

クラスタリングに基づく情報の検索と視覚化

堀田, 政二

<https://doi.org/10.15017/1398256>

出版情報 : 九州芸術工科大学, 2001, 博士 (工学), 課程博士
バージョン :
権利関係 :

第4章

混成グラフからのファジークラスタ抽出

4.1 まえがき

本章では第3章で提案したグラフからのファジークラスタ逐次抽出法を複数のグラフが混成された複雑なグラフに適用できるように拡張する [26, 27].

前章では無向グラフや2部無向グラフ, および有向グラフのクラスタリング法が固有値問題に帰着できることを述べた. しかし, これらのグラフが組合さった複雑なグラフ構造データには, 前章のクラスタリング法は適用できない. 例えば, ウェブページに関する情報として, ページ間のリンク関係と, 各ページのキーワードの出現頻度が与えられたとする. この場合, リンクが有向辺で表され, ページとキーワードとの関係が無向辺で表されるので, 全体は有向グラフと2部無向グラフとを合成したものになる. 本章では, このような複雑なグラフ構造データからファジークラスタを抽出する方法として, べき乗法を一般化した反復法による抽出法を導く. はじめに, 無向グラフと2部無向グラフ, および有向グラフから反復法に基づいてファジークラスタを逐次に抽出する方法を示す. 次に, 有向グラフと2部無向グラフが組合わさった複雑なグラフからファジークラスタを抽出する方法を示す. また, これらの反復法の収束性を証明する. 最後に, 簡単な混成グラフを使った実験例を示し提案手法の有効性を検証する.

4.2 反復法による無向グラフからのファジークラスタ抽出

n 個の点からなる無向グラフを考え, 第 i 点の重みを v_i , 第 i 点と第 j 点の間の辺の重みを w_{ij} とする (無向であるから $w_{ij} = w_{ji}$ である. $w_{ii} = 0$ とする). v_i も w_{ij} も 0 以上 1 以下に規格化しておく. この点集合からクラスタを順番に取り出す. 第 i 点

が第1クラスタに所属する割合を x_{1i} とする. 前章に述べたように, これは

$$\begin{aligned} \max_{\mathbf{x}_1} \quad & \sum_{i=1}^n \sum_{j=1}^n v_i x_{1i} w_{ij} v_j x_{1j} \\ \text{subj.to} \quad & \sum_{i=1}^n v_i x_{1i}^2 = 1 \end{aligned} \quad (4.1)$$

によって求められる. 前章では式 (4.1) の解を固有値分解で求めたが, 本章では後で出てくる混成グラフに固有値分解が適用できないので, 反復法を用いることにし, 解法を本章全体で統一して理解しやすくするために式 (4.1) の反復解法を導いておく. Lagrange 乗数法により式 (4.1) の Lagrange 関数は

$$L = \sum_{i=1}^n \sum_{j=1}^n v_i x_{1i} w_{ij} v_j x_{1j} - \lambda \left(\sum_{i=1}^n v_i x_{1i}^2 - 1 \right) \quad (4.2)$$

となる. λ は Lagrange 乗数である. 式 (4.1) の解は

$$\frac{1}{2} \frac{\partial L}{\partial x_{1i}} = \sum_{j=1}^n v_i w_{ij} v_j x_{1j} - \lambda v_i x_{1i} = 0 \quad (4.3)$$

$$\frac{\partial L}{\partial \lambda} = 1 - \sum_{i=1}^n v_i x_{1i}^2 = 0 \quad (4.4)$$

を満たす. 式 (4.3) から

$$x_{1i} = \sum_{j=1}^n w_{ij} v_j x_{1j} / \lambda \quad (4.5)$$

となるから, これを式 (4.4) に代入して λ を求めると

$$\lambda = \sqrt{\sum_{i=1}^n v_i \left(\sum_{j=1}^n w_{ij} v_j x_{1j} \right)^2} \quad (4.6)$$

となり, これを式 (4.5) に代入すると

$$x_{1i} = \frac{\sum_{j=1}^n w_{ij} v_j x_{1j}}{\sqrt{\sum_{i=1}^n v_i \left(\sum_{j=1}^n w_{ij} v_j x_{1j} \right)^2}} \quad (i = 1, \dots, n) \quad (4.7)$$

となる. そこで $\mathbf{x}_1 = [x_{11}, \dots, x_{1n}]^T$ を任意に初期設定し, それを式 (4.7) の右辺に代入して \mathbf{x}_1 を更新していく. これは固有ベクトルを求める反復法としてよく知られているべき乗法を一般化したものとなっている. すなわち, 後で出てくる反復解法もすべてべき乗法のある種の一般化であるといえる. x_{1i} が最大の i を i_1 とすると $p_{1i} = x_{1i} / x_{1i_1}$ がデータ i の第1クラスタへのメンバシップとなる.

次に第2クラスタ抽出では第1クラスタを取り除いて同じことを行う. すなわち各点の重みを $(1 - p_{1i})v_i$ とする. したがって, 各点が第2クラスタに所属する割合 x_{2i}

は式 (4.1) の v_i を $(1 - p_{1i})v_i$ に変えて同じ反復法を行えば求まる. 第3 クラスタ以降も同様であり, 一般に第 k クラスタでは式 (4.1) の v_i を $\prod_{i=1}^{k-1} (1 - p_{ki})v_i$ として同じことを行えばよい.

4.3 反復法による2部無向グラフからのファジークラスタ抽出

次に2部無向グラフ, すなわち点が2つの部分集合 $i \in S_1: i = 1, \dots, m$ と $j \in S_2: j = 1, \dots, n$ とに別れており, 部分集合間には重み w_{ij} の無向辺 (すなわち $w_{ij} = w_{ji}$) があり, 部分集合内には辺がない場合を考える. 例えば部分集合 S_1 をクラスタリングするとする. 部分集合 S_1 の点 i が第1 クラスタに所属する割合を x_{1i} , S_2 の点 j の所属度を y_{1j} とすると第1 クラスタは

$$\begin{aligned} \max_{\mathbf{x}_1, \mathbf{y}_1} \quad & \sum_{i=1}^m \sum_{j=1}^n v_i x_{1i} w_{ij} v_j y_{1j} \\ \text{subj.to} \quad & \sum_{i=1}^m v_i x_{1i}^2 = 1, \quad \sum_{j=1}^n v_j y_{1j}^2 = 1 \end{aligned} \quad (4.8)$$

で求められる. これも前章 3.4 節で示したように $\tilde{x}_{1i} = \sqrt{v_i} x_{1i}$, $\tilde{y}_{1j} = \sqrt{v_j} y_{1j}$ と変数変換すれば固有値問題に帰着できるが, ここでは式 (4.8) を直接反復法で解く. 式 (4.8) の Lagrange 関数は

$$L = \sum_{i=1}^m \sum_{j=1}^n v_i x_{1i} w_{ij} v_j y_{1j} - \lambda \left(\sum_{i=1}^m v_i x_{1i}^2 - 1 \right) - \mu \left(\sum_{j=1}^n v_j y_{1j}^2 - 1 \right) \quad (4.9)$$

となる. 式 (4.8) の解は

$$\frac{\partial L}{\partial x_{1i}} = \sum_{j=1}^n v_i w_{ij} v_j y_{1j} - 2\lambda v_i x_{1i} = 0 \quad (4.10)$$

$$\frac{\partial L}{\partial y_{1j}} = \sum_{i=1}^m v_i w_{ij} v_i x_{1i} - 2\mu v_j y_{1j} = 0 \quad (4.11)$$

$$\frac{\partial L}{\partial \lambda} = 1 - \sum_{i=1}^m v_i x_{1i}^2 = 0 \quad (4.12)$$

$$\frac{\partial L}{\partial \mu} = 1 - \sum_{j=1}^n v_j y_{1j}^2 = 0 \quad (4.13)$$

を満たす. 前節で式 (4.7) を導いたのと同様にして, 式 (4.10) から

$$x_{1i} = \sum_{j=1}^n w_{ij} v_j y_{1j} / 2\lambda \quad (4.14)$$

となるから、これを式 (4.12) に代入すると λ が求まるのでそれを再び式 (4.14) に代入すると

$$x_{1i} = \frac{\sum_{j=1}^n w_{ij} v_j y_{1j}}{\sqrt{\sum_{i=1}^m v_i \left(\sum_{j=1}^n w_{ij} v_j y_{1j} \right)^2}} \quad (i = 1, \dots, m) \quad (4.15)$$

が得られる。同様にして式 (4.11) と式 (4.13) とから μ を消去すると

$$y_{1j} = \frac{\sum_{i=1}^m v_i w_{ij} x_{1i}}{\sqrt{\sum_{j=1}^n v_j \left(\sum_{i=1}^m v_i w_{ij} x_{1i} \right)^2}} \quad (j = 1, \dots, n) \quad (4.16)$$

が導かれる。そこで $\mathbf{y}_1 = [y_{11}, \dots, y_{1n}]^T$ を任意に初期設定し、それを式 (4.15) の右辺に代入して \mathbf{x}_1 を求め、それを式 (4.16) の右辺に代入して \mathbf{y}_1 を更新していく。この反復法も、後で出てくる混成グラフの反復法の特殊な場合である。 S_1 内で x_{1i} が最大の i を i_1 とすると $p_{1i} = x_{1i}/x_{1i_1}$ がデータ $i \in S_1$ の第1クラスタへのメンバシップ値となり、 S_2 内で y_{1j} が最大の j を j_1 とすると $q_{1j} = y_{1j}/y_{1j_1}$ がデータ $j \in S_2$ の第1クラスタへのメンバシップ値となる。

次に第2クラスタ抽出では、今は部分集合 S_1 をクラスタリングしているのであるから、点 $i \in S_1$ の重みが $(1 - p_{1i})v_i$ と削減され、点 $j \in S_2$ の重みは v_j のままとする。したがって、各点が第2クラスタに所属する割合 x_{2i}, y_{2j} は、式 (4.8) の v_i を $(1 - p_{1i})v_i$ に変えて同じことをすれば求まる。第3クラスタ以降も同様であり、一般に第 k クラスタでは式 (4.8) の v_i を $\prod_{l=1}^{k-1} (1 - p_{li})v_i$ として同じことを行えばよい。

4.4 反復法による有向グラフからのファジークラスタ抽出

有向グラフでは $w_{ij} \neq w_{ji}$ である。この場合、各点 i を2つの点 \hat{i} と \hat{i} に分けて点の数を2倍にし、 \hat{i} から \hat{j} に無向辺を引き、その重みを w_{ij} とすれば有向グラフは2部無向グラフに帰着する(前章図 3.4 参照)。したがって、有向グラフは前節の2部無向グラフと同じ方法でクラスタリングでき、各点について辺の始点としてのメンバシップ p_i と終点としてのメンバシップ q_j が得られることになる。ただしこの2部無向グラフ表現はあくまで形式的であり、 \hat{i} と \hat{i} は元は同じ点 i であるので、これらのメンバシップ p_i と q_j は同時に削減されていく。

4.5 混成グラフからのファジークラスタ抽出

以上、無向グラフや2部無向グラフ、および有向グラフの場合の反復解法に基づくクラスタリング法を示した。これらのグラフ構造データに対しては第3章のクラスタリング法を適用することができるが、与えられる情報が多い場合はこれらを組み合わせた複雑なグラフで表され、固有値問題に帰着できない。例えば、ウェブページに関する情報として、ページ間のリンク関係と、各ページのキーワードの出現頻度が与えられたとする。この場合、リンクが有向辺で表され、ページとキーワードとの関係が無向辺で表されるので、全体は有向グラフと2部無向グラフとを合成したものになる。そこでここでは新たにそのような場合についてのファジークラスタリング法を導く。

4.5.1 混成グラフからのファジークラスタ抽出

例として、向きを無視したリンクとキーワードに基づいてウェブページをクラスタリングする場合について考えてみる。ページを $i = 1, \dots, m$ とし、キーワードを $j = 1, \dots, n$ とする。ページ i からページ i' へのリンクの重み (0 か 1) を $w_{ii'}$ (向きは無視するので $w_{ii'} = w_{i'i}$) とし、キーワード j のページ i での頻度を w_{ij} とする。ページ i が第1クラスタに所属する割合を x_{1i} 、キーワード j の所属度を y_{1j} とすると、式 (4.1) と式 (4.8) を混合した

$$\begin{aligned} \max_{\mathbf{x}_1, \mathbf{y}_1} \quad & \alpha \sum_{i=1}^m \sum_{i'=1}^m v_i x_{1i} w_{ii'} v_{i'} x_{1i'} \\ & + (1 - \alpha) \sum_{i=1}^m \sum_{j=1}^n v_i x_{1i} w_{ij} v_j y_{1j} \\ \text{subj.to} \quad & \sum_{i=1}^m v_i x_{1i}^2 = 1, \quad \sum_{j=1}^n v_j y_{1j}^2 = 1 \end{aligned} \quad (4.17)$$

で第1クラスタが求められる。 $0 \leq \alpha \leq 1$ は重みであり、これが大きいほどリンクを重視したクラスタリングになり、小さいとキーワード、すなわちページの内容を重視したクラスタリングとなる。Lagrange 乗数法より式 (4.17) から

$$x_{1i} = \frac{z_{1i}}{\sqrt{\sum_{i=1}^m v_i z_{1i}^2}} \quad (i = 1, \dots, m) \quad (4.18)$$

$$y_{1j} = \frac{\sum_{i=1}^m v_i w_{ij} x_{1i}}{\sqrt{\sum_{j=1}^n v_j \left(\sum_{i=1}^m v_i w_{ij} x_{1i} \right)^2}} \quad (j = 1, \dots, n) \quad (4.19)$$

が導かれる。ここで

$$z_{1i} = 2\alpha \sum_{i'=1}^m w_{i'i} v_{i'} x_{1i'} + (1 - \alpha) \sum_{j=1}^n w_{ij} v_j y_{1j} \quad (4.20)$$

である。式 (4.18) と式 (4.19) を反復計算して解を求める。まず \mathbf{x}_1 と \mathbf{y}_1 の初期値を任意に設定し、それを式 (4.18) の右辺に代入して \mathbf{x}_1 を求め、それを式 (4.19) に代入して \mathbf{y}_1 を求め、これらの新しい \mathbf{x}_1 と \mathbf{y}_1 を式 (4.18) に代入するというのを \mathbf{x}_1 と \mathbf{y}_1 が収束するまで繰り返す。以上の反復で得られる \mathbf{x}_1 と \mathbf{y}_1 とから $p_{1i} = x_{1i}/\max\{x_{1i}\}$, $q_{1j} = y_{1j}/\max\{y_{1j}\}$ によりページのメンバシップ p_{1i} とキーワードのメンバシップ q_{1j} が得られる。

次に第2クラスタ抽出では、今はページ集合をクラスタリングしているのであるから、ページ i の重みが $(1 - p_{1i})v_i$ と削減され、キーワード j の重みは v_j のままとする。したがって、各ページとキーワードが第2クラスタに所属する割合 x_{2i} と y_{2j} は、式 (4.17) の v_i を $(1 - p_{1i})v_i$ に変えて同じことをすれば求まる。第3クラスタ以降も同様であり、一般に第 k クラスタでは式 (4.17) の v_i を $\prod_{l=1}^{k-1} (1 - p_{li})v_i$ として同じことを行えばよい。この場合も抽出されるクラスタの凝集度は順番が進むに従って単調に減少するので、この凝集度の変化に基づいてクラスタ数を決める。

4.5.2 収束証明

本クラスタリング法の解は大域的に収束する。以下に証明を与える。

まずベクトルのノルムを

$$\|\mathbf{x}\|_v = \sqrt{\sum_{i=1}^m x_i^2/v_i} \quad (4.21)$$

と定義する。 \mathbf{x}_1 と \mathbf{y}_1 の第 ξ 反復解を $\mathbf{x}_1^{(\xi)}$, $\mathbf{y}_1^{(\xi)}$ と記す。また、式 (4.17) の目的関数を $h(\mathbf{x}_1, \mathbf{y}_1)$ とし、その \mathbf{x}_1 に関する勾配ベクトルを

$$\nabla_{x_1}(\mathbf{x}_1, \mathbf{y}_1) = [\partial h/\partial x_{11}, \dots, \partial h/\partial x_{1m}]^T \quad (4.22)$$

とし

$$\nabla^{(\xi)} = \nabla_{x_1}(\mathbf{x}_1^{(\xi)}, \mathbf{y}_1^{(\xi)}) = [\nabla_1^{(\xi)}, \dots, \nabla_m^{(\xi)}]^T \quad (4.23)$$

と書く。すると式 (4.18) による \mathbf{x}_1 の更新は

$$x_{1i}^{(\xi+1)} = \nabla_i^{(\xi)}/v_i \|\nabla^{(\xi)}\| \quad (4.24)$$

と書ける。これを使って $(\nabla^{(\xi)})^T(\mathbf{x}_1^{(\xi+1)} - \mathbf{x}_1^{(\xi)})$ を評価すると

$$(\nabla^{(\xi)})^T(\mathbf{x}_1^{(\xi+1)} - \mathbf{x}_1^{(\xi)}) = (\nabla^{(\xi)} \circ \mathbf{v})^T(\nabla^{(\xi)} - \|\nabla^{(\xi)}\|_v \mathbf{v} \otimes \mathbf{x}) / \|\nabla^{(\xi)}\|_v \quad (4.25)$$

となる. ここで $\mathbf{v} \otimes \mathbf{x} = [v_1 x_1, \dots, v_m x_m]^T$, $\nabla^{(\xi)} \odot \mathbf{v} = [\nabla_1^{(\xi)}/v_1, \dots, \nabla_m^{(\xi)}/v_m]^T$ である. $\|\mathbf{v} \otimes \mathbf{x}\|_c = \sqrt{\sum_{i=1}^m v_i x_i^2} = 1$ であるから

$$\|\|\nabla^{(\xi)}\|_v \|\mathbf{v} \otimes \mathbf{x}\|_c \leq \|\nabla^{(\xi)}\|_c \|\mathbf{v} \otimes \mathbf{x}\|_c = \|\nabla^{(\xi)}\|_c \quad (4.26)$$

となる. したがって,

$$(\nabla^{(\xi)})^T (\mathbf{x}_1^{(\xi+1)} - \mathbf{x}_1^{(\xi)}) \geq 0 \quad (4.27)$$

が成り立つ. このことから式 (4.18) による \mathbf{x}_1 の更新において

$$h(\mathbf{x}_1^{(\xi+1)}, \mathbf{y}_1^{(\xi)}) \geq h(\mathbf{x}_1^{(\xi)}, \mathbf{y}_1^{(\xi)}) \quad (4.28)$$

が成り立つ. 同様に式 (4.19) による \mathbf{y}_1 の更新においても

$$h(\mathbf{x}_1^{(\xi+1)}, \mathbf{y}_1^{(\xi+1)}) \geq h(\mathbf{x}_1^{(\xi+1)}, \mathbf{y}_1^{(\xi)}) \quad (4.29)$$

が成り立つことが示せる. 以上より $\mathbf{x}_1, \mathbf{y}_1$ を 1 回更新すると

$$h(\mathbf{x}_1^{(\xi+1)}, \mathbf{y}_1^{(\xi+1)}) \geq h(\mathbf{x}_1^{(\xi+1)}, \mathbf{y}_1^{(\xi)}) \geq h(\mathbf{x}_1^{(\xi)}, \mathbf{y}_1^{(\xi)}) \quad (4.30)$$

が成り立つ. すなわち, 更新するたびに $h(\mathbf{x}_1, \mathbf{y}_1)$ は単調に増加する. \mathbf{x}_1 も \mathbf{y}_1 も有界なのでこの反復更新は収束する. なお, 式 (4.7) や式 (4.15), 式 (4.16) は式 (4.18), 式 (4.19) の特殊な場合 ($\alpha = 1$ と $\alpha = 0$) であるから, それらの反復も収束する.

4.5.3 実験例

混成グラフからのクラスタ抽出について簡単なデータで実験を行った. データはページ間のリンク関係を無向グラフで表したものと, ページとキーワードの関係を2部無向グラフで表したものを合成させた3部無向グラフを用いた. これらのデータは以下のように人工的に作成した. まず, ページ間のリンク関係を図 4.1 のように作成した. 数字がページ名であり, リンクを張っているページ同士を線で結んでいる. リンクの向きは無視し, 各辺の重みはすべて1としている. ページが3つの独立なコミュニティを生成するようにした. 一方, ページとキーワードの関係は表 4.1 となるように作成した. キーワードは全部で6個あり, 表の要素はそれぞれのページ中のキーワードの出現頻度である. ページ1から3はキーワード1, 2を含み, ページ4から6はキーワード4を, ページ7から9はキーワード5, 6を含むようにしてクラスタを作成した. ただし, ページ3はキーワード3を, ページ6はキーワード5を含むようにして, ページ間の潜在的な関連性が生じるようにした.

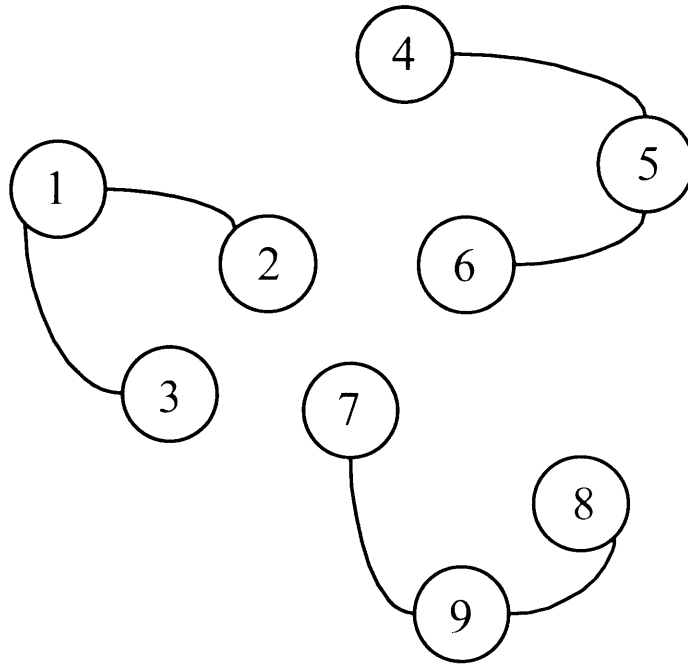


図 4.1: ページのリンク関係

表 4.1: ページとキーワードの関係

	page 1	page 2	page 3	page 4	page 5	page 6	page 7	page 8	page 9
keyword 1	2	3	1	0	0	0	0	0	0
keyword 2	2	2	1	0	0	0	0	0	0
keyword 3	0	0	1	0	1	0	0	0	0
keyword 4	0	0	0	3	1	1	0	0	0
keyword 5	0	0	0	0	0	1	2	3	2
keyword 6	0	0	0	0	0	0	3	3	1

まず最初に、式 (4.17) の α を 0.9 としてページをクラスタリングした。図 4.2 に凝集度の変化を示す。この図からクラスタ数を 3 個とした。図 4.3 に本方法で得られた各ページのクラスタへのメンバシップを示す。横軸はページ名である。クラスタを線の種類で区別している。図 4.3 からわかるように、各ページはほとんど 1 つのクラスタへ所属している。これは α の値が大きい、すなわちページのリンク関係を重視したクラスタリングを行ったため、各ページが独立した 3 つのコミュニティとして抽出されたことを示している。各クラスタの代表ページは、リンクの結合和が最大であるページ 1, 5, 9 となった。

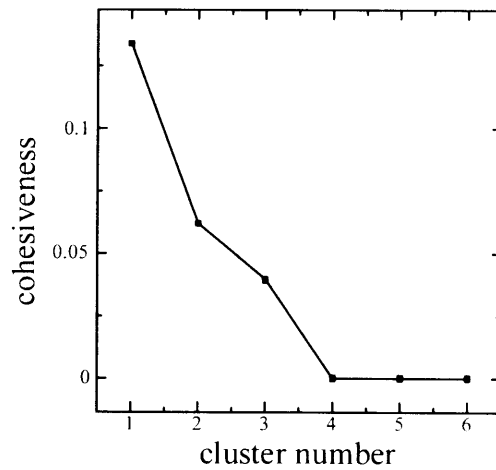


図 4.2: 凝集度の変化 ($\alpha = 0.9$)

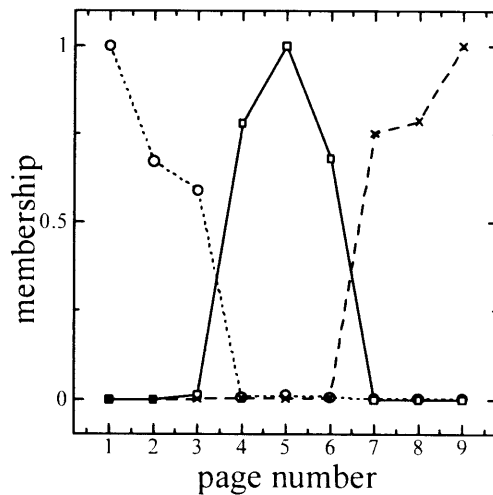


図 4.3: ページのメンバシップ ($\alpha = 0.9$)

次に、式 (4.17) の α を 0.5 としてページをクラスタリングした。図 4.4 に凝集度の変化を示す。3 番目のクラスタから 4 番目にかけての減少が最も大きい。これは主要なクラスタが 3 個であることを示している。図 4.5 に各ページのクラスタへのメンバシップを示す。図 4.5 を見ると、ページ 3 とページ 6 の中央のクラスタへのメンバシップ値が僅かであるが増加している。これはリンク関係だけでは抽出できない潜在的なページ間の関係が、キーワード情報を付加することにより抽出されていることを示している。各クラスタの代表データもキーワードの出現頻度に影響され、ページ 1, 4, 8 となった。

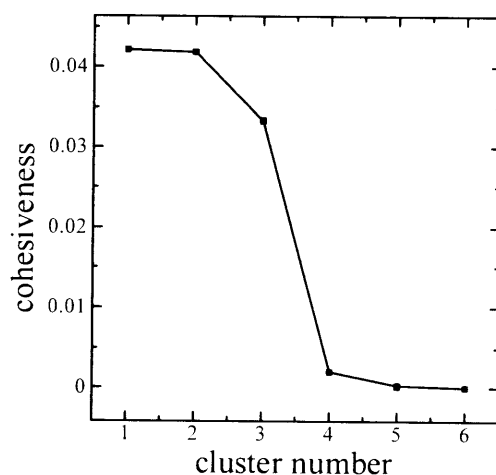


図 4.4: 凝集度の変化 ($\alpha = 0.5$)

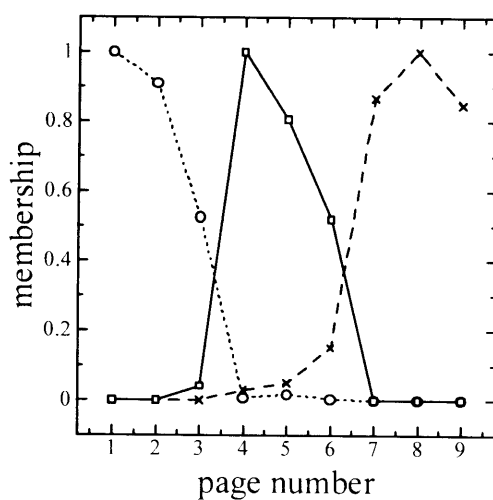


図 4.5: ページのメンバシップ ($\alpha = 0.5$)

最後に、式 (4.17) の α を 0.1 としてページをクラスタリングした。図 4.6 に凝集度の変化を示す。4 番目以降のクラスタから変化が緩やかになっているのでクラスタ数を 3 個とした。図 4.7 に各ページのクラスタへのメンバシップを示す。この図から特にページ 5 とページ 6 のクラスタへのメンバシップ値が極端に減少していることがわかる。これは α の値が小さいために、ページ間のリンク関係がほとんど無視されて、キーワードの出現頻度を重視したクラスタリングが行われたためであり、ページ中のキーワードの出現頻度が少ないページ 5 や 6 のメンバシップ値が小さくなったと考えられる。また、キーワードの出現頻度が多いものが各クラスタの代表データとなり、ページ 2, 4, 8 となった。

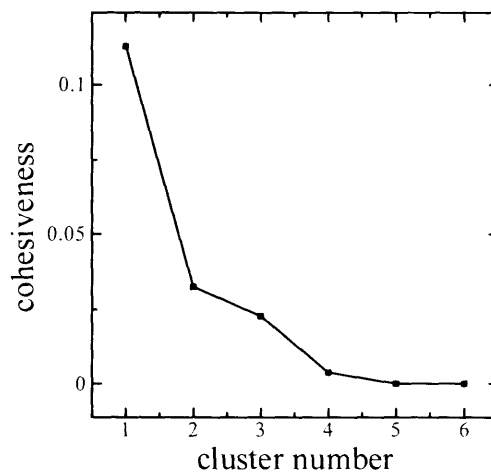


図 4.6: 凝集度の変化 ($\alpha = 0.1$)

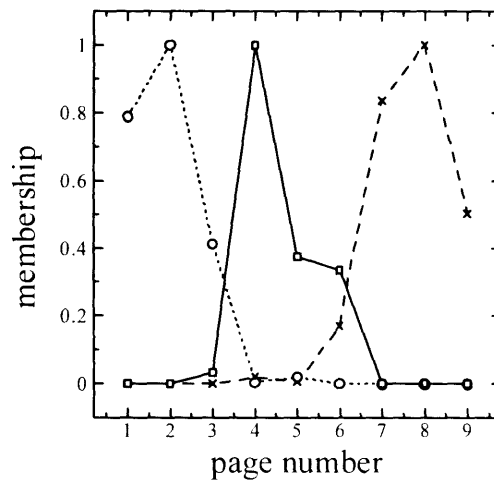


図 4.7: ページのメンバシップ ($\alpha = 0.1$)

4.6 むすび

本章では、第 3 章で提案したグラフからのファジークラスタ逐次抽出法を、複数のグラフが混成された複雑なグラフに適用できるように拡張した。複雑なグラフ構造データからのファジークラスタ抽出問題は固有値問題に帰着できないため、べき乗法を一般化した反復法による抽出法を導いた。はじめに、無向グラフと 2 部無向グラフ、および有向グラフから反復法に基づいてファジークラスタを逐次に抽出する方法を示した。次に、有向グラフと 2 部無向グラフが組合わさったグラフからファジークラスタを抽出する方法を示した。また、これらの反復法が収束することを証明した。最後に、簡単な実験例を示し提案手法の有効性を確認した。本手法のデータ検索への応用は第 6 章で詳しく述べる。