

2007年3月30-31日 行動計量学会 第10回春の合宿セミナー八王子セミナー八
ウス 2007-04-02 修正版

数量化と非線形多変量解析の方法

大学入試センター 研究開発部 大津起夫¹

多変量データの次元縮約法、特に尺度推定を含む手法を概説．

- 主成分分析と特異値分解
- 対応分析（パターン分類の数量化）
- 尺度最適化を伴う主成分分析(OSMOD)
数量化3類（多重対応分析）の問題点を（部分的に）克服
- 非線形因子分析

¹<http://www.rd.dnc.ac.jp/~otsu>

1. 主成分分析法と特異値分解

2

1 主成分分析法と特異値分解

主成分分析法: (principal component analysis, PCA) は複数の連続変数間の相互関係の簡潔な表現を求めるために広く用いられている方法 .

(対応分析: 2重分割表の構造を求める . 共通点 = 特異値分解)

分析対象となる変数: $X_j, (j = 1, \dots, p)$

観測されたデータの値: $x_{ij}, (i = 1, \dots, n; j = 1, \dots, p)$
(i : サンプル(被験者)の番号, j : 観測変数の番号)

1. 主成分分析法と特異値分解

3

いいたいこと

- 対応分析（パターン分析の数量化）+ Bertin グラフィックス
- MCA(数量化 3 類) より尺度最適化つき主成分分析 (OSMOD, PRINQUAL, ほか)

1. 主成分分析法と特異値分解

4

合成変量：

$$Z_m = \mu_m + b_{m1}X_1 + b_{m2}X_2 + \cdots + b_{mp}X_p, \quad (m = 1, \dots, q)$$

残差 2 乗和

$$SSQ_q = \sum_{j=1}^p E(|X_j - \mu_j - \sum_{m=1}^q a_{jm}Z_m|^2) \quad (1)$$

$$E(X_j) = \mu_j + \sum_m a_{jm}E(Z_m)$$

(もし $X_j, (j = 1, \dots, p)$ の平均がゼロなら $\mu_j = 0$)

各変数 X_j の平均がゼロの場合，観測されたデータの要素で表すと

$$SSQ_q = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^p (x_{ij} - \sum_{m=1}^q a_{jm}z_{im})^2$$

1. 主成分分析法と特異値分解

5

X_1, \dots, X_p を要素とする縦ベクトルを x
(資料では大文字の太字になっているが, 行列と間違えやすいので小文字にした)

a_{jm} を要素とする大きさ $p \times q$ の行列を A

Z_1, \dots, Z_q を要素とする縦ベクトルを z (資料は大文字)

Az で x (p 次元) を近似したい. z は x の1次式.

$$z = Bx \text{ なので } Az = ABx = Cx$$

($m > 1$ のときは B は一通りには決まらない)

ただし $\text{rank} C = q$

1. 主成分分析法と特異値分解

6

上の SSQ の式は行列記法を用いると次のように表される .

$$\begin{aligned}SSQ_q &= \text{trace} E\{(\boldsymbol{x} - \boldsymbol{A}z)(\boldsymbol{x} - \boldsymbol{A}z)^T\} \\ &= \text{trace} E\{(\boldsymbol{x} - \boldsymbol{C}x)(\boldsymbol{x} - \boldsymbol{C}x)^T\} \\ &= \text{trace}(\boldsymbol{I}_p - \boldsymbol{C})E(\boldsymbol{x}x^T)(\boldsymbol{I}_p - \boldsymbol{C})^T\end{aligned}$$

記号 trace (トレース) は , 対角要素の和を表す記号 .

確率変数のベクトル x を , $Az = Cx$ で平均的に近似する .

$$x_{ij} \sim \sum_{m=1}^q a_{jm} z_{ij} = \sum_{m=1}^q c_{jm} x_{im}$$

(x_{ij}) の低ランク近似になっている (q が p よりずっと小さいところがミソ)

1. 主成分分析法と特異値分解

7

主成分分析のモデルは理論的には単純（観測変数の最小2乗近似を観測変数の1次式で行う）。

因子分析 (factor analysis) はもっと凝った理論的背景を持っている。主成分分析より大胆な（攻撃的？）モデル推定を行う。

最小の SSQ_q をもたらす $Z_m, m = 1, \dots, q$ は、主成分スコアまたは単に主成分 (principal component) と呼ばれる。

主成分は変数の尺度に依存する。（因子分析の解は尺度に依存しない）

1. 主成分分析法と特異値分解

8

主成分の求め方: 変数の共分散行列 ($p \times p$): $E(\mathbf{x}\mathbf{x}^T) = (E(X_j X_{j'})) = \Sigma_{XX}$

(1) 実対称行列の固有値分解: $\Sigma_{XX} = U\Lambda U^T$

ここで U は直交行列である (異なる列の内積, 異なる行の内積がゼロ). また Λ は固有値 $(\lambda_1, \dots, \lambda_p)$ を成分とする対角行列である.

Λ の対角要素は降順に並んでいるものとする.

U の第1列から第 q 列の $p \times q$ の行列を U_1 とする. ($U = (U_1|U_2)$)

$A = B^T = U_1$ とおくことにより SSQ_q を最小化できる.

1. 主成分分析法と特異値分解

固有値分解の数値例:

共分散行列

	[,1]	[,2]	[,3]	[,4]
[1,]	4	3	2	1
[2,]	3	4	3	2
[3,]	2	3	4	3
[4,]	1	2	3	4

1. 主成分分析法と特異値分解

固有値分解の結果

lambda_1	lambda_2	lambda_3	lambda_4
11.099	3.414	0.901	0.586

	u_1	u_2	u_3	u_4
[1,]	0.448	0.653	0.547	0.271
[2,]	0.547	0.271	-0.448	-0.653
[3,]	0.547	-0.271	-0.448	0.653
[4,]	0.448	-0.653	0.547	-0.271

$$Z_1 = 0.448X_1 + 0.547X_2 + 0.547X_3 + 0.448X_4$$

$$Z_2 = 0.653X_1 + 0.271X_2 - 0.271X_3 - 0.653X_4$$

Z_1 の分散は $\lambda_1 = 11.099$, Z_2 の分散は $\lambda_2 = 3.414$.

これらの共分散はゼロ .

1. 主成分分析法と特異値分解

11

$C = AB = U_1 U_1^T$ ($p \times p$ の行列) は, 上位 q 個の固有値に対応する固有ベクトルによって張られる変数の空間への直交射影子 (対称で巾 (べき) 等な行列, つまり行列積で 2 乗しても元通りであること) になっている.

$Z_m = u_m^T x$ を第 m 主成分と呼ぶ. (u_m は U の第 m 列の縦ベクトルであり, u_m^T はそれを転置した横ベクトルを表す. $u_m^T x$ はこれら 2 つのベクトルの内積を表す.) (x は縦ベクトル)

1. 主成分分析法と特異値分解

12

$\Rightarrow U_1^T$ を p 次元の縦ベクトル x に左からかけることが情報を縮約する (q 個の主成分 Z_1, \dots, Z_q にまとめる) ことになる . ($z = U_1^T x \Rightarrow z^T = x^T U_1$)

これに $U_1(p \times q)$ をさらに左からかけると, q 個の主成分をもとの p 次元のベクトルの世界 (観測変数) に引き戻す . ($\hat{x} = U_1 U_1^T x \Rightarrow \hat{x}^T = x^T U_1 U_1^T$)

1. 主成分分析法と特異値分解

13

もうひとつの求め方（データ行列を直接分解する）

(2) 特異値分解 (singular value decomposition)

一般的に $n \times p$ の行列（正方行列とは限らない） X は

$$X = VDU^T \quad (2)$$

と、直交行列（の一部）と対角行列の積に分解される。
ここで $n \geq p$ とする。

V は $n \times p$ の行列であり、その列ベクトルは各々長さ1であり互いに直交している。

U は、 $p \times p$ の直交行列

D は $p \times p$ の非負要素の対角行列。

1. 主成分分析法と特異値分解

14

各変数の標本平均がゼロであると仮定すると，共分散行列は

$$S = \frac{1}{n} X^T X$$

であり，上の特異値分解を用いると， $S = \frac{1}{n} U D V^T V D U^T = \frac{1}{n} U D^2 U^T$ となる．

この U の各列が S の固有ベクトルであり，対角行列 D の要素の2乗が S の固有値を n 倍したものであることも分かる．

1. 主成分分析法と特異値分解

15

U のうち，大きな特異値に対応する第1列から第 q 列を取り出すと，これが上の U_1 となっている．ここで q 個の主成分は $n \times q$ の行列

$$XU_1 = (\mathbf{v}_1 \times d_1, \dots, \mathbf{v}_q \times d_q)$$

で表され，また主成分を用いた各変数の近似は

$$XU_1U_1^T$$

となる．

1. 主成分分析法と特異値分解

16

V の列ベクトルを v_1, \dots, v_p とし, また U の列ベクトルを u_1, \dots, u_p とすると, 特異値分解の式は

$$X = VDU^T = \sum_m^p d_m v_m u_m^T \quad (3)$$

となる. 上記の主成分による近似は,

$$\sum_m^q d_m v_m u_m^T$$

であり, 行列の階数 (ランク) q であるもののうち, X との残差 2 乗和が最小になることが確認できる.

1. 主成分分析法と特異値分解

特異値分解の数値例：

$X(5 \times 3)$ (各列の平均はゼロ)

	x1	x2	x3
[1,]	-2	2	-6
[2,]	-1	-1	-3
[3,]	0	-2	0
[4,]	1	-1	3
[5,]	2	2	6

1. 主成分分析法と特異値分解

特異値分解の結果 $X = VDU^T$

(注意：R や SPlus の特異値分解では $X = UDV^T$ となっている．ここでは固有値分解の結果と一致させるために U と V を逆にしている)

特異値	d_1	d_2	d_3
	10.000	3.742	0.000

	v_1	v_2	v_3	(d_k × v_k が主成分Z_k)
[1,]	-0.632	-0.535	-0.561	(v_3の値は無意味)
[2,]	-0.316	0.267	0.102	
[3,]	0.000	0.535	-0.510	
[4,]	0.316	0.267	-0.612	
[5,]	0.632	-0.535	-0.204	

1. 主成分分析法と特異値分解

	u_1	u_2	u_3	(u_3も無意味)
[1,]	0.316	0	0.949	
[2,]	0.000	-1	0.000	
[3,]	0.949	0	-0.316	

(u_m は共分散行列の固有ベクトルになる)

1. 主成分分析法と特異値分解

20

主成分分析の（因子分析と比べての）特徴

因子分析と比べて，理論的には洗練されていないが，単純で計算が速い。（確率モデルを想定したあてはめではない）。

変数の尺度が影響を及ぼす（かならずしも欠点ではないが）。

最近の速いPCなら 2000 次元の実対称固有値計算が 30 秒以内で可能（Intel MKL や ACML (AMD Core Math Library) を使うと速い。）

変数の平均がゼロでなくても，特異値分解による行列の近似は可能。

2. 対応分析（パターン分類の数量化）
- 2 対応分析（パターン分類の数量化）
 - 2重分割表の構造を把握したい。

χ^2 検定で，独立性の仮説が棄却されても，それだけではあまり意味がない。

行と列に，データの構造を反映した尺度を与えたい。

2. 対応分析（パターン分類の数量化）

表 1: NLSY79: 1990年における職業と年齢補正済 AFQT
 (クロスセクショナルサンプル) (大津, 2003)

職業	年齢補正済 AFQT					計
	1	2	3	4	5	
事務 (Clerical)	115	194	219	187	132	847
製造 (Craft)	132	129	128	105	56	550
農業 (Farm)	23	10	16	8	7	64
労働 (Labor)	89	65	40	22	13	229
管理職 (Mgr)	56	90	135	164	152	597
熟練工 (Operatv)	149	97	89	59	27	421
専門技術職 (Prof/Tech)	31	80	115	224	428	878
営業 (Sales)	14	37	37	78	74	240
サービス (Service)	197	166	126	99	54	642
運輸 (Transpt)	44	41	27	28	7	147
計	850	909	932	974	950	4615

2. 対応分析（パターン分類の数量化）

2.1 対応分析のモデル

表1のような2重分割表の*i*行*j*列の数値を n_{ij} とする．行数を*I*とし，列数を*J*とする．また，対象（被験者）の総数を

$$N = \sum_i \sum_j n_{ij}$$

とする．

$$P_{ij} = n_{ij}/N, \quad i = 1, \dots, I; \quad j = 1, \dots, J$$

は相対頻度を表す．また，周辺相対頻度を

$$P_{i.} = \sum_{j=1}^J P_{ij}, \quad P_{.j} = \sum_{i=1}^I P_{ij}$$

と表す．

2. 対応分析 (パターン分類の数量化)

24

行スコアを $x_i, i = 1, \dots, I$, 列スコアを $y_j, j = 1, \dots, J$

i 行 j 列のセルの座標 (x_i, y_j)

この2次元上の位置に P_{ij} の確率が与えられているとすると, これによって定まる分布の平均は次のように表される .

$$\mu_x = \sum_i P_{i.} x_i, \quad \mu_y = \sum_j P_{.j} y_j \quad (4)$$

以下では, μ_x と μ_y とが, とともにゼロベクトルに制約する .

2. 対応分析（パターン分類の数量化）

25

この仮定のもとで，行スコア x と列スコア y の分散は，それぞれ

$$\sigma_x^2 = \sum_i P_i x_i^2, \quad \sigma_y^2 = \sum_j P_j y_j^2 \quad (5)$$

となる．以下ではさらに $\sigma_x^2 = \sigma_y^2 = 1$ の制約も，スコアが満たしているとする．

このとき x と y との相関は

$$r_{xy} = \sum_i \sum_j P_{ij} x_i y_j \quad (6)$$

となる．この値は，行と列との関係を表す一つの指標と考えられる．相関 r が大きければ，行と列との間に強い関係があるとみなせる．

2. 対応分析（パターン分類の数量化）

26

r_{xy} を最大にする行スコア，列スコアを求めたい。

これらのスコアは，データに内在する構造をよく表しているだろう（たぶん）。

周辺相対頻度を要素とする対角行列（対角線上の要素のみが非ゼロの行列）を定める。

$$\mathbf{F} = \text{diag}(P_{1.}, \dots, P_{I.}), \quad \mathbf{G} = \text{diag}(P_{.1}, \dots, P_{.J})$$

これらの要素の平方根を要素とする対角行列を，それぞれ

$$\mathbf{F}^{1/2} = \text{diag}(\sqrt{P_{1.}}, \dots, \sqrt{P_{I.}}), \quad \mathbf{G}^{1/2} = \text{diag}(\sqrt{P_{.1}}, \dots, \sqrt{P_{.J}})$$

とする。

2. 対応分析 (パターン分類の数量化)

27

相関(6)を制約のもとで最大化するスコアは、次の特異値分解を利用して得ることができる。

$$F^{-1/2}PG^{-1/2} = UDV^T \quad (7)$$

D の対角要素はすべて非負であり、降順にならんでいるものとする。

変則的だが U と V の列ベクトルを、それぞれ添え字ゼロから始まるものとし、

$$u_k, v_k, (k = 0, 1, \dots, K)$$

とする。また $K = \min(I - 1, J - 1)$ である。

対角行列 D の要素を d_0, d_1, \dots, d_K とする。

2. 対応分析（パターン分類の数量化）

$$F^{-1/2}u_k = x_k, G^{-1/2}v_k = y_k$$

とおくと， x_1 と y_1 とが(6)を最大にし，制約を満たすスコアになり， $d_1 = r_{max}$ となる．

また， x_0 と y_0 とは，要素が全て1のベクトルとなり， $d_0 = 1$ となる．添え字ゼロに対応するスコアは，重みつき平均がゼロの制約を満たさないため，対応分析の解とはならない．

2. 対応分析 (パターン分類の数量化)

29

特異値分解の式 (7) の左から $F^{1/2}$ を乗じ, また右から $G^{1/2}$ を乗ずると

$$\begin{aligned} P &= F^{1/2} U D V^T G^{1/2} \\ &= F F^{-1/2} U D V^T G^{-1/2} G \\ &= F \left(\sum_{k=0}^K d_k \mathbf{x}_k \mathbf{y}_k^T \right) G \end{aligned}$$

となる .

2. 対応分析（パターン分類の数量化）

30

対応分析の解を考えることは，上の式で $k = 0$ と $k = 1$ の部分によって2重分割表 P の重み付きの近似を行っているともみなせる．

$$P \sim F \left(d_0 x_0 y_0^T + d_1 x_1 y_1^T \right) G$$

また，対応分析の利用においては1次元より多くの解，つまり $k = 2, 3, \dots$ に対応する x_k および y_k を解釈の対象とすることがある．

これらの $k = 2, \dots$ に対応する解は，1次元ほど多くの部分を説明しないが，データのより詳細な部分についての情報を与える場合がある（そうでない場合もある）．

2. 対応分析（パターン分類の数量化）

31

U が直交行列であることから，

$$\mathbf{x}_k^T \mathbf{F} \mathbf{x}_l = 0, \quad (k \neq l)$$

であり， V が直交行列であることから，

$$\mathbf{y}_k^T \mathbf{G} \mathbf{y}_l = 0, \quad (k \neq l)$$

が成立する．

スコア x_2 と y_2 とは，平均と分散の制約に加え，上の直交条件を満たすもののうち，(6) を最大にするものである．

2. 対応分析（パターン分類の数量化）

2.2 対応分析の利用例

表1のデータについて，具体的な数値例を示す．
各行の和と平方根は，表に示されているように

行和	847	550	64	229	597	421	878	240	642	147
平方根	29.10	23.45	8.00	15.13	24.43	20.52	29.63	15.49	25.34	12.12

となる．

一方，各列の和と平方根は

列和	850	909	932	974	950
平方根	29.15	30.15	30.53	31.21	30.82

となる．

2. 対応分析 (パターン分類の数量化)

33

$$Q = F^{-1/2} P G^{-1/2} \text{ (各要素は } q_{ij} = p_{ij} / \sqrt{f_i g_j} \text{ である .)}$$

	[,1]	[,2]	[,3]	[,4]	[,5]
[1,]	0.1355	0.2211	0.2465	0.2059	0.1472
[2,]	0.1931	0.1824	0.1788	0.1435	0.0775
[3,]	0.0986	0.0415	0.0655	0.0320	0.0284
[4,]	0.2017	0.1425	0.0866	0.0466	0.0279
[5,]	0.0786	0.1222	0.1810	0.2151	0.2018
[6,]	0.2491	0.1568	0.1421	0.0921	0.0427
[7,]	0.0359	0.0895	0.1271	0.2422	0.4686
[8,]	0.0310	0.0792	0.0782	0.1613	0.1550
[9,]	0.2667	0.2173	0.1629	0.1252	0.0691
[10,]	0.1245	0.1122	0.0729	0.0740	0.0187

2. 対応分析（パターン分類の数量化）

34

行列の特異値分解 $Q = UDV^T$ の，対角行列 D の成分は，

1.0, 0.4625, 0.1457, 0.0514, 0.0456

2番以降が対応分析によって得られる行と列の相関係数になる．第1次元の行と列の相関は $r_{xy} = 0.4625$

2. 対応分析 (パターン分類の数量化)

35

また, U の第 1 列 (制約を満たさない部分) および第 2 , 3 列は , 次のようになる .

	[,1]	[,2]	[,3]
[1,]	0.4284	0.0421	-0.5933
[2,]	0.3452	0.2109	-0.1372
[3,]	0.1178	0.1059	0.1473
[4,]	0.2228	0.2978	0.3540
[5,]	0.3597	-0.1980	-0.3481
[6,]	0.3020	0.3344	0.2738
[7,]	0.4362	-0.7131	0.4678
[8,]	0.2280	-0.2115	-0.1556
[9,]	0.3730	0.3428	0.2068
[10,]	0.1785	0.1794	0.0119

2. 対応分析 (パターン分類の数量化)

36

これに左から $F^{-1/2}$ (相対頻度から求めたもの) をかけた $F^{-1/2}U$ の第 1 ~ 3 列は次のようになる .

	[,1]	[,2]	[,3]
[1,]	1	0.0982	-1.3850
[2,]	1	0.6108	-0.3973
[3,]	1	0.8994	1.2510
[4,]	1	1.3369	1.5892
[5,]	1	-0.5506	-0.9677
[6,]	1	1.1070	0.9065
[7,]	1	-1.6348	1.0725
[8,]	1	-0.9272	-0.6825
[9,]	1	0.9190	0.5545
[10,]	1	1.0050	0.0667

2. 対応分析 (パターン分類の数量化)

一方 V の第 1 列と第 2 列は , 次になる .

	[,1]	[,2]	[,3]
[1,]	0.4292	0.5757	0.6072
[2,]	0.4438	0.2958	-0.1607
[3,]	0.4494	0.1172	-0.4990
[4,]	0.4594	-0.2298	-0.3793
[5,]	0.4537	-0.7173	0.4611

2. 対応分析 (パターン分類の数量化)

38

同様に, これに $G^{-1/2}$ を左からかけると

	[,1]	[,2]	[,3]
[1,]	1	1.3414	1.4149
[2,]	1	0.6666	-0.3621
[3,]	1	0.2607	-1.1103
[4,]	1	-0.5001	-0.8257
[5,]	1	-1.5810	1.0163

となる. この第2列と第3列とが, y_1 および y_2 である.

列についてのスコアは(当然予想されることだが), AFQT
の層別の順を反映したものになっている.

2. 対応分析（パターン分類の数量化）

39

表1の行の順を x_1 の大きさにしたがって昇順に置き換え
ると、次のような表が得られる。

2. 対応分析（パターン分類の数量化）

表 2: 対応分析のスコアによって行を並べ替えた表

	AFQT1	AFQT2	AFQT3	AFQT4	AFQT5
Prof/Tech	31	80	115	224	428
Sales	14	37	37	78	74
Mgr	56	90	135	164	152
Clerical	115	194	219	187	132
Craft	132	129	128	105	56
Farm	23	10	16	8	7
Service	197	166	126	99	54
Transpt	44	41	27	28	7
Operatv	149	97	89	59	27
Labor	89	65	40	22	13

2. 対応分析（パターン分類の数量化）

41

専門技術職(Prof/Tech)ではAFQTの高成績群が相対的に多く、一方、熟練工(Operativ)や労働(Labor)（職種の訳語は大津によるので、公的なものではない。）では、低成績群が多い。中ほどの事務職(Clerical)では、AFQT中成績群が多く、対応分析によるスコアは、すくなくとも表の見かけ上の構造を簡潔に把握する手助けとなることがわかる。

2. 対応分析 (パターン分類の数量化)

42

図1は、横軸に1次元目のスコア、縦軸に2次元目のスコアをとり、 $(d_1 \times x_1, d_2 \times x_2)$ および $(d_1 \times y_1, d_2 \times y_2)$ を、同一のグラフ上にプロットしたものである (文献やソフトウェアによっては d_k で重み付けないものもある。ここでのプロットはRのMASSライブラリ (Venables & Ripley, 1999) を用いている。

2. 対応分析 (パターン分類の数量化)

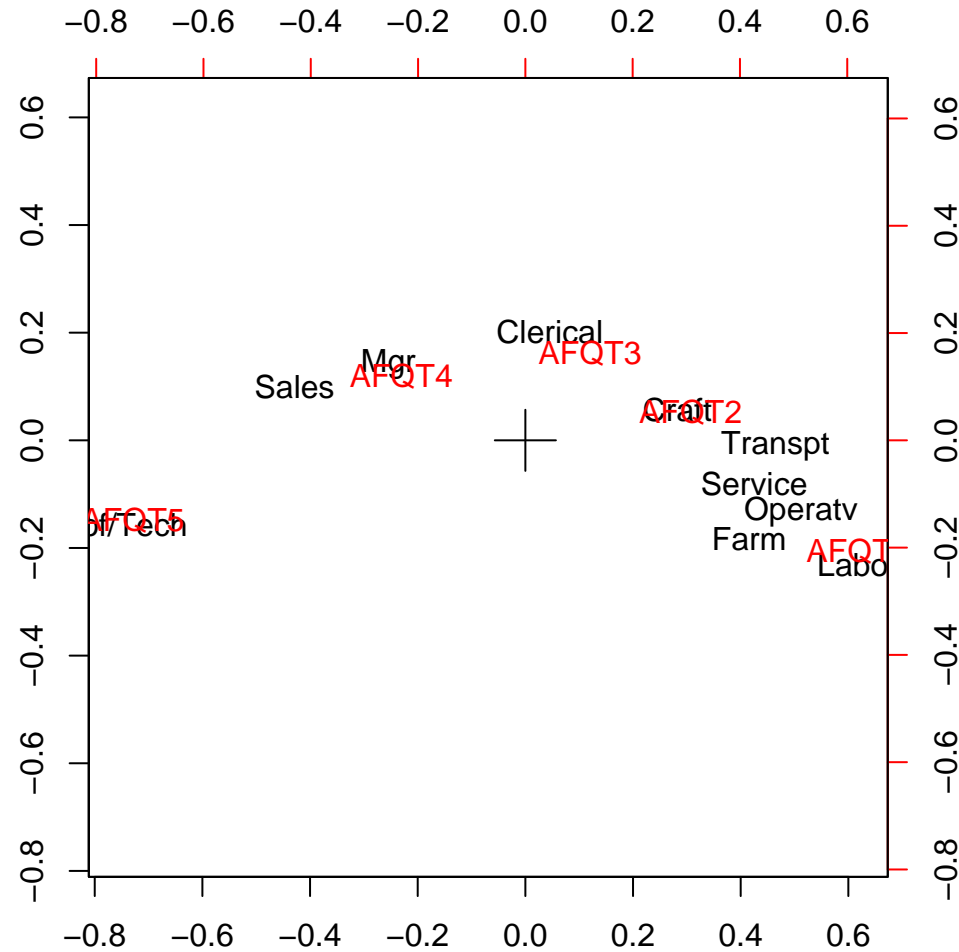


図 1: 対応分析の解，横軸：第 1 次元，縦軸：第 2 次元

2. 対応分析（パターン分類の数量化）

Jaques Bertinの重み付き置換マトリックス．

図2は、表1の表示である．

各行の幅は、対応する周辺分布の大きさに比例．第 i 行の幅（ここでは縦方向の幅）は、 p_{i+} の大きさに比例している．

各セルに対応する棒グラフの長さ（ここでは横方向の長さ）は、各セルの値を周辺度数で割った値 p_{ij}/p_{i+} を表している．小さな長方形の面積が、セルの値を表す．

各列の破線は、 p_{+j} に相当する大きさ（行と列が独立である場合の期待値）を表す．

2. 対応分析 (パターン分類の数量化)

NLSY79 クロスセクショナル (置換後)

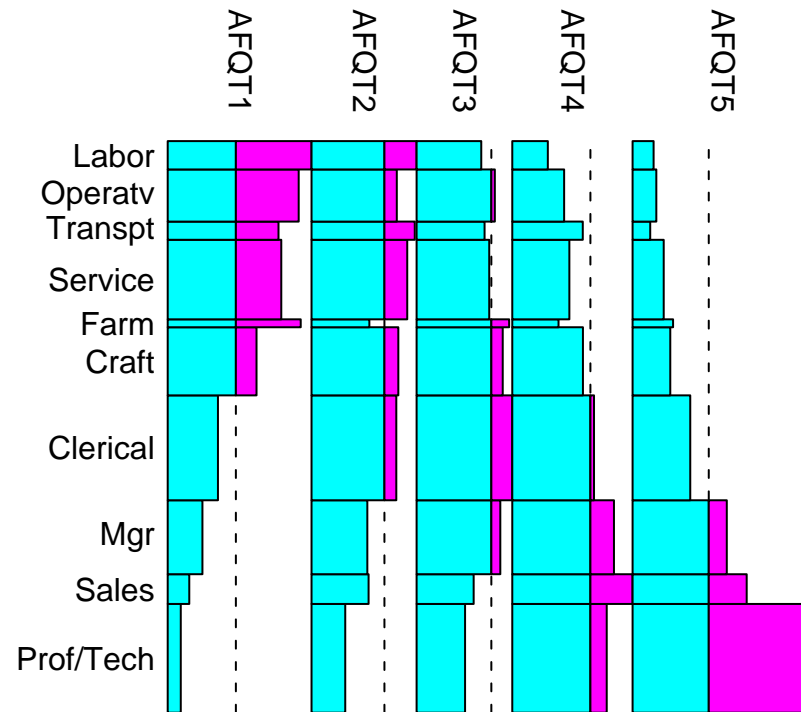


図 2: 重み付き置換マトリックス (第 1 次元スコアによる並べ替え)

2. 対応分析（パターン分類の数量化）

46

2.3 馬蹄形問題

対応分析は，計算が単純で便利な方法であるが，実用上いくつか注意すべき点がある．ひとつは，カテゴリ間の関係を離散変数の間の相関によって捉えるため，周辺度数が偏っている場合（特にカテゴリの水準の数が小さい場合）には，変数間の関係を適切に表現できない場合がありうることである．

さらに，多次元の解を利用する場合には，次のような特徴（馬蹄形問題）にも注意する必要がある．

2. 対応分析（パターン分類の数量化）

47

表3に示す共分散を持つ4変量正規分布を想定する． σ_{13}, σ_{24} の値は，表4に示す2つのケースを考える．次にこの分布に基づく乱数を生成し，次のような手順で2重分割表を作成する．

2. 対応分析（パターン分類の数量化）

48

1. 表3に示す共分散行列に従う4変量正規乱数を5000個生成する．
2. 各変数の値を大きさの順に並べ替え，大きさの順に5つのカテゴリーに分類する．この際，カテゴリー区分は各々の頻度が等しくなるようにする．
3. 次に X_1 と X_2 を離散化して得られたカテゴリーのクロスカテゴリーを求める．同様にして X_3 と X_4 からクロスカテゴリーを求める．各々25のカテゴリを持つ2つの変数 X と Y が得られる(図3)．
4. 離散化された X と Y から 25×25 の2重分割表を求める．

2. 対応分析（パターン分類の数量化）

表 3: 多変量正規分布の共分散行列

	X_1	X_2	X_3	X_4
X_1	1.0			
X_2	0.0	1.0		
X_3	σ_{31}	0.0	1.0	
X_4	0.0	σ_{42}	0.0	1.0

表 4: テストデータの共分散の値

	σ_{13}	σ_{24}
Case 1	0.8	0.5
Case 2	0.8	0.4

2. 対応分析 (パターン分類の数量化)

50

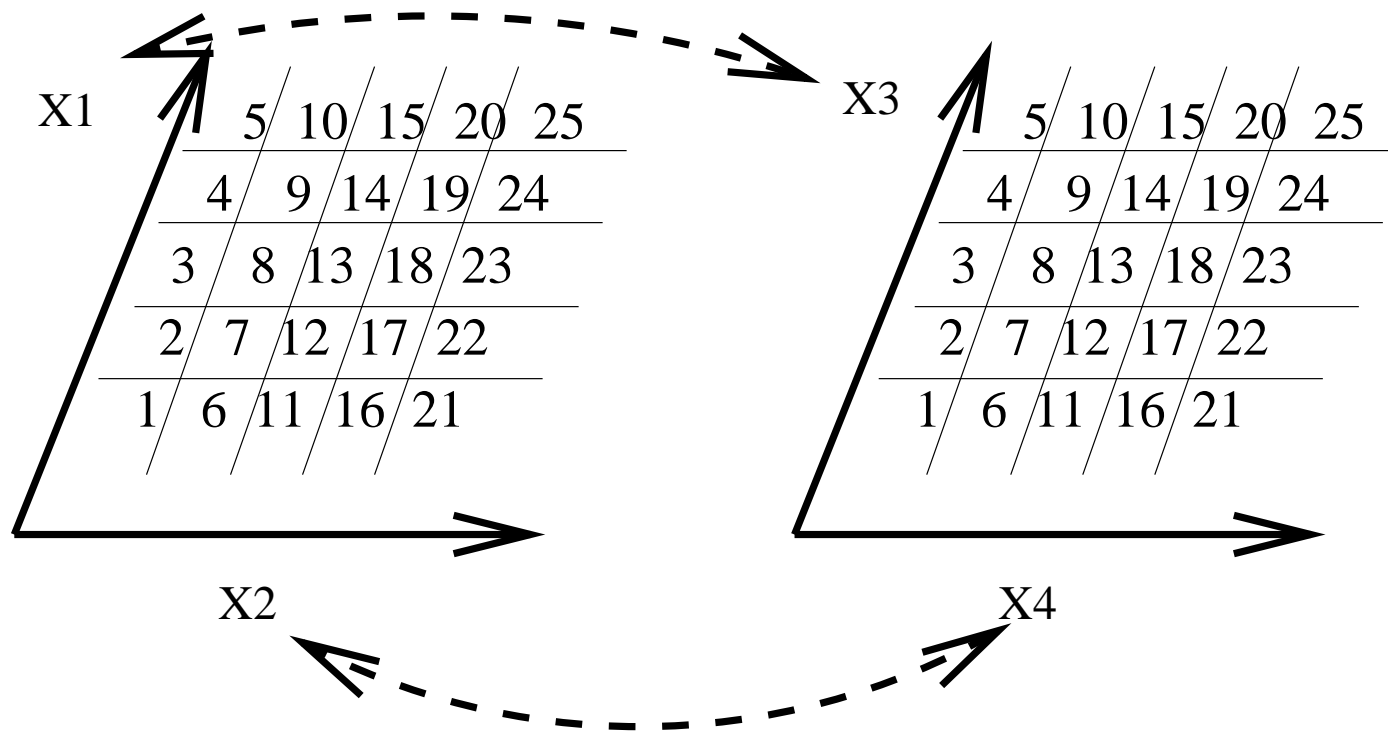


図 3: 2重分割表の作成

2. 対応分析（パターン分類の数量化）

表 5: 対応分析によって推定された相関係数 (2重分割表)

	r_1	r_2	r_3	r_4
Case 1	0.75	0.46	0.42	0.36
Case 2	0.74	0.41	0.36	0.25

この手続きによって生成された2重分割表に対応分析を適用する．表5は，対応分析を適用して得られた相関係数(特異値)を大きさの順に上位4個まで示したものである．ここでは自明な解($r_0 = 1$)は除いてある．

2. 対応分析 (パターン分類の数量化)

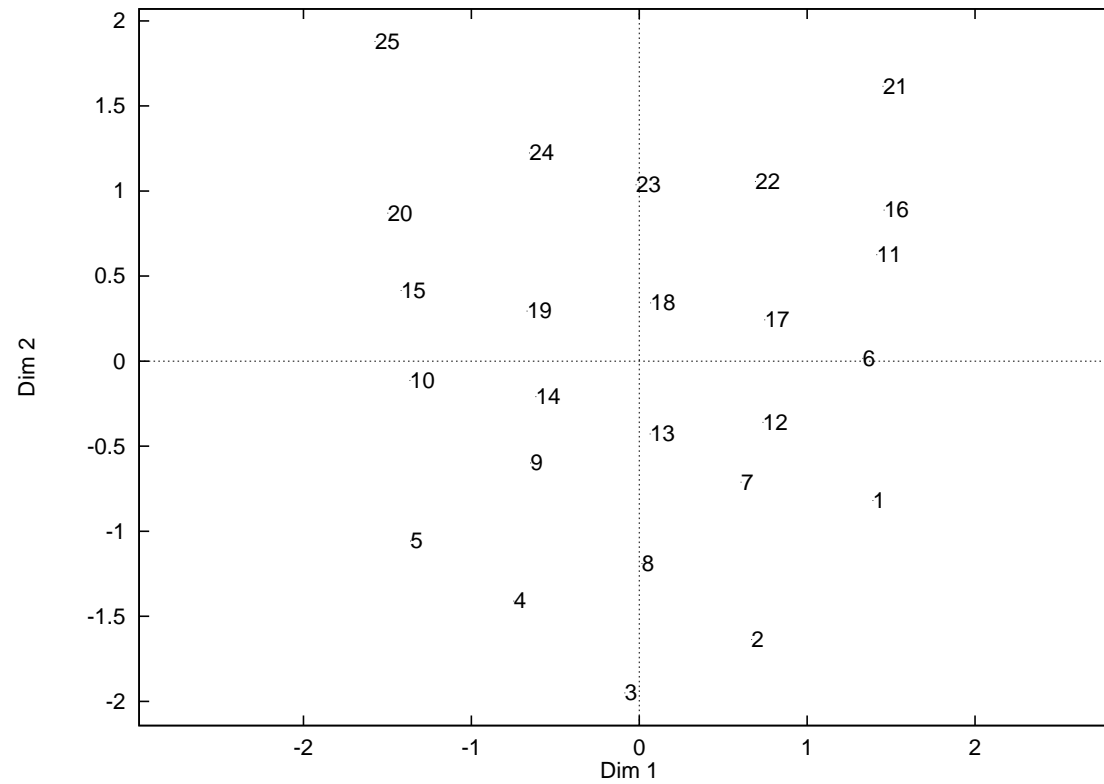


図 4: 対応分析によるスコア (Case1, 第1次元と第2次元)

2. 対応分析 (パターン分類の数量化)

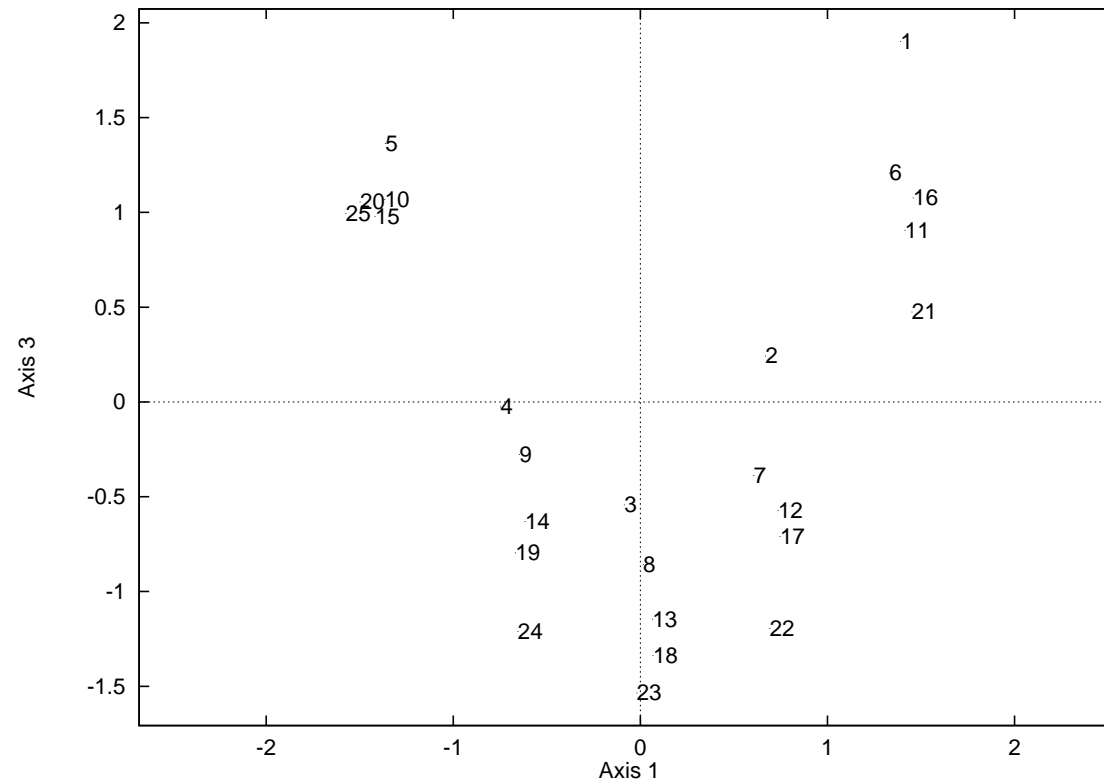


図 5: 対応分析によるスコア (Case1, 第1次元と第3次元)

2. 対応分析（パターン分類の数量化）

54

Case1の計算結果を示す図4においては，データを生成するために設定した格子状の構造が第1次元と第2次元によって再現されている．

Case 1で得られた相関係数は第1次元については0.75であり，また第2次元については0.46となっている．表4に示した正規分布の相関係数よりは幾分小さめではあるが，カテゴリーの背後に設定した分布の構造を推定することに成功している．

2. 対応分析 (パターン分類の数量化)

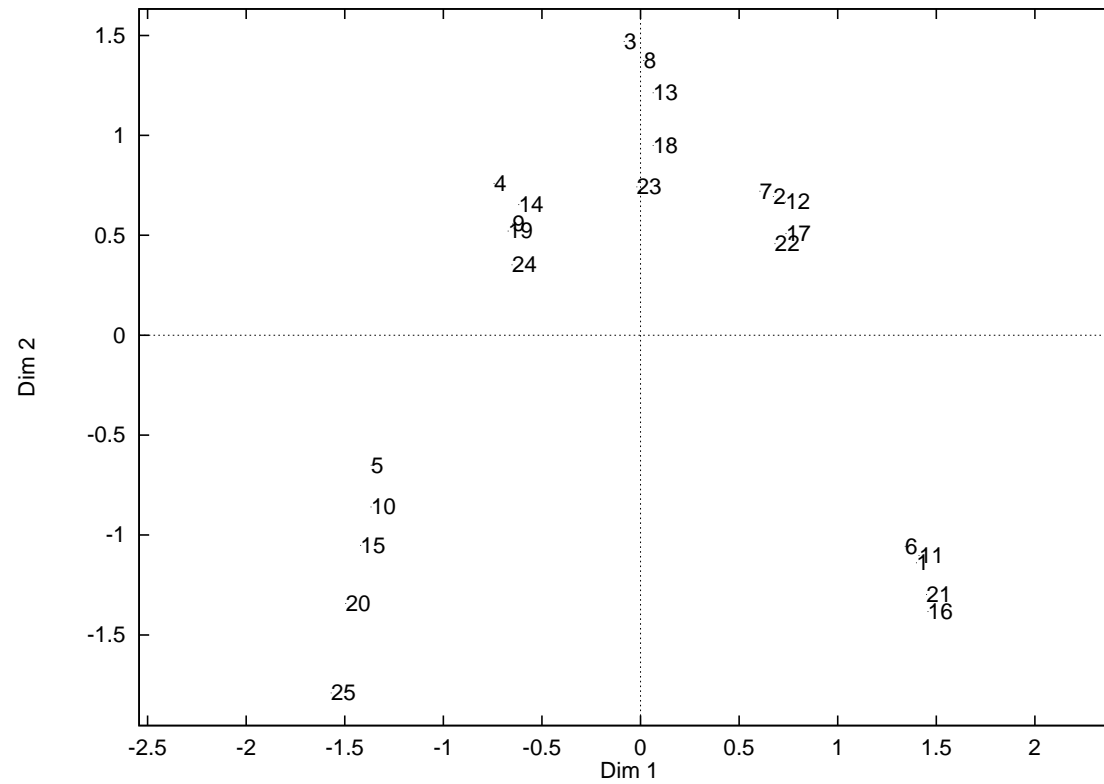


図 6: 対応分析によるスコア (Case2, 第1次元と第2次元)

2. 対応分析 (パターン分類の数量化)

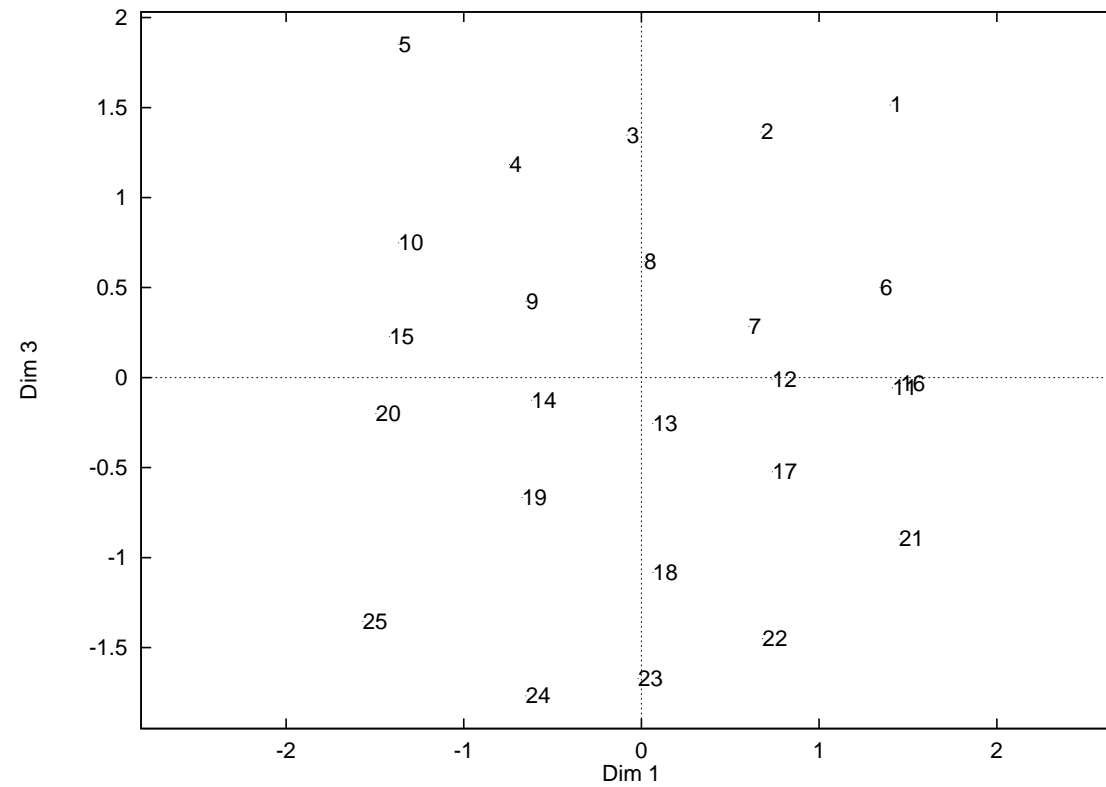


図 7: 対応分析によるスコア (Case2, 第1次元と第3次元)

2. 対応分析（パターン分類の数量化）

57

一方，図6はCase 2における Y のカテゴリースコアである．第2次元のスコアは第1次元のスコアの2次関数になっている．ここには示さなかったが，格子状の構造は第1次元と第3次元のスコアを合わせ見ることによって初めて確認できる．

Case1とCase2の違いは X_2 と X_4 の相関が各々0.5と0.4であることのみであり，わずかな相関構造の違いが大きな結果の違いをもたらしている．

2. 対応分析（パターン分類の数量化）

58

理論的な検討によれば，潜在的な構造が多変量正規分布であり，データ十分多くあり，カテゴリーの分割が小さければ，最大の相関（特異値）が r_1 であるとき，データの実質的な内容に関わらず，ほぼ r_1^2 の相関を持つ解が現れる．

より小さな相関を持つ微妙な成分である場合には，1次元目の2次関数成分が解として得られ，データの内容を表す成分を発見することができない．

ただし，実際の計算においては，カテゴリー数がそれほど多くはないため，2次関数やより高次の関数成分を表す解の相関は， r_1^2 や r_1^3 よりも，かなり小さくなる．

3. 尺度最適化を伴う主成分分析

59

3 尺度最適化を伴う主成分分析

対応分析：2つの離散変数の関係から最適尺度を求める

⇒ r_{XY} の最大化

⇒ $\text{var}(X + Y)$ (主成分の分散) の最大化

3. 尺度最適化を伴う主成分分析

60

多重対応分析(MCA, 数量化3類): 多肢選択の多変量データを対象とする.

X_j は $1, \dots, k_j$ の離散値をとる.

⇒ k_j 個の $0 - 1$ 値をとるダミー変量を対応させる.

$n \times (\sum_j k_j)$ の $0 - 1$ 値のテーブルに対応分析を適用.

⇒ サンプル(被験者)のスコアとアイテムの選択肢のスコアが得られる.

⇒ 第1次元は X_1, X_2, \dots, X_p の尺度最適化による主成分

著名な利用例: 統計数理研究所「日本人の国民性調査」
ピエール・ブルデュー「ディスタンクション」

3. 尺度最適化を伴う主成分分析

61

対応分析と類似の馬蹄形問題が生じる．選択肢のスコアはしばしば不安定．特に第2次元目以降の解に問題がありがち．

- 多重対応分析（数量化3類）の不安定性を回避したい
- 離散値と連続値の混在するデータを分析したい
- 離散値の順序関係も考慮したい

尺度最適化を伴う主成分分析

(SAS:PRINQUAL, Gifi(Heiser et al.) HOMALS, OSMOD)

離散変数については最適スコアを与えつつ主成分分析モデルを当てはめる．

3. 尺度最適化を伴う主成分分析

3.1 データとモデル

分析対象となるデータのサンプル件数（被験者数）を N とし，変数（調査項目）の個数を J とする．第 j 変数が間隔尺度を持つ数値データである場合には，第 i 番目のサンプル（被験者）の第 j 変数の値を $v_i(j)$ と表す．また，第 j 変数が順序尺度または名義尺度のカテゴリー値をとる場合には，それらのカテゴリーを

$$\{C_{j1}, C_{j2}, \dots, C_{jk_j}\}$$

とする．特に，第 j 変数が順序尺度である場合には

$$\{C_{j1} \prec C_{j2} \prec \dots \prec C_{jk_j}\} \quad (8)$$

なる順序関係が定められているものとする．

3. 尺度最適化を伴う主成分分析

63

第 $1 \sim J_1$ 変数は間隔尺度の数値変数

第 $J_1 + 1 \sim J_1 + J_2$ 変数は順序尺度のカテゴリー値を持つ変数

第 $J_1 + J_2 + 1 \sim J_1 + J_2 + J_3$ 変数は名義尺度のカテゴリー値を持つ変数

ここで, $J_1 + J_2 + J_3 = J$ とする.

3. 尺度最適化を伴う主成分分析

64

また，カテゴリー値をとる第 j 変数における第 i サンプルの反応を次のように表記する．

$$\delta_i(jk) = \begin{cases} 1 & (\text{第 } i \text{ サンプルが第 } k \text{ カテゴリーに反応}) \\ 0 & (\text{それ以外}) \end{cases}$$

以降の表記を簡単にするため，反応頻度について次のような記号を導入する．式中， \bar{x} などのように横棒がついているものは， N 件のサンプルについての平均を表す．

$$n_{jk} = \sum_{i=1}^N \delta_i(jk) = N \times \overline{\delta(jk)}$$

$$n_{jk,lm} = \sum_{i=1}^N \delta_i(jk)\delta_i(lm) = N \times \overline{\delta(jk)\delta(lm)}$$

3. 尺度最適化を伴う主成分分析

65

ここで考慮の対象としているカテゴリカルデータは択一式のものである．従って次の等式が成立する．

$$\sum_{k=1}^{k_j} \delta_i(jk) = 1, \quad (i = 1, \dots, N; j = J_1 + 1, \dots, J)$$

ある変数のカテゴリが順序尺度に基づくものである場合には，カテゴリースコアがその順序制約に従うことにする．

3. 尺度最適化を伴う主成分分析

カテゴリースコアのベクトルを

$$w = (w'_{J_1+1}, w'_{J_1+2}, \dots, w'_J)'$$

とする．ここで， $w_j = (w_{j1}, \dots, w_{jk_j})'$ であり，第 j 変数が順序尺度である場合には(1)の順序に基づいて，

$$w_{j1} \leq w_{j2} \leq \dots \leq w_{jk_j}, \quad (j = J_1 + 1, \dots, J_1 + J_2) \quad (9)$$

という制約をスコアに課す．

3. 尺度最適化を伴う主成分分析

67

何らかの基準に基づいて $(J_2 + J_3)$ 個のカテゴリー変数に対応するスコアベクトルが得られたと仮定する．

次の式によって $(J_2 + J_3)$ 個の数値変数が定義される．

$$s_{ij} = \sum_{k=1}^{k_j} w_{jk} \delta_i(jk), \quad (j = J_1 + 1, \dots, J; i = 1, \dots, N) \quad (10)$$

これにより，第 $1 \sim J_1$ 番目の数値変数と併せて J 個の数値変数が得られる．

カテゴリースコアの設定が適切なものであれば，これら J 個の変数の相関行列に基づく主成分分析によってデータの内的構造の簡潔な表現を得ることが期待できる．

3. 尺度最適化を伴う主成分分析

68

どのような基準によってスコアベクトルを定めたら良いか.

ここでは, Saito & Otsu (1983,1988) に採用されている基準 (分散和最大化基準) と一般化分散最小化基準の2つについて検討する.

Maximizing Total Variance (MTV) , と Minimizing Generalized Variance (MGV)

3. 尺度最適化を伴う主成分分析

69

カテゴリースコアを一意に決定するために, (10)の J 個の変数は, 全て平均0・分散1に正規化することにする. 数値変数についてはこれらの制約は

$$\frac{1}{N} \sum_{i=1}^N v_i(j) = 0 \text{ および } \frac{1}{N} \sum_{i=1}^N v_i(j)^2 = 1, \quad (j = 1, \dots, J_1) \quad (11)$$

となる. また, カテゴリー変数については

$$\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^{k_j} w_{jk} \delta_i(jk) = 0 \text{ および } \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^{k_j} w_{jk}^2 \delta_i(jk) = 1, \quad (j = J_1 + 1, \dots, J) \quad (12)$$

と表現される.

3. 尺度最適化を伴う主成分分析

数値変数については

$$s_{ij} = v_i(j), \quad (j = 1, \dots, J_1)$$

とし、またカテゴリー変数については (10) によって正規化条件 (11), (12) を満たす $J_2 + J_3$ 個の変数を定める。

これによって得られる $N \times J$ の行列を $S = (s_{ij})$ とする。 S はカテゴリースコア w の関数であるので、 $S(w)$ とかける。また、このとき J 個の変数の相関行列は

$$R(w) = \frac{1}{N} S'(w) S(w) \quad (13)$$

となる。この相関行列の J 個の固有値を

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_J \geq 0$$

とする。

3. 尺度最適化を伴う主成分分析

71

MTV基準は，あらかじめデータを縮約表現する分析次元数 m を定めておき，次の値をモデル適合度の指標とする．

$$\theta = \lambda_1 + \lambda_2 + \cdots + \lambda_m \quad (14)$$

これを変数の平均と分散についての制約の下で最大化する．

実対称行列の性質より， m 個の J 次元ベクトル $X = (\boldsymbol{x}_1, \dots, \boldsymbol{x}_m)$ が正規直交系であるとの制約のもとで

$$\theta = \max_{X, w} \sum_{i=1}^m \boldsymbol{x}_i' R(\boldsymbol{w}) \boldsymbol{x}_i \quad (15)$$

がいえる．

MTV基準の最大化は w と X についての最大化によって実現される．

3. 尺度最適化を伴う主成分分析

MGV 基準の場合には

$$\det R = \prod_{j=1}^J \lambda_j$$

の最小化を行なえばよい．これは

$$\eta = - \sum_{j=1}^J \log \lambda_j \quad (16)$$

の最大化に等しい．後者のほうが数値的な取扱いが容易であるので，これをMGV基準の目的関数とすることにする．

それぞれの変数 X_j が X_j を除く他の変数によって回帰したときの重相関係数の2乗 R_j^2 が大きければ，MGVの適合度が高まる．

相関行列の行列式 $\det R$ (一般化分散) が最も大きくなるのは，各変数が完全に無相関であるとき，すなわち R

3. 尺度最適化を伴う主成分分析

73

が単位行列である場合である．また， $\det R$ は R が半正定値であることから常に非負であり， 0 となるのは S の列ベクトルが一次従属となる場合である．

3. 尺度最適化を伴う主成分分析

74

推定されたカテゴリースコアベクトルを \hat{w} とし, それに基づく変数間の相関行列を $\hat{R} = R(\hat{w})$ とする. また, \hat{R} の m 個の固有ベクトルを列ベクトルとして持つ $J \times m$ の行列を $\hat{X} = (\hat{x}_{jl})$ とする. これ以降, アイテム(変数) j の布置ベクトルとは, \hat{X} の第 j 行目(長さ m のベクトル)を意味することとする. また, 第 j アイテムの第 k カテゴリーの布置とは, アイテム j の布置ベクトルに, スカラー \hat{w}_{jk} を掛けたものとする. 名義尺度アイテム(変数)について, その布置とカテゴリースコアは符号不定性を持つが, カテゴリー布置は符号不定性を持たない. さらに, $S(\hat{w})$ の第 i 行目のベクトルに \hat{X} を右から掛けて得られる m 次元ベクトルによって, サンプル i の布置を定義する. 因子分析の用語を用いると, アイテム布置は因子負荷量に相当し, サンプル布置は因子得点に相当する.

3. 尺度最適化を伴う主成分分析

75

3.2 OSMOD 分析例: NLSY79 の概要

NLSY79 (BLS による調査)

被験者は1957年から1964年のあいだに生まれた男女であり全部で12,686名。

クロスセクショナルサンプル(6,111名)

サプリメンタルサンプル(5,295名)

軍隊サンプル(1,280名)

このうち分析対象 7,025名 (検討対象の変数の欠測を原則として除く)

3. 尺度最適化を伴う主成分分析

76

NLSY79, 19変数 :

Age 年齢, WithMan 14歳時の同居男性, WithWo
14歳時の同居女性,
MEdu 母親の学歴, FEdu 父親の学歴, Sib 同胞の数(1979),
Reli 被験者の宗教, RelAtt 宗教活動への参加, Race 人種,
Sex 性別, REdu 被験者の学歴(1992), EmpStat 雇用状態
(1996),
LkJob 仕事への好感度(1996), Delinq 非行, Jail 矯正施設
への入所,
Resi 居住地域, Pov 貧困基準に該当, AFQT AFQT得点,
FmInc 世帯収入(1979,1980)

3. 尺度最適化を伴う主成分分析

3.3 NLSY79 分析結果

MTV 基準 ($q = 3$)

表 6: OSMOD による固有値

	1	2	3	4	5	6	7
固有値 (MTV $q = 3$)	3.72	1.72	1.67	1.49	1.29	1.14	1.03
累積寄与率	0.20	0.29	0.37	0.45	0.52	0.58	0.63
固有値 (MGV)	3.63	1.74	1.68	1.48	1.40	1.16	1.06
累積寄与率	0.19	0.28	0.37	0.45	0.52	0.58	0.64

3. 尺度最適化を伴う主成分分析

表 7: 回転後の主成分の分散共分散/相関行列

	V_1	V_2	V_3	V_4	V_5
V_1	2.91	-0.32	0.19	-0.07	-0.18
V_2	-0.78	2.00	-0.16	0.02	0.26
V_3	0.45	-0.31	1.85	-0.03	-0.10
V_4	-0.15	0.04	-0.05	1.50	0.03
V_5	-0.39	0.46	-0.17	0.05	1.62

(上三角部分は相関係数)

3. 尺度最適化を伴う主成分分析

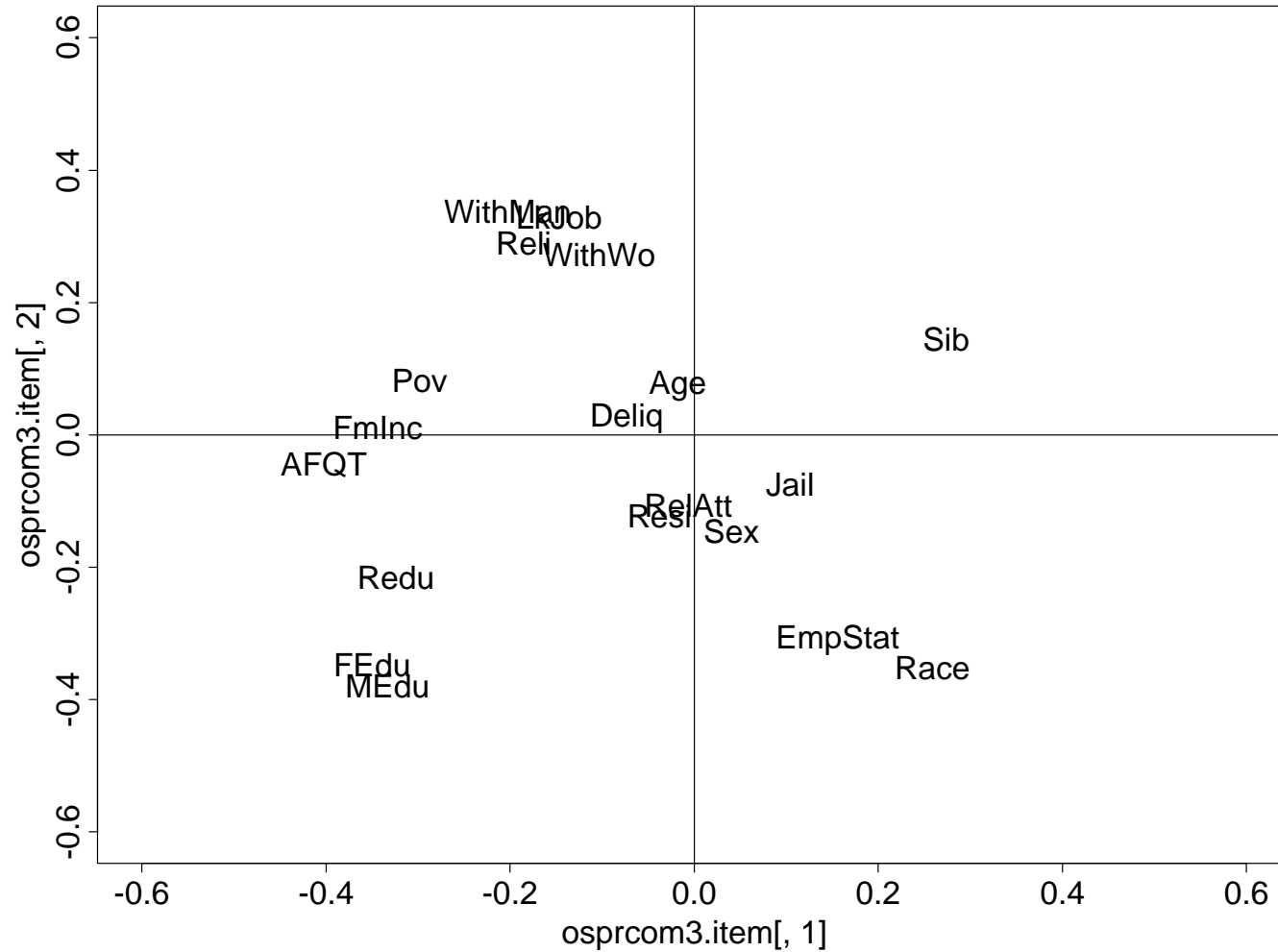


図 8: 回転前固有ベクトル (因子負荷量) のプロット

横軸 1 次元目, 縦軸 2 次元目, MTV, $m = 3$

3. 尺度最適化を伴う主成分分析

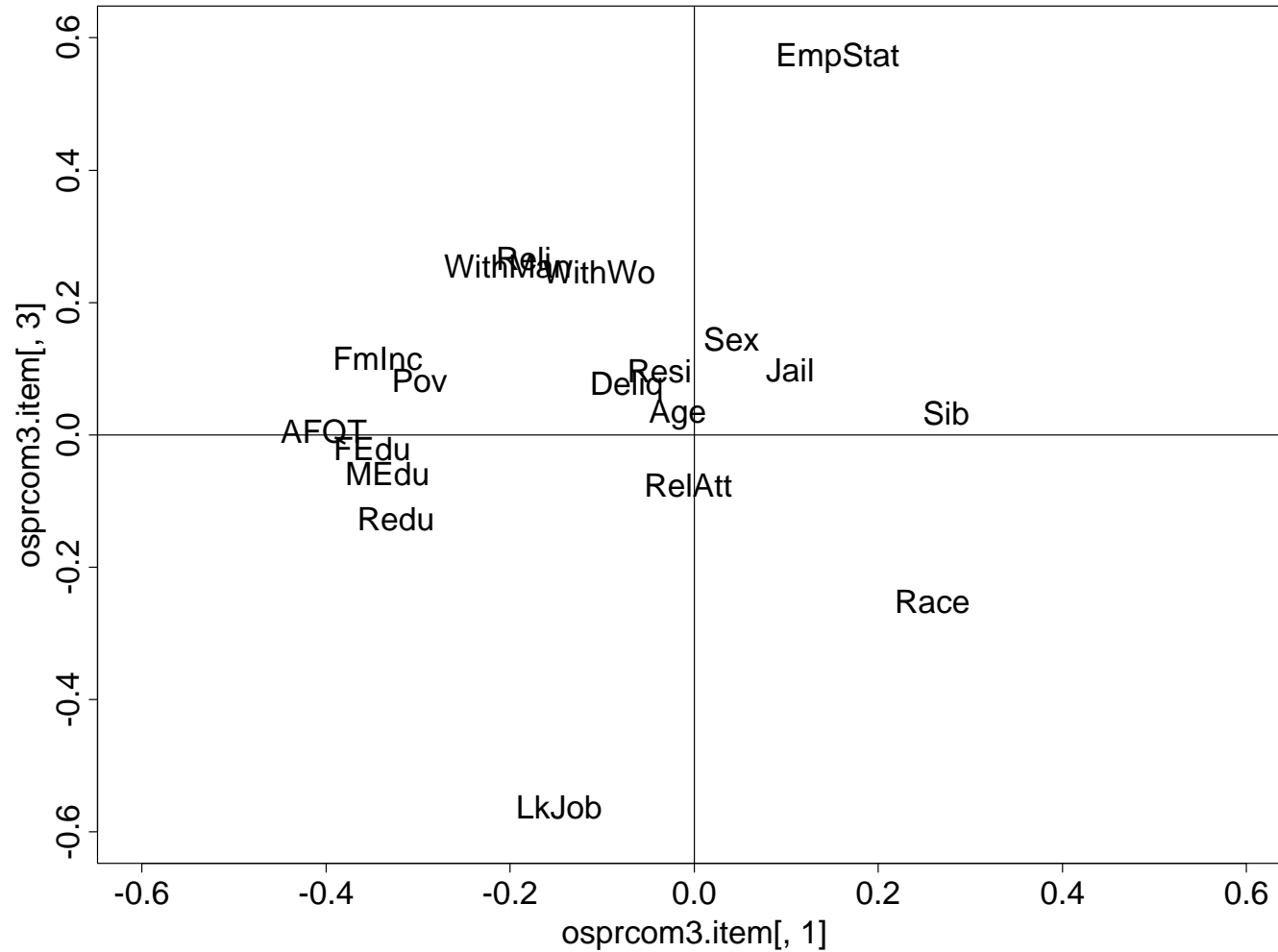


図 9: 回転前固有ベクトル (因子負荷量) のプロット

横軸 1 次元目, 縦軸 3 次元目, MTV, $m = 3$

3. 尺度最適化を伴う主成分分析

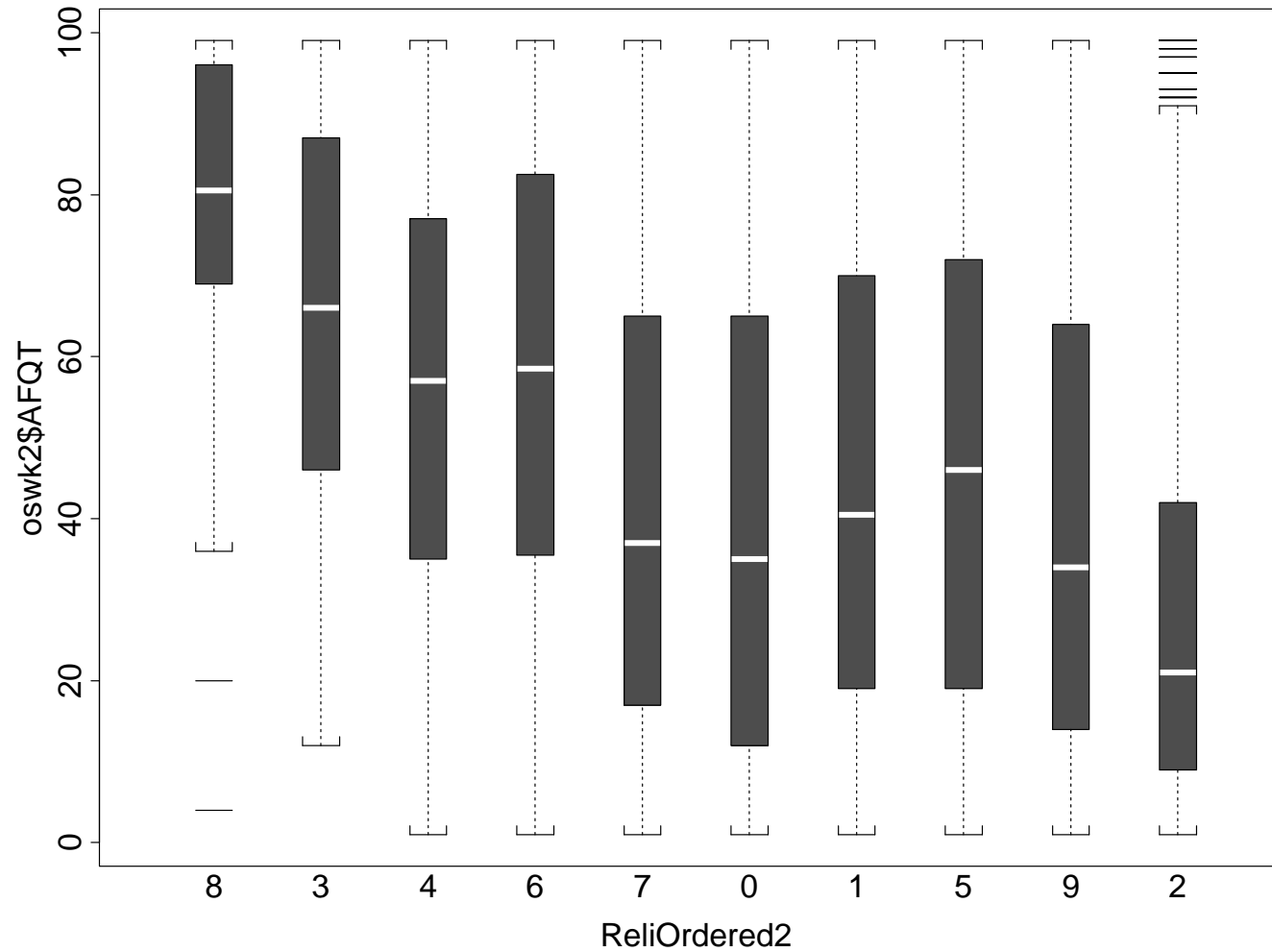


図 10: 宗教 (スコア順) と AFQT

3. 尺度最適化を伴う主成分分析

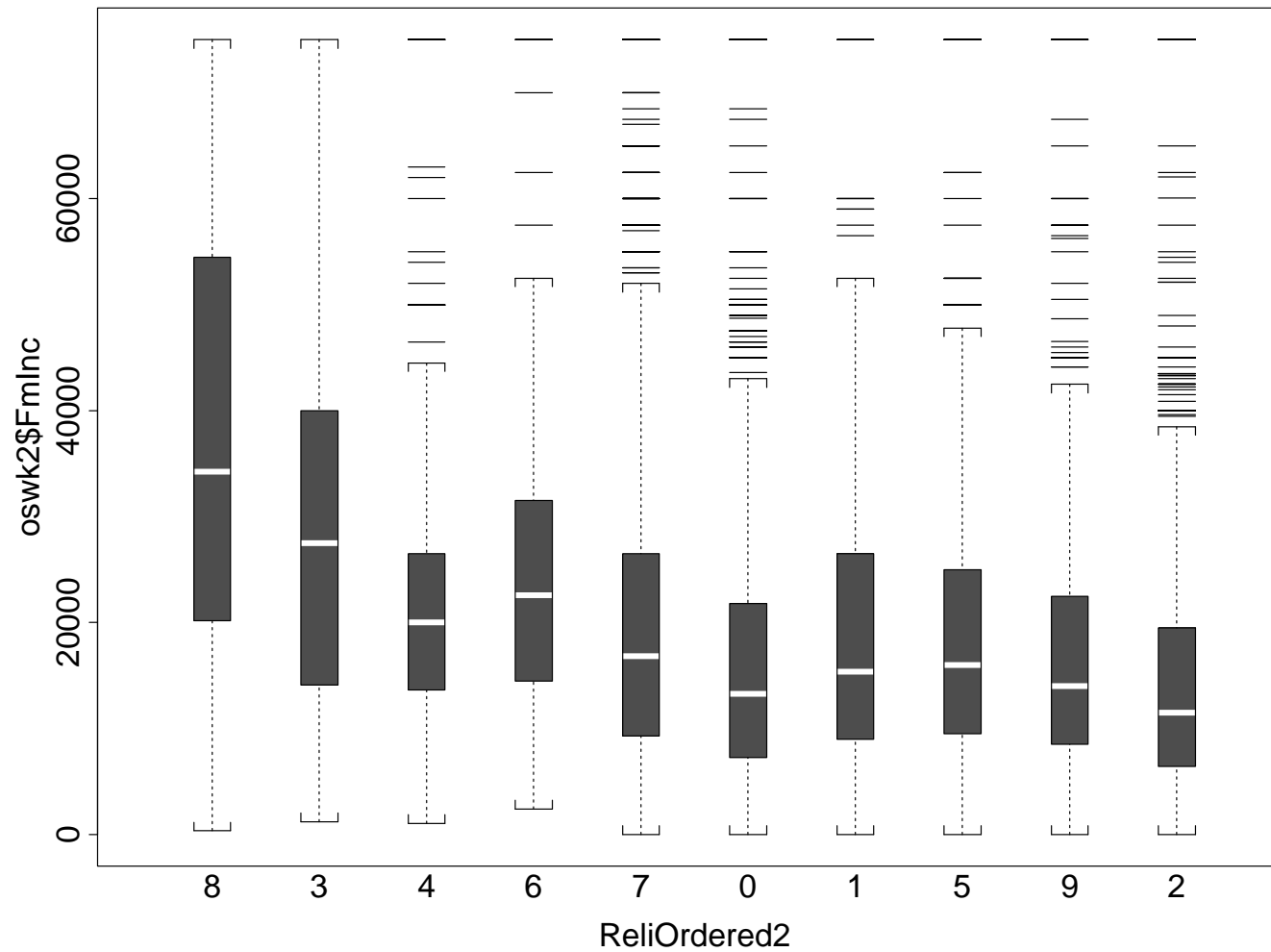


図 11: 宗教 (スコア順) と世帯収入

3. 尺度最適化を伴う主成分分析

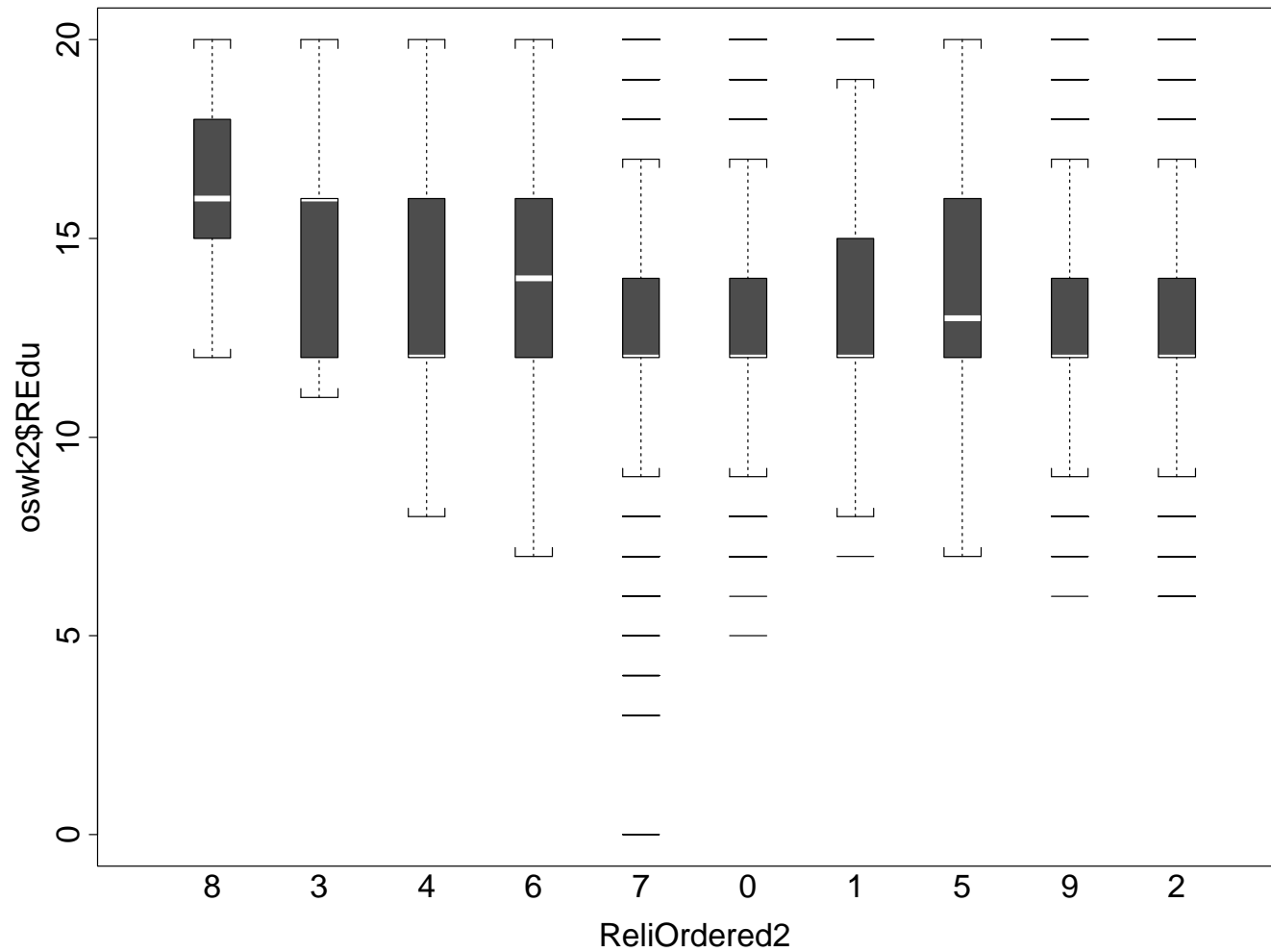


図 12: 宗教 (スコア順) と被験者の教育年数

3. 尺度最適化を伴う主成分分析

84

プログラムについて

- 対応分析

SAS Corresp, SPSS Categories ,
R & Splus (`corresp` in Modern Applied Statistics with
S (MASS). Venables & Ripley, Springer)

- 尺度最適化をともなう主成分分析

SAS PRINQUAL, SPSS Category?, OSMOD (free)

- 大津の OSMOD サイト

<http://www.rd.dnc.ac.jp/~otsu/osmod>

ソースコード + Winバイナリ, MinGW の g77 でコン
パイル.

3. 尺度最適化を伴う主成分分析

85

OSMOD プログラム

1. Fortran 77 で記述（スタイルは古い）、無償、無保証
2. 連続変数の非線形変換の機能はなし
3. 変数の重み付けは一律に固定（相関行列を分析）
4. Fortran の実数データとしてデータを読み込み（文字型データの入力不可）
5. コマンドウィンドウで実行。パラメータを標準入力から読み込み、パラメータで指定されたファイルにバイナリの計算結果を保持。標準出力にログ出力。
6. R & SPlus インタフェースも用意（データはファイル渡し）。

3. 尺度最適化を伴う主成分分析

86

OSMODの実行:

コマンドウィンドウで以下を実行

```
osmod < param.a1 > a1output.txt
```

```
osmodpr < param.b1 > b1output.txt
```

3. 尺度最適化を伴う主成分分析

param.a1の中身 LARGE

```
*-----  
infile      'sample.data'  
format      '*'  
outfile     sample.out  
nitem       5  
* kfunc 1:MTV 2:MGV  
kfunc       2  
* level 0: nominal, 1:ordinal, 2:interval(numeric)  
level       1 1  
level       2 1  
level       3 1  
level       4 0  
level       5 2
```

3. 尺度最適化を伴う主成分分析

88

```
title 'DATA ANDREW HERTZBERG TB63.2'
```

```
* nr:分析次元 MGV のときは推定には関係しない
```

```
nr      2
```

```
* nvec: 固有値計算の反復のための固有ベクトル数
```

```
nvec    3
```

```
* limcat カテゴリ区分総数の上限
```

```
limcat  200
```

```
*-----
```


4. まとめなど

89

4 まとめなど

- 名義尺度変数を含むデータの分析が面白い。
- 数値・順序尺度の場合は、多少の単調変換で結果が大きく影響されるのはまれ。ただし2値はかなり情報が落ちる。
- 数量化3類で分析していたデータを osmod などで再分析すると、新たな発見があるかも。
- 精密なモデル評価には耐えないが、それなりには有用。欠測の扱いは実用上課題。
- Bertin のグラフィックスは商用システムでも欲しい。
Mondrian システム (Java based GUI data analysis)