カーネル主成分分析に基づく高次元データのクラスタリングとチューニング

筑波大学・数理物質科学 中山 優吾 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

本講演では、高次元データのクラスタリングを考える。高次元 PCA によって潜在空間を推定し、そこへデータを射影させることで高精度なクラスタリング手法を提案する。2 つの d 次元分布を Π_1,Π_2 と名付け、それぞれ平均 μ_1 、 μ_2 と、共分散行列 Σ_1 、 Σ_2 をもつと仮定する。いま、データはこの母集団から n (≥ 2) 個のデータを無作為に抽出し、データ行列を $X=[x_1,...,x_n]$ とする。 $n_i=\#\{j|x_j\in\Pi_i,\ j=1,...,n\}$ とする。ここで、#A は集合 A の要素の個数とし、 $n_i\geq 1$ とする。簡単のため、以下を仮定する。

$$tr(\Sigma_1) \le tr(\Sigma_2), \quad x_j \in \Pi_1 \text{ for } j = 1, ..., n_1 \text{ and } x_j \in \Pi_2 \text{ for } j = n_1 + 1, ..., n.$$

本講演では、ガウシアンカーネルを用いたカーネル PCA の漸近的性質を考えることでクラスタリング手法を与える. $n \times n$ のグラム行列 K の (j,j') 成分を

$$k(x_i, x_{i'}) = \exp(-\|x_i - x_{i'}\|^2 / \gamma) \quad (\gamma > 0)$$

とする. \boldsymbol{I}_n を n 次の単位行列, $\boldsymbol{1}_n = (1,...,1)^T$ として, $\boldsymbol{P}_n = \boldsymbol{I}_n - n^{-1} \boldsymbol{1}_n \boldsymbol{1}_n^T$ とする. 中心化グラム行列を $\boldsymbol{K}_0 = \boldsymbol{P}_n \boldsymbol{K} \boldsymbol{P}_n$ とし, \boldsymbol{K}_0 の固有値分解を $\boldsymbol{K}_0 = \sum_{i=1}^{n-1} \hat{\lambda}_i \hat{\boldsymbol{u}}_i \hat{\boldsymbol{u}}_i^T$ ($\hat{\boldsymbol{u}}_i = (\hat{u}_{i1},...,\hat{u}_{in})^T$, $\|\hat{\boldsymbol{u}}_i\|^2 = 1$) とする. \boldsymbol{x}_i の (基準化した) 第 i 主成分スコアを $s_{ij} = \sqrt{n} \hat{u}_{ij}$ とおく.

 $\kappa_{\mu} = \exp(-\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|^2/\gamma), \ \kappa_{\Sigma} = \exp(-|\mathrm{tr}(\boldsymbol{\Sigma}_1) - \mathrm{tr}(\boldsymbol{\Sigma}_2)|/\gamma). \ \Delta_{\kappa} = 1 + \kappa_{\Sigma}^2 - 2\kappa_{\mu}\kappa_{\Sigma}$ とおき, 以下を仮定する.

$$\limsup_{d \to \infty} \frac{(1 - \kappa_{\Sigma}^2)n}{\Delta_{\kappa} n_1 n_2} < 1 \tag{1}$$

第 1 主成分スコア s_{1j} , j=1,...,n について、以下の一致性が成り立つ。

定理 1. 適当な正則条件と (1) のもと, $d \to \infty$ で以下が成り立つ.

$$s_{1j} = \begin{cases} \sqrt{n_2/n_1} + o_P(1) & when \ j = 1, ..., n_1, \\ -\sqrt{n_1/n_2} + o_P(1) & when \ j = n_1 + 1, ..., n. \end{cases}$$
 (2)

定理 1 から第 1 主成分スコアの符号で高次元データを分類できることがわかる. Yata and Aoshima [1] は線形カーネルに対して、平均ベクトル間の距離 $\|\mu_1 - \mu_2\|^2$ が十分大きければ、(1) と同様な主成分スコアの一致性を得られることを示した. ガウシアンカーネルを用いたカーネル PCA は $\mu_1 = \mu_2$ のときでさえも、共分散行列間の距離によって (2) を与えることができる.

当日は、従来の PCA とガウシアンカーネルを用いたカーネル PCA の理論的な比較を与える. さらに、(1) を満たすような γ の選択法を与え、その性能を数値実験と実データ解析を用いて検証する.

[1] Yata, K., Aoshima, M. (2015). Principal component analysis based clustering for high-dimension, low-sample-size data. arXiv:1503.04525.