

# Going further in cluster analysis and classification: Bi-clustering and co-clustering

C Biernacki

► **To cite this version:**

C Biernacki. Going further in cluster analysis and classification: Bi-clustering and co-clustering. Summer School on Clustering, Data Analysis and Visualization of Complex Data, May 2018, Catania, Italy. hal-01810380

**HAL Id: hal-01810380**

**<https://hal.inria.fr/hal-01810380>**

Submitted on 7 Jun 2018

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Going further in cluster analysis and classification: **Bi-clustering and co-clustering**

C. Biernacki

Summer School on Clustering, Data Analysis and Visualization of Complex Data  
May 21-25 2018, University of Catania, Italy



# Outline

**1** HD clustering

2 Modeling

3 Estimating

4 Selecting

5 BlockCluster in MASSICCC

6 To go further

## Motivation

High dimensional (HD) data sets are now frequent:

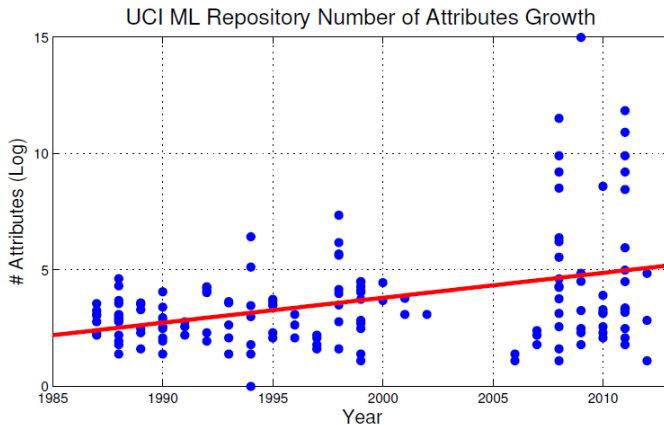
- Marketing:  $d \sim 10^2$
- microarray gene expression:  $d \sim 10^2-10^4$
- SNP data:  $d \sim 10^6$
- Curves: depends on discretization but can be very high
- Text mining
- ...

Clustering has to be applied for HD datasets for the same reasons as the lower dimensional datasets:

- Data summary
- Data exploratory
- Preprocessing for more flexibility of a forthcoming prediction step

But clustering is even more important since visualization in the HD setting can be hazardous. . .

# Today': exponential growing of dimension<sup>1</sup>



<sup>1</sup>S. Alelyani, J. Tang and H. Liu (2013). Feature Selection for Clustering: A Review. *Data Clustering: Algorithms and Applications*, 29

## HD data: definition (1/2)

### An attempt in the non-parametric case

Dataset  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ ,  $\mathbf{x}_j$  described by  $d$  variables, where  $n = o(e^d)$

Justifications:

- To approximate within error  $\epsilon$  a (Lipschitz) function of  $d$  variables, about  $(1/\epsilon)^d$  evaluations on a grid are required [Bellman, 61]
- Approximate a Gaussian distribution with fixed Gaussian kernels and with approximate error of about 10% [Silverman, 86]

$$\log_{10} n(d) \approx 0.6(d - 0.25)$$

For instance,  $n(10) \approx 7.10^5$

## HD data: definition (2/2)

### An attempt in the parametric case

Dataset  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ ,  $\mathbf{x}_j$  described by  $d$  variables and a model  $\mathbf{m}$  with  $\nu$  parameters, where  $n = o(g(\nu))$ , with  $g$  a given function

Justification:

- We consider the heteroscedastic Gaussian mixture with of true parameter  $\theta^*$  with  $K^*$  components. We note  $\hat{\theta}$  the Gaussian MLE with  $K^*$  components. We have  $g$  linear from the following result [Michel, 08]: it exists constants  $\kappa$ ,  $A$  and  $C$  such that

$$E_{\mathbf{x}}[\text{Hellinger}^2(p_{\theta^*}, p_{\hat{\theta}_{\hat{K}}})] \leq C \left[ \kappa \frac{\nu}{n} \left\{ 2A \ln d + 1 - \ln \left( 1 \wedge \left[ \frac{\nu}{n} A \ln d \right] \right) \right\} + \frac{1}{n} \right].$$

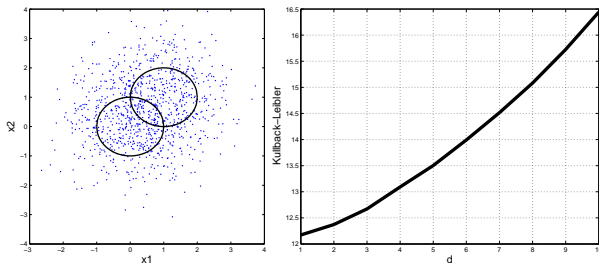
But  $\nu$  can be high since  $\nu \sim d^2/2$ , combined with potentially large constants.

## HD density estimation: curse

A two-component  $d$ -variate Gaussian mixture:

$$\pi_1 = \pi_2 = \frac{1}{2}, \quad \mathbf{X}_1|z_{11} = 1 \sim N_d(\mathbf{0}, \mathbf{I}), \quad \mathbf{X}_1|z_{12} = 1 \sim N_d(\mathbf{1}, \mathbf{I})$$

Components are **more and more separated** when  $d$  grows:  $\|\mu_2 - \mu_1\|_1 = \sqrt{d} \dots$



... but **density estimation quality decreases** with  $d$



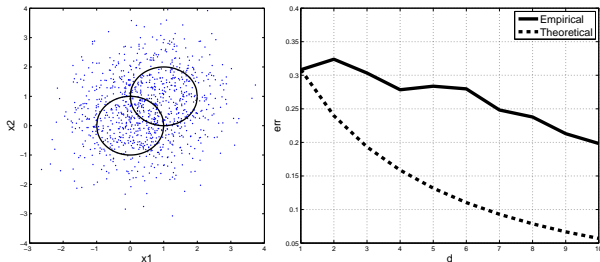
## HD clustering: blessing (1/2)

A two-component  $d$ -variate Gaussian mixture:

$$\pi_1 = \pi_2 = \frac{1}{2}, \quad \mathbf{X}_1|z_{11} = 1 \sim N_d(\mathbf{0}, \mathbf{I}), \quad \mathbf{X}_1|z_{12} = 1 \sim N_d(\mathbf{1}, \mathbf{I})$$

Each variable provides **equal** and **own** separation information

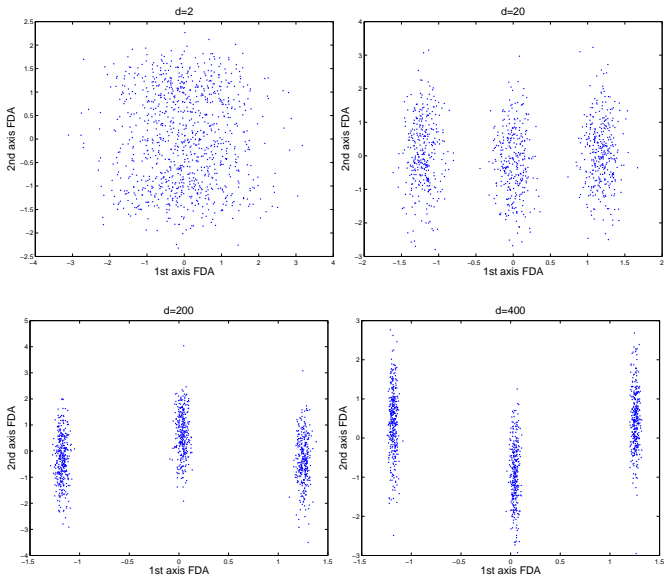
Theoretical error decreases when  $d$  grows:  $err_{theo} = \Phi(-\sqrt{d}/2) \dots$



... and empirical error rate decreases also with  $d$ !

## HD clustering: blessing (2/2)

FDA



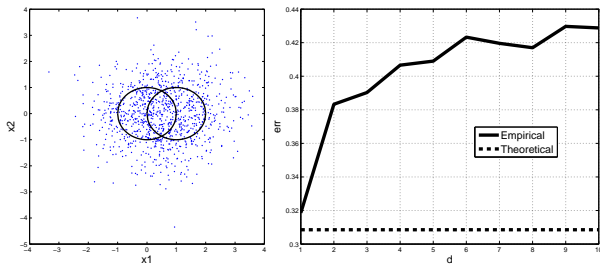
## HD clustering: curse (1/2)

Many variables provide **no separation information**

Same parameter setting except:

$$\mathbf{X}_1 | z_{12} = 1 \sim N_d((1 \ 0 \ \dots \ 0)', \mathbf{I})$$

Groups are **not separated more** when  $d$  grows:  $\|\mu_2 - \mu_1\|_1 = 1 \dots$



... thus **theoretical error is constant** ( $= \Phi(-\frac{1}{2})$ ) and **empirical error increases** with  $d$

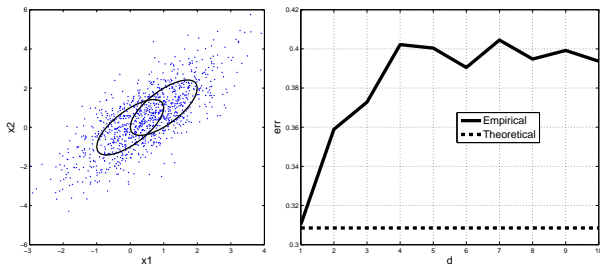
## HD clustering: curse (2/2)

Many variables provide **redundant separation information**

Same parameter setting except:

$$\mathbf{x}_1^j = \mathbf{x}_1^1 + N_1(0, 1) \quad (j = 2, \dots, d)$$

Groups are **not separated more** when  $d$  grows:  $\|\mu_2 - \mu_1\|_{\Sigma} = 1 \dots$



... thus  $err_{theo}$  is constant ( $= \Phi(-\frac{1}{2})$ ) and empirical error increases (less) with  $d$

## The trade-off bias/variance

### The fundamental statistical principle

Always minimize an error  $\text{err}$  between truth ( $\mathbf{z}$ ) and estimate ( $\hat{\mathbf{z}}$ )

- Gap between true ( $\mathbf{z}$ ) and model-based ( $\mathcal{Z}_p$ ) partitions:  $\mathbf{z}^* = \arg \min_{\tilde{\mathbf{z}} \in \mathcal{Z}_p} \Delta(\mathbf{z}, \tilde{\mathbf{z}})$
- Estimation  $\hat{\mathbf{z}}$  of  $\mathbf{z}^*$  in  $\mathcal{Z}_p$ : any relevant method (bias, consistency, efficiency. . .)
- Fundamental decomposition of the observed error  $\text{err}(\mathbf{z}, \hat{\mathbf{z}})$ :

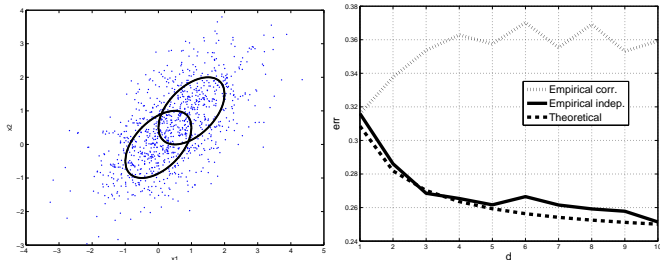
$$\begin{aligned}\text{err}(\mathbf{z}, \hat{\mathbf{z}}) &= \left\{ \text{err}(\mathbf{z}, \mathbf{z}^*) - \text{err}(\mathbf{z}, \mathbf{z}) \right\} + \left\{ \text{err}(\mathbf{z}, \hat{\mathbf{z}}) - \text{err}(\mathbf{z}, \mathbf{z}^*) \right\} \\ &= \left\{ \text{bias} \right\} + \left\{ \text{variance} \right\} \\ &= \left\{ \text{error of approximation} \right\} + \left\{ \text{error of estimation} \right\}\end{aligned}$$

## Bias/variance in HD: reduce variance, accept bias

A two-component  $d$ -variate Gaussian mixture with **intra-dependency**:

$$\pi_1 = \pi_2 = \frac{1}{2}, \quad \mathbf{X}_1 | z_{11} = 1 \sim N_d(\mathbf{0}, \Sigma), \quad \mathbf{X}_1 | z_{12} = 1 \sim N_d(\mathbf{1}, \Sigma)$$

- Each variable provides **equal** and **own** separation information
- Theoretical error decreases** when  $d$  grows:  $\text{err}_{\text{theo}} = \Phi(-\|\mu_2 - \mu_1\|_{\Sigma^{-1}}/2)$
- Empirical error rate with the (true) **intra-correlated model worse** with  $d$
- Empirical error rate with the (false) **intra-independent model better** with  $d$ !



## Some alternatives for reducing variance

- Dimension reduction in non-canonical space (PCA-like typically)
- Dimension reduction in the canonical space (variable selection)
- Model parsimony in the initial HD space (constraints on model parameters)

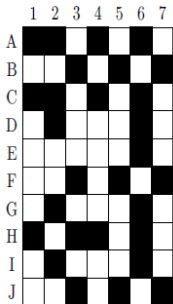
### But which kind of parsimony?

- Remember that clustering is a way for dealing with large  $n$
- Why not reusing this idea for large  $d$ ?

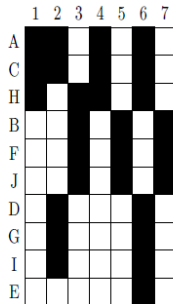
### Co-clustering

It performs parsimony of row clustering through variable clustering

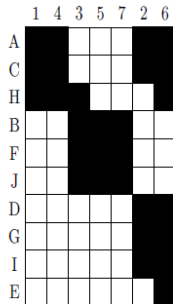
## From clustering to co-clustering



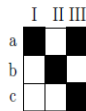
(1)



(2)



(3)



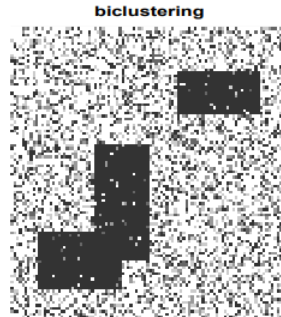
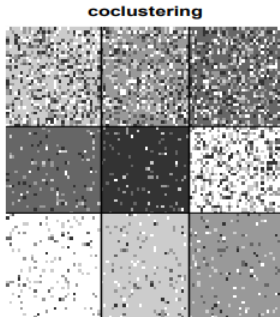
(4)

[Govaert, 2011]



# Bi-clustering

- A generalization of co-clustering
- Look for submatrices of  $x$  which are homogeneous
- We do not consider bi-clustering here



# Outline

1 HD clustering

**2 Modeling**

3 Estimating

4 Selecting

5 BlockCluster in MASSICCC

6 To go further

## Notations

- $\mathbf{z}_i$ : the cluster of the row  $i$
- $\mathbf{w}_j$ : the cluster of the column  $j$
- $(\mathbf{z}_i, \mathbf{w}_j)$ : the **block** of the element  $x_{ij}$  (row  $i$ , column  $j$ )
  
- $\mathbf{z} = (\mathbf{z}_1, \dots, \mathbf{z}_n)$ : partition of individuals in  $K$  clusters of rows
- $\mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_d)$ : partition of variables in  $L$  clusters of columns
- $(\mathbf{z}, \mathbf{w})$ : **bi-partition** of the whole data set  $\mathbf{x}$
- Both space partitions are respectively denoted by  $\mathcal{Z}$  and  $\mathcal{W}$

### Restriction

All variables are of the same kind (see discussion at the end)

## A geometric approach

- Example in the continuous case:  $\mathbf{x} \in \mathbb{R}^{n \times d}$
- It could be possible to define a **within-block** inertia criterion

$$W(\mathbf{z}, \mathbf{w}) = \underbrace{\sum_{k=1}^K \sum_{l=1}^L \sum_{i=1}^n \sum_{j=1}^d}_{\sum_{i,j,k,l}} z_{ik} w_{jl} \|x_{ij} - \mu_{kl}\|^2$$

with  $\mu_{kl}$  the **center of the block**  $(k, l)$

$$\mu_{kl} = \frac{1}{n_{kl}} \sum_{i,j} z_{ik} w_{jl} x_{ij}$$

where  $n_{kl} = \sum_{i,j} z_{ik} w_{jl}$  is the **sample size of the block**  $(k, l)$

But we know now that it hides some model-based assumptions. . .

## The latent block model (LBM)

- Generalization of some existing non-probabilistic methods
- Extend the latent class principle of local (or conditional) independence
- Thus  $x_{ij}$  is assumed to be independent once  $z_i$  and  $w_j$  are fixed ( $\alpha = (\alpha_{kl})$ ):

$$p(\mathbf{x}|\mathbf{z}, \mathbf{w}; \alpha) = \prod_{i,j} p(x_{ij}; \alpha_{z_i w_j})$$

- $\pi = (\pi_k)$  : vectors of proba.  $\pi_k$  that a row belongs to the  $k$ th row cluster
- $\rho = (\rho_l)$  : vectors of proba.  $\rho_l$  that a row belongs to the  $l$ th column cluster
- Independence between all  $z_i$  and  $w_j$
- Extension of the traditional mixture model-based clustering ( $\alpha = (\alpha_{kl})$ ):

$$p(\mathbf{x}; \theta) = \sum_{(\mathbf{z}, \mathbf{w}) \in \mathcal{Z} \times \mathcal{W}} \prod_{i,j} \pi_{z_i} \rho_{w_j} p(x_{ij}; \alpha_{z_i w_j})$$

## Distribution for different kinds of data

[Govaert and Nadif, 2014] The pdf  $p(\cdot; \alpha_{z_i w_j})$  depends on the kind of data  $x_{ij}$ :

- **Binary** data:  $x_{ij} \in \{0, 1\}$ ,  $p(\cdot; \alpha_{kl}) = \mathcal{B}(\alpha_{kl})$
- **Categorical** data with  $m$  levels:  
 $\mathbf{x}_{ij} = \{x_{ijh}\} \in \{0, 1\}^m$  with  $\sum_{h=1}^m x_{ijh} = 1$  and  $p(\cdot; \alpha_{kl}) = \mathcal{M}(\alpha_{kl})$  with  $\alpha_{kl} = \{\alpha_{kjh}\}$
- **Count** data:  $x_i^j \in \mathbb{N}$ ,  $p(\cdot; \alpha_{kl}) = \mathcal{P}(\mu_k \nu_l \gamma_{kl})^2$
- **Continuous** data:  $x_i^j \in \mathbb{R}$ ,  $p(\cdot; \alpha_{kl}) = \mathcal{N}(\mu_{kl}, \sigma_{kl}^2)$

---

<sup>2</sup>The Poisson parameter is here split into  $\mu_k$  and  $\nu_l$  the effects of the row  $k$  and the column  $l$  respectively and  $\gamma_{kl}$  the effect of the block  $kl$ . Unfortunately, this parameterization is not identifiable. It is therefore not possible to estimate simultaneously  $\mu_k$ ,  $\nu_l$  and  $\gamma_{kl}$  without imposing further constraints. Constraints  $\sum_k \pi_k \gamma_{kl} = \sum_l \rho_l \gamma_{kl} = 1$  and  $\sum_k \mu_k = 1, \sum_l \nu_l = 1$  are a possibility.

## Extreme parsimony ability

Model	Number of parameters
Binary	$\dim(\boldsymbol{\pi}) + \dim(\boldsymbol{\rho}) + KL$
Categorical	$\dim(\boldsymbol{\pi}) + \dim(\boldsymbol{\rho}) + KL(m - 1)$
Contingency	$\dim(\boldsymbol{\pi}) + \dim(\boldsymbol{\rho}) + KL$
Continuous	$\dim(\boldsymbol{\pi}) + \dim(\boldsymbol{\rho}) + 2KL$

Very parsimonious so well suitable for the (ultra) HD setting

$$\text{nb. param.}_{\text{HD}} = \text{nb. param.}_{\text{classic}} \times \frac{L}{d}$$

**Other advantage:** stay in the canonical space thus meaningful for the end-user

## Binary illustration: easy interpretation

[Govaert, 2011]

	<i>abcdefghij</i>
y1	1010001101
y2	0101110011
y3	1000001100
y4	1010001100
y5	0111001100
y6	0101110101
y7	0111110111
y8	1100110111
y9	0100110000
y10	1010101101
y11	1010001100
y12	1010000100
y13	1010001101
y14	0010011100
y15	0010010100
y16	1111001100
y17	0101110011
y18	1010011101
y19	1010001000
y20	1100101100

Raw data

	<i>a c g h</i>	<i>b d e f i j</i>
y2	0 0 0 0	1 1 1 1 1 1
y6	0 0 0 1	1 1 1 1 0 1
y7	0 1 0 1	1 1 1 1 1 1
y8	1 0 1 0	1 0 1 1 1 1
y9	0 0 0 0	1 0 1 1 0 0
y17	0 0 0 0	1 1 1 1 1 1
y1	1 1 1 1	0 0 0 0 0 1
y3	1 0 1 1	0 0 0 0 0 0
y4	1 1 1 1	0 0 0 0 0 0
y5	0 1 1 1	1 1 0 0 0 0
y10	1 1 1 1	0 0 1 0 0 1
y11	1 1 1 1	0 0 0 0 0 0
y12	1 1 0 1	0 0 0 0 0 0
y13	1 1 1 1	0 0 0 0 0 1
y14	0 1 1 1	0 0 0 1 0 0
y15	0 1 0 1	0 0 0 1 0 0
y16	1 1 1 1	1 1 0 0 0 0
y18	1 1 1 1	0 0 0 1 0 1
y19	1 1 1 0	0 0 0 0 0 0
y20	1 0 1 1	1 0 1 0 0 0

Permuted data  
(rows/columns)

mode

0	1
1	0

Summary

0.86	0.79
0.83	0.86

Homogeneity

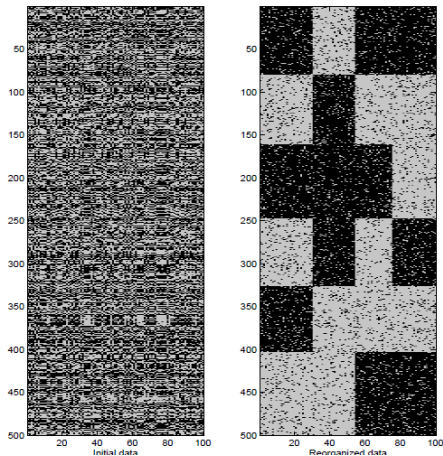
proba=mode

iid Bin(0.83)



## Binary illustration: user-friendly visualization

[Govaert, 2011]



$$n = 500, d = 10, K = 6, L = 4$$

## Other kind of data: ordinal

[Jacques and Biernacki, 2018]

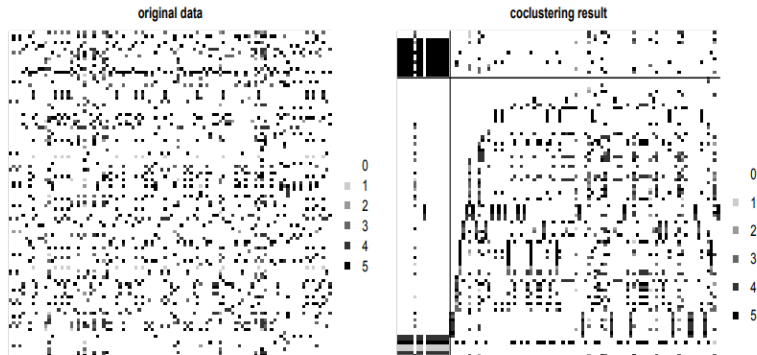
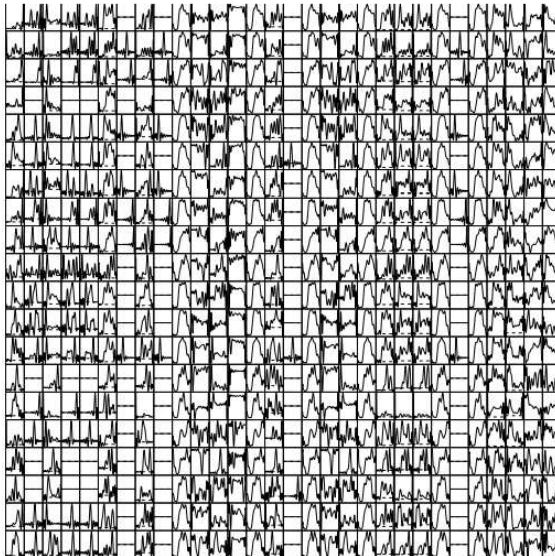


Figure 11: Top 100 Amazon Fine Food Review data (left) and co-clustering result (right).

## Other kind of data: functional

[Jacques, 2016]

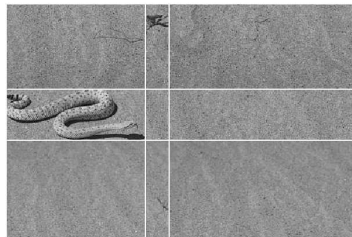


## Other kind of data: image

Original Data



Co-Clustered Data



## Particular case: graph clustering

Stochastic Block Model (SBM): adjacency matrix with  $n = d$  and  $K = L$



# Outline

1 HD clustering

2 Modeling

**3 Estimating**

4 Selecting

5 BlockCluster in MASSICCC

6 To go further

## MLE estimation: log-likelihood(s)

- Remember Lesson 3: first estimate  $\theta$ , then deduce estimate of  $(z, \mathbf{w})$
- Observed log-likelihood:  $\ell(\theta; \mathbf{x}) = \ln p(\mathbf{x}; \theta)$
- MLE:

$$\hat{\theta} = \arg \max_{\theta} \ell(\theta; \mathbf{x})$$

- Complete log-likelihood:

$$\begin{aligned} \ell_c(\theta; \mathbf{x}, \mathbf{z}, \mathbf{w}) &= \ln p(\mathbf{x}, \mathbf{z}, \mathbf{w}; \theta) \\ &= \sum_{i,k} z_{ik} \log \pi_k + \sum_{k,l} w_{jl} \log \rho_l + \sum_{i,j,k,l} z_{ik} w_{jl} \log p(x_i^j; \alpha_{kl}) \end{aligned}$$

Be careful with asymptotics...

If  $\ln(d)/n \rightarrow 0$ ,  $\ln(n)/d \rightarrow 0$  when  $n \rightarrow \infty$  and  $d \rightarrow \infty$ , then the MLE is consistent

[Brault et al., 2017]

## MLE estimation: EM algorithm

- **E-step** of EM (iteration  $q$ ):

$$\begin{aligned}
 Q(\theta, \theta^{(q)}) &= E[\ell_c(\theta; \mathbf{x}, \mathbf{z}, \mathbf{w}) | \mathbf{x}; \theta^{(q)}] \\
 &= \sum_{i,k} \underbrace{p(z_i = k | \mathbf{x}; \theta^{(q)})}_{t_{ik}^{(q)}} \ln \pi_k + \sum_{j,l} \underbrace{p(w_j = l | \mathbf{x}; \theta^{(q)})}_{s_{jl}^{(q)}} \ln \rho_l \\
 &\quad + \sum_{i,j,k,l} \underbrace{p(z_i = k, w_j = l | \mathbf{x}; \theta^{(q)})}_{e_{ijkl}^{(q)}} \ln p(x_{ij}; \alpha_{kl})
 \end{aligned}$$

- **M-step** of EM (iteration  $q$ ): classical. For instance, for the Bernoulli case, it gives

$$\pi_k^{(q+1)} = \frac{\sum_i t_{ik}^{(q)}}{n}, \quad \rho_l^{(q+1)} = \frac{\sum_j s_{jl}^{(q)}}{d}, \quad \alpha_{kl}^{(q+1)} = \frac{\sum_{i,j} e_{ijkl}^{(q)} x_{ij}}{\sum_{i,j} e_{ijkl}^{(q)}}$$



## MLE: intractable E step

$e_{ijkl}^{(q)}$  is usually intractable. . .

- Consequence of dependency between  $x_{ij}$ s (link between rows and columns)
- Involve  $K^n L^d$  calculus (number of possible blocks)
- Example: if  $n = d = 20$  and  $K = L = 2$  then  $10^{12}$  blocks
- Example (cont'd): 33 years with a computer calculating 100,000 blocks/second

### Alternatives to EM

- **Variational EM** (numerical approx.): conditional independence assumption

$$p(\mathbf{z}, \mathbf{w} | \mathbf{x}; \boldsymbol{\theta}) \approx p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}) p(\mathbf{w} | \mathbf{x}; \boldsymbol{\theta})$$

- **SEM-Gibbs** (stochastic approx.): replace E-step by a S-step approx. by Gibbs

$$\mathbf{z} | \mathbf{x}, \mathbf{w}; \boldsymbol{\theta} \quad \text{and} \quad \mathbf{w} | \mathbf{x}, \mathbf{z}; \boldsymbol{\theta}$$

## MLE: variational EM (1/2)

- Use a general variational result from [Hathaway, 1985]
- Maximizing  $\ell(\boldsymbol{\theta}; \mathbf{x})$  on  $\boldsymbol{\theta}$  is equivalent to maximize  $\tilde{\ell}_c(\boldsymbol{\theta}; \mathbf{x}, \mathbf{e})$  on  $(\boldsymbol{\theta}, \mathbf{e})$

$$\tilde{\ell}_c(\boldsymbol{\theta}; \mathbf{x}, \mathbf{e}) = \sum_{i,k} t_{ik} \ln \pi_k + \sum_{j,l} s_{jl} \ln \rho_l + \sum_{i,j,k,l} e_{ijkl} \ln p(x_{ij}; \boldsymbol{\alpha}_{kl})$$

where  $\mathbf{e} = (e_{ijkl})$ ,  $e_{ijkl} \in \{0, 1\}$ ,  $\sum_{k,l} e_{ijkl} = 1$ ,  $t_{ik} = \sum_{j,l} e_{ijkl}$ ,  $s_{jl} = \sum_{i,k} e_{ijkl}$

- Of course maximizing  $\ell(\boldsymbol{\theta}; \mathbf{x})$  or  $\tilde{\ell}_c(\boldsymbol{\theta}; \mathbf{x}, \mathbf{e})$  are both intractable
- Idea: restriction on  $\mathbf{e}$  to obtain tractability  $e_{ijkl} = t_{ik}s_{jl}$
- New variables are thus now  $\mathbf{t} = (t_{ik})$  and  $\mathbf{s} = (s_{jl})$
- As a consequence, it is a maximization of a lower bound of the max. likelihood

$$\max_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathbf{x}) \geq \max_{\boldsymbol{\theta}, \mathbf{t}, \mathbf{s}} \tilde{\ell}_c(\boldsymbol{\theta}; \mathbf{x}, \mathbf{e})$$

## MLE: variational EM (2/2)

### Approximated E-step

$$Q(\theta, \theta^{(q)}) \approx \sum_{i,k} t_{ik}^{(q)} \ln \pi_k + \sum_{j,l} s_{jl}^{(q)} \ln \rho_l + \sum_{i,j,k,l} t_{ik}^{(q)} s_{jl}^{(q)} \ln p(x_{ij}; \alpha_{kl})$$

- We called it now VEM
- Also known as **mean field** approximation
- **Consistency** of the variational estimate [Brault *et al.*, 2017]

## MLE: local maxima

- More local maxima than in classical mixture models
- It is a consequence of many more latent variables (blocks)
- Thus: either many VEM runs, or use the SEM-Gibbs algorithm

## MLE: SEM-Gibbs

- We have already seen the SEM algorithm in Lesson 3 (thus we do not detail more)
- It limits dependency to starting point, so it limits local maxima
- The S-step: a draw  $(\mathbf{z}^{(q)}, \mathbf{w}^{(q)}) \sim p(\mathbf{z}, \mathbf{w} | \mathbf{x}; \boldsymbol{\theta}^{(q)})$  instead an expectation
- But it is still intractable, thus use a Gibbs algorithm to approx. this draw

### Approximated S-step

Two easy draws

$$\mathbf{z}^{(q)} \sim p(\mathbf{z} | \mathbf{w}^{(q-1)}, \mathbf{x}; \boldsymbol{\theta}^{(q)})$$

and

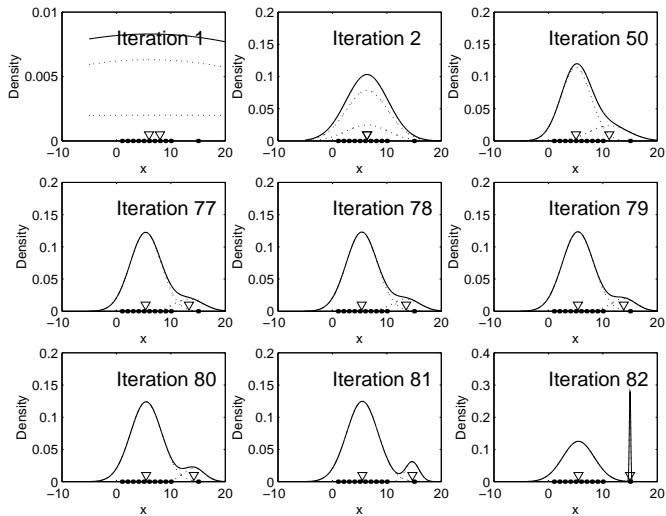
$$\mathbf{w}^{(q)} \sim p(\mathbf{w} | \mathbf{z}^{(q)}, \mathbf{x}; \boldsymbol{\theta}^{(q)})$$

- Rigorously speaking, many draws within the S-step should be performed
- Indeed, Gibbs has to reach a stochastic convergence
- In practice it works well while saving computation time

## MLE: degeneracy

- More degenerate situations than in classical mixture models
- It is again a consequence of many more latent variables (blocks)
- The Bayesian regularization (instead MLE) can be an answer

## Illustration of a degenerate situation



## Bayesian estimation: pitch

- Everything passes by the **posterior distribution of  $\theta$**

$$p(\theta|\mathbf{x}) \propto \underbrace{p(\mathbf{x}|\theta)}_{\text{log-likelihood}} \underbrace{p(\theta)}_{\text{prior}}$$

- Then, take (for instance) the **MAP** as a  $\theta$  estimate (use a VEM like algo...)

$$\hat{\theta} = \arg \max_{\theta} p(\theta|\mathbf{x})$$



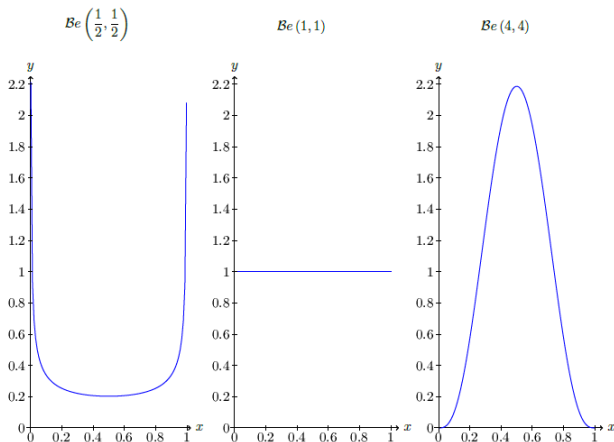
## Bayesian estimation: limiting degeneracy

- Interest for avoiding degeneracy is the prior: it acts as a **penalization** term
- Typical choices are **Dirichlet** for  $\pi$  and  $\rho$  (with independence between  $\pi$ ,  $\rho$ ,  $\alpha$ )

$$p(\theta) = \underbrace{p(\pi)}_{D_K(a, \dots, a)} \times \underbrace{p(\rho)}_{D_L(a, \dots, a)} \times \underbrace{p(\alpha)}_{\text{model dependent}}$$

- The Dirichlet distribution is conjugate, thus easy calculus
- **Control degeneracy frequency with the  $a$  value:**
  - $a = 1$ : uniform prior, so  $\hat{\theta}$  is strictly the MLE (no regularisation)
  - $a = 1/2$ : Jeffreys prior, classical (no informative prior) but may favor degeneracy
  - $a = 4$ : a rule of thumb working well for limiting degeneracy frequency

## Bayesian estimation: prior overview



## Block estimation: estimate

- Once we have a parameter estimate  $\hat{\theta}$ , we need to have an block estimate  $(\hat{z}, \hat{w})$
- But MAP not directly available because of the following maximization difficulty

$$(\hat{z}, \hat{w}) = \arg \max_{(z, w)} \underbrace{p(z, w | x; \hat{\theta})}_{\text{intractable}}$$

- Instead the following (easily, as classical mixtures) estimates are usually retained

$$\hat{z} = \arg \max_z p(z | x; \hat{\theta}) \quad \text{and} \quad \hat{w} = \arg \max_w p(w | x; \hat{\theta})$$

## Block estimation: evaluation

- **Empirical error rate** between blocks:

$$\text{err}_{\text{blocks}}\left(\underbrace{(\mathbf{z}, \mathbf{w})}_{\text{"True" blocks}}, \underbrace{(\hat{\mathbf{z}}, \hat{\mathbf{w}})}_{\text{Estimated blocks}}\right) = \text{err}(\mathbf{z}, \hat{\mathbf{z}}) + \text{err}(\mathbf{w}, \hat{\mathbf{w}}) - \text{err}(\mathbf{z}, \hat{\mathbf{z}}) \times \text{err}(\mathbf{w}, \hat{\mathbf{w}})$$

- **Rand index** between blocks: it exists also a recent definition. . .

## Block estimation: consistency

[Mariadassou and Matias, 12]

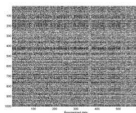
$$\underbrace{\hat{\theta} \xrightarrow{n, d \rightarrow \infty} \theta^*}_{\text{we have seen that...}} \Rightarrow \underbrace{p(\hat{\mathbf{z}} = \mathbf{z}^*, \hat{\mathbf{w}} = \mathbf{w}^* | \mathbf{x}; \hat{\theta}) \xrightarrow{n, d \rightarrow \infty} 1}_{\text{exact bi-partition retrieval!}}$$

Thus we retrieve the HD clustering blessing...

## Block estimation: non asymptotic properties (1/2)

Binary case: marginals seems so **simple mixtures!** [Brault, 14]

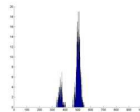
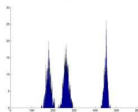
Matrice initiale



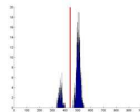
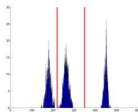
Lignes

Colonnes

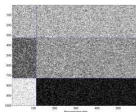
Histogrammes des sommes



Séparations



Matrice réorganisée



## Block estimation: non asymptotic properties (2/2)

[Brault, 14]

- Probability of  $x_{ij}$  with no regard to the column membership is Bernoulli

$$p(x_{ij} = 1 | z_{ik} = 1) = \tau_k = \sum_{l=1}^L \alpha_{kl} \rho_l$$

- Thus marginal distribution of  $x_{ij}$  is a mixture (indep. of  $x_{ij}$  cond.  $z_{ik} = 1$ )

$$\left( \sum_j x_{ij} \right) | z_{ik} = 1 \sim B(d, \tau_k)$$

- Control of error on this partition mixture estimate  $\hat{\mathbf{z}}^{mix}$  of binomial distributions

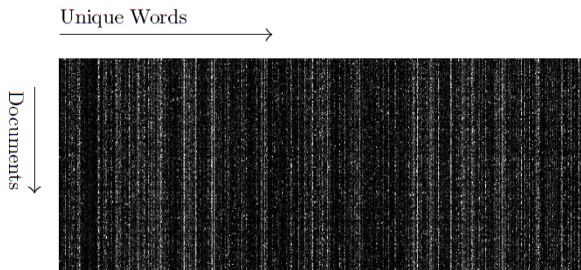
$$p(\hat{\mathbf{z}}^{mix} \neq \mathbf{z}^*) \leq 2n \exp \left\{ - \frac{1}{8} d \underbrace{\left[ \min_{k \neq k'} |\tau_k - \tau_{k'}| \right]}_{\text{overlap}} \right\} + K(1 - \min_k \pi_k)^n$$

- We retrieve also consistency for very high dimension with constraint

$$\ln(n) = o(d)$$

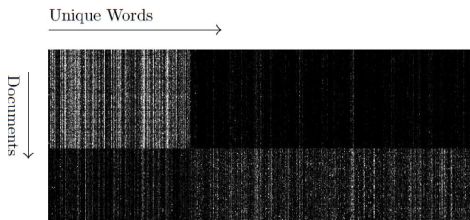
## Illustration: document clustering (1/2)

- Mixture of 1033 medical summaries and 1398 aeronautics summaries
- **Lines:** 2431 documents
- **Columns:** present words (except stop), thus 9275 unique words
- Data matrix: cross counting document  $\times$  words
- Poisson model





## Illustration: document clustering (2/2)



### Results with $2 \times 2$ blocks

	Medline	Cranfield
Medline	1033	0
Cranfield	0	1398

Experiment illustrates previous theory: HD clustering is blessing

# Outline

- 1 HD clustering
- 2 Modeling
- 3 Estimating
- 4 Selecting**
- 5 BlockCluster in MASSICCC
- 6 To go further

## Models in competition

$\mathbf{m} = (K, L)$  typically, but not restricted to

## BIC criterion: two difficulties

- **Difficult 1:** which BIC definition because of the double asymptotic on  $n$  and  $d$ ?
- **Difficult 2:** the observed log-likelihood value is intractable

$$\ell(\boldsymbol{\theta}; \mathbf{x}) = \sum_{(\mathbf{z}, \mathbf{w}) \in \mathcal{Z} \times \mathcal{W}} p(\mathbf{x}, \mathbf{z}, \mathbf{w}; \boldsymbol{\theta})$$

Could be estimated by harmonic mean but time consuming and high variance

## ICL criterion: overcome both difficulties

- ICL uses complete likelihood thus no intractability

$$\text{ICL} = \ln p(\mathbf{x}, \hat{\mathbf{z}}, \hat{\mathbf{w}}) = \ln p(\mathbf{x}|\hat{\mathbf{z}}, \hat{\mathbf{w}}) + \ln p(\hat{\mathbf{z}}) + \ln p(\hat{\mathbf{w}})$$

- Multinomial case ( $r$  levels): [Keribin *et al.*, 2014]

- Derive an exact (non-asymptotic) ICL version
- Deduce an asymptotic approximation of ICL

$$\text{ICLbic} = \ell_c(\hat{\boldsymbol{\theta}}; \mathbf{x}, \hat{\mathbf{z}}, \hat{\mathbf{w}}) - \frac{K-1}{2} \ln(n) - \frac{L-1}{2} \ln(d) - \frac{KL(r-1)}{2} \ln(nd)$$

- We can make a conjecture for the general case

$$\text{ICLbic} = \ell_c(\hat{\boldsymbol{\theta}}; \mathbf{x}, \hat{\mathbf{z}}, \hat{\mathbf{w}}) - \frac{K-1}{2} \ln(n) - \frac{L-1}{2} \ln(d) - \frac{KL\nu_{\alpha_{kl}}}{2} \ln(nd)$$

## ICL criterion: consistency

- We can obtain a BIC expression from ICLbic

$$\begin{aligned} \text{BIC} &= \text{ICLbic} - \ln p(\hat{\mathbf{z}}, \hat{\mathbf{w}} | \mathbf{x}; \hat{\boldsymbol{\theta}}) \\ &= \underbrace{\ell(\hat{\boldsymbol{\theta}}; \mathbf{x})}_{\text{difficult}} - \frac{K-1}{2} \ln(n) - \frac{L-1}{2} \ln(d) - \frac{KL(m-1)}{2} \ln(nd) \end{aligned}$$

- [Brault *et al.*, 2017] establish that asymptotically on  $n$  and  $d$

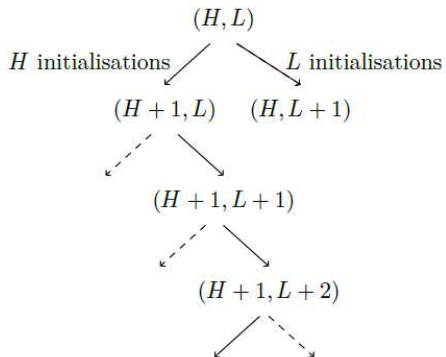
$$“\ell(\hat{\boldsymbol{\theta}}; \mathbf{x}) = \ell_c(\hat{\boldsymbol{\theta}}; \mathbf{x}, \hat{\mathbf{z}}, \hat{\mathbf{w}})”$$

- Thus, since BIC is consistent, ICL is also consistent

Again the HD clustering blessing is here!

## Strategy to smart browsing of $(K, L)$

[Robert, 2017] Algorithm Bi-KM1



## Illustration: discuss the dimension (1/2)

- SPAM E-mail Database<sup>3</sup>
- $n = 4601$  e-mails composed by 1813 “spams” and 2788 “good e-mails”
- $d = 48 + 6 = 54$  continuous descriptors<sup>4</sup>
  - 48 percentages that a given **word** appears in an e-mail (“make”, “you’... )
  - 6 percentages that a given **char** appears in an e-mail (“;”, “\$”... )
- Transformation of continuous descriptors into **binary descriptors**

$$x_{ij} = \begin{cases} 1 & \text{if word/char } j \text{ appears in e-mail } i \\ 0 & \text{otherwise} \end{cases}$$

---

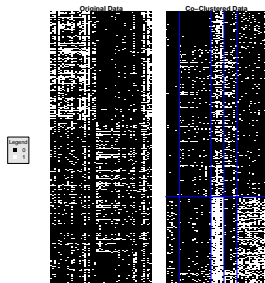
<sup>3</sup><https://archive.ics.uci.edu/ml/machine-learning-databases/spambase/>

<sup>4</sup>There are 3 other continuous descriptors we do not use



## Illustration: discuss the dimension (2/2)

- Perform **co-clustering** with  $K = 2$  and  $L = 5$ : ICL=-92,682, err=0.1984



- Perform **clustering**<sup>5</sup> with  $K = 2$ : ICL=-89,433, err=0.1837

Thus use preferably co-clustering in the HD setting

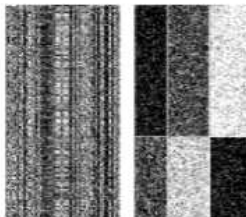
<sup>5</sup>Equivalent to co-clustering with  $L = 54$

# Outline

- 1 HD clustering
- 2 Modeling
- 3 Estimating
- 4 Selecting
- 5 BlockCluster in MASSICCC**
- 6 To go further

# MASSICCC platform for the BLOCKCLUSTER software

<https://massiccc.lille.inria.fr/>

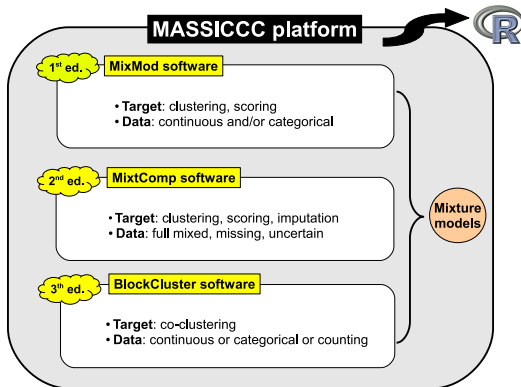


## BlockCluster

BlockCluster can estimate the parameters of co-clustering models for binary, contingency and continuous data. Simply put, when considering a set of data as rows and columns, BlockCluster will make simultaneous permutations of rows and columns in order to organise the data into homogenous blocks.

[Read more about BlockCluster](#)

## MASSICCC?



A high quality and easy to use web platform  
where are transferred mature research clustering (and more) software  
towards (non academic) professionals

Here is the computer you need!



# Running BlockCluster

## Configuration

If you change the configuration of your job and save it, it will start a new process with the updated parameters. This will erase previous results.

Parameters

Title	<input type="text" value="Trial BlockCluster"/>
Data File	<input type="text" value="Blockcluster-Example.csv"/>
Data Type	<input type="text" value="Categorical"/> ⓘ
Rows Cluster Groups	<input type="text" value="1:5"/> ⓘ
Column Cluster Groups	<input type="text" value="1:5"/> ⓘ

# Running BlockCluster

MASSICCC Dashboard Help Profile Logout









RESULTS

DATA FILES

CREATE JOB

RESULTS

Select a job execution from the list below

69		Trial BlockCluster Blockcluster-Example.csv	<div style="width: 42%;"><div style="width: 42%;"></div></div>	23 May 20:47 
68		Genes K1-12 log.cpm.txt		23 May 08:12 
67		Genes log.cpm.txt		22 May 15:38 
65		Genes K1-10 log.cpm.txt		22 May 15:27 

# Running BlockCluster

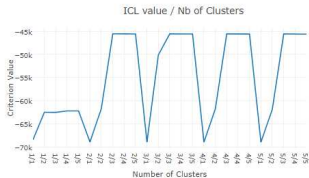
Model	Criterion	Nb Clusters	Error
<i>pik_rho_multi</i>	ICL (-45557.1)	[2,3]	No error
<i>pik_rho_multi</i>	ICL (-45563.3)	[3,3]	No error
<i>pik_rho_multi</i>	ICL (-45566.6)	[2,4]	No error
<i>pik_rho_multi</i>	ICL (-45573.9)	[4,3]	No error
<i>pik_rho_multi</i>	ICL (-45574.6)	[5,3]	No error
<i>pik_rho_multi</i>	ICL (-45577.7)	[3,4]	No error
<i>pik_rho_multi</i>	ICL (-45578.8)	[2,5]	No error

Cluster Plot

Criterion Plot

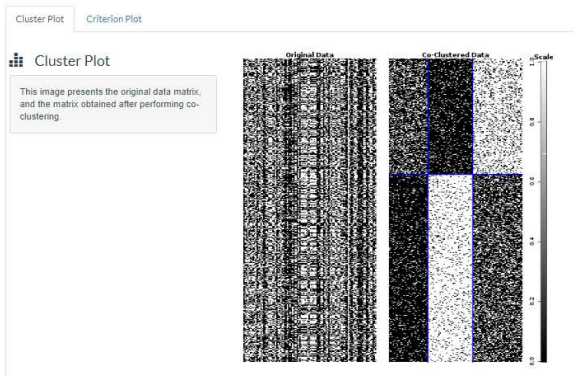
## Model Criterion

This chart represents the criterion value for each model that was built. The higher the value (close to 0) the better the model.





# Running BlockCluster



# Outline

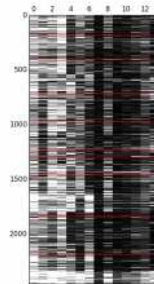
- 1 HD clustering
- 2 Modeling
- 3 Estimating
- 4 Selecting
- 5 BlockCluster in MASSICCC
- 6 To go further**

## Co-clustering of mixed data

- Same partitions in lines, disjoint partitions in columns
- Example: data set TED talks, with talks  $\times$  (terms,scores)



Poisson



Gaussian