高次元現象の統計数理

青嶋 誠 (筑波大数理物質)*

1. はじめに

ゲノム科学・情報工学・金融工学などの現代科学の一つの特徴は,データがもつ次元数の膨大さにある.例えば,DNAマイクロアレイによる遺伝子発現データを見ると,次元数が数万にものぼる一方で,標本数は100にも満たないという事例が多い.高次元データは,高次元になるほど標本を入手することが難しく,一般にデータの次元数pと標本数nには

$p \gg n$ (もしくは, p > n)

といった大小関係がある.これを強調して,高次元小標本データとよぶこともある. 2005年頃から,次世代シーケンサーによる超高次元ゲノムデータも統計解析の対象と なり,例えば Wang et al. [13] は,次元数が300万を超え,標本数は40というゲノム データを扱っている.近年,高次元小標本データの解析は,様々な分野で需要が高まっ ている.本講演は,次世代シーケンサーによる超高次元ゲノムデータをはじめ,天文 宇宙学における天体の分光マッピングによる波長データ([11])や,材料科学における 結晶構造データなど,高次元小標本データの解析例を幾つか紹介する.

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いる ために見つけ難い.伝統的な多変量解析法では,高次元データの解析に精度を保証す ることができない.次元数が標本数を超えると,次元の呪いとして様々な高次元現象 が現れ,多変量解析の理論と方法論の成立を阻むのである.特に非スパース性という 現象に陥ると,データの潜在情報は高次元において膨張するノイズに埋没してしまい, 推測の強不一致という最悪な結果を招くことにもなる.

高次元データの解析には,高次元データ特有の理論と方法論が必要になる.著者の 研究グループは,高次元統計解析と非スパースモデリングという新しい理論と方法論 を考案した.本講演は,様々な高次元現象を説明し,たった数十程度の小標本であって も高速かつ高精度に解析するための,高次元現象の統計数理を紹介する.

2. 高次元データの非スパース性

1990年代後半,高次元データを解析するための方法論として,機械学習が目覚ましく 発展した.スタンフォード大学のDonoho教授やTibshirani教授らは,スパースモデリ ングを発展させた.これは,高次元データの母共分散行列にスパース性を仮定して,高 次元データの特徴を少数の変数や標本で捉えようとするもので,L1正則化法や圧縮セ ンシングがよく知られる([7]).近年,多層ニューラルネットワークを用いた深層学習 が盛んに研究されている.これもスパース性を仮定して,データは大標本であること

2010 Mathematics Subject Classification: 62H25, 62H30

本研究は,筑波大学の矢田和善准教授との共同研究である.本研究は,科学研究費補助金基盤研究(A) 20H00576 研究代表者:青嶋誠「大規模複雑データの理論と方法論の革新的展開」,および,学術研究 助成基金助成金挑戦的研究(萌芽)22K19769 研究代表者:青嶋誠「テンソル構造をもつ巨大データの 統計的圧縮技術の開発」から助成を受けたものである.

キーワード:高次元統計解析,非スパースモデリング

^{*}e-mail: aoshima@math.tsukuba.ac.jp web: https://www.math.tsukuba.ac.jp/~aoshima-lab/jp/

が前提になっている.ところで,高次元データの母共分散行列に,スパース性は仮定 できるのだろうか.

2.1. 高次元共分散行列の非スパース性

p次非負定値対称行列 Σ をp次元データの母共分散行列とし,その固有値を $\lambda_1 \ge \cdots \ge \lambda_p$ (≥ 0)とする.次の3つの高次元データセットについて,それぞれの固有値を上から 10個求めてみる.第1のデータセットは,Takeuchi et al. [11]による天体の分光マッピン グによる波長データで,標本数は231地点,スペクトルは1971次元である.第2のデー タセットは,von Roemeling et al. [12]による2群からなる淡明細胞型腎細胞癌 (ccRCC) のマイクロアレイデータで,ccRCC群は71サンプル,健常群 (Normal)は72サンプル で,2群とも遺伝子数(プロープ数)が54675次元である.第3のデータセットは,Wang et al. [13]の次世代シーケンサーによる2群からなる加齢黄斑変性 (AMD)のゲノムデー タで,AMD群は21サンプル,健常群 (Normal)は19サンプルで,2群ともゲノム領域 数が3095656次元である.図1は,寄与率 $\lambda_j/\text{tr}(\Sigma)$ について,Yata and Aoshima [15]の ノイズ掃き出し法による推定値 $\tilde{\lambda}_j/\text{tr}(S)$ をプロットしたものである¹.どの高次元デー タも,最初の幾つかの固有値が飛び抜けて大きく,他の固有値は小さなものであること が見てとれる.特に,第1固有値は強くスパイクしており,一般に $\text{tr}(\Sigma)/p \rightarrow c$ (>0) となることを考えれば, λ_1 はO(p)のスピードで発散することが分かる.



図 *I*. 天体の分光マッピングによる波長データ (p = 1971), 2つの群のマイクロアレイ データ (p = 54675), 2つの群の次世代シーケンサーによるゲノムデータ (p = 3095656). それぞれの高次元データについて,第10寄与率まで推定値をプロットした.

高次元データの母共分散行列の固有値に対して, Johnstone [9] や Paul [10] 等は, ランダム行列理論から標本固有値の漸近的性質を導出した².特に,スタンフォード大学の Johnstone 教授は, Tracy-Widom 分布に基づく標本固有値の極限分布を与えた([9]).

¹ $\tilde{\lambda}_j$ は(4.2)式で与えられるもの,*S*は標本共分散行列である.後述の(A-i)と(A-ii)のもと, $\tilde{\lambda}_j$ /tr(*S*) = $\lambda_j \{1 + o_P(1)\}/$ tr(Σ) ($p \to \infty, n \to \infty$)なる一致性をもつ.

²これらの研究は、次元数
 pと標本数
 nが $n/p \to c \ (>0)$ なる場合を考えている
 . 高次元小標本データ $(n/p \to 0)$ には対応していないことに注意する .

そのとき,∑の固有値には,次のようなスパイクモデルが仮定されていた.

 $\lambda_1, ..., \lambda_t$ は1よりも大きく,次元数pには依存しない定数

$$\lambda_{t+1} = \dots = \lambda_p = 1 \tag{2.1}$$

(2.1)のもとでは, $\|\Sigma - I_p\|_F^2 = \sum_{s=1}^t (\lambda_s - 1)^2 = O(1) (p \to \infty)$ となり³, $p \times p$ 行列 $\Sigma \geq p$ 次単位行列 I_p の距離が有界になる.つまり, (2.1)は, Σ の非対角成分の殆どが 零である,といった非常にスパースな共分散構造を想定していることになる.スパー スモデリングでは, (2.1)と同様な「 λ_1 は次元数pに関して有界である」という条件が 仮定される.図1でも見たが,固有値が次元数pに依存しないという仮定は,実データ がもつ固有値構造から大きく乖離する.それでは,実データの固有値構造に飛び抜け て大きな固有値が現れる原因は何か.主たる原因は,高次元データを構成する変数間 の相関にある.実データに対する Σ は,非対角成分の多くが非零となり,非スパース な共分散構造をもつのである.困ったことに,非スパースな共分散構造は高次元にお いて巨大なノイズを生み,潜在情報を埋もれさせてしまう.これが,高次元データの 解析を困難にしている.このような高次元データに,ランダム行列理論・スパースモ デリング・深層学習などを,そのまま使用することは望ましくない.そこで,高次元 データが非スパースな共分散構造をもつことを考慮して,新たに開発された方法論が 非スパースモデリングである.



図2. 高次元データの解析: スパース性と非スパース性.

2.2. 強スパイク固有値モデル

高次元データの母共分散行列 Σ は,次元数の増加とともに行列のサイズが大きくなり,図1で見たように最初の幾つかの固有値は次元数に依存して大きなものとなる.特に,最大固有値 λ_1 は,次元数の冪関数として捉えることが自然である. Aoshima and Yata [4] は,高次元データの固有値モデルを2つに分類した.一つ目は,強スパイク固有値モデル (strongly spiked eigenvalue (SSE) モデル)とよばれ,次のように定義される.

$$\liminf_{p \to \infty} \left\{ \frac{\lambda_1^2}{\operatorname{tr}(\boldsymbol{\Sigma}^2)} \right\} > 0 \tag{2.2}$$

いま一つは,弱スパイク固有値モデル(NSSEモデル)とよばれ,次のように定義される.

$$\frac{\lambda_1^2}{\operatorname{tr}(\boldsymbol{\Sigma}^2)} \to 0 \ (p \to \infty) \tag{2.3}$$

 $||^{3} || \cdot ||_{F}$ はフロベニウスノルムを表す.

任意の高次元データは,上記2つの固有値モデルの何れかに分類される.簡単のため, 以下のような固有値構造を考えてみる.

$$\lambda_s = c_s p^{\alpha_s} \ (s = 1, ..., t), \quad \lambda_s = c_s \ (s = t + 1, ..., p) \tag{2.4}$$

ここで, $c_s > 0$, $\alpha_1 \ge \cdots \ge \alpha_t > 0$ は次元数 p に依存しない未知の実数, t は次元数 p に依存しない未知の自然数である. (2.4) は, $\alpha_1 \ge 0.5$ のとき強スパイク固有値モデル (2.2) の一つとなり, $\alpha_1 < 0.5$ のとき弱スパイク固有値モデル(2.3) の一つとなる.

後述するが,弱スパイク固有値モデルにおいては種々の統計量に高次元漸近正規性 という望ましい性質が得られるが,強スパイク固有値モデルでは高次元漸近正規性が 成り立たない.それどころか,強くスパイクする固有値が一つでもあると,それが巨 大なノイズとなって潜在情報を覆い隠してしまう.そのとき,強スパイク固有値モデ ルは,非対角成分の多くが非零といった非スパースなΣを想定したものになる.(2.2) は母共分散行列に関する境界条件になっており,SSEモデルは非スパース性に対応し, NSSEモデルはスパース性に対応する.

3. 高次元漸近正規性

平均がp次ベクトル μ ,共分散行列が Σ の母集団から, $n \geq 3$)個のデータベクトル $x_1, ..., x_n$ を無作為に抽出する.それらを並べて, $p \times n$ のデータ行列を $X = [x_1, ..., x_n]$ とする.適当な直交行列 $H = [h_1, ..., h_p]$ で $\Sigma = H\Lambda H^T$, $\Lambda = \text{diag}(\lambda_1, ..., \lambda_p)$ と分解する. そのとき, $X - [\mu, ..., \mu] = H\Lambda^{1/2}Z$ となるZを $Z = [z_1, ..., z_p]^T$, $z_s = (z_{1s}, ..., z_{ns})^T$ と表記する.Zの成分は4次モーメントが一様有界であると仮定し,さらに,必要な箇所で,適宜,次を仮定する.

(A-i)
$$E(z_{\ell s}^2 z_{\ell t}^2) = 1$$
, $E(z_{\ell s} z_{\ell t} z_{\ell s'}) = 0$, $E(z_{\ell s} z_{\ell t} z_{\ell s'} z_{\ell t'}) = 0$, $s \neq t, s', t'$

(A-i)は,母集団の分布について,正規分布を緩めた仮定になっている.

母平均 μ に関する推測を考える.そのとき,共分散行列 Σ はノイズとなり,前述の SSE モデルとNSSE モデルはノイズモデルになる.標本平均ベクトル $\overline{x} = \sum_{\ell=1}^{n} x_{\ell}/n$ について,平均2乗誤差は $E(\|\overline{x} - \mu\|^2) = \operatorname{tr}(\Sigma)/n$ となる.ここで, $\|\cdot\|$ はユーク リッドノルムである.一般に $\operatorname{tr}(\Sigma)/p \to c$ (> 0)となることを考えれば,高次元小標本 $(p/n \to \infty)$ において

$$E(\|\overline{\boldsymbol{x}} - \boldsymbol{\mu}\|^2) \to \infty$$

となり,平均2乗誤差は発散してしまう.すなわち,次元数の増加とともに標本平均に 含まれるノイズが巨大化し, \overline{x} で μ を推定することは困難になる.ところで, $\|\overline{x} - \mu\|^2$ の分散は,(A-i)のもと $p \to \infty$ のとき

$$\operatorname{Var}(\|\overline{\boldsymbol{x}} - \boldsymbol{\mu}\|^2) = 2\frac{\operatorname{tr}(\boldsymbol{\Sigma}^2)(n-1)}{n^3} + \sum_{s=1}^p \frac{\lambda_s^2 \operatorname{Var}(z_{\ell s}^2)}{n^3} = 2\frac{\operatorname{tr}(\boldsymbol{\Sigma}^2)}{n^2} + O\left(\frac{\operatorname{tr}(\boldsymbol{\Sigma}^2)}{n^3}\right)$$

となる.さらに,球形条件 ⁴"tr(Σ^2)/{tr(Σ)}² → 0 ($p \to \infty$)"を満たすと仮定すると, $p \to \infty$ のとき (nは固定したままでも)次が成り立つ.

$$\|\overline{\boldsymbol{x}} - \boldsymbol{\mu}\| = \sqrt{\operatorname{tr}(\boldsymbol{\Sigma})/n} \{1 + o_P(1)\}$$

 ${}^4 \operatorname{tr}(\Sigma^2)/{\operatorname{tr}(\Sigma)}^2 \leq \lambda_1/\operatorname{tr}(\Sigma) \leq \lambda_1/{\operatorname{tr}(\Sigma^2)}^{1/2}$ に注意すれば, NSSE モデル(2.3)は球形条件を満たす.

すなわち,高次元小標本 $(p/n \to \infty)$ において, \overline{x} は中心 μ ,半径 $\sqrt{\operatorname{tr}(\Sigma)/n}$ の球面に収束する.図3は,単位行列 I_p を共分散行列にもつp次元正規分布 $N_p(\mu, I_p)$ について,大きさn = 3の無作為標本による標本平均ベクトルを200回発生させ, $\overline{x} - \mu$ を固有空間 $(\hat{h}_1, \hat{h}_2, \hat{h}_3$ で張られる空間)にプロットしたものである.ここで, \hat{h}_j は標本共分散行列の第j固有ベクトルを表す.次元数がp = 4のとき(左図)は,半径 $\sqrt{p/n} = \sqrt{4/3}$ の球の周辺に点が散らばっている.一方,次元数がp = 1000のとき(右図)は,ノイズが膨張して半径 $\sqrt{p/n} = \sqrt{1000/3}$ の球の表面に点が集中し,球面集中現象が見てとれる.



図 3. 高次元小標本における 家の球面集中現象.

球面での *x* の確率変動について, Aoshima and Yata [1, 3] は, 次のような高次元漸近 正規性を証明した.

定理 1 ([3]). (*A*-*i*)を仮定する. *NSSE* モデルのもとで, $p \to \infty, n \to \infty$ のとき次が成り 立つ⁵.

$$\frac{\|\overline{\boldsymbol{x}} - \boldsymbol{\mu}\|^2 - tr(\boldsymbol{\Sigma})/n}{\sqrt{2tr(\boldsymbol{\Sigma}^2)/n^2}} \xrightarrow{\mathcal{L}} N(0, 1)$$

ここで, $\stackrel{\mathcal{L}}{\rightarrow}$ は分布収束を表す.

定理1は,高次元において, \overline{x} が中心 μ ,半径 $\sqrt{\operatorname{tr}(\Sigma)/n}$ の球面の周りに正規分布で変動することを意味する.このことから,NSSEモデルであれば,球面付近の微小な確率変動を捉えることで統計的推測の精度を保証することが可能となる.しかしながら,SSEモデルの場合,これが上手くいかない.強くスパイクする(非スパースな)ノイズが悪さをして,高次元漸近正規性は成立しない.いま,p = 1000のとき,NSSEモデルとして $\Sigma = \operatorname{diag}(p^{1/3}, 1, ..., 1)$ を,SSEモデルとして $\Sigma = \operatorname{diag}(p, 1, ..., 1)$ を考えてみる. 図4は,p次元正規分布 $N_p(\mu, \Sigma)$ について,大きさn = 10の無作為標本による標本平均ベクトルをもとに $U = \{||\overline{x} - \mu||^2 - \operatorname{tr}(\Sigma)/n\}/\sqrt{2\operatorname{tr}(\Sigma^2)/n^2}$ を計算し,これを2000回発生させてヒストグラムを作成したものである.定理1の通り,NSSEモデルにおけるヒストグラムの形状はN(0,1)に近いことが確認できる.一方,SSEモデルにおけるとへがうんの形状はN(0,1)に近いことが確認できる.一方,SSEモデルにおいては,高次元漸近正規性が成立しない.さらに,SSEモデルにおいては $\operatorname{tr}(\Sigma^2) = O(\lambda_1^2)$ となるので,図1のように最大固有値が飛び抜けている場合, $||\overline{x} - \mu||^2$ の分散が非常に大きくなる.つまり,SSEモデルが当てはまる状況では,高次元漸近正規性が壊れるだけでなく,推測の精度も著しく悪くなるのである.このことから,NSSEモデルを想定して(つまり,ノイズのスパース性を想定して)構築された統計的推測法は,SSE

⁵本講演で $n \to \infty$ と表記するとき, $n = \log p$ 程度で十分である.つまり,高次元小標本であっても, 結果は保証される.



図 4. NSSE モデル $\Sigma = \operatorname{diag}(p^{1/3}, 1, ..., 1)$ とSSE モデル $\Sigma = \operatorname{diag}(p, 1, ..., 1)$ の場合の, 球面付近の \overline{x} の漸近的挙動. U(p = 1000, n = 10) を 2000 回発生させ, ヒストグラム を作成した.実線は N(0, 1) の確率密度関数を表す.

モデルに対しては精度を保証することができず,それゆえSSEモデルに使うわけにはいかない.それでは,SSEモデルも考慮した統計的推測法を構築することはできないか.Aoshima and Yata [4] は,SSEモデルの場合は自動的にNSSEモデルに変換して推測を行うデータ変換法を考案した.次節でこれを説明するが,表記を簡単にするために,ここでSSEモデルを次のように簡略化しておく.

(A-ii) (次元数 p に依存しない) ある自然数 k に対して,

(i)
$$1 \le s < s' \le k$$
のとき , $\liminf_{p \to \infty} (\lambda_s / \lambda_{s'} - 1) > 0$
(ii) $\liminf_{p \to \infty} \frac{\lambda_k^2}{\Psi_k} > 0$ かつ $\frac{\lambda_{k+1}^2}{\Psi_{k+1}} \to 0, p \to \infty$ (ただし , $\Psi_j = \sum_{s=j}^p \lambda_s^2$ である)

(2.4)の固有値構造であれば, k(< t)に対して $\alpha_k \ge 0.5 > \alpha_{k+1}$ 且つ $1 \le s < s' \le k$ に対して $c_s \ne c_{s'}$ ならば,条件(A-ii)を満たす.(A-ii)は,図1の実データ解析をヒントに,最初のk個の固有値が飛び抜けて大きいことを表現したSSEモデルになっている.

4. データ変換法

Aoshima and Yata [4] は,次のようなデータ変換を考えた⁶.

$$\boldsymbol{x}_{\ell*} = \boldsymbol{A}\boldsymbol{x}_{\ell} \ (\ell = 1, ..., n) \quad ($$
ただし , $\boldsymbol{A} = \sum_{s=k+1}^{p} \boldsymbol{h}_{s} \boldsymbol{h}_{s}^{T} \ (k \ge 0)$ である) (4.1)

変換後のデータについて, $Var(\boldsymbol{x}_{\ell*}) = \sum_{s=k+1}^{p} \lambda_s \boldsymbol{h}_s \boldsymbol{h}_s^T (= \Sigma_* とおく) となるので, \lambda_{\max}(\Sigma_*) を \Sigma_* の最大固有値とすれば, SSE モデル(A-ii)の条件(ii)から$

$$\frac{\{\lambda_{\max}(\boldsymbol{\Sigma}_*)\}^2}{\operatorname{tr}(\boldsymbol{\Sigma}_*^2)} = \frac{\lambda_{k+1}^2}{\Psi_{k+1}} \to 0, \ p \to \infty$$

が成立する.つまり,変換後のデータ $x_{\ell*}$ には,NSSEモデルが当てはまることになる. データ変換に使う正射影行列は $A = I_p - \sum_{s=1}^k h_s h_s^T$ と書けるので,強くスパイクする 固有空間を如何に高精度に推定するかが鍵となる.本節では,高次元主成分分析(PCA) として知られるノイズ掃き出し法を用いて固有値・固有ベクトルを推定する.

標本共分散行列

$$m{S}=(n-1)^{-1}(m{X}-\overline{m{X}})(m{X}-\overline{m{X}})^T$$
 (ただし, $\overline{m{X}}=[\overline{m{x}},...,\overline{m{x}}]$ である)

⁶強くスパイクする固有空間の個数 k は未知である . k の推定は , Aoshima and Yata [4] の S2.2 節を参照 のこと.なお , k = 0 の場合は恒等変換となり , データは変換せずとも NSSE モデルが当てはまる .

について, Sの固有値を $\hat{\lambda}_1 \geq \cdots \geq \hat{\lambda}_p$ (≥ 0), 各固有値 $\hat{\lambda}_j$ に対する固有ベクトルを \hat{h}_j と する.ただし, $\hat{h}_1, ..., \hat{h}_p$ はp次元実ベクトル空間の正規直交基底をなすとする.Sの固 有値分解は $S = \sum_{s=1}^p \hat{\lambda}_s \hat{h}_s \hat{h}_s^T$ と書ける. $S_D = (n-1)^{-1} (X - \overline{X})^T (X - \overline{X})$ とおくと,n次の対称行列 S_D はSと正の固有値を共有する. S_D の固有値を $\hat{\lambda}_1 \geq \cdots \geq \hat{\lambda}_{n-1}$ (≥ 0), 各 $\hat{\lambda}_j$ に対する固有ベクトルを $\hat{u}_j = (\hat{u}_{1j}, ..., \hat{u}_{nj})^T$ とする.ただし, $\hat{u}_1, ..., \hat{u}_{n-1}$ は互い に正規直交とする. S_D の固有値分解は $S_D = \sum_{s=1}^{n-1} \hat{\lambda}_s \hat{u}_s \hat{u}_s^T$ と書ける. $S \geq S_D$ の固有 ベクトルには,次の関係がある.

$$\hat{h}_j = \{(n-1)\hat{\lambda}_j\}^{-1/2} (X - \overline{X})\hat{u}_j \quad (j = 1, ..., n-1)$$

Yata and Aoshima [14] は,固有値構造が(2.4)の場合に, $\hat{\lambda}_j$, \hat{h}_j , さらに主成分スコアについて,推定量の一致性(と不一致性)を明らかにした. Aoshima and Yata [4] と Yata and Aoshima [16] は,一般のSSEモデルの場合にその理論を拡張している. SSEモデル(A-ii)を使って $\hat{\lambda}_j$ と \hat{h}_j について簡単に述べると,次のようになる.

定理 2 ([4]). (*A-i*) を仮定する . *SSE* モデル (*A-ii*) のもとで , 各 $j (\leq k)$ について , $p \rightarrow \infty, n \rightarrow \infty$ のとき次が成り立つ .

$$\hat{\lambda}_j / \lambda_j = 1 + \delta_j + O_P(n^{-1/2}), \quad (\hat{\boldsymbol{h}}_j^T \boldsymbol{h}_j)^2 = (1 + \delta_j)^{-1} + O_P(n^{-1/2})$$

ただし , $\delta_j = \lambda_j^{-1} \sum_{s=k+1}^p \lambda_s / (n-1)$ である .

上記の通り, $\hat{\lambda}_j \geq \hat{h}_j$ は, $\delta_j \geq N$ うノイズをもつ.(2.4)の固有値構造であれば, $\alpha_j = 0.5$ に対して $n = o(p^{0.5})$ だとすると $\delta_j \rightarrow \infty \geq 0$ より, $\lambda_j / \hat{\lambda}_j = o_P(1) \lor (\hat{h}_j^T h_j)^2 = o_P(1) \geq 0$ いった強不一致性を起こす.この問題を解決するために, Yata and Aoshima [15, 16]はノイズ掃き出し法という高次元におけるPCAを考案し,次のような固有値の推定量を与えた.

$$\tilde{\lambda}_j = \hat{\lambda}_j - \frac{\operatorname{tr}(\boldsymbol{S}_D) - \sum_{s=1}^j \hat{\lambda}_s}{n - 1 - j} \quad (j = 1, ..., n - 2)$$
(4.2)

固有ベクトルの推定量は次のようになる.

$$\tilde{\boldsymbol{h}}_j = \{(n-1)\tilde{\lambda}_j\}^{-1/2} (\boldsymbol{X} - \overline{\boldsymbol{X}}) \hat{\boldsymbol{u}}_j \quad (j = 1, ..., n-2)$$
(4.3)

Aoshima and Yata [4] と Yata and Aoshima [16] による次の定理は, $\delta_j \rightarrow \infty$ となる場合にも, ノイズ掃き出し法による推定量が一致性をもつことを示している.

定理 3 ([4]). (*A-i*) を仮定する. *SSE* モデル (*A-ii*) のもとで, 各 $j (\leq k)$ について, $p \rightarrow \infty, n \rightarrow \infty$ のとき次が成り立つ.

$$\tilde{\lambda}_j / \lambda_j = 1 + O_P(n^{-1/2}), \quad (\tilde{\boldsymbol{h}}_j^T \boldsymbol{h}_j)^2 = 1 + O_P(n^{-1})$$

データ変換(4.1)を使うためには, Ax_{ℓ} を構成する $h_j^T x_{\ell}$ (= $x_{\ell j}$ とおく),j = 1, ..., kを推定する必要がある. Aoshima and Yata [4] は, $\hat{h}_j^T x_{\ell}$ や $\tilde{h}_j^T x_{\ell}$ といった推定は非常に大きなバイアスが生むことを指摘し,次のような推定を考えた.

$$\tilde{x}_{\ell j} = \tilde{\boldsymbol{h}}_{\ell j}^{T} \boldsymbol{x}_{\ell} \quad (j = 1, ..., k)$$
(4.4)

ただし,

$$\tilde{\boldsymbol{h}}_{\ell j} = (n-1)^{1/2} (\boldsymbol{X} - \overline{\boldsymbol{X}}) \hat{\boldsymbol{u}}_{\ell j} / \{ (n-2) \tilde{\lambda}_{j}^{1/2} \}, \\ \hat{\boldsymbol{u}}_{\ell j} = (\hat{u}_{1j}, ..., \hat{u}_{\ell-1j}, -\hat{u}_{\ell j} / (n-1), \hat{u}_{\ell+1j}, ..., \hat{u}_{nj})^{T}$$

ここで, $\sum_{\ell=1}^n ilde{m{h}}_{\ell j}/n = ilde{m{h}}_j$ となることに注意する.

5. 高次元2標本検定とデータ解析

母集団が2つあり,各母集団 π_i (i = 1, 2) は平均に p 次ベクトル μ_i , 共分散行列に p 次 非負定値対称行列 Σ_i $(= \sum_{s=1}^p \lambda_{is} \mathbf{h}_{is} \mathbf{h}_{is}^T)$ をもつとする.次の検定を考える⁷.

$$H_0: \mu_1 = \mu_2$$
 vs. $H_1: \mu_1 \neq \mu_2$ (5.1)

各母集団から, n_i (\geq 3) 個のp次データベクトル $x_{i1}, ..., x_{in_i}$ を無作為に抽出し,データ 行列を $X_i = [x_{i1}, ..., x_{in_i}]$ とする.高次元小標本では標本共分散行列の逆行列が存在し ないので,従来のホテリングの T^2 -統計量は使えない.高次元2標本検定では,次の統 計量が使われる.

$$T = \|\overline{\boldsymbol{x}}_{1} - \overline{\boldsymbol{x}}_{2}\|^{2} - \sum_{i=1}^{2} \frac{\operatorname{tr}(\boldsymbol{S}_{i})}{n_{i}} = 2\sum_{i=1}^{2} \frac{\sum_{\ell < \ell'}^{n_{i}} \boldsymbol{x}_{i\ell'}^{T} \boldsymbol{x}_{i\ell'}}{n_{i}(n_{i}-1)} - 2\frac{\sum_{\ell=1}^{n_{1}} \sum_{\ell'=1}^{n_{2}} \boldsymbol{x}_{1\ell}^{T} \boldsymbol{x}_{2\ell'}}{n_{1}n_{2}}$$
(5.2)

ただし, $\overline{x}_i = \sum_{\ell=1}^{n_i} x_{i\ell}/n_i$, $\overline{X}_i = [\overline{x}_i, ..., \overline{x}_i]$, $S_i = (n_i - 1)^{-1} (X_i - \overline{X}_i) (X_i - \overline{X}_i)^T$ である.ここで, $E(T) = \|\mu_1 - \mu_2\|^2$ となり, H_0 のもとで

$$\operatorname{Var}(T) = 2\sum_{i=1}^{2} \frac{\operatorname{tr}(\boldsymbol{\Sigma}_{i}^{2})}{n_{i}(n_{i}-1)} + \frac{4}{n_{1}n_{2}}\operatorname{tr}(\boldsymbol{\Sigma}_{1}\boldsymbol{\Sigma}_{2}) \quad (= K \succeq \mathfrak{s} \triangleleft)$$

となることに注意する.2つの母集団がともにNSSEモデルのとき,*T*に関する高次元 漸近正規性がChen and Qin [6] と Aoshima and Yata [1,3] によって示され,それぞれ別 の検定手法が与えられた.ところが,2つの母集団のどちらか一方でもSSEモデルの とき,3節の図4のように,*T*の高次元漸近正規性は成立しない.さらに,2つの母集 団がともにSSEモデルのとき, $K = O(\lambda_{11}^2/n_1^2 + \lambda_{21}^2/n_2^2)$ となり,各母集団の最大固有 値の影響で*T*の分散は非常に大きくなる.つまり,SSEモデルが当てはまる状況では, *T*の高次元漸近正規性が壊れるだけでなく,検定の精度も著しく悪くなる.

Aoshima and Yata [4] は, SSE モデルも考慮し, データ変換法を使って高次元2標本 検定を次のように与えた.各母集団のデータを, (4.1)と同様にデータ変換する.

$$\boldsymbol{x}_{i\ell*} = \boldsymbol{A}_i \boldsymbol{x}_{i\ell} \ (\ell = 1, ..., n_i) \ (ただし, \boldsymbol{A}_i = \sum_{s=k_i+1}^p \boldsymbol{h}_{is} \boldsymbol{h}_{is}^T \ (k_i \ge 0)$$
である)

変換後のデータを使うと,(5.2)式のTは次のようになる.

$$T_{*} = 2 \sum_{i=1}^{2} \sum_{\ell < \ell'}^{n_{i}} \frac{\boldsymbol{x}_{i\ell}^{T} \boldsymbol{A}_{i} \boldsymbol{x}_{i\ell'}}{n_{i}(n_{i}-1)} - 2 \frac{\sum_{\ell=1}^{n_{1}} \sum_{\ell'=1}^{n_{2}} \boldsymbol{x}_{1\ell}^{T} \boldsymbol{A}_{1} \boldsymbol{A}_{2} \boldsymbol{x}_{2\ell'}}{n_{1}n_{2}}$$
$$= 2 \sum_{i=1}^{2} \frac{\sum_{\ell < \ell'}^{n_{i}} (\boldsymbol{x}_{i\ell}^{T} \boldsymbol{x}_{i\ell'} - \sum_{j=1}^{k_{i}} x_{i\ell j} x_{i\ell' j})}{n_{i}(n_{i}-1)}$$
$$- 2 \frac{\sum_{\ell=1}^{n_{1}} \sum_{\ell'=1}^{n_{2}} (\boldsymbol{x}_{1\ell} - \sum_{j=1}^{k_{1}} x_{1\ell j} \boldsymbol{h}_{1j})^{T} (\boldsymbol{x}_{2\ell'} - \sum_{j=1}^{k_{2}} x_{2\ell' j} \boldsymbol{h}_{2j})}{n_{1}n_{2}}$$

⁷本節は,SSEモデルも考慮した高次元統計的推測を説明するために,基本的で本質的な高次元2標本 検定を扱う.高次元判別分析などは,[5,8]を参照のこと.なお,いわば基礎編となるNSSEモデルの 高次元統計的推測については,青嶋・矢田[2]が入門書となる.

ただし, $x_{i\ell j} = h_{ij}^T x_{i\ell}$ である.ここで, $E(T_*) = \|A_1\mu_1 - A_2\mu_2\|^2$ (= Δ_* とおく)となり, H_0 のもとで

$$\operatorname{Var}(T_*) = 2\sum_{i=1}^{2} \frac{\operatorname{tr}(\boldsymbol{\Sigma}_{i*}^2)}{n_i(n_i - 1)} + \frac{4}{n_1 n_2} \operatorname{tr}(\boldsymbol{\Sigma}_{1*} \boldsymbol{\Sigma}_{2*}) \quad (= K_* \boldsymbol{\Sigma} \boldsymbol{\mathfrak{s}} \boldsymbol{\zeta})$$

となることに注意する.ただし, $\Sigma_{i*} = \sum_{s=k_i+1}^{p} \lambda_{is} h_{is} h_{is}^{T}$ である. Aoshima and Yata [4] は, T_* に関して高次元漸近正規性を証明した.さらに, Aoshima and Yata [4] は, T_* を 構成する $h_{ij} \ge x_{i\ell j}$ を(4.3)式の $\tilde{h}_{ij} \ge$ (4.4)式の $\tilde{x}_{i\ell j}$ で置き換えた T_* の推定量 \hat{T}_* について, 次の結果を与えた.

定理 4 ([4]). 2つの母集団に (A-i) を仮定する.SSE モデルに (A-ii) を仮定する.さらに, $v = \min\{p, n_1, n_2\}$ とおき, $\limsup_{v \to \infty} \Delta_*^2 / K_* < \infty$ を仮定する.適当な正則条件⁸の もと, $v \to \infty$ のとき次が成り立つ.

$$(\widehat{T}_* - \Delta_*) / \sqrt{\widehat{K}_*} \xrightarrow{\mathcal{L}} N(0, 1)$$

ただし, \hat{K}_* は*Aoshima and Yata* [4]の5節で与えられる K_* の一致推定量である. 仮説(5.1)は,次のように検定される.

$$\widehat{T}_*/\sqrt{\widehat{K}_*} > z_{\alpha} \Longrightarrow$$
有意水準 α で H_0 を棄却 (5.3)

ただし, z_{α} は N(0,1) の上側 α 点を表す.検定方式 (5.3) の第1種の過誤と検出力については, Aoshima and Yata [4] の定理6を参照のこと.

図1で紹介した Wang et al. [13]の次世代シーケンサーによる AMD ゲノムデータを 解析する.ゲノム領域数が 3095656 (= p)次元, AMD 群が 21(= n_1) サンプル, 健常群 (Normal) が 19(= n_2) サンプルであった.図1 で見たように,両群とも SSE モデルが当 てはまる.仮説 (5.1)を有意水準 $\alpha = 0.05$ で検定する.検定方式 (5.3)の検定統計量を 計算すると⁹, $\hat{T}_*/\sqrt{\hat{K}_*} = 23.44 > z_{0.05} = 1.65$ となり,有意水準 5% で H_0 は棄却され る.SSE モデル(非スパース性)を考慮し,データ変換法を使って非スパースなノイズ を除去したことで潜在情報が浮き彫りとなり,2群間の平均の差異が検出できた.

ちなみに,両群にスパース性を仮定してNSSEモデルのもとで検定統計量を計算すると, $T/\sqrt{\hat{K}} = 1.01 < z_{0.05}$ となり, H_0 は棄却されない.しかし,この結論は統計的に保証されるものではない.なぜならば,NSSEモデル(スパース性の仮定)が妥当ではなく,検定統計量に高次元漸近正規性が成立していないからである. $K = O(\lambda_{11}^2/n_1^2 + \lambda_{21}^2/n_2^2)$ に注意すれば,Tの潜在情報である $E(T) = \|\mu_1 - \mu_2\|^2$ が非スパースなノイズとなる最大固有値に押しつぶされ,2群間の平均の差異を検出できなかったことが分かる.

次に,AMD 群と健常群 (Normal)の差異の原因となる変数 (ゲノム領域)を絞り込む. いま, $\mu_1 - \mu_2 = (\mu_{o1}, ..., \mu_{op})^T$, $\overline{x}_1 - \overline{x}_2 = (\overline{x}_{o1}, ..., \overline{x}_{op})^T$ とおく.有意な変数の集合を $D = \{j \mid |\mu_{oj}| > 0\}$ とおく. $\overline{x}_{o1}, ..., \overline{x}_{op}$ を $|\overline{x}_{o(1)}| \ge \cdots \ge |\overline{x}_{o(p)}|$ と絶対値の大きい順に並 び替え,m個の変数からなる $\hat{D}_m = \{j \mid |\overline{x}_{oj}| \ge |\overline{x}_{o(m)}|\}$ でDを推定する.適当なmと

⁸ Aoshima and Yata [4] の定理 5 を参照のこと.

 $^{{}^9}k_i$ の推定については, [4]のS2.2節を参照のこと.このデータセットは, $(\hat{k}_1,\hat{k}_2)=(7,7)$ であった.

適当な正則条件のもと, $P(D \subseteq \widehat{D}_m) \to 1 \ (p \to \infty)$ となる¹⁰.しかし,この変数選択はスパース性を当てにしたものである.

そこで,非スパース性を考慮し,データ変換法を使って $\overline{x}_{i*} = A_i \overline{x}_i$ とする.(4.3)式の \tilde{h}_{ij} と(4.4)式の $\tilde{x}_{i\ell j}$ を用いて, \overline{x}_{i*} を

$$ar{ar{x}}_{i*} = \overline{m{x}}_i - \sum_{j=1}^{\kappa_i} ar{ar{x}}_{ij} ar{m{h}}_{ij}$$
 (ただし, $ar{ar{x}}_{ij} = \sum_{\ell=1}^{n_i} ar{x}_{i\ell j}/n_i$ である)

で推定する. $\overline{\tilde{x}}_{1*} - \overline{\tilde{x}}_{2*} = (\overline{\tilde{x}}_{o1}, ..., \overline{\tilde{x}}_{op})^T と \mathbf{b}$, $\overline{\tilde{x}}_{o1}, ..., \overline{\tilde{x}}_{op} \mathbf{\epsilon} |\overline{\tilde{x}}_{o(1)}| \ge \cdots \ge |\overline{\tilde{x}}_{o(p)}| と並び$ 替え, m個の変数からなる $\widetilde{D}_m = \{j | |\overline{\tilde{x}}_{oj}| \ge |\overline{\tilde{x}}_{o(m)}|\}$ でDを推定する¹¹.

結果を見てみよう. $\hat{D}_m \geq \tilde{D}_m$ を用いて, 3095656 個の変数からm = 10, 100, 1000 個の変数をそれぞれ選択した. 変数選択の精度を高次元 PCA で確認する. x_{ij} において, \hat{D}_m の変数のみを利用したm次元データを \hat{x}_{ij} , \tilde{D}_m の変数のみを利用したm次元データを \hat{x}_{ij} , \hat{D}_m の変数のみを利用したm次元データを \hat{x}_{ij} とし, AMD 群と健常群 (Normal)を混合させたデータセットを $\hat{X} = (\hat{x}_{11}, ..., \hat{x}_{1n_1}, \hat{x}_{21}, ..., \hat{x}_{2n_2})$, $\tilde{X} = (\tilde{x}_{11}, ..., \tilde{x}_{1n_1}, \tilde{x}_{2n_2})$ とおく. Yata and Aoshima [17] は, 規準化した主成分スコアを用いることで,高次元混合データのクラスタリングが可能であることを証明した.2 群間の差異に有意な変数を正しく選択できていれば, 理論的には,主成分スコアを使って2 群を精度良く分離できるはずである.そこで, \hat{X} と \tilde{X} の各々について, $n_1 + n_2 = 40$ 個のサンプルに対する(規準化した)第1主成分スコアと第2主成分スコアを計算し,図5 にプロットした.データ変換法を使って \tilde{D}_m による変数選択は,2群を精度よく分離できていることが見てとれる.データ変換法を使って非スパースなノイズを除去することで,3095656 個の変数から2群の差異に有意な100 個の変数を絞り込めたのである.

非スパースモデリングの一幕をお見せした.データ解析に要した時間は,モバイル PCで5分程度である.高い計算処理能力と大標本を前提とする深層学習が,経済的に も技術的にも環境的にも敷居の高いテクノロジーであるのに対して,高次元統計解析 による非スパースモデリングは,いわば庶民的な町工場の技術である.非スパースモ デリングを使うことで,非スパースなノイズの影響を取り除くことができ,超高次元 データもたった数十程度の小標本で処理でき,モバイルPCでも高速かつ高精度に解析 することができるのである.

参考文献

- [1] Aoshima, M. and Yata, K. (2011). Two-stage procedures for high-dimensional data, *Sequential Anal*. (*Editor's special invited paper*), **30**, 356–399.
- [2] 青嶋 誠, 矢田和善 (2019). 高次元の統計学, 共立出版, 東京.
- [3] Aoshima, M. and Yata, K. (2015). Asymptotic normality for inference on multisample, high-dimensional mean vectors under mild conditions, *Meth. Comput. Appl. Probab.*, **17**, 419–439.
- [4] Aoshima, M. and Yata, K. (2018). Two-sample tests for high-dimension, strongly spiked eigenvalue models, *Statist. Sinica*, **28**, 43–62.
- [5] Aoshima, M. and Yata, K. (2019). Distance-based classifier by data transformation for highdimension, strongly spiked eigenvalue models, Ann. Inst. Statist. Math., 71, 473–503.
- [6] Chen, S.X. and Qin, Y.-L. (2010). A two-sample test for high-dimensional data with applications to gene-set testing, *Ann. Statist.*, 38, 808–835.

¹⁰ 例えば, Aoshima and Yata [1] の5節を参照のこと.

 $^{{}^{11}\}widetilde{D}_m$ は \widehat{D}_m よりも緩い条件で一致性をもつ.



図 5. m = 10, 100, 1000 個の変数を選択した AMD 群と健常群 (Normal)の混合データに 対する第1と第2の主成分スコア. \hat{D}_m : データ変換なし(左), \tilde{D}_m : データ変換法(右).

- [7] Donoho, D.L. (2006). Compressed sensing, IEEE Trans. Information Theory, 52, 1289–1306.
- [8] Ishii, A., Yata, K. and Aoshima, M. (2022). Geometric classifiers for high-dimensional noisy data, J. Multivariate Anal. (Editor's invited paper), 188, 104850.
- [9] Johnstone, I.M. (2001). On the distribution of the largest eigenvalue in principal components analysis, *Ann. Statist.*, **29**, 295–327.
- [10] Paul, D. (2007). Asymptotics of sample eigenstructure for a large dimensional spiked covariance model, *Statistica Sinica*, **17**, 1617–1642.
- [11] Takeuchi, T, Yata, K., Aoshima, M., et al. (2022). High dimensional statistical analysis and its application to ALMA map of NGC 253, *arXiv preprint*, arXiv:2203.04535.
- [12] von Roemeling, C., Radisky, D., Marlow, L., et al. (2014). Neuronal pentraxin 2 supports clear cell renal cell carcinoma by activating the AMPA-selective glutamate receptor-4, *Cancer Res.*, 74, 4796–4810.
- [13] Wang, J., Zibetti, C., Shang, P., et al. (2018). ATAC-seq analysis reveals a widespread decrease of chromatin accessibility in age-related macular degeneration, *Nat. Commun.*, **9**, 1364.
- [14] Yata, K. and Aoshima, M. (2009). PCA consistency for non-Gaussian data in high dimension, low sample size context, *Comm. Statist. Theory Methods, Special Issue Honoring Zacks, S. (ed. Mukhopadhyay, N.)*, **38**, 2634–2652.
- [15] Yata, K. and Aoshima, M. (2012). Effective PCA for high-dimension, low-sample-size data with noise reduction via geometric representations, J. Multivariate Anal., 105, 193–215.
- [16] Yata, K. and Aoshima, M. (2013). PCA consistency for the power spiked model in highdimensional settings, J. Multivariate Anal., 122, 334-354.
- [17] Yata, K. and Aoshima, M. (2020). Geometric consistency of principal component scores for high-dimensional mixture models and its application, *Scand. J. Stat.*, 47, 899–921.