

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

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

情報検索における適合性フィードバックの手法としては Rocchio の方法がよく知られており、幅広く利用されている。しかし、この方法は、適合／不適合の 2 値でなされた適合判定にしか適用できず、利用者が自由に適合度の値を回答した場合には、このフィードバック情報をそのまま生かすことができない。本稿では、利用者から返された適合度の値を目標値として、その値に近い文献スコアをシステムが算出できるように、テイラー展開を用いて検索質問ベクトルを修正する手法を提案する。そして、NTCIR-1 のテストコレクションを用い、Rocchio の方法と本稿の方法との実証的な性能比較を試みる。

Feedback Method for Document Retrieval using Numerical Values on Relevance Given by Users

Kazuaki KISHIDA

Faculty of Cultural Information Resources, Surugadai University

698 Azu, Hanno, Saitama 357-8555

kishida@surugadai.ac.jp

The Rocchio method, a well-known approach for relevance feedback, is based on dichotomous relevance judgments, i.e., relevant or non-relevant. However, if numerical values representing the degree of relevance are given by users, this method is unable to make use of the relevance information effectively. This paper explores an alternative feedback method to solve this problem. The basic idea is to adjust the query vector by using Taylor formula so that the system returns retrieval status values near the numerical values on relevance. It is empirically shown that our method outperforms the Rocchio method, by using NTCIR-1 test collection.

1 はじめに

対話的な情報検索環境における検索性能の向上を目的とした適合性フィードバック (relevance feedback) の研究がこれまで数多く積み重ねられ、いくつかの手法が提案されてきた。なかでも、ベクトル空間型モデルに基づ

く Rocchio の方法[1][2]は、約 30 年前に開発されたものであるにも関わらず、その性能の高さと方法の簡明さから、現在でも文献検索(あるいは文書検索; document retrieval)における重要な手法として広く認められている。実際、最近の数多くの文献がこの手法を直接的・間接的に取り上げている[3][4][5][6][7]。

その手順の概略は次のとおりである。

- ① 利用者自身が作成した検索質問を使って文献を検索する（第1次の検索）。
- ② その検索結果に含まれる何件かの文献に対して利用者が適合判定（適合／不適合）をおこなう。
- ③ 判定結果に基づいて、検索質問ベクトルを修正する。具体的には、語句ごとに、適合文献群における重みの平均から不適合文献群における重みの平均を差し引いたものを、元の検索質問ベクトルに加える。
- ④ 修正された検索質問ベクトルを使って検索を再実行する（第2次の検索）。

段階③における検索質問ベクトルの修正方法は、より正確には、修正前のベクトルを \mathbf{q} 、修正後を $\tilde{\mathbf{q}}$ とかくと、

$$\tilde{\mathbf{q}} = \alpha \mathbf{q} + \frac{\beta}{|D_1|} \sum_{i:d_i \in D_1} \mathbf{d}_i - \frac{\gamma}{|D_0|} \sum_{i:d_i \in D_0} \mathbf{d}_i \quad (1)$$

である。ここで、 \mathbf{d}_i は文献 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、 $\mathbf{d}_i = (w_{i1}, \dots, w_{iM})^T$ である (M はデータベースに含まれる語句の異なり総数、 T は転置を示す記号)。また、 D_1 は適合文献の集合、 D_0 は不適合文献の集合、 α 、 β 、 γ はパラメータである。

Rocchio の方法の長所はその簡便さにある。その計算は適合文献および不適合文献における語句の重みの平均の加減にすぎない。そのため、システムへの実装は容易であり、また他の場面にも応用しやすい。事実、Rocchio の方法は、文献検索だけでなく、画像検索[8]やテキスト分類[9]の問題などにも広く応用されている。

しかしながら、(1)式から明らかなように、Rocchio の方法には、適合／不適合の2値での適合判定にしか対応できないという問題があり、この点からのさらなる検討の余地が残されている。このことは Rocchio の方法だけでなく、確率的なフィードバック手法[10]にもあてはまる。

本稿の目的はこの問題を解決するための新たなフィードバック手法の検討を試みることに

にある。そして本稿で考案した方法と Rocchio の方法の2つの手法の性能を、日本語のテストコレクションである NTCIR-1 を使って実証的に比較評価する。

2 適合度の目標値に基づくフィードバック手法

2.1 線形関数による検索モデル

代表的な情報検索モデルとして、ベクトル空間型モデル (vector space model) と確率型モデルとが一般によく知られている。前者については、最近の TREC での検索実験では、文献ベクトルと検索質問ベクトルの重みはそれぞれ

$$w_{ij} = \log x_{ij} + 1.0 \quad (2)$$

$$w_{qj} = (\log x_{qj} + 1.0) \log(N/n_j) \quad (3)$$

と設定されることが多い[11]。ここで、 x_{ij} は文献 d_i における語 t_j の出現頻度、 w_{qj} は検索質問ベクトルの要素 ($\mathbf{q} = (w_{q1}, \dots, w_{qM})^T$)、 x_{qj} は検索質問における語 t_j の出現頻度、 N はデータベース中の文献総数、 n_j は語 t_j が出現する文献数である。そして、検索質問に対する文献 d_i の類似度 s_i を、

$$s_i = \sum_{j=1}^M w_{ij} w_{qj} / \sqrt{\sum_{j=1}^M w_{ij}^2 \sum_{j=1}^M w_{qj}^2} \quad (4)$$

で計算する。

一方、確率型モデルとしては Okapi システムにおける方法[12]が著名であり、その方法の1つでは、 s_i を

$$s_i = \sum_{j=1}^M \left(\frac{2.2 x_{ij}}{1.2(0.25 + 0.75 l_i / \bar{l}) + x_{ij}} \times x_{qj} \log \frac{N - n_j + 0.5}{n_j + 0.5} \right) \quad (5)$$

とする。ここで、

$$l_i = \sum_{j=1}^M x_{ij}, \quad \bar{l} = N^{-1} \sum_{i=1}^N l_i$$

である (すなわち文献の長さ、およびデータベース全体でのその平均)。なお、(5)式は Okapi の方法の数多くのバリエーションのうちの1つにすぎない。

これらの2つの代表的な検索モデルは、

$$\mathbf{s} = f(\mathbf{b}) = \mathbf{A}\mathbf{b} \quad (6)$$

と表現できる。ここで、 $\mathbf{s} = (s_1, \dots, s_N)^T$ であり、 f は M 次元ベクトルを定義域、 N 次元ベクトルを値域とする関数 ($f: R^{M \times 1} \rightarrow R^{N \times 1}$)、 \mathbf{A} は $N \times M$ 行列で、その要素 a_{ij} は、ベクトル空間型(2)~(4)式の場合には、

$$a_{ij} = (\log x_{ij} + 1.0) / \sqrt{\sum_{j=1}^M (\log x_{ij} + 1.0)^2} \quad (7)$$

確率型の場合には、例えば、

$$a_{ij} = \frac{2.2x_{ij}}{1.2(0.25 + 0.75l_i/l_j) + x_{ij}}$$

である。また \mathbf{b} は M 次元列ベクトルであって、要素 b_j ($j = 1, \dots, M$) はベクトル空間型の場合には

$$b_j = w_{qj} / \sqrt{\sum_{j=1}^M w_{qj}^2} \quad (8)$$

ただし $w_{qj} = (\log x_{qj} + 1.0) \log(N/n_j)$ 、確率型

の場合には $b_j = x_{qj} \log \frac{N - n_j + 0.5}{n_j + 0.5}$ である。

以上のように、現在の代表的な2つの検索モデルはともに検索質問ベクトル \mathbf{b} についての線形関数として捉えることができる。

2.2 検索システムの目的とフィードバック

実際には、(6)式のベクトル \mathbf{s} によって文献の順位が決められて出力されることになる。これは $\mathbf{s} = f(\mathbf{b})$ を「真の」適合度ベクトル $\mathbf{r} = (r_1, \dots, r_N)^T$ の推計値として考えていることにほかならない。とすれば、検索システムの目的は、与えられた検索質問に対して、「真の」ベクトル \mathbf{r} に最も近い \mathbf{s} を求めるということになる。これを \mathbf{s}^* とかき、 \mathbf{s}^* を求めることのできる最適な質問ベクトルを \mathbf{b}^* とする（すなわち $\mathbf{s}^* = f(\mathbf{b}^*)$ ）。

もちろん \mathbf{r} は未知である。しかし、適合性フィードバックによって部分的に知ることができる。例えば、第1次の検索結果の上位 n 件に対して適合度の値を利用者に回答してもらえばよい。この n 件の文献集合を X と表記し、この集合に限定された \mathbf{r} を \mathbf{r}_X とすれば、この \mathbf{r}_X だけは利用者からのフィードバックによっ

て知ることが可能である。

これに対応して、

$$\mathbf{s}_X = f_X(\mathbf{b}) = \mathbf{A}_X \mathbf{b} \quad (9)$$

と定義する (\mathbf{A}_X は $n \times M$ 行列で、 \mathbf{A} から X に含まれる文献の行のみを取り出したもの、 \mathbf{s}_X は n 次元列ベクトル、 $f_X: R^{M \times 1} \rightarrow R^{n \times 1}$)。また、 \mathbf{r}_X と \mathbf{s}_X との距離を $\phi(\mathbf{r}_X, \mathbf{s}_X)$ とかく。これらの記号を使えば、 \mathbf{r}_X が与えられた場合の適合性フィードバックの目的は、結局、

$$\tilde{\mathbf{b}} = \arg \min_{\mathbf{b}} \phi(\mathbf{r}_X, f_X(\mathbf{b})) \quad (10)$$

を求めることに帰着する（ただし、 $\tilde{\mathbf{b}}$ が \mathbf{b}^* に一致する保証はない）。

この(10)式の $\tilde{\mathbf{b}}$ を求めることのできるフィードバック手法は、各文献の適合度が数値で与えられた場合にその情報を十分に生かすことのできる方法にほかならない。したがって、そのような方法こそが本稿で求めているものである。

2.3 テイラー展開による解法

ここでは(10)式を解くための近似的な方法を提案する。まず、関数 $f_X(\tilde{\mathbf{b}})$ は一般的にはテイラー展開によって、

$$f_X(\tilde{\mathbf{b}}) = f_X(\mathbf{b}) + \frac{\partial f_X(\mathbf{b})}{\partial \mathbf{b}^T} (\tilde{\mathbf{b}} - \mathbf{b}) + K \quad (11)$$

とかける (K は2次以上の項) [13]。ここで考える検索モデルを線形関数(6)式に限れば、

$$\frac{\partial f_X(\mathbf{b})}{\partial \mathbf{b}^T} = \frac{\partial (\mathbf{A}_X \mathbf{b})}{\partial \mathbf{b}^T} = \mathbf{A}_X$$

と計算されるから、(11)式は、

$$f_X(\tilde{\mathbf{b}}) = f_X(\mathbf{b}) + \mathbf{A}_X (\tilde{\mathbf{b}} - \mathbf{b}) \quad (12)$$

となる ($K=0$ に注意)。

次に、(11)式中の \mathbf{b} を第1次検索で用いたベクトルとし、また $\mathbf{r}_X = f_X(\tilde{\mathbf{b}})$ が成立すると仮定する。第1次の検索結果は $f_X(\mathbf{b}) = \mathbf{s}_X$ であるから、(12)式は、

$$\mathbf{r}_X = \mathbf{s}_X + \mathbf{A}_X (\tilde{\mathbf{b}} - \mathbf{b})$$

となり、これを变形すれば、

$$\mathbf{A}_X (\tilde{\mathbf{b}} - \mathbf{b}) = \mathbf{r}_X - \mathbf{s}_X \quad (13)$$

$$\therefore \tilde{\mathbf{b}} = \mathbf{b} + \mathbf{A}_X^{-1} (\mathbf{r}_X - \mathbf{s}_X) \quad (14)$$

を得る。 \mathbf{A}_X^{-1} は $M \times n$ 行列で、 $\mathbf{A}_X^{-1} \mathbf{A}_X = \mathbf{I}_M$ が

成り立つものとする (\mathbf{I}_M はすべての対角要素が1, その他は0である $M \times M$ 行列)。この $\tilde{\mathbf{b}}$ は文献集合 X による \mathbf{b}^* の近似値にすぎないが, X が適合度に関してデータベース全体の状況を正しく反映しているとすれば, フィードバック後に(14)式で計算される $\tilde{\mathbf{b}}$ を用いた第2次検索の性能は向上するはずである。

(14)式から明らかなように, この方法では, n 件の文献に対する適合度 \mathbf{r}_x と第1次の検索結果 \mathbf{s}_x との差に \mathbf{A}_x^{-1} を掛けたものを元のベクトル \mathbf{b} に加えることによって, \mathbf{b} を修正する。元のベクトルに何らかの値を加算 (または減算) することによって修正する点は Rocchio の方法(1)式と同様である。しかし, 本稿のテイラー展開による方法では「修正後のベクトルによる再計算結果 $f_x(\tilde{\mathbf{b}})$ が適合度の目標値 \mathbf{r}_x に等しくなるように修正する」という基準が設定されている点が異なっているわけである。

2.4 特異値分解の応用

(14)式中の \mathbf{A}_x^{-1} は一種の逆行列であるが, \mathbf{A}_x は正方行列ではなく $n \times M$ 行列であるため, 通常の手順では \mathbf{A}_x^{-1} を求めることはできない。そこで特異値分解 (SVD) を利用する。ここでは特異値分解の計算アルゴリズム[14]の都合上, 「縦長」の \mathbf{A}_x^T に対して特異値分解を施す (文献検索の状況では文献数よりもそれに含まれる語数のほうが常に大きくなることに注意)。結果として $\mathbf{A}_x^T = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T$ と分解され (\mathbf{U} は $M \times n$ の直交行列, $\mathbf{\Lambda}$ は $n \times n$ の対角行列, \mathbf{V} は $n \times n$ の直交行列), 最終的に $\mathbf{A}_x = (\mathbf{A}_x^T)^T = \mathbf{V}\mathbf{\Lambda}\mathbf{U}^T$ を得る。これを(13)式に代入すれば,

$$\mathbf{V}\mathbf{\Lambda}\mathbf{U}^T(\tilde{\mathbf{b}} - \mathbf{b}) = \mathbf{r}_x - \mathbf{s}_x$$

となるが, \mathbf{V} と \mathbf{U} は直交行列であるから,

$$\tilde{\mathbf{b}} - \mathbf{b} = \mathbf{U}\mathbf{\Lambda}^{-1}\mathbf{V}^T(\mathbf{r}_x - \mathbf{s}_x)$$

$$\therefore \tilde{\mathbf{b}} = \mathbf{b} + \mathbf{U}\mathbf{\Lambda}^{-1}\mathbf{V}^T(\mathbf{r}_x - \mathbf{s}_x) \quad (15)$$

と計算できる。(15)式を使えば, 本稿のフィードバック手法の実装が可能になる (なお, 実際には, M の代わりに, X における語の異なり

総数 M' を使って特異値分解をおこなえば十分である。当然, X に含まれない語についての修正は何ら施されない)。

3 検索実験による性能比較

ここでは, NTCIR のテストコレクションを使って, ベクトル空間型モデルに基づく Rocchio の方法とテイラー展開による方法 (14)あるいは(15)式の性能を比較評価する。

3.1 テストコレクションの概要

本稿で用いるテストコレクションは国立情報学研究所 (NII) の検索実験プロジェクトによる NTCIR-1 の Jコレクションである[15]。これは, さまざまな主題領域の国内学会が主催する全国大会や研究大会などでの発表要旨 (抄録) から成るデータベースから作成されたものである (Jコレクションには日本語の標題・抄録等から成るレコードのみが含まれる)。検索対象フィールドは標題・抄録・著者キーワードとする。

検索質問 (topic) については 031-083 までを使う。ただし, 後述するように今回の実験には適さないものが何件もあり, 最終的に, 48 件の検索質問を使用する。検索質問ベクトルの自動生成に用いるフィールドは <TITLE>, <DESCRIPTION>, <NARRATIVE>, <CONCEPT> である (すなわち「長い」検索質問とする)。

3.2 実験の手順と方法の詳細

実験手順の概略は以下のとおりである。

- ① 第1次の検索をベクトル空間型モデルでおこなう。具体的には, (2), (3)および(4)式を用いる。この検索の実行を, 便宜上, VECORG と表記する。
- ② 上記①の検索結果の上位 10 件の文献について, テストコレクションの適合判定ファイルから適合/不適合の情報を抽出し, Rocchio の方法(1)式, 本稿でのテイラー展開に基づく方法(15)式で, それぞれ検索質問ベクトルを修正する。したがって, 今回は $n = 10$ の大きさの文献集合

X を使ってフィードバックを試みることになる。

- ③ 修正された検索質問ベクトルを用いて検索を再実行する(それぞれ, ROCCHIO, TAYLOR と表記する)。2つの結果を, 平均精度 (average precision) と再現率-精度グラフを使って比較する。

なお, (2)式において $x_{ij} = 0$ ならば $w_{ij} = 0$ とおく ((3)式も同様)。また, ②の段階で, 上位10件中に適当文献がまったく含まれない検索質問, あるいは, 不適当文献がまったく含まれない検索質問は, 今回はいちおう実験の対象外とした (このような検索質問が5件あった)。

より具体的には, Rocchio の方法では, (1)式の \mathbf{d}_i の要素は(2)式をそのまま使い, パラメータは先行研究[6]に従って $\alpha = 8$, $\beta = 16$, $\gamma = 4$ と設定する。そしてこの結果計算された(1)式の \mathbf{q} をそのまま(4)式に投入して各文献の s_i を算出する (③の段階)。

テイラー展開による方法では, \mathbf{A} の要素は(7)式, \mathbf{b} の要素は(8)式を使う。そして(15)式で \mathbf{b} を求め, ③の段階で $\mathbf{s} = \mathbf{A}\mathbf{b}$ によって各文献の s_i を再計算する。この過程で問題となるのは \mathbf{r}_X である。本来ならば, n 件の文献に対して利用者が回答した適合度の数値を \mathbf{r}_X とすべきである。しかし, テストコレクションにはこのような数値は用意されておらず, 実際には適合/不適合の2値のみしか利用できない。そこで, 今回の実験では, 次の方法によって, 適合/不適合の2値の情報から「自動的に」 \mathbf{r}_X を設定することにする。

- (a) 適合文献 $d_i (\in X)$ に対しては, 第1次の検索で(4)式によって計算された s_i についての適合文献中での最大値 s_{\max} を使って, $s_i + (1.0 - s_{\max})$ を \mathbf{r}_X の要素とする。
- (b) 不適合文献 $d_i (\in X)$ に対しては, 第1次の検索で(4)式によって計算された s_i についての不適合文献中での最小値 s_{\min} を使って, $s_i - s_{\min}$ を \mathbf{r}_X の要素とする。

この方法によれば, 適合文献中で最大の s_i を持つ文献 d_i については $r_i = 1.0$, 不適合文献中で最小の s_i を持つ文献 d_i については $r_i = 0.0$ と

なる。利用者による初期の検索質問ベクトルに負の重みが含まれていないとすれば, (4)式によって計算される s_i の最大値は 1.0, 最小値は 0.0 である。したがって, 上記の方法は, X 中で最も適合している文献の適合度を第1次検索での理論的上限值, 最も適合していない文献の適合度を理論的下限值として, 各文献の第1次検索の値 s_i に応じて, \mathbf{r}_X を設定していることになる。

この設定方法は, テイラー展開による(14)式が本来持っている応用可能性を制限している。しかし, 現実の利用者にとっては, 適合度の数値を割りふるよりも, 適合か不適合かだけを回答したほうが負担がかかって少ないという状況も考える。したがって, 上記の方法を実験で試すことには, 単なる検索性能の確認だけでなく, 実用的な意味での価値をも見出すことができよう。

3.3 日本語テキストからの語句の抽出

実験のためには, テストコレクション中の文献と検索質問に含まれる日本語テキストから語 t_j を抽出しなければならない。このためには一般的な機械可読辞書との照合による手法を用いる。つまり, レコード中のテキストと辞書とを突き合わせていき, テキスト中の文字列と一致する辞書の見出し語のうち最長のものを語句として採用する (いわゆる最長一致法)。また, テキスト中の一致しない部分に関しては, 同一文字種の連続を語句として見なすこととする。例えば, 「情報検索システムの研究」という文があつて「情報」のみが辞書に含まれる場合には, 「情報」がまず切り出され, 次に, 漢字の連続として「検索」が識別され, さらに, カタカタの連続として「システム」, 同様に「の」と「研究」がそれぞれ抽出される。ここで, ひらがなは落とすこととする。そこで, この段階での抽出結果は「情報」「検索」「システム」「研究」となる。

機械可読辞書としては, 形態素解析システム「茶筌」[16]で使われているものを利用させていただいた。ただし, この辞書に含まれている品詞情報などは使用せず, この辞書の見出し語のみを抽出して, 照合用の辞書とする。した

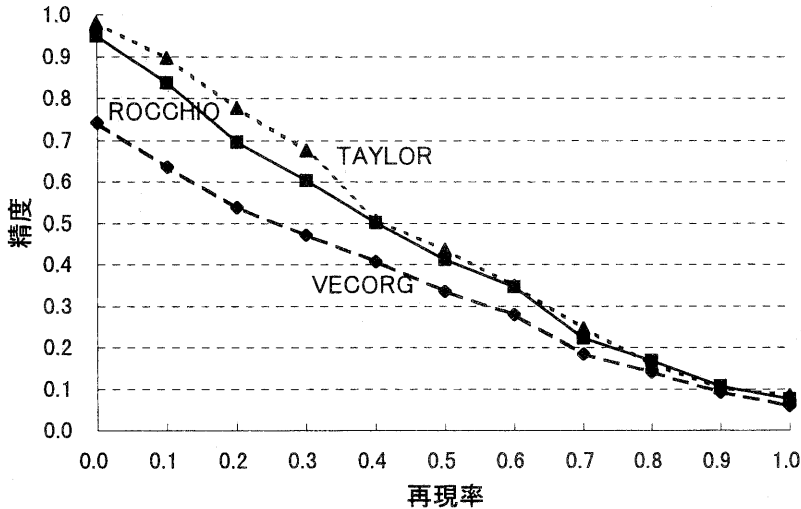


図1 再現率—精度グラフ

がって、本稿で使用するのは辞書というよりも単なる語のリストである。

問題は「茶筌」の辞書には各主題領域の専門用語が含まれていないことである。したがって、上記の語句の識別方法だけを用いると、一般性の高い語のみしか抽出されない可能性がある。そこで、上記の方法で切り出された語句のうち、隣接する2語を機械的に組み合わせ、複合語を自動生成する。ただし、2つの語句の間にひらがなや記号が挟まる場合には、複合語を生成しない。したがって、上記の例では結果的に、「情報」「検索」「システム」「研究」「情報検索」「検索システム」の6語が切り出されることになる。そして、これらを語 t_i として用いることとする。

3.4 実験結果

NTCIR-1のJコレクションに含まれるレコード総数(=N)は332,918件、語の出現延べ総数は39,284,817回であった(1件あたり118.0回)。

表1 各手法の平均精度

VECORG (第1次検索)	ROCCHIO (第2次検索)	TAYLOR (第2次検索)
.342	.431	.461

再現率—精度グラフを図1に示す。また、平均精度の48検索質問での平均値(mean average precision)を表1に示す。平均精度については、本稿の方法(TAYLOR)がわずかにRocchioの方法(ROCCHIO)を上回った。第1次の検索(VECORG)と比較して、テイラー展開による方法は約38%の性能の増加、Rocchioの方法は26%の増加である。

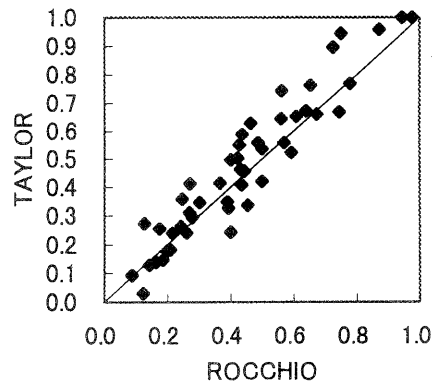


図2 検索質問ごとの平均精度

次に、各検索質問ごとに両手法の性能を比較する。図2は、ROCCHIOによる平均精度をx軸の値、TAYLORによる平均精度をy軸の値

として、48 の検索質問をプロットした図である。したがって、対角線よりも上側の点は TAYLOR が優れている検索質問を示し、逆に下側は ROCCHIO が優っている検索質問である。この図からは、検索質問によって、TAYLOR が優れている場合と ROCCHI が優れている場合とに分かれていることがわかる。つまり、表 1 で示された本稿の方法 (TAYLOR) の優位性はすべての検索質問に対して個別的に成立するわけではなく、その優位性は検索質問で全体平均した場合に観察されるにすぎない。

すでに述べたように、テストコレクションでは 2 値での適合判定の情報が利用できるのみであるため、今回の実験では両手法が使用する「学習のための情報」に大差がなく、その結果、検索性能に大きな相違が生じないのかもしれない。しかし、そのような条件で、本稿の方法の性能がマクロ的にも Rocchio の方法を上回ったということから、適合度に関するより豊富な情報が与えられた場合での本稿の方法の有効性には十分な期待が持てると考えられるだろう。

4 おわりに

本稿では、Rocchio の方法に代わりうる新しい適合性フィードバックの手法を提案した。この手法は、フィードバックの結果として利用者から与えられた各文献の適合度の値を目標値として設定し、その目標値になるべく近い値がシステムによって計算されるように、Taylor 展開を利用して検索質問ベクトルに修正を施す方法である。そして、NTCIR-1 の J コレクションを使って、この手法の性能が Rocchio による方法のそれを平均的には上回ることを実証的に確認した。ただし、その差は顕著なものではなく、検索質問によっては Rocchio の方法の性能が優れている場合もあった。なお、Rocchio の方法はクラスター化などを併用すると性能がより向上することなどが報告されているが[7]、今回の本稿の比較はそのような工夫なしでの単純な比較に留まるものである。

本稿の手法は、適合/不適合の 2 値だけでな

く、自由に設定された適合度の値に対応できる点が大きな長所である。しかし、テストコレクションの制約から、今回の実験では、適合/不適合の 2 値でのフィードバック情報が与えられた状況に限定した。フィードバックによって適合度の値が実数として与えられたときに、本稿での手法がどの程度検索性能を改善するかは今後の研究課題である。

謝辞

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

参考文献

- [1] Rocchio, J.J. Jr.: Relevance feedback in information retrieval, in G. Salton ed.: *The SMART Retrieval System: Experiments in Automatic Document Processing*, Englewood Cliffs, New Jersey, Prentice-Hall, 1971. p313-323.
- [2] 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).
- [3] Buckley, C., Allan, J. and Salton, G.: Automatic routing and retrieval using SMART: TREC-2, *Information Processing and Management*, Vol.31, No.3, p.315-326 (1995).
- [4] Sarinivasan, P.: Query expansion and MEDLINE, *Information Processing and Management*, Vol.32, No.4, p.431-443 (1996).
- [5] Lee, J.H.: Combining the evidence of different relevance feedback methods for information retrieval, *Information Processing and Management*, Vol.34, No.6, p.681-691 (1998).
- [6] Mandala, R., Tokunaga, T. and Tanaka, H.: Query expansion using heterogeneous thesauri, *Information Processing and Management*, Vol.36, p.361-378 (2000).
- [7] 岩山真: 適合性フィードバックの効率化について, 情報処理学会研究報告, Vol.2000, No.29, p.1-8 (2000).
- [8] Ciocca, G. and Schettini, R.: A relevance feedback mechanism for content-based image retrieval, *Information Processing and Management*, Vol.35, p.605-632 (1999).
- [9] Moens, Marie-Francine and Dumortier,

- Jos.: Text categorization: the assignment of subject descriptors to magazine articles, *Information Processing and Management*, Vol.36, No.6, p.841-861 (2000).
- [10]Harper, D.J. and van Rijsbergen, C.J.: An evaluation of feedback in Document Retrieval using co-occurrence data, *Journal of Documentation*, Vol.34, No.3, p.189-216 (1978).
- [11]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. <http://trec.nist.gov/>
- [12]Robertson, et al.: Okapi at TREC-3, D.K. Harman ed.: Overview of *the Third Text Retrieval Conference (TREC-3)*, Gaithersburg, MD, National Institute of Standards and Technology, 1995. p.109- 126.
- [13]Harville, D.A.: *Matrix Algebra from a Statistician's Perspective*, New York, Springer, 1997.
- [14]Press, W.H. et al.: *Numerical Recipes in C*, 2nd ed., New York, Cambridge Univ. Press, 1992.
- [15]<http://research.nii.ac.jp/ntcir/>
- [16]<http://cl.aist-nara.ac.jp/lab/nlt/chasen/>