

## 数量化 2 類

グループ(群)別に項目の選択データが与えられているとき、それらのグループを区別(判別)するのに最適な項目選択枝の数量化を行うのが数量化 2 類である。データは、各項目(アイテム) $j$ 、 $j=1,\dots,n$ 、に対する選択枝(カテゴリ) $1 \sim N_j$ から 1 つ  $k$ 、 $k=1,\dots,N_j$ 、を選ぶ形式のアイテム・カテゴリ型である。岩坪(1987)では、次の表 1 の例があげられている。1 はその選択枝が選ばれたことを、0 は選ばれなかったことを表す。

表 1 グループ別に与えられたアイテム(項目)・カテゴリ(選択枝)型データ

群	個体	項目 1		項目 2		
		選択枝 1	選択枝 2	選択枝 1	選択枝 2	選択枝 3
A	1	1	0	1	0	0
	2	1	0	1	0	0
	3	1	0	0	1	0
B	1	1	0	0	1	0
	2	0	1	0	1	0
C	1	0	1	0	1	0
	2	0	1	0	0	1
	3	0	1	0	0	1

表 1 の形式のデータを表す変数  $n_i(v, jk)$  を次のように定義する。

個体  $i$  がグループ  $v$  に属しており、かつアイテム  $j$  のカテゴリ  $k$  が選ばれているとき

$$n_i(v, jk) = 1$$

それ以外るとき

$$n_i(v, jk) = 0$$

このとき、アイテム  $j$  のカテゴリ  $k$  に数値  $x(jk)$  を与えると、グループ  $v$  の  $i$  番目の個体の得点  $y_i(v)$  を次式で与えることができる

$$y_i(v) = \sum_j \sum_k x(jk) \cdot n_i(v, jk) \quad (1)$$

$y_i(\nu)$ の値のグループ内での変動は小さく、グループ間の変動は大きくなるように  $x(jk)$ の値を定める。すなわち、相関比の2乗

$$\eta^2 = \frac{V_B}{V_T}$$

が最大になるように  $x(jk)$ の値を定める。ここで、 $V_B$ は  $y_i(\nu)$ の群間分散を、 $V_T$ は全体の分散を表す。

いま、 $n_i(\nu, jk)$ の項目選択肢間の分散共分散行列を  $A$ 、群間分散共分散行列を  $B$  とおくと、 $\eta^2$ を極大にする  $x(jk)$ は次式を満たすことが導かれる(岩坪、1987)。

$$B\mathbf{x} = \eta^2 A\mathbf{x} \quad (2)$$

ここで、 $\mathbf{x}$ は  $x(jk)$ を要素とするベクトルである。

$\eta^2$ は(2)式の固有値となっているので、相関比の2乗を最大にするためには(2)式の最大固有値に対応する固有ベクトルを  $x(jk)$ の値とする。複数個の解を求めるときは、(2)式の固有値のうち大きさの順に必要な個数を取り、対応する固有ベクトルを用いる。項目選択肢(アイテム・カテゴリ)の値  $x(jk)$ が決まると、(1)式より個体の得点  $y_i(\nu)$ も定まる。

#### 解法

(2)式を満たす  $\mathbf{x}$ を岩坪(1987)の式を用いて以下のように求める。

$A$ の固有値分解を次式で表す。

$$A = [T \ T_0] \begin{bmatrix} \Lambda_A & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} T' \\ T_0' \end{bmatrix}$$

このとき、(2)式は次のように変形できる(岩坪、(3.31)式)。

$$T'BTz = \lambda \Lambda_A z \quad (3)$$

ここで、

$$z = T(1k_1, \dots, nk_n)' \cdot U(1k_1, \dots, nk_n)' \cdot x$$

である(岩坪、1987、(3.30)式)。 $T(1k_1, \dots, nk_n)$ は、 $T$ の  $jk_j$ 行(項目  $j$ のカテゴリ  $k_j$ に

対応する行)を除いて縦に詰めた行列である。 $U(1k_1, \dots, nk_n)$ は、 $\sum_{j=1}^n N_j$  次単位行列から  $jk_j$

列を取り除いて横に詰めた行列において  $jk_j$ 行の項目  $j$ に対応する要素を  $-1$ とおいたもの

である。項目  $j$ における  $k_j$ は適当に選ぶ。

( 3 ) 式を次のように変形する。

$$\Lambda_A^{-1/2} \mathbf{T}' \mathbf{B} \mathbf{T} \Lambda_A^{-1/2} \Lambda_A^{1/2} \mathbf{z} = \lambda \Lambda_A^{1/2} \mathbf{z}$$

いま

$$\mathbf{u} = \Lambda_A^{1/2} \mathbf{z} \quad (4)$$

とおけば

$$\Lambda_A^{-1/2} \mathbf{T}' \mathbf{B} \mathbf{T} \Lambda_A^{-1/2} \mathbf{u} = \lambda \mathbf{u}$$

と書けるので、 $\lambda$  と  $\mathbf{u}$  は  $\Lambda_A^{-1/2} \mathbf{T}' \mathbf{B} \mathbf{T} \Lambda_A^{-1/2}$  の固有値と固有ベクトルとして求めることができる。 $\mathbf{u}$  が求まると、岩坪の ( 3.35 ) 式

$$\mathbf{w} = \mathbf{T} \mathbf{z} \quad (5)$$

および ( 3.39 ) 式

$$x(jk) = w(jk) - w(jk_j) \quad (6)$$

を用いて、 $\mathbf{x}$  を求めることができる。すなわち、式 ( 4 ) と式 ( 5 ) から  $\mathbf{w}$  が

$$\mathbf{w} = \mathbf{T} \Lambda_A^{-1/2} \mathbf{u}$$

と求まる。この  $\mathbf{w}$  から  $\mathbf{x} = (x(jk))$  が ( 6 ) 式より算出できる。

$\mathbf{x}$  は、項目内のカテゴリに同じ値を加えても ( 2 ) 式の解であるという性質がある ( 付録参照 )。この性質によって、 $x(jk)$  の項目内の平均が 0 になるように基準化する。すなわち、

$$x^*(jk) = w(jk) - \sum_{k=1}^{N_j} \bar{n}(jk) w(jk) \quad (7)$$

として基準化した最適解  $\mathbf{x}^* = (x^*(jk))$  を求める。( 7 ) 式で与えられる解が平均値が 0 になるように基準化されていることは以下のように確かめられる。

$$\begin{aligned} \sum_k \bar{n}(jk) x^*(jk) &= \sum_k \bar{n}(jk) w(jk) - \sum_k \left[ \bar{n}(jk) \left\{ \sum_{k'} \bar{n}(jk') w(jk') \right\} \right] \\ &= \sum_k \bar{n}(jk) w(jk) - \sum_k \left[ \bar{n}(jk) \left\{ \sum_{k'} \bar{n}(jk') w(jk') \right\} \right] \\ &= \sum_k \bar{n}(jk) w(jk) - \left\{ \sum_{k'} \bar{n}(jk') w(jk') \right\} \sum_k \bar{n}(jk) \\ &= \sum_k \bar{n}(jk) w(jk) - \left\{ \sum_{k'} \bar{n}(jk') w(jk') \right\} \sum_k \left\{ \frac{1}{m} \sum_v \sum_i n_i(v, jk) \right\} \\ &= \sum_k \bar{n}(jk) w(jk) - \left\{ \sum_{k'} \bar{n}(jk') w(jk') \right\} \left[ \frac{1}{m} \sum_v \sum_i \left\{ \sum_k n_i(v, jk) \right\} \right] \\ &= \sum_k \bar{n}(jk) w(jk) - \left\{ \sum_{k'} \bar{n}(jk') w(jk') \right\} \left[ \frac{1}{m} \sum_v \sum_i 1 \right] \end{aligned}$$

$$= \sum_k \bar{n}(jk)w(jk) - \left\{ \sum_{k'} \bar{n}(jk')w(jk') \right\} \times 1$$

$$= 0$$

## プログラム

プログラム PQ2a.dpr は、 $x^*(jk)$ 、 $y_i(v)$ 、および  $y_i(v)$  のグループごとの平均値  $\overline{y(v)}$  を求めるものである。このプログラムを実行すると、図 1 のフォームが表示される。

図 1 起動時に表示されるフォーム

「Add Data」ボタン、および「Add Var」ボタンをクリックして行数および項目（アイテム）数をデータに合わせる。「Add Data」ボタンのクリックでアクティブなセルの下に行が追加・挿入され、「Add Var」ボタンのクリックでアクティブなセルの右側に列（項目欄）が追加・挿入される。セルはそのセルのクリックでアクティブになる。「Del Data」ボタンあるいは「Del Var」ボタンをクリックするとアクティブなセルを含む行あるいは列が削除される。追加・挿入あるいは削除の操作は、データの設定の途中においてもできる。図 2 はデータの設定例である。

	グループID	個体ID	項目 1	項目 2
ラベル	*****	*****	a	b
項目数	*****	*****	2	3
個体 1	A	1	1	1
個体 2	A	2	1	1
個体 3	A	3	1	2
個体 4	B	1	1	2
個体 5	B	2	2	2
個体 6	C	1	2	2
個体 7	C	2	2	3
個体 8	C	3	2	3

図2 データの設定

第 2 行目（ラベル行）に各項目のラベルを設定する。設定された文字列の 1 バイト目の文字がマップの描画に用いられる。グループ ID と個体 ID の欄はラベルを設定しない（データとしては無視される）。第 3 行目（項目数の行）に各項目のカテゴリ数を設定する。この行もグループ ID と個体 ID の欄は数値を設定しない（無視される）。第 4 行目から各個体のデータを設定していく。グループ ID 欄に所属するグループを表す文字（1 バイト）を設定する。この文字はマップの描画においてもグループを表すのに用いられる。個体 ID に個体のラベルを設定する。項目欄には各項目における個体の該当するカテゴリ番号を設定する。

項目  $j$  におけるカテゴリ数が  $N_j$  であるとき、カテゴリ番号は 1 から  $N_j$  までの数値が当てられているものとする。図 2 は表 1 のデータを設定したものである。

図 2 のように設定したデータは、「Save(CSV)」ボタンのクリックで保存することができる。保存したデータは「Open(CSV)」ボタンのクリックで読み込むことができる。データは CSV 形式のファイルで保存されるので、Excel で開くこともできる。図 2 のデータを保存したファイルを Excel で開くと図 3 のように表示される。

	A	B	C	D	E
1	*****	*****	a	b	
2	*****	*****	2	3	
3	A	1	1	1	
4	A	2	1	1	
5	A	3	1	2	
6	B	1	1	2	
7	B	2	2	2	
8	C	1	2	2	
9	C	2	2	3	
10	C	3	2	3	
11					

図3 Excel で開いたデータ

逆に、Excel で図3の形式によって作成したデータは、CSV 形式のファイルで保存する（拡張子を.csv として保存する）と「Open(CSV)」のクリックで読み込むことができる。

データの設定後、「Calc」ボタンをクリックすると計算が始まる。「Calc」ボタンのクリックで、先ず図4のダイアログボックスが表示される。

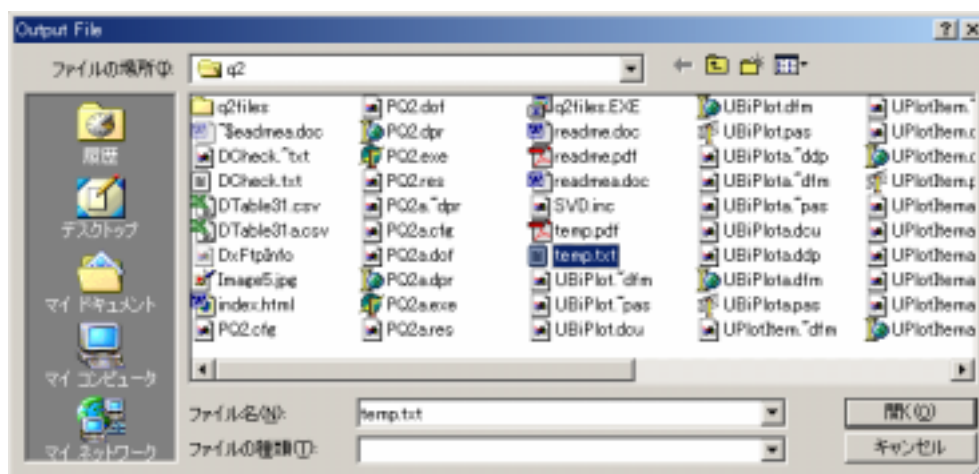


図4 計算結果の書き出し用ファイル名の設定

図4のダイアログボックスで設定した名前のファイルに計算結果がテキストファイルとして書き出される。このファイルは、プログラムの実行終了後エディタなどで開いて見ることができる。

出力ファイル名の設定後、ダイアログボックスの「開く」ボタンをクリックすると計算が始まる。計算が終了すると図5のフォームになる。

	グループID	個体ID	項目 1	項目 2
ラベル	*****	*****	a	b
項目数	*****	*****	2	3
個体 1	A	1	1	1
個体 2	A	2	1	1
個体 3	A	3	1	2
個体 4	B	1	1	2
個体 5	B	2	2	2
個体 6	C	1	2	2
個体 7	C	2	2	3
個体 8	C	3	2	3

図 5 計算終了時のフォーム

「Plot(Subj)」ボタンは個体の値  $y_i(\nu)$  とグループの平均値  $\overline{y(\nu)}$  を図示するものであり、  
「Plot(Item)」ボタンはアイテム・カテゴリ（項目・選択肢）の値  $x^*(jk)$  とグループの平均値  $\overline{y(\nu)}$  を図示するためのものである。

「Plot(Subj)」ボタンをクリックすると図 6 のフォームが表示される。

図 6 描画次元の選択

個体の値  $y_i(\nu)$  とグループの平均値  $\overline{y(\nu)}$  のプロットにおいて使用する解（相関比の 2 乗の

大きさの順に対応する解を 1 軸（第 1 次元）、2 軸（第 2 次元）・・・と数える）の軸（次元）を選ぶ。軸の選択後、「OK」ボタンをクリックすると、選んだ軸の解がプロットされる（図 7）。

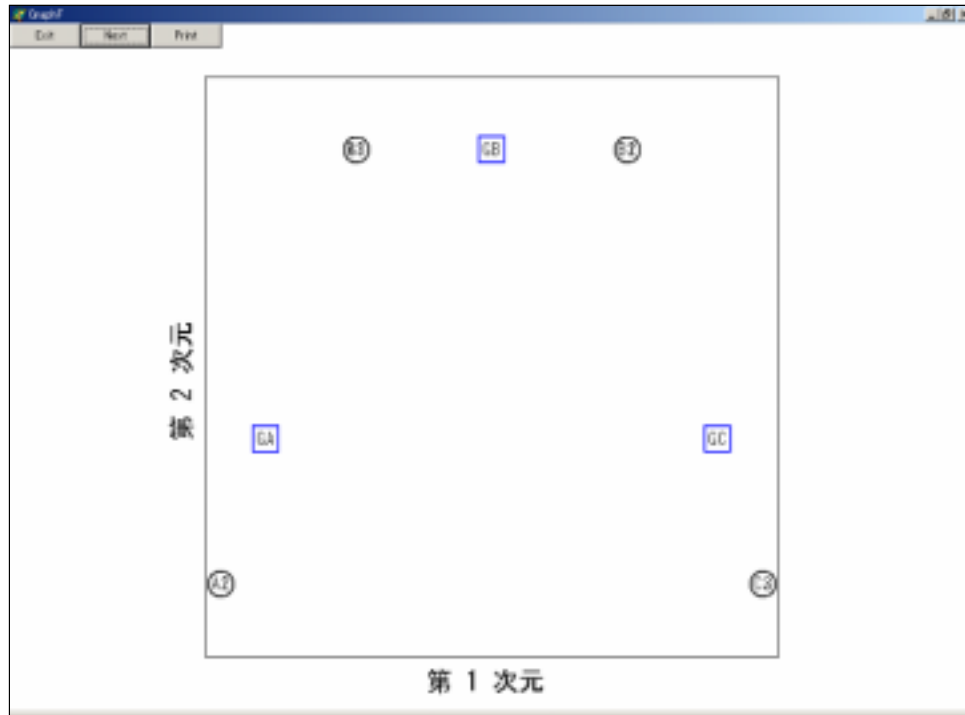


図 7 個体の値とグループ平均値の描画

個体の値は小円によって表されている。小円内の先頭の文字は、グループのラベルとして入力されたものが表示されている。2 番目の数値はグループ内の個人の番号である。入力データファイルに書かれている順番を表す。グループ内の個人の数が 9 人を超えるときは、10 人目以降は ? で表される。図 7 の表示では、グループ A に属する個人は左側、グループ B に属する個人は上方、グループ C に属する個人は右側に位置している。グループの平均値は、小正方形で表されている。小正方形内の 2 文字目がそのグループを表すラベルである。それぞれのグループ内の個人を表す位置の平均位置にグループを表す小正方形が位置していることがわかる。ディスプレイに描画されている図は、図 7 の画面の左上のボタン「Print」をクリックするとプリンタに出力される。「Next」ボタンをクリックすると図 6 の描画次元の選択に戻る。描画したい次元を再設定した後、「OK」ボタンをクリックすると再設定した次元での描画が行われる。「Exit」ボタンのクリックで図 5 の計算終了時のフォームに戻る。

アイテム・カテゴリの位置を表示するときは、図 5 のフォームにおいて「Plot(Item)」ボタンをクリックする。「Plot(Item)」ボタンのクリックで図 8 のフォームが表示される。



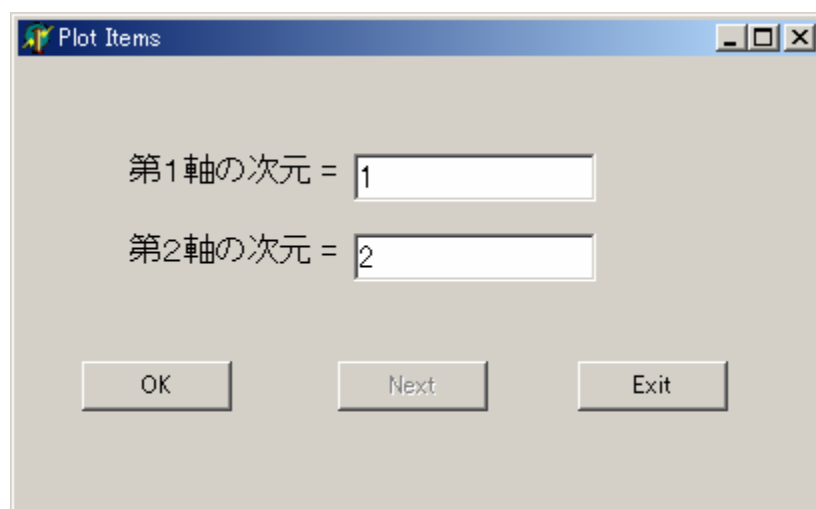


図8 描画次元の選択

図6の場合と同様に、描画する軸（次元）を選び、「OK」ボタンをクリックすると選択した軸のアイテム・カテゴリ値  $x^*(jk)$  とグループの平均値  $\overline{y(\nu)}$  がプロットされる（図9）。

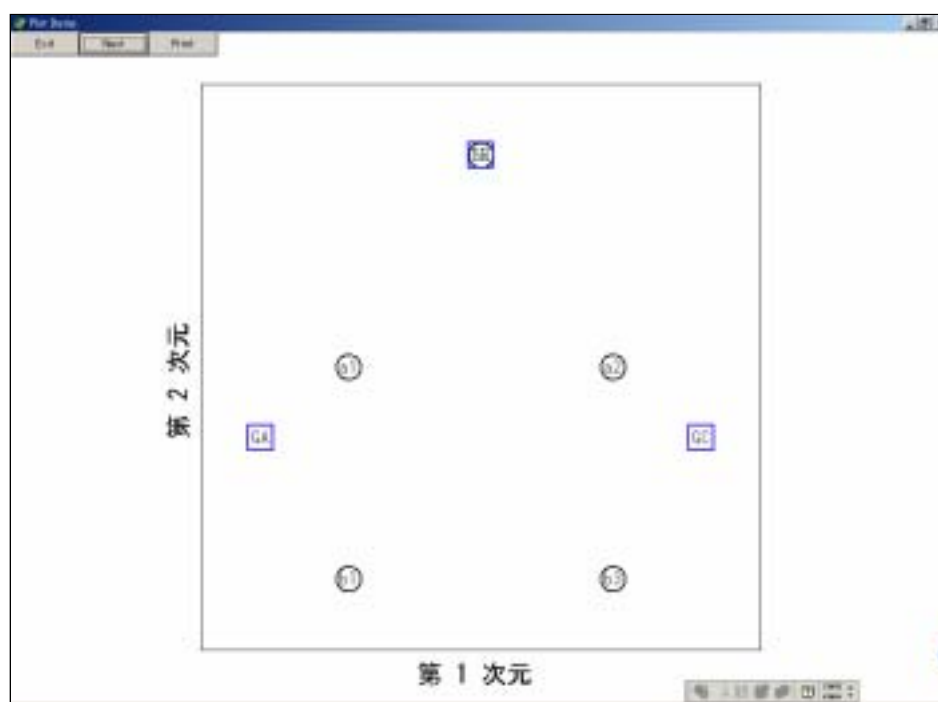


図9 アイテム・カテゴリ値とグループ平均値のプロット

グループ平均値は図7の個人の値のプロットのとおり描かれている。アイテム・カテゴリの値は小円によって描かれ、小円内の最初の文字がそのアイテム（項目）を表す。アイテムは、図2のデータ設定において項目のラベルとして設定された文字列の先頭

の1バイトが用いられる。アイテムを表す文字の右側の数値はカテゴリを表す。カテゴリはアイテム内の通し番号で1、2、・・・、 $N_j$ と数値で表している。図9から、アイテムaのカテゴリ1およびアイテムbのカテゴリ1を選ぶ者はグループAに、アイテムaのカテゴリ2およびアイテムbのカテゴリ3を選ぶ者はグループCに属する傾向のあることがわかる。アイテムbのカテゴリ2を選ぶ者はグループBに属する傾向にある。これらの傾向は表1のデータにおいて確認できる。

「Exit」、「Print」、「Next」ボタンの使い方は、図7の場合と同じである。

図5のフォームにおいて「Exit」ボタンをクリックすると、プログラムの実行終了となる。プログラムの実行終了後、図4で設定した名前の計算結果の出力ファイルを開いて見ることができる。2図のデータ設定の場合、出力ファイルの内容はリスト1のようになっている。

リスト1 計算結果の出力例

*****	*****	a	b
*****	*****	2	3
A	1	1	1
A	2	1	1
A	3	1	2
B	1	1	2
B	2	2	2
C	1	2	2
C	2	2	3
C	3	2	3
入力データ...			
グループ: A			
	1	10	100
	2	10	100
	3	10	010
グループ: B			
	1	10	010
	2	01	010
グループ: C			
	1	01	010
	2	01	001
	3	01	001
個体総数 = 8			
固有値 (相関比の2乗) =			
		0.833333333333333	
		0.333333333333333	
基準化された最適スコア			
a1 ->	-0.63246	0.00000	
a2 ->	0.63246	-0.00000	
b1 ->	-0.63246	-1.00000	
b2 ->	0.00000	1.00000	
b3 ->	0.63246	-1.00000	

個体得点		
グループ A		
1	<1>	-1.265    -1.000
2	<2>	-1.265    -1.000
3	<3>	-0.632    1.000
グループ B		
1	<1>	-0.632    1.000
2	<2>	0.632    1.000
グループ C		
1	<1>	0.632    1.000
2	<2>	1.265    -1.000
3	<3>	1.265    -1.000
グループ平均得点		
グループ A		
		-1.054    -0.333
グループ B		
		0.000    1.000
グループ C		
		1.054    -0.333

リスト 1 を見ると、まず、グリッドのセルに設定されていたデータがそのまま書き出されている。続いて、入力データがグループ別に書き出されるが、このときは項目のカテゴリの選択は 0 - 1 のパターンで書き出されている。選択されたカテゴリの値が 1、それ以外のカテゴリの値は 0 である。グループごとのデータの出力後、個体の総数が出力される。固有値（相関比の 2 乗）の書き出しに続いて、図 9 で表示されているアイテム・カテゴリの値が基準化された最適スコアとして書き出され、図 7 で表示されている個体の値が個体得点として書き出されている。最後にグループごとの平均得点がグループ平均得点として書き出されている。

# 付 録

$\sum_{j=1}^n N_j$  次列ベクトルにおいて項目  $j$  に対応する要素を  $c$ 、それ以外の要素を 0 とおいたものを次式のように  $\mathbf{d}$  とおく。

$$\mathbf{d}' = \begin{bmatrix} 0 & \cdots & 0 & \underbrace{c \cdots c}_{\text{項目 } j \text{ のカテゴリ}} & 0 & \cdots & 0 \end{bmatrix}$$

このとき  $\mathbf{Ad}$  の  $(j_1 k_1)$  行目の要素は以下ようになる。

$$\begin{aligned} \sum_{k=1}^{N_j} a(j_1 k_1) c &= \frac{c}{m} \sum_k \left[ \sum_{v=1}^g \sum_{i=1}^{m_v} \{n_i(v, j_1 k_1) - \bar{n}(j_1 k_1)\} \{n_i(v, jk) - \bar{n}(jk)\} \right] \\ &= \frac{c}{m} \sum_v \sum_i \left[ \{n_i(v, j_1 k_1) - \bar{n}(j_1 k_1)\} \left\{ \sum_k n_i(v, jk) - \sum_k \bar{n}(jk) \right\} \right] \quad (\text{A1}) \end{aligned}$$

ここで、

$$\sum_k n_i(v, jk) = 1$$

および

$$\begin{aligned} \sum_k \bar{n}(jk) &= \sum_k \left\{ \frac{1}{m} \sum_{v'} \sum_{i'} n_{i'}(v', jk) \right\} \\ &= \frac{1}{m} \sum_{v'} \sum_{i'} \sum_k n_{i'}(v', jk) \\ &= \frac{1}{m} \sum_{v'} \sum_{i'} 1 \\ &= 1 \end{aligned}$$

であることに注意すると、(A1) 式は

$$\frac{c}{m} \sum_v \sum_i [\{n_i(v, j_1 k_1) - \bar{n}(j_1 k_1)\} \{1 - 1\}] = 0$$

となる。すなわち、次式が成り立つ。

$$\sum_{k=1}^{N_j} a(j_1 k_1) c = 0$$

上式は  $\mathbf{Ad}$  の任意の  $(j_1 k_1)$  行について成り立つ。したがって、

$$\mathbf{Ad} = \mathbf{0} \quad (\text{A2})$$

である。

同様に、 $\mathbf{Bd}$  の  $(j_1 k_1)$  行目の要素は次のように計算できる。

$$\begin{aligned} \sum_k b(j_1 k_1, jk)c &= \frac{c}{m} \sum_k \left[ \sum_v m_v \{ \bar{n}(v, j_1 k_1) - \bar{n}(j_1 k_1) \} \{ \bar{n}(v, jk) - \bar{n}(jk) \} \right] \\ &= \frac{c}{m} \sum_v \left[ m_v \{ \bar{n}(v, j_1 k_1) - \bar{n}(j_1 k_1) \} \left[ \sum_k \{ \bar{n}(v, jk) - \bar{n}(jk) \} \right] \right] \quad (\text{A3}) \end{aligned}$$

ここで、次式が成り立っている。

$$\begin{aligned} \sum_k \{ \bar{n}(v, jk) - \bar{n}(jk) \} &= \sum_k \left\{ \frac{1}{m_v} \sum_{i'} n_{i'}(v, jk) \right\} - \sum_k \left\{ \frac{1}{m} \sum_{v'} \sum_{i'} n_{i'}(v', jk) \right\} \\ &= \frac{1}{m_v} \sum_{i'} \sum_k n_{i'}(v, jk) - \frac{1}{m} \sum_{v'} \sum_{i'} \sum_k n_{i'}(v', jk) \\ &= \frac{1}{m_v} \sum_{i'} 1 - \frac{1}{m} \sum_{v'} \sum_{i'} 1 \\ &= 1 - 1 \\ &= 0 \end{aligned}$$

故に (A3) 式は次のようになる。

$$\frac{c}{m} \sum_v [m_v \{ \bar{n}(v, j_1 k_1) - \bar{n}(j_1 k_1) \} \times 0] = 0$$

すなわち、

$$\sum_k b(j_1 k_1, jk)c = 0$$

が成り立つ。

上式は  $\mathbf{Bd}$  の任意の  $(j_1 k_1)$  行について成り立つ。したがって、

$$\mathbf{Bd} = \mathbf{0} \quad (\text{A4})$$

が成り立つ。

(A2) 式と (A4) 式より、 $\mathbf{x}$  が

$$\mathbf{Bx} = \lambda \mathbf{Ax} \quad (\text{A5})$$

を満たすとき、

$$\mathbf{B}(\mathbf{x} + \mathbf{d}) = \mathbf{Bx} + \mathbf{Bd} = \mathbf{Bx} = \lambda \mathbf{Ax} = \lambda \mathbf{Ax} + \lambda \mathbf{Ad} = \lambda \mathbf{A}(\mathbf{x} + \mathbf{d})$$

が成り立つ。すなわち、 $\mathbf{x} + \mathbf{d}$  も (A5) 式を満たす。

## 参考文献

岩坪秀一「数量化法の基礎」朝倉書店、1987