

2N-2

## ファジィ文書検索システム(1) ～実験システムと評価～

森田哲也 小川泰嗣 小林清彦  
(株)リコー 中央研究所

### 1. はじめに

光ディスクファイリングシステムの普及により大容量データベースがオフィス・個人規模で利用可能となっている。しかしこれらのシステムは、ファイリング時には大量文書に対する分類やキーワード付け作業が煩雑であり、検索時には利用者の要求主題に適した文書がなかなか探し出せない等の問題点があった。これらを解決するためのファイリング支援方式として、自動分類方式・キーワード自動抽出方式の研究が、また検索支援方式としてファジィシソーラス[1]・キーワード関係知識ベース[2,3]等を用いた情報検索の研究が行なわれている。

我々は、[2,3]で提案されたキーワードコネクションマトリックスを用いたファジィ文書検索および学習方式を拡張し、また文書数10,000件についてキーワード自動抽出を行なうことによりキーワード数約2,000・コネクション数約72,000の知識ベースを持つ知的ファイリングシステムを開発した。本論文では、その構成および性能評価結果について述べる。

### 2. キーワードコネクションマトリックス

[2,3]では語彙関係を用いた概念空間として、キーワードコネクションマトリックス(以下、KCM)を提案した。これは任意の2キーワード間の概念的類似度を定量的に表現したキーワード関連度を、知識として持つキーワード関係知識ベースである。KCMは、キーワード数をKとすると、 $K \times K$ の行列Wとして与えられる。キーワードiとjの関連度の初期値 $W_{ij}$ は、

$$W_{ij} = \begin{cases} \frac{N_{ij}}{N_i + N_j - N_{ij}} & : i \neq j \\ 1.0 & : i = j \end{cases} \quad \dots (1)$$

$N_i, N_j$ : キーワードi, jそれぞれの出現頻度

$N_{ij}$ : キーワードiとjの共出現頻度

KCMは一種のシソーラスであり、その初期値は[1]におけるRT(Related Term)の計算式と等価であるが、後述の学習によりキーワード関連度が変化する点で異なっている。

### 3. 機能および処理概要

#### ①自動ファイリング機能:

まず、入力文書を形態素解析した結果に語接続規則が適用され、キーワード候補が抽出される。候補から不用語を除去した後、既存キーワードとの比較・出現頻度処理が施され、結果がユーザに提示される。ユーザは文書内容に適したキーワードを抽出結果から選択し文書に付与することができる。

#### ②あいまい連想検索(ファジィ文書検索)機能:

一般にキーワードと論理演算子(AND/OR/NOT)からなる検索式Queryは(2)式のような積標準形に変形できるため、ファジィ文書検索における結果集合は、検索式中の任意のキーワードKに対して各文書が[0,1]の値をメンバーシップ値としてもつファジィ集合D(K)を用いて(4)(5)式のように表現される。即ち積標準形化されたQueryに対するファジィ検索の結果集合は、それぞれの副検索式SubQueryに対する結果を用いて表せる。

ここで、Nおよびhは副検索式の総数および番号、nおよびmはそれぞれ否定演算子( $\neg$ )のないキーワードと否定演算子付きのキーワードの数を表す。

$$\text{Query} = \text{SubQuery}(1) \wedge \dots \wedge \text{SubQuery}(N) \quad \dots (2)$$

$$\text{SubQuery}(h) = K_1 \vee \dots \vee K_{n_h} \vee \neg K_{n_h+1} \vee \dots \vee \neg K_{m_h} \quad \dots (3)$$

$$\text{Result} = \text{SubResult}(1) \cap \dots \cap \text{SubResult}(n) \quad \dots (4)$$

$$\text{SubResult}(h) = D(K_1) \cup \dots \cup D(K_{n_h}) \cup (\neg D(K_{n_h+1})) \cup \dots \cup (\neg D(K_{m_h})) \quad \dots (5)$$

集合演算として代数和・代数積を使用すると、副結果集合SubResult(h)の文書iのメンバーシップ値 $r_i(h)$ は、

$$r_i(h) = 1 - \left( \prod_{K_j \in Q(h)+} S_{ij} \right) \left( \prod_{K_j \in Q(h)-} R_{ij} \right) \quad \dots (6)$$

$$\text{where, } R_{ij} = 1 - \prod_{K_k \in A_i} (1 - W_{jk}), \quad S_{ij} = 1 - R_{ij}$$

$Q(h)+, Q(h)-$ : 否定なし、付きの検索キー集合

$A_i$ : 文書i中に含まれるキー集合

各文書確度 $r_i$ は、全ての副検索結果の結合であり、

$$r_i = \prod_{h=1}^N r_i(h) \quad \dots (7)$$

#### ③適応化学習機能

ユーザは検索結果に対する適切さの評価値を入力することにより、キーワードコネクション中のキーワード関連度を変更できる。これを学習機能と呼び、検索結果の文書確度と適切さの評価値によって定義される評価関数に最急降下法を適用することによって実現した[4]。

Fuzzy Document Retrieval System.(1)

--Experimental System and Results--

Tetsuya MORITA, Yasushi OGAWA, Kiyohiko KOBAYASHI  
Ricoh, Co. Ltd. E-mail: morita@takezo.rdc.ricoh.junet

4. システム構成

本システムでは、文書画像を取り込み文字コード化するOCRユニットおよび文書・KCM等を保持する光ディスクユニットが、Unixワークステーションに接続されている (Fig.1)。またソフトウェアはC言語で記述し、文書ファイリング支援を行なうキーワード自動抽出部と、あいまい連想検索および学習を行なうファジィ検索・学習部の2ブロックからなる。ユーザI/FはXウィンドウバージョン10を用いて開発した。

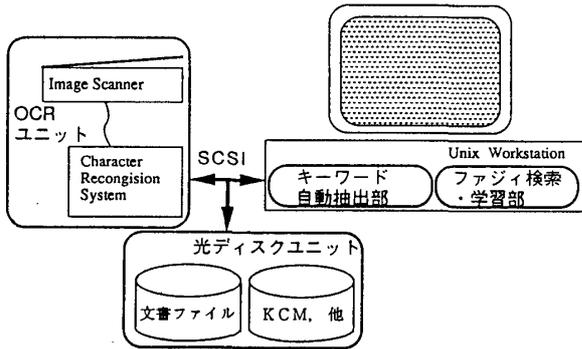


Figure.1 実験システムの構成

4. 実験方法

技術文書10,000件からキーワード自動抽出機能により約2,000 キーワードを抽出、またこれらを基に72,000リンクからなるキーワードコネクションを作成した。

4.1 検索機能の性能測定方法

検索条件式として4つのキーワードの単体および論理和/論理積/否定の組合せによって構成した。条件式に対して4人の実験者が個々に正解と見做した文書集合を、統計処理して正解集合を作成し、ファジィ検索の結果集合との関係を再現率・適合率で測定した。同様の検索条件に対する従来のクリスプな検索方式を用いた場合の再現率・適合率を測定し、上記結果と比較した。

4.2 適応化学習機能の性能測定方法

キーワード単体および論理和を用いた検索式の2種類について学習前の再現率・適合率と、学習を経過した後のそれらを測定した。学習は実験者によって作成された正解集合を適切さの評価値=1.0 (正解)として30回の学習を試行した。

5. 結果および評価

本ファジィ文書検索方式と従来方式について3種類の検索条件に対する性能比較をTable.1に示す。

検索条件	ファジィ文書検索		従来方式	
	再現率	適合率	再現率	適合率
2キ-論理積	52 %	5 %	22 %	6 %
2キ-否定	72 %	19 %	23 %	7 %
1キ-単体	66 %	40 %	44 %	62 %
平均	63 %	21 %	30 %	25 %

Table.1 各検索条件における性能比較

検索式として {CAD or LSI} を用いた時の再現率 (Recall) ・適合率 (Precision) vs. 閾値の関係を図. 2に示す。閾値は検索結果集合を決定する文書確度の閾値である。また、検索式 {CAD} で正解集合を用いて学習を繰り返した結果を図. 3に示す。

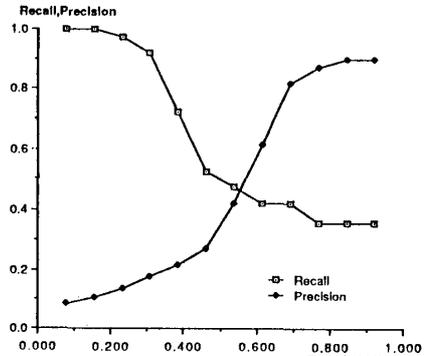


Fig.2 Recall, Precision vs. Thresh. (query={CAD or LSI})

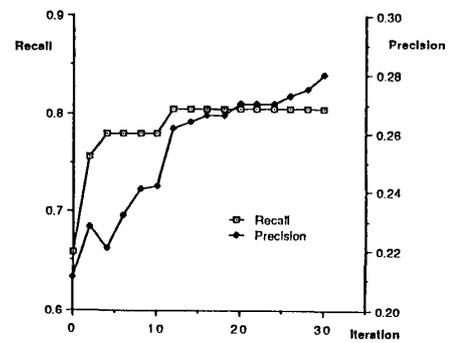


Fig.3 Learning Effect. (query={CAD})

6. まとめ

キーワード自動抽出・連想検索・適応化学習機能を持つファジィ文書検索システムをWS上に構築し、検索および学習機能について再現率・適合率を測定することにより評価を行なった。従来のクリスプ検索に比べ連想による再現率の向上・および学習による効果が認められた。

7. 参考文献

[1] S.Miyamoto, T.Miyake, and K.Nakayama. "Generation of a Pseudthesaurus for Information Retrieval Based on Co-occurrences and Fuzzy Set Operations". IEEE Trans. of SMC, vol.SMC-13, No.1, pp.62-70,1983.  
 [2] 森田, 小川, 小林. "動的キーワードコネクション方式を用いた文書検索システム". 昭63信学全大. vol.D1, pp.132, 1988.  
 [3] Y.Ogawa, T.Morita, and K.Kobayashi. "Fuzzy Document Retrieval System and Its Learning Method Based on the Keyword Connection". Proc.of Int.Workshop on Fuzzy System Applications, pp. 143-144, 1988.  
 [4] 小川, 森田, 小林. "ファジィ文書検索システム(2) ~キーワードコネクションマトリックスの学習方式~" 情報処理学会第39回全国大会予稿集, 1989.