

カテゴリカルなデータにおける正準相関分析

—特性グループ間の相関に基く数量化—

統計数理研究所 仁 木 直 人

(1980年8月 受付)

Canonical Correlation Analysis of Categorical Data —Quantification Based on the Correlation between Two Sets of Categorical Variates—

Naoto Niki

(The Institute of Statistical Mathematics)

A method of "quantification" based on the correlation between two sets of categorical variates is described. We consider "scores" x and y of categories in the first and second sets of variates and consider a particular pair of variates as "scores" of samples given by $u = D^{(1)}x$ and $v = D^{(2)}y$, where $D^{(t)}$ ($t=1, 2$) are the response patterns of samples to the categories in the t -th set. Under the conditions $\bar{u} = \bar{v} = 0$ the correlation coefficient between u and v attains its maximum ρ_1 for $x = x_1$ and $y = y_1$ provided that ρ_1^2 is the second largest eigenvalue of the equations $F + GH + G'x = \rho^2 x$ and $H + G'F + Gy = \rho^2 y$ with the eigenvectors x_1 and y_1 , where $F = D^{(1)'}D^{(1)}$, $G = D^{(1)'}D^{(2)}$, $H = D^{(2)'}D^{(2)}$ and X^+ signifies a generalized inverse of X . The largest eigenvalue is always one but the associated eigenvectors do not satisfy the conditions $\bar{u} = \bar{v} = 0$. We can also determine the optimum multi-dimensional "scores" $x_1 \dots x_S$ and $y_1 \dots y_S$ which are the successive eigenvectors associated with the second to $(S+1)$ -th largest eigenvalues.

1. 序

カテゴリカルなデータの解析においては、2元クロス表を取扱うことが非常に多く、また数多くの解析手法が知られている。

この中で、Guttman [2] や林 [3] に始まる数量化は、「分析目的に応じて操作的に測度を定義し、その測度に基づいて各区分（カテゴリ）に数値を割りつける」という非常に柔軟な考え方をするもので、クロス表の取扱いに関しても応用範囲は広大である。

いま、多くの特性項目（カテゴリカルな変数）について全ての組合せでの2元クロス表が与えられているとしよう。そして、これらの特性項目を2群に分け、2つの群間の関係を解析することを考える。典型的な群の分け方は「説明項目群」と「基準項目群」の2群とすることであるが、とくにそのような意識がなくとも2群に分けて考える方が合理的な場合も多いであろう。たとえば、「政治意識に関する項目群」と「教育問題に関する項目群」との相関を見たいとする場合などである。

項目を数群に分けて解析する試みは林 [4] によってすでに行なわれているが、これは分割された各群間の分離を一般化相関比を測度として最大化しようとするものである。

本稿では、これと異なり、2群に分かれた各群ごとにある一次結合を作り、その一次結合間

の相関を最も高くするように各区分に与える数値を決めようとするものである。いわば、カテゴリカルなデータにおける正準相関分析である。

なおごく最近、岩坪 [9] が「正準相関係数最大化の数量化法」を発表している。これは本稿の方法と基本的に同一のものである。しかし、根本となる発想には異なる面があり、解くべき方程式も本稿の形の方が多次元化した場合や既存の数量化法との関連をみる上で理解しやすいだろう。また、Okamoto and Endo [11] がまとめを行なった「正準相関を測度とした数量化」は普通「クロス表の数量化」と呼ばれるもの（形式的に通常の正準相関分析に一致する）であって、のちに述べるように本稿の特別な場合に相当する。

第2節では、数量化の基準となる測度の定義を行なったのち、解の満たすべき固有方程式を導く。第3節では、その固有方程式の基本的な性質について考察する。次に第4節では、各区分に与える数値を多次元化する拡張について述べる。第5節では、本稿の特別な場合に相当する、数量化 II 類、「クロス表の数量化」、数量化 III 類との関連について議論する。最後に第6節では、簡単な実例を用いて具体的な解法の説明を行なう。

2. カテゴリカルな変数における正準相関分析

n 個の対象について、ある 1 群の特性に関する $m^{(1)}$ 項目の測定データ、および、別の 1 群の特性に関する $m^{(2)}$ 項目の測定データがあるものとする。各測定項目はいくつかの区分に分割されていて、各対象は必ずそのうちのどれかひとつだけに該当することとする。

すなわち、データとしては、

$$\delta_{ijk}^{(t)} = \delta_i^{(t)}(l) = \begin{cases} 1 & \text{第 } i \text{ 対象が第 } t \text{ 群の第 } j \text{ 項目第 } k \text{ 区分に該当するとき,} \\ 0 & \text{該当しないとき,} \end{cases}$$

$$\left[\begin{array}{l} t = 1, 2; i = 1, 2, \dots, n; \\ j = 1, 2, \dots, m^{(t)}; k = 1, 2, \dots, c_j^{(t)} \\ l = 1, 2, \dots, p^{(t)} \end{array} \right]$$

が与えられる。ただし、 $c_j^{(t)}$ は第 t 群の第 j 項目の区分数であり、

$$(1) \quad \sum_{k=1}^{c_j^{(t)}} \delta_{ijk}^{(t)} = 1$$

を満たす。また l は各群における区分の通し番号で、

$$l = k \quad (j = 1), \quad l = k + \sum_{h=1}^{j-1} c_h^{(t)} \quad (j > 1),$$

$$p^{(t)} = \sum_{h=1}^{m^{(t)}} c_h^{(t)} \quad (\text{各群での総区分数})$$

である。

ここで、2 群の特性の間になんらかの関係が想定される場合を考えてみよう。

いま、第 1 群の第 j 項目第 k 区分（通し番号で第 l 区分）に数値 $x_{jk} = x(l)$ を与え、第 2 群の第 j 項目第 k 区分（通し番号第 l 区分）に $y_{jk} = y(l)$ を与えることにする。そして、第 i 対象に対し、その該当する区分に与えた数値の合計をそれぞれ u_i, v_i すなわち、

$$u_i = \sum_{l=1}^{p^{(1)}} \delta_i^{(1)}(l) x(l),$$

$$v_i = \sum_{l=1}^{p^{(2)}} \delta_i^{(2)}(l) y(l)$$

としたとき、この u_i と v_i の相関係数を測度として、両者の相関が最も高くなるように、逆に $x(l), y(l)$ の値をもとめようとする。

これは、カテゴリカルな変数における正準相関分析とみることができよう。

計算の便利のため、各変数をマトリクス表示しておこう。

$x(l) (l=1, 2, \dots, p^{(1)})$, $y(l) (l=1, 2, \dots, p^{(2)})$; $u_i, v_i (i=1, 2, \dots, n)$ を要素とする縦ベクトルをそれぞれ $\mathbf{x}, \mathbf{y}, \mathbf{u}, \mathbf{v}$ で表わし、 $\delta_i^{(t)}(l) (t=1, 2; i=1, 2, \dots, n; l=1, 2, \dots, p^{(t)})$ を要素とする n 行 $p^{(t)}$ 列行列を $\mathbf{D}^{(t)}$ で表わす。明らかに

$$\mathbf{u} = \mathbf{D}^{(1)} \mathbf{x}$$

$$\mathbf{v} = \mathbf{D}^{(2)} \mathbf{y}$$

である。

いま、 \mathbf{x}, \mathbf{y} は u_i, v_i の平均がともに 0、すなわち、要素が全て 1 の n 次の縦ベクトルを $\mathbf{1}_n$ と書くことにすれば、 $\mathbf{1}_n' = (1 \ 1 \ \dots \ 1)$ で ($'$ は転置を表わす)、

$$(2) \quad n\bar{u} = \mathbf{1}_n' \mathbf{u} = \mathbf{1}_n' \mathbf{D}^{(1)} \mathbf{x} = 0$$

$$n\bar{v} = \mathbf{1}_n' \mathbf{v} = \mathbf{1}_n' \mathbf{D}^{(2)} \mathbf{y} = 0$$

となるように選ぶことにする。

すると、 u_i および v_i の分散 $V^{(1)}, V^{(2)}$ および共分散 C は、

$$\mathbf{F} = \mathbf{D}^{(1)'} \mathbf{D}^{(1)}$$

$$\mathbf{H} = \mathbf{D}^{(2)'} \mathbf{D}^{(2)}$$

$$\mathbf{G} = \mathbf{D}^{(1)'} \mathbf{D}^{(2)}$$

とおけば、

$$V^{(1)} = \frac{1}{n} \mathbf{u}' \mathbf{u} = \frac{1}{n} \mathbf{x}' \mathbf{F} \mathbf{x}$$

$$V^{(2)} = \frac{1}{n} \mathbf{v}' \mathbf{v} = \frac{1}{n} \mathbf{y}' \mathbf{H} \mathbf{y}$$

$$C = \frac{1}{n} \mathbf{u}' \mathbf{v} = \frac{1}{n} \mathbf{x}' \mathbf{G} \mathbf{y}$$

よって、 u_i と v_i の相関係数 ρ は

$$\rho = \frac{\frac{1}{n} \mathbf{x}' \mathbf{G} \mathbf{y}}{\sqrt{\frac{1}{n} \mathbf{x}' \mathbf{F} \mathbf{x} \cdot \frac{1}{n} \mathbf{y}' \mathbf{H} \mathbf{y}}}$$

と表わされる。なお、 \mathbf{F} は第 1 群の特性相互の、 \mathbf{H} は第 2 群の特性相互の、また \mathbf{G} は第 1 群と第 2 群の特性間のクロス表になっていることを注意しておく。

さて、 u_i, v_i のスケールは自由にとることができるが、それぞれの分散が

$$V^{(1)} = V^{(2)} = 1$$

となるようにするのが便利であろう。すなわち、

$$(3) \quad \begin{aligned} \frac{1}{n} \mathbf{x}' \mathbf{F} \mathbf{x} &= 1, \\ \frac{1}{n} \mathbf{y}' \mathbf{H} \mathbf{y} &= 1 \end{aligned}$$

なる条件のもとで、

$$\rho = \frac{1}{n} \mathbf{x}' \mathbf{G} \mathbf{y}$$

を最大とするような \mathbf{x} , \mathbf{y} を決めよう。(2)の条件についてはのちに考慮する。

ラグランジェの未定数法を適用する。未定数 λ, μ を導入して、

$$Q = \frac{1}{n} \mathbf{x}' \mathbf{G} \mathbf{y} - \frac{1}{2} \lambda \left(\frac{1}{n} \mathbf{x}' \mathbf{F} \mathbf{x} - 1 \right) - \frac{1}{2} \mu \left(\frac{1}{n} \mathbf{y}' \mathbf{H} \mathbf{y} - 1 \right)$$

とし、この Q を \mathbf{x}, \mathbf{y} の各要素で偏微分したものをマトリクス表示すれば、形式的に

$$\begin{aligned} \frac{\partial Q}{\partial \mathbf{x}} &= \frac{1}{n} \mathbf{G} \mathbf{y} - \frac{1}{n} \lambda \mathbf{F} \mathbf{x} \\ \frac{\partial Q}{\partial \mathbf{y}} &= \frac{1}{n} \mathbf{G}' \mathbf{x} - \frac{1}{n} \mu \mathbf{H} \mathbf{y} \end{aligned}$$

と書ける。ゆえに、

$$\begin{aligned} \mathbf{G} \mathbf{y} &= \lambda \mathbf{F} \mathbf{x} \\ \mathbf{G}' \mathbf{x} &= \mu \mathbf{H} \mathbf{y} \end{aligned}$$

を解けば求める解が得られる。

これらの式にそれぞれ \mathbf{x}' および \mathbf{y}' を左からかければ、

$$\begin{aligned} \mathbf{x}' \mathbf{G} \mathbf{y} &= \lambda \mathbf{x}' \mathbf{F} \mathbf{x} = n\lambda \\ \mathbf{y}' \mathbf{G}' \mathbf{x} &= \mu \mathbf{y}' \mathbf{H} \mathbf{y} = n\mu \end{aligned}$$

となり、未定数 λ, μ はともに u_i と v_i の相関係数

$$\rho = \frac{1}{n} \mathbf{x}' \mathbf{G} \mathbf{y} = \frac{1}{n} \mathbf{y}' \mathbf{G}' \mathbf{x}$$

に等しいことがわかる。よって

$$(4) \quad \begin{aligned} \mathbf{G} \mathbf{y} &= \rho \mathbf{F} \mathbf{x} \\ \mathbf{G}' \mathbf{x} &= \rho \mathbf{H} \mathbf{y} \end{aligned}$$

を解くことになる。

いま議論を簡単にするため、2つの対称行列 \mathbf{F} および \mathbf{H} の一般化逆行列のうち、各群の第1項目を除き、他の項目の第1区分に対応する行および列の要素が全て0であるものを、それぞれ \mathbf{F}^+ , \mathbf{H}^+ と書くことにする。一般にこれらは一意に定まり、

$$\begin{aligned} \text{rank}(\mathbf{F}^+) &= \text{rank}(\mathbf{F}) \\ \text{rank}(\mathbf{H}^+) &= \text{rank}(\mathbf{H}) \end{aligned}$$

を満たす反射型一般化逆行列である。

このような \mathbf{F}^+ , \mathbf{H}^+ を選んだことは、

$$\begin{aligned} x_{j1} &= 0 & (j = 2, 3, \dots, m^{(1)}) \\ y_{j1} &= 0 & (j = 2, 3, \dots, m^{(2)}) \end{aligned}$$

なる条件をつけることにより, (4) の特解が

$$(5) \quad \begin{aligned} \mathbf{F} + \mathbf{G}\mathbf{y} &= \rho\mathbf{x} \\ \mathbf{H} + \mathbf{G}'\mathbf{x} &= \rho\mathbf{y} \end{aligned}$$

により一意に定まることを意味する.

(5) を互いに代入しあうことにより, 結局

$$\begin{aligned} \mathbf{F} + \mathbf{G}\mathbf{H} + \mathbf{G}'\mathbf{x} &= \rho^2\mathbf{x} \\ \mathbf{H} + \mathbf{G}'\mathbf{F} + \mathbf{G}\mathbf{y} &= \rho^2\mathbf{y} \end{aligned}$$

なる固有値問題に帰着する.

3. 固有値問題の基本的性質

$$\begin{aligned} \mathbf{A} &:= \mathbf{F} + \mathbf{G}\mathbf{H} + \mathbf{G}' \\ \mathbf{B} &:= \mathbf{H} + \mathbf{G}'\mathbf{F} + \mathbf{G} \end{aligned}$$

と書くことにする. \mathbf{A} と \mathbf{B} の固有値は共通であり, 一般に $\text{rank}(\mathbf{G}) = \min(p^{(1)} - m^{(1)} + 1, p^{(2)} - m^{(2)} + 1)$ 個の 0 でない値を持つ.

このうち, $\rho^2 = 1$ なる固有値が含まれていることを述べる.

いま, 始めから $c_1^{(1)}, c_1^{(2)}$ 個の要素が 1 で他の全要素が 0 である $p^{(1)}, p^{(2)}$ 次の縦ベクトルをそれぞれ $\mathbf{x}_*, \mathbf{y}_*$ とする. すなわち,

$$\mathbf{x}_* = \left(\begin{array}{c} 1 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{array} \right) \left. \vphantom{\begin{array}{c} 1 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{array}} \right\} c_1^{(1)}, \quad \mathbf{y}_* = \left(\begin{array}{c} 1 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{array} \right) \left. \vphantom{\begin{array}{c} 1 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{array}} \right\} c_1^{(2)}$$

である.

(1) より明らかに

$$\begin{aligned} \mathbf{F}\mathbf{x}_* &= n\mathbf{1}_{p^{(1)}} \\ \mathbf{G}'\mathbf{x}_* &= n\mathbf{1}_{p^{(2)}} \\ \mathbf{G}\mathbf{y}_* &= n\mathbf{1}_{p^{(1)}} \\ \mathbf{H}\mathbf{y}_* &= n\mathbf{1}_{p^{(2)}} \end{aligned}$$

であるから,

$$(6) \quad \begin{aligned} \mathbf{A}\mathbf{x} &= \rho^2\mathbf{x} \\ \mathbf{B}\mathbf{y} &= \rho^2\mathbf{y} \end{aligned}$$

は, $\rho^2 = 1, \mathbf{x} = \mathbf{x}_*, \mathbf{y} = \mathbf{y}_*$ なる解を持つ.

しかし, この解は (3) の条件

$$\frac{1}{n} \mathbf{x}_* \mathbf{F} \mathbf{x}_* = 1$$

$$\frac{1}{n} \mathbf{y}_* \mathbf{H} \mathbf{y}_* = 1$$

を満たすが、(2)の条件は

$$\begin{aligned} \mathbf{1}_n' \mathbf{D}^{(1)} \mathbf{x}_* &= n \neq 0 \\ \mathbf{1}_n' \mathbf{D}^{(2)} \mathbf{y}_* &= n \neq 0 \end{aligned}$$

より満たすことができない。すなわち、 $\rho^2=1$, $\mathbf{x}=\mathbf{x}_*$, $\mathbf{y}=\mathbf{y}_*$ は排除されるべき形式解である。

次に $\rho^2 \neq 1$ なる固有値に対応する固有ベクトル \mathbf{x} , \mathbf{y} が条件 (2) を満たすことを述べる。

$$\begin{aligned} \mathbf{1}_{p^{(1)}}' \mathbf{F} \mathbf{x} &= \mathbf{1}_{p^{(1)}}' \mathbf{D}^{(1)'} \mathbf{D}^{(1)} \mathbf{x} \\ &= m_1 \mathbf{1}_n' \mathbf{D}^{(1)} \mathbf{x} \\ &= m_1 n \bar{u} \\ \mathbf{1}_{p^{(1)}}' \mathbf{G} \mathbf{y} &= m_1 \mathbf{1}_n' \mathbf{D}^{(2)} \mathbf{y} \\ &= m_1 n \bar{v} \\ \mathbf{1}_{p^{(2)}}' \mathbf{G}' \mathbf{x} &= m_2 \mathbf{1}_n' \mathbf{D}^{(1)} \mathbf{x} \\ &= m_2 n \bar{u} \\ \mathbf{1}_{p^{(2)}}' \mathbf{H} \mathbf{y} &= m_2 \mathbf{1}_n' \mathbf{D}^{(2)} \mathbf{y} \\ &= m_2 n \bar{v} \end{aligned}$$

であるから、(4)の両辺にそれぞれ $\mathbf{1}_{p^{(1)'}}$ および $\mathbf{1}_{p^{(2)'}}$ を左からかけることにより、

$$\begin{aligned} m_1 n \bar{v} &= \rho m_1 n \bar{u} \\ m_2 n \bar{u} &= \rho m_2 n \bar{v} \end{aligned}$$

を得る。ゆえに

$$\begin{aligned} (1-\rho^2) \bar{u} &= 0 \\ (1-\rho^2) \bar{v} &= 0 \end{aligned}$$

より、 $\rho^2 \neq 1$ ならば $\bar{u}=0$ かつ $\bar{v}=0$ となって条件 (2) を満たす。

したがって \mathbf{A} , \mathbf{B} の最大固有値 1 に次ぐ、2 番目に大きい固有値 ρ_1^2 を求めれば、それが u_i , v_i の相関係数の最大値 (正準相関係数に相当する) を与え、また対応する固有ベクトル \mathbf{x}_1 , \mathbf{y}_1 はそれぞれ第 1 群、第 2 群の各区分に付する最適な数値を与える。

実際に (6) を解く際には、 \mathbf{x} , \mathbf{y} の 0 でない要素数 $(p^{(1)}-m^{(1)}+1)$, $(p^{(2)}-m^{(2)}+1)$ の小さい方について先に解を求め、次いで (5) によりもう一方を求めるという手続きに従う。

なお、区分に与えた数値を各項目ごとの平均が 0 となるように修正しておくことと結果が見やすい。それには、各項目ごとにある値 ξ_j ($j=1, 2, \dots, m^{(1)}$), η_j ($j=1, 2, \dots, m^{(2)}$) を各区分に与えた数値 x_{jk} , y_{jk} に加え、

$$\begin{aligned} \sum_{i=1}^n \sum_{k=1}^{c_j^{(1)}} \delta_{ijk}^{(1)} (x_{jk} + \xi_j) &= 0 \\ \sum_{i=1}^n \sum_{k=1}^{c_j^{(2)}} \delta_{ijk}^{(2)} (y_{jk} + \eta_j) &= 0 \end{aligned}$$

が満足されれば、 x_{jk} , y_{jk} の代わりに $x_{jk} + \xi_j$, $y_{jk} + \eta_j$ を採用することにする。このような修正が可能なのは、 ξ_j , η_j を $c_j^{(1)}$, $c_j^{(2)}$ 個ずつ順に並べた $p^{(1)}$, $p^{(2)}$ 次の縦ベクトル $\boldsymbol{\xi}$, $\boldsymbol{\eta}$ が、

$$\begin{aligned} \mathbf{F} \mathbf{x} &= \mathbf{0}_{p^{(1)}} = \mathbf{0} \cdot \mathbf{1}_{p^{(1)}} \\ \mathbf{H} \mathbf{y} &= \mathbf{0}_{p^{(2)}} = \mathbf{0} \cdot \mathbf{1}_{p^{(2)}} \end{aligned}$$

の解であることから明らかで、 \mathbf{F} , \mathbf{H} の一般化逆行列を適当に選び直すことに対応する。

4. 多次元への拡張

次に各区分に多次元の数値を付すことを考えよう。

第 1 群および第 2 群の各区分に付す S 次元の数値をベクトルとして、

$$\begin{aligned} & {}^1\mathbf{x}, {}^2\mathbf{x}, \dots, {}^S\mathbf{x} \\ & {}^1\mathbf{y}, {}^2\mathbf{y}, \dots, {}^S\mathbf{y} \end{aligned} \quad (S \leq S_0 = \text{rank}(\mathbf{G}) - 1)$$

で表わす。またそれらによって第 i 対象に与えられる数値 ${}^s u_i, {}^s v_i$ ($s=1, 2, \dots, S; i=1, 2, \dots, n$) をベクトル化して、

$$\begin{aligned} {}^s\mathbf{u} &= \mathbf{D}^{(1)} {}^s\mathbf{x} \\ {}^s\mathbf{v} &= \mathbf{D}^{(2)} {}^s\mathbf{y} \end{aligned} \quad (s = 1, 2, \dots, S)$$

で表わす。

ここで、1次元の場合と同じく、 ${}^s\mathbf{u}, {}^s\mathbf{v}$ の平均に

$$(7) \quad \begin{aligned} n {}^s\bar{u} &= \mathbf{1}_n' {}^s\mathbf{u} = 0 \\ n {}^s\bar{v} &= \mathbf{1}_n' {}^s\mathbf{v} = 0 \end{aligned}$$

なる条件を、さらに、分散および共分散について

$$(8) \quad \begin{aligned} \frac{1}{n} \mathbf{r} {}^r\mathbf{u}' {}^s\mathbf{u} &= \frac{1}{n} \mathbf{r} \mathbf{x}' \mathbf{F} {}^s\mathbf{x} = \begin{cases} 1 & (r = s) \\ 0 & (r \neq s) \end{cases} \\ \frac{1}{n} \mathbf{r} {}^r\mathbf{v}' {}^s\mathbf{v} &= \frac{1}{n} \mathbf{r} \mathbf{y}' \mathbf{H} {}^s\mathbf{y} = \begin{cases} 1 & (r = s) \\ 0 & (r \neq s) \end{cases} \\ \frac{1}{n} \mathbf{r} {}^r\mathbf{u}' {}^s\mathbf{v} &= \frac{1}{n} \mathbf{r} \mathbf{x}' \mathbf{G} {}^s\mathbf{y} = 0 \quad (r \neq s) \end{aligned}$$

なる条件を課すことにする。すなわち、 ${}^s\mathbf{u}, {}^s\mathbf{v}$ は次元間で互いに直交する。

このとき、各対象に与えた 2 つの S 次元ベクトル

$$\begin{pmatrix} {}^1u_i & {}^2u_i & \dots & {}^Su_i \\ {}^1v_i & {}^2v_i & \dots & {}^Sv_i \end{pmatrix}$$

のベクトル相関を考える。すなわち、 ${}^s\mathbf{u}, {}^s\mathbf{v}$ ($s=1, 2, \dots, S$) を列ベクトルとする $n \times S$ 行列をそれぞれ \mathbf{U}, \mathbf{V} と書き、回帰モデル

$$(9) \quad \mathbf{V} = \mathbf{U}\mathbf{C} + \mathbf{E}$$

を想定する。誤差の一般化分散

$$V_\epsilon = \left| \frac{1}{n} \mathbf{E}' \mathbf{E} \right|$$

を最小化するような係数行列は、(7) (8) より

$$\begin{aligned} \hat{\mathbf{C}} &= (\mathbf{U}'\mathbf{U})^{-1} (\mathbf{U}'\mathbf{V}) \\ &= \begin{pmatrix} {}^1\rho & & 0 \\ & {}^2\rho & \\ 0 & & \ddots \\ & & & {}^S\rho \end{pmatrix} \quad \left(\text{ただし, } {}^s\rho = \frac{1}{n} \mathbf{x}' {}^s\mathbf{G} {}^s\mathbf{y} \right) \end{aligned}$$

となり、また、 $\hat{\mathbf{C}}$ を (9) に与えて最小となった V_ϵ は

$$\begin{aligned} \min V_\epsilon &= \left| \frac{1}{n} (\mathbf{V} - \mathbf{U}\hat{\mathbf{C}})' (\mathbf{V} - \mathbf{U}\hat{\mathbf{C}}) \right| \\ &= |\mathbf{I} - \hat{\mathbf{C}}^2| \quad (\mathbf{I} \text{ は単位行列}) \\ &= \prod_{s=1}^S (1 - \rho_s^2) \end{aligned}$$

と表わすことができる。よって、 \mathbf{U} と \mathbf{V} の相関係数 ${}^S R$ としては、

$$\begin{aligned} {}^S R^2 &= 1 - (\min V_\epsilon) \left/ \left(\frac{1}{n} \mathbf{V}'\mathbf{V} \right) \right. \\ &= 1 - \prod_{s=1}^S (1 - \rho_s^2) \end{aligned}$$

で定義するのが自然である。

この ${}^S R$ あるいは ${}^S R^2$ を最大化するような ${}^s \mathbf{x}$, ${}^s \mathbf{y}$ は、(6) において \mathbf{A} , \mathbf{B} の 1 を除く固有値のうち大きいものから S 個とって

$$\rho_1^2, \rho_2^2, \dots, \rho_S^2$$

とし、それらに対応する固有ベクトルを

$$\begin{aligned} \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_S \\ \mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_S \end{aligned}$$

としたとき、

$$\begin{aligned} {}^s \mathbf{x} &= \mathbf{x}_s, \\ {}^s \mathbf{y} &= \mathbf{y}_s, \end{aligned} \quad (s = 1, 2, \dots, S)$$

とすればよいことはただちにわかる。このとき ${}^s \rho^2 = \rho_s^2$ ($s = 1, 2, \dots, S$) であって、

$${}^S R^2 = 1 - \prod_{s=1}^S (1 - \rho_s^2)$$

である。

なお、 ${}^S R^2$ の代わりの測度として、Wilks [13] の A ,

$$\begin{aligned} A &= 1 - {}^S R^2 \\ &= \prod_{s=1}^S (1 - \rho_s^2) \end{aligned} \quad (\text{最小化})$$

や Hotelling [6] の一般化 T_0^2 ,

$$T_0^2 = \sum_{s=1}^S \frac{{}^s \rho^2}{1 - {}^s \rho^2} \quad (\text{最大化})$$

Pillai [12] の V ,

$$V = \sum_{s=1}^S {}^s \rho^2 \quad (\text{最大化})$$

などを用いても同じ解を得る。

5. 特別な場合

ここでは、既存の数量化手法と本稿の方法との関連について主に議論する。

a. 数量化 II 類

いま、2群に分かれた項目群の一方が、ただひとつの項目から成る場合を考える。すなわち $m^{(2)}=1$ とする。(もちろん $m^{(1)}=1$ としても全く同じ議論ができる。)

このとき H は対角要素以外は全て 0 である対角行列となり、逆行列 H^{-1} をもつ。よって y について解くべき式は、

$$H^{-1}G'F + Gy = \rho^2 y$$

となるが、これは数量化 II 類の場合と全く一致する ([7], [10])。

小さな違いとしては、固有値が相関係数の 2 乗の代わりに相関比と意味づけされていること x や y に与えるスケールが (3) に従わないこと、があげられるがともに本質的なものではない

また岩坪 [8] の「重相関係数による数量化」は $G'F + G$ のひとつを具体的に与える方法を述べたものと解される。

b. クロス表の数量化 ([1], [11] など)

第 1 群、第 2 群のどちらもがただひとつの項目から成る $m^{(1)}=m^{(2)}=1$ の場合を考える。もちろん数量化 II 類の特別な場合でもある。

このとき、 x および y の満たす固有方程式は

$$\begin{aligned} F^{-1}GH^{-1}G'x &= \rho^2 x \\ H^{-1}G'F^{-1}Gy &= \rho^2 y \end{aligned}$$

であって、いわゆる「クロス表の数量化」の式に一致する。

c. 数量化 III 類

数量化 III 類 [3] は、形式的には G の要素が 0 または 1 のみから成る「クロス表の数量化」と考えることができる。しかし、問題の設定やデータの与え方には独特なものがあり、一概にこの系列に入れてしまうのは少し無理があるかもしれない。

一方、数量化 III 類の適用例には、各測定対象についての数種の特性項目からなる第 1 群と「測定対象」という単一項目からなる第 2 群との間の相関を扱っているとみられるものが多いことも事実である。この場合、第 1 群の特性項目間の直接の相関は無視されていることに注意を払う必要がある。

6. 簡単な計算例

表 1 はある本 [5] の 11 人の共著者のプロフィールを彼らの神秘感に関する設問に対する解答パターンとして示したものである。いずれも、肯定的な解答を 1、否定的な解答を 2 で表現してある。

各解答 (すなわち区分) 間のクロス集計を行なうと、次の表 2 を得る。

ここで「信仰」「神」を第 1 群の項目、「UFO」「ネッシー」を第 2 群の項目とし、両群の相関関係を最も高くするように数量化を行なってみよう。すなわち、この 11 人の共著者に関して、第 1 群の「宗教観」と第 2 群の「科学的側面を持つ神秘感」との関わりを知ろうとする。

表1 共著者のプロフィール

	宗教を信仰	神の存在	UFOの存在	ネッシーの存在
1	1	1	1	2
2	1	1	2	1
3	2	1	1	1
4	2	1	1	1
5	2	1	1	2
6	2	1	2	1
7	2	2	1	1
8	2	2	1	1
9	2	2	2	1
10	2	2	2	1
11	2	2	2	2

表2 クロス表

		信 仰		神		UFO		ネッシー	
		1	2	1	2	1	2	1	2
信 仰	1	2	0	2	0	1	1	1	1
	2	0	9	4	5	5	4	7	2
神	1	2	4	6	0	4	2	4	2
	2	0	5	0	5	2	3	4	1
UFO	1	1	5	4	2	6	0	4	2
	2	1	4	2	3	0	5	4	1
ネッシー	1	1	7	4	4	4	4	8	0
	2	1	2	2	1	2	1	0	3

本稿の記述法に従い,

$$\mathbf{F} = \begin{pmatrix} 2 & 0 & 2 & 0 \\ 0 & 9 & 4 & 5 \\ 2 & 4 & 6 & 0 \\ 0 & 5 & 0 & 5 \end{pmatrix}, \quad \mathbf{H} = \begin{pmatrix} 6 & 0 & 4 & 2 \\ 0 & 5 & 4 & 1 \\ 4 & 4 & 8 & 0 \\ 2 & 1 & 0 & 3 \end{pmatrix}, \quad \mathbf{G} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 5 & 4 & 7 & 2 \\ 4 & 2 & 4 & 2 \\ 2 & 3 & 4 & 1 \end{pmatrix}$$

とする. \mathbf{F} および \mathbf{H} の一般化逆行列 \mathbf{F}^+ , \mathbf{H}^+ としては, どちらも第2行の要素が全て0であるものを選べば,

$$\mathbf{F}^+ \mathbf{G} = \begin{pmatrix} -1/4 & 1/4 & -1/4 & 1/4 \\ 0 & 0 & 0 & 0 \\ 3/4 & 1/4 & 3/4 & 1/4 \\ 2/5 & 3/5 & 4/5 & 1/5 \end{pmatrix},$$

$$\mathbf{H}^+ \mathbf{G}' = \begin{pmatrix} -1/16 & 1/16 & 1/4 & -1/4 \\ 0 & 0 & 0 & 0 \\ 5/32 & 27/32 & 3/8 & 5/8 \\ 3/8 & 5/8 & 1/2 & 1/2 \end{pmatrix}$$

となるから、

$$\begin{aligned} \mathbf{A} &= \mathbf{F} + \mathbf{GH} + \mathbf{G}' \\ &= \begin{pmatrix} 9/128 & -9/128 & -1/32 & 1/32 \\ 0 & 0 & 0 & 0 \\ 21/128 & 107/128 & 19/32 & 13/32 \\ 7/40 & 33/40 & 1/2 & 1/2 \end{pmatrix} \end{aligned}$$

が得られる。

\mathbf{A} が固有値 1 を持つことは

$$\mathbf{A} \begin{pmatrix} 0 \\ 0 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 1 \end{pmatrix}$$

より確かめられる。

さて、

$$(10) \quad \mathbf{Ax} = \rho^2 \mathbf{x}$$

なる固有値問題を解くことになるが、 $x_{12} = x(2) = 0$ は明らかであるから、 \mathbf{A} の第 2 行および第 2 列を取り除いた行列を \mathbf{A}_1 として、

$$\mathbf{A}_1 \mathbf{z} = \rho^2 \mathbf{z}$$

なる固有値問題を解けば充分である。さらに、この式は

$$\rho^2 = 1, \quad \mathbf{z}_* = (0 \ 1 \ 1)$$

なる形式解を持っていることから、単位行列の第 3 列を \mathbf{z}_* で置きかえた行列を \mathbf{T} とし、 \mathbf{T} による \mathbf{A}_1 の相似変換

$$\mathbf{A}_2 = \mathbf{T}^{-1} \mathbf{A}_1 \mathbf{T}$$

を考える。ここに

$$\mathbf{T} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{pmatrix}, \quad \mathbf{T}^{-1} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & -1 \\ 0 & 0 & 1 \end{pmatrix}$$

である。

すると、 \mathbf{A}_2 の第 3 列は $(0 \ 0 \ 1)'$ となり、

$$\mathbf{A}_2 = \begin{pmatrix} 9/128 & -1/32 & 0 \\ 7/640 & 3/32 & 0 \\ 7/40 & 1/2 & 1 \end{pmatrix}$$

である。ゆえに \mathbf{A}_2 の第 3 行および第 3 列をとり除いた行列を

$$\mathbf{A}_3 = \begin{pmatrix} 9/128 & -1/32 \\ 7/640 & 3/32 \end{pmatrix}$$

とおけば、(10) の固有値問題は $\rho^2 = 1$ の解を含まない

$$(11) \quad \mathbf{A}_3 \mathbf{w} = \rho^2 \mathbf{w}$$

まで簡単化される. $\mathbf{x}, \mathbf{z}, \mathbf{w}$ の関係を整理しておけば,

$$\begin{aligned}\mathbf{x}' &= (x(1) \ 0 \ x(3) \ x(4)) \\ \mathbf{z}' &= (x(1) \ x(3) \ x(4)) \\ \mathbf{w}' &= (x(1) \ x(3) - x(4))\end{aligned}$$

である.

(11) の固有値 ρ^2 は 2 次方程式の 2 根として与えられるから, 計算により,

$$\begin{aligned}\rho_1^2 &= 0.104, & \rho_2^2 &= 0.0601 \\ \rho_1 &= 0.322, & \rho_2 &= 0.245\end{aligned}$$

またこれらに対する固有ベクトルは,

$${}^1\mathbf{w} = w_1 \begin{pmatrix} 1 \\ -1.075 \end{pmatrix}, \quad {}^2\mathbf{w} = w_2 \begin{pmatrix} 1 \\ 0.325 \end{pmatrix}$$

と求まる (w_1, w_2 は 0 でない任意定数).

\mathbf{w} から \mathbf{x} を求める際には, (2) の $\bar{u}=0$ の条件より ${}^s x(4)/w_s (s=1, 2)$ を決定し, 次いで (3) の $V^{(1)}=1$ の条件より w_1, w_2 を定める. このようにして

$${}^1\mathbf{x} = \begin{pmatrix} 1.970 \\ 0 \\ -1.321 \\ 0.797 \end{pmatrix}, \quad {}^2\mathbf{x} = \begin{pmatrix} 2.091 \\ 0 \\ -0.071 \\ -0.751 \end{pmatrix}$$

を得る. 次いで (5) より $\mathbf{y} = (1/\rho) \mathbf{H}^+ \mathbf{G}' \mathbf{x}$ として

$${}^1\mathbf{y} = \begin{pmatrix} -2.025 \\ 0 \\ 0.964 \\ 1.479 \end{pmatrix}, \quad {}^2\mathbf{y} = \begin{pmatrix} 0.161 \\ 0 \\ -0.691 \\ 1.521 \end{pmatrix}$$

が求まる. このとき $\mathbf{x} = (1/\rho) \mathbf{F}^+ \mathbf{G} \mathbf{y}$ となっていることも確かめられる.

さらに結果を見やすくするために, 各項目ごとに平均を 0 としておくと,

$$\begin{aligned}{}^1\mathbf{x} &= \begin{pmatrix} 1.612 \\ -0.358 \\ -0.963 \\ 1.155 \end{pmatrix}, & {}^2\mathbf{x} &= \begin{pmatrix} 1.710 \\ -0.380 \\ 0.309 \\ -0.371 \end{pmatrix}, \\ {}^1\mathbf{y} &= \begin{pmatrix} -0.920 \\ 1.104 \\ -0.141 \\ 0.375 \end{pmatrix}, & {}^2\mathbf{y} &= \begin{pmatrix} 0.073 \\ -0.088 \\ -0.603 \\ 1.608 \end{pmatrix}\end{aligned}$$

となる. 結果は図示すると見やすいだろう. 図 1~3 はその例で, 図 1 は ${}^1\mathbf{x}, {}^1\mathbf{y}$ を横軸, ${}^2\mathbf{x}, {}^2\mathbf{y}$ を縦軸として各区分に与えた数値をプロットしたものであり, 図 2 および図 3 は各共著者に与えた ${}^1\mathbf{u}$ と ${}^1\mathbf{v}$ の組合せおよび ${}^2\mathbf{u}$ と ${}^2\mathbf{v}$ の組合せをプロットしたものである.

なお以上の計算には数量化 II 類について述べた [10] が参考になるだろう. ただし, 計算機のプログラムを組む上からは, 上記のような手続きは必ずしも得策とはいえない. \mathbf{H}^+ および $\mathbf{G}'\mathbf{F}^+\mathbf{G}$ などの対称性を積極的に使った方法を用いるべきであろう.

執筆にあたり, 水野欽司氏から貴重な助言をいただいた. ここに記して感謝いたします.

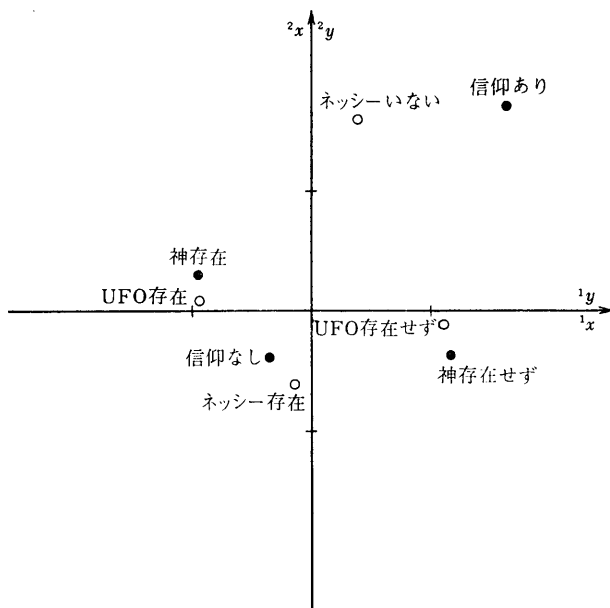


図 1

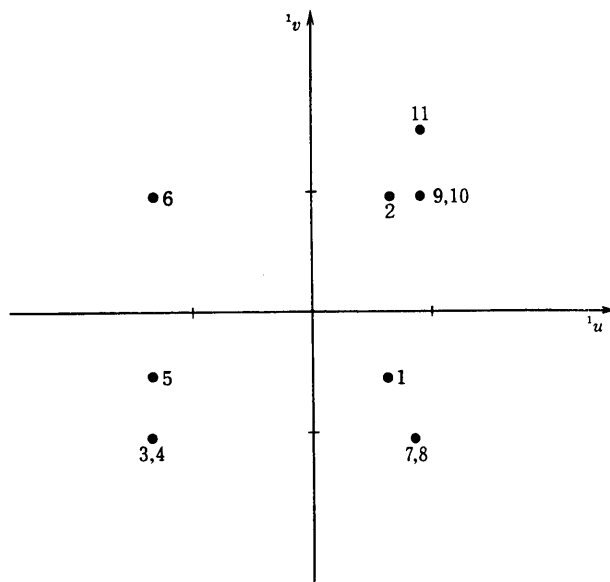


図 2

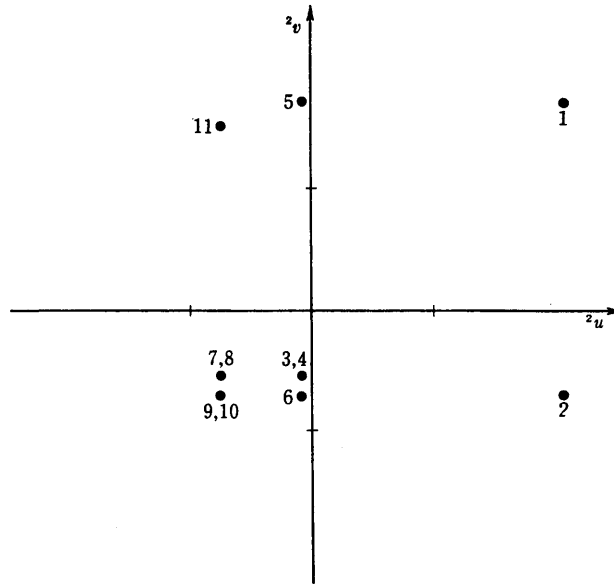


図 3

参 考 文 献

- [1] Benzécri, J.P. (1973). *L'Analyse des Données*, Tom 2 *L'Analyse des Correspondances*, Dunod.
- [2] Guttman, I. (1946). An approach for quantifying paired comparison and rank order, *Ann. Math. Statist.*, **17**, 144-163.
- [3] Hayashi, C. (1952). On the prediction of phenomena from qualitative data and the quantification of qualitative data from the mathematico-statistical point of view, *Ann. Inst. Statist. Math.*, **3**, 69-98.
- [4] Hayashi, C. (1954). Multidimensional quantification — with the applications to analysis of social phenomena —, *Ann. Inst. Statist. Math.*, **5**, 121-143.
- [5] 林知己夫, 飽戸 弘 (編) (1976). 多次元尺度解析法, サイエンス社.
- [6] Hotelling, H. (1951). A generalized *T*-test and measure of multivariate dispersion, *Proc. 2nd Berkeley Symp. Math. Statist. Prob.*, 23-41.
- [7] 石田正次 (1979). データ解析の基礎, 森北出版.
- [8] 岩坪秀一 (1971). 多次元質的データ構造の解析 —— 重相関係数による数量化 ——, 電子技術総合研究所集報, **35**, 622-632.
- [9] 岩坪秀一 (1979). *n*-way 質的データの多変量解析の手法の研究, 電子技術総合研究所研究報告, No 801.
- [10] 仁木直人 (1979). 数量化II類の解法について, 統計数理研究所集報, **26**, 115-123.
- [11] Okamoto, M. and Endo, H. (1973). Basic properties of categorical canonical correlation analysis, *J. Japan Statist. Soc.* **4**, 15-23.
- [12] Pillai, K.C.S. (1955). Some new test criteria in multivariate analysis, *Ann. Math. Statist.*, **26**, 117-121.
- [13] Wilks, S.S. (1932). Certain generalizations in the analysis of variance, *Biometrika*, **24**, 471-494.