

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)

○荒木啓介, 金子明夫, 高野文雄, 日夏健一
(日本科学技術情報センター)

1. はじめに

文献のインテグリングの補助機能として、日本語のタイトルおよび/または抄録文からのキーワード自動抽出は重要な課題であり、相当な精度のものが実現すれば大幅な省力化になるが、技術的にはおつかしい問題があった。^(文献2)つまり

- ①. 連続した漢字が混り日本語を適切な単位に切断する。
- ②. 切断された文字列がキーワードとして重要か否かを判定する。
- ③. 取り出されたキーワードにフリガナを与える。

の三点である。

③は、著者ら外既に開発し、JICSTにおいて実用化している漢字-カナ変換システム(K-KACS)により、99.94%の精度が保証されている。

②は、タイトルについては、英語文などで確立した手法である。不要語(of, with, that など)を除いた、残りをキーワードとする方法により容易であるが、長い抄録文については、構文解析や意味解析を要するため、英語文等についても解決されていない。

①は日本語特有の問題であり、今までいくつかの研究があるが完璧とは行かないようである。^(文献3)

以上の状況から、当面論文タイトルを対象に、日本語文の効率良い切断方式を検討した。

著者らは既に、カナに変換された日本語文が、適度に分かり書きされている必要性から、K-KACSに分かり書き機能を与えていたが、これだけでは適切なキーワードを取り出すのに十分ではない(精度78%, 文献4)

そこで、日本語文中における漢字の役割に注目し、以下のアルゴリズムを考察し、それに基づき9万タイトル以上の論文タイトルを人手で分析したから、一見バカバカしい作業を繰り返して用語を収集し、辞書を作成した。並行してプログラムを作成し、未知のタイトルについてテストしたところ、予想を上回る性能が得られた。

2. 切断アルゴリズム

2.1. Cut 処理 : 日本語文の漢字列のうち、接頭辞、接尾辞、接中辞になることのある文字(文字列)を、辞書により切断する。

(例) 水溶液中触媒反応 → 水溶液中 触媒反 応

2.2. Pass 処理 : 切断対象となる文字(文字列)も、意味のある熟語を作るものについては、その網羅的な辞書により切断させない。

(例) 中央 . P. 中毒 . P. 卒中 . P. 死

以上の処理は、特に接辞のみに適用するのではなく、接辞を含めた用語の長い列についても適用される。その場合は、用語のタイポとしてその数少ないものを主として辞書化し、Cut 処理し、全体をハローハローにいく。

(例) カブトムシ幼虫ミトコンドリアATPアーゼ分解性物質

↑ 無限にあるタイポの用語 ↑ 有限なタイポの用語 ↑ 無限にあるタイポの用語

2.3. イタリック文字のプログラムによる切断 : 動植物名などは、イタリックリフト記号に囲まれているので、これを単にスペースに置き換えて切断する。半スペースは入すする。

(例) *Escherichia coli*

↑ ↑
半スペース

2.4. Quench 処理 : 通常の単純な切断では、切口に不要な文字列が付着することがあるので、これを消去(切断)する。

(例) 濃厚 trans-テカリン - テリミクターにより切断した。前後に逆の向き、文字の変わり目でも切断する。

↓

濃厚 Trans テカリン

2.5. Alter 処理 : 以上の切断機能があるが適切なおらずる場合に自由に対処し、具体的には切断結果を指示する。

(例) 上 Cut, 陸上 Pass とすると、大陸上高気圧が切れるので、大陸上 Alter 大陸上 高気圧 と指示していく。

3.2. システム規模 : 辞書 (Version 1)、切断辞書 4,305 語
 不要語辞書, 2161 語.

プログラム, アセンブラ, 1本, 1200 ステップ
 (切断処理用).
 マスプロ 20本, 3500 ステップ
 コボル 1本, 670 ステップ
 (共に, I/O, ファイル作成
 等用)

4. 性能

4.1. 処理速度

使用計算機 : 日立 M-170
 対象 : JICST 科学技術文献タイトル, 3,500 件
 速度 (処理時間) : 40分16秒
 (内訳) 分かち書き 35秒 切断 34分42秒 抽出 1分50秒 不要語除去, 重複語除去 2分32秒, 1分16秒

これは, 1万件当りに換算すると, 1.91 時間 (設計時予測, 3.3 時間)

4.2. 精度

出力の 1/10 サンプリングによる推算

タイトル数 283
 正しく抽出されるキーワード, 1288 (平均タイトルキーワード4.6個)
 抽出されるべき全キーワード, 1321

$$\text{精度} = \frac{1288}{1321} \times 100 = 97.5\% \quad (\text{設計時予測, } 95\%)$$

正しく抽出される部分は, 切りすぎと切り不足である。

5. 出力例,

(1) (原文タイトル)

G010 「地震対策特集」 大規模地震の情報伝達と応急復旧
 (切断済タイトル)

M000 「地震対策 特集」 大規模 地震 の 情報伝達 と 応急
 復旧

(自動抽出キーワード)

M020 地震対策

M050 大規模

M060 地震

M080 情報伝達

M100 応急復旧

- (2) (原文タイトル)
 G010 音声信号に対する適応型反響阻止
 (切断済タイトル)
 M000 音声信号 に対する 適応型 反響阻止
 (自動抽出キーワード)
 M010 音声信号
 M050 適応型
 M060 反響阻止
- (3) (原文タイトル)
 G010 改善された予測法による言語のエントロピー推定とそのヘブライ語への応用
 (切断済タイトル)
 M000 改善 された 予測法 による 言語 の エントロピー 推定
 と その ヘブライ語 への 応用
 (自動抽出キーワード)
 M030 予測法
 M060 言語
 M080 エントロピー
 M120 ヘブライ語
 M140 応用
- (4) (原文タイトル)
 G010 スピノーダルCu-10wt%Ni-6wt%Snにおける低サイクル疲労
 (切断済タイトル)
 M000 スピノーダル Cu - 10wt % Ni - 6wt % Sn に
 おける 低サイクル疲労
 (自動抽出キーワード)
 M010 スピノーダル
 M020 Cu
 M040 10wt
 M060 Ni
 M080 6wt
 M100 Sn
 M130 低サイクル疲労
- (5) (原文タイトル)
 G010 鉾山総労働者数の変動状況調査
 (切断済タイトル)
 M000 鉾山 総 労働者数 の 変動 状況調査
 (自動抽出キーワード)
 M010 鉾山
 M030 労働者数
 M050 変動

6. JICSTにおける応用

- (1) 今年10月に向け辞書メンテナンスを重ねて精度向上を図り、57年4月からのサービスに使用する他、56年分の蓄積ファイルにも適用し、タイトルキーワードのJOLISオンラインにおける一次タグとして使用できるようにする。
- (2) 55年以前の旧ファイル、約200万件については、同様にタイトルキーワードが使用できるように検討する(主として計算機時間の問題)
- (3) 切断のスピードアップのため、全切断辞書が大幅に増加したのを見通した上で、オンコア処理することを検討(1/10の割合に短縮される見込)。

7. 文献

- (1) 菊池敏典. 自動索引法と自動分類法・展望. 情報管理 Vol 8 No 13 pp 3-9 (1965), Vol 8 No 5 pp 9-14 (1965) Vol 9 No 1 pp 21-28 (1966), Vol 9 No 4 pp 185-189 (1966)
- (2) 長尾真, 水谷幹男, 池田浩之. 日本語文献における重要語の自動抽出. 情報処理 Vol 17 No 2. pp 110-117 (1976)
- (3) 植村俊亮. 日本語KWIC索引の試み. 第5回情報科学技術研究会発表論文集 pp 143-152. (1968)
- (4) 荒木啓介. 漢字-カナ変換システムのその後の展開. 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)の開発. 計算言語学 24-4 (190. 8. 21)