

# 高次元小標本における固有値の推定とその応用

矢田和善（筑波大学・数理物質科学研究科・院）

青嶋 誠（筑波大学・数理物質科学研究科）

## 1 はじめに

マイクロアレイデータやMRIデータに見られるように、情報化の進展に伴い、データの次元数  $d$  が標本数  $n$  よりも遥かに大きな高次元小標本 (HDLSS) データが、解析対象になる場面が増えてきている。このような HDLSS データに対して従来の統計手法を用いると、次元の呪いによって解析が上手くいかない。低次元空間への次元縮約の最も一般的な手法の一つに、主成分分析 (PCA) がある。しかしながら、HDLSS データに対して従来型の PCA の使用には限界があることが、Muller et al. (2008), Jung and Marron (2008) 等によって報告されている。本発表では、HDLSS データに対して、従来型の固有値や固有ベクトルの推定が有効になるための、標本数  $n$  の  $d$  に関するオーダーを導いた。さらに、Yata and Aoshima (2008) のアプローチを用いて新しい固有値の推定法が提案され、これが、HDLSS データに対して PCA の適用範囲を広げ、推定に必要な標本数  $n$  のオーダーも緩めることを理論的に示した。この方法論は、固有ベクトルの推定にも応用された。

## 2 固有値と固有ベクトルの推定

平均が 0 の  $d$  次元分布をもつ母集団から、 $n$  個のデータベクトルを無作為に抽出して、データ行列  $X: d \times n$  を定義する。ただし、 $d > n$  と仮定する。母共分散行列  $\Sigma$  の固有値を  $\lambda_1 \geq \dots \geq \lambda_d \geq 0$  とし、適当な直交行列  $H = [h_1, \dots, h_d]$  で  $\Sigma = H\Lambda H^T$ ,  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$  と分解する。ここで  $Z = \Lambda^{-1/2} H^T X$  を定義する。 $Z$  の成分は、4 次モーメントが一様有界になることを仮定する。いま、 $\Sigma$  の固有値に、Jung and Marron (2008) 等と同等な次のモデルを仮定する。

$$\lambda_i = a_i d^{\alpha_i} \quad (i = 1, \dots, m), \quad \lambda_j = c_j \quad (j = m+1, \dots, d).$$

ここで、 $a_i (> 0)$ ,  $c_j (\geq 0)$ ,  $\alpha_i (\alpha_1 \geq \dots \geq \alpha_m > 0)$  は未知の実数、 $m$  は未知の整数とする。いま、 $\alpha_i > 1/2$  なる最大の  $i$  を  $k_1$  とし、 $\lambda_1 > \dots > \lambda_{k_1}$  と仮定する。標本共分散行列  $S = n^{-1} X X^T$  と同じ固有値を有する Dual な標本共分散行列を  $S_D = n^{-1} X^T X$  として、 $S_D$  の固有値を  $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_n$  とする。そのとき、 $\hat{\lambda}_i$  ( $i = 1, \dots, k_1$ ) について、以下の結果が導かれた。

定理 1. (i)  $\alpha_i > 1$  ならば、 $d \rightarrow \infty$ ,  $n \rightarrow \infty$ , (ii)  $\alpha_i \in (1/2, 1]$  のとき、 $d \rightarrow \infty$ ,  $d^{1-\alpha_i}/n \rightarrow 0$  のもと、 $\hat{\lambda}_i/\lambda_i = 1 + o_p(1)$ .

HDLSS データに対して,  $S_D$  に基づく推定法では,  $\alpha_i \leq 1/2$  の場合には固有値を推定できず,  $\alpha_i > 1/2$  の場合であっても固有値の推定に多くの標本数が必要になる. Yata and Aoshima (2008) のアプローチを用いて, 2つの独立な  $d \times n$  データ行列  $X_1, X_2$  を使って,  $S^2 = n^{-2} X_1 X_1^T X_2 X_2^T$  を定義する. (ここでの  $n$  は,  $n' = n/2$  を意味する.)  $S^2$  と同じ固有値を有する Dual な 2 乗行列を  $S_D^2 = n^{-2} X_1^T X_2 X_2^T X_1$  とし,  $S_D^2$  の固有値を  $\tilde{\lambda}_1^2 \geq \dots \geq \tilde{\lambda}_n^2$  とする. そのとき, 以下の結果が導かれた.

**定理 2.**  $\alpha_i > 1/4$  なる最大の  $i$  を  $k_2$  とし,  $\lambda_1 > \dots > \lambda_{k_2}$  と仮定する. そのとき,  $\tilde{\lambda}_i^2$  ( $i = 1, \dots, k_2$ ) について, (i)  $\alpha_i > 1/2$  ならば,  $d \rightarrow \infty, n \rightarrow \infty$ , (ii)  $\alpha_i \in (1/4, 1/2]$  ならば,  $d \rightarrow \infty, d^{2-4\alpha_i}/n \rightarrow 0$  のもと,  $\sqrt{\tilde{\lambda}_i^2}/\lambda_i = 1 + o_p(1)$ .

この推定量の良さは, シミュレーション実験において, 推定量の分散も考慮して報告された.

本研究で得られたアプローチを用いれば, 標本数  $n$  を  $d$  の適当なオーダーで定めることで, HDLSS データに対して固有ベクトルを有効に推定することが可能である. いま,  $S$  に対して  $\hat{H}^T S \hat{H} = \hat{\Lambda}$ ,  $\hat{\Lambda} = \text{diag}(\hat{\lambda}_1, \dots, \hat{\lambda}_d)$  となる直交行列を,  $\hat{H} = [\hat{h}_1, \dots, \hat{h}_d]$  とする. そのとき,  $\hat{h}_i$  ( $i = 1, \dots, k_1$ ) について, 以下の結果が導かれた.

**定理 3.** (i)  $\alpha_i > 1$  ならば,  $d \rightarrow \infty, n \rightarrow \infty$ , (ii)  $\alpha_i \in (1/2, 1]$  のとき,  $d \rightarrow \infty, d^{1-\alpha_i}/n \rightarrow 0$  のもと,  $\text{Angle}(\hat{h}_i, h_i) \xrightarrow{p} 0$ .

十分な大きさの標本が得られないとき,  $S_D$  に基づく固有値の推定量  $\hat{\lambda}_i$  は不安定になる. それに伴い, 固有ベクトルの推定となる  $\hat{h}_i$  に誤差が生じ, 主成分スコアにも誤差が生じることとなる.  $S_D^2$  に基づく固有値の推定量  $\sqrt{\tilde{\lambda}_i^2}$  を使えば, HDLSS データにおいて, より精度の高い固有値の推定が可能になる. 本発表では,  $\sqrt{\tilde{\lambda}_i^2}$  に基づくベイズ的なアプローチを使って, 固有ベクトル  $\hat{h}_i$  を修正する方法も報告した.

Jung, S. and Marron, J. S. (2008). PCA consistency in high dimension, low sample size context. *Ann. Statist.*

Muller, K. E., Chi, Y.-Y., Ahn, J. and Marron, J. S. (2008). Limitations of high dimension, low sample size principal components for gaussian data. *J. Amer. Statist. Assoc.*

Yata, K. and Aoshima, M. (2008). Intrinsic dimensionality estimation of high dimension, low sample size data with  $d$ -asymptotics, submitted.