

不均一構造・不完全情報下での主成分分析

01300450	日本大学	*高橋 磐郎	TAKAHASHI Iwaro
01011500	日本大学	大澤 慶吉	OSAWA Keikichi
01404360	日本大学	西澤 一友	NISHIZAWA Kazutomo
02991410	日本大学	王 克義	Keyi WANG

§ 1. はじめに

主成分分析[1]のデータは、一般に p 種の項目 A_1, \dots, A_p に関する N 個の個体を対象とする観測データから構成される。個体 n の A_j に関する観測値を a_{nj} とすると、 $n=1 \sim N$, $j=1 \sim p$ のすべての組み合わせに対して a_{nj} が揃っていない表1のような場合がしばしば起こる。

n	A ₁	A ₂	A ₃	A ₄	p=4
1	a ₁₁	a ₁₂		a ₁₄	
2	a ₂₁	a ₂₂	a ₂₃	a ₂₄	
3	a ₃₁		a ₃₃	a ₃₄	
4		a ₄₂	a ₄₃	a ₄₄	
5	a ₅₁	a ₅₂	a ₅₃		
6	a ₆₁		a ₆₃	a ₆₄	
7		a ₇₂		a ₇₄	
8	a ₈₁	a ₈₂	a ₈₃	a ₈₄	

N=8

データの欠けている部分を欠測部、データが存在する部分を実測部ということにする。

このような事が起こる原因として2通りの場合が考えられる。1つは欠測部の (n, j) の組合せは、構造上もともと起こり得ないものである場合、もう1つは構造的にはすべてのデータを採り得るが、情報が不完全のため、たまたまデータが欠損して、欠測部が生じた場合である。前者を不均一構造の場合、後者を不完全情報の場合と呼ぶことにしよう。両者は見かけ上は表1のように同じ姿をしているが、データ構造としては異なるものと考えねばならない。

前回[2]では(標題は「不完全情報」

となっているが)じつは不均一構造の場合を取り扱ったことになるが、今回は主として不完全情報の場合について考察しよう。

§ 2. 不完全情報の主成分分析のデータモデル

不完全情報の欠測部は観測すればデータを採れる可能性はあるが、たまたま欠落したものである。そこで対象個体 n 、項目 j が欠測部である場合そのデータ a_{nj} のデータ構造を、ここでは次のように仮定することにする。

$$(1) \quad a_{nj} = \beta p_n + \gamma q_j$$

つまり(1)式によって欠測部のデータを予測しようということである。ここで p_n は個体 n の個体平均、 q_j は項目 A_j の項目平均(→(3))で、 α 、 β 、 γ は未知パラメータで、以下の方法で推定されるものである。

たとえば、個体を学生と考え、項目を授業科目と考え、選択科目方式の問題を不完全情報とみなすと、学生 n がもし科目 A_j を受けたとすると、どの程度の得点をとるか予測する式が(1)である。

(1)は最も簡単なモデルであるが、問題によって各項目あるいは各個体の特殊性について何らかの情報があれば、それを考慮したモデルを作ることもし得る。(1)を用いて a_{nj} を予測すれば、あとは完全情報の場合と同様に(第一)主成分のスコアを求め、それによって、各個体の総合評価を行うことができる。第一主成分が個体の総合評価として最も適切であることは[2]にも触れた。第二主成分以下については今回は特に考慮しない。

§ 3. 不完全情報の主成分分析の原則

個体 n 、項目 A_j に対し、欠測部では $\delta_{nj}=0$ 、 $\varepsilon_{nj}=1$ 、実測部では $\delta_{nj}=1$ 、 $\varepsilon_{nj}=0$ とすると、 $n=1\sim N$ 、 $j=1\sim p$ の全体に対して

(2) $\alpha_{nj} = \delta_{nj} a_{nj} + \varepsilon_{nj} (\alpha + \beta p_n + \gamma q_j)$ と書ける。ここで

$$(3) \begin{cases} p_n = \frac{\sum_j \delta_{nj} a_{nj}}{\sum_j \delta_{nj}} \quad \dots \text{個体平均} \\ q_j = \frac{\sum_n \delta_{nj} a_{nj}}{\sum_n \delta_{nj}} \quad \dots \text{項目平均} \end{cases}$$

個体 n の第一主成分スコアを z_n 、項目のウェイトを u_1, \dots, u_p とすると

$$(4) z_n = \sum_j \alpha_{nj} u_j$$

となるが z_n の分散 $V_n = \sum (z_n - \bar{z})^2 / N$ 、 $\bar{z} = \sum z_n / N$ を $u_1^2 + \dots + u_p^2 = 1$ の下に最大にするように $u_1, \dots, u_p, \alpha, \beta, \gamma$ を決めようという考えが、ここで用いる原則である。この分散最大という原則は簡単ではあるが最も合理的な情報縮約を与えることは [2] にも述べたところである。

§ 4. 解析および計算方式

上記の原則に基づく解の候補は、ラグランジュ関数

$$L = Vz - \lambda (u_1^2 + \dots + u_p^2 - 1)$$

の停留点の中に求められる。つまり

$$(5) \partial L / \partial u_i = 0 \quad (i=1\sim p)$$

$$(6) \begin{cases} \partial L / \partial \alpha = 0, & \partial L / \partial \beta = 0, \\ \partial L / \partial \gamma = 0 \end{cases}$$

の解となるが、このうち (5), (6) からそれぞれ次の式 (7), (8) が得られる。

$$(7) Ru = \lambda u, \quad u^T = (u_1, \dots, u_p)$$

$$(8) \begin{cases} (EE)\alpha + (PE)\beta + (QE)\gamma + (AE) = 0 \\ (PE)\alpha + (PP)\beta + (QP)\gamma + (AP) = 0 \\ (QE)\alpha + (QP)\beta + (QQ)\gamma + (AQ) = 0 \end{cases}$$

ここで行列 R 、スカラー $(EE), (PE), \dots$ 等は次のように与えられる；

$$(9) R = [r_{ij}], \quad r_{ij} = \sum_n \dot{\alpha}_{ni} \dot{\alpha}_{nj} / N$$

$$(10) (EE) = \sum_i \sum_j \sum_n \dot{\varepsilon}_{ni} \dot{\varepsilon}_{nj} u_i u_j,$$

$$(PE) = \sum_i \sum_j \sum_n \dot{p}_{ni} \dot{\varepsilon}_{nj} u_i u_j,$$

$$(QE) = \sum_i \sum_j \sum_n \dot{q}_{ni} \dot{\varepsilon}_{nj} u_i u_j,$$

$$(AE) = \sum_i \sum_j \sum_n \dot{a}_{ni} \dot{\varepsilon}_{nj} u_i u_j, \dots$$

$$(11) \dot{\alpha}_{nj} = \dot{a}_{nj} + \dot{\varepsilon}_{nj} \alpha + \dot{p}_{nj} \beta + \dot{q}_{nj} \gamma$$

$$\dot{a}_{nj} = a_{nj} - \bar{a}_j \quad (\bar{a}_j = \sum_n \delta_{nj} a_{nj} / N)$$

$$\dot{\varepsilon}_{nj} = \varepsilon_{nj} - \bar{\varepsilon}_j \quad (\bar{\varepsilon}_j = \sum_n \varepsilon_{nj} / N)$$

$$\dot{p}_{nj} = \varepsilon_{nj} p_n - \bar{p}_j \quad (\bar{p}_j = \sum_n \varepsilon_{nj} p_j / N)$$

$$\dot{q}_{nj} = \varepsilon_{nj} q_n - \bar{q}_j \quad (\bar{q}_j = \sum_n \varepsilon_{nj} q_j / N)$$

(7) の R の中には、 α, β, γ が含まれているし、(8) の $(EE), (PE), \dots$ 等には u_1, \dots, u_p が含まれているので、(7), (8) を連立させて解くことは殆ど不可能だから、次のような逐次近似計算を用いる；

u_1, \dots, u_p の初期値として $u_1 = \dots = u_p = 1/\sqrt{p}$ を定め、(10) より $(EE), (PE), \dots$ 等を計算して (8) を解き α, β, γ を求める。これを用いて (11), (9) より R を求め (7) の主 (最大固有値に対する) 固有ベクトル u_1, \dots, u_p を求める。これを新たな初期値として以上の手順を繰り返す。

§ 5. 応用例および今後の課題

選択授業科目の問題が不完全情報に属するか不均一構造に属するかは微妙な問題であるが、ある大学のある学科についてこの方式を適用してみた。このような場合、第二主成分の問題をどのように扱うかまたどのように解釈すべきか等今後の問題として検討したい。

参考文献

- [1] 奥野忠一他：「多変量解析法」日科技連、1971
 [2] 高橋磐郎他：不完全情報下での主成分分析、1995年日本OR学会秋季研究発表会アブストラクト集、pp60-61(1995).