

スピアマンランク行列による主成分分析のロバストネス

千葉大・融合理工学府 渡邊 宏大
千葉大・理学研究院 内藤 貫太

はじめに: ロバストな主成分分析については、これまで多くの手法が提案されている。本発表では、Marden (1999) の Rank から派生したスピアマンランク行列に基づく主成分分析のロバスト性について報告する。Han and Liu (2018) で議論されている Kendall's tau との比較についても考察する。

設定と定義: $\mathbf{y} \in \mathbb{R}^d$ の Spatial sign を

$$S(\mathbf{y}) = \begin{cases} \frac{\mathbf{y}}{\|\mathbf{y}\|_2} & , \mathbf{y} \neq \mathbf{0} \\ \mathbf{0} & , \mathbf{y} = \mathbf{0} \end{cases}$$

で定義する。 $\mathbf{X}, \mathbf{Y}, \mathbf{Z}, \mathbf{X}_1, \dots, \mathbf{X}_n$ を有限な分散共分散行列 Σ を持つ d 次元確率分布 F にしたがう互いに独立な確率ベクトルとする。Spatial sign に基づくロバストな主成分分析の手法が、Marden (1999), Han and Liu (2018) で議論されている。

母集団での Kendall's tau matrix $\mathbf{K} \in \mathbb{R}^{d \times d}$ は

$$\mathbf{K} = E_{\mathbf{X}, \mathbf{Y}} [S(\mathbf{X} - \mathbf{Y})S(\mathbf{X} - \mathbf{Y})^T]$$

で定義され、その推定量は

$$\widehat{\mathbf{K}} = \frac{2}{n(n-1)} \sum_{j < i} S(\mathbf{X}_i - \mathbf{X}_j)S(\mathbf{X}_i - \mathbf{X}_j)^T$$

で定義される。一方、母集団でのスピアマンランク行列 $\Sigma_R \in \mathbb{R}^{d \times d}$ は

$$\Sigma_R = E_{\mathbf{X}, \mathbf{Y}, \mathbf{Z}} [S(\mathbf{X} - \mathbf{Y})S(\mathbf{X} - \mathbf{Z})^T] = E_{\mathbf{X}} [E_{\mathbf{Y}} [S(\mathbf{X} - \mathbf{Y})] E_{\mathbf{Z}} [S(\mathbf{X} - \mathbf{Z})]^T]$$

で定義され、その推定量は

$$\widehat{\Sigma}_R = \frac{1}{n(n-1)(n-2)} \sum_{i=1}^n \sum_{j \neq i, k \neq i, j \neq k} S(\mathbf{X}_i - \mathbf{X}_j)S(\mathbf{X}_i - \mathbf{X}_k)^T$$

で定義される。

問題: 重要な事実として、 F が楕円分布のとき、 Σ , \mathbf{K} および Σ_R は同じ直交行列で対角化されることが知られている (Marden (1999), Han and Liu (2018) を参照)。Marden (1999) においても、 Σ_R の推定量 $\widetilde{\Sigma}_R$ によるロバストな主成分分析が議論されているが、その $\widetilde{\Sigma}_R$ は

$$\widetilde{\Sigma}_R = \frac{1}{n} \widehat{\mathbf{K}} + \left(1 - \frac{2}{n}\right) \widehat{\Sigma}_R$$

と分解される. $n \rightarrow \infty$ のとき, $\widehat{\mathbf{K}}$ は \mathbf{K} に, $\widehat{\Sigma}_R$ は Σ_R に確率収束することが示される. このことから, $\widetilde{\Sigma}_R$ と $\widehat{\mathbf{K}}$ の比較は, 漸近的には $\widehat{\Sigma}_R$ と $\widehat{\mathbf{K}}$ の比較となる. 我々の問題は特に, $\widehat{\Sigma}_R$ と $\widehat{\mathbf{K}}$ はどちらがロバストな主成分分析を提供するのか? ということになる.

影響関数: ロバストネスの指標として影響関数を考える (影響関数については, Hampel et al. (1986) を参照). F を楕円分布とする. 比較の対象として, \mathbf{K} と Σ_R の最大固有値に関する影響関数および固有ベクトルに関する影響関数を与える.

\mathbf{K} は直交行列 $\Gamma = [\gamma_1 \cdots \gamma_d]$ と対角行列 $A = \text{diag}(a_1, \dots, a_d)$, $a_1 > \cdots > a_d > 0$ によって $\mathbf{K} = \Gamma A \Gamma^T$ とスペクトル分解されるとする. 上で述べた重要な事実から, Σ_R も同じ直交行列によって $\Gamma \Lambda \Gamma^T$, $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$ とスペクトル分解できる. ただし, 対応する固有値については, $\lambda_1 > \cdots > \lambda_d > 0$ となるかは定かではない. したがって, γ_1 が属している固有値を λ_M , $M \in \{1, \dots, d\}$, Σ_R の最大固有値と対応する固有ベクトルをそれぞれ λ_{max} , γ_L , $L \in \{1, \dots, d\}$ とする. このとき, \mathbf{K} と Σ_R の最大固有値に関する影響関数 $IF_{\mathbf{K}}(\mathbf{x}, a_1, F)$, $IF_{\Sigma_R}(\mathbf{x}, \lambda_{max}, F)$ はそれぞれ

$$\begin{aligned} IF_{\mathbf{K}}(\mathbf{x}, a_1, F) &= \gamma_1^T A(\mathbf{x}) \gamma_1 - 2a_1, \\ IF_{\Sigma_R}(\mathbf{x}, \lambda_{max}, F) &= \gamma_L^T B(\mathbf{x}) \gamma_L - 3\lambda_{max} \end{aligned}$$

で与えられる. ここで,

$$\begin{aligned} A(\mathbf{x}) &= 2E_{\mathbf{X}} [S(\mathbf{x} - \mathbf{X})S(\mathbf{x} - \mathbf{X})^T], \\ B(\mathbf{x}) &= E_{\mathbf{X}, \mathbf{Y}} [S(\mathbf{X} - \mathbf{Y})S(\mathbf{X} - \mathbf{x})^T + S(\mathbf{X} - \mathbf{x})S(\mathbf{X} - \mathbf{Y})^T + S(\mathbf{x} - \mathbf{X})S(\mathbf{x} - \mathbf{Y})^T] \end{aligned}$$

である. また, \mathbf{K} と Σ_R の固有ベクトル γ_1 に関する影響関数 $IF_{\mathbf{K}}(\mathbf{x}, \gamma_1, F)$, $IF_{\Sigma_R}(\mathbf{x}, \gamma_1, F)$ はそれぞれ

$$\begin{aligned} IF_{\mathbf{K}}(\mathbf{x}, \gamma_1, F) &= \sum_{k=2}^d \frac{1}{a_1 - a_k} \gamma_k^T A(\mathbf{x}) \gamma_1 \cdot \gamma_k, \\ IF_{\Sigma_R}(\mathbf{x}, \gamma_1, F) &= \sum_{k \neq M} \frac{1}{\lambda_M - \lambda_k} \gamma_k^T B(\mathbf{x}) \gamma_1 \cdot \gamma_k \end{aligned}$$

で与えられる.

適用例: 影響関数を用いた固有ベクトルのロバストネスの比較, シミュレーション結果と実データへの適用結果についてはシンポジウムにて報告する.

参考文献

- Hampel, F. R., Ronchetti, E. M., Rousseeuw, P. J., and Stahel, W. A. (1986). *Robust Statistics : The Approach Based on Influence Functions*. Wiley.
- Han, F. and Liu, H. (2018). Eca: High-dimensional elliptical component analysis in non-gaussian distributions. *Journal of the American Statistical Association*, 113:252–268.
- Marden, J. I. (1999). Some robust estimates of principal components. *Statistics & Probability Letters*, 43:349–359.