

MT（マハラノビスタグチ）システム

1. MTシステムとは

MTシステムとは、多次元の情報で構成されたデータベースに対して、一つの尺度を導入し、パターン認識をする情報処理の技術である。

この多次元空間の中に、普通・一様・中間の集団で単位空間を定義し、尺度を作り、単位空間に属さない対象について、単位空間と異なる集団かどうかとか、この先どうなるかを診断・予測する。

近年注目され、事例も急増しているが、初期の方法であるMT法では多重共線性や分散ゼロについての問題があり、その対応によりMTA法・TS法・T(1)(2)(3)法と多くの対応策が提案されてきた。

今回は最初に提案されたMT法と最新のT(1)法を中心に解説する。

2. MT法

1970年代に田口先生から提案があり、通信病院の兼高先生との協同研究がスタートしたが、80年代になって、ようやく成果がでるようになった。

インドの統計学者P. C. マハラノビスが提唱したマハラノビスの汎距離を使用する方法で、最初にMTシステムとして提案されたものである。

表. 1のような単位空間データをk個集める。単位空間データは普通・一様・中間なもので十分な数があることが重要である。単位空間データの数kは項目の数nより大きい必要がある ($k > n$)。

表. 1の各項目の平均値をm、標準偏差をσとする。

表. 2のように、表. 1と同じ項目で、単位空間に属さない信号データをl個集める。信号データの数lは1以上であればよい。

信号データは一つひとつパターンが異なるが、単位空間内では同一パターンとし、その中に、パターン内距離としてマハラノビスによる複合距離（マハラノビスの汎距離）を求める。

表. 1 単位空間データ

	項目 1	項目 2	...	項目 n
1	X_{11}	X_{21}	...	X_{n1}
2	X_{12}	X_{22}	...	X_{n2}
.
.
.
k	X_{1k}	X_{2k}	...	X_{nk}

表. 2 信号空間データ

	項目 1	項目 2	...	項目 n
1	Y_{11}	Y_{21}	...	Y_{n1}
2	Y_{12}	Y_{22}	...	Y_{n2}
.
.
.
l	Y_{1l}	Y_{2l}	...	Y_{nl}

平均値 m_1 m_2 ... m_n

標準偏差 $\sigma_1 \quad \sigma_2 \quad \cdots \quad \sigma_n$

MT法ではデータの規準化を行う。単位空間データでの各項目の平均値 m と標準偏差 σ により

$$x_{ij} = \frac{X_{ij} - m_i}{\sigma_i}$$
$$y_{ij} = \frac{Y_{ij} - m_i}{\sigma_i}$$

とする。

単位空間データを用い、各項目間の相関係数 r を求め、相関行列 R を作成する。

$$r_{ij} = \frac{S(x_i, x_j)}{\sqrt{S(x_i, x_i) S(x_j, x_j)}}$$

ここで

$$S(x_i, x_i) = \sum x_{ik}^2 / k$$
$$S(x_i, x_j) = \sum x_{ik} x_{jk} / k$$

とする。

相関行列 R は

$$R = \begin{pmatrix} r_{11} & r_{21} & \cdots & r_{n1} \\ r_{12} & r_{22} & \cdots & r_{n2} \\ r_{13} & r_{23} & \cdots & r_{n3} \\ \cdot & \cdot & \cdots & \cdot \\ \cdot & \cdot & \cdots & \cdot \\ \cdot & \cdot & \cdots & \cdot \\ r_{1n} & r_{2n} & \cdots & r_{nn} \end{pmatrix}$$

である。

相関行列の逆行列 A を求める。

$$A = R^{-1} = \begin{pmatrix} a_{11} & a_{21} & \cdots & a_{n1} \\ a_{12} & a_{22} & \cdots & a_{n2} \\ a_{13} & a_{23} & \cdots & a_{n3} \\ \cdot & \cdot & \cdots & \cdot \\ \cdot & \cdot & \cdots & \cdot \\ \cdot & \cdot & \cdots & \cdot \\ a_{1n} & a_{2n} & \cdots & a_{nn} \end{pmatrix}$$

逆行列 A が求まらない場合があり、これを多重共線性という。

多重共線性については後で詳細を述べる。

逆行列Aと規準化した単位空間データよりマハラノビスの距離を次式で計算する。

$$D^2 = (x_{1i} \cdots x_{ni}) \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{pmatrix} \begin{pmatrix} x_{1i} \\ \vdots \\ x_{ni} \end{pmatrix}$$

単位空間でのサンプルに対して2乗距離を求めたことになる。

単位空間のマハラノビス距離の平均値は1となり、1とならない場合は逆行列が正常に求まっていないので結果が信用できない。

信号データについても単位空間で求めた逆行列Aと規準化したデータを用い、信号データのサンプルに対するマハラノビス距離を計算する。

$$D^2 = (y_{1i} \cdots y_{ni}) \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{pmatrix} \begin{pmatrix} y_{1i} \\ \vdots \\ y_{ni} \end{pmatrix}$$

信号データで求めたマハラノビス距離により、信号データのサンプル一つひとつが単位空間からどれくらい離れているかをみて、単位空間に属するものか否かを判別する。信号データで真値がわかる場合は真値を信号Mにし、マハラノビス距離の平方根Dを特性値としてSN比を計算する。

計測項目の中に無駄なものがあると、予測や判別の精度が落ちるので、計測項目を2水準系の直交表にわりつけ、水準1をその計測項目を用いる、水準2をその計測項目を用いないとして、マハラノビス距離を計算する。

このとき、直交表はL12、L20等の4×素数のものが望ましい。

直交表の組み合わせ数だけマハラノビス空間をつくり、SN比を計算して、利得を求めて、計測項目の取捨選択を行う。

信号データで真値がわからない場合は望大特性のSN比を計算し、直交表を利用して、計測項目の取捨選択を行う。

適用事例に健康診断や火災報知機のもの、外観検査技術の研究等に成果が報告されている。

何をしたいのか目的を明確にし、目的に合った計測項目を選び、普通・均質・正常・一様な計測項目で単位空間とするデータベースを定義することからMT法をスタートさせる。

例えば、偽札の判別をする場合、単位空間は本物の札のデータとし、いろいろな状態の本物の札（新品、使い古し、汚れたもの、折り目の付いたもの等）を用意する。偽札には多くの種類・方式があり、いくら偽札の研究をしても意味がなく、ひとつの偽札に対応できても、すぐに新しい種類・方式の偽札が作られ、対応にはきりが無い。本物の札を研究し、単位空間以外のサンプル（信号データ）に対するマハラノビス距離を計算し、その距離で本物か偽札かを判別するのがMT法である。

MT法には多くの問題点があり、取り扱うデータの吟味が必要である。

- ①単位空間の項目で標準偏差がゼロ（その項目のデータが全て同一）の場合、規準化をしようとするときゼロ除算となり、計算できない。
- ②多重共線性により計算が進まない、計算結果が信用できない。
多重共線性とは逆行列が求まらない状態をいい、下記が原因である。
 - 1) 相関係数の絶対値が1または1に近い
 - 2) 自由度以上の項目がある
 - 3) サンプル数が項目数より少ない

これらの問題点をクリアするため、色々な対応策が提案され、現在では、MT法以外に、MTA法、TS法、T(1)法(両側T法)、T(2)法(片側T法)、T(3)法(RT法)が存在する。

3. T(1)法(両側T法)

MT法では普通・正常等の単位空間はゼロ点近傍で、単位空間に属さない信号データでは正の距離となる。

MT法の単位空間は端にあったが、T(1)法は単位空間が真ん中にあり、正負の符号が付いた距離を計算する予測を主体とする方法である。

T(1)法では計測項目以外に必ず真値が必要で、真値とは予測したい項目をいう。企業の売上予測、販売数の予測、不動産価格の予測などでは、真中付近に単位空間をとるのがよく、売上高・販売数・価格等が真値で、T(1)法での予測が使える。昨年に比べて今年の売上の伸びが普通(例えば±1%)の企業を集めて単位空間として、売上が5%増加した企業、3%減少した企業などは信号とする。

表. 3のように単位空間データを収集する。

この時、真値については真ん中で均質なものとする必要がある。

表. 3 単位空間データ

	項目 1	項目 2	...	項目 n	真値
1	X_{11}	X_{21}	...	X_{n1}	XM_1
2	X_{12}	X_{22}	...	X_{n2}	XM_2
.
.
.
k	X_{1k}	X_{2k}	...	X_{nk}	XM_k

平均値 m_1 m_2 ... m_n x_m

単位空間データの規準化を項目毎の平均値、真値の平均値を引いて行う。

$$x_{ij} = X_{ij} - m_i$$

$$x_{m_i} = XM_i - x_m$$

単位空間には異なるデータ群で信号データを収集する。

信号データにも真値（予測したい項目）が必ず必要である。

表. 4 信号データ

	項目 1	項目 2	...	項目 n	真値
1	Y_{11}	Y_{21}	...	Y_{n1}	YM_1
2	Y_{12}	Y_{22}	...	Y_{n2}	YM_2
.
.
.
1	Y_{11}	Y_{21}	...	Y_{n1}	YM_1

信号データを単位空間の項目毎平均値、真値の平均値を引き規準化する。

$$y_{ij} = Y_{ij} - m_i$$

$$y_{m_i} = YM_i - x_m$$

規準化された信号空間の真値と項目 1 とで SN 比 η （真数）と比例定数 β を求める。

$$\eta = (S_\beta - Ve) / r / Ve \quad S\beta > Ve \text{ の時}$$

$$\eta = 0 \quad S\beta \leq Ve \text{ の時}$$

$$\beta = L / r$$

ここで

$$S_T = \sum y_{1j}^2$$

$$r = \sum y_{mj}^2$$

$$L = \sum y_{1j} y_{mj}$$

$$S_\beta = L^2 / r$$

$$Se = S_T - S_\beta$$

$$Ve = Se / (1 - \eta)$$

以下、項目 2 ~ n についても同様に SN 比と比例定数を求める。

この時、 i 番目の項目のみで真値を予測すると、 $y = \beta M$ より

$$y_{m_{ij}} = y_{ij} / \beta_i$$

となる。

j 番目のサンプルについて各項目のみの予測値を用い、項目毎 η の加重平均を総合予

測値 $e m_j$ とする。

$$e m_j = \frac{n_1 y_{1j} / \beta_1 + n_2 y_{2j} / \beta_2 + \dots + n_n y_{nj} / \beta_n}{n_1 + n_2 + \dots + n_n}$$

総合予測値は規準化されているので、真値の予測値は

$$E M_j = e m_j + x m$$

となる。

ここで予測の精度を表す総合SN比を求める。

信号データの真値を信号に、予測値を特性値としてゼロ点比例式のSN比を

$$\eta = 10 \log \frac{(S_\beta - V e) / r}{V e}$$

とする。

ここで

$$S_T = \sum E M_j^2$$

$$r = \sum Y M_j^2$$

$$L = Y M_1 \cdot E M_1 + Y M_2 \cdot E M_2 + \dots + Y M_1 \cdot E M_1$$

$$S_\beta = L^2 / r$$

$$S e = S_T - S_\beta$$

$$V e = S e / (1 - 1)$$

である。

総合SN比は大きければ大きい程よい。

計測項目の中に無駄なものがあると、予測の精度が落ちるので、計測項目を2水準系の直交表にわりつけ、水準1をその計測項目を用いる、水準2をその計測項目を用いないとして、総合SN比を計算する。

このとき、直交表はL 12、L 20等の4×素数のものが望ましい。

直交表の組み合わせ数だけ総合SN比を計算して、利得を求めて、計測項目の取捨選択を行う。

4. その他のMTシステム

MTシステムでは分散ゼロ対策、多重共線性対策のため、色々な提案がされ、現在に至っている。

以下、簡単に解説する。

(1) MTA法

MT法からの変更点のみを述べる。

①基準化を単位空間の各項目の平均値を引くことで行う。

②相関行列でなく分散・共分散行列から余因子行列を用い、直接、逆行列を求める。

以下、MT法と同様で、単位空間は端にある。

行列演算が繰り返し行われるので、計算に時間が掛かる。

(2) TS法

予測を主体としたもので、正負の符号付き予測ができる。

真値が必ずあること。

項目間にグラム・シュミットの直交展開を採用した。

直交展開は逐次計算であるので、計測項目が多いと時間が掛かる。

項目の順序が結果に影響する。

順序の決め方には真値との相関係数の絶対値が大きい順、項目を主成分分析にかけ、因子付加量の大きい順等がある。

T(1)法同様、単位空間は真ん中にある。

(3) T(2)法(片側T法)

T(1)法の単位空間が真ん中であるのに対し、T(2)法では端にある。

真値が必ずあり、ランク値でもよい。

各項目毎の予測値を2乗の形でいき、項目毎 η の加重平均で2乗距離を予測する。

(4) T(3)法(RT法)

画像処理を念頭に置いたもので、真値がない。

単位空間の各項目の平均値 m と単位空間での各サンプル毎の各項目値に対し、ゼロ点比例式の η と β を求める。

$$y_1 = \beta \quad y_2 = 1 / \sqrt{\eta}$$

としてMTA法により2乗距離を求める。

同様に単位空間の各項目の平均値 m と信号データでの各サンプル毎の各項目値に対し、ゼロ点比例式の η と β を求める。

$$y_1 = \beta \quad y_2 = 1 / \sqrt{\eta}$$

としてMTA法により2乗距離を求め、単位空間に属するか否かを判別する。