

# 双対平坦性に基づく正準相関分析と 混合分布を用いた項目反応理論

数理情報学専攻 48-196235

渡邊 元和

指導教員

駒木 文保 教授

## 1 はじめに

本論文は2部構成となっている。

第I部では、指数型分布族の持つ双対平坦性を利用した正準相関分析の手法を提案する。本研究では、指数型分布族の持つ双対平坦性を利用して与えられる角度を考え、この角度を最小化する合成ベクトルを求めるものとして、指数型分布族における正準相関分析を定めている。第II部では、項目反応理論と呼ばれるテスト理論に用いられる統計モデルを扱う。本研究では、被験者能力の分布に混合正規分布を用いたモデル [3] を、問題の選択肢が順序を考えられる多肢となっている場合に用いられる段階反応モデル [6] へと拡張し、パラメータ推定手法を与えている。

本発表では主に第I部を扱い、第II部は概要を紹介する。

## 2 正準相関分析

正準相関分析 [4,5] は、2つの多変量変数群が与えられたときに、変数群の間の相関構造を分析するために用いられる多変量解析手法である。

正準相関分析は、行列  $Y \in \mathbb{R}^{n \times k}$ ,  $X \in \mathbb{R}^{n \times \ell}$  が与えられたとき、それぞれの行列を  $\alpha \in \mathbb{R}^k$ ,  $\beta \in \mathbb{R}^\ell$  を用いた線形変換をすることで得られる2本のベクトル  $r := Y\alpha$ ,  $u := X\beta$  の相関係数を最大化するような  $\alpha, \beta$  を求める手法であり、以下の制約付き最適化問題

$$\begin{aligned} (\alpha, \beta) &:= \operatorname{argmax}_{\alpha, \beta} \beta^\top X^\top Y \alpha \\ \text{s.t. } &\|Y\alpha\| = \|X\beta\| = 1 \end{aligned}$$

を解くことで、その解が与えられる。これを幾何的に捉えると、正準相関分析は、各列ベクトルを基底ベクトルとする超平面

$$S_1 := \{r : r = Y\alpha\}, S_2 := \{u : u = X\beta\}$$

を考え、2本のベクトルのなす角度が最小となるような  $r \in S_1$ ,  $u \in S_2$  を求める手法であると解釈することができる。

## 3 指数型分布族の双対平坦性

指数型分布族は、確率変数  $y = (y_1, \dots, y_n)$  の確率密度 (質量) 関数が、パラメータ  $\theta = (\theta^1, \dots, \theta^n)$  を用いて、 $p(y|\theta)d\mu(y) = \exp(\theta^i y_i - \psi(\theta))d\mu(y)$  で与えられるような確率分布族のことである。

情報幾何 [1,2] では、分布全体の集合  $\mathcal{S}$  を座標系  $\theta$  を持つ  $n$  次元多様体と考え、フィッシャー情報行列を計量として用いる。ここで、e-接続、m-接続と呼ばれる双対な接続に対して、座標系  $\theta$  と、 $\eta_i = E_\theta[y_i]$  によって定められる座標系  $\eta = (\eta_1, \dots, \eta_n)$  がそれぞれアファイン座標系となり、このようにe-接続、m-接続に対するアファイン座標系が存在する空間はそれぞれe-平坦、m-平坦と呼ばれる。また、平坦な座標系における測地線は、その座標系で直線として表され、それぞれe-測地線、m-測地線と呼ばれる。

指数型分布族は、e-平坦かつm-平坦な空間であることより双対平坦空間と呼ばれる。双対平坦空間では、双対な座標系を  $\theta, \eta$  として、 $\theta$  座標系に対する計量を  $g_{ij}$ ,  $\eta$  座標系に対する計量を  $g^{ij}$  とすれば、 $\frac{\partial \eta_i}{\partial \theta^j} = g_{ij}$ ,  $\frac{\partial \theta^i}{\partial \eta_j} = g^{ij}$  が成立し、 $g_{ij}$  と  $g^{ij}$  は逆行列の関係にある。また、同じ接空間における  $\theta$  座標系の基底  $e_i$  と、 $\eta$  座標系の基底  $e^j$  に関して、双対直交性と呼ばれる関係  $\langle e_i, e^j \rangle = \delta_i^j$  が成り立っている。

## 4 指数型分布族に対する正準相関分析

指数型分布族のなす空間における、双対平坦性を利用した角度を定め、その角度を最小化する変数を求めるものとして指数型分布族における正準相関分析を与える。

### 4.1 指数型分布族における角度

指数型分布族上の3点  $p, q, r \in \mathcal{S}$  が与えられたとき、 $q$  を頂点として  $p$  と  $q$  を結ぶe-測地線と、 $p$  と  $q$  を結ぶm-測地線を利用して与えられる角度  $\phi$  ( $0 \leq \phi \leq \pi$ ) を以下のように定める。

$$\phi := \arccos \left[ \frac{\left\langle \frac{d}{dt} \tilde{p}(0), \frac{d}{dt} \tilde{r}(0) \right\rangle}{\left\| \frac{d}{dt} \tilde{p}(0) \right\| \left\| \frac{d}{dt} \tilde{r}(0) \right\|} \right]$$

ただし,  $\frac{d}{dt}\tilde{p}(0)$  は,  $p$  と  $q$  を結ぶ e-測地線  $\tilde{p}(t) = t\theta_p + (1-t)\theta_q$  の  $t=0$  での接ベクトル,  $\frac{d}{dt}\tilde{r}(0)$  は,  $r$  と  $q$  を結ぶ m-測地線  $\tilde{r}(t) = t\eta_r + (1-t)\eta_q$  の  $t=0$  での接ベクトルである. また,  $\theta_p, \theta_q$  はそれぞれ  $p, q$  の  $\theta$  座標系における座標,  $\eta_q, \eta_r$  はそれぞれ  $q, r$  の  $\eta$  座標系における座標を表している.

#### 4.2 指数型分布族に対する正準相関分析の定式化

指数型分布族を考えたときの正準相関分析を, 以下の手順で与える.

まず, 頂点  $v \in S$  を 1 点与え, e-平坦な部分空間  $e\text{-span}(X)$  と, m-平坦な部分空間  $m\text{-span}(Y)$  を

$$e\text{-span}(X) := \left\{ \theta : \theta = \sum_{i=1}^{\ell} (x_i - \theta_v) \beta^i + \theta_v \right\}$$

$$m\text{-span}(Y) := \left\{ \eta : \eta = \sum_{i=1}^k (y^i - \eta_v) \alpha_i + \eta_v \right\}$$

と定める, このとき,  $u \in e\text{-span}(X), r \in m\text{-span}(Y)$  として,  $u$  と  $v$  を結ぶ e-測地線と,  $r$  と  $v$  を結ぶ m-測地線を利用して与えられる角度を最小とするような,  $u, r$  を求めるものとして正準相関分析を与える.

#### 4.3 指数型分布族に対する正準相関分析の解

指数型分布族に対する正準相関分析の解となる  $\alpha, \beta$  は, 一般化固有値問題

$$\begin{pmatrix} O & \tilde{Y}^\top \tilde{X} \\ \tilde{X}^\top \tilde{Y} & O \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \lambda \begin{pmatrix} \tilde{Y}^\top g^{ij}(\eta_v) \tilde{Y} & O \\ O & \tilde{X}^\top g_{ij}(\theta_v) \tilde{X} \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$

の固有値が最大となるような固有ベクトルを選ぶことで得られる. ただし,  $\tilde{X}, \tilde{Y}$  は, それぞれ,

$$\tilde{X} := (x_1 - \theta_v, \dots, x_\ell - \theta_v)$$

$$\tilde{Y} := (y^1 - \eta_v, \dots, y^k - \eta_v)$$

を表し,  $\theta_v, \eta_v$  は, それぞれ  $v$  の  $\theta$  座標系での座標と  $\eta$  座標系での座標を表すものとする.

### 5 数値実験

指数型分布族に対する正準相関分析と, 従来の正準相関分析の比較をするために数値実験を行った.

指数型分布族として, 平均パラメータと分散パラメータを持つ正規分布を用いて,  $n = 2000, k = \{60, 100\}, \ell = \{20, 40, 60, 80, 100\}$  で実験を行った. 正準変数のスケージングは接ベクトルの長さを利用し

て与え, 評価指標は, 正準変数  $u, r$  に対する KL ダイバージェンス  $KL(r : u)$  の  $n$  変数での平均値を用いた.

表 1. KL ダイバージェンスの平均値 ( $\times 10^{-4}$ )

$\ell$	$k = 60$		$k = 100$	
	提案手法	従来手法	提案手法	従来手法
20	4.0370	4.0402	3.9986	4.0032
40	3.7753	3.7778	3.7343	3.7381
60	3.7020	3.7041	3.6616	3.6649
80	3.6211	3.6235	3.5848	3.5883
100	3.4893	3.5504	3.4479	3.5122

結果は表 1 のようになった. ただし, 提案手法は, 指数型分布族に対する正準相関分析, 従来手法は従来の正準相関分析をそれぞれ表す. 実験の結果, 提案手法のほうが, より変数間の KL ダイバージェンスが小さくなるような正準変数の組が得られることがわかった.

## 6 第 II 部: 混合分布を用いた項目反応理論

項目反応理論とは, テストやアンケートの回答結果から, 被験者の潜在的な能力と各項目の持つ特性とを同時に推定する統計モデルである. 項目反応理論では, 被験者能力の分布を正規分布に仮定してパラメータ推定を行うことが多いが, その仮定が妥当とはいえない場合があることが指摘されているため, より柔軟な被験者の能力の分布を考えた研究が行われている.

本研究では, 被験者の潜在能力に混合正規分布を用いたモデル [3] を, 項目が順序のある多肢選択の場合に用いられる段階反応モデル [6] の枠組みに拡張し, マルコフ連鎖モンテカルロ法によるパラメータ推定手法を与えた. また, 数値実験によりその性能を確かめた.

### 参考文献

- [1] S. Amari. *Differential-geometrical Methods in Statistics*. Lecture notes in statistics. Springer-Verlag, 1985.
- [2] S. Amari and H. Nagaoka. *Methods of information geometry*, Vol. 191. American Mathematical Soc., 2007.
- [3] F. B. Gonçalves, B. C. C. Dias, and T. M. Soares. Bayesian item response model: a generalized approach for the abilities' distribution using mixtures. *Journal of Statistical Computation and Simulation*, Vol. 88, No. 5, pp. 967–981, 2018.
- [4] H. Hotelling. Relations between two sets of variates. *Biometrika*, Vol. 28, No. 3–4, pp. 321–377, 1936.
- [5] C. Jordan. Essai sur la géométrie à  $n$  dimensions. *Bulletin de la Société Mathématique de France*, Vol. 3, pp. 103–174, 1875.
- [6] F. Samejima. Estimation of latent ability using a response pattern of graded scores. *Psychometrika monograph supplement*, 1969.