

# データマイニングにおける INDSCALの活用プロセス

奥 喜正

## 1. はじめに

IT時代ではデータを大量に容易に集計することが可能になってきたので、蓄積された大量データを効率的に圧縮して視覚化する新たなデータマイニング技術の開発が要請されている(柴田, 2001)。よって、従来よりもデータマイニングメソッドを利用する頻度が増してきた(福田・森本・徳山, 2001)。それゆえ、統計的データ解析もデータマイニングへの新たな貢献が求められてきているが、このような状況に貢献しているデータ解析メソッド(Computational Methods for Data Mining)には機械学習(Machine Learning)などから発展した、非確率構造モデル(Nonparametric Model)のクラスター分析、決定木(Tree Models)、ニューラルネットなどが、現時点では挙げられる(Giudici, 2003)。しかしながら、同様に非確率構造モデルに属する多次元尺度法(Multi Dimensional Scaling: MDS)は、データマイニングの実行時に利用される頻度が今までのところ少なかったように窺える(Kruskal, 1964)。

そこで、本稿では多次元尺度法のなかでも立体データの2相3元データ解析メソッドのINDSCAL(Individual Difference SCALing: 個人差多次元尺度法)を数年間の日本のプロ野球成績データに適用して(Carroll & Chang, 1970)、3元多次元尺度法(3元MDS)のINDSCALをデータマイニングに有効活用する方法を検討する。INDSCALの応用例には、新製品の購入数予測に利用されたものがある(岡太・宮内, 1996)。本稿では、3元データ解析でデータの経時的変動を適確に捉えることが可能なことや、データ圧縮過程で重要課題のデータ濃縮(Data Compression)をINDSCALがうまく遂行する様子と有効活用を示す。さらに、データマイニングの特徴を鑑みると、統計的データ解析メソッドの利用形態に多少の工夫を施してから、それらメソッドを利用する必要がある。福田ら(2001)は、データマイニングでクラスター分析を使用する場合には、留意点として例外的なデータや不要属性の排除を指摘した。本稿では実際のデータ解析プロセスを通じて3元MDS利用時の留意点を検討し、データマイニングへのINDSCALの有用な利用形

態を模索する。

## 2. 多次元尺度法とINDSCAL

最初に、多次元尺度法をクルスカルの方法(Kruskal, 1964)を例に説明する。多次元尺度法(MDS)とは、対象間の心理的距離など非類似性データから、潜在する少数の次元や背後にある構造を明らかにしたいときに利用される多変量データ解析手法の一つである。対象*i, j, k, l*における非類似性データ $\delta_{ij}$ ,  $\delta_{kl}$ が得られたとすると、それらをユーザが視覚的に認識できる距離量 $d_{ij}$ ,  $d_{kl}$ にできるだけその大小関係(順序関係)が保存されるように単調変換する。距離量とは三角不等式

$$d_{ik} + \leq d_{ij} + d_{jk} \quad \text{for } \forall i, j, k,$$

を満足する非負の量である。そして

$$\delta_{ij} < \delta_{kl} \Rightarrow d_{ij} \leq d_{kl} \quad \text{for } \forall i, j, k, l \quad (1)$$

のように「クルスカルの単調回帰」を利用して、非類似性データ量を距離量に変換する(Kruskal, 1965)。各々の対象を点として表現する空間を布置と呼び、布置の次元を*T*とすると対象*i, j*の次元*t*における座標を $x_{it}$ ,  $x_{jt}$ のように求められれば

$$\delta_{ij} \stackrel{m}{=} \hat{d}_{ij} \approx d_{ij} = \sqrt{\sum_{t=1}^T (x_{it} - x_{jt})^2}$$

が成立することになる。ここで、記号 $\stackrel{m}{=}$ は単調増加関係を示し、 $d_{ij}$ は条件式(1)を満足する擬似距離量で、得られた布置の $d_{ij}$ から計算される量でディスパリティとも呼ばれる。ところで、ある次元*T*(通常、*T* = 2, 3)で完全に式(1)を満足するような布置を得ることは一般には困難である。そこで、与えられた非類似性データ $\delta_{ij}$ の順序関係と、求められた布置から計算される $d_{ij}$ の順序関係がどの程度、一致しているかの評価基準にストレス*S*と呼ばれる、データに対するモデル適合度の評価尺度をMDSでは導入する。

$$S = \sqrt{\frac{\sum_{i < j} \sum (d_{ij} - \hat{d}_{ij})^2}{\sum_{i < j} \sum d_{ij}^2}}$$

順序関係が完全に保存されていれば、すなわち式(1)が成立すれば*S* = 0となり、ストレス*S*の値が0.2未満であることが最悪でも要求される。そうでない場合には、該当する次元ではデータ行列 $[\delta_{ij}]$ に十分に適合する布置は得られないとみなして、通常は解の次元を一次元上げて、改めて布置を求める。

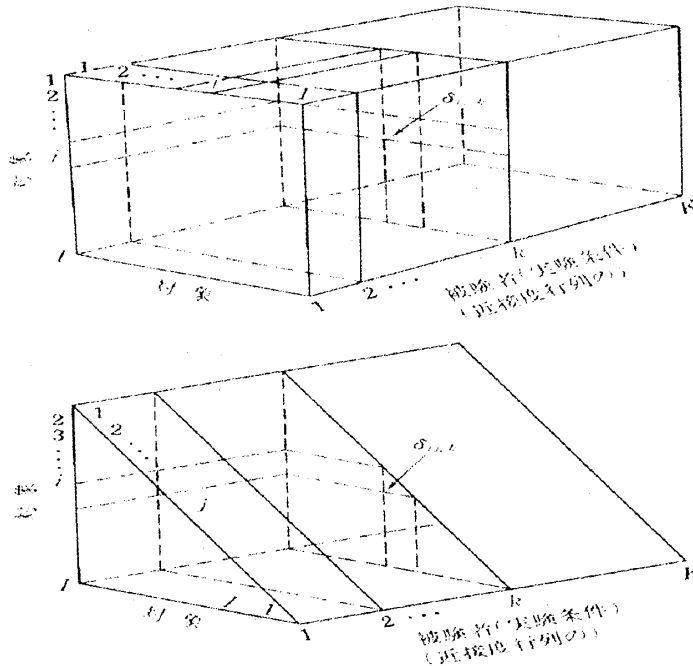


図1. 3元データ (岡太・今泉, 1990)

つぎに、3元MDSの代表モデル、INDSCAL (INDividual Difference SCALing : 個人差多次元尺度法) について説明する。INDSCALは2相3元データを解析して(図1)、全体の被験者に共通する対象を布置する「共通対象空間X」と、被験者間の違い、すなわち、個人差を表示する「被験者空間W」を同時に出力するモデルである(図2)。被験者間の認知の相違を、布置次元に与える重みづけを被験者ごとに変えることで表現するモデルである。被験者kの対象iとjの距離  $d_{ij,k}$  は

$$d_{ij,k} \approx \left[ \sum_{t=1}^T w_{kt} (x_{it} - x_{jt})^2 \right]^{1/2} \quad (2)$$

と定式化する。 $w_{kt}$ は被験者k ( $k = 1, \dots, N$ )が次元t ( $t = 1, \dots, T$ )を重要視する程度を示す重み係数である。 $x_{it}$ と $x_{jt}$ は共通対象空間における、第t次元の対象iと対象jの座標である。すなわち、 $x_{it}$ とは被験者全体に共通する共通対象空間での対象iのt座標であるのに対して、特定の個人kから対象空間を眺めた場合、対象iの布置は $\sqrt{w_{kt}} x_{it}$  となって、被験者kの認知空間では座標軸tが $\sqrt{w_{kt}}$ だけINDSCALでは伸縮されることになる。INDSCALは、被験者ごとに次元tで、 $\sqrt{w_{kt}}$ だけt座標軸を引き伸ばしたり縮めたりすることで被験者間の認知個人差を表現するモデルである(図2)。ここにINDSCALの3元データ解析モデルの使い易さがあると同時に、個人差を単に軸の伸縮だけで説明するという当該モデルの持つ限界も存在する。

ところで、式(2)の形では $d_{ij,k}$ から、 $w_{kt}$ 、 $x_{it}$ 、 $x_{jt}$ の値を得ることはできない。

そこで、距離 $d_{ij,k}$ を、対象 $i$ と対象 $j$ を表すベクトルの内積 $b_{ij,k}$ にダブルセンタリング法(double centering)を用いて変換すると

$$b_{ij,k} = \sum_{t=1}^T w_{kt} x_{it} x_{jt} + e_{ij,k}$$

と式(2)は書ける。ここで、 $e_{ij,k}$ は誤差項である。この表現式にすると、 $b_{ij,k}$ が与えられたとき $T$ の値を固定すると、 $b_{ij,k}$ に最小2乗的にあてはまりのよい $w_{kt}$ ,  $x_{it}$ ,  $x_{jt}$ の値が正準分解(Canonical Decomposition Analysis)を利用して得られることが知られている。つまり

$$\sum_i \sum_j \sum_k (b_{ij,k} - \sum_{t=1}^T w_{kt} x_{it} x_{jt})^2 \rightarrow \min$$

が成立するようにNILS (Nonlinear Iterative Least Squares)を使用して計算できる。このようにINDSCALはモデル構築過程で内積を利用する。INDSCALの場合には、データに対するモデル適合度の尺度、ストレス $S$ は

$$S = \frac{1}{m} \sum_k \left( \frac{\sum_i \sum_j [(d_{ij,k})^2 - (\hat{d}_{ij,k})^2]}{\sum_i \sum_j (\hat{d}_{ij,k})^4} \right)^{1/2}$$

のように定義される。

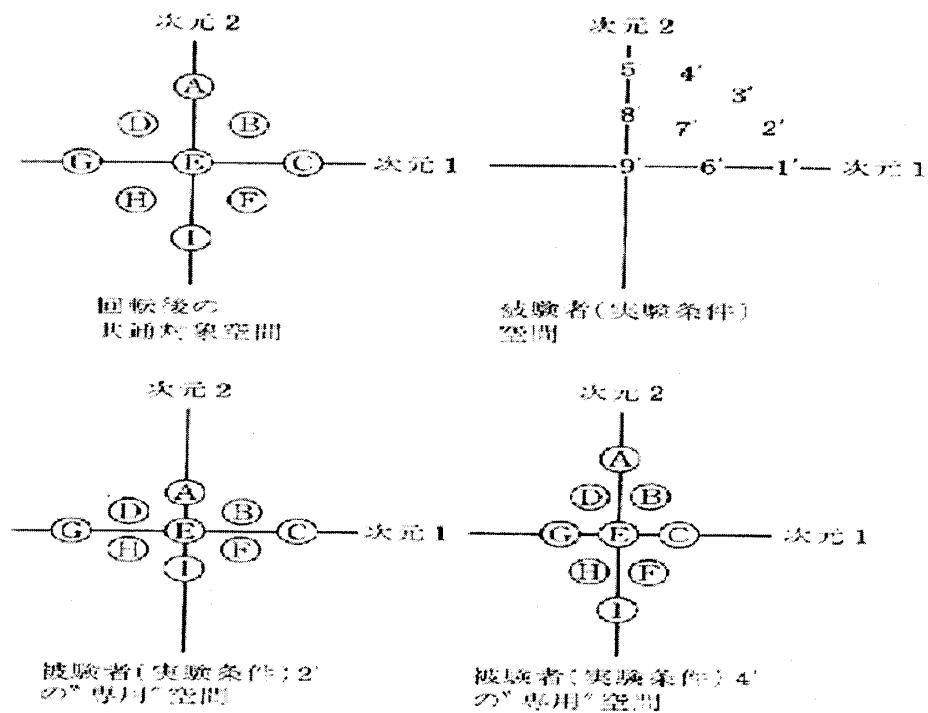


図2. INDSCALによる共通対象空間 $X$ と被験者空間 $W$  (岡太・今泉, 1990)

### 3. 分析方法

1990年度から2001年度までの12年間にわたるプロ野球チーム12球団の成績データ，すなわち，6個の対象である成績，勝率，打率，得点，失点，本塁打数，防御率に関するデータをINDSCALで解析して対象6個の布置と年度格差を示す重みを得る。

一年間の成績データでは， $12 \times 6$ 型の2相2元データ行列 $Z_i$  ( $i = 1, \dots, 12$ )が形成される。対象間の非類似性はユークリッド距離を採用して計算すると，6個の対象間の ${}_6C_2$ 個の非類似性データが計算されて単相2元データの $6 \times 6$ 型非類似性行列 $\Delta$ を生成する。12年間にわたってデータ集計したので12個の非類似性行列 $\Delta_1, \Delta_2, \dots, \Delta_{12}$ が作成される。さらに，これら12個の非類似性行列 $\Delta_k$  ( $k = 1, \dots, 12$ )を縦に並べると，2相3元データの $72 \times 6$ 型の非類似性行列 $D = [d_{ij,k}]$ が形成される。

この非類似性行列 $D$ をINDSCALで解析して，2次元の共通対象空間 $X$ と被験者空間 $W$ ，そして，モデル適合度の指標であるストレス値を得る。今回のデータ解析では，布置の次元は一様に2に固定して，データ圧縮過程における情報損失に留意する。

### 4. 結果

ストレス値が0.2を超える1992，1995，1996年度のデータは次元数2でINDSCALモデルで分析するのは不適切であると見做して3元データ解析対象から外し，残り9年間のデータのみで改めてINDSCALによる解析を実行した。このようなデータマイニング的な対処によって(福田他，2001)，平均ストレス値を0.166から0.119にまで低下させることに成功した。一年間のデータだけでは安定した対象布置は得られないことは実証されていて(奥・前鶴，2002)，9年間の成績データから計算された非類似性行列 $D$ に基づいて，対象に関する安定した布置を得たが，特にストレス値が大きいデータは当面の解析対象から除いて，これらデータは別途，他のメソッドで分析する必要がある。

INDSCALによる共通対象布置 $X_1$ (12年間)， $X_2$ (9年間)，被験者空間布置 $W$ を図3，図4および図5に示す。INDSCALの被験者空間 $W$ という出力は，プロ野球成績データの時系列的变化を捉えたものと理解できて(図5)，3元データ圧縮過程でのデータ濃縮という課題をINDSCALが首尾よく遂行していることが窺えた。12年間のデータによる布置 $X_1$ (図3)とモデル適合度が良い9年間データの布置 $X_2$ (図4)が，それ程異なっていないことが確認でき，3年間のデータ排除は総括的なデータ解析に大きな影響を与えなかったものといえる。

表1. 12年間のINDSCAL解析結果における各年度の布置のストレス値と平均ストレス値.

年度	Stress	RSQ	年度	Stress	RSQ
1990	.116	.920	1991	.099	.940
1992	<b>.279</b>	.553	1993	.141	.877
1994	.172	.838	1995	<b>.227</b>	.694
1996	<b>.246</b>	.660	1997	.137	.885
1998	.097	.957	1999	.115	.922
2000	.138	.904	2001	.093	.960

Averaged (rms) over matrices

$$\text{StrPPess} = .16614$$

表2. 9年間のINDSCAL解析結果における各年度の布置のストレス値と平均ストレス値.

年度	Stress	RSQ	年度	Stress	RSQ
1990	.100	.943	1991	.094	.957
1993	.110	.929	1994	.188	.818
1997	.130	.912	1998	.105	.953
1999	.087	.958	2000	.124	.930
2001	.099	.957			

Averaged (rms) over matrices

$$\text{Stress} = .11871$$

個人差 (重み付き) ユークリッド距離モデル

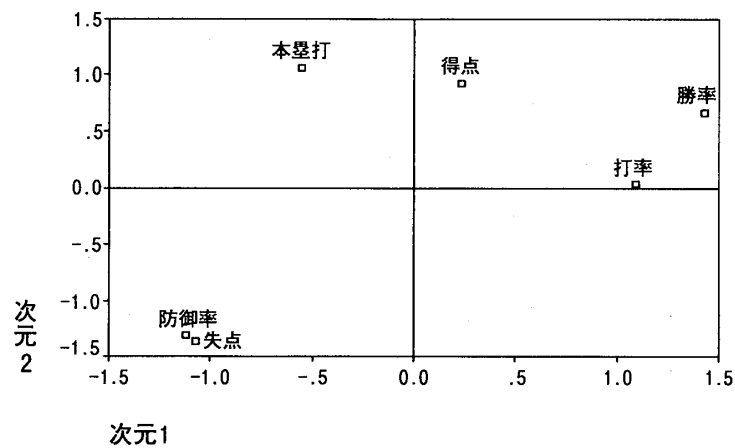


図3. INDSCALによる12年間の共通対象空間の2次元布置.

個人差（重み付き）ユークリッド距離モデル

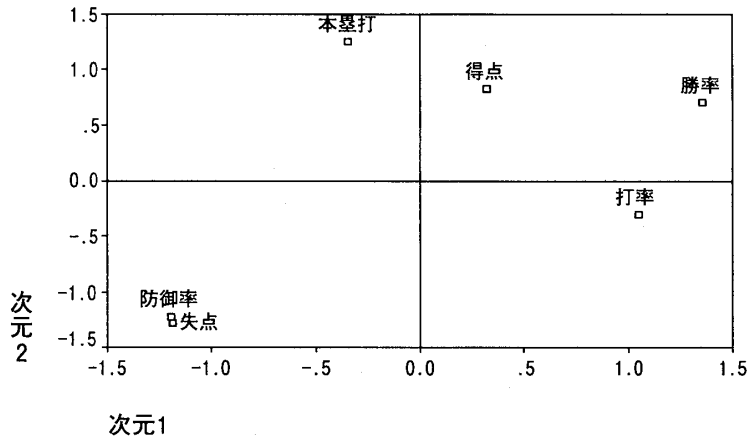


図4. INDSCALによる9年間の共通対象空間の2次元布置.

個人差（重み付き）ユークリッド距離モデル

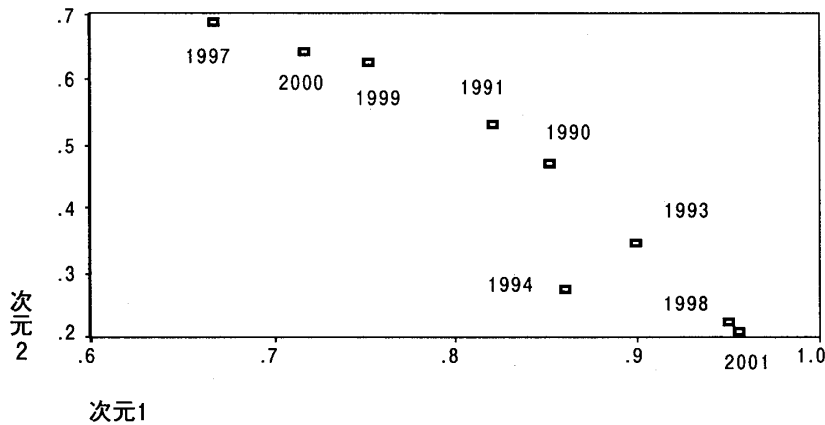


図5. INDSCALによる9年間の被験者空間 W.

5. 考察

2相3元データをクルスカルの方法で解析した場合に比べてINDSCALを利用した場合のほうが、被験者空間Wが出力されるだけデータ圧縮過程における情報損失を減らすことができる(奥, 2002)。すなわち、2相3元データのもつ経時的変動をINDSCALは被験者空間Wで捉えることができ、INDSCALが時系列的な3元データを解析するのに適切な方法であることが窺える。データ圧縮過程でデータサイエンスが目指すデータ濃縮、データ集約という作業に(柴田, 2001)、INDSCALが参画して有用な情報を形成している様子が理解できた。すなわち、9年間のデータから成る2相3元データを

INDSCALのよって解析した結果の共通対象空間  $X$  は、6 個の対象の安定した布置となっており、同時に、被験者空間出力  $W$  は年毎の布置変動を示すものと解釈できる。

INDSCALによるモデリングは年毎の布置変動を軸の重みの変化として捉えて、プロ野球成績の年度毎の特徴とみなして個人差、このケースでは年度差を被験者空間  $W$  で捉えるモデル構築の立場を採っている。これに対して成績の年度差を偶然的変動とみなす立場では、9 年間のデータを全体として総括してから、クルスカルの方法を一回利用したデータ解析を実行すればよいことになるが、この場合には成績の年度差は考慮しないから、データの経時的変動を無視されることを覚悟しなければならない。集計されたデータが経時的性質を伴う場合には、経時変化を無視せずに時系列的変化を捉えることが可能なINDSCALのような3元データ解析モデルを使用してデータ圧縮時の情報損失を少なくすべきであろう。

被験者ごとに次元  $t$  において、 $\sqrt{w_{kt}}$  だけ  $t$  座標軸を引き伸ばしたり縮めたりすることで被験者の認知空間個人差をINDSCALは表現するモデルであって、ここにINDSCALが単に軸の伸縮だけで個人差を説明するという、当該モデルのもつ限界が存在する。しかるに、特に、ストレス値が0.2より大きい年度データはINDSCALで無理に分析しないことが合理的対処であろう。これらの年度については、そのデータ構造が共通対象空間の軸の伸縮だけでは捉えられない独特の布置構造を持っていると見做すべきで、3年間のデータは別途、Kruskalの方法で解析すべきである。データマイニングの例外的データは排除すべきであるという原則に従うより(福田他, 2001), これらのデータを注意深く解析処理すべきであるという立場を筆者はとる。

そこで、1992年度、1995年度、1996年度のデータをそれぞれKruskalの方法で解析して得た布置をつぎに示す。これらのKruskalの方法による布置をINDSCALによる共通布置空間  $X$  と比べると(図7, 図8), 共通布置  $X$  からの軸の伸縮だけでは得られない独特の布置構造を、これらの布置が実際に持っていることが視覚的に理解できる。



個人差（重み付き）ユークリッド距離モデル

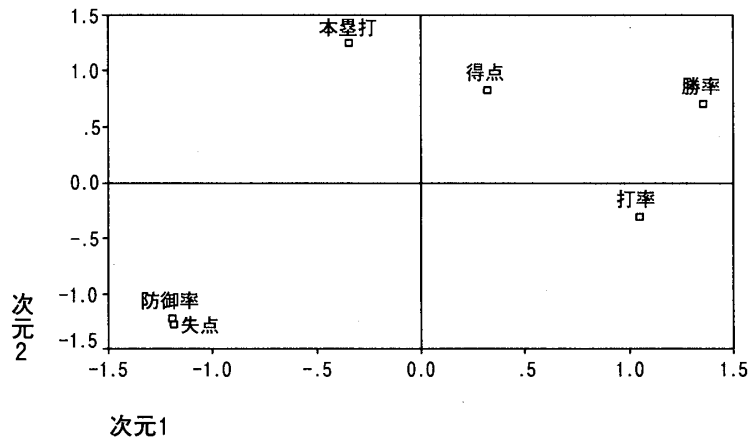


図6. INDSCALの9年間の共通布置（参考）

ユークリッド距離モデル

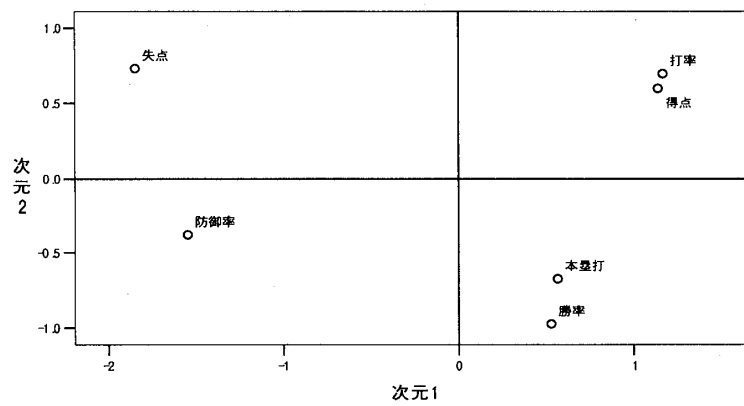


図7. 1992年度の布置 ストレス値=0.182

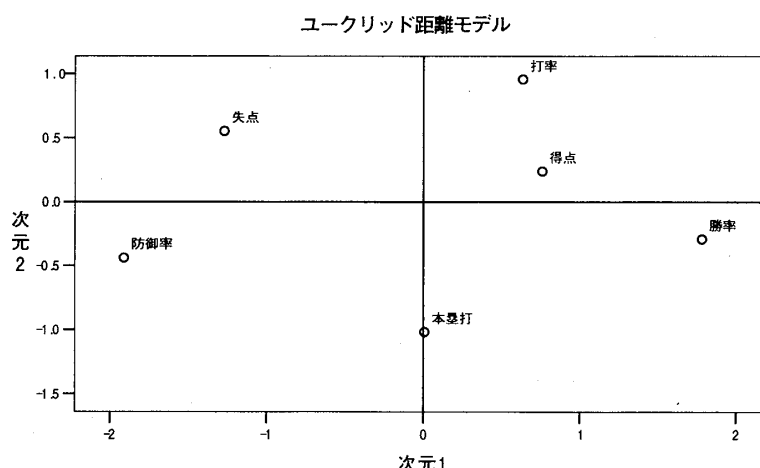


図8. 1995年度布置 ストレス値=0.108

ところで、従来の統計的データ解析メソッドをデータマイニングの枠組みで使用する際には、使用形態に工夫を加える必要が一般的にある。例えば、データマイニングでクラスター分析を実行するときは、完全にk個のクラスターに分別する必要はなく、おおよそその数のクラスターに分別すれば十分であり、例外的なデータを排除することで、その目的を達成することができる。福田らは(2001)データマイニングの実行時には例外的なデータや不要属性の排除を、データ解析メソッドを使用する場合の注意点として要求した。同様にデータマイニングで3元多次元尺度法を利用する際にも、例外的なデータは別途扱うことを本稿では実践した。結果として総括的なストレス値を0.166から0.119へ低下できたことで、効率的なデータ解析メソッドの有効利用を例証できた。

本稿は、従来はそれほど利用されなかった3元多次元尺度法INDSCALのデータマイニングへの適用が、データマイニング手法のクラスター分析や決定木と同様に有効なメソッドになり得ることを具体的なデータにINDSCALを適用して指摘した。特に当該モデルが3元データに対して安定した布置を供給するとともに、データ濃縮という事項の時系列変動を適切に捉えることを示し、また、データマイニングでのINDSCALの有効利用には例外的なデータは別途に取り扱うことが必要であることを強調した。

## 謝辞

本研究において流通経済大学のバンドン・エリ君、于永慶君には協力を頂きましたのでここに感謝します。

参考文献

- Arabie, P., Carroll, J. D. & Desarbo, W. S. (1987). Three-way Scaling and Clustering. Newbery Park, CA: Sage.
- 岡太・今泉（共訳）（1990）. 3元データの分析. 共立出版
- Carroll, J.D. and Chang, J. J. (1970). Analysis of Individual Differences in Multi- Dimensional Scaling via an N-way Generalization of Eckart-Young Decomposition. *Psychometrika*, 35, 283-319.
- Everitt, B.S. & Rabe-Hesketh, S. (1997). The Analysis of Proximity Data. Kendall's Library of Statistics 4. Arnold, London.
- Giudici, P. (2003). Applied Data Mining. Wiley, Chichester.
- Kruskal, J. B. (1964). Nonmetric Multidimensional Scaling: Numerical Method. *Psychometrika*, 29, 28-42.
- Kruskal, J. B. (1965). Analysis of Factorial Experiments by Estimating Transformations of the Data. *Journal of the Royal Statistical Society Series B*, 27, 251-263.
- 岡太彬訓, 宮内綾子 (1996). INDSCALを用いた予測の一方法：新製品の購入予測心理学評論, 30, 4, 439-458.
- 奥喜正, 前鶴政和 (2002). INDSCALのデータサイエンスへの適用と展望. 日本経営数学会第24回全国研究大会報告要旨集, 27-31.
- 柴田 里程 (2001). データリテラシー. 共立出版
- 福田 剛志, 森本 康彦, 徳山 豪 (2001). データマイニング. 共立出版