

Classification non-supervisée de données de grande dimension et de graphes à l'aide de modèles à variables latentes discrètes

Soutenance de Thèse

Nicolas Jouvin

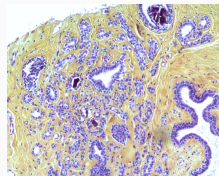
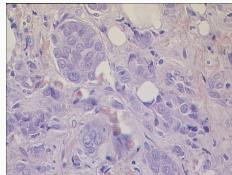
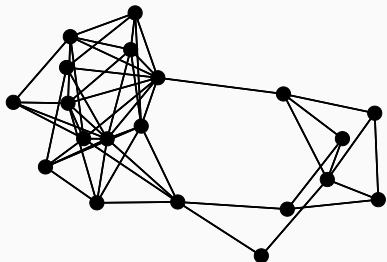
Directeurs : P. Latouche, C. Bouveyron & A. Livartowski

Vendredi 11 décembre 2020



Clustering in a nutshell

clustering
dimension
algorithm
model
data
latent
analysis
matrix
clusters
mixture



Clustering is the task of **grouping objects** together into classes or *clusters*, in an **unsupervised** fashion based on some **criterion**.

Example 1: document clustering

Grouping similar texts together based on their topics.

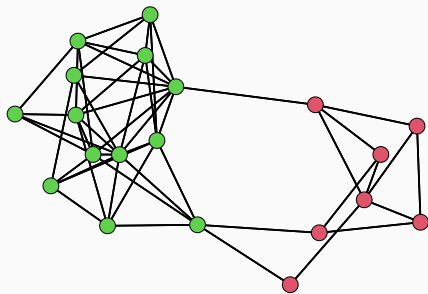
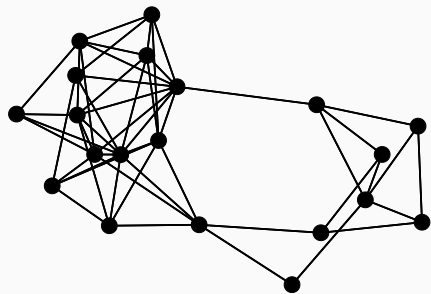
MICROBIOPSIE SOUS ECHOGRAPHIE DU SEIN DROIT	MICROBIOPSIE SOUS ECHOGRAPHIE DU SEIN DROIT	MACROBIOPSIE DU SEIN GAUCHE
MACROSCOPIE Cinq fragments de 5 à 15 mm	MACROSCOPIE Cinq fragments de 5 à 15 mm	MACROSCOPIE 3 fragments de 7 à 15 mm
MICROSCOPIE Les prélèvements examinés correspondent à des fragments de tissu mammaire remanié par une prolifération tumorale dont les caractères morphologiques sont ceux d'un adénocarcinome canalaire infiltrant. Cette lésion est peu différenciée, d'architecture essentiellement trabéculaire. Les cellules néoplasiques comportent des atypies nucléaires marquées. L'index mitotique est élevé (22 mitoses sur 10 champs au grossissement 400). Deux fragments de 8 et 15 mm. Adénocarcinome mammaire de type canalaire infiltrant, peu différencié. Grade histo-pronostique (EE) : III Index mitotique élevé.	MICROSCOPIE L'examen histologique met en évidence des lésions tumorales dont les caractères morphologiques sont ceux d'un carcinome canalaire infiltrant moyennement différencié. La lésion est d'architecture trabéculaire et glandulaire. Les cellules sont caractérisées par des atypies cytonucléaires modérées. L'activité mitotique est faible : deux mitoses ont été dénombrées sur dix champs au grossissement 400. Ces lésions sont associées à un stroma dense fibreux. Elles infiltrant le tissu adipeux. Deux séries de prélèvements ont été confondues : A - 1er tour : onze cylindres biopsiques mesurant 10 à 30 mm de long. B - 2ème tour : onze cylindres biopsiques mesurant 5 à 30 m de long. Adénocarcinome mammaire de type canalaire infiltrant. Grade histopronostique (EE) I. Index mitotique faible.	MICROSCOPIE Tous les prélèvements ont un aspect histologique similaire. Ils correspondent à des fragments de tissu mammaire remanié par des lésions de mastose fibreuse commune. Présence d'un discret infiltrat inflammatoire. On retrouve également quelques microcalcifications. L'un des prélèvements cryo-prélevés sera analysé histologiquement et un compte rendu complémentaire adressé ultérieurement. Trois fragments de 7 à 15 mm. Lésions de mastose fibreuse. Le prélèvement paraît peu significatif. Une analyse complémentaire sur le prélèvement cryo-prélevé sera réalisée.

Doc 1	"Lésions cancéreuses (...) carcinome canalaire"
Doc 2	"Lésions cancéreuses (...) carcinome lobulaire"
...	...
Doc n	"Lésions bénignes (...) métaplasie"

Clustering is the task of **grouping objects** together into classes or *clusters*, in an **unsupervised** fashion based on some **criterion**.

Example 2: Network clustering

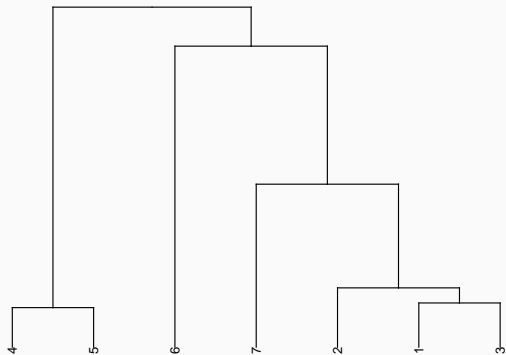
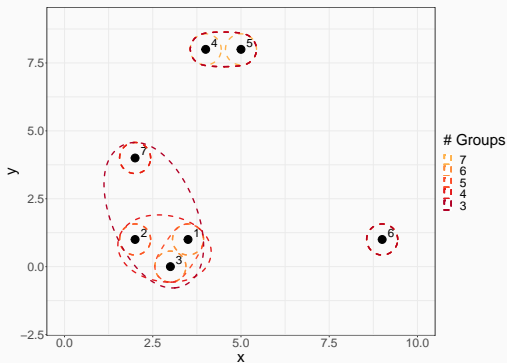
Group nodes of a network based on their connections with respect to others



Clustering is the task of **grouping objects** together into classes or *clusters*, in an **unsupervised** fashion based on some **criterion**.

Example 3: Hierarchical clustering

Build a hierarchy of *nested* clusters



Probabilistic approach for three types of data

- Count data, e.g. text documents
- Continuous data, e.g. images
- Graph data

Handle high-dimensionality: large number p of variables

Joint clustering & model selection: select the K number of clusters

Model-based clustering

The rationale of model-based clustering

Observe \mathbf{Y} related to n objects

Search for $z_i \in \{0, 1\}^K$ the cluster assignment of object i

Assume $\mathbf{Z} = \{z_i\}$ contains independent and identically distributed (*i.i.d.*) discrete latent variables

$$p(\mathbf{Z} \mid \boldsymbol{\pi}) = \prod_{i=1}^n \mathcal{M}_K(z_i \mid 1, \boldsymbol{\pi})$$

Posit a statistical model on $\mathbf{Y} \mid \mathbf{Z}, \boldsymbol{\theta}$

- perform inference, e.g. maximum-likelihood, to get $(\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\theta}})$
- use the posterior $p(\mathbf{Z} \mid \mathbf{Y}, \hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\theta}})$

A fundamental assumption: conditional independence

Discrete Latent Variable Models (DLVMs)

$$p(\mathbf{Y} \mid \mathbf{Z}, \boldsymbol{\theta}) = \prod_{\mathbf{y} \in \mathbf{Y}} p(\mathbf{y} \mid \mathbf{Z}, \boldsymbol{\theta}) \quad (1)$$

Example 1: Finite Mixture Models (FMM)

Observations $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_n\}$ are *i.i.d.* inside a cluster

$$\forall i, \quad \mathbf{y}_i \mid \{z_{ik} = 1\} \sim p(\cdot \mid \boldsymbol{\theta}_k)$$

- Gaussian mixture model: $p(\mathbf{y}_i \mid \boldsymbol{\theta}_k) = \mathcal{N}_p(\mathbf{y}_i \mid \mathbf{m}_k, \mathbf{S}_k)$
- Mixture of multinomials: $p(\mathbf{y}_i \mid \boldsymbol{\theta}_k) = \mathcal{M}_p(\mathbf{y}_i \mid \boldsymbol{\theta}_k)$

$$p(\mathbf{Y} \mid \boldsymbol{\pi}, \boldsymbol{\theta}) = \prod_{i=1}^n p(\mathbf{y}_i \mid \boldsymbol{\pi}, \boldsymbol{\theta}) = \prod_{i=1}^n \sum_{k=1}^K \pi_k p(\mathbf{y}_i \mid \boldsymbol{\theta}_k)$$

A fundamental assumption: conditional independence

Discrete Latent Variable Models (DLVMs)

$$p(\mathbf{Y} \mid \mathbf{Z}, \boldsymbol{\theta}) = \prod_{\mathbf{y} \in \mathbf{Y}} p(\mathbf{y} \mid \mathbf{Z}, \boldsymbol{\theta}) \quad (1)$$

Example 2: Stochastic Block Model (SBM)

Observe n^2 edges $\mathbf{Y} = \{y_{ij}\}_{ij}$, cluster n nodes

$$\forall (i, j), \quad y_{ij} \mid \{z_{ik}z_{jl} = 1\} \sim p(\cdot \mid \boldsymbol{\theta}_{kl})$$

Edges are *i.i.d.* inside a *block* of clusters, **not marginally**

- Binary SBM: $p(y_{ij} \mid \boldsymbol{\theta}_{kl}) = \mathcal{B}(y_{ij} \mid \boldsymbol{\theta}_{kl})$
- Poisson SBM: $p(y_{ij} \mid \boldsymbol{\theta}_{kl}) = \mathcal{P}(y_{ij} \mid \boldsymbol{\theta}_{kl})$

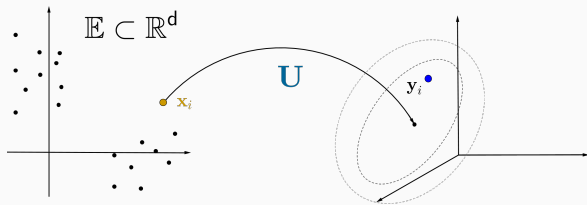
Main challenges tackled in this thesis

High-dimensional clustering: mixture estimation is cumbersome

- Gaussian mixtures: $\mathcal{O}(p^2)$ free parameters
- Small-sample scenario $n \ll p$

Probabilistic dimension reduction

$$\mathbf{x}_i \in \mathbb{E} \subset \mathbb{R}^d \longrightarrow \mathbf{y}_i \approx \mathbf{U} \mathbf{x}_i$$



Greedy clustering: discrete optimization w.r.t. \mathbf{Z}

- Joint inference and clustering
- Joint model selection and clustering

High-dimensional count data clustering

The Bayesian Fisher-EM algorithm

Hierarchical model-based clustering in DLVMs

Conclusion

High-dimensional count data clustering

MICROBIOPSIE SOUS ECHOGRAPHIE DU SEIN DROIT

MACROSCOPIE

Cinq fragments de 5 à 15 mm

MICROSCOPIE

Les prélèvements examinés correspondent à des fragments de tissu mammaire remanié par une prolifération tumorale dont les caractères morphologiques sont ceux d'un adénocarcinome canalaire infiltrant. Cette lésion est peu différenciée, d'architecture essentiellement trabéculaire. Les cellules néoplasiques comportent des atypies nucléaires marquées. L'index mitotique est élevé (22 mitoses sur 10 champs au grossissement 400). Deux fragments de 8 et 15 mm. Adénocarcinome mammaire de type canalaire infiltrant peu différencié. Grade histo-pronostique (EE) : III Index mitotique élevé.

MACROBIOPSIE DU SEIN GAUCHE

MACROSCOPIE

3 fragments de 7 à 15 mm

MICROSCOPIE

Tous les prélèvements ont un aspect histologique similaire. Ils correspondent à des fragments de tissu mammaire remanié par des lésions de mastose fibreuse commune. Présence d'un discret infiltrat inflammatoire. On retrouve également quelques microcalcifications. L'un des prélèvements cryo-préservés sera analysé histologiquement et un compte rendu complémentaire adressé ultérieurement. Trois fragments de 7 à 15 mm. Lésions de mastose fibreuse. Le prélèvement paraît peu significatif. Une analyse complémentaire sur le prélèvement cryo-préservé sera réalisée.

...

Document-term matrix					
Documents \ Terms	lésions	canalaire	...	lobulaire	métaplasie
“Lésions (...) carcino- me canalaire”	2	1	...	0	0
“Lésions bénignes (...) métaplasie”	3	0	...	0	1

Probabilistic dimension reduction for count data

Count data: $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_n\}$ with $\mathbf{y}_i \in \mathbb{N}^p$

- e.g. word count in a document, *read* count in a gene
- Total count $c_i := \sum_j y_{ij}$
- Zero inflated data, high-dimensional problems $n \ll p$

Multinomial PCA (MPCA, Buntine 2002)

$$\mathbf{x}_i \sim \mathcal{D}_d(\boldsymbol{\alpha}) \quad (\text{latent space: } \Delta_d)$$

$$\mathbf{y}_i \mid \mathbf{x}_i \sim \mathcal{M}_p(c_i, \mathbf{U} \mathbf{x}_i) \quad (\text{observation space})$$

- $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_d] \in (\Delta_p)^d$ is called the *topic* matrix
- also known as Latent Dirichlet Allocation (Blei et al. 2003)

Mixture of multinomial PCA (Yu et al. 2005)

One latent variable per cluster:

$$\mathbf{x} = (\mathbf{x}_k)_k, \quad \mathbf{x}_k \stackrel{i.i.d.}{\sim} \mathcal{D}_d(\boldsymbol{\alpha})$$
$$\forall i, \quad \mathbf{y}_i \mid \mathbf{x} \sim \sum_{k=1}^K \pi_k \mathcal{M}_p(c_i, U \mathbf{x}_k)$$

Constrained multinomial model: $\boldsymbol{\theta}_k = U \mathbf{x}_k$ (Carel et al. 2017)

Property

Suppose \mathbf{Z} known and fixed, construct K meta-observations

$$\tilde{\mathbf{Y}}_k(\mathbf{Z}) = \sum_{i=1}^n z_{ik} \mathbf{y}_i$$

Then, $\mathbf{Y} \mid \mathbf{Z}$ follows a MPCA model on $\tilde{\mathbf{Y}}(\mathbf{Z})$

Classification likelihood approach

$$\arg \max_{\mathbf{Z}, \mathbf{U}, \boldsymbol{\pi}} \left\{ \log p(\mathbf{Y}, \mathbf{Z} \mid \boldsymbol{\pi}, \mathbf{U}) = \log p(\mathbf{Z} \mid \boldsymbol{\pi}) + \underbrace{\log p(\mathbf{Y} \mid \mathbf{Z}, \mathbf{U})}_{(*) \text{ MPCA on } \tilde{\mathbf{Y}}(\mathbf{Z})} \right\}$$

Problems:

1. Combinatorics: number of partitions exponential with n
2. $(*)$ is intractable because of marginal over \mathbf{x}

Solutions:

1. Variational inference layer on \mathbf{x}

$$\mathbf{x} \sim q, \quad \log p(\mathbf{Y}, \mathbf{Z} \mid \boldsymbol{\pi}, \mathbf{U}) \geq \mathcal{J}(\mathbf{Z}, \boldsymbol{\pi}, \mathbf{U}, q)$$

2. Greedy algorithm for joint inference and clustering

Branch & bound C-VEM algorithm

Algorithm: Explore partition space using \mathcal{J} as a surrogate objective

Input: $K, d, \mathbf{Z}^{(0)}, \boldsymbol{\pi}^{(0)}, \mathbf{U}$

while \mathbf{Z} has not converged **do**

For all $i = 1, \dots, n$, try individual swaps: $z_{ik}^{(t)} = 1 \rightarrow z_{il}^{(tmp)} = 1$

// Difference with standard greedy approaches

Use variational inference to update q

$$(\mathcal{J}_l, q_l) = \arg \max_q \mathcal{J}(\mathbf{Z}^{(tmp)}, \boldsymbol{\pi}^{(t)}, \mathbf{U}, q)$$

Select $l^* = \arg \max_l \mathcal{J}_l$

$$z_{il^*}^{(t+1)} = 1, \quad q^{(t+1)} = q_{l^*} \quad \boldsymbol{\pi}^{(t+1)} = \sum_i \mathbf{z}_i^{(t+1)} / n$$

end

How to choose the pair (K, d) ?

Integrated Classification Likelihood (ICL, Biernacki et al. 2000)

$$\log p(\mathbf{Y}, \mathbf{Z}) = \int_{\pi} \int_U \log p(\mathbf{Y}, \mathbf{Z}, \pi, \mathbf{U}) d\mathbf{U} d\pi$$

ICL criterion for MMPCA

Laplace and Stirling approximations combined with a variational approximation on $p(\mathbf{Y} \mid \mathbf{Z})$ lead to

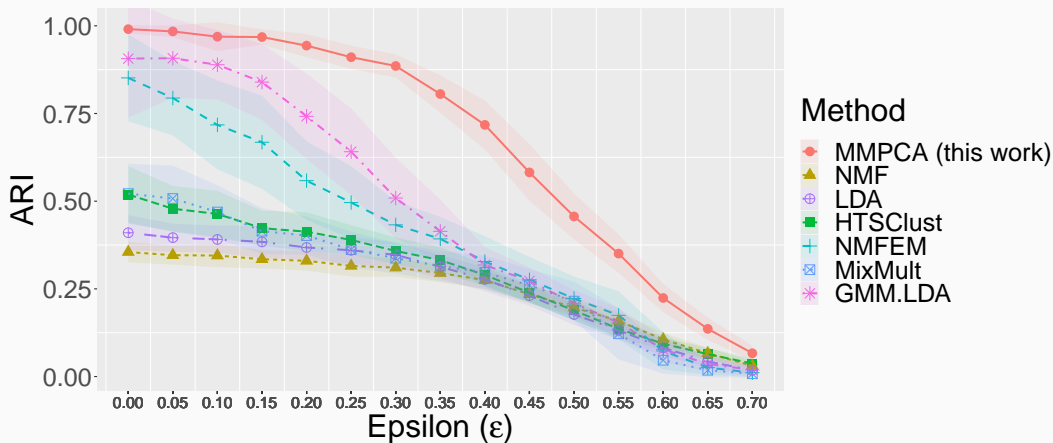
$$\begin{aligned} \text{ICL}_{MMPCA}(K, d) = & \mathcal{J}(\hat{\mathbf{Z}}, \hat{\pi}, \hat{\mathbf{U}}, \hat{q}) \\ & - \frac{d(p-1)}{2} \log(K) - \frac{K-1}{2} \log(n) \end{aligned}$$

Scenario 1: noisy setting

$n = 400,$

$p = 1000,$

Noise level: $\epsilon \in [0, 1]$

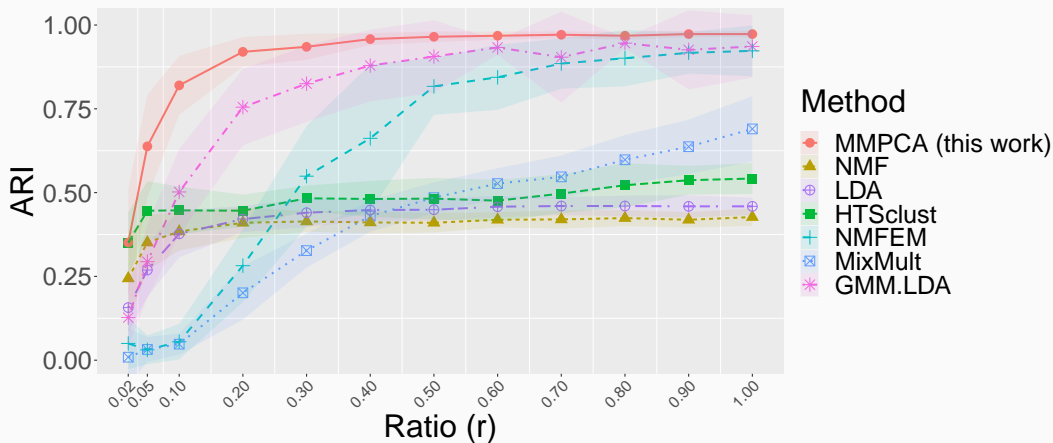


Scenario 2: small-sample sizes

$p = 1000,$

$\epsilon = 0.2,$

$n = r \times p, r \in [0, 1]$



Application: clustering of anatomopathological reports

Context: textual reports describing histopathological slides

- Benign
- Lobular carcinoma
- Non Special Type (NST) carcinoma, e.g. ductal

Unsupervised analysis: select $K = 7$ and $d = 5$

	Benign	NST carcinoma	Lobular carcinoma
1	0	0	43
2	1	31	1
3	0	106	0
4	231	3	0
5	0	211	0
6	0	126	0
7	0	113	0

Application: clustering of anatomopathological reports

Context: textual reports describing histopathological slides

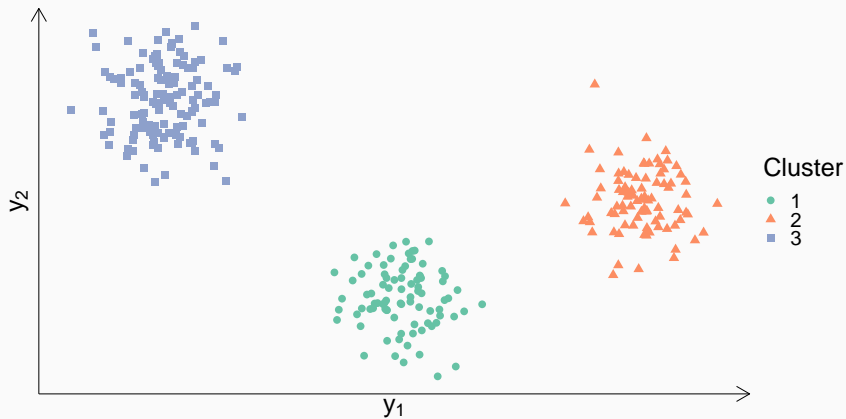
- Benign
- Lobular carcinoma
- Non Special Type (NST) carcinoma, e.g. ductal

Unsupervised analysis: select $K = 7$ and $d = 5$

	Benign	NST carcinoma	Lobular carcinoma
1	0	0	43
2	1	31	1
3	0	106	0
4	231	3	0
5	0	211	0
6	0	126	0
7	0	113	0

The Bayesian Fisher-EM algorithm

Low-dimensional mixture



What if p is large ?

Gaussian subspace clustering

Continuous data: $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_n\}$, $\mathbf{y}_i \in \mathbb{R}^p$

$$\forall i, \quad \mathbf{y}_i \stackrel{i.i.d.}{\sim} \sum_{k=1}^K \pi_k \mathcal{N}_p(\mathbf{m}_k, \mathbf{S}_k) \quad (\text{GMM})$$

Problem when p is large: over-parameterized \mathbf{S}_k

Gaussian subspace clustering

Continuous data: $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_n\}$, $\mathbf{y}_i \in \mathbb{R}^p$

$$\forall i, \quad \mathbf{y}_i \stackrel{i.i.d.}{\sim} \sum_{k=1}^K \pi_k \mathcal{N}_p(\mathbf{m}_k, \mathbf{S}_k) \quad (\text{GMM})$$

Problem when p is large: over-parameterized \mathbf{S}_k

Constrained GMM: low-rank covariance

$$\mathbf{m}_k = \mathbf{U} \boldsymbol{\mu}_k \quad \mathbf{S}_k = \mathbf{U} \boldsymbol{\Sigma}_k \mathbf{U} + \boldsymbol{\Psi}_k, \quad \mathbf{U}^\top \mathbf{U} = \mathbf{I}_d$$

Factor analysis formulation: low-dimensional embedding

$$\mathbf{x}_i \stackrel{i.i.d.}{\sim} \sum_{k=1}^K \pi_k \mathcal{N}_d(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (\text{Latent space: } \mathbb{R}^d)$$

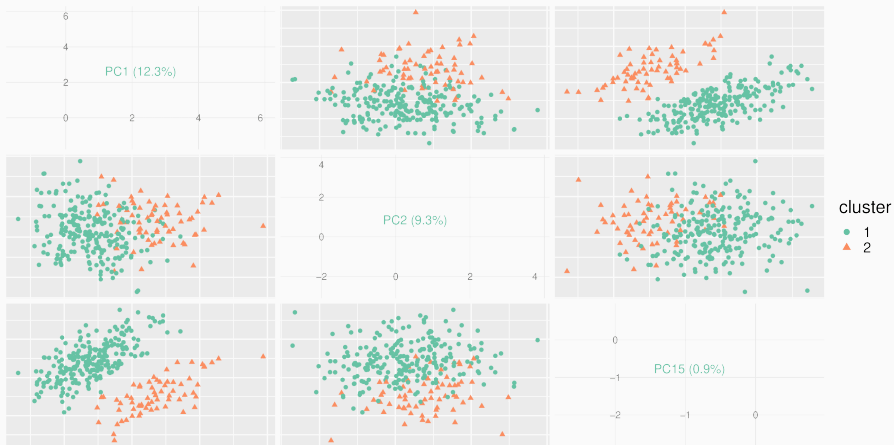
$$\mathbf{y}_i = \mathbf{U} \mathbf{x}_i + \boldsymbol{\epsilon}_{ik}, \quad \boldsymbol{\epsilon}_{ik} \sim \mathcal{N}_p(\mathbf{0}_p, \boldsymbol{\Psi}_k) \quad (\text{Observation space})$$

The tension between clustering and density estimation

Maximum-likelihood estimation (MLE):

$$(\hat{\pi}, \hat{\mu}, \hat{\Sigma}, \hat{U}) \in \arg \max_{\pi, \mu, \Sigma, U} \log p(Y \mid \pi, \mu, \Sigma, U)$$

PCA-like objective: preserves variance of the signal



Supervised: Fisher Discriminant Analysis (FDA, Fisher 1936)

1. \mathbf{Z} is known: construct scatter matrices

$$\mathbf{S}_B = \sum_k n_k (\mathbf{m}_k - \bar{\mathbf{y}})(\mathbf{m}_k - \bar{\mathbf{y}})^\top \quad (\text{between-class})$$

$$\mathbf{S}_W = \sum_k \sum_i z_{ik} (\mathbf{y}_i - \mathbf{m}_k)(\mathbf{y}_i - \mathbf{m}_k)^\top \quad (\text{within-class})$$

2. maximize their ratio in the latent space, $d \leq K - 1$

$$\hat{\mathbf{U}} = \arg \max_{\mathbf{U}} \left\{ F(\mathbf{U}) := \text{Tr} \left[(\mathbf{U}^\top \mathbf{S}_W \mathbf{U})^{-1} \mathbf{U}^\top \mathbf{S}_B \mathbf{U} \right] \right\}$$

Discriminative subspace

Supervised: Fisher Discriminant Analysis (FDA, Fisher 1936)

1. \mathbf{Z} is known: construct scatter matrices

$$\mathbf{S}_B = \sum_k n_k (\mathbf{m}_k - \bar{\mathbf{y}})(\mathbf{m}_k - \bar{\mathbf{y}})^\top \quad (\text{between-class})$$

$$\mathbf{S}_W = \sum_k \sum_i z_{ik} (\mathbf{y}_i - \mathbf{m}_k)(\mathbf{y}_i - \mathbf{m}_k)^\top \quad (\text{within-class})$$

2. maximize their ratio in the latent space, $d \leq K - 1$

$$\hat{\mathbf{U}} = \arg \max_{\mathbf{U}} \left\{ F(\mathbf{U}) := \text{Tr} \left[(\mathbf{U}^\top \mathbf{S}_W \mathbf{U})^{-1} \mathbf{U}^\top \mathbf{S}_B \mathbf{U} \right] \right\}$$

Unsupervised: use posterior (Bouveyron et al. 2012)

$$\tau_{ik} = p(z_{ik} = 1 \mid \mathbf{y}_i, \boldsymbol{\pi}, \boldsymbol{\Sigma}_k, \beta_k, \mathbf{U})$$

Bayesian discriminative latent mixture model (BDLM)

Problem: algorithmic instability, bad conditioning of the scatter matrices

New Gaussian subspace clustering model:

$$\begin{aligned}\boldsymbol{\mu} &= (\boldsymbol{\mu}_k), \quad \boldsymbol{\mu}_k \stackrel{i.i.d.}{\sim} \mathcal{N}_d(\boldsymbol{\nu}, \lambda \mathbf{I}_d) \\ \mathbf{y}_i \mid \boldsymbol{\mu} &\stackrel{i.i.d.}{\sim} \sum_k \pi_k \mathcal{N}_p(\mathbf{U} \boldsymbol{\mu}_k, \underbrace{\mathbf{U} \boldsymbol{\Sigma}_k \mathbf{U} + \boldsymbol{\Psi}_k}_{\mathbf{S}_k})\end{aligned}$$

- λ controls the separation in the latent space
- \mathbf{U} is supposed to be **discriminative**
- Block diagonal hypothesis on $\mathbf{S}_k = \mathbf{D} \boldsymbol{\Delta}_k \mathbf{D}^\top$:

$$\boldsymbol{\Delta}_k = \begin{pmatrix} \boldsymbol{\Sigma}_k & 0 \\ 0 & \beta_k \mathbf{I}_{p-d} \end{pmatrix}, \mathbf{D} = [\mathbf{U}, \mathbf{U}^\perp]$$

- Possible constraints on $\boldsymbol{\Sigma}_k$, 12 submodels

Joint inference and clustering

Objective: joint MLE and FDA

- Maximize likelihood with respect to π, Σ, β
- Update τ_{ik} and maximize $F(\mathbf{U})$ w.r.t $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_d$

Problems:

1. orthonormal FDA \rightarrow no closed-form solution
2. intractable likelihood because of marginalization on (\mathbf{Z}, μ)

Solutions:

1. Use variational inference layer on (\mathbf{Z}, μ)

$$(\mathbf{Z}, \mu) \sim q \quad \log p(\mathbf{Y} \mid \pi, \Sigma, \beta, \mathbf{U}) \geq \mathcal{J}(\pi, \Sigma, \beta, \mathbf{U}, q)$$

2. Use iterative procedure solving 1-D FDA problems

Bayesian Fisher-EM algorithm (BFEM)

Fix $\mathbf{U}^{(0)}$ and $(\boldsymbol{\pi}, \boldsymbol{\Sigma}, \boldsymbol{\beta})^{(0)}$ and iterate over

- **VE-step:** Find

$$q^{(t+1)} = \arg \max_q \mathcal{J}(\boldsymbol{\pi}^{(t)}, \boldsymbol{\Sigma}^{(t)}, \boldsymbol{\beta}^{(t)}, \mathbf{U}^{(t)}, q)$$

- **M-step:** Find

$$(\boldsymbol{\pi}, \boldsymbol{\Sigma}, \boldsymbol{\beta})^{(t+1)} = \arg \max_{\boldsymbol{\pi}, \boldsymbol{\Sigma}, \boldsymbol{\beta}} \mathcal{J}(\boldsymbol{\pi}, \boldsymbol{\Sigma}, \boldsymbol{\beta}, \mathbf{U}^{(t)}, q^{(t+1)})$$

- **F-step:** use $q^{(t+1)}(\mathbf{Z})$ to construct \mathbf{S}_W and \mathbf{S}_B , then

$$\begin{aligned} \mathbf{U}^{(t+1)} &:= [\mathbf{u}_1 \mid \dots \mid \mathbf{u}_d], \text{ with } \forall h = 1, \dots, d, \\ \mathbf{u}_h &= \arg \max_{\mathbf{u} \in \mathbb{R}^p} F(\mathbf{u}) \text{ s.t. } \mathbf{u} \in \left\{ \forall r < h, \mathbf{u}^\top \mathbf{u}_r = 0 \right\} \end{aligned}$$

Questions:

1. How to set (ν, λ) ? $\lambda \rightarrow +\infty \implies$ frequentist setting
2. How to choose K and a submodel \mathcal{M} ?

Empirical Bayes:

$$(\hat{\nu}, \hat{\lambda}) = \arg \max_{\nu, \lambda} \mathcal{J}(\nu, \lambda)$$

ICL criterion for BFEM

Denote $\gamma_{\mathcal{M}, K}$ the number of free parameters in model \mathcal{M} with K clusters

$$\text{ICL}_{BIC}(\mathcal{M}, K) = \log p(\mathbf{Y}, \hat{\mathbf{Z}} \mid \hat{\boldsymbol{\vartheta}}, \mathcal{M}, K) - \frac{\gamma_{\mathcal{M}, K}}{2} \log(n),$$

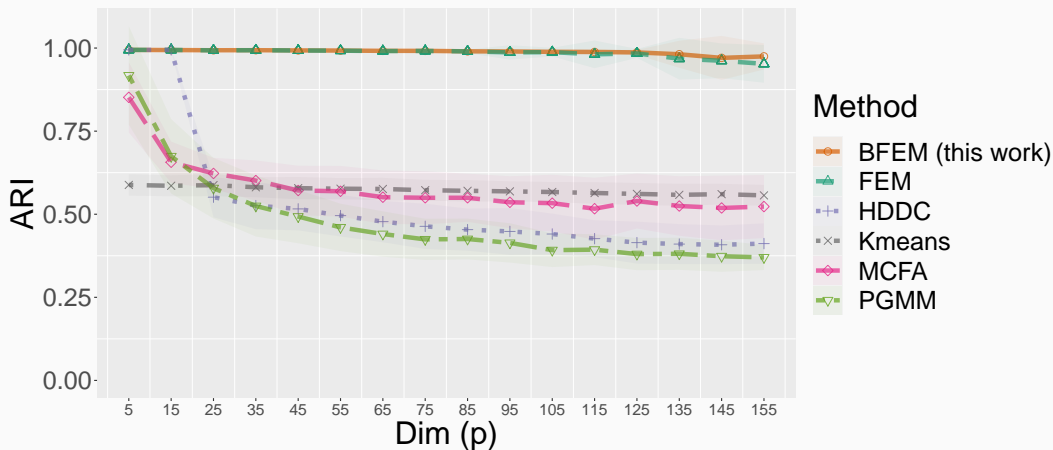
Scenario 1: increasing dimension

$n = 900,$

$d = 2,$

$\beta_k = 1,$

$p \in \{5, 15, \dots, 155\}$



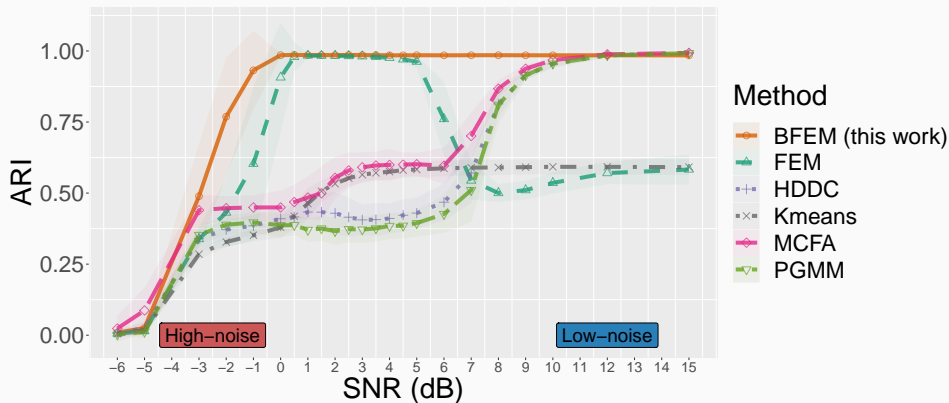
Scenario 2: signal-to-noise ratio

$n = 900,$

$d = 2,$

$p = 150,$

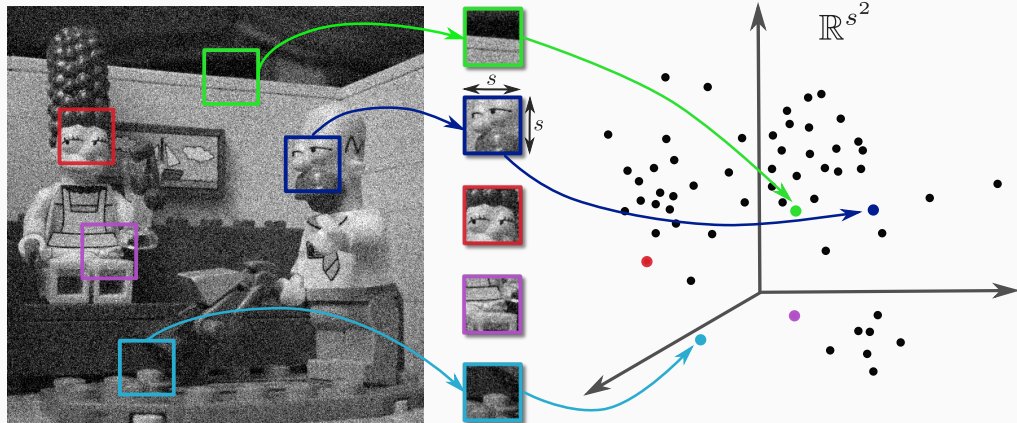
$\beta \in \{8, \dots, 0.8, \dots, 0.08\}$



Application: patch-based image denoising (Houdard et al. 2018)

$$I = I_0 + N,$$

$$N \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

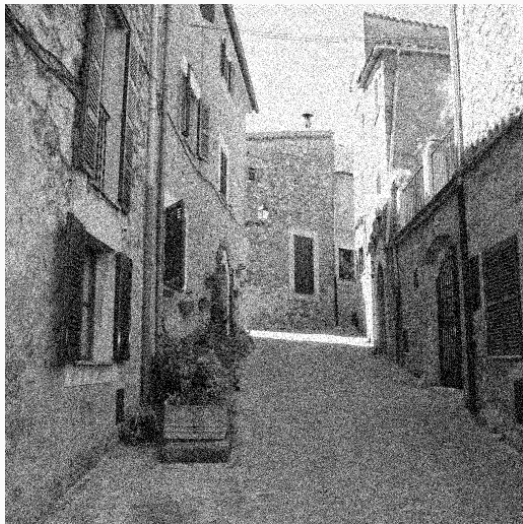


Alley image

Original



Noisy, $\sigma = 30$



Alley image

S-PLE, $PSNR = 28.22dB$



BFEM, $PSNR = 28.95dB$



Hierarchical model-based clustering in DLVMs

Overview of contributions

- Generic approach: applies in the framework of DLVMs
- Model selection criterion as a clustering objective

$$\text{ICL}_{\text{ex}}(\mathbf{Z}) = \log \int_{\pi} \int_{\theta} p(\mathbf{Y}, \mathbf{Z}, \theta, \pi) d\theta d\pi$$

Two contributions:

1. **Genetic algorithm**: greedy maximization w.r.t \mathbf{Z}
 - Based on selection mechanisms: *mutation* and *cross-over* operators
 - Perform clustering and model selection, return $\mathbf{Z}^{(K^*)}$
 - Bypass inference
2. **Hierarchical algorithm**: start from $\mathbf{Z}^{(K^*)}$ and merge clusters

$$\mathbf{Z}^{(K^*)} \leq \dots \leq \mathbf{Z}^{(1)}$$

Overview of contributions

- Generic approach: applies in the framework of DLVMs
- Model selection criterion as a clustering objective

$$\text{ICL}_{\text{ex}}(\mathbf{Z}) = \log \int_{\boldsymbol{\pi}} \int_{\boldsymbol{\theta}} p(\mathbf{Y}, \mathbf{Z}, \boldsymbol{\theta}, \boldsymbol{\pi}) \, \mathrm{d}\boldsymbol{\theta} \, \mathrm{d}\boldsymbol{\pi}$$

Two contributions:

1. **Genetic algorithm**: greedy maximization w.r.t \mathbf{Z}
 - Based on selection mechanisms: *mutation* and *cross-over* operators
 - Perform clustering and model selection, return $\mathbf{Z}^{(K^*)}$
 - Bypass inference
2. **Hierarchical algorithm**: start from $\mathbf{Z}^{(K^*)}$ and merge clusters

$$\mathbf{Z}^{(K^*)} \leq \dots \leq \mathbf{Z}^{(1)}$$

Exact integrated classification likelihood

Proposition (Fubini)

With a factorized prior: $p(\boldsymbol{\theta}, \boldsymbol{\pi}) = p(\boldsymbol{\theta} \mid \boldsymbol{\beta}) p(\boldsymbol{\pi} \mid \boldsymbol{\alpha})$

$$\text{ICL}_{\text{ex}}(\mathbf{Z}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \underbrace{\log p(\mathbf{Y} \mid \mathbf{Z}, \boldsymbol{\beta})}_{(1)} + \underbrace{\log p(\mathbf{Z} \mid \boldsymbol{\alpha})}_{(2)}$$

Conjugate prior for exact (1) available in standard DLVMs

- MoM (Biernacki et al. 2010)
- Binary SBM (Côme et al. 2015)
- GMM (Bertoletti et al. 2015), modulo *informative* prior

Suppose $\boldsymbol{\beta}$ fixed and denote

$$D(\mathbf{Z}) := \log p(\mathbf{Y} \mid \mathbf{Z}, \boldsymbol{\beta})$$

Proposition (Fubini)

With a factorized prior: $p(\boldsymbol{\theta}, \boldsymbol{\pi}) = p(\boldsymbol{\theta} \mid \boldsymbol{\beta}) p(\boldsymbol{\pi} \mid \boldsymbol{\alpha})$

$$\text{ICL}_{\text{ex}}(\mathbf{Z}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \underbrace{\log p(\mathbf{Y} \mid \mathbf{Z}, \boldsymbol{\beta})}_{(1)} + \underbrace{\log p(\mathbf{Z} \mid \boldsymbol{\alpha})}_{(2)}$$

Exact expression of (2) with universal prior

$$p(\boldsymbol{\pi} \mid \boldsymbol{\alpha}) = \mathcal{D}_K(\boldsymbol{\pi} \mid \boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_K))$$

Set $\forall k, \alpha_k = \alpha$,

- $\alpha = 1$: uniform prior on the simplex
- $\alpha = 1/2$: Jeffreys prior

A sketch of the hierarchical algorithm

Standard agglomerative method

- Starts from $\mathcal{Z}^{(K)}$ with $K \leq n$ cluster
- At each stage, find the best fusion w.r.t ICL_{ex}

Problem: a fusion is not always possible

Solution:

- Use α as a regularization parameter
- Extract a set of dominating *nested* partitions

$$\log p(\mathbf{Z} \mid \alpha) = \log \frac{\Gamma(K\alpha) \prod_k \Gamma(\alpha + n_k)}{\Gamma(\alpha)^K \Gamma(n + \alpha K)}, \quad n_k = \sum_i z_{ik}$$

Our proposition: asymptotic of $\log \Gamma$ near 0

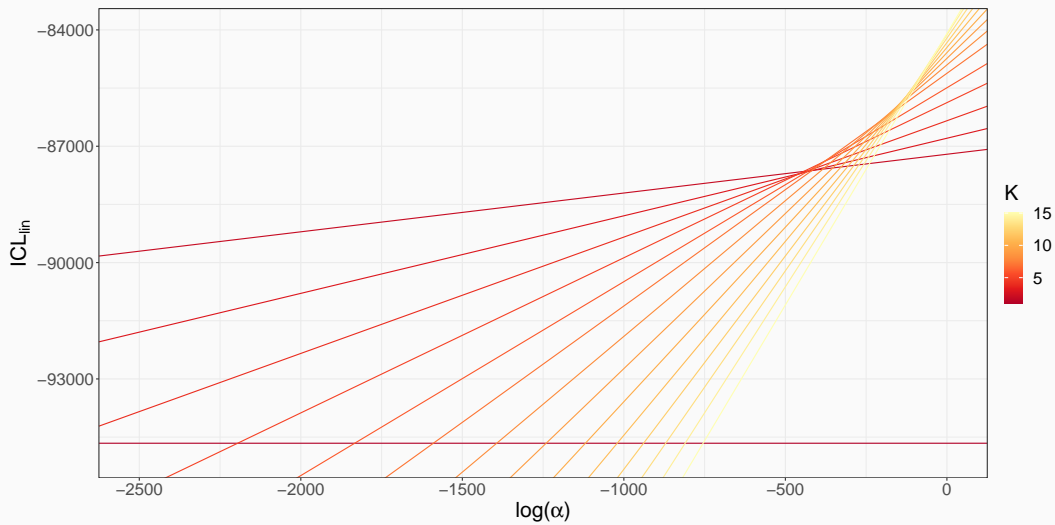
$$\log \Gamma(\alpha) \underset{\alpha \rightarrow 0}{\sim} -\log(\alpha)$$

Log-linear ICL

$$\text{ICL}_{\text{lin}}(\mathbf{Z}, \alpha) := (K - 1) \log(\alpha) + I(\mathbf{Z})$$

$$I(\mathbf{Z}) = D(\mathbf{Z}) + \sum_k \log \Gamma(n_k) - \log \Gamma(n) - \log(K)$$

ICL_{lin} increasing slope with K



Fusion opportunity at stage (k)

Fixed partition $\mathbf{Z}^{(k)}$ with k clusters

Two clusters (g, h): ICL_{lin} change for $g \cup h$?

$$\Delta_{g \cup h}(\alpha) = \text{ICL}_{lin}(\mathbf{Z}_{g \cup h}, \alpha) - \text{ICL}_{lin}(\mathbf{Z}^{(k)}, \alpha)$$

Proposition

$$\forall g \neq h, \Delta_{g \cup h}(\alpha) > 0 \iff \log(\alpha) < I(\mathbf{Z}_{g \cup h}) - I(\mathbf{Z}^{(k)})$$

Regularization parameter: α unlocks fusions

Question: $k(k-1)/2$ fusions, which one is the best ?

$$(g^*, h^*) = \arg \max_{g, h} I(\mathbf{Z}_{g \cup h})$$

Hierarchy construction and dendrogram representation

Repeat procedure at each stage $\mathbf{Z}^{(k)}$

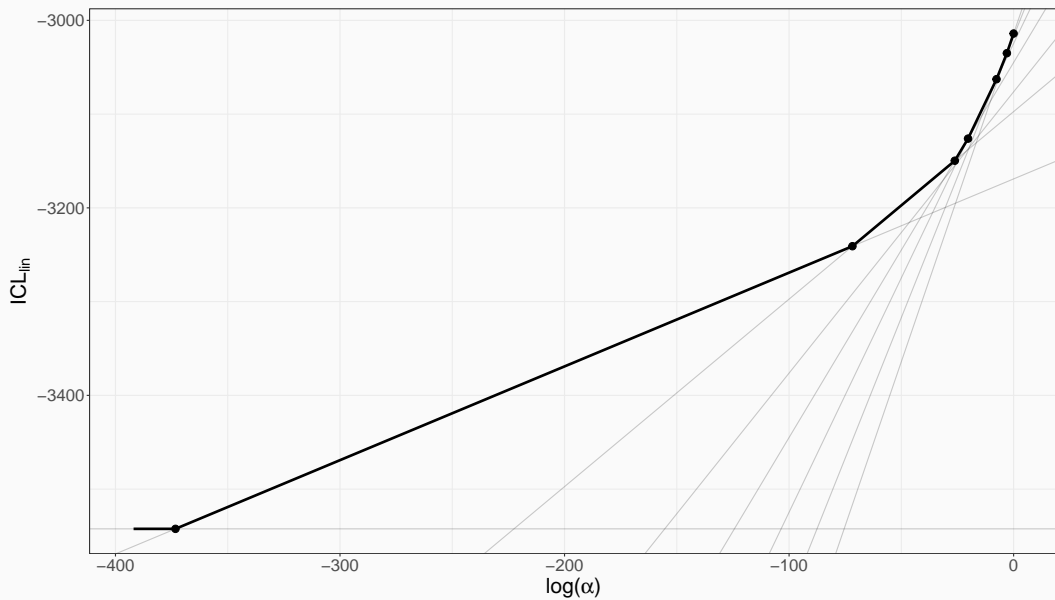
$$\log \alpha^{(k)} := I(\mathbf{Z}_{g^* \cup h^*}) - I(\mathbf{Z}^{(k)})$$

Outputs a hierarchy of partitions

Dendrogram representation:

- $\alpha^{(k)}$ is the amount of regularization needed for the fusion
- Extract a **front** of dominating partitions on range $[\alpha^{(k-1)}, \alpha^{(k)}]$

A discrete Pareto frontier



Simulation scenario: graph clustering

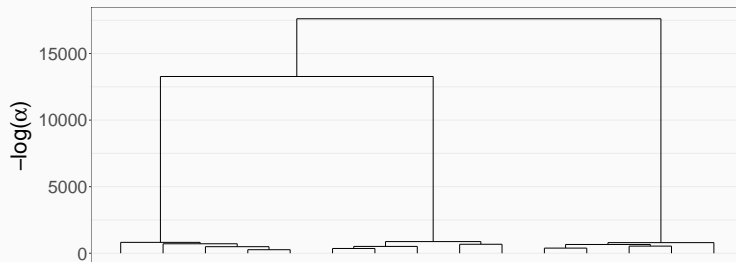
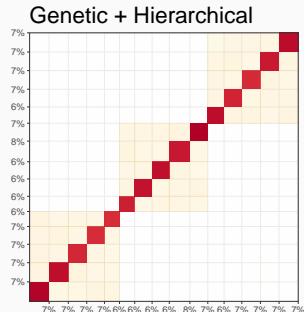
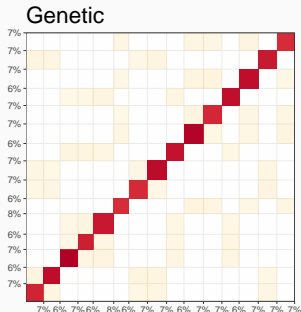
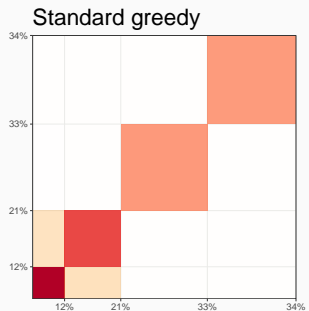
Simulate according SBM with nested structure

- $n = 1500$ nodes, $K = 15$
- 3 clusters composed of 5 smaller ones

Two-step methodology

1. Maximization of ICL_{ex} (genetic algorithm): find $\mathbf{Z}^{(K)}$
2. Hierarchy construction using ICL_{lin} and α

Link density 0.00 0.02 0.04 0.06



Conclusion

Clustering algorithms using DLVMs

- ▶ High-dimensional count data (Jouvin et al. 2020)
- ▶ High-dimensional continuous data, *preprint*, submitted to journal
- ▶ Graph data (and co-clustering), *preprint*, submitted to journal

Applications: medical data, image denoising, graph clustering

Reproducible research: R packages available

- MoMPCA
- FisherEM
- **greed** (E. Côme)

Clustering categorical data

- survey, census
- Mixture of Multinomial multiple correspondence analysis

$$\mathbf{y}_i \sim \mathcal{M}_p(1, \boldsymbol{\theta}_i), \quad \boldsymbol{\theta}_i = \text{softmax}(\boldsymbol{\beta} + \mathbf{U}\mathbf{x}_i), \quad \mathbf{x}_i \sim \text{GMM}(\mathbb{R}^d)$$

Extensions to BFEM

- Sparsity on \mathbf{U} through l_1 penalty \rightarrow **variable selection**
- Other formulations of orthonormal FDA

Exact ICL for GMM \rightarrow handling informative prior

Thank you for your attention !

Questions

1. Appendix MoMPCA [Go to](#)

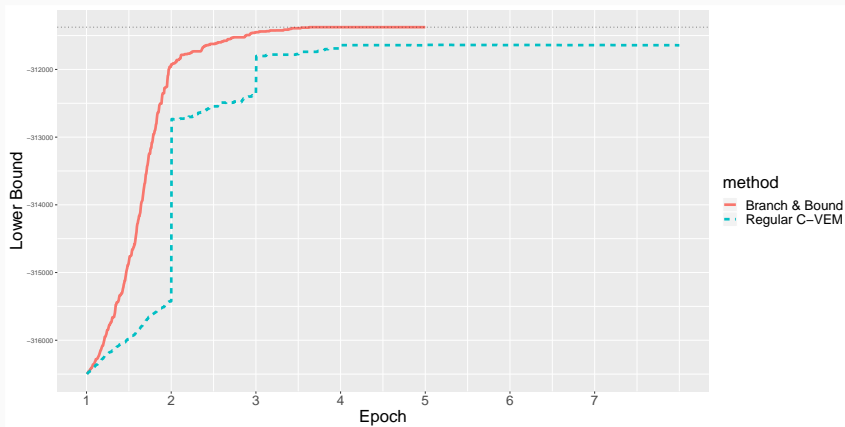
2. Appendix BFEM [Go to](#)

3. Appendix HC-ICL [Go to](#)

MoMPCA (Appendix)

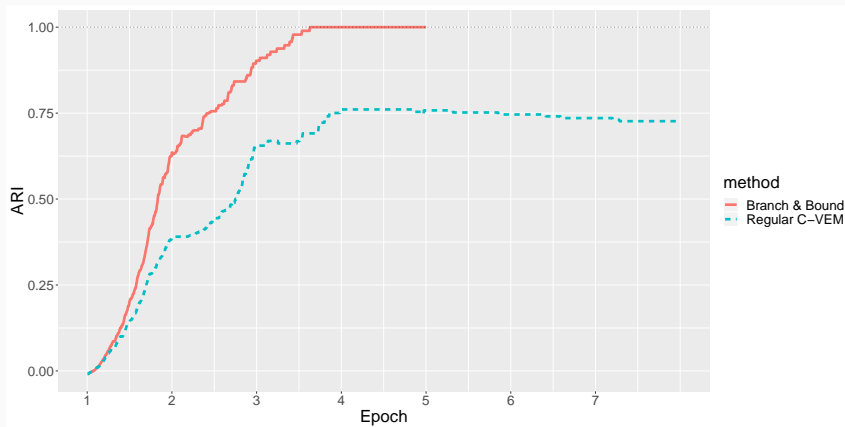
Branch & Bound VS standard C-VE

Standard C-VE: no variational inference step after a swap



Branch & Bound VS standard C-DEM

Standard C-DEM: no variational inference step after a swap



MoMPCA: two simulation scenarios

Fixed setting:

$$p = 1000, \quad K = 6, \quad d = 4, \quad \mathbf{U}^*, \quad \mathbf{x}^*, \quad \forall i, c_i = 400$$

Scenario 1: noisy structure [Goto](#) $n = 400$

$$\mathbf{x}_{\epsilon,k} = (1 - \epsilon)\mathbf{x}_k^* + \frac{\epsilon}{d} \underbrace{(1, \dots, 1)}_d^\top, \quad \epsilon \in [0, 1]$$

- $\epsilon = 0 \rightarrow \mathbf{x}_{0,k} = \mathbf{x}_k^*$ distribution across topics
- $\epsilon = 1 \rightarrow \mathbf{x}_{1,k}$ uniform across topics (no cluster structure)

Scenario 2: small-sample size [Goto](#) $\epsilon = 0.2$

$$n = r \times p, \quad r \in [0, 1]$$

Metric: *adjusted Rand index* (ARI) the higher, the better

	Topic1	Topic2	Topic3	Topic4	Topic5
x_1	0.00	0.01	0.98	0.00	0.00
x_2	0.19	0.11	0.04	0.38	0.29
x_3	0.13	0.09	0.01	0.76	0.00
x_4	0.01	0.00	0.01	0.01	0.97
x_5	0.00	1.00	0.00	0.00	0.00
x_6	0.05	0.65	0.03	0.26	0.01
x_7	0.74	0.12	0.03	0.11	0.00

Cluster 2 contains micro-calcifications and peaked towards

- Topic4: vocabulary of *in-situ* lesions
- Topic5: vocabulary of benign lesions

Posterior explanation: all samples came from macro-biopsy exams

Bayesian Fisher EM (appendix)

BFEM: two simulation scenarios

Fix $K^* = 3$, $d^* = 2$, $n = 900$, $\boldsymbol{\pi}^*$, $\boldsymbol{\Sigma}_k^* = \boldsymbol{\Sigma}^*$

Scenario 1: increasing dimension [Goto](#) Fix $\beta_k = \beta = 1$, increase p

Scenario 2: signal-to-noise [Goto](#) Fix $p = 150$, increase β

$$SNR = 10 \times \log_{10} \left(\frac{\text{Tr} [\boldsymbol{\Sigma}]}{\beta} \right)$$

- $SNR \gg 0$: Noiseless regime
- $SNR \ll 0$: No signal

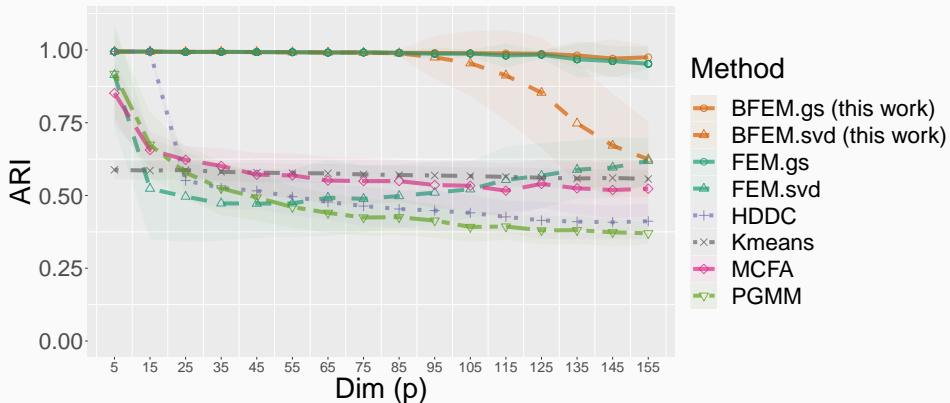
BFEM scenario 1: graph with SVD method

$$n = 900,$$

$$d = 2,$$

$$\beta_k = 1,$$

$$p \in \{5, 15, \dots, 155\}$$



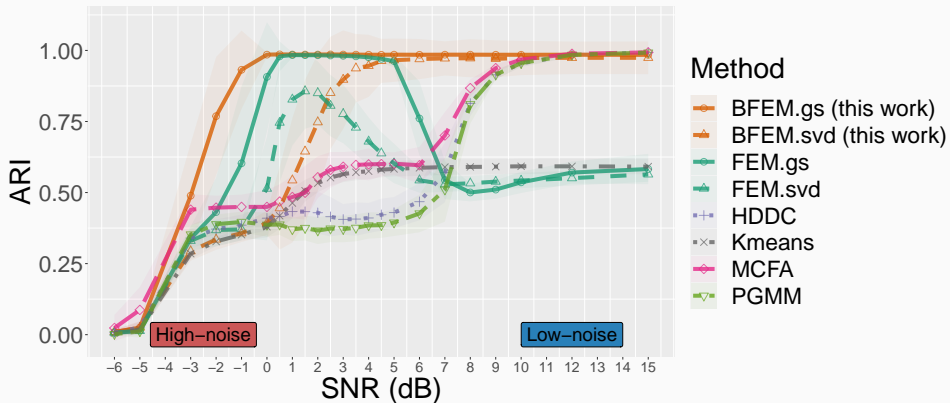
BFEM scenario 2: graph with SVD method

$n = 900,$

$d = 2,$

$p = 150,$

$\beta \in \{8, \dots, 0.8, \dots, 0.08\}$



BFEM performances on real clustering datasets

Dataset	p	Non-HD models		HD models				
		k -means	Mclust	HDDC	MCFA	PGMM	FEM	BFEM
<i>Iris</i>	4	0.73	0.90	0.90	0.92	0.94	0.88	0.90
<i>Wine 27</i>	27	0.90	0.93	0.95	0.96	0.98	0.93	0.93
<i>Satellite</i>	36	0.53	0.36	0.45	0.43	0.56	0.53	0.64
<i>USPS358</i>	256	0.64	0	0.35	0.28	0.38	0.66	0.76

Application to single image denoising

Context: Observe I_0 blurred with Gaussian white noise

$$I = I_0 + N, \quad N \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

Patch denoising: decompose I into $f \times f$ sub-images $\mathbf{y}_i \in \mathbb{R}^{f^2}$

- GMM prior on the true patch \mathbf{t}_i

$$\mathbf{y}_i = \mathbf{t}_i + \epsilon_i \sim \sum_k \pi_k \mathcal{N}_{f^2}(\mathbf{m}_k, \boldsymbol{\phi}_k + \sigma^2 \mathbf{I}_{f^2})$$

- Optimal estimate of \mathbf{t}_i w.r.t. quadratic risk

$$\hat{\mathbf{t}}_i = \mathbb{E}[\mathbf{t}_i \mid \mathbf{y}_i] = \sum_{k=1}^K \tau_{ik} \mathbb{E}[\mathbf{t}_i \mid \mathbf{y}, z_{ik} = 1, \pi_k, \boldsymbol{\phi}_k, \sigma^2]$$

Use **BFEM** to estimate $\boldsymbol{\pi}$, τ_{ik} and $\boldsymbol{\phi}_k = \mathbf{U} \boldsymbol{\Sigma}_k \mathbf{U}^\top$

Clustering and hierarchical clustering in DLVMs (appendix)

Genetic algorithm

Works with ICL_{ex} , α is fixed e.g. to 1 or 1/2

Existing works: find local maxima of ICL_{ex} with greedy swaps and merge

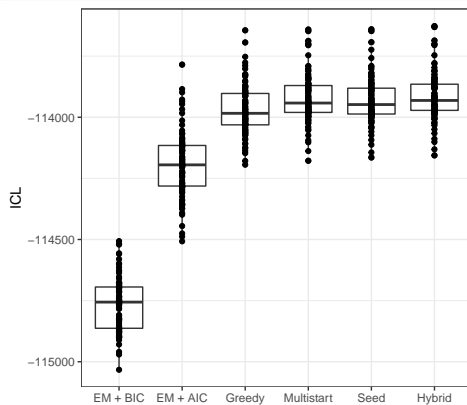
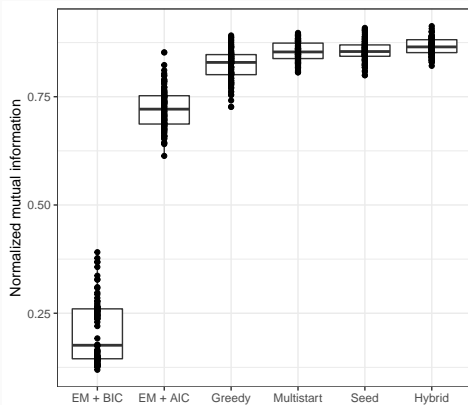
Problem: too many spurious local maxima

Contribution: genetic algorithm to improve exploration

- Recombination: cross-partition operator
- Mutation: cluster split

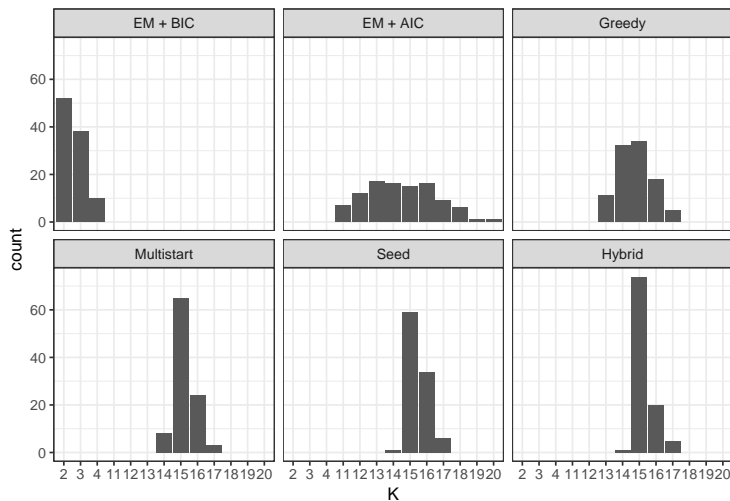
Genetic algorithm: medium-scale MoM

$n = 500$, $p = 100$, $K = 15$, θ_k peaked toward 10 variables



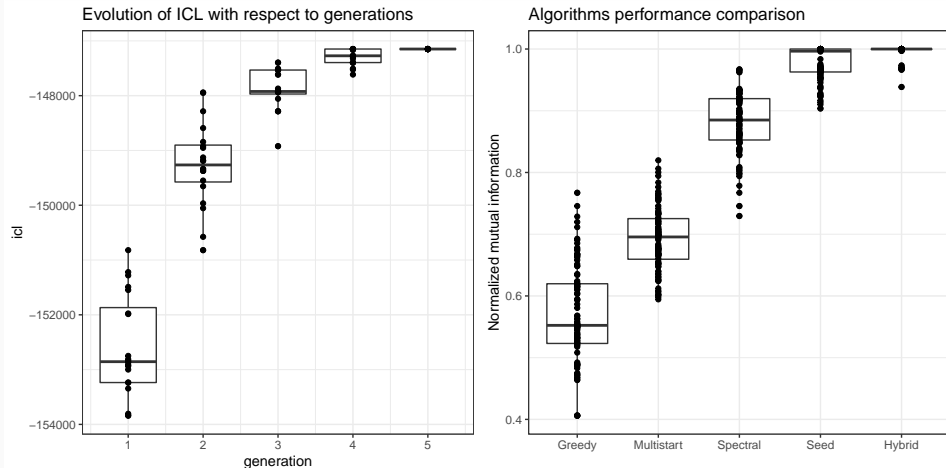
Genetic algorithm: medium-scale MoM

$n = 500$, $p = 100$, $K = 15$, θ_k peaked toward 10 variables



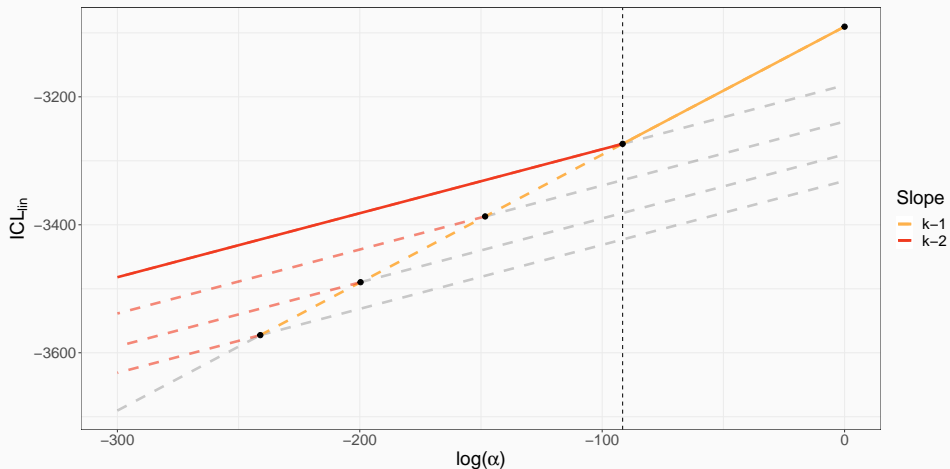
Genetic algorithm: medium-scale SBM

Hierarchical nested SBM with $K = 15$ and $n = 1500$



Choosing best fusion at stage (k)

$$\forall g \neq h, \Delta_{g \cup h}(\alpha) > 0 \iff \log(\alpha) < I(\mathbf{Z}_{g \cup h}) - I(\mathbf{Z}^{(k)})$$



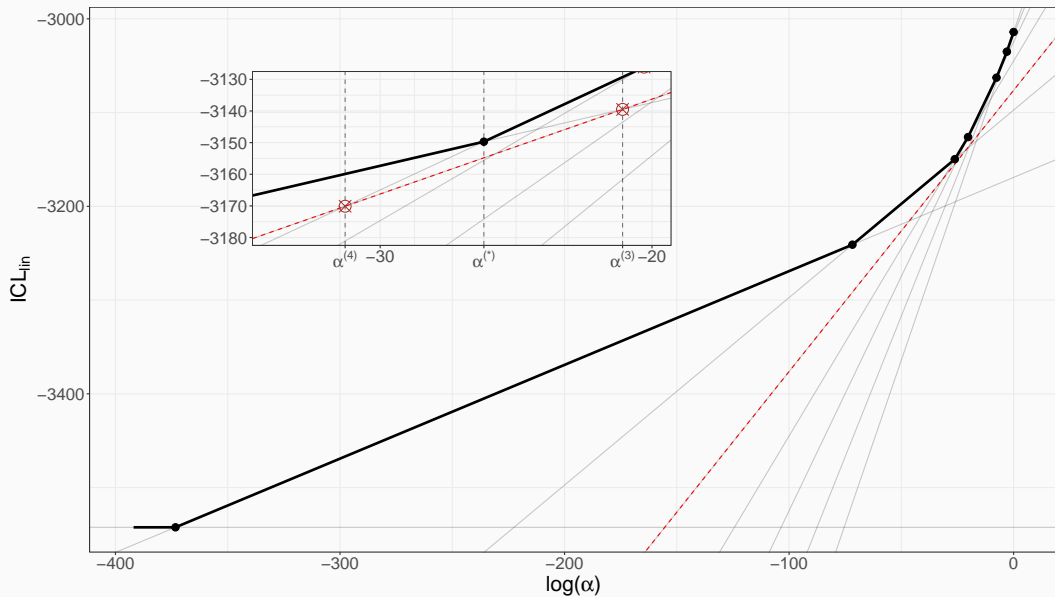
Backtracking step

Question: do we have $\alpha^{(1)} \leq \dots \leq \alpha^{(K)}$?

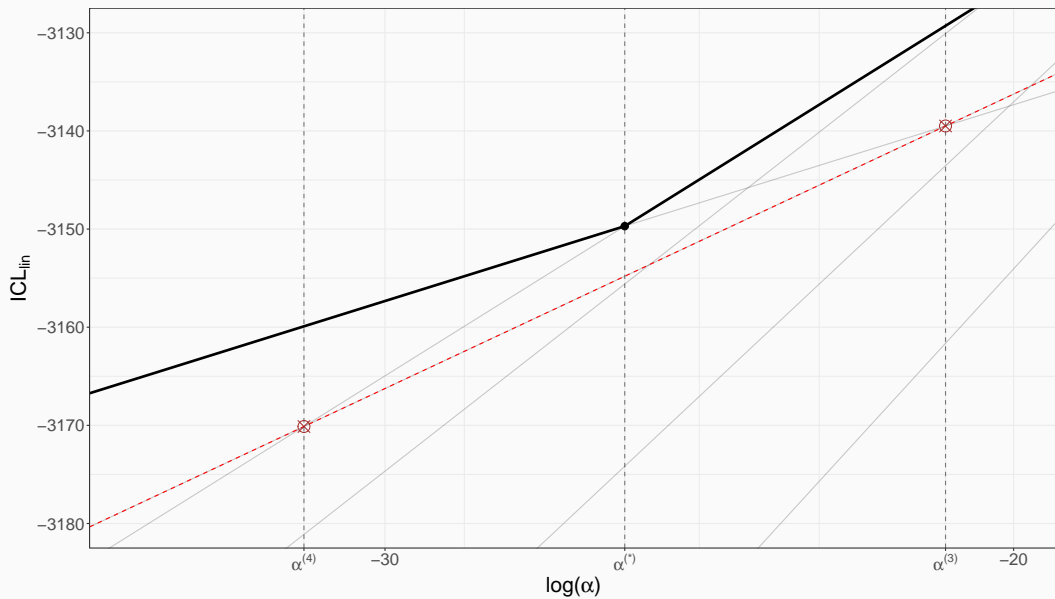
Answer:

- not necessarily, some $Z^{(k)}$ can be nowhere dominant
- easily tracked: corresponds to $\alpha^{(k-1)} \geq \alpha^{(k)}$
- happens when several fusions are possible at once
- compute the new sequence to get the **dendrogram** representation

Backtrack nowhere dominant partitions



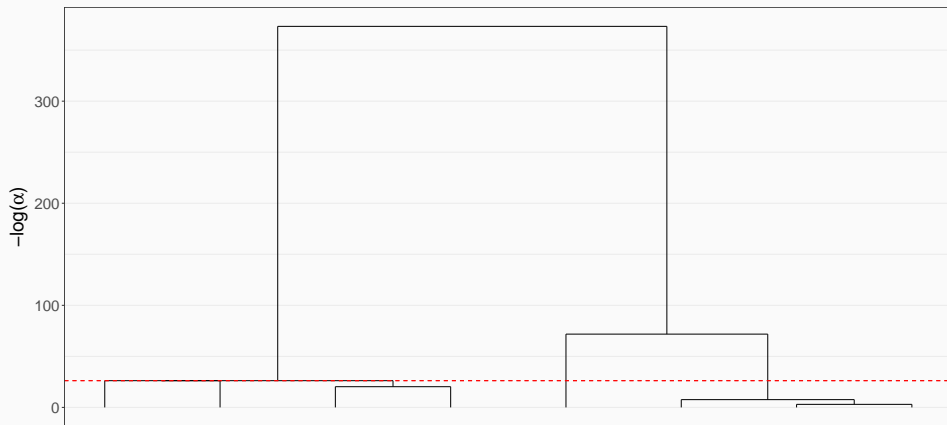
Backtrack nowhere dominant partitions: fusion $5 \rightarrow 3$







Dendrogram representation: several fusions at α^*







The partition $\mathcal{Z}^{(4)}$ is not in the Pareto front



Remove and compute the intersection α^* between $\mathcal{Z}^{(3)}$ and $\mathcal{Z}^{(5)}$.



References

-  Bertoletti, Marco, Nial Friel, and Riccardo Rastelli (Aug. 2015). “Choosing the number of clusters in a finite mixture model using an exact integrated completed likelihood criterion”. In: *METRON* 73.2, pp. 177–199.
-  Biernacki, Christophe, Gilles Celeux, and Gérard Govaert (2000). “Assessing a mixture model for clustering with the integrated completed likelihood”. In: *IEEE transactions on pattern analysis and machine intelligence* 22.7, pp. 719–725.
-  Biernacki, Christophe, Gilles Celeux, and Gerard Govaert (2010). “Exact and monte carlo calculations of integrated likelihoods for the latent class model”. In: *Journal of Statistical Planning and Inference* 140, pp. 2991–3002.
-  Blei, David M, Andrew Y Ng, and Michael I Jordan (2003). “Latent dirichlet allocation”. In: *Journal of machine Learning research* 3.Jan, pp. 993–1022.

-  Bouveyron, Charles and Camille Brunet (2012). "Simultaneous model-based clustering and visualization in the Fisher discriminative subspace". In: *Statistics and Computing* 22.1, pp. 301–324.
-  Buntine, Wray (2002). "Variational extensions to EM and multinomial PCA". In: *European Conference on Machine Learning*. Springer, pp. 23–34.
-  Carel, Léna and Pierre Alquier (2017). "Simultaneous dimension reduction and clustering via the NMF-EM algorithm". In: *Advances in Data Analysis and Classification*, pp. 1–30.
-  Côme, Etienne and Pierre Latouche (2015). "Model selection and clustering in stochastic block models based on the exact integrated complete data likelihood". In: *Statistical Modelling* 15.6, pp. 564–589.
-  Fisher, Ronald A (1936). "The use of multiple measurements in taxonomic problems". In: *Annals of eugenics* 7.2, pp. 179–188.
-  Houdard, Antoine, Charles Bouveyron, and Julie Delon (2018). "High-dimensional mixture models for unsupervised image denoising (HDMI)". In: *SIAM Journal on Imaging Sciences* 11.4, pp. 2815–2846.

-  Jouvin, Nicolas et al. (2020). “Greedy clustering of count data through a mixture of multinomial PCA”. In: *Computational Statistics*.
-  Yu, Shipeng et al. (2005). “A probabilistic clustering-projection model for discrete data”. In: *European Conference on Principles of Data Mining and Knowledge Discovery*. Springer, pp. 417–428.