

電子文書参照支援システム BENLI における内容検索機能

3 F - 8

— 前処理としてのキーワード検索 —

石橋 由希夫, 和田 浩一, 小川 洋一

セコム (株) IS 研究所

1 はじめに

我々は、電子文書参照支援システム BENLI の研究/作成を進めている [1]。BENLI における内容検索とは、ユーザが入力した自然言語文に対して、その文と検索対象の文の意味的な比較を行なう [2] ことで類似する部分を検索する機能である。しかし、対象文書の全てについて比較を行なうとシステム全体の負荷が大きくなり過ぎると考えられる。そこで、前処理としてキーワード検索を行なうことにより比較すべき部分を絞り込むという方法をとることとした。そのために、あらかじめ設定した検索件数が得られるような検索式を生成しなければならない。本稿では、一般のサーチャーがとっている戦略 [3, 4, 5] を取り入れた、検索式生成のためのアルゴリズムを提案する。

2 キーワード検索の概要

検索対象となる文書はハイパーテキストとして格納される。文書は構造的なまとまりごとにノード (以後、“文書ノード”と呼ぶことにする) として表される。この文書ノードを検索の単位とする。また、各文書ノードには ID 番号が与えられている。

キーワード検索では、ユーザが検索キーとして入力した文を解析して得られた意味構造中の自立語をキーワードとみなして、それらを論理結合した検索式で検索する。その検索式は、検索件数が前もって設定しておいた範囲になるように生成される (以後、このような検索式を“適切な検索式”と呼ぶことにする)。そして、適切な検索式による検索の結果得られた文書ノードの ID 番号のリストをキーワード検索の結果とする。検索には、キーワードからそのキーワードが出現する文書ノードの ID のリストを得ることができるキーワード索引を用いる。

3 概念とその特定性の高さ

概念とその特定性の高さを以下のように定義する。

(1) 概念

ある自立語に対して、その自立語に意味的に類似している自立語の中で、上位概念語または下位概念語にあたるものを類義語、それ以外の語を同義語と呼ぶことにする。

このとき、自立語とその自立語の同義語や類義語を和結合したものを概念と呼ぶことにする。

あるキーワード k について、 k とその同義語の集合を和結合したものを K 、また、類義語の集合を和結合したものを K^+ と表す。

(2) 概念の特定性の高さ

キーワードの論理結合式 E による検索件数を返す関数を、 $docfreq(E)$ としたとき、概念 K の特定性の高さを、

$$S = \frac{1}{docfreq(K)}$$

(ただし、 $docfreq(K) = 0$ のとき $S = 0$)

で表すことにする。

4 適切な検索式を生成するためのアルゴリズム

入力の N 個の自立語から得られた概念を、

$$K_1, K_2, \dots, K_N$$

とする。この N 個の概念の中に特定性の高いものがある場合、その特定性の高い概念から検索式を生成しても検索候補の絞り込みが十分にできる [3, 5]。そこで、 K_i の特定性の高さを S_i としたとき、与えられた閾値 θ に対して、

$$S_i > \theta$$

となる i があれば特定性の高い概念から検索式の生成を行なう (4.1 節)。そうでなければ、全ての概念の積から検索式の生成を行なう (4.2 節)。

なお、以下の説明で、前もって設定しておいた検索件数の範囲の上限と下限を各々 MAX , MIN とする。

4.1 特定性の高い概念からの検索式の生成

生成した検索式の検索件数が MAX より多い場合は、その検索式に新たな概念を積結合して、検索件数を減らす。逆に MIN より少ない場合は、検索式中の同義語からなる概念に対する類義語の和結合などにより、検索件数を増やす。具体的なアルゴリズムは以下に示す通り。

要素数 r 個の配列 $record$ について、

$$\bigvee_{i=1}^r (record[i])$$

が適切な検索式となるような *record* を求める。ただし、 $i < j$ ならば、 $S_i \geq S_j$ 。

```

rlen ← 1; {record の要素数}
i ← 1; record[1] ← (K1);
loop:
  if i > N then exit;
  if MAX < doqfreq(∏j=1rlen(record[j])) then
    begin “record[rlen]と(Ki+1)を積結合したものを
           新たにrecord[rlen]に代入する;
           i ← i + 1; goto loop end
  else if MIN > doqfreq(∏j=1rlen(record[j])) then
    begin rlen ← rlen + 1;
           if “record[rlen - 1]に(Ki)がある” then
             “record[rlen - 1]の(Ki)の部分を(Ki+)に置
             き換えられた式をrecord[rlen]に代入する”
           else begin
             “record[rlen - 1]の(Ki+)の部分を(Ki+1)に
             置き換えられた式をrecord[rlen]に代入する”;
             i ← i + 1 end;
           goto loop end

```

4.2 全ての概念の積からの検索式の生成

生成した検索式の検索件数が *MAX* より多い場合は、その検索式中の同義語と類義語からなる概念から類義語を除くことで、検索件数を減らす。逆に *MIN* より少ない場合は、積結合している概念の除去などにより、検索件数を増やす。適切な検索式を *record* に求める具体的なアルゴリズムを以下に示す。ただし、配列 *exp*, *under*, *over* には検索式が代入される。また *docfreq(nil)* = 0 とする。

```

exp[1] ← ∏i=1N(Ki ∨ Ki+);
record ← nil; elen ← 1; {exp の要素数}
olen ← 0; ulen ← 0; {over, under の要素数}
loop:
  {適切な検索式が生成できない場合}
  if elen = 0 then
    if ulen ≠ 0 then begin
      “配列underの中で検索件数が最も多いものを
       recordに格納する”; exit end
    else begin
      “配列overの中で検索件数が最も少ないものを
       recordに格納する”; exit end;
  {以下でexpの中に適切な検索式を探す}
  olen ← 0; ulen ← 0;
  for i = 1 to elen do
    if MAX < docfreq(exp[i]) then
      begin olen ← olen + 1;
            over[olen] ← exp[i] end
    else if MIN > docfreq(exp[i]) then
      begin ulen ← ulen + 1;
            under[ulen] ← exp[i] end
    else if docfreq(exp[i]) > docfreq(record) then
      record ← exp[i];
  if record = nil then

```

```

“underの組を調べる”-(1);
if record = nil then
  begin “別の検索式を新たにexpに代入する”-(2);
        goto loop end

```

(1),(2)の手順は以下の通りである。

(1) *under* の組を調べる。

```

if ulen > 1 then
  for i = 2 to ulen do
    “under の要素からi個取り出して和結合したもので
    中で適切な検索式になるものがあるならば、その中で
    検索件数が最も多い式をrecordに代入する”

```

(2) *exp* に適切な検索式がなかった場合に別の検索式を *exp* に代入する。

```

elen ← 0;
for i = 1 to olen do
  for j = 1 to N do
    if “over[i]に(Kj ∨ Kj+)がある” then
      begin elen ← elen + 1;
            “over[i]の(Kj ∨ Kj+)の部分をKjに
            置き換えられた式をexp[elen]に代入する”
      end;
  for i = 1 to elen do
    for j = 1 to N do
      if “under[i]に(Kj)または(Kj ∨ Kj+)がある” then
        begin elen ← elen + 1;
              “under[i]から(Kj), (Kj ∨ Kj+)の部分を
              除いた式をexp[elen]に代入する”
        end

```

5 おわりに

検索件数が適当な値になるような検索式を生成するためのアルゴリズムを提案した。現在、このアルゴリズムを実現するプログラムを作成中である。また、並行して対象文書の技術用語などを登録したソースの作成も進めている。今後、検索実験によりこのアルゴリズムの評価を行なう予定である。

参考文献

- [1] 小川他. 電子文書参照支援システム BENLI-システムの構成と機能の概要-. 第45回情報学大会 3F-07, 1992-10.
- [2] 和田他. 電子文書参照支援システム BENLIにおける内容検索機能-意味構造の一致度の判定手法-. 第45回情報学大会 3F-09, 1992-10.
- [3] Oldroyd B.K. Citroen C.L. Study of strategies used in on-line searching. *On-Line Review*, Vol. 1, No. 4, pp. 295-309, 1977.
- [4] Marcia J.Bates. Information Search Tactics. *Journal of the American Society for Information Science*, pp. 205-214, 1979.
- [5] 三輪眞木子. サーチャーの時代. 丸善, 1986.