

文書検索におけるクエリーの拡張方法 —大域的分析和局所的分析の実証比較—

岸田和明

駿河台大学文化情報学部
〒357-8555 埼玉県飯能市阿須 698
kishida@surugadai.ac.jp

利用者が作成した検索のための質問文をシステムが自動的に拡張することによって、検索性能をより向上させる研究が進められている。その方法としては、データベース中の語の出現統計を用いて自動的に作成されたシソーラスを利用する大域的な分析方法と、利用者自身の質問文によって第1次の検索をおこない、その上位文書を使って拡張する局所的な分析方法とがある。本研究の目的はこれらの2種類の分析方法の性能を実証的に比較評価することにある。データとしてはNTCIR-1のテストコレクションを用いた。今回は、大域的分析の例としてダイス係数に基づいて自動構築されたシソーラスによる拡張を取り上げ、また局所的分析の例として自動適合フィードバックを用いた。その結果、局所的な分析は元の質問文による性能を改善し、それに対して、大域的な分析ではその性能は低下する結果となった。なお、検索モデルとしてはベクトル空間型を用い、フィードバックはRocchioの方法を利用した。

Global and Local Analysis for Automatic Query Expansion in Document Retrieval: An Empirical Study

Kazuaki KISHIDA

Faculty of Cultural Information Resources, Surugadai University
698 Azu, Hanno, Saitama 357-8555
kishida@surugadai.ac.jp

Automatic query expansion is widely known as one of the useful methods for enhancing retrieval performance. There are two ways for the expansion: (1) global analysis and (2) local analysis. An example of the global analysis is automatically constructed thesaurus based on term co-occurrence statistics. On the other hand, the local analysis means that a set of documents retrieved by the original query is examined for the expansion. The purpose of the paper is to compare empirically the retrieval performance between the global and local methods. A Japanese test collection NTCIR-1 is used for the retrieval experiment. As a result, automatic relevance feedback as a local analysis outperforms global analysis based on statistical thesaurus constructed from statistics of co-occurrence. In the experiment, the vector space model is used as a retrieval algorithm, and the automatic relevance feedback is based on the Rocchio method.

1 はじめに

データベースやインターネットを検索する際に、利用者が自分の情報要求を的確かつ十分に表

現できるとは限らない。例えば、検索したいこと
がらに利用者が精通していない、あるいは、その
データベースを検索した経験が乏しいなどの理
由から、適切な検索語がわからないという状況は

頻繁に起こりうる。もともと、何らかの不明確なことがらに対して検索要求が生じるのであるから、むしろ、その情報要求を適切に記述できる場合のほうが少ないとも考えられる。この結果、多くの利用者の質問文（またはクエリー）は十分な長さを持たず、適合文書を検索するための十分な検索語を含まない可能性がある。

この問題を解決するためのひとつの方法は、何らかの別の情報源から新たな検索語を選び出し、利用者自身による元のクエリーに自動的または半自動的に追加することである。これは一般にクエリーの拡張 (query expansion) と呼ばれ、これまでに数多くの研究が積み重ねられてきた [1]。

その情報源としては、まず、伝統的なオンライン情報検索におけるシソーラスや辞書・事典などの外部的な資源を活用することが考えられる。しかし、適切なシソーラス等が常に利用可能であるとは限らない。そこで、検索対象であるデータベース自体から、いわば bootstrap 的に、新規追加する語を識別・抽出する方法が重要となる。

このための方法には、大域的分析と局所的分析の2つがある [2]。本稿では、これらの分析方法の検索性能の相違に焦点をあて、その実証的な比較を試みる。具体的には、大域的分析の例として、データベース全体から自動構築されたシソーラスによる拡張の性能を調査する。一方、局所的分析としては、ここではいわゆる自動適合フィードバックを取り上げる。実証比較に使用するデータベースは、国立情報学研究所 (NII) による日本語検索用テストコレクション NTCIR-1 である [3]。これは日本の学会発表の標題や抄録のレコード約 33 万件から成るデータベースに基づいている。

以下、第2節では、これまでのクエリー拡張についての研究の概観を展望し、第3節では、本稿における実証比較の実験方法について述べる。第4節では実験結果について議論する。

2 クエリーを拡張するための大域的分析和局所的分析

2.1 シソーラスの自動構築

検索対象のデータベースにシソーラスが存在しない場合、それを自動構築することが考えられる。これには伝統的に、語の出現頻度についての統計情報が利用されてきた。すなわち、語 t_j と t_k

の出現文書数をそれぞれ n_j 、 n_k とし、また両者が共起する文書数を n_{jk} とすれば、これらの2つの語間の関連度を、例えばダイス係数 (Dice coefficient) によって、

$$\rho = 2n_{jk} / (n_j + n_k) \quad (2.1)$$

として計算できる。そして、密接な関連を持つ語のペアをデータベース中から自動識別することが可能になる。

この統計的なシソーラスを使えば、クエリーの拡張は容易である。つまり、元の検索語と各語の関連度をシソーラスから探索し、ある閾値を超えたものを新たな検索語として追加すればよい。この方法に対しては、語の統計的共起関係は検索性能の向上には役立たないという批判や実験結果がある [4][5]。そのため、語や文書のクラスタリングを応用する手法 [6][7] や、非対称の関連尺度を導入する試み [8]、ベイズ推論ネットワークを用いる研究 [9] などがなされている。しかしその一方で、(2.1)式 (あるいはコサイン係数などの他の尺度) に基づく方法は実装が容易などの利点を持ち、研究が続けられているのが現状である [10][11][12]。

2.2 自動適合フィードバック

語の共起関係に基づくシソーラスの自動構築では、基本的には、データベース全体を対象に出現文書数等が計算される。この点に関して、さまざまなトピックが含まれる「異質な」データベース全体を大域的 (global) に分析しても有用な情報は得られないという批判がある。この批判からは、検索質問に関連した部分のみを局所的 (local) に分析するという考え方が導かれる。例えば、元のクエリーによって検索された何件かの上位の文書のみを対象として、局所的なシソーラスを自動構築することなどが考えられる [2]。

また、フィードバック手法の適用も局所的分析の範疇に含めることができる。いわゆる自動適合フィードバック (automatic relevance feedback, あるいは pseudo relevance feedback) は元のクエリーによって検索された上位何件かの文書を、強制的に適合文書と見なして、Rocchio の方法などの適合フィードバックの手法を応用するものであり、これも一種のクエリー拡張のための局所分析と見なすことが可能である。

Rocchio の方法はベクトル空間型の検索モデルに基づいており、修正前のベクトルを q

```

<REC>
<ACCN>gakkai-j-0000441590</ACCN>
<TITL TYPE="kanji">大規模テストコレクション NTCIR-1 の構築 (1) -プーリングと正解判定
の分析-</TITL>
<AUPK TYPE="kanji">栗山 和子 / 江口 浩二 / 野末 俊比古 / 神門 典子</AUPK>
<CONF>全国大会</CONF>
<CNFD>1999. 09. 28 - 1999. 09. 30</CNFD>
<ABST TYPE="kanji"><ABST.P>本研究の目的は, (1) 大規模テストコレクションを構築する手
法としてのプーリングの有効性を検証し, (2) プーリング件数が検索システムの評価に関連が
あるかどうか調べ, (3) 正解判定の際の判定のゆれがシステムの評価に関係してくるかどうか
を明らかにすることである。</ABST.P>... (中略) ...<ABST.P> (3) のために, NTCIR-1 の評価
用セットを構築する際に行なった, 異なる判定者による3種類の正解判定結果 (判定者 A, B
それぞれによる判定, 両者の協議による最終判定) を用いて評価実験を行なった。結果として,
53 件の検索課題を用いて検索結果を評価したとき, 検索精度の平均は異なる正解判定リスト間
においてほとんど差がなくなり, 他数の検索課題を用いて評価を行なえば, 判定者間の判定の
ゆれは評価においては問題ではないということがわかった。</ABST.P></ABST>
<KYWD TYPE="kanji">テストコレクション // プーリング // 情報検索 // NTCIR // Move-to-
Front 法 // 網羅性</KYWD>
<SOCN TYPE="kanji">情報処理学会</SOCN>
</REC>

```

図 3.1 NTCIR-1 における文書レコードの例

($= (w_{q1}, \dots, w_{qM})$, w_{qj} は語 t_j の検索質問中の重み, M はデータベースに含まれる語の異なり総数), フィードバックによる修正後を \tilde{q} とかくと,

$$\tilde{q} = \alpha q + \frac{\beta}{|D_1|} \sum_{i: d_i \in D_1} d_i - \frac{\gamma}{|D_0|} \sum_{i: d_i \in D_0} d_i \quad (2.2)$$

である[13]。ここで d_i は文書 d_i の主題表現ベクトルであり, w_{ij} を文献 d_i における語 t_j の重みとすれば, $d_i = (w_{i1}, \dots, w_{iM})$ である。また, D_1 は適合文書の集合, D_0 は不適合文書の集合, α , β , γ はパラメータである。自動適合フィードバックの場合には, 各文書の適合/不適合の情報が存在しない状況を考え, 第1次の検索結果の上位何件かの文書のベクトルをそのまま(2.2)式の右辺第2項 d_i の計算に利用することになる(第3項は考えない)。

3 大域的解析と局所的解析との比較実験

3.1 実験目的

日本語テストコレクション NTCIR-1 を使って, 大域的解析の例としての統計的ソーラスによる拡張と, 局所的解析の例としての自動適合フィ

ードバックによる拡張とを, その検索性能に焦点をあて, 実証的に比較することがここでの目的である。

3.2 実験用テストコレクション: NTCIR-1

検索実験用の日本語テストコレクション NTCIR-1 の文書レコードの具体例を図 3.1 に示す。このうち, 検索に用いるのは標題<TITL>と抄録<ABST>のみである。後述するように著者キーワード<KYWD>はソースの自動構築のみに使用する。なお, 本稿で利用する文書レコードは, NTCIR-1 に含まれる日本語レコード 332,918 件である。

また, NTCIR-1 に含まれる検索質問(トピック)のうち, No.31~No.83までの53個を用い, そのうちの<DESCRIPTION>フィールドだけを抜き出して質問文と見なすことにする。これはいわば「短い質問文」に相当する。NTCIR ではより詳しくトピックを表現した<NARRATIVE>フィールドも用意されているが, 今回の研究目的からすればこれは使用する必要はない。

3.3 使用する索引作成法と検索手法

日本語テキストからの索引語および検索語の抽出方法は辞書との最長一致に基づく単純な方

法とする。すなわち、一般的な機械可読辞書の見出し語とテキストとを突き合わせ、最も長い見出し語と一致する部分を「語」とみなす。ただし専門的な用語を拾うために、さらに隣接する2つの「語」を複合語として機械的に組み合わせる方法を併用した(詳細は岸田[14]を参照)。

また、検索手法は、今回はベクトル空間型モデル[15]を用いる。すなわち、

$$w_{ij} = \log x_{ij} + 1.0 \quad (3.1)$$

$$w_{qj} = (\log x_{qj} + 1.0) \log(N/n_j) \quad (3.2)$$

である。ここで、 x_{ij} は文書 d_i 中の語 t_j の出現頻度、 x_{qj} は質問文中の語 t_j の出現頻度、 N はデータベース中の文書総数である。なお、 $x_{ij} = 0$ の場合には $w_{ij} = 0$ と定義する (w_{qj} についても同様)。そして、質問文に対する文献 d_i の類似度 s_i を、

$$s_i = \sum_{j=1}^M w_{ij} w_{qj} / \sqrt{\sum_{j=1}^M w_{ij}^2 \sum_{j=1}^M w_{qj}^2}$$

で計算する。

3.4 統計的シソーラスによる拡張

すでに述べたように、共出現統計から自動構築されるシソーラスはその方法の簡便性などの利点から、まだ探究する価値があると考ええる。そこで、本稿では複雑な工夫は加えずに、ダイス係数(2.1)式に基づく素朴なシソーラスの自動構築を考える。

ただし、検索性能を向上させるための若干の工夫は加える。すなわち、図 3.1 の著者キーワード <KYWD> を一種の「疑似ディスクリプタ」と見なして、この著者キーワードのみを新規の検索語として追加することにする。この方法は MeSH のディスクリプタを新規追加した Srinivasan の研究[10]と似ている。もちろん、NTCIR-1 の著者キーワードは統制された語彙ではない。しかし、著者キーワードは各文書の著者がその内容を的確に表現しようと努力して付与したものであることを考えれば、これを疑似ディスクリプタとして扱うことができると考えられる。

Title: シソーラスの自動構築

Author Keyword: シソーラス / 質問文拡張

図 3.2 語の関連の抽出方法

本稿で作成する統計的シソーラスにおける、語間の関連の抽出方法を図 3.2 に示す。まず、抄録は使わずに、標題だけを利用する。これは、標題中には著者キーワードと同様に内容を比較的よく表わす語が含まれているのに対して、抄録中の語は、文書の主題を表現する以外のさまざまな目的で使用されていると想定されるからである。次に、著者キーワードはそのまま使わずに、標題や抄録に対する索引作成と同じ方法で分解する。これは、著者キーワードは統制語彙ではなく、表記のゆれが統一されていないわけではないという理由による。

この結果、本稿のシソーラスでは、その見出し語は標題中に出現する語のみであり、その見出し語の下に、著者キーワード(の断片)と(2.1)式で計算される関連度がそれぞれ記録されることになる。

そして、実際の質問文の拡張の手順は次のようになる。

- ①元の質問文に対して索引作成をおこなって検索語を識別する
- ②検索語ごとにシソーラスを探索し、追加の候補となる著者キーワード(の断片)を抽出する
- ③そのうち閾値を超えるもののみを質問文に追加する。なおこの際に次の(a)~(b)のルールを導入する。
 - (a)シソーラスから新規に追加する語は元の検索語1語あたり20語までとし、それを超える場合には関連度上位20語のみを追加する。
 - (b)新規に追加された語の重みは1.0とする。
 - (c)複数の検索語が同じ語をシソーラスから新規追加した場合でも、その重みは1.0とする。

3.5 自動適合フィードバックの手法

本稿では、検索手法としてベクトル空間型(3.1)式および(3.2)式を使うため、Rocchioの方法(2.2)式をそのまま利用してフィードバックをおこなうことができる。なお、パラメータは $\alpha = 8$ 、 $\beta = 16$ とする。また、適合と見なす上位の文書件数は1, 2, 5, 10の4種類を試してみることにする。

3.6 評価基準としての検索実行

以下の実験では、基本的には、上記の統計的シソーラスによる拡張方法と自動適合フィードバ

ックによる拡張方法とを比較するが、その他の評価基準として、

- ① 元の質問文（検索語）のみを使用した、(3.1)と(3.2)式での検索
- ② 適合・不適合情報を利用した、本来的な適合フィードバック ((2.2)式を利用、ただし $\gamma = 4$)。

の実行も併せておこなう。①は何の工夫も加えない検索であり、いわば評価のためのベースラインとして機能する。一方、②は他の方法と比較して、より多くの情報（適合判定の情報）を使用しており、当然、検索性能は高くなる。したがって、これは一種の「上限値」を与える方法として利用できる。なお、今回の実験では、②のフィードバックの場合に使用する判定文書の数として、上位5件と上位10件の2種類を試してみることにする。

4 比較実験の結果

4.1 索引作成の結果

3.3 節で説明した索引作成の方法を使って、すべての日本語文書レコード 332,918 件から語を自動識別したところ、1 文書あたりの平均で 118.0 語（異なり数）が抽出された。一方、53 件のトピック（質問文）に対して同様の自動抽出を試みた結果、表 4.1 のような数値が得られた。

表 4.1 元の質問文中の語数

平均	標準偏差	最大	最小
9.452	3.569	23	2

表 4.2 シソーラスから追加された語数

	ダイス係数の閾値			
	0.01	0.05	0.1	0.2
平均	127.98	45.90	14.82	3.77
標準偏差	48.73	23.13	9.58	3.08
最大値	299	123	49	13
最小値	20	3	2	0

表 4.3 自動適合フィードバックでの追加語数

	適合と見なす上位文書数			
	1	2	5	10
平均	48.77	99.34	225.02	400.94
標準偏差	18.27	32.33	73.46	118.64
最大値	110	205	514	763
最小値	10	26	102	274

4.2 クエリーの拡張の状況

以上の索引作成結果を使って統計的シソーラスを実際に作成したところ、標題中の語と著者キーワードの断片のペアの総数は約 16,870,000 件、シソーラスの見出し語は約 449,000 語となった。このシソーラスを使って 3.4 節の手順およびルールに従って質問文（クエリー）を拡張した結果を表 4.2 に示す（53 トピックの平均値および標準偏差など）。なお、ダイス係数の閾値は 0.01, 0.05, 0.1, 0.2 の 4 段階とした。また、自動適合フィードバックにより追加された語数を表 4.3 に示す。これらの表が示すように、自動適合フィードバックによる追加語数のほうがシソーラスからの追加語数よりも多い。

4.3 検索性能の比較

次に、平均精度の平均（Mean Average Precision : MAP）と再現率-精度グラフを使って、各実行の検索性能を比較する。まず、MAP を表 4.4 に示す。

表 4.4 平均精度の平均 (MAP)

実行	略称	MAP
(1)ベースライン		
元の質問文のみ	ORG	0.228
(2)シソーラスでの拡張		
閾値 0.01	TH1	0.123
閾値 0.05	TH2	0.158
閾値 0.1	TH3	0.181
閾値 0.2	TH4	0.197
(3)自動適合フィードバック		
上位 1 件を適合	AT1	0.242
上位 2 件を適合	AT2	0.261
上位 5 件を適合	AT3	0.265
上位 10 件を適合	AT4	0.244
(4)適合フィードバック		
判定文書数 5 件	FB1	0.348
判定文書数 10 件	FB2	0.376

表 4.4 が示すように、統計的なシソーラスから新規に語を追加した場合、元の質問文のみによる検索よりも、性能が低下してしまっていることがわかる。逆に、自動適合フィードバックは検索性能の改善をもたらしている。

すでに述べたように、追加語数はシソーラスを使った場合のほうが少ないが、検索性能の低下は、

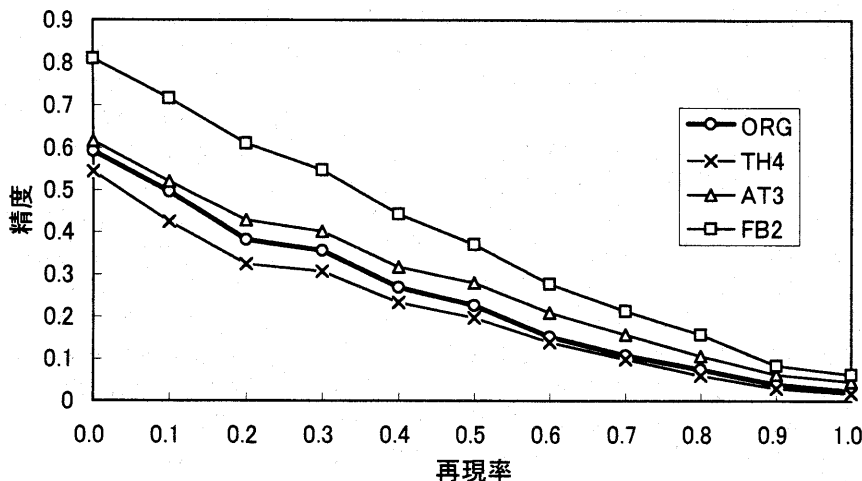


図4.1 再現率-精度グラフ(補間による)

明らかに、語数の少なさに起因するものではない。当然、ダイス係数の閾値を下げれば、追加語数は増加するが、それに伴って、検索性能は低下している(表4.4参照)。逆に、閾値を上げると検索性能は向上するが、これは単に追加語数が少なくなり、元の質問文に近づいていることを意味しているにすぎない。今回の実験では、ダイス係数に基づくマクロ的な統計的シソーラスによる拡張は負の効果をもたらす結果となった。

一方、自動適合フィードバックの場合、正解と見なす文書数を1件から増やしていくと、検索性能は増加していく。しかし、10件までを適合文書としてしまうと、5件の場合よりも検索性能が低下する。表4.4からはその文書数の最適値は読み取れないが、5件~9件の間にそれが存在していると考えられる。

表4.4の(1)~(4)のそれぞれから1つの実行を選んで再現率-精度グラフを描いた結果を図4.1として示す。(2)~(4)についてはそれぞれ最も性能の高いものを採用した(すなわち、TH4, AT3, FB2)。

表4.5 各実行のMAPの増加率(ORGに対する比率)

TH4	AT3	FB1	FB2
-13.6%	16.2%	52.6%	64.9%

元の質問文による実行ORGに対しての、TH4, AT3, FB1, FB2の各実行のMAPの増加率を表4.5に示す。最高でFB2の64.9%の増加である。文書数5の場合の真の適合フィードバック(FB1)

による増加率は52.6%であり、それに対して、文書数5の自動適合フィードバックは16.2%であるから、文書数5の場合の「上限値」をFB1であると仮定すれば、自動適合フィードバックはこの場合、上限値の約30%(=16.2/52.6)の改善をもたらしたことになる。

4.3 実験結果に対する考察

上で述べたように、統計的シソーラスによる拡張は検索性能の低下をもたらし、逆に、自動適合フィードバックは性能を改善した。この結果を素直に受け入れれば、大域的分析よりも局所的分析のほうが優れているということになる。

しかし、今回用いた大域的分析の方法はかなり素朴なものであることに注意する必要がある。文書レコード中の著者キーワードを使用するという工夫は加えられているものの、共起関係から計算されるダイス係数を単純に用いているにすぎない。しかも、新規に追加される語の重みに対する工夫はなく、単に1.0とするのに留まっている。それに対して、今回の自動適合フィードバックは長年研究されてきた(2.2)式のモデルに基づいており、この意味では、もう1段階上のレベルの工夫がなされているといつてよい。

この点を考慮すると、今回の結果から、ただちに大域的分析を無用なものとして即断するのは早すぎるかもしれない。単にダイス係数やコサイン係数を適用するだけでなく、もう少し「きめこまかい」何らかの確率的モデルを導入してみる余地は

ある。

ただ、幅広く有用性が認められている、大域的な統計的シソーラスのより高度な構築法・活用法がないという現状においては、これまで一般的に利用されてきた素朴な方法を試した結果として検索性能の向上が観察されなかったという事実は、大域的分析の有用性に対する疑義を裏づけたといつてよいだろう。この点、今回の結果は、共起関係に基づく統計的シソーラスの有用性に否定的な先行研究[4][5]の主張を再確認したということになる。

自動適合フィードバックにも、もちろん、改善の余地はある。上で述べたように、この手法は、今回は上限値(理想)の約30%の改善を達成したのにすぎない。今後はこの比率を上げる工夫が必要になる。

ひとつの問題は、第1次の検索結果がうまくいかなかった場合の処置であろう。例えば、上位5文書を適合と見なすことにした場合に、第1次検索による上位5件に実際の適合文書がまったく含まれていなければ、自動適合フィードバックはむしろ第1次検索の結果を改悪する可能性がある。ORGの平均精度を横軸に、AT3とFB2の平均精度を縦軸として、今回使用した53件のトピックをプロットしたものを図4.2として示す。この図において、対角線よりも下に位置づけられている検索質問は、フィードバックの適用によって検索性能が低下したものである。当然、ORGの平均精度が低い場合にこのようなケースが多い(図4.2参照)。この点の改善がひとつの課題である。

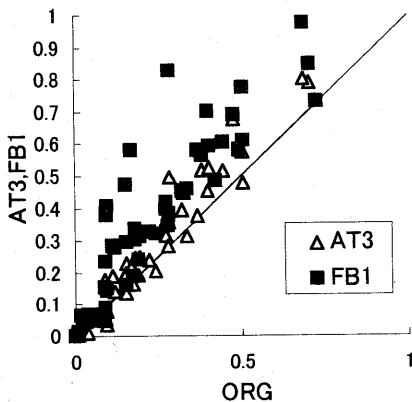


図4.2 ORGとAT3, FB1のプロット
(53質問)

また実際的な問題として、自動適合フィードバツ

クは2度検索を繰り返さねばならず、この点、応答速度に問題が生じる可能性がある。もちろん、統計的シソーラスを使う場合でも、シソーラスを参照する分、実行に余分な時間がかかるが、この点の検討も今後必要である。

5 おわりに

本稿では、クエリーの拡張のための大域的分析の例としての語の共起関係に基づく統計的シソーラスを活用する方法と、局所的分析の例としての自動適合フィードバックの検索性能を実証的に比較した。今回の検索実験の範囲内では、統計的シソーラスによる拡張は検索性能を低下させるのに対して、自動適合フィードバックは元の検索質問による実行に比べて、約16%の性能向上を示した。しかし、この数値は、理想値である真の適合フィードバックの増加率の約30%にすぎず、さらなる改善の余地のあることが示された。

謝辞

貴重なテストコレクションを準備され、研究目的での使用を認めていただいた、国立情報学研究所の皆様には感謝いたします。

参考文献

- [1] Efthimiadis, E. N.: Query expansion, *Annual Review of Information Science and Technology*, Vol.31, p.121-187 (1996).
- [2] Baeza-Yates, R. and Ribeiro-Neto, B.: *Modern Information Retrieval*, Harlow, England, Addison-Wesley, 1999.
- [3] Kando, N. and Nozue, T. eds: *Proceedings of the First NTCIR Workshop on Research in Japanese Text Retrieval and Term Recognition*, Tokyo, National Center for Scientific Information Systems, 1999.
- [4] Lesk, M.E.: Word-word associations in document retrieval, *Journal of Documentation*, Vol.20, No.1, p.8-36 (1969).
- [5] Peat, H.J. and Willett, P.: The limitation of term co-occurrence data for query expansion in document retrieval system, *Journal of the American Society for Information Science*, Vol.42, No.5, p.378-383 (1991).
- [6] Crouch, C.J.: An approach to the automatic construction of global thesauri, *Information Processing & Management*, Vol.26, No.5, p.629-640 (1990).
- [7] Schutze, Hinrich and Pedersen, Jan O.: A

- cooccurrence-based thesaurus and two applications to information retrieval, *Information Processing & Management*, Vol.33, No.3, p.307-318 (1997)
- [8]Chen, H. et al.: Automatic thesaurus generation for an electronic community system, *Journal of the American Society for Information Science*, Vol.46, No.3, p.175-193 (1995)
- [9]Park, Young C. and Choi, Key-Sun.: Automatic thesaurus construction using Bayesian networks, *Information Processing & Management*, Vol.32, No.5, p.543-553 (1996)
- [10]Srinivasan, Padmini: Query expansion and MEDLINE, *Information Processing & Management*, Vol.32, No.4, p.431-443 (1996)
- [11]Kim, Myoung-Cheol and Choi, Key-Sun.: A comparison of collocation-based similarity measures in query expansion, *Information Processing & Management*, Vol.35, p.19-30 (1999)
- [12]Mandala, Rila, Tokunaga Takenobu and Hozumi Tanaka: Query expansion using heterogeneous thesauri, *Information Processing & Management*, Vol.36, p.361-378 (2000)
- [13]Salton, G. and Buckley, C.: Improving retrieval performance by relevance feedback, *Journal of the American Society for Information Science*, Vol.41, No.4, p.288-297 (1990).
- [14]岸田和明:文献の適合度に関する目標値に基づくフィードバック手法, 情報処理学会研究報告, Vol.2001, No.21, p.189-196 (2001)
- [15]Buckley, C., Allan, J. and Salton, G. : Automatic routing and ad-hoc retrieval using SMART: TREC2, in D.K. Harman ed.: *The Second Text Retrieval Conference (TREC2)*, Gaithersburg, MD, National Institute of Standards and Technology, 1994. p.45- 55.