

ランダム行列理論によるノイズ推定を用いた統計的データ解析

小林由香 (指導教員：吉田裕亮)

1 はじめに

ランダム行列とは、20世紀初めの Wishart らによる数理統計学の研究に起源をもっている。その後、Wigner が原始核物理学に応用してから理論物理学の研究対象としても注目されるようになった。数学の分野では最近、非可換確率論における確率変数を漸近的に表現するなど、発展を続けている。また、金融工学における金融相関行列のモデルになることや、ランダムグラフの研究との合流により、複雑ネットワークのモデルとしても現れるなど、ランダム行列の応用分野は広がっている。本研究では、ノイズを含む二点間の距離が推定されたデータにおいて、データ距離行列 (三角行列) に含まれる雑音量の推定を行うために、ランダム行列理論を用いた応用を考える。特に MDS (Multi Dimensional Scaling) 法において用いられる内積行列から、雑音量の推定を行う手法を提案する。

2 半円分布

ランダム行列とは、一般に確率変数を要素にもつ行列である。特に成分が、 $a_{ii} = 0$ かつ $a_{ij} = a_{ji}$ が成り立ち、確率変数の族 $a_{ij} (1 \leq i < j \leq N)$ が独立に $N(0, \sigma^2)$ に従うランダム行列は、ガウス型直交ランダム行列とも呼ばれ、その固有値経験分布は行列のサイズ N を $N \rightarrow \infty$ とするとき、以下の中心 0 の半円分布に収束されることが知られている。

一般に、中心が $m \in \mathbb{R}$ 、半径 $r > 0$ の半円分布 $\omega_{m,r}$ とは、密度関数 $P(t)$ が

$$P(t) = \frac{2}{\pi r^2} \sqrt{r^2 - (t - m)^2} \quad (1)$$

で与えられる分布である。

3 モーメント

$f(x)$ を連続確率変数 X の密度関数とすると、原点まわりの k 次モーメント m_k は、

$$m_k = \int_{-\infty}^{\infty} x^k f(x) dx \quad (2)$$

と表せ、これらの値は確率分布の特徴を与える。また、先に述べた Marcenko-Pastur 分布のモーメントは、

$$m_k = \frac{2k!}{k!(k+1)!} m_1^k \quad (3)$$

で与えられる。

本研究では、観測データからの標本モーメント列を理論値と比較することにより、最適なノイズ部の推定を行うことにする。これにより、目視でスペクトルを比較するより、定量的な推定が行える。

半円分布 $\omega_{m,r}$ の k 次モーメント m_k は

$$\begin{aligned} m_k &= \frac{2}{\pi r^2} \int_{-r}^r x^k \sqrt{r^2 - x^2} dx \\ &= \begin{cases} 0, & (k = 2m + 1), \\ \frac{(2m)!}{m!(m+1)!} \left(\frac{r^2}{4}\right)^m, & (k = 2m), \end{cases} \end{aligned} \quad (4)$$

で与えられる。ただし、 $C_m = \frac{(2m)!}{m!(m+1)!}$ は m 次 Catalan 数と呼ばれる組合せ論でよく現れる数である。

4 MDS (Multi Dimensional Scaling) 法

MDS 法とは幾つかの対象があり、任意の 2 つの対象間の距離が Euclid 距離として推定されている場合、Young-Householder の定理をもとに対象を Euclid 空間の点として布置 (位置付ける) 方法である。 n 個の対称 O_1, O_2, \dots, O_n のうち、 O_n を原点とし、残りの $n-1$ 個の対称を終点とする (列) ベクトルを x_1, x_2, \dots, x_{n-1} とする。さらに、対象 O_i の第 j 軸の座標値を x_{ij} とし、行列 $X = x_{ij}$ を定義する。行列 X は、 $n-1$ 個の対象の埋め込まれる空間の次元数分の座標値を対象ごと各行に並べたものである。この時、この空間の次元数 r は、

$$B = (b_{ij}) = XX^t \quad (5)$$

と置けば、 $r = \text{rank} X = \text{rank} B$ となる。行列 B は、 $n-1$ 個の対象間の内積を要素とする行列で、内積行列である。

$$B = \begin{pmatrix} (\vec{x}_1 | \vec{x}_1) & (\vec{x}_1 | \vec{x}_2) & \dots & (\vec{x}_1 | \vec{x}_{n-1}) \\ (\vec{x}_2 | \vec{x}_1) & (\vec{x}_2 | \vec{x}_2) & \dots & (\vec{x}_2 | \vec{x}_{n-1}) \\ \vdots & \vdots & \ddots & \vdots \\ (\vec{x}_{n-1} | \vec{x}_1) & (\vec{x}_{n-1} | \vec{x}_2) & \dots & (\vec{x}_{n-1} | \vec{x}_{n-1}) \end{pmatrix} \quad (6)$$

とする。つぎに、ベクトル $v_{ij} = x_j - x_i$ を定義すると、対象 O_i と O_j 間のユークリッド距離 d_{ij} を Euclid ノルムを用いて表現すれば

$$\begin{aligned} \|v_{ij}\|^2 &= d_{ij}^2 = v_{ij}^t v_{ij} = (x_j - x_i)^t (x_j - x_i) \\ &= \|x_i\|^2 + \|x_j\|^2 - 2b_{ij} \end{aligned} \quad (7)$$

と書ける。 $\|x_i\|^2 = d_{ni}^2$ より

$$b_{ij} = \frac{1}{2}(d_{ni}^2 + d_{nj}^2 - d_{ij}^2) \quad (8)$$

である。つまり、対象間の距離情報 (距離の二乗) から内積行列が構成されることになる。行列 B は対象 O_n を原点とした場合の内積行列であり、どの対象を原点

に取るかで行列 B は一般に異なってくる。誤差を含む場合には特に、原点の取り方には工夫が必要であり、どの対象 (点) も対等に扱うことなので、 n 点の重心を原点に選ぶことが多い。 n 個の対象行列の重心を原点に選んだときの各点の位置ベクトルの内積行列 B_c を求めるには、対象間 2 乗距離を要素とする $n \times n$ の行列 $T = \{d_{ij}^2\}$ を作る。次に、中心化行列

$$J_n = I - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^t \quad (9)$$

を作り、ただし、 I_n は n 次単位行列 $\mathbf{1}_n = (1, 1, \dots, 1)$ 、 J_n は

$$J_n = \begin{pmatrix} 1 - \frac{1}{n} & -\frac{1}{n} & \dots & -\frac{1}{n} \\ -\frac{1}{n} & 1 - \frac{1}{n} & \dots & -\frac{1}{n} \\ \vdots & \vdots & \ddots & \vdots \\ -\frac{1}{n} & -\frac{1}{n} & \dots & 1 - \frac{1}{n} \end{pmatrix} \quad (10)$$

となる。これを用いて、次のような計算により、重心を原点とする位置ベクトルから構成される。内積行列 B_c は

$$B_c = -\frac{1}{2} J_n T J_n^t \quad (11)$$

と求められる。これは Young-Householder 変換と呼ばれている。MDS では、 B_c スペクトル分解を行うことにより、布置座標を求めることになる。

本研究では、この内積行列の対象の配置構造以外のスペクトル (固有値) として雑音部を推定する手法を提案する。すなわち、2 点間の距離が推定された r 次元に埋め込まれたデータより内積行列 B_c を求める。 r 個の優固有値を除いたスペクトルからモーメント法を用いて、半円分布と比較することにより雑音量の推定を行う。これは、半円分布の密度関数の半径を推定することと同等である。

5 実験例

	2 乗	4 乗	6 乗	8 乗	10 乗
理論値	1	2	5	14	42
モーメント	1.02	2.148	5.657	16.671	52.537

表 1: モーメント

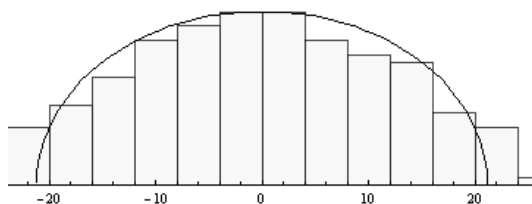


図 1: B_n のスペクトルと Wigner 分布

まず、ノイズを含み 2 点間の距離が推定されている、2 次元の構造データを 200 点用意する。この点を用い

て、距離行列を生成する。距離行列には、四捨五入を行うことにより、ノイズを混入する。これらの距離行列から、先に述べた MDS 法により内積行列 B_c を求める。 B_c のスペクトルから、2 個の優固有値を除き、各次数のモーメントを計算する (モーメントの各値は表 1 のようになる)。この場合は、 $r = 21.2$ と推定され、推定された半径を元に半円分布と B_c のヒストグラムの密度と重ね合わせると、図 1 のようになり、モーメント法により半径、すなわち雑音量、の推定ができたと考えられる。

6 実データによる実験例

実データとして、京王電鉄の運賃表を用いて実験を行った。構造データは、 70×70 行列である。本手法を用いて実験を行った結果は、表 2、図 2 で表わされる。結果、半径は $r = 250$ と推定された。また、標準偏差として 2 点間の距離に 14.9 円程度ずつノイズが含まれていることが分かった。

	2 乗	4 乗	6 乗	8 乗	10 乗
理論値	1	2	5	14	42
モーメント	0.966	1.937	4.843	13.489	40.202

表 2: 実データのモーメント

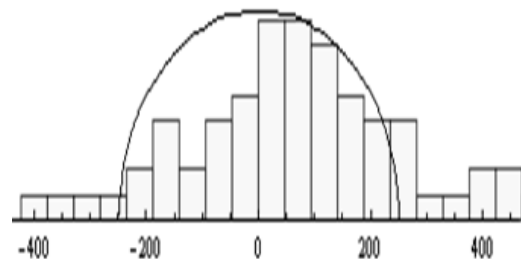


図 2: 実データ B_n のスペクトルと Wigner 分布

7 まとめ

データ距離行列とランダム行列理論の Wigner 半円分布を用いて、雑音部の推定を行った。雑音部の推定は目視ではなくモーメント法を用いて、定量的に行う手法を提案した。今後の課題として、距離行列から生成される様々なタイプの行列とランダム行列との関連によりノイズ推定の精度を上げることも可能と考えられる。

参考文献

- [1] 広尾太郎, 「ランダム行列の広がり」, 数理科学, 第 45 巻 2 号, サイエンス社, (2007 年)
- [2] 赤穂昭太郎:カーネル多変量解析 非線形データ解析の新しい展開 (2008 年)