化】【一般化】【応答】【応用】確率過程高次平均【周知】順序キュミュラント順序指数積【定理】【定理】【導入】非交換確率過程【物理系】【平均】【利用】量子系</p>	<p>の定理) を非交換確率過程の場合に一般化。を最大限利用できる。1962, 17) の自然な一般化である順の順序指数 (J. Phys. Soc. Japでの応用を示す。の積の平均に関する周知の定理 (特に Fur の平均に関する周知の定理 (特に Furut を導入。である順序付けられたはん関数を導入。の高次平均が順序キュミュラントで表現されを示す。の2つのはん関数の積の平均に関する周知のが順序キュミュラントで表現され、物理系のの定理 (特に Furutsu - Novikoで表現され、物理系の Gauss 特性を最大 (J. Phys. Soc. Japan, 19の平均に関する周知の定理 (特に Furut (特に Furutsu - Novikov の定) を非交換確率過程の場合に一般化。の場合に一般化。の Gauss 特性を最大限利用できる。に関する周知の定理 (特に Furutsu - できる。のじょう乱解析での応用を示す。</p>

図-5 一般用語の除去と切断誤り回復の結果

と変更して、切断誤りの回復処理を行なわない場合、処理時間は100文献で平均して一文献あたり601.9msecになるが、部分的に抽出した専門用語の個数は43個になり、

抽出率  $(a-b-d-e)/a = 95.9\%$

抽出の精度  $(a-b-d-e)/(a-b+c) = 94.7\%$

となる。従って、切断誤りの回復処理による抽出率の改善率は  $(4.1-1.7)/4.1$  より58.5%である。

本研究は昭和58年度文部省科学研究費特定研究「言語の標準化」によった。

## 7. 参考文献

- (1) 荒木啓介・金子明夫・高野文雄・日夏健一：“日本語論文タイトルからのキーワード自動抽出システム(JAKAS)”，情報処理学会自然言語処理研究会資料26-3、1981。
- (2) 吉村賢治・日高達・吉田将：“未登録語を含む日本語文の形態素解析アルゴリズム”，九州大学工学集報、第55巻、第6号、1982。
- (3) 吉村賢治・日高達・吉田将：“文節数最小法を用いたべた書き日本語文の形態素解析”，情報処理学会論文誌、第24巻、第1号、1983。
- (4) 吉田将・日高達・稲永敏之・田中武美・吉村賢治：“公用データベース日本語単語辞書の使用について”，九州大学大型計算機センター広報、第16巻、第4号、1983。
- (5) 首藤公昭：“文節構造モデルによる日本語の機械処理に関する研究”，福岡大学研究所報、第45号、1980。
- (6) 日高達・吉田将・稲永敏之：“拡張B-treeによる日本語単語辞書の作成”，情報処理学会自然言語処理研究会資料33-8、1982。