

高次元空間における特徴ベクトルの一致推定について

筑波大学・数理物質系 矢田 和善
筑波大学・数理物質系 青嶋 誠

1. はじめに. ゲノム科学などの現代科学の1つの特徴は、データがもつ次元数の膨大さにある。こういった高次元データの第1の特徴は、次元数が標本数を遥かに超えることである。第2の特徴は、高次元データは豊富な情報を有するものの、それが巨大なノイズに埋もれて見つけ難いことである。これらの理由から、通常の多変量解析法では高次元データの推測に精度を保証することができず、間違った解析結果を導くことさえある。Yata and Aoshima (2009,CSTM) は、高次元データにおけるPCAの性質を研究し、PCAが一致性をもつための標本数 n の d に関するオーダー条件を導き、高次元小標本においてPCAが不適解を起こすことを示した。この問題を解決する策として、Yata and Aoshima (2012,JMA) は、高次元データ空間の幾何学的表現を研究し、それに基づいて“ノイズ掃き出し法”とよばれる方法論を考案した。最近、Aoshima and Yata (2017,SS) は、ノイズ掃き出し法による新たな高次元平均ベクトルの検定法を考案した。

本講演では、高次元空間における特徴ベクトルの一致推定について論じる。具体的には、高次元固有ベクトルや高次元平均ベクトルの一致推定を考える。

2. 高次元固有ベクトルの一致性. 共分散行列に d 次の半正定値行列 Σ をもつ母集団を考える。母集団から n (≥ 3) 個の d 次データベクトル $\mathbf{x}_1, \dots, \mathbf{x}_n$ を無作為に抽出する。 Σ の固有値を $\lambda_1 \geq \dots \geq \lambda_d (\geq 0)$ とし、適当な直交行列 $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_d]$ で Σ を $\Sigma = \mathbf{H}\Lambda\mathbf{H}^T$, $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$ と分解する。標本共分散行列 \mathbf{S} のスペクトル分解を $\mathbf{S} = \sum_{i=1}^d \hat{\lambda}_i \hat{\mathbf{h}}_i \hat{\mathbf{h}}_i^T$ とする。そのとき、Yata and Aoshima (2009) は、適当な正則条件のもと、

$$\hat{\lambda}_j / \lambda_j = 1 + \delta_j + o_p(1) \quad \text{and} \quad \mathbf{h}_j^T \hat{\mathbf{h}}_j = (1 + \delta_j)^{-1/2} + o_p(1) \quad \text{as } d, n \rightarrow \infty \quad (1)$$

が成り立つことを示した。ただし、 $\delta_j = \lambda_j^{-1} \sum_{i=m+1}^d \lambda_i / (n-1)$, $m (\geq j)$ はある定数である。一方で、Yata and Aoshima (2012) は、高次元データ空間の幾何学的表現を研究し、それに基づいて“ノイズ掃き出し法”とよばれる方法論を考案し、 $\tilde{\lambda}_j = \hat{\lambda}_j - (\text{tr}(\mathbf{S}) - \sum_{i=1}^j \hat{\lambda}_i) / (n-1-j)$ ($j = 1, \dots, n-2$) なる固有値の推定量を提案した。さらに、 Σ の固有ベクトル \mathbf{h}_j を $\tilde{\mathbf{h}}_j = (\hat{\lambda}_j / \tilde{\lambda}_j)^{1/2} \hat{\mathbf{h}}_j$ で推定する。そのとき、Yata and Aoshima (2012) は、適当な正則条件のもと、

$$\tilde{\lambda}_j / \lambda_j = 1 + o_p(1) \quad \text{and} \quad \mathbf{h}_j^T \tilde{\mathbf{h}}_j = 1 + o_p(1) \quad \text{as } d, n \rightarrow \infty \quad (2)$$

が成り立つことを示した。それゆえ、 \mathbf{h}_j の内積に関する一致性をもつ。

ここで、 $\|\tilde{\mathbf{h}}_j\|^2 = \hat{\lambda}_j / \tilde{\lambda}_j$ であることに注意し、(1) と (2) より、適当な正則条件のもと、
 $\|\tilde{\mathbf{h}}_j - \mathbf{h}_j\|^2 = 2\{1 - (1 + \delta_j)^{-1/2}\} + o_p(1)$ and $\|\tilde{\mathbf{h}}_j - \mathbf{h}_j\|^2 = \delta_j + o_p(1)$ as $d, n \rightarrow \infty$
となる。すなわち、 $\liminf_{d, n \rightarrow \infty} \delta_j > 0$ のとき、 $\hat{\mathbf{h}}_j$ と $\tilde{\mathbf{h}}_j$ はノルムに関する一致性をもたない。

本講演では、閾値を用いて $\tilde{\mathbf{h}}_j$ を補正することで、 $\delta_j \rightarrow \infty$ のもとでも、

$$\|\tilde{\mathbf{h}}_j - \mathbf{h}_j\|^2 = o_p(1)$$

が成り立つような新たな固有ベクトルの推定量 $\tilde{\mathbf{h}}_j$ を提案し、理論的かつ数値的に既存の推定量と比較する。さらに、上記のアイデアを高次元平均ベクトルの推測にも応用し、新たな高次元平均ベクトルの一致推定量も考案する。