Clustering bounds for hidden Partition Models

Nicolas Verzelen INRAE, Montpellier

Mathematical Statistics in the Information Age

Clustering arises in various contexts



Clustering features



Clustering graphs



Topic of the talk

- Introduce several hidden partition problems.
- fresh view on some classical clustering algorithms.
- Some recovery error bounds for Gaussian Mixtures

Main Message

K-means and its relaxations (and corrections) are versatile and near optimal tools.



2 Omnibus Clustering Algorithms

- K-means criterion
- Lloyd's Algorithm
- Peng and Wei's SDP
- Spectral Clustering

3 Recovery bounds for (sub)Gaussian Mixtures

- Loss functions
- Minimax lower bound
- Recovery bounds for Peng and Wei' SDP

Mixture of Gaussian vectors [Pearson('1895)]

Partition

 $G^* = \{G^*_1, \ldots, G^*_K\}$ of [n]

Mixture of Gaussian variables (conditional)

 $X_1, \ldots, X_n \in \mathbb{R}^p$ are independent with $X_a \sim \mathcal{N}(\theta_k, \mathbf{\Sigma}_k)$ if $a \in G_k^*$

The observations are gathered in
$$\mathbf{X} = \begin{bmatrix} X_1 \\ \cdots \\ X_n \end{bmatrix} \in \mathbb{R}^{n \times p}$$



Mixture of Gaussian vectors [Pearson('1895)]

Partition

 $G^* = \{G^*_1, \ldots, G^*_K\}$ of [n]

Mixture of Gaussian variables (conditional)

 $X_1,\ldots,X_n\in\mathbb{R}^p$ are independent with $X_a\sim\mathcal{N}(heta_k,\mathbf{\Sigma}_k)$ if $a\in G_k^*$



Objective : recovering G^* from **X** (θ and **\Sigma** are unknown but K is known)

Mixture of Gaussian vectors [Pearson('1895)]

Partition

 $G^* = \{G^*_1, \ldots, G^*_K\}$ of [n]

Mixture of Gaussian variables (conditional)

 $X_1,\ldots,X_n\in\mathbb{R}^p$ are independent with $X_a\sim\mathcal{N}(heta_k,\mathbf{\Sigma}_k)$ if $a\in G_k^*$



Objective : recovering G^* from **X** (θ and **\Sigma** are unknown but K is known)

Clustering Problem \neq Parameter Estimation

$$\begin{array}{l} \text{Membership Matrix } \mathbf{A} \in \mathbb{R}^{n \times K} \text{ defined by } \mathbf{A}_{ak} = \mathbf{1}_{a \in G_k} = \mathbf{\Pi} \begin{bmatrix} 1 & 0 & 0 \\ & \dots & \\ 0 & 1 & 0 \\ & \dots & \\ 0 & 0 & 1 \end{bmatrix} \\ \text{Mean component Matrix } \mathbf{\Theta} = \begin{bmatrix} \theta_1 \\ \dots \\ \theta_K \end{bmatrix}.$$

Population Version : $\mathbb{E}[\mathbf{X}] = \mathbf{A} \boldsymbol{\Theta}$

$$\mathbf{X} = \mathbf{A}\mathbf{\Theta} + \mathbf{E} =$$
"Signal" + "Noise"

Holland et al. ('83), Abbe('17),...,

 $\mathbf{X}=\mathsf{adjacency}$ matrix of an undirected graph $\in \{0,1\}^{n imes n}$.

Let $\mathbf{Q} \in [0, 1]_{sym}^{K \times K}$

(conditional) SBM

The graph is generated by a SBM with partition G^* and matrix ${\bf Q}$ if ${\bf X}_{ab}$ with a < b are independent and

$$\mathbb{P}[\mathbf{X}_{ab} = 1] = \mathbf{Q}_{jk}$$
 for any $a \in G_j^*$ and $b \in G_k^*$,

Holland et al. ('83), Abbe('17),...,

 $\mathbf{X}=\mathsf{adjacency}$ matrix of an undirected graph $\in \{0,1\}^{n imes n}$.

Let $\mathbf{Q} \in [0, 1]_{sym}^{K \times K}$

(conditional) SBM

The graph is generated by a SBM with partition G^* and matrix ${\bf Q}$ if ${\bf X}_{ab}$ with a < b are independent and

$$\mathbb{P}[\mathbf{X}_{ab} = 1] = \mathbf{Q}_{jk}$$
 for any $a \in G_j^*$ and $b \in G_k^*$



Holland et al. ('83), Abbe('17),...,

 $\mathbf{X}=\mathsf{adjacency}$ matrix of an undirected graph $\in \{0,1\}^{n imes n}$.

Let $\mathbf{Q} \in [0, 1]_{sym}^{K \times K}$

(conditional) SBM

The graph is generated by a SBM with partition G^* and matrix ${\bf Q}$ if ${\bf X}_{ab}$ with a < b are independent and

$$\mathbb{P}[\mathbf{X}_{ab} = 1] = \mathbf{Q}_{jk}$$
 for any $a \in G_j^*$ and $b \in G_k^*$



Holland et al. ('83), Abbe('17),...,

 $\mathbf{X}=\mathsf{adjacency}\ \mathsf{matrix}\ \mathsf{of}\ \mathsf{an}\ \mathsf{undirected}\ \mathsf{graph}\in\{0,1\}^{n imes n}$.

Let $\mathbf{Q} \in [0, 1]_{sym}^{K \times K}$

(conditional) SBM

The graph is generated by a SBM with partition G^* and matrix ${\bf Q}$ if ${\bf X}_{ab}$ with a < b are independent and

$$\mathbb{P}[\mathbf{X}_{ab} = 1] = \mathbf{Q}_{jk}$$
 for any $a \in G_j^*$ and $b \in G_k^*$.



Objective recovering G^* from **X** (**Q** is unknown)

Population Version : $\mathbb{E}[\mathbf{X}] = \mathbf{A}\mathbf{Q}\mathbf{A}^T - \text{Diag}(\mathbf{A}\mathbf{Q}\mathbf{A}^T) \approx \mathbf{A}[\mathbf{Q}\mathbf{A}^T]$

Outline

1 Hidden Partition problems



2 Omnibus Clustering Algorithms

- K-means criterion
- Lloyd's Algorithm
- Peng and Wei's SDP
- Spectral Clustering

3 Recovery bounds for (sub)Gaussian Mixtures

- Loss functions
- Minimax lower bound
- Recovery bounds for Peng and Wei' SDP

Algorithms do not depend on a particular model. To provide some intuition :

Mixture of Gaussian variables (conditional)

 $X_1,\ldots,X_n\in\mathbb{R}^p$ are independent with $X_a\sim\mathcal{N}(heta_k,\mathbf{\Sigma}_k)$ if $a\in G_k^*$

Signal + Noise decomposition

 $\mathbf{X} = \mathbf{A} \boldsymbol{\Theta} + \mathbf{E}$



2 Omnibus Clustering Algorithms K-means criterion

- Lloyd's Algorithm
- Peng and Wei's SDP
- Spectral Clustering

- Loss functions
- Minimax lower bound
- Recovery bounds for Peng and Wei' SDP

Maximum Likelihood Estimation

For the Gaussian Model, MLE

$$\hat{G}^{MLE} \in \arg\min_{G} \sum_{k=1}^{K} \min_{\boldsymbol{\Sigma}_{k} \in \mathcal{S}_{p}^{+}} \min_{\boldsymbol{\theta}_{k} \in \mathbb{R}^{p}} \sum_{a \in G_{k}} \left((X_{i} - \boldsymbol{\theta}_{k})^{T} \boldsymbol{\Sigma}_{k}^{-1} (X_{i} - \boldsymbol{\theta}_{k}) + \log(\det(\boldsymbol{\Sigma}_{k})) \right)$$

Two Difficulties :

- Computational : (In principle) requires to scan over the space of partitions (size of order Kⁿ/K!)
- **Statistical** : Estimation of Σ_k unstable for large p.

Maximum Likelihood Estimation

For the Gaussian Model, MLE

$$\hat{G}^{MLE} \in \arg\min_{G} \sum_{k=1}^{K} \min_{\boldsymbol{\Sigma}_{k} \in \mathcal{S}_{p}^{+}} \min_{\boldsymbol{\theta}_{k} \in \mathbb{R}^{p}} \sum_{a \in G_{k}} \left((X_{i} - \boldsymbol{\theta}_{k})^{T} \boldsymbol{\Sigma}_{k}^{-1} (X_{i} - \boldsymbol{\theta}_{k}) + \log(\det(\boldsymbol{\Sigma}_{k})) \right)$$

Two Difficulties :

- Computational : (In principle) requires to scan over the space of partitions (size of order Kⁿ/K!)
- **Statistical** : Estimation of Σ_k unstable for large p.

ightarrow For the latter, assume in the criterion that $\mathbf{\Sigma}_k = \sigma^2 \mathbf{I}_p$

$$\hat{G} \in \arg\min_{G} \sum_{k=1}^{K} \min_{\theta_k \in \mathbb{R}^p} \sum_{a \in G_k} \|X_a - \theta_k\|_2^2$$

 $\widehat{G}\in \arg\min_{G}\mathsf{Crit}(\mathbf{X},G)$ where

$$\mathsf{Crit}(\mathbf{X},G) = \sum_{k=1}^{K} \sum_{a \in G_k} \|X_a - \overline{X}_{G_k}\|^2 = \frac{1}{2} \sum_{k=1}^{K} \frac{1}{|G_k|} \sum_{a,b \in G_k} \|X_a - X_b\|^2 ,$$

where $\overline{X}_{G_k} = \frac{1}{|G_k|} \sum_{a \in G_k} X_a$

 $\widehat{G}\in \arg\min_{G}\mathsf{Crit}(\mathbf{X},G)$ where

$$\mathsf{Crit}(\mathbf{X},G) = \sum_{k=1}^{K} \sum_{a \in G_k} \|X_a - \overline{X}_{G_k}\|^2 = \frac{1}{2} \sum_{k=1}^{K} \frac{1}{|G_k|} \sum_{a,b \in G_k} \|X_a - X_b\|^2 ,$$

where
$$\overline{X}_{G_k} = \frac{1}{|G_k|} \sum_{a \in G_k} X_a$$

We did not address the computational issues

- There can be many local optima
- In worst-case solving K-means is NP-hard (Mahajan et al. ('09))
- $(1 + \epsilon)$ -approximation is also NP-hard : Finding \hat{G} s.t. $Crit(\mathbf{X}, \hat{G}) \leq (1 + \epsilon) \min_{G} Crit(\mathbf{X}, G)$
- 8-approximation is possible in polynomial time (Kanungo et al. ('04))

Encoding partition learning as a Matrix Problem.

Membership Matrix A $\in \{0,1\}^{n \times K}$ defined by $\mathbf{A}_{ak} = \mathbf{1}_{a \in G_k}$ (or equivalently function $k : [n] \mapsto [K]$) is NOT Identifiable. Why?

Encoding partition learning as a Matrix Problem.

Membership Matrix $\mathbf{A} \in \{0,1\}^{n \times K}$ defined by $\mathbf{A}_{ak} = \mathbf{1}_{a \in G_k}$ (or equivalently function $k : [n] \mapsto [K]$) is at best identifiable up to permutation

A more suitable object : The $n\times n$ partnership matrix ${\bf B}^*={\bf A}({\bf A}^T{\bf A})^{-1}{\bf A}^T$

$$\mathbf{B}^*_{ab} = \left\{ egin{array}{cc} rac{1}{|G^*_k|} & ext{if a and b belong to the same G^*_k} \ 0 & ext{else} \end{array}
ight.$$

Invariant with respect to the group labeling.

Encoding partition learning as a Matrix Problem.

Membership Matrix $\mathbf{A} \in \{0,1\}^{n \times K}$ defined by $\mathbf{A}_{ak} = \mathbf{1}_{a \in G_k}$ (or equivalently function $k : [n] \mapsto [K]$) is at best identifiable up to permutation

A more suitable object : The $n\times n$ partnership matrix ${\bf B}^*={\bf A}({\bf A}^T{\bf A})^{-1}{\bf A}^T$

$$\mathbf{B}^*_{ab} = \left\{ \begin{array}{ll} \frac{1}{|G^*_k|} & \text{ if } a \text{ and } b \text{ belong to the same } G^*_k \\ 0 & \text{ else} \end{array} \right.$$

Invariant with respect to the group labeling.

$$\begin{split} \frac{\text{Properties}:}{\text{tr}[\mathbf{B}^*] = K} \\ \mathbf{B}^* = \mathbf{\Pi}^T \begin{pmatrix} 1 & \mathbf{J}_{|G_1^*|} & 0 & 0 \\ 0 & 1 & \mathbf{J}_{|G_2^*|} & 0 \\ 0 & 0 & 1 & \mathbf{J}_{|G_2^*|} & \mathbf{J}_{|G_1^*|} \\ \end{pmatrix} \mathbf{\Pi} \end{split}$$

Rewriting K-means as a matrix estimation problem

$$\begin{aligned} \mathsf{Crit}(\mathbf{X}, G) &= \frac{1}{2} \sum_{k=1}^{K} \frac{1}{|G_k|} \sum_{a, b \in G_k} \|X_a - X_b\|^2 \\ &= -\sum_k \sum_{a, b \in G_k} \langle X_a, X_b \rangle \frac{1}{|G_k|} + \sum_{a=1}^n \|X_a\|_2^2 \\ &= -\langle \mathbf{X} \mathbf{X}^T, \mathbf{B} \rangle + \dots \end{aligned}$$

Rewriting K-means as a matrix estimation problem

$$\begin{aligned} \mathsf{Crit}(\mathbf{X}, G) &= \frac{1}{2} \sum_{k=1}^{K} \frac{1}{|G_k|} \sum_{a, b \in G_k} \|X_a - X_b\|^2 \\ &= -\sum_k \sum_{a, b \in G_k} \langle X_a, X_b \rangle \frac{1}{|G_k|} + \sum_{a=1}^n \|X_a\|_2^2 \\ &= -\langle \mathbf{X} \mathbf{X}^T, \mathbf{B} \rangle + \dots \end{aligned}$$

K-means : linear minimization problem over the space of partnership matrices

$$\widehat{G} = \arg\min_{\mathbf{B}} \langle -\mathbf{X}\mathbf{X}^T, \mathbf{B} \rangle$$

One caveat of K-means criterion



Bias of K-means :

Tends to split groups with large variance and favors small groups

Outline



2 Omnibus Clustering Algorithms

- K-means criterion
- Lloyd's Algorithm
- Peng and Wei's SDP
- Spectral Clustering

- Loss functions
- Minimax lower bound
- Recovery bounds for Peng and Wei' SDP

K-means as a least-square problem :

$$\hat{G} \in \arg\min_{G} \sum_{k=1}^{K} \min_{\theta_k \in \mathbb{R}^p} \sum_{a \in G_k} \|X_a - \theta_k\|_2^2$$

K-means as a least-square problem :

$$\hat{G} \in \arg\min_{G} \sum_{k=1}^{K} \min_{\theta_k \in \mathbb{R}^p} \sum_{a \in G_k} \|X_a - \theta_k\|_2^2$$

Alternate Minimization between estimation of the centroids and of the partition

- 1 Compute the centroids
- 2 Update the partition



K-means as a least-square problem :

$$\hat{G} \in \arg\min_{G} \sum_{k=1}^{K} \min_{\theta_k \in \mathbb{R}^p} \sum_{a \in G_k} \|X_a - \theta_k\|_2^2$$

Alternate Minimization between estimation of the centroids and of the partition

- 1 Compute the centroids
- 2 Update the partition

K-means as a least-square problem :

$$\hat{G} \in \arg\min_{G} \sum_{k=1}^{K} \min_{\theta_k \in \mathbb{R}^p} \sum_{a \in G_k} \|X_a - \theta_k\|_2^2$$

(

Alternate Minimization between estimation of the centroids and of the partition

- 1 Compute the centroids
- 2 Update the partition

K-means as a least-square problem :

$$\hat{G} \in \arg\min_{G} \sum_{k=1}^{K} \min_{\theta_k \in \mathbb{R}^p} \sum_{a \in G_k} \|X_a - \theta_k\|_2^2$$

Alternate Minimization between estimation of the centroids and of the partition

- 1 Compute the centroids
- 2 Update the partition



K-means as a least-square problem :

$$\hat{G} \in \arg\min_{G} \sum_{k=1}^{K} \min_{\theta_k \in \mathbb{R}^p} \sum_{a \in G_k} \|X_a - \theta_k\|_2^2$$

Alternate Minimization between estimation of the centroids and of the partition

- 1 Compute the centroids
- 2 Update the partition



K-means as a least-square problem :

$$\hat{G} \in \arg\min_{G} \sum_{k=1}^{K} \min_{\theta_k \in \mathbb{R}^p} \sum_{a \in G_k} \|X_a - \theta_k\|_2^2$$

Alternate Minimization between estimation of the centroids and of the partition

- 1 Compute the centroids
- 2 Update the partition



K-means as a least-square problem :

$$\hat{G} \in \arg\min_{G} \sum_{k=1}^{K} \min_{\theta_k \in \mathbb{R}^p} \sum_{a \in G_k} \|X_a - \theta_k\|_2^2$$

Alternate Minimization between estimation of the centroids and of the partition



K-means as a least-square problem :

$$\hat{G} \in \arg\min_{G} \sum_{k=1}^{K} \min_{\theta_k \in \mathbb{R}^p} \sum_{a \in G_k} \|X_a - \theta_k\|_2^2$$

Alternate Minimization between estimation of the centroids and of the partition



K-means++ achieves a $\log(K)$ approximation of K-means criterion (in worst case).

Outline

2 Omnibus Clustering Algorithms

- K-means criterion
- Lloyd's Algorithm
- Peng and Wei's SDP
- Spectral Clustering

- Loss functions
- Minimax lower bound
- Recovery bounds for Peng and Wei' SDP
Matrix Characterization of the Partnership matrix.

The $n \times n$ Partnership matrix

$$\mathbf{B}_{ab} = \left\{ \begin{array}{cc} \frac{1}{|G_k|} & \text{ if } a \text{ and } b \text{ belong to the same } G_k \\ 0 & \text{ else} \end{array} \right.$$

Lemma (Peng & Wei(07))

```
The K-means minimizer \widehat{G} satisfies
```

$$\begin{split} \widehat{\mathbf{B}} &\in \arg\min_{\mathbf{B}\in\mathcal{D}} \langle -\mathbf{X}\mathbf{X}^T, \mathbf{B} \rangle \ , \\ \mathcal{D} &:= \left\{ \begin{aligned} \mathbf{B} &\in \mathbb{R}^{n \times n} : & \mathbf{B} \succcurlyeq 0 \\ \mathbf{B} &\in \mathbb{R}^{n \times n} : & \mathbf{B}_{ab} \geqslant 0, \ \forall a, b \\ &\bullet \operatorname{Tr}(\mathbf{B}) = K \\ &\bullet \mathbf{B}^2 = \mathbf{B} \end{aligned} \right\}$$

<u>**Proof</u>** : **B** is a bistochastic matrix with K eigenvalues equal to 1.</u>

Matrix Characterization of the Partnership matrix.

The $n \times n$ Partnership matrix

$$\mathbf{B}_{ab} = \left\{ \begin{array}{cc} \frac{1}{|G_k|} & \text{ if } a \text{ and } b \text{ belong to the same } G_k \\ 0 & \text{ else} \end{array} \right.$$

Lemma (Peng & Wei(07))

The K-means minimizer \widehat{G} satisfies

$$\widehat{\mathbf{B}} \in \arg\min_{\mathbf{B}\in\mathcal{D}} \langle -\mathbf{X}\mathbf{X}^T, \mathbf{B} \rangle ,$$

$$\mathcal{D} := \left\{ \begin{aligned} \mathbf{0} & \mathbf{B} \succcurlyeq \mathbf{0} \\ \mathbf{0} & \sum_{a} \mathbf{B}_{ab} = 1, \forall b \\ \mathbf{B} \in \mathbb{R}^{n \times n} : & \mathbf{0} \mathbf{B}_{ab} \geqslant \mathbf{0}, \forall a, b \\ \mathbf{0} & \operatorname{Tr}(\mathbf{B}) = K \\ \mathbf{0} & \mathbf{B}^2 = \mathbf{B} \end{aligned} \right\}$$

<u>Proof</u>: **B** is a bistochastic matrix with K eigenvalues equal to 1. Perron-Frobenius theorem \rightsquigarrow support of **B**= adjacency matrix of a graph with K cc each block has rank 1 and is bistochastic.

Relaxed *K*-means

Idea : drop the $\mathbf{B}^2 = \mathbf{B}$ condition.

1 Estimate \mathbf{B}^* using the semi-definite program (SDP)

$$\widehat{\mathbf{B}} = \operatorname*{arg\,min}_{\mathbf{B} \in \mathcal{C}} \langle -\mathbf{X}\mathbf{X}^T, \mathbf{B} \rangle$$

where

$$\mathcal{C} := \begin{cases} \mathbf{B} \in \mathbb{R}^{n \times n} : & \mathbf{B} \succcurlyeq 0 \\ \mathbf{\Phi} \sum_{a} \mathbf{B}_{ab} = 1, \ \forall b \\ \mathbf{B}_{ab} \geqslant 0, \ \forall a, b \\ \mathbf{\Phi} \text{ Tr}(\mathbf{B}) = K \end{cases}$$

2 (Compute \widehat{G} by applying any clustering algorithm on $\widehat{\mathbf{B}}$)

- Convex optimization but many constraints: https://www.ams.jhu.edu/villar/research/ (n ≈ a few hundreds) Iguchi et al.('15), Mixon et al.('17)
- No information of the group sizes is needed.

Relaxed *K*-means

Idea : drop the $\mathbf{B}^2 = \mathbf{B}$ condition.

1 Estimate \mathbf{B}^* using the semi-definite program (SDP)

$$\widehat{\mathbf{B}} = \operatorname*{arg\,min}_{\mathbf{B} \in \mathcal{C}} \langle -\mathbf{X}\mathbf{X}^T, \mathbf{B} \rangle$$

where

$$\mathcal{C} := \begin{cases} \mathbf{B} \in \mathbb{R}^{n \times n} : & \mathbf{B} \succcurlyeq 0 \\ \mathbf{\Phi} \sum_{a} \mathbf{B}_{ab} = 1, \ \forall b \\ \mathbf{B}_{ab} \geqslant 0, \ \forall a, b \\ \mathbf{\Phi} \text{ Tr}(\mathbf{B}) = K \end{cases}$$

2 (Compute \widehat{G} by applying any clustering algorithm on $\widehat{\mathbf{B}}$)

Remark :

- Convex optimization but many constraints : https://www.ams.jhu.edu/villar/research/ (n ≈ a few hundreds) lguchi et al.('15), Mixon et al.('17)
- No information of the group sizes is needed.

Differs from the SDP relaxation of Max-Cut Problem Goeman and Williamson('95) (see e.g. Hajek et al.('16)) where the size of the communities has to be known...

Outline



2 Omnibus Clustering Algorithms

- K-means criterion
- Lloyd's Algorithm
- Peng and Wei's SDP
- Spectral Clustering

- Loss functions
- Minimax lower bound
- Recovery bounds for Peng and Wei' SDP

Heuristic

$$\mathbf{E}[\mathbf{X}\mathbf{X}^T] = \mathbf{A}\boldsymbol{\Theta}\boldsymbol{\Theta}^T\mathbf{A}^T \left(+ \mathbb{E}[\mathbf{E}\mathbf{E}^T] \right)$$

Lemma (e.g. Lei and Rinaldo('15))

Assume that $\Theta \Theta^T$ has full rank. Let $\mathbf{A} \Theta \Theta^T \mathbf{A}^T = \sum_{k=1}^{K} d_k u_k u_k^T$ and set $\mathbf{U} = [u_1, \dots, u_k] \in \mathbb{R}^{n \times K}$. Then, there exist $Z_1, \dots, Z_K \in \mathbb{R}^K$, such that

$$U_{i:} = Z_k \text{ for all } i \in G_k, \text{ and } ||Z_k - Z_\ell||^2 = \frac{1}{|G_k|} + \frac{1}{|G_\ell|}.$$

$$\mathbf{U} = \begin{pmatrix} Z_{k(1)} \\ Z_{k(2)} \\ \cdots \\ Z_{k(n)} \end{pmatrix}$$

Spectral Clustering as relaxed K-means

Spectral Clustering

- 1 Compute the matrix $\widehat{\mathbf{U}}$ made of the K-leading eigenvectors of $\mathbf{X}^T \mathbf{X}$
- 2 Estimate \widehat{G} by distance clustering on the rows of $\widehat{\mathbf{U}}$.

(e.g. Apply an approximate K-means algorithm to the rows of the matrix $\widehat{\mathbf{U}}$)

Spectral Clustering as relaxed K-means

Spectral Clustering

- 1 Compute the matrix $\widehat{\mathbf{U}}$ made of the K-leading eigenvectors of $\mathbf{X}^T \mathbf{X}$
- 2 Estimate \widehat{G} by distance clustering on the rows of $\widehat{\mathbf{U}}$.
- (e.g. Apply an approximate K-means algorithm to the rows of the matrix $\widehat{\mathbf{U}})$



 \implies it amounts to dropping the constraints ${f B}1=1,\;{f B}_{ab}\geqslant 0$ in the former relaxation

Spectral Clustering as relaxed K-means

Spectral Clustering

- 1 Compute the matrix $\widehat{\mathbf{U}}$ made of the K-leading eigenvectors of $\mathbf{X}^T \mathbf{X}$
- 2 Estimate \widehat{G} by distance clustering on the rows of $\widehat{\mathbf{U}}$.
- (e.g. Apply an approximate K-means algorithm to the rows of the matrix $\widehat{\mathbf{U}}$)

 $\begin{array}{l} \Longrightarrow \text{ it amounts to dropping the constraints } \mathbf{B}1=1, \ \mathbf{B}_{ab} \ge 0 \text{ in the former relaxation} \\ \underline{\text{Proof}}:1) \ \mathbf{\overline{B}}=\widehat{\mathbf{U}}\widehat{\mathbf{U}}^T \\ 2) \ (\widehat{\mathbf{U}}\widehat{\mathbf{U}}^T)_{a:} \text{ is some orthogonal transformation of } \widehat{\mathbf{U}}_{a:}. \end{array}$

Outline

1 Hidden Partition problems

2

- Omnibus Clustering Algorithms
- K-means criterion
- Lloyd's Algorithm
- Peng and Wei's SDP
- Spectral Clustering

3 Recovery bounds for (sub)Gaussian Mixtures

Loss functions

- Minimax lower bound
- Recovery bounds for Peng and Wei' SDP

Partial recovery bounds

Proportion of misclustered points

$$err(\widehat{G}, G^*) = \min_{\pi \in \mathcal{S}_K} \frac{1}{2n} \sum_{k=1}^K \left| G_k^* \triangle \widehat{G}_{\pi(k)} \right|$$

Our goal

Prove that with high-probability, when s^2 is large

prop. misclustered =
$$err(\widehat{G}, G^*) \le e^{-cs^2}$$

where s^2 is an appropriate SNR.

Other related goals :

- **partial recovery** : Find the minimal s^2 such that $err(\hat{G}, G^*)$ is smaller than random guess whp.
- Almost full recovery : Find the minimal s^2 such that $err(\widehat{G}, G^*) \to 0$
- Perfect recovery : Find the minimal s^2 such that $err(\widehat{G}, G^*) = 0$ whp.

Mixture of Gaussian variables

 $X_1,\ldots,X_n\in\mathbb{R}^p$ are independent with

 $X_a \sim \mathcal{N}(\theta_k, \mathbf{\Sigma}) \ a \in G_k^*$

many results discussed later readily extend to subGaussian random variables and to different covariances

Set
$$\Delta^2 = \min_{j \neq k} \|\theta_k - \theta_j\|^2$$
, $\sigma^2 = \|\mathbf{\Sigma}\|_{op}$ and $R_{\mathbf{\Sigma}} = \frac{\|\mathbf{\Sigma}\|_F^2}{\|\mathbf{\Sigma}\|_{op}^2}$,

Specific case of Isovolumetric spherical Gaussians

 $\Sigma = \sigma^2 \mathbf{I}_p$

Outline

1 Hidden Partition problems

2

Omnibus Clustering Algorithms

- K-means criterion
- Lloyd's Algorithm
- Peng and Wei's SDP
- Spectral Clustering

3 Recovery bounds for (sub)Gaussian Mixtures

- Loss functions
- Minimax lower bound
- Recovery bounds for Peng and Wei' SDP

What should we expect? What is the correct SNR?

Toy example : K = 2, $|G_1^*| = |G_2^*| = n/2$, $\Sigma = \sigma^2 \mathbf{I}_p$, $\theta_2 = -\theta_1$.

Simpler Problem 1 (known parameters) : θ_1 is known.



What should we expect? What is the correct SNR?

Toy example : K = 2, $|G_1^*| = |G_2^*| = n/2$, $\Sigma = \sigma^2 \mathbf{I}_p$, $\theta_2 = -\theta_1$.

Simpler Problem 1 (known parameters) : θ_1 is known.



Bayes Classifier achieves :

$$\mathbb{E}[err(\widehat{G}, G^*)] = \mathbb{P}[\mathcal{N}(0, \sigma^2) > \|\theta_1\|] \le \exp\left[-\frac{\Delta^2}{8\sigma^2}\right]$$

Benchmark 2 : Supervised case

Parameter : θ_1 is sampled uniformly on the sphere of radius $\Delta/2$. **Supervised observations** : $\mathcal{L} = (X_a, Z_a)_{a=1,...,n}$ where $Z_a \in \{1, 2\}$ is the label of X_a and is sampled uniformly in $\{1, 2\}$

Objective : classify a new observation X.

Benchmark 2 : Supervised case

Parameter : θ_1 is sampled uniformly on the sphere of radius $\Delta/2$. **Supervised observations** : $\mathcal{L} = (X_a, Z_a)_{a=1,...,n}$ where $Z_a \in \{1, 2\}$ is the label of X_a and is sampled uniformly in $\{1, 2\}$

Objective classify a new observation X.

 $\begin{array}{l} \text{Optimal Classifier is achieved by LDA} \\ \widehat{h}(x) = \left\{ \begin{array}{ll} 2 & \text{if } \mathbb{P}[Z=2|X=x,\mathcal{L}] > \mathbb{P}[Z=1|X=x,\mathcal{L}] \rangle \\ 1 & \text{if } \mathbb{P}[Z=2|X=x,\mathcal{L}] \leq \mathbb{P}[Z=1|X=x,\mathcal{L}] \rangle \end{array} \right. \end{array}$

$$\widehat{h}(x) = \frac{3}{2} + \frac{1}{2}\operatorname{sign}\left(\left\langle\frac{1}{n}\sum_{a=1}^{n}(2Z_a - 3)X_a, x\right\rangle\right)$$

Benchmark 2 : Supervised case

Parameter : θ_1 is sampled uniformly on the sphere of radius $\Delta/2$. **Supervised observations** : $\mathcal{L} = (X_a, Z_a)_{a=1,...,n}$ where $Z_a \in \{1, 2\}$ is the label of X_a and is sampled uniformly in $\{1, 2\}$

Objective classify a new observation X.

 $\begin{array}{l} \text{Optimal Classifier is achieved by LDA} \\ \widehat{h}(x) = \left\{ \begin{array}{ll} 2 & \text{if } \mathbb{P}[Z=2|X=x,\mathcal{L}] > \mathbb{P}[Z=1|X=x,\mathcal{L}] \rangle \\ 1 & \text{if } \mathbb{P}[Z=2|X=x,\mathcal{L}] \leq \mathbb{P}[Z=1|X=x,\mathcal{L}] \rangle \end{array} \right. \end{array}$

$$\widehat{h}(x) = \frac{3}{2} + \frac{1}{2}\operatorname{sign}\left(\left\langle\frac{1}{n}\sum_{a=1}^{n}(2Z_a - 3)X_a, x\right\rangle\right)$$

$$\begin{split} \mathbb{P}[\hat{h}(X) \neq Z] &= & \mathbb{E}_{\theta} \, \mathbb{P}\left[\left\langle \theta_{1} + \frac{\sigma}{\sqrt{n}} \epsilon, \theta_{1} + \sigma \epsilon' \right\rangle < 0 | \theta \right] \\ &= & \mathbb{P}\left[\frac{\Delta^{2}}{4\sigma^{2}} < \frac{\Delta}{2\sigma} \left(\frac{\epsilon_{1}}{\sqrt{n}} + \epsilon'_{1} \right) - \frac{1}{\sqrt{n}} \left\langle \epsilon, \epsilon' \right\rangle \right] \\ &\sim_{\log} & \begin{cases} & \exp\left(- \frac{\Delta^{2}}{8\sigma^{2}} \right) & \text{if } \frac{\Delta^{2}}{\sigma^{2}} \gg \left[1 \lor \frac{p}{n} \right] \\ & \exp\left(- \frac{n\Delta^{4}}{32p\sigma^{2}} \right) & \text{if } \left[1 \lor \sqrt{\frac{p}{n}} \right] \ll \frac{\Delta^{2}}{\sigma^{2}} \ll \frac{p}{n} \end{split}$$

Toy Model : K = 2, $|G_1^*| = |G_2^*| = n/2$, $\Sigma = \sigma^2 \mathbf{I}_p$



For more communities, define the SNR

$$s^2 = \frac{\Delta^2}{\sigma^2} \wedge \frac{m\Delta^4}{R_{\Sigma}\sigma^4},$$

where

$$\blacksquare R_{\Sigma} = \frac{\|\Sigma\|_{F}^{2}}{\|\Sigma\|_{op}^{2}} \text{ is the effective rank of } \Sigma$$

Outline

1 Hidden Partition problems

- 2
 - Omnibus Clustering Algorithms
 - K-means criterion
 - Lloyd's Algorithm
 - Peng and Wei's SDP
 - Spectral Clustering

3 Recovery bounds for (sub)Gaussian Mixtures

- Loss functions
- Minimax lower bound
- Recovery bounds for Peng and Wei' SDP

relaxed K-means

Solve the SDP

$$\widehat{B} \in \underset{\mathbf{B} \in \mathcal{C}}{\operatorname{argmin}} \langle -\mathbf{X}^T \mathbf{X}, \mathbf{B} \rangle \ ,$$

with

$$\mathcal{C} := \left\{ \begin{aligned} \mathbf{B} \in \mathbb{R}^{n \times n} : & \stackrel{\bullet}{\bullet} \stackrel{\mathbf{B} \succcurlyeq 0}{\underset{ab}{\bullet} \stackrel{\bullet}{=} 1, \forall b} \\ \mathbf{B} \in \mathbb{R}^{n \times n} : & \stackrel{\bullet}{\bullet} \stackrel{\mathbf{B}_{ab}}{\underset{ab}{\bullet} 0, \forall a, b} \\ & \stackrel{\bullet}{\bullet} \operatorname{Tr}(\mathbf{B}) = K \end{aligned} \right\}$$

Step 2 : Apply approximate K-medoid method Charikar et al. ('02).

$$|\widehat{\mathbf{A}}\widehat{\mathbf{M}} - \widehat{\mathbf{B}}|_1 \le \rho \min_{A, \operatorname{Rows}(M) \subset \operatorname{Rows}(\widehat{\mathbf{B}})} |\mathbf{A}\mathbf{M} - \widehat{\mathbf{B}}|_1$$

$$s^2 = \frac{\Delta^2}{\sigma^2} \wedge \frac{m\Delta^4}{R_{\Sigma}\sigma^4},$$

Theorem (Giraud and V. ('18))

If
$$s^2 \gtrsim n/m$$
, then $\mathbb{P}\left[err(\widehat{G}, G^*) > e^{-cs^2}\right] \lesssim \frac{1}{n^2}$.

$$s^2 = \frac{\Delta^2}{\sigma^2} \wedge \frac{m\Delta^4}{R_{\Sigma}\sigma^4},$$

Theorem (Giraud and V. ('18))

If $s^2\gtrsim n/m,$ then $\mathbb{P}\left[err(\widehat{G},G^*)>e^{-cs^2}\right]\lesssim \frac{1}{n^2}.$

$$s^2 \gtrsim n/m = K$$
 is equivalent to $\Delta^2 \gtrsim \sigma^2 \frac{n}{m} \left(1 \lor \sqrt{\frac{R_{\Sigma}}{n}} \right) = \sigma^2 K \left(1 \lor \sqrt{\frac{R_{\Sigma}}{n}} \right).$

$$s^2 = \frac{\Delta^2}{\sigma^2} \wedge \frac{m\Delta^4}{R_{\Sigma}\sigma^4},$$

Theorem (Giraud and V. ('18))

If
$$s^2\gtrsim n/m$$
, then $\mathbb{P}\left[err(\widehat{G},G^*)>e^{-cs^2}\right]\lesssim rac{1}{n^2}.$

$$s^2 \gtrsim n/m = K$$
 is equivalent to $\Delta^2 \gtrsim \sigma^2 \frac{n}{m} \left(1 \lor \sqrt{\frac{R_{\Sigma}}{n}} \right) = \sigma^2 K \left(1 \lor \sqrt{\frac{R_{\Sigma}}{n}} \right).$

- 1 Optimal convergence rate provided the SNR is large enough.
- 2 No restriction on the dimension

$$s^2 = \frac{\Delta^2}{\sigma^2} \wedge \frac{m\Delta^4}{R_{\Sigma}\sigma^4},$$

Theorem (Giraud and V. ('18))

If
$$s^2\gtrsim n/m,$$
 then $\mathbb{P}\left[err(\widehat{G},G^*)>e^{-cs^2}\right]\lesssim \frac{1}{n^2}.$

$$s^2 \gtrsim n/m = K$$
 is equivalent to $\Delta^2 \gtrsim \sigma^2 \frac{n}{m} \left(1 \lor \sqrt{\frac{R_{\Sigma}}{n}} \right) = \sigma^2 K \left(1 \lor \sqrt{\frac{R_{\Sigma}}{n}} \right).$

- 1 Optimal convergence rate provided the SNR is large enough.
- 2 No restriction on the dimension
- 3 Ideas come from Fei and Chen('17) for SBMs. See also Mixon et al.('16), Fei and Chen('18).

$$s^2 = \frac{\Delta^2}{\sigma^2} \wedge \frac{m\Delta^4}{R_{\Sigma}\sigma^4},$$

Theorem (Giraud and V. ('18))

If
$$s^2\gtrsim n/m,$$
 then $\mathbb{P}\left[err(\widehat{G},G^*)>e^{-cs^2}\right]\lesssim \frac{1}{n^2}.$

$$s^2 \gtrsim n/m = K$$
 is equivalent to $\Delta^2 \gtrsim \sigma^2 \frac{n}{m} \left(1 \lor \sqrt{\frac{R_{\Sigma}}{n}} \right) = \sigma^2 K \left(1 \lor \sqrt{\frac{R_{\Sigma}}{n}} \right).$

- 1 Optimal convergence rate provided the SNR is large enough.
- 2 No restriction on the dimension
- 3 Ideas come from Fei and Chen('17) for SBMs. See also Mixon et al.('16), Fei and Chen('18).
- 4 It does not recover the tight constant inside the exponential. This is possible in specific situations; see Fei and Chen('19) for SBM.

$$s^2 = \frac{\Delta^2}{\sigma^2} \wedge \frac{m\Delta^4}{R_{\Sigma}\sigma^4},$$

Theorem (Giraud and V. ('18))

$$\text{If } s^2 \gtrsim n/m \text{, then } \mathbb{P}\left[err(\widehat{G}, G^*) > e^{-cs^2}\right] \lesssim \frac{1}{n^2}.$$

$$s^2 \gtrsim n/m = K$$
 is equivalent to $\Delta^2 \gtrsim \sigma^2 \frac{n}{m} \left(1 \lor \sqrt{\frac{R_{\Sigma}}{n}} \right) = \sigma^2 K \left(1 \lor \sqrt{\frac{R_{\Sigma}}{n}} \right).$

- 1 Optimal convergence rate provided the SNR is large enough.
- 2 No restriction on the dimension
- 3 Ideas come from Fei and Chen('17) for SBMs. See also Mixon et al.('16), Fei and Chen('18).
- 4 It does not recover the tight constant inside the exponential. This is possible in specific situations; see Fei and Chen('19) for SBM.
- 5 perfect recovery for $s^2 \gtrsim \log(n) \lor (n/m) = \log(n) \lor K$

Large number K of clusters

Simplification : balanced partition + Isovolumetric spherical Gaussians.

- SNR condition $s^2 \ge K$ is needed for Peng and Wei's SDP or for spectral clustering.
- From Regev and Vijayaraghavan('17), the condition $s^2 \gtrsim \log(K)$ is needed for partial recovery to be information theoretically possible.

Large number K of clusters

Simplification : balanced partition + Isovolumetric spherical Gaussians.

- SNR condition $s^2 \ge K$ is needed for Peng and Wei's SDP or for spectral clustering.
- From Regev and Vijayaraghavan('17), the condition $s^2 \gtrsim \log(K)$ is needed for partial recovery to be information theoretically possible.
- In low-dimensional settings (e.g. $n\widetilde{\gg}p^{3}K^{2}),~{\rm Vempala}$ and ${\rm Wang}('04)$ achieve exact recovery if

 $s^2 \gtrsim \sqrt{K \log(n)} + \log(n)$.

Large number K of clusters

Simplification : balanced partition + Isovolumetric spherical Gaussians.

- SNR condition $s^2 \ge K$ is needed for Peng and Wei's SDP or for spectral clustering.
- From Regev and Vijayaraghavan('17), the condition $s^2 \gtrsim \log(K)$ is needed for partial recovery to be information theoretically possible.
- In low-dimensional settings (e.g. $n \gg p^3 K^2$), Vempala and Wang('04) achieve exact recovery if

 $s^2 \gtrsim \sqrt{K \log(n)} + \log(n)$.

In high-dimension (p ≥ n), it is conjectured that no polynomial-time estimator can beat the condition s² ≳ K − Banks et al.('18) With a proper correction (see e.g. Bunea et al.('16)), general covariances $\Sigma_1, \ldots, \Sigma_K$ can be handled by K-means and its relaxations.

- With a proper correction (see e.g. Bunea et al. ('16)), general covariances $\Sigma_1, \ldots, \Sigma_K$ can be handled by *K*-means and its relaxations.
- ... However, those methods do not build upon general covariances. Optimal convergence rate with respect to the **Mahalanobis** distance $(\mu_k - \mu_l)^T \Sigma^{-1} (\mu_k - \mu_l)$
 - \rightsquigarrow For unknown $\Sigma,$ this is an active and challenging research direction Dabis et al.('21)

First Message

Various convex relaxations of K-means seem to work well...

First Message

Various convex relaxations of K-means seem to work well...

Second Message

... for a variety of models. (SBM, GMM, ...)

First Message

Various convex relaxations of K-means seem to work well...

Second Message

... for a variety of models (SBM, GMM, ...)

Third Message

Large K asymptotic is still not completely understood.
First Message

Various convex relaxations of K-means seem to work well...

Second Message

... for a variety of models (SBM, GMM,)

Third Message

Large K asymptotic is still not completely understood.

Danke für Ihre Aufmerksamkeit!

References I



E. Abbe.

Community detection and stochastic block models : recent developments. ArXiv e-prints, March 2017.



Pranjal Awasthi, Moses Charikar, Ravishankar Krishnaswamy, and Ali Kemal Sinop

The Hardness of Approximation of Euclidean k-means.

arXiv preprint arXiv :1502.03316, 2015.



Emmanuel Abbe, Jianqing Fan, Kaizheng Wang, and Yiqiao Zhong. Entrywise eigenvector analysis of random matrices with low expected rank. arXiv preprint arXiv :1709.09565, 2017.



Matthew Brennan and Guy Bresler.

Average-case lower bounds for learning sparse mixtures, robust estimation and semirandom adversaries.

arXiv preprint arXiv :1908.06130, 2019.



Optimal Average-Case Reductions to Sparse PCA : From Weak Assumptions to Strong Hardness.

arXiv preprint arXiv :1902.07380, 2019.

References II



Matthew Brennan, Guy Bresler, and Wasim Huleihel.

Reducibility and computational lower bounds for problems with planted sparse structure.

arXiv preprint arXiv :1806.07508, 2018.



Matthew Brennan, Guy Bresler, and Wasim Huleihel. Universality of computational lower bounds for submatrix detection. arXiv preprint arXiv :1902.06916, 2019.





Jess Banks, Cristopher Moore, Roman Vershynin, Nicolas Verzelen, and Jiaming Xu.

Information-theoretic bounds and phase transitions in clustering, sparse PCA, and submatrix localization.

IEEE Transactions on Information Theory, 64(7) :4872-4894, 2018.

Moses Charikar, Sudipto Guha, Éva Tardos, and David B. Shmoys. A Constant-Factor Approximation Algorithm for the k-Median Problem. Journal of Computer and System Sciences, 65(1) :129–149, 2002.



📄 Peter Chin, Anup Rao, and Van Vu.

Stochastic Block Model and Community Detection in Sparse Graphs : A spectral algorithm with optimal rate of recovery.

In Peter Grünwald, Elad Hazan, and Satyen Kale, editors, Proceedings of The 28th Conference on Learning Theory, volume 40 of Proceedings of Machine Learning Research, pages 391-423, Paris, France, 03-06 Jul 2015 PMLR



Yudong Chen and Jiaming Xu.

Statistical-computational tradeoffs in planted problems and submatrix localization with a growing number of clusters and submatrices. Journal of Machine Learning Research, 17(27) :1-57, 2016.



Ilias Diakonikolas, Daniel M. Kane, and Alistair Stewart. List-Decodable Robust Mean Estimation and Learning Mixtures of Spherical Gaussians

CoRR, abs/1711.07211, 2017.



Yingjie Fei and Yudong Chen.

Hidden integrality of sdp relaxation for sub-gaussian mixture models. arXiv preprint arXiv :1803.06510, 2018.



Paul W Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. Stochastic blockmodels : First steps. Social networks, 5(2) :109-137, 1983.

References V



References VI

Y. Lu and H. H. Zhou. Statistical and Computational Guarantees of Lloyd's Algorithm and its Variants. ArXiv e-prints, December 2016.



Matthias Löffler, Anderson Y Zhang, and Harrison H Zhou. Optimality of spectral clustering for gaussian mixture model. arXiv preprint arXiv :1911.00538, 2019.

Meena Mahajan, Prajakta Nimbhorkar, and Kasturi Varadarajan. The planar k-means problem is NP-hard. In International Workshop on Algorithms and Computation, pages 274–285. Springer, 2009.

Ankur Moitra, William Perry, and Alexander S Wein. How robust are reconstruction thresholds for community detection ? In

Proceedings of the forty-eighth annual ACM symposium on Theory of Computing, pages 828-841. ACM, 2016.

Dustin G Mixon, Soledad Villar, and Rachel Ward. Clustering subgaussian mixtures by semidefinite programming. Information and Inference : A Journal of the IMA, 6(4) :389–415, 2017.

References VII



Mohamed Ndaoud.

Sharp optimal recovery in the Two Gaussian Mixture Model. arXiv preprint arXiv :1812.08078, 2018.



Jiming Peng and Yu Wei.

Approximating K-means-type Clustering via Semidefinite Programming. SIAM J. on Optimization, 18(1) :186–205, February 2007.

M. Royer. Adaptive Clustering through Semidefinite Programming. Advances in Neural Information Processing Systems (NIPS), 2017.



Oded Regev and Aravindan Vijayaraghavan.

On learning mixtures of well-separated gaussians.

ln

2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS), pages 85–96. IEEE, 2017.



Nathan Srebro, Gregory Shakhnarovich, and Sam Roweis.

An investigation of computational and informational limits in gaussian mixture clustering.

In Proceedings of the 23rd international conference on Machine learning, pages 865–872. ACM, 2006.



Yihong Wu and Harrison H Zhou.

Randomly initialized EM algorithm for two-component Gaussian mixture achieves near optimality in $O(n\hat{1}/2)$ iterations. arXiv preprint arXiv :1908.10935, 2019.



Anru Zhang, T Tony Cai, and Yihong Wu. Heteroskedastic PCA : Algorithm, optimality, and applications. arXiv preprint arXiv :1810.08316, 2018.