

論文 / 著書情報  
Article / Book Information

論題(和文)	学習用Web情報の検索支援システムの開発
Title(English)	A development of information retrieval support system for learning material on the Web
著者(和文)	森本容介, 中山実, 清水康敬
Authors(English)	YOUSUKE MORIMOTO, MINORU NAKAYAMA, YASUTAKA SHIMIZU
出典(和文)	教育システム情報学会誌, Vol. 17, No. 3, pp. 231-240
Citation(English)	Transactions of Japanese Society for Information and Systems in Education, Vol. 17, No. 3, pp. 231-240
発行日 / Pub. date	2000,

## 学習用Web情報の検索支援システムの開発

森本容介\*, 中山実\*, 清水康敬\*

## A development of information retrieval support system for learning material on the Web

Yousuke MORIMOTO\*, Minoru NAKAYAMA\*, Yasutaka SHIMIZU\*

This paper describes development of information retrieval support system for educational materials on the web, and its performance of the retrieval. The feature vectors of the term and the subject and grade were extracted from the teaching guide line for the primary school by singular value decomposition (SVD). The self-organizing maps for the terms and the subject and grade were trained by both feature vectors. The developed information retrieval support system displays the terms map and the subject map, and terms in high similarity according to a term which is assigned by users. The performance of this system is evaluated by variances of precision score for a subset of educational documents on the web. The variance increase with addition of several terms which show higher similarity to user assigned term. It suggests this system extracts intended documents by assisting terms selection for searching.

キーワード：学習情報，WWW，情報検索，学習指導要領，自己組織化マップ

## 1. はじめに

小中高等学校へのインターネット接続が急速に進んでいる。授業におけるインターネット利用の一つとして、Web情報を教材として利用することがあげられる。児童・生徒が主体となってWeb情報を参照したり、教員が教材研究のためにWeb情報を参考にするなどである。このようなWeb情報の検索には、gooなどの全文検索型サーチエンジンを使って情報検索が行われるが、適切な検索結果を得ることは容易ではない<sup>(1)</sup>。

これは、検索語が含まれるページが非常に多く得られて、検索者が求める教科学年に関連する情報がごく

わずかであるためである。この問題は、検索者が与える検索語が、検索目的の教科や学年に関連するWeb情報の検索語として適切でないことが原因の一つと考えられる。本研究では、検索者は対象とする教科や学年に関連した学習に利用可能な情報を求めていると考え、対象とする教科や学年に着目した検索を検討することにした。

情報検索の効率化を図る方法として、検索される情報の特徴を基にクラスタリングし、類似情報の収集を容易にしたシステムが検討されている (WEBSOM<sup>(2)</sup>)。ただし、これはインターネットニュースを対象としており、学習情報は対象にしていない。そこで本研究では、学習利用できるWeb情報の検索を対象として、対象とする教科や学年に着目した検索支援を検討した。具体的には、検索に必要な検索語の選定が容易かつ適

\* 東京工業大学教育工学開発センター  
CRADLE, Tokyo Institute of Technology

切な行えるシステムを開発し、その有効性を明らかにすることを目的とした。この目的のために、以下の手順で研究を行った。

- (1) 学習情報の特徴を抽出するために、学習指導要領の記述から各教科学年や用語の特徴ベクトルを算出し、自己組織化マップを用いて、クラスタリングを行った。
- (2) 自己組織化マップと用語間類似度を用いて、検索対象に適した検索語を明示できる検索支援システムを開発した。
- (3) 開発した検索支援システムによる効果を評価、検討した。

## 2. 学習用語の抽出と自己組織化マップの作成

検索支援に必要な用語辞書<sup>(3)</sup>と自己組織化マップを以下の要領で作成し、これを検索支援システムに組み込んだ。

### 2.1 用語辞書の作成

#### 2.1.1 用語の頻度行列の作成

現行の小学校学習指導要領<sup>(4)</sup>の全文を電子化し、学年・教科と用語の特徴を抽出した。電子化された全文を形態素解析して語句に分け、その品詞情報も得た。形態素解析には、「茶筌Ver1.5」<sup>(5)</sup>を用いた。その結果、用語として1919語が抽出され、これを分析対象とした。これらの用語の頻度について、教科の目標や学年ごとの記述を53項目に、教科ごとに再分類し、1919×53の出現頻度表を作成し、これを分析に用いた。図1は、抽出された語句が各学年・教科の記述に出現する頻度の一部を示している。

この表において、縦方向には1919の用語を列挙し、横方向には学年・教科として53項目を示している。ここで、学年・教科は、以下に列挙する53項目である。国語、算数、音楽、図工の24項目(1年～6年×4教科)、社会、理科の8項目(3年～6年×2教科)、生活(1年～2年)の2項目、家庭(5年～6年)の2項目、道徳、体育の6項目(1年～6年までの2学年ごと×2教科)、特別活動(全学年共通)の1項目、教科の目標と指導計画(特別活動を除く10教科)の10項目、合計53項

目である。

なお、教科の目標と指導計画は、本論文では両者をまとめて指導目標と呼ぶことにする。

### 2.1.2 特徴ベクトルの生成

前項で得られた出現頻度表から、用語や教科などの特徴ベクトルを得るために、特異値分解<sup>(6)</sup>を行った。ここでは、用語の頻度を考慮して分析するために、各頻度を平方変換<sup>(7)</sup>して分析に用いた。特異値分解の概要を、図2に示す。用語の出現表は、特異値分解によ

学年 教科 用語	1年 国工	2年 国工	3年 国工	4年 国工	5年 国工	6年 国工	図工 指導 目標	1,2年 道徳	...
あいさつ	0	0	0	0	0	0	0	1	...
...	...	...	...	...	...	...	...	...	...
水上	0	0	0	0	0	0	0	0	...
表	0	0	0	0	0	0	0	0	...
表す	3	2	2	2	3	3	0	0	...
表する	1	3	5	4	12	12	5	0	...
表れる	0	0	0	0	4	4	0	0	...
表記	0	0	0	0	0	0	0	0	...
表現	1	1	2	2	3	3	9	0	...
表情	0	0	0	0	0	0	0	0	...
...	...	...	...	...	...	...	...	...	...

図1 出現頻度表

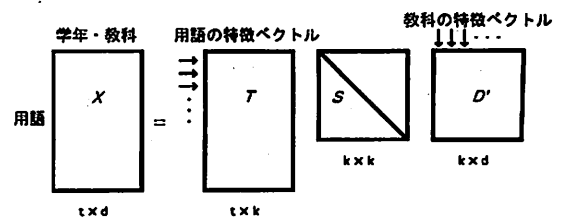


図2 特異値分解の概略

って、用語、学年・教科それぞれの特徴を示す行列2つと、固有値を示す行列の3つに分解できる。

すなわち、1919の用語と53の学年・教科の出現表(行列X)を特異値分解すると、3つの行列T, S, D'の行列の積で表される。

$$X = TSD' \quad (1)$$

このうち、行列Sは対角行列で、その成分 $\mu_i$ は行列の積 $X'X$ の固有値である。行列Tは1919×53型、行列D'は53×53型である。これらの行列の成分を、それぞれの特徴とすることができる。

例えば、行列Tは1919の用語に対して、それぞれ53成分の表現が行われている。これを用語の特徴ベクトルとして扱う。同様に、行列D'は53の教科学年について、53成分の表現が得られ、教科学年の特徴ベクトルとする。

Deerwesterらは、この様に導出した文書の特徴ベクトルと単語の特徴ベクトルを用いて、文献検索が行えることを示した<sup>(6)</sup>。具体的には、特徴ベクトルのベクトル内積が、コサイン類似度として用いられた。すなわち、2つのベクトルを $v_u, v_v$ とすると、コサイン類似度 $\cos \theta = (v_u, v_v) / |v_u| |v_v|$ である。 $|v_u|, |v_v|$ は各ベクトル長である。ここで、 $v_u, v_v$ の大きさを1とすると、コサイン類似度はベクトル間の内積値に一致する。

## 2.2 自己組織化マップの作成

学習指導要領から抽出された用語は、各教科で用いられていることから、前節で求めたそれぞれの特徴ベクトルによって学習内容が構成されていると考えることができる。これらの特徴ベクトルがどのような空間を構成しているかを視覚的に捉えるために、自己組織化マップによって非階層クラスタリングを行い、ベクトル間の関係を2次元空間に次元圧縮表示した。

自己組織化マップは、対象とする特徴ベクトルのパターンをコードブックベクトルとして学習する方法で、パターンの教師なし学習の一種である。すなわち、2次元上の格子点に配置されたベクトルであるコードブックベクトルは、特徴ベクトルの傾向を学習する。これを用いて特徴ベクトルをマップ上に配置すると、近傍空間に類似する特徴ベクトルが集まるように表現できる。特徴ベクトルによるコードブックベクトルの学

習には、特徴ベクトル間の非類似度としてユークリッド距離を用いるのが一般的である。

中山ら<sup>(8)</sup>は、ユークリッド距離によって学習指導要領に基づいてマップを作成する方法を提案した。また、マップの作成は初期状態を乱数で与えて、学習を行わせた。しかし、用語の頻度に基づく特徴ベクトルでは、前項で述べたようにベクトル間内積による類似度が用いられる。また、後述のように用語間や教科学年との関連性は、すべてベクトル間内積による類似度によって検討される。ユークリッド距離とベクトル間内積は異なる尺度であることから、自己組織化の学習においても、ベクトル間内積による類似度を用いて行った方が対応関係が明確になると考えられる。また、教科学年の関係を初期状態として与えた方が、より実際の関係を反映したマップが作成されると考えられる。すなわち、学習のアルゴリズムを先行研究の方法に変更を加えてマップを作成した方が、本研究の目的に合致したマップが得られると期待される。

自己組織化の学習には、KohonenらのSOM\_PAK3.1<sup>(9)</sup>を用いた。ただし、上述のように類似度としてベクトル間内積を用いるため、コードブックベクトル $m_i$ の学習則が以下の式で行われるようにプログラムを修正して用いた。この学習にあたっては、学習データである特徴ベクトルは、常にそれぞれベクトル長を1に規格化した。これによって、ベクトル間の内積値がコサ

$$m_i(t+1) = \begin{cases} \frac{m_i(t) + \alpha(t)T_j(t)}{\|m_i(t) + \alpha(t)T_j(t)\|} & i \in N_c(t) \\ m_i(t) & i \notin N_c(t) \end{cases} \quad (2)$$

イン類似度になる。

ここで、 $T_j$ は用語の特徴ベクトル、 $N_c$ は近傍集合、 $\alpha$ は学習係数である。

特徴ベクトルの学習は、次の手順で行った。

- (1) 教科学年の特徴ベクトルの2成分( $d_{1i}, d_{2i}$ )を用いて、2次元空間に各教科学年の特徴ベクトルを配置する。すなわち、53次元空間で表現される特徴ベクトルについて、2成分で張られた空間での各教科学年の関係を初期値とする。
- (2) 用語の特徴ベクトルを用いて学習する。ある特徴ベクトルの配置される近傍範囲を比較的広くして、そこにあるコードブックベクトルの学習をするため



る。CGIプログラムはPerlで作成し、後述するシステム上でのチェックボックスによる用語の選択や自己組織化マップの操作の機能はJava Scriptを用いて実現した。図3は、「宇宙」をキーワードとして本システムで検索を行ったときに表示される画面の一例である。システムは、自己組織化マップの図2つと用語リストおよび検索用語リストの入力行を表示する。

### 3. 3 自己組織化マップの表示

ユーザが検索語1語を入力すると、本システムは、自己組織化マップを2枚表示する。

この2枚の自己組織化マップは、同じ領域を拡大標記したものであるが、マップ上に配置したラベルが異なる。1枚は、用語の類似関係を示すために用語のラベルをマップ上に配置した。もう1枚は、教科学年の類似関係を示すために教科学年のラベルをマップ上に配置した。すなわち、検索しようとしている用語に類似する用語をマップ上で表示するとともに、もう一方のマップでは、それらの用語と類似関係にある教科学年に対応していることを示している。

ユーザはブラウザ上のボタンをクリックすることにより、両方のマップの拡大、縮小や表示範囲の移動を自由に行うことができる。これによって、入力語に近い関係用語群を知ることができる。また、入力語がどの学年・教科と関連が高いかを視覚的に見ることができる。

### 3. 4 類似用語の表示

検索語の入力を受け付けると、CGIプログラムは類似度が高い用語のリストを、類似度の数値とともに表示する。

本システムで表示する自己組織化マップは、Ultschら<sup>(10)</sup>の方法による表示を用いた。この表示方法では、マップのコードブックベクトル間の類似度を、グレースケールで表現している。すなわち、色の濃い部分は類似度が低く、色の薄い部分は類似度が高いことを示している。このため、検索語の近傍にある用語間でも、配置される位置のマップ上でのグレースケールの濃淡によって、実際の類似度が異なる。そこで、正確に類似度を示すために、図中の表示だけでなく類似度が高い用語のリストも示している。

図3の類似用語の表示を、分かりやすくするために、表1に類似用語とその類似度の一覧を示す。

### 3. 5 検索語の選択と検索オプションの付加

ユーザは、表示された用語リストやマップの情報を基に、既存のサーチエンジンで検索する用語群を決定する。抽出された用語に関しては、各用語の先頭に付けられたチェックボックスをクリックすることによって、簡単に選択や追加することができる。抽出された用語以外にも、必要に応じて手入力によって任意の用語を、検索の用語群に加えることができる。最後に、ユーザは検索オプションを選択し、既存のサーチエンジンで実際に検索を行う。検索オプションは、サーチ

表1 「宇宙」の近傍用語

順位	抽出された用語	コサイン類似度
1	地球	1.000
2	エネルギー	0.964
3	物質	0.964
4	問題	0.951
5	見いだす	0.815
6	水	0.723
7	空気	0.657
8	違う	0.647
9	作れる	0.643
10	食べ物	0.643
11	動物	0.638
12	生物	0.616
13	植物	0.599
14	種子	0.588
15	出る	0.588

表2 交差行列

	検索された文書	検索されなかった文書
適合文書	$w$	$x$
非適合文書	$y$	$z$

エンジンの種類や, AND・OR検索の選択である。

#### 4. 検索支援システムの評価

##### 4.1 検索性能の指標

一般に検索システムの性能評価には, 精度と再現率が用いられる。精度と再現率を簡単に説明する。有限個の文書情報を検索し, その検索結果は表2の様に分類される。

この時, 再現率 (recall)  $R$ と精度 (precision)  $P$ は次のように定義される。

$$\text{再現率: } R = \frac{w}{w+x} \quad (4)$$

$$\text{精度: } P = \frac{w}{w+y} \quad (5)$$

すなわち, 再現率は検索された適合文書数 ( $w$ ) と検索対象となる文書集合中の適合文書数 ( $w+x$ ) の比であり, これが大きい方が検索漏れが少ないよいシステムであるといえる。一方, 精度は検索された適合文書数 ( $w$ ) と検索された文書数 ( $w+y$ ) の比であり, こちらも大きい方が不要な検索結果が少ないよいシステムであるといえる。

検索システムの総合的な性能を評価するための指標として, マクロ平均がある。マクロ平均は, 検索質問ごとに再現率, 精度を計算し, それらの値をすべての検索質問にわたって平均したものである。再現率, 精度のマクロ平均  $\bar{R}$  と  $\bar{P}$  はそれぞれ式6, 式7で計算できる。ただし,  $Q$  は検索要求の総数,  $w_i, x_i, y_i$  は表2と同じ意味で, 添字は検索質問を表している。

$$\bar{R} = \frac{1}{Q} \sum_{i=1}^Q \frac{w_i}{w_i + x_i} \quad (6)$$

$$\bar{P} = \frac{1}{Q} \sum_{i=1}^Q \frac{w_i}{w_i + y_i} \quad (7)$$

これらの指標による評価を行うためには, 検索語に対する適合文書と非適合文書が定められていることが前提になる。この目的のために, いくつかのテストコレクションが作成されている。しかし, 教育利用の

Web情報に関しては, テストコレクションが整備されていない。そこで, 本システムを評価するために, 評価実験用の文書集合を作成した。

##### 4.2 文書集合の作成

本システムの評価実験用の文書集合を作成するために, 学習用語として任意の検索用語10語を用意し, インターネット上でヒットしたHTML文書上位1000件を得た。検索には, 全文検索型サーチエンジンであるgooを用い, gooの最大検索件数である1000件を対象とした。しかし, 1000件のうちの58件は取得することができなかったため, 942件が検索実験用の文書である。検索語には, 用語辞書中から無作為に選んだ10語 (遺跡, 宇宙, 栄養, 音符, 茎, 主語, 水泳, 打楽器, 地図, 立方体) を使用し, OR検索を行った。本研究では, これを評価用の文書集合として用いた。

取得した942件をそれぞれの用語について, 適合文書と非適合文書とに分類した。これらの文書は, 現職教員の協力を得て, 学習に利用可能な内容を適合文書, それ以外を非適合文書に人手で分類した。分類結果は, 表3のようになった。

次に, 取得した942件のHTML文書を, パブリックドメインソフトウェアであるNamazu<sup>(11)</sup>でインデックス化した。この際に必要となる, 日本語のわかち書きには茶筌<sup>(5)</sup>を用いた。

表3 942件の分類結果

	適合文書数
遺跡	16
宇宙	12
地図	9
栄養	5
打楽器	4
立方体	4
音符	2
水泳	2
茎	1
主語	1

### 4. 3 文書集合に対する検索結果の評価

インデックス化した942の文書集合に対して、実際に本システムを利用して、適合文書数が比較的多い上位6語の検索語について検索実験を行った。6語は、遺跡、宇宙、栄養、打楽器、地図、立方体である。残りの4語（音符、茎、主語、水泳）は適合文書数が1~2と少ないため用いなかった。

検索語を本システムに入力し、近傍に配置される2用語を抽出した。実際のサーチエンジンにおける検索では、ユーザは2~3語を検索語に用いている<sup>(12)</sup>とされている。そこで、検索語1語のみで検索した場合と、

抽出された2用語を加えて3用語でAND、OR検索を行った場合との計3群で、再現率、精度の変化について調べた。

再現率、精度は次のようにして求めた<sup>(12)</sup>。まず、検索語1語のみ、3用語でAND検索、OR検索のそれぞれについて、取得した942の文書に対して検索を行った。この検索結果を順位の高い順に調べ、適合文書がヒットした時点で、それまでにヒットした適合文書数を全適合文書数で割ったものを再現率、それまでにヒットした適合文書数をそれまでにヒットした文書数で割ったものを精度とした。なお、検索結果の順位付けはNamazuのスコアに基づいて行った。システムの検索

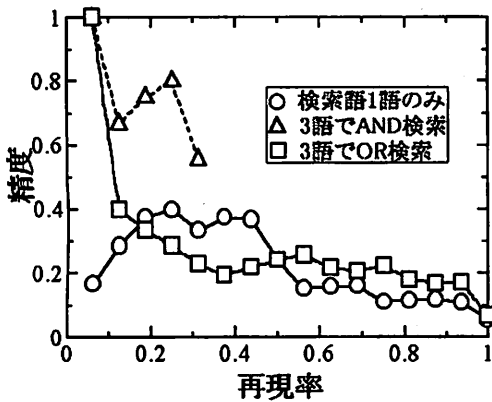


図4 再現率—精度グラフ (遺跡)

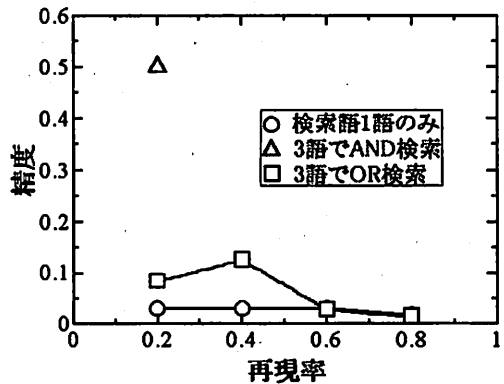


図6 再現率—精度グラフ (栄養)

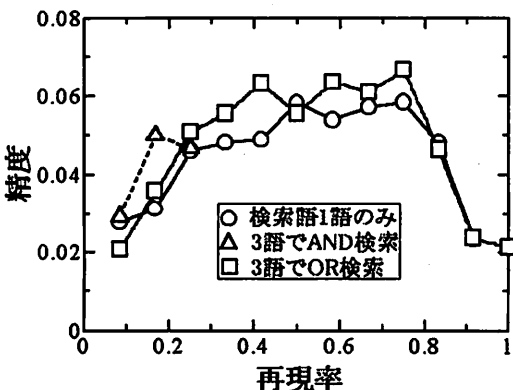


図5 再現率—精度グラフ (宇宙)

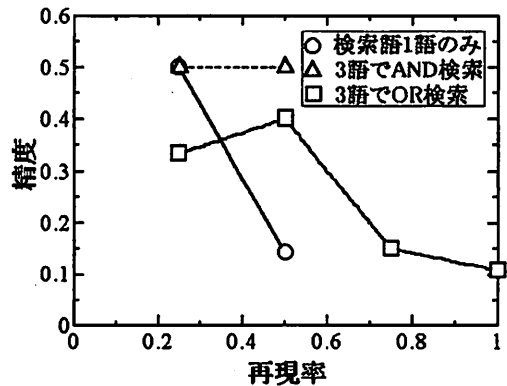


図7 再現率—精度グラフ (打楽器)

性能の評価方法として、再現率-精度グラフによる評価方法がある。各検索語について、再現率-精度グラフで、3群間の比較を行った。図4～図9に6用語の再現率-精度グラフを示す。

上記の定義により、再現率は下位の文書までを多く対象とするにしたがって高くなる。よって再現率は、検索結果の順位として見ることもできる。

図4は「遺跡」の結果を示している。再現率が0.15以下と、0.5以上、すなわち順位が比較的上位と中位以下においては、本システムの支援による3語で検索した方が精度が高い。ただし、再現率0.15～0.5では検索語1語のみで検索した結果が、3語で検索した場合よりも精度が高かった。しかし、その他5用語については、1語の場合よりも3語で検索した場合の方がすべて精度が高かった。「宇宙」では、AND検索と再現率0.2～0.8におけるOR検索によって精度が向上した。「栄養」や「打楽器」では広い範囲でAND検索とOR検索による精度が向上した。「地図」、「立方体」では再現率に関わらず、AND検索とOR検索によって精度が向上した。本結果は、本システムで検索語を追加した方が、1語のみで検索する場合よりも精度の高い検索結果が得られることを示している。

全体の結果をまとめるため、再現率と精度のマクロ平均を求めた。結果を表4に示す。

AND検索では再現率は他条件よりも劣る。一方、精度については、AND検索、OR検索、1語のみで検索の

順によくなった。

AND検索とOR検索について比較したところ、一般的に言われるように、AND検索は精度に優れるが、再現率は劣る。逆に、OR検索は、精度はAND検索よりも劣るが、再現率は優れている。これより、少数の適合ページのみ検索できればよい場合にはAND検索を、多くのページを検索したい場合にはOR検索を用いればよいことが分かる。

ただし、本検討は1つの文書集合に対する結果であることから、複数の文書集合についても同様な評価を行う必要がある。インターネット上の適合文書が少ないことを考慮した評価も必要である。これらは今後の検討課題である。

5. まとめ

本研究では、インターネット上にある学習情報を効

表4 マクロ平均

	再現率 $\bar{R}$	精度 $\bar{P}$
検索語1語のみ	0.527	0.130
3語でOR検索	0.551	0.175
3語でAND検索	0.211	0.531

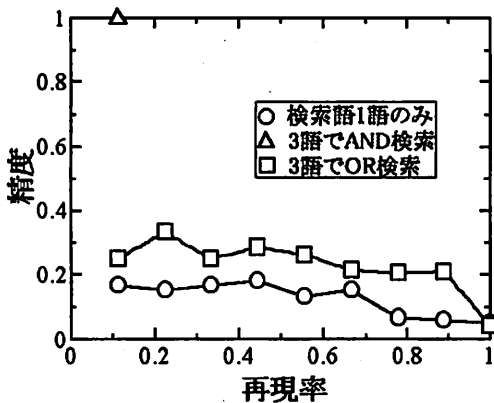


図8 再現率-精度グラフ (地図)

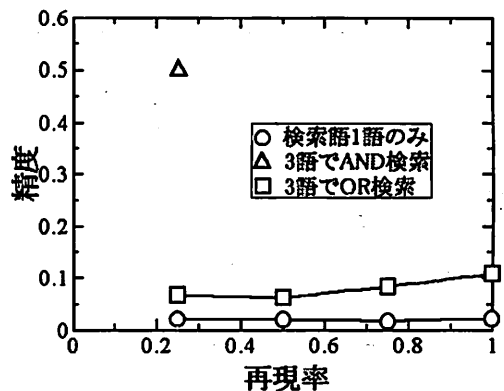


図9 再現率-精度グラフ (立方体)

率的に検索できる検索支援システムを開発した。また、本システムを用いて、実際にインターネット上にある文書を用いて検索性能の評価を行った。以下に得られた結果を示す。

- (1) 学習指導要領を分析して、学習情報に用いられる用語辞書を作成し、これを視覚的に表示する自己組織化マップを作成した。これらの情報をWebブラウザ上に表示し、検索支援を行うシステムを開発した。
- (2) 検索支援システムでは、1つの検索用語が与えられると類似度が高い15の用語を類似度とともに表示し、これを簡便な操作で検索用語リストに加えて、既存のサーチエンジンを利用できるようにした。
- (3) 本システムの検索支援の効果を検討するために、類似用語を加えた場合の検索性能を検討した。その結果、本システムを用いて検索語を追加した方が、検索の精度が向上することを示した。

なお、本論文における検索用語は辞書に含まれる用語に限られることから、学習情報に関連するWebサイトの検索に限られる。これを任意の分野に適用するためには、学習用語としての辞書語の拡充が必要である。また、本システムのユーザインターフェースを含めた、利用者による性能評価は残された課題である。

## 謝 辞

本研究は、文部省科学研究費基盤研究(B)(10400015)、基盤研究(C)(11680210)の補助を受けた。

(2000年3月30日受付)

## 参 考 文 献

- (1) 山本朋弘, 清水康敬, 環境教育に関連するキーワードをインターネットで検索した場合の学習情報に関する一検討, 日本教育工学会誌, Vol.23, Suppl., pp.63-66, 1999
- (2) T.Honkela, S.Kaski, K.Lagus and T.Kohonen, WEB-SOM, URL: <http://websom.hut.fi/websom/>
- (3) 伊藤哲郎, 情報検索, 昭晃堂, 1986
- (4) 文部省, 小学校学習指導要領, 1989
- (5) 日本語形態素解析システム「茶筌」, 奈良先端科学技術大学松本研究室, URL: <http://cl.aist-nara.ac.jp/>

- lab/ntl/chasen/
- (6) S.Deerwester, S.T.Dumais, G.W.Furnus, T.K.Landauer and R.Harshman, "Indexing by Latent Analysis", J. of the Am. Soc. for Info. Sci., Vol.41, No.6, pp.391-407, 1990
- (7) 高山泰博, R.Flournoy, S.Kaufman, S.Peters, 用語の連想関係に基づく情報検索システムInfoMAP, 情報処理学会情報学基礎研究会, 53-1, pp.1-8, 1999
- (8) 中山実, 實松北斗, 清水康敬, 学習指導要領に基づいた学習情報のキーワード検索のための用語辞書に関する検討, 信学論, J83-D-I, No.1, pp.225-233, 2000
- (9) T.Kohonen, J.Hynninen, J.Kangas, J.Laaksonen, SOM\_PACK (1992-1995), 1996, FTP: [ftp://cochlea.hut.fi/pub/som\\_pak/](ftp://cochlea.hut.fi/pub/som_pak/)
- (10) T.Kohonen, "Self-Organizing Maps", Springer Series in Information Science, Volume 30, Springer-Verlag Berlin Heidelberg, 1995 (訳本: 徳高平蔵, 岸田悟, 藤村喜久朗訳, 自己組織化マップ)
- (11) 全文検索システム「Namazu」, Namazu Project, URL: <http://www.namazu.org/>
- (12) 徳永健伸, 情報検索と言語処理, 東京大学出版会, 1999

## 著 者 略 歴

### 森本 容介

2000年, 東京工業大学工学部制御システム工学科卒業。現在, 同大学大学院社会理工学研究科人間行動システム専攻修士課程に在学中。教育工学に関する研究に従事。



### 中山 実

1983年, 東京学芸大学教育学部理科学卒業。1985年同大学院修士課程修了。同年, 東京工業大学教育工学開発センター研究生。1989年東京工業大学大学院博士課程単位取得退学。同年,



