

画像データベースのためのデータマイニング法の拡張

仲川亜希 片山幸治 金山智一 小西 修 菊地時夫
高知大学人文学部社会経済学科
高知大学理学部数理情報科学科
780 高知県高知市曙町 2-5-1
{akinakag} @cc.kochi-u.ac.jp
{katayama, kanayama, konishi, tkukuchi} @is.kochi-u.ac.jp

概要

我々は以前に、ドキュメントの中の共出現キーワードに注目した自己組織化マップによる自動クラスタリング法を提案した。これは、大規模なドキュメント集合における重要キーワード間の関係を表した概念マップを創出できる。今回、我々は、このアルゴリズムを画像に適応し、画像の特徴パターンに注目したクラスタリングを行った。対象画像は、気象衛星による雲の濃淡画像であり、冬型の高気圧での”吹き出し”や台風の目などの雲の動きの特徴によるクラスタリングに良い結果を得た。

キーワード: データマイニング、クラスタリング、自己組織化マップ、相関ルール、共出現情報、気象衛星画像

Enhancement of a Data Mining Algorithm - Application to the Weather Satellite Images -

Aki NAKAGAWA, Koji KATAYAMA, Tomokazu KANAYAMA, Osamu KONISHI, Tokio KIKUCHI

Dept. of Information Science, Faculty of Science, Kochi University

2-5-1 Akebono-cho Kochi 780 Japan

{akinakag, katayama, kanayama, konishi, tkukuchi} @is.kochi-u.ac.jp

Abstract

We have previously shown how a Self-Organization Map (SOM) of automatic clustering using co-occurrence term pairs can be used to perform the data mining of large-scale text databases, the discovery of important data within large document sets by finding optimal data clustering. We report here on an extension of our previous work, substituting the co-occurring relationships among the feature patterns in images for term-pairs in documents. The experimental data is the huge volume of the weather satellite images. We performed the good results for the automatic clustering of these image data.

Key words: Data Mining, Clustering, Self-Organization Map, Association Rules, Co-Occurrence Data, Weather Satellite Image

1 はじめに

近年、データベースをはじめデジタル技術の進歩と相俟って種々の大量の情報が蓄積されてきており、その大量情報の有効活用が求められている。その重要な課題の一つとして、“データマイニング”や“データベースからの知識発見”という研究が盛んになり、その手法の実用化も始まっている。そこで用いられるアルゴリズムは、分類、クラスタリング、決定木、相関ルールなどが主要なものである。これらは、従来、文献情報を対象にした情報検索の研究分野や、機械学習などの人工知能分野で研究開発されたものであるが、データマイニングは、これらを超大規模なデータに適用して、そこから有用な知識を発見しようとする試みである。[1], [2], [3]

我々は、大規模な(約10万件)文献情報を対象に、文献の中に共出現する重要ターム(キーワード)間の情報に注目した共出現ルール抽出アルゴリズム(共出現ルール法と呼ぶ)を開発してきた。これは、共出現語対集合を自己組織化マップによってクラスタリングし、そのクラスから共出現語対ルールを抽出することによって、与えられた文献集合の特徴を表す概念マップを創出できるものである。[7], [8]

本稿では、このアルゴリズムを、大量の画像データに適用し、有用なデータマイニングを行った実験の試みを報告する。上記のアルゴリズムは、文献集合のキーワードが要素であるが、画像データの場合は、画像の中のある単位のブロック(例えば、 20×20 ピクセル)がその要素となり、それは、画像の部分を表す特徴パターンのセルと見なすことができる。

今回の実験のデータは、気象衛星の雲の動きを表す濃淡画像であり、それらの画像の雲の動きと分布による自動クラスタリングが目標である。

2章では、我々の開発した共出現ルール法について、3章では、気象衛星画像の特徴について述べる。つづいて4章では、実験結果を述べ、最後の章でまとめを行っている。

2 共出現ルール法

我々の提案している共出現ルール法について説明する。[9], [10]ここでは対象を文献集合とし、次のような3段階の処理からなる。

1. 共出現語対の抽出
2. 自己組織化特徴マップ
3. ルールの導出

(1) 共出現語対 (co-occurrence term pairs) の抽出

共出現語関係

ドキュメントの中のキーワード(以降KWとする)間の概念の階層関係の情報を得るために同じドキュメント中に共出現するKWとKWの対(共出現語対)とその順序関係を求める。

定義1 ドキュメント集合 D を

$$D = (D_1, D_2, \dots, D_n) \text{ とする。}$$

ドキュメント集合から抽出されるKWを

$$\text{TERM}_k = (t_{1k}, t_{2k}, \dots, t_{nk}),$$

(t_{ik} はドキュメント D_i のKW)とすると、共出現語対は

$$C(\text{TERM}_k, \text{TERM}_h) \\ = \{[t_{ik}, t_{ih}] | t_{ik} \in D_i, t_{ih} \in D_i\} \\ \text{となる。}$$

定義2 $C(\text{TERM}_k, \text{TERM}_h)$ に重みを付けるために、 TERM_k と TERM_h の間の距離(結合度)を次のような関数で与える。

$$f(\text{TERM}_k, \text{TERM}_h) \\ = \frac{\text{freq.of}C(\text{TERM}_k, \text{TERM}_h)}{\sqrt{[\text{freq.of} \text{TERM}_k \times \text{freq.of} \text{TERM}_h]}}$$

$$\text{ここで、} \text{freq.of} \text{TERM}_k = \sum t_{ik} \\ \text{freq.of} \text{TERM}_h = \sum t_{ih}$$

定義 3 $C(\text{TERM}_k, \text{TERM}_h)$ において、 $\text{freq.ofTERM}_k > \text{freq.ofTERM}_h$ ならば、そのとき TERM_k は TERM_h よりも概念の上位関係にあるとする。
 $\text{freq.ofTERM}_k \geq \text{freq.ofTERM}_h$

このように順序を有する共出現語対の二項関係を共出現語関係と呼ぶ。

共出現語対の抽出アルゴリズム

次にこの定義に基づいた共出現語関係の抽出手順を示す。

共出現語対の抽出では、KW から KW のマトリックス上の組み合わせが考えられ、ドキュメント内の共出現回数を求めるために、 $n(n-1) \times m$ 回の計算を必要とする。ここで、 n はドキュメント中から抽出された KW 数、 m はドキュメント数を表す。

step1 ドキュメント集合から KW を抽出し、KW の出現頻度の降順にソートした用語候補リストを準備する。

step2 KW 候補リストを得た元のドキュメント集合を対象に、KW 候補リストの各 KW を検索語とした検索を行なう。

step3 その検索結果の集合から KW 候補リストと同様に KW を切り出し、頻度の降順にソートする。ここで、検索語となった KW(頻度統計の第 1 位の KW) とそれ以外の KW(第 2 位以降からある頻度以上のものまで) との組み合わせが、共出現語対である。このとき、第 2 位以降の KW の頻度は第 1 位の KW との共出現回数を示している。

step4 KW 候補リストの全ての KW について、step2,3 を繰り返す。

step5 得られた共出現語対に対して、定義 2 による結合度を計算する。このようにして得られた共出現語対のデータは表 1 に示すような属性をもったリレーションとして表される。

表 1: 共出現語対データの例

T_1	T_2	f_1	f_2	cof	$cohesion$
learning	training	206	65	37	0.42

ここで T_1, T_2 はキーワード、 f_1, f_2 はそれぞれのキーワードの頻度、 cof は共出現する頻度、 $cohesion$ は結合度を表す。

そこで、表 2 のような関係データベースを考える。

表 2: 文献とキーワードのデータベース例

文献 No.	キーワード
1	A, C, D
2	B, C, E
3	A, B, C, E
4	B, E

これに共出現語対抽出アルゴリズムを用いると図 1 のようになる。ここに、# は文献 No. を、 T_1, T_2 はキーワードを表し、 f はそのキーワードの頻度、 cof は共出現の頻度を表す。

(2) Kohonen の自己組織化マップによるクラスタリング

ドキュメント集合をそのコンテンツに従って分類するために、Kohonen の自己組織化マップを用いた。[6]

Kohonen の自己組織化(特徴)マップ(Self-Organizing (Feature) Map) は、1990 年に T.Kohonen によって提案されたパラダイムであり、ベクトルで表される入力パターン間の位相関係を、学習アルゴリズムにより発見、分類して位相地図を組織化する 2 層のネットワークである。このときベクトルの各成分はパターンの要素に対応している。

この結果得られた地図は、ネットワークに与えられたパターン間の自然な関係構造を表

#	K	#	K	#	f	K
1	A	A	1	1	2	A
1	C	A	3	3	2	A
1	D	B	2	2	3	B
2	B	B	3	3	3	B
2	C	B	4	4	3	B
2	E	C	1	1	3	C
3	A	C	2	2	3	C
3	B	C	3	3	3	C
3	C	D	1	1	1	D
3	E	E	2	2	3	E
4	B	E	3	3	3	E
4	E	E	4	4	3	E

T1	T2	f1	f2	T1	T2	f1	f2	Cf
A	D	2	1	A	C	2	3	2
A	E	2	3	B	C	3	3	2
B	A	3	2	B	E	3	3	3
B	C	3	3	C	E	3	3	2
B	E	3	3					
B	E	3	3					
B	E	3	3					
C	A	3	2					
C	A	3	2					
C	B	3	3					
C	D	3	1					
C	E	3	3					
C	E	3	3					

図1: 共出現語対抽出の過程

#	f	K	T1	T2	f1	f2
1	3	C	C	A	3	2
1	2	A	C	D	3	1
1	1	D	A	D	2	1
2	3	B	B	C	3	3
2	3	C	B	E	3	3
2	3	E	C	E	3	3
3	3	C	C	B	3	3
3	3	B	C	A	3	2
3	2	A	C	E	3	3
3	3	E	B	A	3	2
4	3	B	B	E	3	3
4	3	E	A	E	2	3
			B	E	3	3

自己組織化マップアルゴリズム

step1 入力パターンを与える。

$$E = [e_1, e_2, e_3, \dots, e_n]$$

step2 この入力から競合層の各ユニット i への結合の重みを与える。

$$U_i = [u_{i1}, u_{i2}, \dots, u_{in}]$$

step3 その重みが入力パターンと最もよく一致する競合層のユニット c を定める。すなわち、ベクトル E と U_i の間の距離が最小となるものを探す。

$$\|E - U_c\| = \min_j \|E - U_j\|$$

$$= \sqrt{\sum_j (e_j - u_{ij})^2}$$

step4 このユニット i とその近傍 N_c で重みを調整して一致を増大させる。

$$\Delta u_{ij} = \begin{cases} \alpha(e_j - u_{ij}) & (i \in N_c) \\ 0 & (i \notin N_c) \end{cases}$$

また

$$u_{ij}^{new} = u_{ij}^{old} + \Delta u_{ij}$$

$$\alpha_t = \alpha_0 \left(1 - \frac{t}{T}\right)$$

している。ネットワークは処理ユニットの入力層と競合層の組み合わせであり、教師なし学習により訓練される。入力パターンは競合層で活性化されるユニットにより分類される。パターン間の類似は競合層のグリッド上の近さの関数に写される。

訓練が終了した後、パターン関係やパターングループが競合層で観察される。

Kohonen の自己組織化マップのアルゴリズムは以下の通りである。

ここで、 α は学習率でその値は訓練が進むにつれて0へと減少していく。また、 t は現在の訓練回数であり、 T は行われるべき訓練の全回数である。

step5 学習反復が進むに連れて近傍のサイズと重みの変化の量を次第に減少させる。

ドキュメントをそのドキュメントのコンテンツによって分類された Kohonen マップを生成することができる。こうして得られたマップでは自動的に関連の強い KW が近くにまとめられる。マップ上の KW は、その専門分野の概念体系を表している主要な KW でありこれらの代表的な KW に連なって他の多くの KW があると考えられる。そこで、これらの KW 間の関係からその専門分野の知識構造を把握するために、マップ上の KW をクラスタリングする。

(3) クラスタからのルール抽出

いくつかの出力ノードをひとつのあるクラスタにグループ化し、ルールを定義して概念関係を識別する。クラスタは、それぞれの出力ノードに関するルールの条件文により決定する。すなわち、条件文と同じ属性群を含むルールの出力ノードは同じクラスタにグループ化される。そして、次節で述べる方法により、これらの概念(クラスタ)は階層化される。

ルール導出のアルゴリズムは、以下の通りである。

ルール導出のアルゴリズム

step1 すべての入力から競合層のあるユニット b_k への結合の重みの中で、最大のものを探す。

$$W_{max} = (w_{1k}, w_{2k}, \dots, w_{nk})$$

ただし $W_{max} \neq 0$

step2 $W_{ik} \geq \beta W_{max}$ となる入力 a_i をすべて選ぶ。ここで β は 0 と 1 の間の定数とする。

step3 **step2** で選んだすべての入力を AND でつなぎ、ルールの条件文とする。例えば、**step2** で選ばれた入力を a_{i1}, a_{i2}, a_{i3} とするとルールは

IF (a_{i1} AND a_{i2} AND a_{i3}) THEN(b_k)
となり

(a_{i1} AND a_{i2} AND a_{i3}) \Rightarrow (b_k)
と表す。

step4 **step3** をすべての出力に対して行い、初期ルール集合をつくる。

step5 条件文の中の入力属性が最も少ないルールを選ぶ。

step6 初期ルール集合に、**step5** で選んだ概念を代入する。

step7 代入がそれ以上できなくなるまで、**step5,6** を繰り返す。

step8 最終ルール集合から、概念階層をつくる。

以上の手順によって我々は、ドキュメントの中の共出現キーワードに注目し、大規模なドキュメント集合における重要キーワード間の関係を表した概念マップを創出した。今回、このアルゴリズムを画像に適用し、大量の画像集合における共出現ブロックの関係をを用いる。つまり、画像集合におけるブロック値(パターンセル値)の共出現関係に着目し、一つの画像の中に共出現するブロックの対を取り出す。このブロックの対に頻度情報による順序関係を持たせ、ブロックの階層関係を抽出し、ルールを導出する。

3 気象衛星画像の特徴

日本の静止軌道気象衛星 GMS-5 は、空間的観測密度は衛星直下で 1.25km (可視) ないし 5km (赤外) と極軌道衛星に比べて粗いが、毎時の観測を行っているため、一般気象観測網と極軌道衛星との間を補うものとして重要

である。また、赤外チャンネルは NOAA 極軌道衛星と同様に 2 チャンネルに分かれていたほか、水蒸気チャンネルによる観測も行っており、地表面積の正確な評価や蒸発量の計算などに有効に使える可能性が高い。このような GMS-5 のデータは重要性が高いが、現在のところ、研究者が誰でも、いつでも、過去のデータだけでなく最新のデータについても、使い易い形で直ちに入手することは困難である。東京大学生産技術研究所においては、GMS-5 の運用開始から継続的に高解像度でのデータ受信を行っているが、幾何補正などの処理を施していない生のデータで保存されている。また、現在のところ GMS-5 のデータについてネットワークから利用できるようになってきているのは、間引きを施した、いわゆるクイックルックデータのみである。従って、現在までに受信された GMS データに補正を施した上で、データベース化し、インターネットを通じて研究者に供給する体制が整えば、飛躍的に各種の研究が推進される可能性がある。

また、一般向けの気象情報画像としての利用では、必ずしも画像の全ての情報を保存する必要はなく、一部の情報に欠落があっても、おおまかな気圧配置を推測することができる場合などが考えられる。いずれにしろ、これらの大量の情報を効率よく利用するためには、画像の特徴による自動クラスタリングが必要である。[5]

4 実験

本実験では、高知大学理学部情報科学科気象ページ (<http://weather.is.kochi-u.ac.jp>) で提供されている横 640 ピクセル縦 480 ピクセルのグレースケールの pgm フォーマットの雲の画像情報を使用した。

4.1 ブロック DCT→SOM 方式

- 実験用データ
1997 年の 362 枚 (一日一枚) の画像を対象とした。

- 実験方法

この画像を横 40 ピクセル縦 30 ピクセルの縦横 16 ブロックずつの 256 ブロックに分け、各ブロックごとに以下の式により DCT 変換を行なった。

$$F(u, v) = \frac{2}{\sqrt{MN}} k(u)k(v) \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} f(m, n) \cos\left\{\frac{\pi u}{2M}(2m+1)\right\} \cos\left\{\frac{\pi v}{2N}(2n+1)\right\}$$

$$k(u), k(v) = \begin{cases} \frac{1}{\sqrt{2}} & \text{if } u, v = 0 \\ 1 & \text{otherwise} \end{cases}$$

DCT 変換された画像を SOM にかけてクラスタリングを行なう。SOM に入力する各画像のベクトルは、各ブロックから三つずつである。

$$DC = F(0,0)$$

$$AC1 = F(0,1)$$

$$AC2 = F(1,0)$$

これを 64 ブロック分 768 (3 × 64) 個の要素をを入力し、5 × 5 のマップに出力する。SOM の学習回数は、362 (画像枚数) × 20 回の 7240 回である。

- 実験結果

典型的な冬型の雲の配置となっているものの集まり図 2 に示す。これらの画像は、出力されるマップの座標 (0,4) に出力されたもので、19 枚 (19 日分) クラスタリングされている。

4.2 共出現ルール方式

2 章で述べた共出現ルール抽出アルゴリズムの有用性を検証するための実験を行なった。

- 実験用データ

1996 年から 1997 年までの夏から秋にかけて 900 枚 (6 月から 10 月:1 日約 3 から 5 時間おき) の画像を対象とした。

- 実験方法

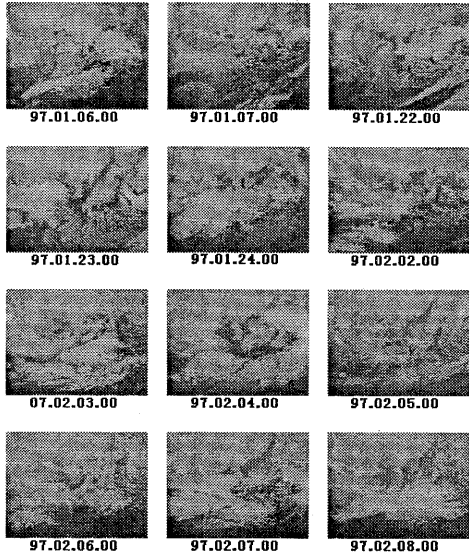


図2 冬型高気圧の張り出しによる日本列島上空の”吹き出し”を含む画像

1. 元画像 (480x640) を 8x8 単位のブロックで表す

- このブロックをパターンセル (patterncell) とする。これは、ドキュメントに出現する k w と同じ、つまり、1 画像のなかの特徴パターンの分布とみなす
- SOM : 入力ベクトル (60 × 80) → 出力 (マップ) 5 × 5
- 元画像を上記のマップクラスタのアドレスで表現 (図3)

24	24	24	23	13	2	24	24
21	11	24	14	3	12	21	11
1	7	11	19	9	0	1	7
2	1	1	16	13	3	2	1
9	9	9	9	9	24	9	9
20	24	24	18	5	15	20	24
6	24	24	12	0	3	6	24
3	20	18	1	0	8	3	20
17	24	24	24	2	1	17	24
16	24	24	24	14	13	16	24
13	24	24	24	9	5	13	24
2	15	18	23	5	2	2	15
13	14	13	2	13	8	13	14
7	9	9	13	5	15	7	9
3	19	24	23	0	8	3	19
9	11	24	12	0	17	9	11
19	24	24	9	14	3	19	24
19	24	24	24	9	17	19	24
2	19	24	24	9	18	2	19
2	13	9	24	24	12	2	13
20	14	9	24	23	0	13	14
14	16	9	24	1	6	20	16
5	8	19	9	12	6	14	8
5	14	9	23	1	17	9	14

図3 マッピング例

2. これらの集合に共出現語対抽出アルゴリズムを適用

- 各パターンセルの出現頻度 → 共出現関係を求める

3. 共出現語対集合を自己組織化マップによりクラスタする。

- SOM : 入力ベクトル (5 × 5: マップクラスタ) → 出力 (3 × 3: マップ)

	0	1	2	3	4
0	1.00	0.64	0.53	0.54	0.70
1	0.64	1.00	0.23	0.35	0.29
2	0.53	0.23	1.00	0.32	0.54
3	0.54	0.35	0.32	1.00	0.38
4	0.70	0.29	0.54	0.38	1.00

図4 入力パターン例 (結合度を用いる)

4. 得られたクラスタからルールを導出する。(図5)

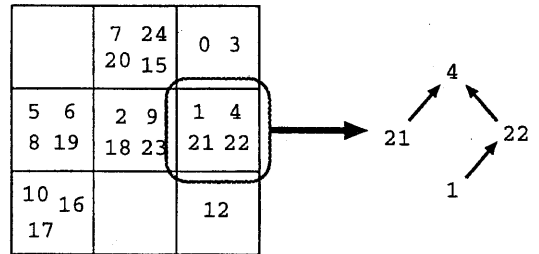


図5 共出現パターンセルからのルール導出の例

5. このパターンによって画像の分類を行う

実験結果

ルール導出によって獲得された台風の中心を含むようなパターンを用いて、実際にクラスタリングを行った結果を図6に示す。

考察

(1) の手法でも図2のように、似かよった雲の分布を持つ画像は同じクラスにクラスタリングされている。しかし、この手法では、画像のパターンセルの順序がそのまま入力ベクトルになるためクラスタは画像のパターンの

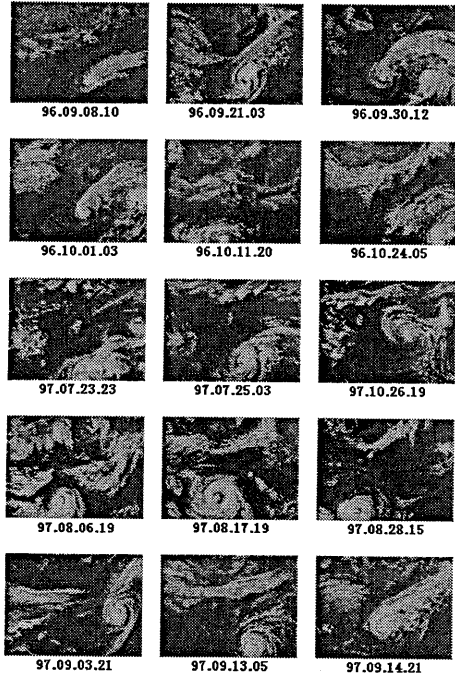


図6 "台風を中心"を含む画像クラスタの例

位置に影響を受ける。このため、"前線による雲"や"台風を中心"など移動するパターンを含む画像については、うまくクラスタリングされない場合がある。(2)の方法では、特徴パターンの分布に基づいてクラスタリングされるため、図6のように"台風を中心"を含む画像が同じクラスにクラスタリングされる。

5 おわりに

画像のパターンセルの共出現関係に注目した共出現ルールアルゴリズムを提案し、大量の気象画像に適用した実験を行った。「日本近辺」の雲の動きと分布の特徴によるクラスタリングに良い結果を示した。

今後の課題として、

1. 今回は、雲の動きの濃淡情報のみであるが、他の気象画像の属性情報(風の状況、海面温度、年月日の季節情報など)と組み合わせた手法の開発

2. 美術画など RGB画像への適用

3. シーン分割など映像情報への適用

などがある。

参考文献

- [1] Alex, A.F., and Simon, H.L., "Mining Very Large Databases with Parallel Processing", Luwer Academic Publishers, 1998.
- [2] Fayyas, U.M., Djorgovski, S.G., and Weir, N., "Automatic teh Analysis and Cataloging of Sky Surveys", In Advances in Knowledge Discovery and Data Mining, pp.471-493, AAAI Press/MIT Press, 1996.
- [3] Fomg, J. (Edt.), "Data Mining, Data Warehousing & Client / Server Databases", Proc. 8 the Int. Database Workshop, Springer, 1997.
- [4] 保木大典, 片山幸治, 仲川亜希, 小西 修: 協調マルチメディア情報収集法, 情処研報, Vol.97, No.56, 97-DBS-113, pp.335-340, 1997.7.
- [5] 菊地時夫, "気象衛星GMS画像の処理とデータベース化体制の確立による長時間の熱収支の解明(序報)", Mem.Fac.Sci.Kochi Univ. (Inform.Sci), 19, March 1998.
- [6] Kohonen. T., Self - Organizing Maps, Springer, 1995.
- [7] 小西 修: アクティブ・メディエーション・システムの設計と実装 - エージェント型データベースの研究 -, 平成9年度科学研究費重点領域研究「高度データベース」東京ワークショップ講演論文集, pp. 260-275, 1997年6月.
- [8] 小西 修, "異種情報源統合のためのアクティブメディエーション・システム-HI-AMS: High Intelligent - Active Mediation system-", Mem.Fac.Sci.Kochi Univ. (Inform.Sci), 17, March 1997.
- [9] 仲川亜希, 小西 修: 自己組織化マップを用いたテキスト情報からの知識獲得, 情処研報, Vol.96, No.68, pp31-36, 96-DBS-109, 1996.7.
- [10] 仲川亜希, 小西 修: 情報探索のための自己組織化アプローチ, 情処研報, Vol.96, No.103, pp39-46, 96-DBS-110, 1996.10.