

## 計量MDSのストレスによる評価とグラフ表現の高度化

渡辺 洋平 三浦 孝夫  
法政大学工学部電気電子工学科

## 1 前書き

多次元尺度構成 (MDS) は、一対の高次元データの非類似度を収集し、この間の距離を直接に低次元空間で表現する技法である。MDSの主たるねらいは、機械学習や統計・推定とは異なって、データオブジェクト間の関連性の把握にある。しかも大づかみな傾向を捕らえることが狙いであり、多くの場合、視覚的な表現を伴う [2]。各次元は解釈可能でなければならないため、結果の評価は”心理的”であり客観的な評価は難しい。反面、実用的な観点から広範囲な応用を有する [1]。

計量MDS(Metric MDS) は、データ間の距離が数値として得られるときにその比を保ったままで表現する手法である。順位尺度データを扱う場合には非計量MDS手法があり、Kruskal手法などが知られる。非計量データは、(非)類似性を相関性などの統計値で与え、これをもとに類似度を推定するため、MDSによる結果の”よさ”を客観的に評価するストレス (Stress) 関数が必要となる。

本研究では、計量MDS (Torgerson法) を用い、次元選択とカラー値 (RGB) による6次元までの低次元化を行う。評価関数 (ここでは単純にストレスを利用) と次元縮小を利用者が判断できるシステムを試作する。次元縮小によって最大6次元までの可視化を行う。次元は明示的に利用者が与えてもよいが、自動的に次元縮小を行える。これにより、縮小次元、ストレス、表示次元の選択を入力として可視化状況を制御することができる。

## 2 多次元尺度構成とストレス関数

非計量MDSではデータ間  $i, j$  の類似度は  $\delta_{ij}$  と表現されこれが次元縮小されて距離  $d_{ij}$  をもつ2点に配置される。 $F(\delta_{ij}) = d_{ij}$  を与える変換  $F$  を多次元尺度構成MDSと定式化する。しかし実際には  $d_{ij}$  ではなくて、誤差  $e_{ij}$  を伴う。現実には、この誤差を小さくする変換  $f$  を考えることになる:  $f(\delta_{ij}) = d'_{ij} = F(\delta_{ij}) + e_{ij}$

”Stress Function on Metric MDS for Advanced Visualization”: Yohei Watanabe, Takao Miura: Hosei University, Dept. of Elec. and Elec. Eng. Kajino-cho 3-7-2, Koganei, Tokyo, JAPAN

すべての誤差の二乗和はデータと、空間上の距離の適合度を表すと考えられるのでストレス評価関数  $S$  を次のように定義する:  $S = \sqrt{\frac{\sum_{i < j} (d_{ij} - d'_{ij})^2}{\sum_{i < j} d_{ij}^2}}$

計量MDSでは、変換  $f$  の計算は単純な線形代数計算に基づくものがある。このうちTorgerson法は、すべてのデータオブジェクト間に距離  $\delta$  が与えられている (したがって対称関係) と仮定する [1]。  $n$  次実対称正定行列  $A = (a_{ij})$  を次式により求める:

$$a_{ij} = \frac{1}{2} \left( \frac{1}{n} \sum_k \delta_{ik}^2 + \frac{1}{n} \sum_k \delta_{kj}^2 - \frac{1}{n^2} \sum_k \sum_h \delta_{kh}^2 - \delta_{ij}^2 \right)$$

この要素は実は全オブジェクトの重心から  $i, j$  に至るベクトルの内積を意味し、固有値・固有ベクトルから因子分解して対応する点の座標を求めることができる。すなわち、 $A$  は直交行列  $X$  を用いて固有値からなる対角行列  $\Lambda$  に変換できる:  $X^t \cdot A \cdot X = \Lambda$

このとき固有値は大きいものから並べ、直交行列  $X = (\mathbf{x}_1 \dots \mathbf{x}_n)$  は固有ベクトル列となる。すべての固有値が非負であるとき、 $\Lambda$  の各要素  $\lambda_k$  を  $\sqrt{\lambda_k}$  に置き換えた対角行列を  $\Lambda_2$  と表せば、行列  $P = X \cdot \Lambda_2$  の第  $i$  行はデータオブジェクト  $i$  の座標を実ユークリッド空間で表現できる。

この因子分解の結果に対して縮小次元の決定が必要である。大きい固有値のみを取り出し (この個数が縮小する次元に対応する)、残りを無視することにより、次元縮小が可能となる。作為的に無視する固有値があるならば、縮小する次元に応じて”無視される”部分が增加し、類似性行列の近似的な復元が困難となる。この誤差は上述のストレス関数を用いて計算できる。

ストレスを計算するとき、直交行列 (固有ベクトル列)  $X$  で規定される座標系を用いる。本試作システムでは通常3次元表現に加え、RGBカラーを用いて最高6次元オブジェクトの可視化を行うが、視点の向きを変更する機能も有する。

本試作システムでは、固有値が正である限り表示次元 (1-6) を動的に選択しストレス最小を選ぶ機能を有する。あるいは、利用者の指定した閾値を下回るものがあればその時点で停止する。そのようなものがな

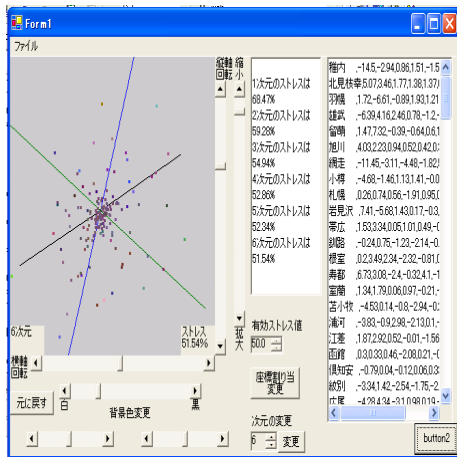
れば、エラーとする。

### 3 試作システム

本章では、前述の機能を実行する試作システムを構築し、その有効性を確認する。入力データはファイル形式で与えられ、Torgerson 法による計量 MDS 計算を行った結果を(最高 6 次元に)可視化する。表示空間上で 3 つの直線が X, Y, Z 軸を表しデータオブジェクトが配置され、これらを射影変換した結果が表示する。表示次元は自動的に選定され、正の固有値により決定される。即ち、正の固有値が 2 つならば、2 次元の座標、6 つならば 3 次元座標と RGB カラー値に変換される。これを超える場合は、固有値の大きいほうから 6 個を選び、同様の処理を実行する。

しかし、自動的に選定された次元値は、利用者が任意に変更できる。本試作システムでは、正の固有値が 3 個以下はそのまま、4 つ以上では大きいものから 3 つをそれぞれ X,Y,Z 座標に対応させ、それ以降を RGB 値に対応させる。各 RGB 値は、0-256 に比例配分しあてはめる。

ストレス基準(この評価値を上限とする)を設定することにより、計算結果がこれを満たさないときには、改善の方法をとる。例えば 4 次元でのストレス評価値が基準を満たさない場合、5 次元に設定しなおして再計算する。6 次元を超えた場合は、表示不能と判断する。次の例は 6 次元までのストレス計算をすべて実行し 51.54% の最小値になるまで繰り返している。



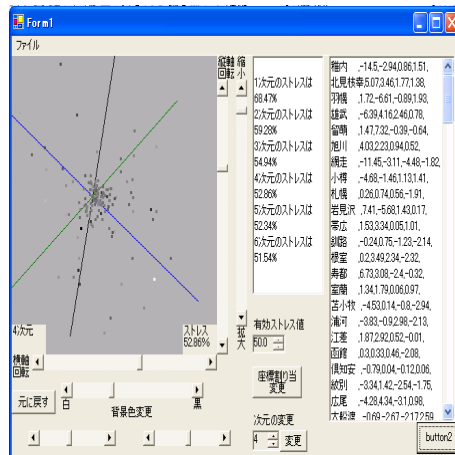
### 4 実験

以下では、試作システムを用いて 16 属性からなる全国の 6 月の気象データ 150 件を用いる。実験データ、得られたグラフを以下に示す。縮小次元は自動的に選定したものであり、ストレス基準(この評価値を上限とする)を設定していないため、6 個以上の正の固有値により 6 次元で表示する。先の例はこれを表示したものである。

である。

地点	平均気温	最高気温	最低気温	平均湿度	最小湿度	日照時間	不照日数	降水量	日降水mm	降雪
箱内	3.5,14.1,-1.5,82,30,139.7,8,109.5,28.5,24.4,2,7,1,1,166.5,49.1,4.7									
北見枝幸	2.7,12.3,-3.5,79,25,133.9,8,124.5,22.5,25.4,7,0.6,169.7,66.5,3.1									
羽幌	5.1,18.5,-2.3,73,25,144.1,7,94.5,30.5,23.5,3,9,1,3,172.2,64.2,3.4									
雄武	2.5,13.8,-7.7,78,20,138.4,8,128.34,26.3,8.8,-0.6,175.9,50.3,3.8									
留萌	4.9,14.8,-1.6,79,33,137.2,8,81.5,20.24,5.3,9.2,1.5,168.7,51.8,4.7									
旭川	4.2,15.7,-6.2,76,23,114.5,7,97.5,35.5,24.5,2,10.5,0.1,169.5,8.2,1									
網走	3.15,-5.1,79,21,145.7,8,160.53,22.4,1.8,5.0,2,174.5,54.7,3.5									
小樽	5.6,15.8,-0.3,69,20,137.1,5,81.5,12.5,24.6,3,10.5,2.5,177.4,61.7,2.7									
札幌	8.1,15.2,-1.1,73,19,139.4,4,123.5,38.24,6.7,11.1,2.7,178.4,60.9,3.3									
岩見沢	5.2,16.5,-2.7,78,26,127.7,8,103.5,32.5,23.5,6,10.5,1,175.58,3.7									
帯広	4.9,19.1,-4.4,70,16,163.1,6,260.5,89,21.5,4,11.3,0.2,197.5,60.5,2.3									
釧路	3.9,12.5,-4.8,78,29,165.1,7,186.5,49,21.3,5.7,5.0,181.1,78.8,3.8									
根室	3.1,14.4,-3.4,79,25,162.5,5,166,45,20.3,2,6.9,0.2,175.4,77.8,5.5									
帯広	5.3,12.8,-2.5,80,18,122.3,10,117.5,31.5,24.6,1.9,6.2,6.173.6,61.9,3.6									
室蘭	5.4,12.1,0.3,84,24,158.4,5,122.5,42.5,23.5,5.9,1.2,7,193.5,76.8,4.6									
苫小牧	5.2,13.4,-2.7,77,18,148.8,4,143,37.5,22.4,8,8,1,174.8,5,3.2									
浦河	5.2,15.4,-1.78,37,193.5,4,169.5,44.5,20.4,8,8,5.1,4,190,78.7,4.4									
江津	7.2,17.5,-0.1,75,21,140.5,5,53.5,9.5,23,7,10.6,3,9,169.4,77,4.3									

各データオブジェクトは観測点を示し、RGB の順で近接する。もとの数値と比べると、右側が北の地方に、左側が南に対応する。実際に計算した結果、4 次元空間を指定し表示させた例である。この結果では上で 52.86% のストレス値を得る。



### 5 結論

本稿では次元数を与え、これを満す解があるとして可視化・表示した。表示方法で、4 次元以上ではどのデータを選ぶかで印象が異なる。表示軸の設定を変更することで、(次元数やストレス値は同じであっても)可視化結果の印象が変化する。しかし、"印象の計算"は、新たな評価の観点を必要とする。

### 参考文献

- [1] 林知己夫, 飽戸弘: 多次元尺度解析法, サイエンス者, 1976
- [2] 細田博史: 高次元の利用による対話的グラフ配置法, WISS, 2003
- [3] Young, F.W. and Hamer, R.M.: Multidimensional Scaling: History, Theory and Applications, Erlbaum, 1986