

ランダム行列を用いたスペクトラルクラスタリング

茨木 志織^{†1}

カーネル法を用いたクラスタリングの1つにスペクトラルクラスタリングがある。本研究では、ノイズを含んだデータに対し、そのノイズをランダム行列の手法を用いて取り除くことにより、スペクトラルクラスタリングの精度を上げる手法を提示する。カーネルにはガウスカーネルを用いて、Wishart 行列の固有値分布とガウスカーネルで写像した特徴空間における内積行列の固有値分布が等価であることを利用し、ノイズを推定する。

Spectral clustering using random matrices

SHIORI IBARAKI^{†1}

The spectral clustering is known as one of methods for clustering by using kernel technique. In this study, we shall show a method of improving spectral clustering by removing the noise from data with the theory of random matrices. We will use the Gaussian kernel and estimate the noise since the spectral distribution of Wishart matrix is equivalent to one of the matrix constituted from inner products of Gaussian kernel.

1. はじめに

データをまとまりごとに集めてグループ分けをすることは、クラスタリングと呼ばれている。クラスタリングは、階層的手法と非階層的手法の2つに大きく分類され、非階層的手法の代表例に K -平均法がある。 K -平均法は非常に有効なクラスタリング手法ではあるが、反復演算を必要とする点と、収束解が必ずしも目的関数を最適にするものではないという欠

点がある。スペクトラルクラスタリングでは、クラスタリングの問題を固有値問題として定式化することによって、これらの問題点のないアルゴリズムを構成することが出来る。本研究ではノイズを含んだデータに対し、そのノイズをランダム行列の手法を用いて除去することにより、スペクトラルクラスタリングの精度を上げる手法を提示する。

2. ランダム行列理論

一般に、ランダム行列とは確率変数を要素に持つ行列であり、その代表例として Wishart 行列が挙げられる。

2.1 Wishart 行列

各成分が独立に $N(0, 1)$ の標準正規分布に従う変数をもつ $n \times p$ の行列を C とする。このランダム行列 C から

$$S = \frac{1}{n} C^T C \quad (1)$$

で求められる $n \times n$ 対称ランダム行列 S を Wishart 行列という。 $p/n = \lambda$ を保ちながら、 $n \rightarrow \infty, p \rightarrow \infty$ の極限をとると、Wishart 行列 S の固有値の経験分布は、 $\lambda_{min} \leq t \leq \lambda_{max}$ のときに以下の確率密度関数に収束することが知られている。

$$p(t) = \frac{1}{2\pi} \frac{\sqrt{-(t - \lambda_{max})(t - \lambda_{min})}}{\lambda t}, \quad \lambda_{min}^{max} = (1 \pm \sqrt{\lambda})^2 \quad (2)$$

また、このような確率密度関数を持つ分布は Marcenko-Pastur 分布と呼ばれている。

2.2 ガウスカーネル

変数の集合の二つの要素 x, x' に対し、カーネル関数 $k(x, x')$ は x, x' それぞれの特徴ベクトルどうしの内積

$$k(x, x') = \phi(x)^T \phi(x') \quad (3)$$

として定義される。カーネルには様々なものがあるが、その中でもガウスカーネル

$$k(x, x') = \exp(-\beta \|x - x'\|^2) \quad (4)$$

を用いて特徴空間に写像した行列は、相関行列と同じような振る舞いをする事が知られている。ここで、 $\|\cdot\|^2$ は通常のユークリッド2乗距離で、 $\beta \in R$ は適当なパラメータである。

また、Wishart 行列の固有値分布と、ガウスカーネルで写像した特徴空間における内積行列の固有値分布は等価であると知られており、これによりガウスカーネル行列におけるノイズ部に相当する固有値分布も Marcenko-Pastur 分布と同様の性質を持つことがわかる。

^{†1} お茶の水女子大学大学院 人間文化創成科学研究科
Graduate School of Humanities and Sciences, Ochanomizu University

3. スペクトラルクラスタリング

スペクトラルクラスタリングでは、サンプル点をグラフ構造として考え、各頂点がサンプル点で、枝にはサンプル点同士の近さを表す重みがついているとする。例えば、サンプル点を2つのグループに分けると、それに伴いグラフも2分割される。分割されたグループ間を結ぶ枝のことを分割のカットと呼ぶ。グループ内は出来るだけ近いものどうしが集まり、グループ間は遠く離れていることが望ましいので、カットの重みの合計が小さくなるようにグループ分けを行う。式で表すと、以下ようになる。

$$\min_{\beta} \sum_{i,j} K_{ij}(\beta_i - \beta_j)^2 = 2^T \beta P \beta, \quad \beta_i = \pm 1 \quad (5)$$

ここで、 P は対角行列 Λ を $\Lambda_{ii} = \sum_{j=1}^n K_{ij}$ として、 $P = \Lambda - K$ と書ける。 β は2値ベクトルという制約がある。それは整数計画問題と呼ばれ、一般には解くのが困難である。そこで、整数という制約を取り払って任意の実数ベクトルに、 ${}^T \beta \Lambda \beta = 1$ という条件の下、制約を緩めることにより推定を行う。この場合、最小固有値0が存在するが、これはすべてのサンプルを1つにまとめてしまうという意味のない解のため、実際には2番目以降の固有ベクトルの成分符号に基づいてクラスタリングを行う。

4. 提案手法

本研究ではランダム行列の理論を用いてスペクトラルクラスタリングを以下のように行うことを提案する。

- (1) サンプルデータのガウスカーネル行列 K を構成し、そのスペクトルを求める。
- (2) 求めたスペクトル分布と Marcenko-Pastur 分布を比較し、ノイズに相当すると考えられる部分を取り除く。その残ったスペクトルを用いて、カーネル行列 K' を再構成する。
- (3) この新たな K' を与えられたカーネル行列と見なし、スペクトラルクラスタリングを行う。

5. 実験例

図1のような、線形で分けることの出来ない2群データを用意する。1群が300点、2群が250点の、合計550点のサンプルデータとなっている。これらをスペクトラルクラスタリングで2つにクラス分けする。

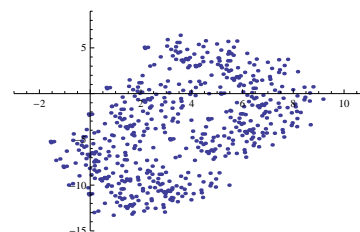


図1 サンプルデータ
 Fig.1 sample data

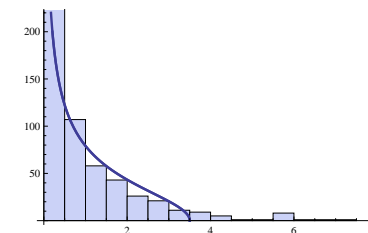


図2 ヒストグラムと Marcenko-Pastur 分布
 Fig.2 histogram and Marcenko-Pastur distribution

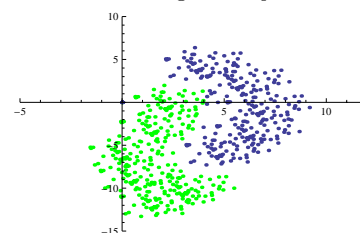


図3 ノイズを除いた場合
 Fig.3 with reduction

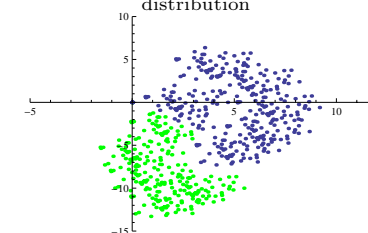


図4 除いていない場合
 Fig.4 without reduction

5.1 ノイズ部スペクトルの推定

サンプルデータのガウスカーネル行列のヒストグラムと Marcenko-Pastur 分布を重ね合わせると、図2のようになった。この分布から外れた部分がこのデータにおける取り出すべき情報、つまりデータからノイズを除いた構造部であると考えた。今回の実験では、この構造部のスペクトルは27個であり、このスペクトルのみを用いてカーネル行列を再編成し、スペクトラルクラスタリングを行った。

5.2 結果

図3の結果が得られた。比較の為、右に通常のスペクトラルクラスタリングの結果(図4)を並べた。図3は、ほぼ設定した通りにクラス分けが出来ていることが読み取れる。

6. まとめ

ランダム行列理論を用いて、ノイズを除去することでスペクトラルクラスタリングの精度を高めることが可能である場合が確認された。また今後の課題として、Marcenko-Pastur 分布の適切なパラメータの推定を行いたいと考えている。