

文字認識と形態素解析を用いた
類似文書検索の試み

6M-7

間瀬久雄 小山幸子 木山忠博 辻 洋 絹川博之
(株)日立製作所 システム開発研究所 関西システムラボラトリ

1. はじめに

新聞や雑誌、特許、論文、報告書、議事録などの大量文書から必要な文書だけを検索するための技術が要求されている。我々は、印刷文書紙面を入力として、この文書に類似した文書を全自動で検索する方式を検討している。すなわち、印刷文書を文字認識し、認識結果を形態素解析し、解析結果を基にキーワードを抽出し、抽出したキーワードを基に検索式を生成し、検索式に基づいて文書DBを全文検索し、検索結果文書に対して類似度を判定し、類似文書を抽出する。本方式は、特許の公知例調査などにとくに有効である。

本稿では、印刷文書の文字認識精度がキーワード抽出精度に及ぼす影響について考察する。

2. 処理の流れ

図1に本方式の処理の流れを示す。

(1) 文字認識

LBP出力のA4横書き印刷文書を入力とする。印刷紙面やスキャナのノイズによる精度劣化を含めた文字認識精度が95%前後のものを採用した。

(2) 形態素解析¹⁾

5万語の語彙を持つ辞書を参照して最長一致による単語分割を行い、品詞情報を取得する。

(3) キーワード抽出

品詞が名詞である単語と、辞書に未登録の単

語(未知語)を対象として、出現頻度を算出し、頻度の高い単語をキーワードとする。ただし、ほとんどの文書に共通に出現する単語(「場合」「とき」「こと」「もの」等)は除去する。

(4) 検索式生成

キーワードを論理演算子(AND/OR)で結合して検索式を生成する。

(5) 全文検索

全文検索システム Bibliotheca/TS²⁾により全文検索する。検索式のキーワードを含む文書を漏れなく抽出する。同義語や異表記も吸収できる。

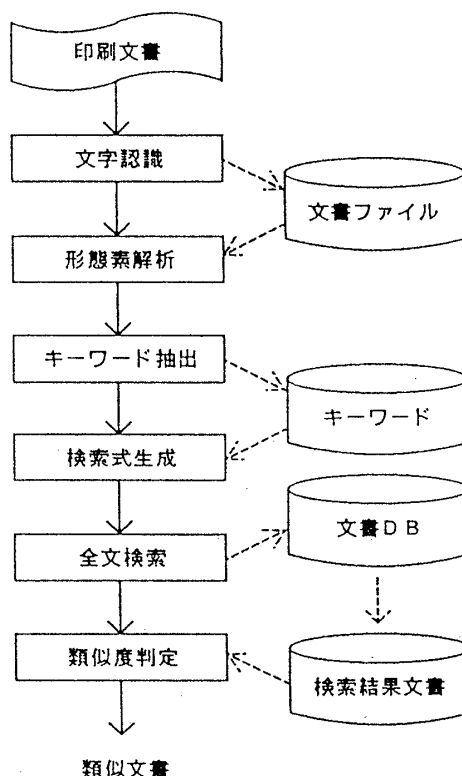


図1 処理の流れ

Text Retrieval using OCR and Morph Analysis
Hisao MASE, Sachiko KOYAMA, Tadahiro KIYAMA,
Hiroshi TSUJI, Hiroshi KINUKAWA
Systems Development Laboratory, Hitachi Ltd.

表1 文字認識結果

正解	4938字 (93.6%)
不正	かな 232字 (4.4%)
	記号 53字 (1.0%)
正解	漢字 21字 (0.4%)
	数字 16字 (0.3%)
解	力ナ 14字 (0.3%)
	英字 3字 (0.1%)
合計	339字 (6.4%)

表2 字種別認識結果

字種	正解文字数
かな	1720字 / 1952字 (88%)
記号	397字 / 450字 (88%)
漢字	2477字 / 2498字 (99%)
数字	202字 / 218字 (93%)
力ナ	126字 / 140字 (90%)
英字	16字 / 19字 (84%)
合計	4938字 / 5277字 (94%)

表2 キーワード抽出結果

評価項目(%)	頻度					
	7以上	6以上	5以上	4以上	3以上	2以上
再現率	9/9 (100)	12/12 (100)	27/27 (100)	43/47 (91)	79/80 (99)	164/174 (94)
適合率	9/11 (82)	12/14 (86)	27/30 (90)	43/50 (86)	79/88 (90)	164/189 (87)
正当率	9/11 (82)	12/14 (86)	27/30 (90)	43/54 (80)	79/89 (89)	164/199 (82)

※ 再現率(%)=(両方の文書から抽出されたキーワードの種類数)/(認識率100%の文書から抽出されたキーワードの種類数)

※ 適合率(%)=(両方の文書から抽出されたキーワードの種類数)/(文字認識された文書から抽出されたキーワードの種類数)

※ 正当率(%)=(両方の文書から抽出されたキーワードの種類数)/(どちらか一方の文書から抽出されたキーワードの種類数)

(6) 類似度判定

抽出された各文書に、何種類のキーワードが何回出現したかによって類似度を求め、類似度の高い文書から順に表示する。

3. 本方式における仮説

本方式では、次の仮説を立てている。

[仮説1] 文書に頻繁に現れる語句は、その文書の「キーワード」となりうる。

[仮説2] 文書1の「キーワード」を多く含む文書nは、文書1の「類似文書」である。

[仮説3] 「キーワード」を仮説1のように定義し、「類似文書」を仮説2のように定義すると、文字認識精度がキーワード抽出精度および類似文書検索精度に及ぼす影響は小さい。

4. 文字認識精度とキーワード抽出精度の関係

文字認識精度がキーワード抽出に影響を及ぼすかを新聞記事10文書を対象に実験評価した。

表1、表2に文字認識結果を示す。精度は93.6%であった。また、表3にキーワード抽出結果を示す。ここでは文字認識精度を100%とした場合に抽出されるキーワードを正解と

した。

再現率については、誤認識結果が上位のキーワード(頻度5以上)には影響がなかった。適合率については、「た」を「左」と誤認識した結果、キーワード「左」がノイズとなった文書が多かった(10文書中8文書)。適合率の改善については、全文検索方法(検索式の生成方法)の検討などがあり、今後の課題である。

5. おわりに

文字認識精度がキーワード抽出精度に及ぼす影響についての評価実験結果について述べた。文字の誤認識は、キーワード抽出の再現率にはほとんど影響を及ぼさなかった。

今後は、検索式の生成方法について検討するとともに、類似文書の検索精度について実験評価していく予定である。

参考文献

- 1) 西森ほか：汎用日本語形態素解析ツールの開発、情処学会第44回全国大会
- 2) 浅川ほか：フルテキストサーチシステム Bibliotheca/TSの開発、情処学会第45回全国大会