

# 欠測値を含むデータからの局所的な主成分の抽出法\*

(Extraction of Local Principal Components  
from Data with Missing Values)

本多 克宏<sup>†</sup>・神田 章裕<sup>†</sup>・市橋 秀友<sup>†</sup>・山川あす香<sup>‡</sup>

<sup>†</sup>大阪府立大学 大学院 工学研究科

Graduate School of Engineering,

Osaka Prefecture University;

1-1 Gakuen-cho, Sakai city, Osaka 599-8531, JAPAN

<sup>‡</sup>福島学院短期大学 情報ビジネス科

Information and Business Communication,

Fukushima College;

1-1 Chigoike Miyashiro, Fukushima 960-0181, JAPAN

## 概要

大規模なデータベースからの知識発見の手法として、ファジィクラスタリングと多変量解析法を組み合わせることにより局所的な特徴量を抽出する研究が近年盛んに行われているが、実データを扱う際にはしばしば欠測値の処理方法が問題となる。本論文では、観測値のみを用いる主成分分析法に Fuzzy  $c$ -Means 法を融合することにより、欠測値を含むデータから局所的な主成分を抽出する手法を提案する。ファジィ散布行列の固有値問題に帰着される Fuzzy  $c$ -Varieties 法が分散共分散行列を用いた主成分分析とクラスタリングの同時適用法ととらえられるのに対して、相関係数行列の固有値問題に帰着される提案法は相関係数行列を用いた主成分分析とクラスタリングの同時適用法であるといえる。

(In many real world applications data sets with missing values are quite common. In this paper, we propose a new approach which extracts local principal components for the feature extraction from a large scale database. The new method is based on a simultaneous approach to principal component analysis and fuzzy clustering with an incomplete data set including missing values. The simultaneous approach extracts local principal components by using the eigenvectors of the correlation coefficient matrix, while Fuzzy  $c$ -Varieties(FCV) proposed by Bezdek *et al.* uses the eigenvectors of the fuzzy scatter matrix.)

*keywords:* fuzzy clustering, principal component analysis, missing value, eigen value problem.

## 1 はじめに

近年、大規模なデータベースからの知識発見の手法として、ファジィクラスタリングと多変量解析法を組み合わせることにより局所的な特徴量を抽出する研究が盛んに行われている。Bezdek らによる Fuzzy  $c$ -Varieties (FCV) 法および Fuzzy  $c$ -Elliptotypes (FCE) 法 [1, 2] は、局所的な主成分ベクトルによりいくつかの線形多様体状あるいは楕円体状のクラスタに分類する手法であり、クラスタ分析と主成分分析の同時適用法であるといえる。また、その改良として、クラスタごとの分散共分散行列の固有値をトレードオフパラメータに採用した Adaptive Fuzzy  $c$ -Elliptotypes (AFC) 法 [3] や、次元の異なる多様体の発見法 [4, 5] など提案されている。その他、Bezdek らの Fuzzy  $c$ -Means

---

\*システム制御情報学会論文誌, 15, 12, 663-672 (2002)

(FCM) 法 [1] の目的関数とその他の多変量解析手法の目的関数を組み合わせることにより、様々な局所的特徴量を抽出する試みもなされている [6, 7]。

しかし、実データを分析するには、データの中に観測されなかった部分、すなわち欠測値を含み、直接、これらの手法を適用できないことがある。多変量解析法の適用において、このような不完全データを処理する方法がいくつか考えられる [8]。最も簡単な手法として、欠測値を含むデータを削除した後に分析するということが考えられるが、削除することで多くの情報を切り捨ててしまうという問題点がある。また、欠測値そのものを何らかの方法によって推定し、データの欠落部分を補完してから分析する手法なども考えられるが、補完の際に用いられるモデルの妥当性などの問題がある。そこで、モデルの仮定を設けず、欠測部分の補完という手続きを避けながら、なおかつデータ行列のすべての観測値を直接利用して、多変量データの分析を行うことが望まれる。

欠測値を含む多変量データに対する主成分分析法として、Ruhe [9] や Wiberg [10]、柴山 [11] は、データ行列の最小 2 乗近似において観測値に対応する要素については誤差が小さくなるようにするものの、欠測値に対応する要素については変化するに任せることにより、多変量データの主たる変動成分を抽出する方法を提案している。また、高根 [12] は柴山ら [13] の欠測値がある場合の線形等化法を多次元に拡張することにより、標準化したデータの主成分を抽出する手法を提案している。

一方、ファジィクラスタリングの適用においても、同様に欠測値の処理法が問題となる。宮本ら [14] や Timm ら [15] は、重み付き平均を欠測値の推定値として用いたり、欠測値を無視したりすることにより、FCM 法において欠測値を処理する手法を提案している。また、本多ら [16] は柴山 [11] の最小 2 乗近似に基づく主成分分析法に宮本らの手法を融合することにより FCV 法において欠測値を処理する手法を提案し、協調フィルタリングへの応用を行っている [17] が、データ行列と同じ要素数の近似行列を繰り返しごとに求める必要があり、データ数が増えるにつれてメモリの所要量や計算量が増加するという欠点がある。

そこで本論文では、高根の欠測値を含むデータの主成分分析法にファジィクラスタリングで用いるメンバシップを導入することにより、欠測値を含むデータに対応できるクラスタリング法を提案する。提案法はラグランジュ乗数法を FCM 法と高根による主成分分析法の同時分析に応用したもので、K-L 情報量正則化による FCM クラスタリング [18] により複数のクラスタに分類して分析する。FCV 法が分散共分散行列を用いた主成分分析とファジィクラスタリングの同時適用法とみなされるのに対して、提案法は標準化されたデータによる相関係数行列を用いた主成分分析にクラスタリングを融合した手法である。

数値例では、人工データやスーパーマーケットで収集された POS データを用いて、提案法の特徴を検証する。

## 2 欠測値を含むデータの主成分分析とクラスタリングの同時適用法

$n$  次元の  $N$  個の標本データからなる  $(N \times n)$  データ行列  $X = (x_{ij})$  が与えられたときに、 $N$  個の標本データを  $C$  個のクラスタに分割しながら、クラスタごとに局所的な主成分を抽出する問題を考える。

Bezdek らの Fuzzy  $c$ -Elliptotypes (FCE) 法 [1, 2] は、線形のクラスタを抽出するためにクラスタのプロトタイプとして線形多様体を用いる Fuzzy  $c$ -Varieties (FCV) 法の目的関数と、Fuzzy  $c$ -Means (FCM) 法の目的関数を重み付きで足し合わせることで、楕円体状のクラスタを得る手法であり、以下の目的関数の最小化問題に定式化される。

$$L = \sum_{i=1}^N \sum_{c=1}^C u_{ci}^m \left( (x_i - b_c)^T (x_i - b_c) - \alpha \sum_{k=1}^t v_{ck}^T B_{ci} v_{ck} \right) \quad (1)$$

ただし、 $x_i = (x_{i1}, \dots, x_{in})^T$  は  $n$  次元の第  $i$  標本データを表す。 $u_{ci}$  は第  $i$  データがクラスタ  $c$  に属する度合いを表すメンバシップで、

$$u_c = \left\{ (u_{ci}) \mid \sum_{c=1}^C u_{ci} = 1, u_{ci} \in [0, 1] \right\}$$

を満たすものとする。 $b_c = (b_{c1}, \dots, b_{cn})^T$  はクラスタ  $c$  の中心である。また、 $v_{ck}$  はクラスタ  $c$  のプロトタイプを表す線形多様体を張るベクトルであり、 $t$  は線形多様体の次元を表す。

ここで,

$$B_{ci} = (\mathbf{x}_i - \mathbf{b}_c)(\mathbf{x}_i - \mathbf{b}_c)^T \quad (2)$$

であり,  $\mathbf{v}_{ck}^T B_{ci} \mathbf{v}_{ck}$  は,

$$\begin{aligned} \mathbf{v}_{ck}^T B_{ci} \mathbf{v}_{ck} &= \mathbf{v}_{ck}^T (\mathbf{x}_i - \mathbf{b}_c)(\mathbf{x}_i - \mathbf{b}_c)^T \mathbf{v}_{ck} \\ &= |(\mathbf{x}_i - \mathbf{b}_c)^T \mathbf{v}_{ck}|^2 \end{aligned} \quad (3)$$

のように表すことができ,  $\mathbf{v}_{ck}$  により定まる射影軸に射影されたデータとクラスタ中心 (平均) との 2 乗距離となる.  $\mathbf{v}_{ck}$  を主成分ベクトルと考えると  $\mathbf{v}_{ck}^T B_{ci} \mathbf{v}_{ck}$  はクラスタごとにメンバシップを考慮した主成分分析のための目的関数であるととらえられるので, 最大化することにより局所的な主成分を抽出できる.  $\alpha$  は主成分分析に対する重み係数で,  $L$  は  $\alpha = 0$  のときは FCM 法の目的関数に,  $\alpha = 1$  のときは線形多様体とデータ点との距離の総和を表す FCV 法の目的関数に等しくなる. このように, FCV 法および FCE 法の目的関数は, クラスタごとのメンバシップを考慮した主成分分析の目的関数に FCM 法の目的関数を重み付きで足し合わせることで, データの局所的な構造を考慮しながら主成分を抽出する手法であると考えられることから, 多変量解析手法と FCM 法の目的関数を重み付きで足し合わせることで, 大規模なデータベースから部分構造を考慮した局所的な特徴量を抽出する研究が行われている [6, 7].

しかし, データ行列に欠測値が含まれる場合には, FCE 法の目的関数は定義できない. そこで本章では, 高根 [12] による欠測値を含む場合の主成分分析法の目的関数に, FCM 法の目的関数を融合することにより, 欠測値を含むデータから局所的な主成分を抽出する手法を提案する.

## 2.1 欠測値を含むデータからの局所的な主成分の抽出法

高根は柴山ら [13] の欠測値が含まれる場合の線形等化法を多次元に拡張することにより, 欠測値を含むデータから主成分を抽出する手法を提案している. この手法は標準化されたデータの主成分分析にしか適用できないが, 解が解析的に求まるという特長を持っている. 本節では, 高根の手法をメンバシップが与えられた場合に拡張し, クラスタごとにメンバシップを考慮しながら局所的な主成分を抽出する方法を提案する.

欠測値を含む ( $N \times n$ ) データ行列  $X = (x_{ij})$  の  $i$  行目の要素を対角要素とする対角行列を  $D_{xi}$  で表し, クラスタ ( $c = 1, \dots, C$ ) ごとに次のモデルを定義する.

$$Y_{ci} = D_{xi} V_c + V_{0c} \quad (i = 1, \dots, N) \quad (4)$$

ただし, 求める主成分ベクトルの数を  $t$  とし,  $Y_{ci}$  は  $(n \times t)$  行列,  $V_c$ ,  $V_{0c}$  は  $(n \times t)$  の重み行列であり,  $V_c$  は次の制約条件を満たすものとする.

$$V_c^T S_c V_c = I \quad (5)$$

ここで,  $S_c$  は観測された値のみから計算されたクラスタ  $c$  における各変量の変動を対角要素とする対角行列で,  $x_{ij}$  が観測されていれば 1, 欠測値であれば 0 をとる 2 値変量  $d_{ij}$  を用いると,  $j$  番目の対角要素  $s_{cj}$  は,

$$s_{cj} = \sum_{i=1}^N d_{ij} u_{ci} (x_{ij} - b_{cj})^2$$

と表される. クラスタごとの局所的な主成分分析のための目的関数を次のように定義し, その最小化問題を考える.

$$\begin{aligned} J_c &= \sum_{i=1}^N u_{ci} \text{tr} \left( (Y_{ci} - \mathbf{1}_n \mathbf{w}_{ci}^T)^T D_{wi} \right. \\ &\quad \left. \times (Y_{ci} - \mathbf{1}_n \mathbf{w}_{ci}^T) \right) \end{aligned} \quad (6)$$

ただし,  $\text{tr}$  は行列の対角要素の和 (トレース) を表す.  $\mathbf{1}_n$  は要素がすべて 1 の  $n$  次元ベクトル,  $D_{wi}$  は第  $j$  対角要素が  $d_{ij}$  である対角行列とし,  $\mathbf{w}_{ci}$  はデータごとに与えられる  $t$  次元ベクトルとする.  $t = 1$  のときは, (6) 式の最小化は (4) 式で表されるベクトル  $Y_{ci}$  の要素をすべて  $\mathbf{w}_{ci}$  に等しくすることを表すので, クラスタごとにデータを線形等化することに相当する.

まず,  $V_c, V_{0c}$  を暫時既知として (6) 式を  $w_{ci}$  についてのみ最小化すると, 求める  $\hat{w}_{ci}^T$  は,

$$\hat{w}_{ci}^T = \mathbf{1}_n^T D_{wi} Y_{ci} / n_i \quad (7)$$

$$n_i = \mathbf{1}_n^T D_{wi} \mathbf{1}_n \quad (8)$$

となる. したがって  $w_{ci}$  について最小化した (6) 式は,

$$J_c^*(V_c, V_{0c}) = \sum_{i=1}^N u_{ci} \text{tr}(Y_{ci}^T C_i Y_{ci}) \quad (9)$$

と表せる. ただし

$$C_i = (Q_{n/Dwi})^T D_{wi} Q_{n/Dwi} \quad (10)$$

$$Q_{n/Dwi} = I - \mathbf{1}_n \mathbf{1}_n^T D_{wi} / n_i \quad (11)$$

とおいた. ここで, (4) 式を (9) 式に代入すると,

$$\begin{aligned} J_c^*(V_c, V_{0c}) = & \text{tr}(V_c^T A_{c1} V_c) + 2\text{tr}(V_c^T A_{c2} V_{0c}) \\ & + \text{tr}(V_{0c}^T A_{c3} V_{0c}) \end{aligned} \quad (12)$$

と書き換えることができる. ただし,

$$A_{c1} = \sum_{i=1}^N u_{ci} D_{xi} C_i D_{xi} \quad (13)$$

$$A_{c2} = \sum_{i=1}^N u_{ci} D_{xi} C_i \quad (14)$$

$$A_{c3} = \sum_{i=1}^N u_{ci} C_i \quad (15)$$

とする. さらに,  $J_c^*(V_c, V_{0c})$  を  $V_{0c}$  に関して最小化する.  $\partial J_c^* / \partial V_{0c} = O$  から, 次の関係が得られる.

$$A_{c3} V_{0c} = -A_{c2}^T V_c \quad (16)$$

これを解いて最適解  $\hat{V}_{0c}$  を求めればよい. しかし, データ行列に欠測値が含まれる場合には,  $A_{c3}$  がランク落ちした行列となるため, 通常の逆行列は存在しない. そのため,  $A_{c3}$  のムーア・ペンローズ逆行列  $A_{c3}^+$  を使って (16) 式から  $\hat{V}_{0c}$  を求めると,

$$\hat{V}_{0c} = -A_{c3}^+ A_{c2}^T V_c \quad (17)$$

となる. よって,  $V_{0c}$  について最小化した (12) 式は次のようになる.

$$\begin{aligned} J_c^{**}(V_c) = & \text{tr}(V_c^T (A_{c1} - A_{c2} A_{c3}^+ A_{c2}^T) V_c) \\ = & \text{tr}(V_c^T A_c V_c) \end{aligned} \quad (18)$$

ただし  $A_c = A_{c1} - A_{c2} A_{c3}^+ A_{c2}^T$  とおいた. ここで,

$$V_c^T S_c V_c = V_c^T (S_c^{1/2})^T (S_c^{1/2}) V_c = I \quad (19)$$

と書き換えられるので  $\tilde{V}_c = S_c^{1/2} V_c$  とおくと制約条件は,

$$\tilde{V}_c^T \tilde{V}_c = I \quad (20)$$

となり, (18) 式は,

$$J_c^{**}(\tilde{V}_c) = \text{tr}(\tilde{V}_c^T (S_c^{-1/2})^T A_c (S_c^{-1/2}) \tilde{V}_c) \quad (21)$$

と表すことができ, (20) 式の制約条件のもとで  $J_c^{**}$  を  $\tilde{V}_c$  について最小化すると, 次の一般化固有方程式が導かれる.

$$S_c^{-1/2} A_c S_c^{-1/2} \tilde{V}_c = \tilde{V}_c \Delta_c \quad (22)$$

ただし,  $\Delta_c$  は固有値を対角要素に持つ対角行列である. (22) 式の固有値問題を解くことにより,  $J_c$  を最小とする  $V_c$  を求めることができる.

ここで、(22) 式の固有値問題と主成分分析との関係を調べるために、データ行列に欠測値が含まれない場合について考察する。いま、クラス  $c$  において 1 本の主成分ベクトル  $\mathbf{v}_c = (v_{c1}, \dots, v_{cn})$  を推定するとし、

$$y_{cij} = v_{cj}x_{ij} + v_{0cj} \quad ; i = 1, \dots, N, j = 1, \dots, n \quad (23)$$

なるモデルを仮定すると、目的関数  $J_c$  は、

$$J_c = \sum_{i=1}^N \sum_{j=1}^n u_{ci} (y_{cij} - w_{ci})^2 \quad (24)$$

と定義され、最適性の必要条件から、

$$\hat{w}_{ci} = \frac{1}{n} \sum_{j=1}^n y_{cij} = \bar{y}_{ci} \quad (25)$$

$$\hat{\mathbf{v}}_{0c} = - \frac{\sum_{i=1}^N u_{ci} D_{xi}}{\sum_{i=1}^N u_{ci}} \mathbf{v}_c \quad (26)$$

となり、

$$J_c^{**}(\mathbf{v}_c) = \mathbf{v}_c^T \left( \sum_{i=1}^N u_{ci} D_{xi} C_i D_{xi} - \frac{\left( \sum_{i=1}^N u_{ci} D_{xi} \right) C_i \left( \sum_{i=1}^N u_{ci} D_{xi} \right)}{\sum_{i=1}^N u_{ci}} \right) \mathbf{v}_c \quad (27)$$

と変形される。ここで、

$$C_i = \left( I - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T \right)^T \left( I - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T \right) \quad (28)$$

であり、 $\mathbf{b}_c$  の要素を対角要素にもつ対角行列を  $D_{bc}$ 、

$$\mathbf{b}_c = \frac{\sum_{i=1}^N u_{ci} \mathbf{x}_i}{\sum_{i=1}^N u_{ci}} \quad (29)$$

とすると、

$$D_{xi} C_i D_{xi} = D_{xi}^2 - \frac{1}{n} \mathbf{x}_i \mathbf{x}_i^T \quad (30)$$

$$\begin{aligned} & \frac{\left( \sum_{i=1}^N u_{ci} D_{xi} \right) C_i \left( \sum_{i=1}^N u_{ci} D_{xi} \right)}{\sum_{i=1}^N u_{ci}} \\ &= \sum_{i=1}^N u_{ci} D_{bc}^2 - \frac{1}{n} \sum_{i=1}^N u_{ci} \mathbf{b}_c \mathbf{b}_c^T \end{aligned} \quad (31)$$

となるので,  $J_c^{**} = \mathbf{v}_c^T A_c \mathbf{v}_c$  とおくと,

$$\begin{aligned} A_c &= \sum_{i=1}^N u_{ci} (D_{xi}^2 - D_{bc}^2) \\ &\quad - \frac{1}{n} \sum_{i=1}^N u_{ci} (\mathbf{x}_i - \mathbf{b}_c)(\mathbf{x}_i - \mathbf{b}_c)^T \\ &= \text{diag}(S_{fc}) - \frac{1}{n} S_{fc} \end{aligned} \quad (32)$$

となる．ここで,  $S_{fc}$  はクラス  $c$  におけるファジィ散布行列であり,  $\text{diag}$  は対角要素のみからなる対角行列を表す．したがって,  $\text{diag}(S_{fc}) = S_c$  から,  $S_c^{-1/2} A_c S_c^{-1/2}$  は,

$$S_c^{-1/2} A_c S_c^{-1/2} = I - S_c^{-1/2} S_{fc} S_c^{-1/2} / n \quad (33)$$

$$= I - R_c / n \quad (34)$$

とできる．ここで,  $R_c$  はファジィクラス  $c$  における相関係数行列であり, その大きい固有値に対応する固有ベクトルはメンバシップを考慮しながら標準化したデータの主成分ベクトルに相当する [19]．(22) 式の固有値問題は,  $t = 1$  のとき, 固有値を  $\lambda$  とおくと,

$$S_c^{-1/2} A_c S_c^{-1/2} \tilde{\mathbf{v}}_c = (I - R_c / n) \tilde{\mathbf{v}}_c = \lambda \tilde{\mathbf{v}}_c \quad (35)$$

$$R_c \tilde{\mathbf{v}}_c = (1 - \lambda) n \tilde{\mathbf{v}}_c \quad (36)$$

となることから, 相関係数行列の最大固有値に対応する固有ベクトルが,  $S_c^{-1/2} A_c S_c^{-1/2}$  の最小固有値に対応する固有ベクトルに等しいことがわかる．同様に, 2 本以上の主成分ベクトルを求める際にも,  $S_c^{-1/2} A_c S_c^{-1/2}$  の小さい固有値に対応する固有ベクトルを求めることにより, 相関係数行列の大きい固有値に対応する固有ベクトルを求めることができる．

## 2.2 欠測値を含むデータの主成分分析とファジィクラスタリングの同時適用法

本節では, 前節で提案した局所的な主成分分析のための目的関数に欠測値を含む場合に拡張した FCM 法の目的関数を融合することにより, 欠測値を含むデータの主成分分析とファジィクラスタリングの同時適用法を提案する．

FCM 法における欠測値の処理法としては, 宮本ら [14] がメンバシップのべき乗を用いる標準的な方法とエントロピー正則化 [20] を用いる方法の二つのバリエーションごとに, 欠測値を無視して残りの座標で距離を定義する手法や欠測値に重み付きの平均値を代入する手法を提案している．また, Timm ら [15] も標準的な方法について同様の考察を行い, 欠測値を無視する手法の方が, 代入する場合よりもよりあいまいなメンバシップの割り当てを行う傾向があると報告している．

欠測値を無視することにより, 標準的な FCM 法の目的関数は, 次のように書き換えられる．

$$\psi = \sum_{c=1}^C \sum_{i=1}^N u_{ci}^m \sum_{j=1}^n d_{ij} (x_{ij} - b_{cj})^2 \quad (37)$$

前節の局所的な主成分分析の目的関数と欠測値を含むデータの FCM 法の目的関数を融合することにより, 欠測値を含むデータのファジィクラスタリングと主成分分析の同時適用法の定式化を考える．ただし, 前節の局所的な主成分分析ではデータをクラス  $c$  ごとに標準化したデータの主成分の抽出を考えていたのに対し, 前述の標準的な FCM 法もしくはエントロピー正則化を用いる FCM 法では, クラス  $c$  ごとくデータの標準化を考慮していない．そこで, クラス  $c$  ごとく容量の最適化も考慮に入れながらファジィクラス  $c$  を得ることができる K-L 情報量正則化 [18] を導入することにより, ラグランジュ乗数法による目的関数を以下のように定義する．

$$\begin{aligned} L &= \sum_{c=1}^C \sum_{i=1}^N u_{ci} \sum_{j=1}^n d_{ij} \left( \alpha \sum_{k=1}^t (y_{cijk} - w_{cik})^2 \right. \\ &\quad \left. + (1 - \alpha) (x_{ij} - b_{cj})^2 \right) \end{aligned}$$

$$\begin{aligned}
& +\beta \sum_{c=1}^C \sum_{i=1}^N u_{ci} \log \frac{u_{ci}}{\pi_c} \\
& - \sum_{i=1}^N \gamma_i \left( \sum_{c=1}^C u_{ci} - 1 \right) - \tau \left( \sum_{c=1}^C \pi_c - 1 \right)
\end{aligned} \tag{38}$$

ここで、 $\pi_c$  はクラスタ  $c$  の容量を表す変数で、 $\gamma_i$  および  $\tau$  はラグランジュ乗数である。第 1 項目の  $(y_{cij} - w_{cik})^2$  は前節で述べた主成分分析のための項であり、これを最小化することによって局所的な主成分の抽出を行う。 $(x_{ij} - b_{cj})^2$  は FCM 法の目的関数を表す。 $\alpha$  は主成分分析とクラスタリングのトレードオフを表すパラメータであり、 $\alpha = 0$  のときには FCM 法と同様の球状のクラスタが得られ、 $\alpha$  を 1 に近づけるにつれて、主成分の抽出を重視した線形多様体状のクラスタが得られるようになる。第 2 項目はファジィクラスタを得るための K-L 情報量を表す項であり、メンバシップ  $u_{ci}$  とクラスタ容量  $\pi_c$  を近づけることにより、クラスタ容量を考慮しながらクラスタリングが行われる。重み  $\beta$  はファジィ度を制御する係数で、標準的な FCM 法における  $m$  と同様に、大きくするにしたがってよりあいまいなデータ分割が得られるようになる。第 3, 4 項目はそれぞれ、メンバシップの和が 1 という制約条件、 $\pi_c$  の和が 1 という制約条件を表す項である。文献 [18] ではマハラノビス距離を用い、変量間の分散共分散行列の要素も決定変数としているが、ここでは変量間の相関関係は主成分によって表されるという観点から、クラスタ容量のみを決定変数として用いている。

$L$  の最適性の必要条件  $\partial L / \partial u_{ci} = 0$  より、

$$u_{ci} = \frac{\pi_c \exp(E_{ci})}{\sum_{l=1}^C \pi_l \exp(E_{li})} \tag{39}$$

が得られる。ただし、

$$\begin{aligned}
E_{li} = & -\frac{1}{\beta} \sum_{j=1}^n d_{ij} \left( \alpha \sum_{k=1}^t (y_{lij} - w_{lik})^2 \right. \\
& \left. + (1 - \alpha)(x_{ij} - b_{lj})^2 \right)
\end{aligned} \tag{40}$$

である。また、 $\partial L / \partial \pi_c = 0$ 、 $\partial L / \partial b_{cj} = 0$  から、それぞれ、

$$\pi_c = \frac{1}{N} \sum_{i=1}^N u_{ci} \tag{41}$$

$$b_{cj} = \frac{\sum_{i=1}^N u_{ci} d_{ij} x_{ij}}{\sum_{i=1}^N u_{ci} d_{ij}} \tag{42}$$

である。

クラスタリングアルゴリズムを以下に示す。

step:1 メンバシップ  $u_{ci}$ ,  $i = 1, \dots, N$ ,  $c = 1, \dots, C$  を乱数により  $\sum_{c=1}^C u_{ci} = 1$  となるように定める。

step:2 (22) 式の固有値問題を解き、局所的な主成分ベクトル、 $y_{cij}$ ,  $w_{ci}$  を求める。

step:3 (41) 式より、 $\pi_c$  を更新する。

step:4 (42) 式より、 $b_{cj}$  を更新する。

step:5 (39) 式より、 $u_{ci}$  を更新する。

step:6 小さな正数  $\varepsilon$  に対して収束判定条件

$$\max_{c,i} |u_{ci}^{NEW} - u_{ci}^{OLD}| < \varepsilon$$

を満たせば終了。そうでなければ、ステップ 2 へ。

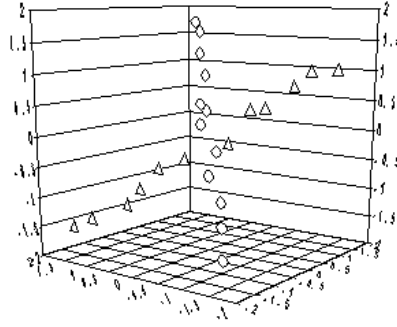


図 1: Clustering result with no missing value

### 2.3 因子負荷量の求め方

クラスタごとに抽出された主成分と各変量の関係を定量的に示すためには、第  $j$  変量  $x_{.j}$  のクラスタ  $c$  における第  $k$  主成分  $z_{c.k}$  に対するファジィ因子負荷量  $\eta_{ckj}$  [19] を用いることができる。

$$\eta_{ckj} = \frac{Cov_c\{z_{c.k}, x_{.j}\}}{\sqrt{V_c\{z_{c.k}\}V_c\{x_{.j}\}}} \quad (43)$$

ただし、 $Cov_c\{z_{c.k}, x_{.j}\}$ 、 $V_c\{z_{c.k}\}$ 、 $V_c\{x_{.j}\}$  はそれぞれ、クラスタ  $c$  における主成分  $z_{c.k}$  と変量  $x_{.j}$  の共分散、主成分の分散、および変量の分散である。

第  $i$  データ点の主成分得点を要素とするベクトル  $z_{ci}$  は、欠測値を含まない場合については、主成分行列  $V_c$  とクラスタ中心  $b_c$  を用いて、

$$z_{ci} = V_c^T (x_i - b_c) \quad (44)$$

により求められるが、欠測値を含むデータについては適用することができない。そこで、欠測値を含むデータについては、データの低階数近似の立場 [11] から、

$$x_i = V_c z_{ci} + b_c + e_i \quad (45)$$

とおいたときに、観測値の 2 乗誤差の和  $e_i^T D_{wi} e_i$  が最小となるように、

$$z_{ci} = (V_c^T D_{wi} V_c)^{-1} V_c^T D_{wi} (x_i - b_c) \quad (46)$$

と求める。これは、欠測値を含むデータが主成分ベクトルにより張られる線形多様体上、もしくはその線形多様体からの距離が最も近くなる点に存在すると仮定して欠測値を補完し、(44) 式の値を求めることに等しい。

## 3 数値例

本章では、提案手法の有効性を検証するために、以下の二つの数値例を示す。

### 3.1 人工データを用いた数値実験

まず、表 1 に示す 3 次元 24 個の人工データを用いて、実験を行った。データ集合は 3 次元空間内で二つの直線状に布置されており、欠測値がない場合には図 1 に示されるように提案法を用いて 2 で表される二つの線形クラスタに分割され、表 2 の因子負荷量を持つ主成分とクラスタ中心が得られた。ただし、パラメータは  $\alpha = 0.99$ 、 $\beta = 0.1$  とした。この人工データではクラスタごとに各変量がほぼ等しい分散を有しているために、得られるクラスタリング結果はデータの標準化を考慮しない FCV 法と同様なものとなっている。

つぎに、このデータ集合に対して、24 個すべてのサンプルからランダムに一つずつ、表 1 で太字で示された値を欠落させることにより、データ行列の要素の 1/3 が欠測値となるテストデータ集合を作成した。すべてのサンプルデータに欠測値が含まれるために、欠測値を含むサンプルを削除すると全く分析することができず、回帰分析などによりデータを補完して分析することも困難な例となっている。作成した不完全データを欠測値が無い場合と同じパラメータを用いて提案法により二つのクラ



表 1: Artificial data set

—	$x_1$	$x_2$	$x_3$
1	-1.43	-1.28	<b>-1.56</b>
2	<b>-0.96</b>	-0.93	-1.21
3	<b>-0.80</b>	-0.77	-0.87
4	<b>-0.63</b>	-0.41	-0.52
5	-0.36	-0.35	<b>-0.17</b>
6	0.18	0.26	<b>0.38</b>
7	<b>0.40</b>	0.55	0.43
8	0.82	0.77	<b>0.88</b>
9	0.99	<b>1.02</b>	1.23
10	<b>1.33</b>	1.28	1.58
11	1.53	<b>1.54</b>	1.73
12	<b>-0.06</b>	0.20	0.18
13	<b>-1.71</b>	1.54	-1.46
14	-1.54	<b>1.28</b>	-1.33
15	-1.17	<b>0.83</b>	-1.10
16	<b>-0.90</b>	0.77	-0.87
17	-0.73	<b>0.51</b>	-0.53
18	-0.36	<b>0.25</b>	-0.40
19	0.18	-0.26	<b>-0.16</b>
20	0.45	-0.51	<b>0.39</b>
21	0.52	-0.77	<b>0.43</b>
22	0.99	<b>-1.03</b>	0.76
23	1.26	<b>-1.18</b>	0.99
24	1.53	-1.54	<b>1.03</b>

表 2: Result of analysis with no missing value

cluster center			
—	$x_1$	$x_2$	$x_3$
$c = 1$	-0.13	-0.01	-0.17
$c = 2$	0.11	0.16	0.16

factor loading			
—	$x_1$	$x_2$	$x_3$
$c = 1$	0.94	-0.93	0.94
$c = 2$	0.94	0.94	0.95

スタに分割した結果を表 3 に示す．欠測値の有無に関わらず，ほぼ等しいクラスタ中心および主成分が得られており，提案法が欠測値を含む不完全データから局所的な構造を考慮しながら主成分を抽出するのに有効であることが分かる．

表 3: Result of analysis with missing values

cluster center			
—	$x_1$	$x_2$	$x_3$
$c = 1$	0.38	-0.07	-0.22
$c = 2$	-0.02	-0.11	0.01

factor loading			
—	$x_1$	$x_2$	$x_3$
$c = 1$	0.85	-0.57	0.84
$c = 2$	0.82	0.65	0.84

表 4: Result of FCV

variable	cluster center		factor loading( $c = 1$ )		factor loading( $c = 2$ )	
	$c = 1$	$c=2$	1st PC	2nd PC	1st PC	2nd PC
1.	-0.197	0.198	0.000	-0.003	-0.080	0.117
2.	-0.022	-0.022	0.223	<b>-0.383</b>	-0.227	<b>-0.449</b>
3.	0.000	0.000	-0.043	<b>0.456</b>	0.087	0.338
4.	0.021	0.021	-0.273	0.289	0.332	<b>0.428</b>
5.	<b>-0.786</b>	<b>0.790</b>	0.252	<b>0.471</b>	0.257	<b>-0.404</b>
6.	<b>-0.775</b>	<b>0.779</b>	0.291	<b>0.481</b>	0.255	<b>-0.378</b>
7.	<b>-0.765</b>	<b>0.769</b>	0.205	<b>0.468</b>	<b>0.397</b>	<b>-0.443</b>
8.	<b>-0.766</b>	<b>0.770</b>	0.112	<b>0.413</b>	<b>0.381</b>	<b>-0.472</b>
9.	<b>-0.768</b>	<b>0.771</b>	0.164	<b>0.412</b>	0.326	<b>-0.444</b>
10.	-0.344	0.345	<b>0.568</b>	0.262	<b>-0.647</b>	0.335
11.	-0.177	0.178	<b>0.687</b>	0.246	<b>-0.750</b>	<b>0.453</b>
12.	-0.105	0.105	<b>0.391</b>	0.034	<b>-0.422</b>	<b>0.421</b>
13.	-0.140	0.140	<b>0.830</b>	<b>0.449</b>	<b>-0.669</b>	<b>0.393</b>
14.	0.058	-0.059	<b>0.426</b>	0.040	<b>-0.693</b>	0.260
15.	0.033	-0.034	<b>0.698</b>	<b>0.435</b>	<b>-0.522</b>	0.182
16.	0.033	-0.034	<b>0.641</b>	<b>0.387</b>	-0.306	0.258
17.	-0.140	0.140	<b>-0.601</b>	<b>0.552</b>	<b>0.608</b>	<b>0.490</b>
18.	-0.162	0.163	<b>-0.577</b>	<b>0.687</b>	<b>0.630</b>	<b>0.608</b>
19.	-0.197	0.198	<b>-0.628</b>	<b>0.709</b>	<b>0.645</b>	<b>0.658</b>
20.	-0.089	0.089	<b>-0.554</b>	<b>0.712</b>	<b>0.562</b>	<b>0.740</b>

### 3.2 POS データの分析

つぎに、文献 [7] で用いられた POS(Point of Sales) データによる数値実験を通して、提案法の特徴を示す。データは 2 店舗のスーパーマーケットで 1 年間にわたって収集・蓄積されたもので、12 月 31 日などの特異な傾向を示す日を除いた 333 日分について、各変量ごとに平均 0、分散 1 に基準化したものを用いた。用いた変量は、2 店舗の来客数、日配品点数に暦、気象情報を加えた以下の 20 個の変量である。

1. 祝日, 2. 金曜, 3. 土曜, 4. 日曜, 5. 平均気温,
- 6-9. 気温 (9 時, 12 時, 15 時, 18 時), 10. 相対湿度,
- 11, 12. 天気概況 (昼, 夜), 13. 日降水量,
- 14-16. 降水量 (9-12 時, 12-15 時, 15-18 時),

表 5: Result of proposed method with no missing value

variable	cluster center		factor loading( $c = 1$ )		factor loading( $c = 2$ )	
	$c = 1$	$c=2$	1st PC	2nd PC	1st PC	2nd PC
1.	-0.021	0.020	0.056	0.005	0.036	0.014
2.	0.033	-0.031	0.192	-0.163	-0.120	-0.203
3.	-0.070	0.065	-0.072	0.128	0.106	0.081
4.	-0.038	0.035	-0.120	0.144	0.077	<b>0.352</b>
5.	<b>-0.829</b>	<b>0.777</b>	<b>0.669</b>	<b>0.660</b>	<b>-0.841</b>	-0.022
6.	<b>-0.823</b>	<b>0.771</b>	<b>0.666</b>	<b>0.646</b>	<b>-0.822</b>	-0.033
7.	<b>-0.820</b>	<b>0.768</b>	<b>0.641</b>	<b>0.674</b>	<b>-0.854</b>	0.047
8.	<b>-0.812</b>	<b>0.761</b>	<b>0.624</b>	<b>0.678</b>	<b>-0.866</b>	0.065
9.	<b>-0.814</b>	<b>0.762</b>	<b>0.634</b>	<b>0.661</b>	<b>-0.857</b>	0.028
10.	-0.279	0.261	<b>0.488</b>	-0.034	0.035	<b>-0.497</b>
11.	-0.119	0.112	0.339	-0.200	0.231	<b>-0.544</b>
12.	-0.107	0.100	0.330	-0.103	0.180	-0.246
13.	-0.155	0.145	<b>0.488</b>	-0.105	0.165	<b>-0.525</b>
14.	-0.093	0.087	<b>0.387</b>	-0.094	0.113	<b>-0.411</b>
15.	-0.089	0.083	<b>0.430</b>	-0.083	0.109	<b>-0.408</b>
16.	-0.088	0.082	0.296	-0.181	0.097	-0.290
17.	-0.194	0.182	-0.195	<b>0.442</b>	-0.008	<b>0.588</b>
18.	-0.227	0.213	-0.188	<b>0.457</b>	0.025	<b>0.607</b>
19.	-0.257	0.241	-0.168	<b>0.536</b>	0.024	<b>0.639</b>
20.	-0.138	0.130	-0.252	<b>0.383</b>	0.143	<b>0.588</b>

17. 来客数 (A 店), 18. 日配品点数 (A 店),

19. 来客数 (B 店), 20. 日配品点数 (B 店)

ただし、日配品とは工場で生産される食品のうち、数日中に消費されるものを指し、その総売上点数が日配品点数である。また、暦に関する変数はおのおのあてはまるか否かを表すダミー変数であり、天気概況は降水の状況により 0 (晴れまたは曇り) から 3 (大雨) の値を与えた。

まず、相関係数行列を用いる提案法と分散共分散行列を用いる FCV 法の違いを調べるために、333 日分の POS データのすべてを用いて提案法で分析した結果と文献 [7] で得られた FCV 法による分析結果の比較を行った。おのおのの手法で得られたクラスタ中心と二つの主成分のファジィ因子負荷量を表 4、表 5 に示す。表中、絶対値が 0.35 以上のものは太字で記載している。ただし、提案法におけるパラメータは  $\alpha = 0.9$ ,  $\beta = 1.0$  とした。

表 4、表 5 の 2~3 列目から、いずれの手法においても太字で示された気温に関する変数でクラスタ中心が分かれており、温暖期と寒冷期の二つのクラスタに分類されている。一方、表中の 4~7 列目に示されたファジィ因子負荷量を比較すると、FCV 法ではいずれのクラスタにおいても第 1 主成分が降水量や来店客数・日配品点数といった変数と相関が大きいものに対して、提案法では第 1 主成分が気温の影響を強く反映している。これは、主に気温の高低を加味したクラスタリングの結果、おのおののクラスタに属するデータの気温の分散が小さくなっているために、FCV 法ではその影響が弱まっているものに対して、提案法ではメンバシップを考慮した標準化を行っているので FCV 法では抽出できない特徴量が抽出できていると考えられる。このように、相関係数行列を用いる提案法は、クラスタリングによる層別の影響で FCV 法が無視してしまう変数についても、標準化して考慮することができるという利点があることが分かる。また、この特徴は主成分分析において解が変量の測定単位により影響を受けるという性質、すなわち尺度不変 (scale invariant) でないという性質に対応するものである。

つぎに、観測データに欠測値が含まれる場合を考える。333 日分のサンプルデータのすべてについて、気象に関する変数 (5~16) からランダムに二つずつの値を欠落させることにより、10%の欠測値が含まれるデータ行列を作成した。前節と同様に、すべてのサンプルが欠測値を有するために、

表 6: Result of FCV with missing values

variable	cluster center		factor loading( $c = 1$ )		factor loading( $c = 2$ )	
	$c = 1$	$c=2$	1st PC	2nd PC	1st PC	2nd PC
1.	-0.018	0.021	0.042	0.032	0.001	0.020
2.	0.016	-0.026	0.164	<b>-0.395</b>	-0.327	0.114
3.	-0.040	0.052	0.002	<b>0.372</b>	0.169	0.333
4.	-0.054	0.056	0.233	0.333	<b>0.379</b>	0.270
5.	<b>-0.791</b>	<b>0.748</b>	0.300	<b>0.366</b>	-0.010	-0.176
6.	<b>-0.761</b>	<b>0.720</b>	0.302	<b>0.375</b>	-0.067	-0.166
7.	<b>-0.771</b>	<b>0.760</b>	0.190	<b>0.391</b>	0.051	-0.219
8.	<b>-0.775</b>	<b>0.707</b>	0.224	<b>0.369</b>	-0.007	-0.239
9.	<b>-0.797</b>	<b>0.738</b>	0.279	<b>0.359</b>	0.012	-0.210
10.	-0.274	0.208	<b>0.580</b>	0.195	<b>-0.418</b>	<b>0.452</b>
11.	-0.054	0.112	<b>0.611</b>	0.204	<b>-0.445</b>	<b>0.541</b>
12.	-0.107	0.099	<b>0.391</b>	0.068	-0.305	<b>0.397</b>
13.	-0.036	0.054	<b>0.659</b>	0.195	<b>-0.466</b>	<b>0.564</b>
14.	-0.067	0.104	<b>0.472</b>	0.028	<b>-0.391</b>	<b>0.398</b>
15.	0.093	-0.066	<b>0.505</b>	0.200	-0.269	0.280
16.	0.006	0.016	<b>0.493</b>	0.151	-0.269	0.280
17.	-0.206	0.229	<b>-0.454</b>	<b>0.611</b>	<b>0.677</b>	0.249
18.	-0.263	0.261	<b>-0.411</b>	<b>0.718</b>	<b>0.749</b>	<b>0.379</b>
19.	-0.263	0.261	<b>-0.445</b>	<b>0.739</b>	<b>0.749</b>	<b>0.366</b>
20.	-0.141	0.148	<b>-0.410</b>	<b>0.725</b>	<b>0.705</b>	<b>0.491</b>

欠測値を含むサンプルを削除すると全く分析することができない例となっている．文献 [16] の FCV 法を不完全データに拡張した手法および提案法を用いて分析した結果を表 6，表 7 に示す．ただし，パラメータは欠測値が無い場合と同じものを用い，表中，絶対値が 0.35 以上のものを太字で示した．

表 6，表 7 の 2～3 列目から，欠測値の有無に関わらず，温暖期と寒冷期に分割されていることが分かる．また，表中の 4～7 列目に示されたファジィ因子負荷量には，以下のような特徴が見られる．欠測値がない場合には目的関数が FCV 法と等しくなる文献 [16] の手法では，気温に関する変量の因子負荷量がさらに小さくなったものの，欠測値がない場合の FCV 法と類似した結果が得られている．一方，提案法では，第 1 主成分において欠測値を含む気象関連の変量の因子負荷量が小さくなり，FCV 法の結果に近づく傾向が見られたものの，データの層別に用いられた気温に関する変数の影響も無視されておらず，欠測値が無い場合と類似した相関関係がとらえられている．これらの違いは，文献 [16] の手法がファジィ散布行列を用いる FCV 法を欠測値が含まれるデータに拡張したものであるのに対して，提案手法が相関係数行列を用いることによりクラスごとにデータの基準化を考慮していることを示している．以上のように，提案法はファジィ分散共分散行列の固有値問題に帰着される FCV 法とは若干異なる主成分を抽出するものの，欠測値を含む場合にも妥当な結果を得られる手法であることが分かる．

## 4 おわりに

欠測値の補完などの手続きをせずに，欠測値を含むデータから局所的な主成分を抽出する手法を提案した．Bezdek らの FCV 法がデータの分散共分散行列を用いた主成分分析とファジィクラスタリングの同時適用法ととらえられるのに対して，提案法は相関係数行列を用いた主成分分析とファジィクラスタリングの同時適用法である．FCV 法を不完全データに拡張した本多ら [16] の手法ではデータ行列の低階数近似行列を保持するために，サンプルデータ数が多くなるにつれてメモリの消費量が増すのに対して，提案法では欠測値があるにもかかわらず固有値問題を解くことにより解析的に主成

表 7: Result of proposed method with missing values

variable	cluster center		factor loading( $c = 1$ )		factor loading( $c = 2$ )	
	$c = 1$	$c=2$	1st PC	2nd PC	1st PC	2nd PC
1.	-0.016	0.016	0.046	0.040	0.038	-0.022
2.	0.049	-0.047	0.280	-0.073	-0.247	-0.114
3.	-0.083	0.080	-0.166	0.100	0.159	-0.001
4.	-0.062	0.059	-0.211	0.082	0.266	0.268
5.	<b>-0.793</b>	<b>0.742</b>	0.286	<b>0.783</b>	<b>-0.680</b>	<b>0.354</b>
6.	<b>-0.764</b>	<b>0.718</b>	0.295	<b>0.751</b>	<b>-0.664</b>	0.304
7.	<b>-0.759</b>	<b>0.755</b>	0.227	<b>0.785</b>	<b>-0.630</b>	<b>0.454</b>
8.	<b>-0.762</b>	<b>0.698</b>	0.269	<b>0.755</b>	<b>-0.647</b>	<b>0.411</b>
9.	<b>-0.791</b>	<b>0.736</b>	0.271	<b>0.769</b>	<b>-0.674</b>	<b>0.391</b>
10.	-0.314	0.229	<b>0.377</b>	0.163	-0.174	<b>-0.444</b>
11.	-0.106	0.138	<b>0.354</b>	0.038	-0.022	<b>-0.538</b>
12.	-0.109	0.091	0.336	0.056	0.017	-0.277
13.	-0.132	0.114	<b>0.421</b>	0.097	-0.068	<b>-0.504</b>
14.	-0.075	0.104	<b>0.352</b>	0.044	-0.073	<b>-0.363</b>
15.	-0.098	0.077	0.323	0.111	-0.059	-0.336
16.	-0.084	0.079	0.311	0.058	-0.055	-0.288
17.	-0.228	0.219	<b>-0.433</b>	0.335	0.332	<b>0.512</b>
18.	-0.263	0.252	<b>-0.448</b>	<b>0.359</b>	<b>0.380</b>	<b>0.505</b>
19.	-0.294	0.283	<b>-0.463</b>	<b>0.431</b>	<b>0.392</b>	<b>0.525</b>
20.	-0.185	0.177	<b>-0.479</b>	0.261	<b>0.476</b>	<b>0.421</b>

分ベクトルが求まるため、データ数が多い場合でも効率的に分析できるという特長がある。ただし、主成分分析の重みやファジィ度を制御するパラメータなどを試行錯誤により定めなければならない点は従来法と同様に問題点として残っており、適応的にクラスタ形状を決定する手法の開発が課題である。数値例では、FCV 法と提案法での欠測値がない場合の結果の比較をとおして、クラスタごとのデータの標準化を考慮するか否かで結果に違いを生じることを示した。また、欠測値がある場合についても妥当な主成分が抽出されることを示した。提案法は欠測値がランダムに発生しているという仮定に基づいて定式化がなされているが、実データを扱う際には欠測値の発生に偏りがある場合も少なくない。欠測値の発生機構を考慮に入れたモデルの作成が今後の課題である。

## 参考文献

- [1] J. C. Bezdek: *Pattern recognition with fuzzy objective function algorithms*, Plenum press (1981)
- [2] J. C. Bezdek, C. Coray, R. Gunderson and J. Watson: Detection and characterization of cluster substructure II. Fuzzy  $c$ -varieties and convex combinations thereof; *SIAM Jour. of Appl. Math.*, Vol.40, No.2, pp.358–372 (1981)
- [3] R. N. Dave: An adaptive fuzzy  $c$ -elliptotype clustering algorithm; *Proc. of the North American Fuzzy Information Processing Society :Quarter Century of Fuzziness*, Vol.1, pp.9–12 (1990)
- [4] 馬屋原, 中森: 線形多様体クラスタリングと楕円形ファジィモデル; 日本ファジィ学会誌, Vol.10, No.1, pp.142–149 (1998)
- [5] 馬屋原, 宮本: 次元係数の正則化による線形ファジィクラスタリング; 日本ファジィ学会誌, Vol.12, No.4, pp.552–561 (2000)
- [6] 本多, 山川, 市橋, 三好, 奥山: ファジィクラスタリングと回帰と主成分の同時分析法; システム制御情報学会論文誌, Vol. 13, No. 5, pp. 236–243 (2000)

- [7] 呉, 本多, 新山, 市橋: ファジィクラスタリングを用いる外的基準に独立な局所的主成分の抽出法; 日本ファジィ学会誌, Vol. 12, No. 6, pp. 826–834 (2000)
- [8] 竹内: 統計学辞典, 東洋経済新報社 (1989)
- [9] A. Ruhe: Numerical computation of principal components when several observations are missing; *Tech Rep. UMINF-48-74, Dept. Information Processing*, Umea Univ. (1974)
- [10] T. Wiberg: Computation of principal components when data are missing; *Proc. of Second Symp. Computational Statistics*, pp. 229–236 (1976)
- [11] 柴山: 欠測値を含む多変量データのための主成分分析的方法; 教育心理学研究, Vol. 40, No. 3, pp. 257–265 (1992)
- [12] 高根: 制約つき主成分分析法, 朝倉書店 (1995)
- [13] 柴山, 芝: 欠測値がある場合の線形等化法; 教育心理学研究, Vol.35, No.1, pp.86–89 (1987)
- [14] S. Miyamoto, O. Takata and K. Umayahara: Handling missing values in fuzzy c-means; *Proc. of Third Asian Fuzzy Systems Symp.*, pp. 139–142 (1998)
- [15] H. Timm and R. Kruse: Fuzzy cluster analysis with missing values; *Proc. of 17th International Conf. of the North American Fuzzy Information Processing Society*, pp. 242–246 (1998)
- [16] 本多, 杉浦, 市橋, 荒木, 久津見: 最小 2 乗基準を用いた Fuzzy  $c$ -Varieties 法における欠測値の処理法; 日本ファジィ学会誌, Vol. 13, No. 6, pp. 680–688 (2001)
- [17] K. Honda, N. Sugiura, H. Ichihashi and S. Araki: Collaborative filtering using principal component analysis and fuzzy clustering; *Web Intelligence: Research and Development*, Lecture Notes in Artificial Intelligence 2198, Springer, pp. 394–402 (2001)
- [18] 宮岸, 市橋, 本多: K-L 情報量正則化 FCM クラスタリング法; 日本ファジィ学会誌, Vol. 13, No. 4, pp. 406–417 (2001)
- [19] Y. Yabuuchi and J. Watada: Fuzzy principal component analysis and its application; *Biomedical Fuzzy and Human Sciences*, Vol.3, No.1, pp.83–92 (1997)
- [20] 宮本: クラスタ分析入門, 森北出版 (1999)