

カテゴリカルデータの単純2段主成分分析

カテゴリ型と連続型の同時主成分分析

○岡本安晴

日本女子大学人間社会学部

1 はじめに

質問項目がカテゴリカル変数である場合は多くみられるが、その分析法として村上（1999）は非計量的主成分分析を提案している。この分析法では、 k 番目（ $k = 1, \dots, m$ ）のカテゴリカル変数に対するデータを指標行列 G_k 、 G_k に対する重み（数量化）行列を V_k 、 $G_k V_k$ の因子モデルにおける得点行列を F 、負荷行列を A_k で表すとき、次式

$$\sum_{k=1}^m \|G_k V_k - F A'_k\|^2 \quad (1)$$

を制約条件

$$V'_k D_k V_k = I_{q_k}, \quad F' \mathbf{1}_n = \mathbf{0}_r, \quad n^{-1} F' F = I_r$$

のもとで最小化する解が求められる。しかし、(1) の最小化では、 $G_k V_k$ が G_k の情報を最大限に反映しているという保証はない。また、 G_k を数量化するときの次元数、すなわち V_k の列数 q_k は、分析者があらかじめ指定するもので、必ずしも最大数 $p_k - 1$ （ p_k は G_k の列数）が設定されるとは限らない。ここでは、 G_k の情報を最大限反映した数量化に基づいた主成分分析を行うため、次の単純2段主成分分析を考える。

2 単純2段主成分分析

各カテゴリカル変数 G_k （ $n \times p_k$ ）の数量化を行い（第1段）、その後、数量化されたカテゴリカルデータと他の数量データがあればそれを加えて主成分分析を行う。

(1) カテゴリカルデータの数量化

数量化のための重み ω_k を $G_k \omega_k$ の分散が最大（極大）になるように選ぶ。

$$\mathbf{1}'_n G_k = \mathbf{f}'_k = \left(f_1^{(k)} \quad \dots \quad f_{p_k}^{(k)} \right), \quad H_k = G_k - \frac{1}{n} \mathbf{1}_n \mathbf{f}'_k$$

とおき、 $\omega'_k \omega_k = 1$ の制約のもとで $\text{tr}((H_k \omega_k)'(H_k \omega_k))$ を最大にする ω_k を求めるため、目的関数を

$$\begin{aligned} Q(v_k) &= \text{tr}((H_k v_k)'(H_k v_k)) - \lambda(v'_k v_k - 1) \\ &= \text{tr}(v'_k H'_k H_k v_k) - \lambda(v'_k v_k - 1) \end{aligned}$$

とおく。

$$\frac{\partial}{\partial v_k} Q(v_k) = 2H'_k H_k v_k - 2\lambda v_k$$

なので、求める ω_k は(2)式を満たす。

$$H'_k H_k \omega_k = \lambda \omega_k \quad (2)$$

(2)式に対して、例えば $\lambda = 0$ 、 $\omega_k = \frac{1}{\sqrt{p_k}} \mathbf{1}_{pk}$ は1組の解である。

(2)式を満たす2組の解、 λ_i と ω_i 、および λ_j と ω_j 、に対して

$$(H_k \omega_i)' (H_k \omega_j) = \lambda_j \omega'_i \omega_j = 0$$

が成り立つ。したがって、

$$\text{Cov}(H_k \omega_i, H_k \omega_j) = 0$$

である。また、

$$(H_k \omega_i)' (H_k \omega_i) = \lambda_j \omega'_i \omega_i = \lambda_j$$

より、

$$\text{Var}(H_k \omega_i) = \frac{1}{n} \lambda_i$$

となる。これより、標準化得点

$$z_i = H_k \omega_i / \sqrt{\lambda_i / n}$$

を得る。

$$G_k^z = (z_1 \quad \cdots \quad z_{pk-1})$$

とおく。

(2) 主成分分析

G_k すなわち H_k の数量化による標準得点 G_k^z に対して

$$G^z = (G_1^z \quad \cdots \quad G_m^z)$$

とおく。

(2a) データがカテゴリカル変量 G_k のみからなる場合

G^z の主成分分析を行う。

(2b) データがカテゴリカル変量 G_k と数量データ X とからなる場合

$(G^z \quad X)$ に対して主成分分析を行う。

3 分析例

10人の作曲家についての属性、生国 (a:ドイツ・オーストリア、b:フランス、c:ロシア)、作品の性格 (a:絶対音楽、b:中間的、c:表題・劇音楽)、時期 (a:19世紀前半、b:19世紀後半、c:20世紀前半)、による評価データ (表1; 村上、1999、Table 1 から転載) を分析してみる。各カテゴリカル変量 (属性) に対する固有値 λ_i と固有ベクトル ω_i を求めると表2のようになった。いずれの変量も3カテゴリであるので、0でない固有値の数は2個である。各固有値と固有ベクトルから算出される標準化得点 z_i の名前を表2の対応する固有値と固有ベクトルの欄の下に示した。

表1 10人の作曲家の評価 (村上 (1999) より)

作曲家	生国	作品の性格	時期
Beethoven	a	a	a
⋮	⋮	⋮	⋮
Stravinsky	c	b	c

これらの数量化変数の標準化得点から構成されるデータ $G^z = (G_1^z \ \dots \ G_m^z)$ に対して主成分分析を行ったときの特異値を表3に示す。2次元の主成分分析を行ったときの各主成分と数量化変数との相関係数を座標値として数量化変数の布置を

表2 各カテゴリカル変量の固有値と固有ベクトル

	生国		作品		時期	
固有値	3.6	3	4	2.4	3.6	3
固有ベクトル	0.816	0.000	0.000	0.816	-0.408	0.707
	-0.408	0.707	0.707	-0.408	0.816	0.000
	-0.408	-0.707	-0.707	-0.408	-0.408	-0.707
数量化変数名	Country-0	Country-1	Charactr-0	Charactr-1	Period-0	Period-1

表3 主成分分析における特異値

1.393	1.252	1.057	0.895	0.614	0.444
-------	-------	-------	-------	-------	-------

表したものを図1に示す。第1主成分 (Comp.1) は、数量化変数 Country-0、Charactr-1 と正の相関があり、period-0 と負の相関があることがわかる。第2主成分 (Comp.2) は、数量化変数 Charactr-0 と正の相関があり、Country-1、period-1 と負の相関がある。すなわち、第1主成分は、「ドイツ・オーストリア → フランス、ロシア」および「絶対音楽 → 中間、表題・劇」と正の相関があり、「19c 後半 → 19c 前半、20c 前半」と負の相関がある。第2主成分は、「中間的 → 絶対音楽 → 表題・劇」と正の相関があり、「フランス → ドイツ・オーストリア → ロシア」および「19c 前半 → 19c 後半 → 20c 前半」と負の相関がある。各作曲家の主成分値を座標として布置を表したものが図2である。作曲家名の先頭の英字は評定値を表す。例えば、「acb: Wagner」は、「生国」、「作品」、「時期」の評価がそれぞれ a、c、b であることを表している。小円が作曲家の位置であり、長方形の枠内に主成分と対応する数量化変数が示されている。図1の表わす主成分に対応して作曲家がグループ分けされていることがわかる。

主成分分析は、データの散布を低次元空間に正射影するものと考えられる (Okamoto, 2005)。このとき、主成分は特殊な座標系に対する座標値である。図1および図2では、2つの主成分と数量化変数との対応は明快であるが、データによっては正射影された空間において解釈の容易な座標系を探すことも必要になると考えられる。すなわち、射影空間における回転である。回転の一般式は例えば Okamoto (2005) で取り上げられているが、具体的な回転方法は因子分析法で開発された種々の回転を用いることができる。

引用文献

村上隆 (1999). 科学研究費補助金 (09610114) 研究成果報告書

Okamoto, Y. (2005). 日本女子大学人間社会学部紀要 (<http://ci.nii.ac.jp/naid/110006223025>)

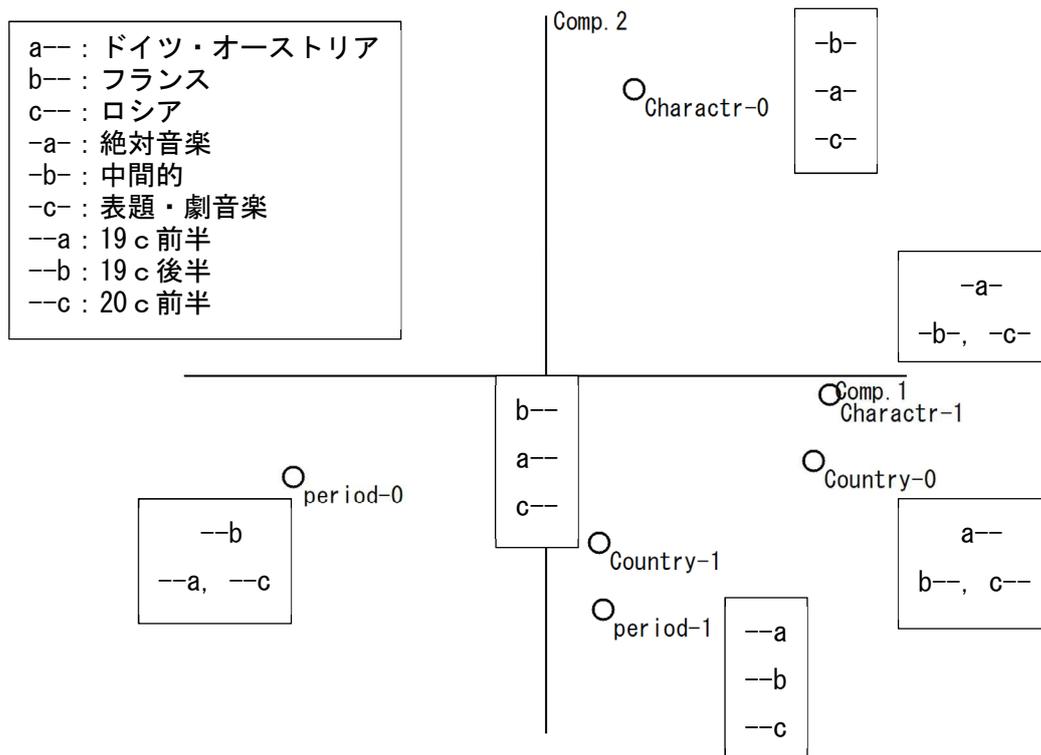


図1 数量化変数と主成分との相関係数を座標値とする布置

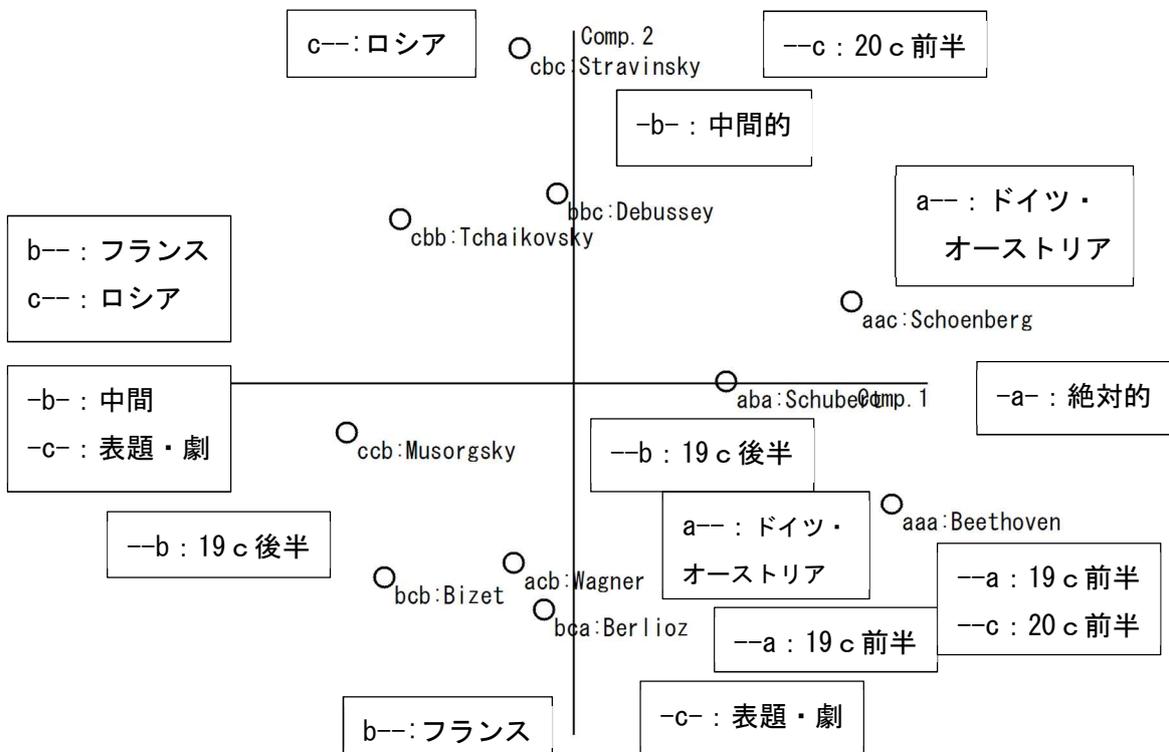


図2 主成分値を座標とする作曲家の布置