

Variable selection for model-based clustering

Matthieu Marbac (Ensai - Crest)

Joint works with:

M. Sedki (Univ. Paris-sud) and V. Vandewalle (Univ. Lille 2)

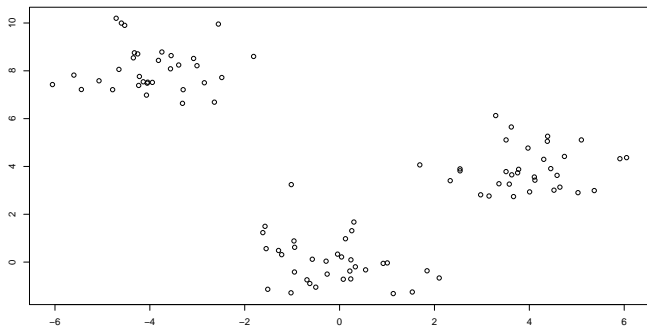
The problem

Objective:

- Estimation of a partition \mathbf{z} among n observations \mathbf{x} .

Notations:

- g : number of clusters.
- $\mathbf{x} = (\mathbf{x}_i; i = 1, \dots, n)$: observed sample of n iid observations.
- $\mathbf{x}_i = (x_{i1}, \dots, x_{id})$: vector of the d features for observation i .
- $\mathbf{z} = (\mathbf{z}_i; i = 1, \dots, n)$: unobserved partition.
- $\mathbf{z}_i = (z_{i1}, \dots, z_{ig})$ with $z_{ik} = 1$ means that observation i belongs to cluster k .



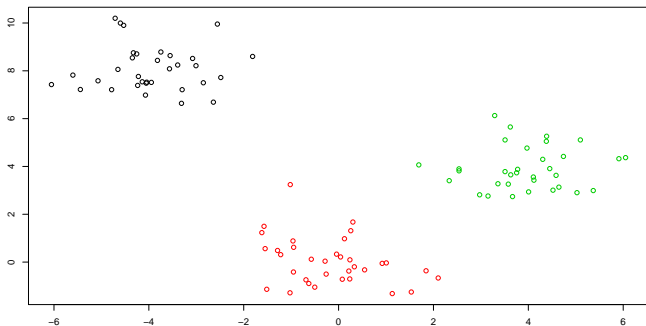
The problem

Objective:

- Estimation of a partition \mathbf{z} among n observations \mathbf{x} .

Notations:

- g : number of clusters.
- $\mathbf{x} = (\mathbf{x}_i; i = 1, \dots, n)$: observed sample of n iid observations.
- $\mathbf{x}_i = (x_{i1}, \dots, x_{id})$: vector of the d features for observation i .
- $\mathbf{z} = (\mathbf{z}_i; i = 1, \dots, n)$: unobserved partition.
- $\mathbf{z}_i = (z_{i1}, \dots, z_{ig})$ with $z_{ik} = 1$ means that observation i belongs to cluster k .



The problem

The objectives

- Estimation of a classification rule (*i.e.*, estimator of z_i given x_i).
- Evaluation of the risk of misclassification.
- Interpretation of the clusters.
- Estimation of the number of clusters g .
- Deal with complex data (mixed-type data, missing values, ...).

1 Model-based clustering

2 Variable selection

3 Multiple partitions

Model-based clustering

Main idea:

model the distribution of the observed data \mathbf{X} .

Generative model:

- $\mathbf{Z}_i \sim \mathcal{M}(\pi_1, \dots, \pi_g)$
- $\mathbf{X}_i | Z_{ik} = 1 \sim \mathcal{L}_k(\cdot)$, e.g., $\mathcal{L}_k(\cdot) = \mathcal{N}(\mu_k, \Sigma_k)$.

Mixture model:

The pdf of the observed data is

$$f(\mathbf{x}_i; \boldsymbol{\theta}) = \sum_{k=1}^g \pi_k f_k(\mathbf{x}_i; \boldsymbol{\theta}_k)$$

where $\boldsymbol{\theta}$ groups all the parameters.

Fuzzy and hard partition:

$$\mathbb{P}(Z_{ik} = 1 | \mathbf{X}_i = \mathbf{x}_i) = \frac{\pi_k f_k(\mathbf{x}_i; \boldsymbol{\theta}_k)}{\sum_{\ell=1}^g \pi_\ell f_\ell(\mathbf{x}_i; \boldsymbol{\theta}_\ell)}$$

The component membership of observation \mathbf{x}_i is obtained by

$$\hat{z}_{ik^*} = 1 \text{ if } k^* = \arg \max_k \mathbb{P}(Z_{ik} = 1 | \mathbf{X}_i = \mathbf{x}_i)$$

Model-based clustering

Main idea:

model the distribution of the observed data \mathbf{X} .

Generative model:

- $\mathbf{Z}_i \sim \mathcal{M}(\pi_1, \dots, \pi_g)$
- $\mathbf{X}_i | Z_{ik} = 1 \sim \mathcal{L}_k(\cdot)$, e.g., $\mathcal{L}_k(\cdot) = \mathcal{N}(\mu_k, \Sigma_k)$.

Mixture model:

The pdf of the observed data is

$$f(\mathbf{x}_i; \boldsymbol{\theta}) = \sum_{k=1}^g \pi_k f_k(\mathbf{x}_i; \boldsymbol{\theta}_k)$$

where $\boldsymbol{\theta}$ groups all the parameters.

Fuzzy and hard partition:

$$\mathbb{P}(Z_{ik} = 1 | \mathbf{X}_i = \mathbf{x}_i) = \frac{\pi_k f_k(\mathbf{x}_i; \boldsymbol{\theta}_k)}{\sum_{\ell=1}^g \pi_\ell f_\ell(\mathbf{x}_i; \boldsymbol{\theta}_\ell)}$$

The component membership of observation \mathbf{x}_i is obtained by

$$\hat{z}_{ik^*} = 1 \text{ if } k^* = \arg \max_k \mathbb{P}(Z_{ik} = 1 | \mathbf{X}_i = \mathbf{x}_i)$$

Model-based clustering

Main idea:

model the distribution of the observed data \mathbf{X} .

Generative model:

- $\mathbf{Z}_i \sim \mathcal{M}(\pi_1, \dots, \pi_g)$
- $\mathbf{X}_i | Z_{ik} = 1 \sim \mathcal{L}_k(\cdot)$, e.g., $\mathcal{L}_k(\cdot) = \mathcal{N}(\mu_k, \Sigma_k)$.

Mixture model:

The pdf of the observed data is

$$f(\mathbf{x}_i; \boldsymbol{\theta}) = \sum_{k=1}^g \pi_k f_k(\mathbf{x}_i; \boldsymbol{\theta}_k)$$

where $\boldsymbol{\theta}$ groups all the parameters.

Fuzzy and hard partition:

$$\mathbb{P}(Z_{ik} = 1 | \mathbf{X}_i = \mathbf{x}_i) = \frac{\pi_k f_k(\mathbf{x}_i; \boldsymbol{\theta}_k)}{\sum_{\ell=1}^g \pi_\ell f_\ell(\mathbf{x}_i; \boldsymbol{\theta}_\ell)}$$

The component membership of observation \mathbf{x}_i is obtained by

$$\hat{z}_{ik^*} = 1 \text{ if } k^* = \arg \max_k \mathbb{P}(Z_{ik} = 1 | \mathbf{X}_i = \mathbf{x}_i)$$

Model-based clustering: inference

Maximum likelihood inference

From the sample $\mathbf{x} = (\mathbf{x}_i; i = 1, \dots, n)$, we want $\hat{\boldsymbol{\theta}} = \operatorname{argmax} \ell(\boldsymbol{\theta}; \mathbf{x})$ where

$$\ell(\boldsymbol{\theta}; \mathbf{x}) = \ln p(\mathbf{x}; \boldsymbol{\theta}) = \sum_{i=1}^n \ln \left(\sum_{k=1}^g \pi_k f_k(\mathbf{x}_i; \boldsymbol{\theta}_k) \right).$$

We consider the complete-data log-likelihood

$$\ell(\boldsymbol{\theta}; \mathbf{x}, \mathbf{z}) = \ln p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = \sum_{i=1}^n \sum_{k=1}^g z_{ik} \ln \left(\pi_k f_k(\mathbf{x}_i; \boldsymbol{\theta}_k) \right).$$

EM algorithm

- Introduced to deal with missing values (here \mathbf{z} is missing).
- Iterative.
- Log-likelihood increases at each iteration.
- At iteration $[r]$, two steps:
 - E-step: computation of

$$t_{ik}^{[r-1]} := \mathbb{P}(Z_{ik} = 1 | \mathbf{x}_i, \boldsymbol{\theta}^{[r-1]}).$$

- M-step: $\boldsymbol{\theta}^{[r]}$ maximizes the complete-data log-likelihood

$$\ln p(\mathbf{x}, \mathbf{t}^{[r-1]} | \boldsymbol{\theta}) = \sum_{i=1}^n \sum_{k=1}^g t_{ik}^{[r-1]} \ln \left(\pi_k f_k(\mathbf{x}_i; \boldsymbol{\theta}_k) \right).$$

Model-based clustering: inference

Maximum likelihood inference

From the sample $\mathbf{x} = (\mathbf{x}_i; i = 1, \dots, n)$, we want $\hat{\boldsymbol{\theta}} = \operatorname{argmax} \ell(\boldsymbol{\theta}; \mathbf{x})$ where

$$\ell(\boldsymbol{\theta}; \mathbf{x}) = \ln p(\mathbf{x}; \boldsymbol{\theta}) = \sum_{i=1}^n \ln \left(\sum_{k=1}^g \pi_k f_k(\mathbf{x}_i; \boldsymbol{\theta}_k) \right).$$

We consider the complete-data log-likelihood

$$\ell(\boldsymbol{\theta}; \mathbf{x}, \mathbf{z}) = \ln p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = \sum_{i=1}^n \sum_{k=1}^g z_{ik} \ln \left(\pi_k f_k(\mathbf{x}_i; \boldsymbol{\theta}_k) \right).$$

EM algorithm

- Introduced to deal with missing values (here \mathbf{z} is missing).
- Iterative.
- Log-likelihood increases at each iteration.
- At iteration $[r]$, two steps:
 - E-step: computation of

$$t_{ik}^{[r-1]} := \mathbb{P}(Z_{ik} = 1 | \mathbf{x}_i, \boldsymbol{\theta}^{[r-1]}).$$

- M-step: $\boldsymbol{\theta}^{[r]}$ maximizes the complete-data log-likelihood

$$\ln p(\mathbf{x}, \mathbf{t}^{[r-1]} | \boldsymbol{\theta}) = \sum_{i=1}^n \sum_{k=1}^g t_{ik}^{[r-1]} \ln \left(\pi_k f_k(\mathbf{x}_i; \boldsymbol{\theta}_k) \right).$$

Model-based clustering: model selection

Model

It is defined by the number of components, the family of the components, constraints among parameters...

Question

How can we perform model selection?

Standard approach

Define the set of competing models \mathcal{M} by considering the maximum number of components g_{\max} . The selected model maximizes an information criterion.

Exhaustive approach

Computation of an information criterion for each model in \mathcal{M} .

Tools

Information criteria (BIC, ICL,...) which penalize the log-likelihood.

Model-based clustering: model selection

Model

It is defined by the number of components, the family of the components, constraints among parameters...

Question

How can we perform model selection?

Standard approach

Define the set of competing models \mathcal{M} by considering the maximum number of components g_{\max} . The selected model maximizes an information criterion.

Exhaustive approach

Computation of an information criterion for each model in \mathcal{M} .

Tools

Information criteria (BIC, ICL,...) which penalize the log-likelihood.

Bayesian Information Criterion: a consistent criterion

With a uniform prior for $\omega \in \mathcal{M}$:

$$p(\omega|\mathbf{x}) \propto p(\mathbf{x}|\omega) \text{ o\`u } p(\mathbf{x}|\omega) = \int p(\mathbf{x}|\omega, \boldsymbol{\theta})p(\boldsymbol{\theta}|\omega)d\boldsymbol{\theta}.$$

To assess $\ln p(\mathbf{x}|\omega)$, the BIC (Schwarz G., 1978) uses a Laplace approximation:

$$\text{BIC}(\omega) = \ell(\hat{\boldsymbol{\theta}}_{\omega}; \omega, \mathbf{x}) - \frac{\nu_{\omega}}{2} \ln n,$$

where $\hat{\boldsymbol{\theta}}_{\omega}$ is the MLE and where ν_{ω} denotes the number of parameters of model ω .

To summarize

- Trade-off: Accuracy / Complexity.
- Consistent criterion (Keribin C., 2000).
- The MLE is needed.
- The clustering purpose is not considered.

Integrated Complete-data Likelihood: a criterion for clustering

The ICL (Biernacki C., Celeux G., Govaert G., 2000) considers the clustering purpose:

$$\text{ICL}(\omega) = \ln p(\mathbf{x}, \hat{\mathbf{z}}|\omega) \text{ où } p(\mathbf{x}, \mathbf{z}|\omega) = \int p(\mathbf{x}, \mathbf{z}|\omega, \boldsymbol{\theta})p(\boldsymbol{\theta}|\omega)d\boldsymbol{\theta},$$

where $\hat{\mathbf{z}}$ is the partition given by the MAP rule with $\hat{\boldsymbol{\theta}}_\omega$.

$p(\mathbf{x}, \mathbf{z}|\omega)$ has closed-form when the mixture components belong to the exponential family (+ conjugate prior).

Other, we use a Laplace approximation

$$\text{ICL}(\omega) \simeq \text{BIC}(\omega) + \sum_{i=1}^n \sum_{k=1}^g \hat{z}_{ik} \ln t_{ik}(\hat{\boldsymbol{\theta}}_\omega).$$

To summarize

- Trade-off: Accuracy / Complexity / Overlap of clusters.
- The clustering purpose is considered.
- The MLE is needed.
- Non-consistent criterion.
- Robustness for model misspecification.

The problem

What variables should be considered in clustering?

- Well-posed problem in the supervised classification setting with objective criteria: error rate, AUC, ...
- Ill-posed problem in clustering since the class variable is not known by advance. Thus, what are the most relevant variables with respect to this unknown variable?
- Pragmatic solution 1: Prior choice of the practitioner among available variables (according to some focus).
- Pragmatic solution 2: Posterior analysis of the correlation between the predicted cluster (based on all the variables) and each variable.

Variable selection for model-based clustering

If a mixture of two isotropic Gaussians is considered then, the minimax expected loss R_n is, up to log factors (Azizyan, Singh and Wasserman, 2013),

$$R_n \approx \kappa^2 \sqrt{\frac{d}{n}}$$

where κ^{-1} is the signal-to-noise ratio.

Now, consider the case where there are $s < d$ relevant features, then

$$\kappa \sqrt{\frac{s \sqrt{\log d}}{\sqrt{n}}} \succ R_n \succ \kappa^2 \sqrt{\frac{s \log d}{n}}$$

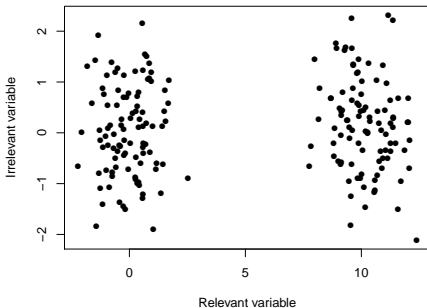
Variable selection in model-based clustering

Main idea

Only a subset of variables explains the partition (but the partition is not observed).

Advantages of variable selection

- Improve the accuracy of the clustering by decreasing the variance of the estimators.
- Allow some specific interpretation of the classifying variables.



Variable selection in model-based clustering

Variable selection is mainly studied for Gaussian mixture

- Tadesse et al. (2005)

$$p(\mathbf{x}_i) = p(\mathbf{x}_i^I) \sum_{k=1}^g \pi_k p(\mathbf{x}_i^R | Z_{ik} = 1)$$

\mathbf{x}_i^I irrelevant variables and \mathbf{x}_i^R relevant variables.

- Raftery and Dean (2006)

$$p(\mathbf{x}_i) = p(\mathbf{x}_i^I) \underbrace{p(\mathbf{x}_i^U | \mathbf{x}_i^R)}_{\text{linear regression}} \sum_{k=1}^g \pi_k p(\mathbf{x}_i^R | Z_{ik} = 1)$$

\mathbf{x}_i^R redundant variables (independent to the partition conditionally on \mathbf{x}_i^R).

- Maugis et al. (2009)

$$p(\mathbf{x}_i) = p(\mathbf{x}_i^I) \underbrace{p(\mathbf{x}_i^U | \mathbf{x}_i^R)}_{\text{sparse linear regression}} \sum_{k=1}^g \pi_k p(\mathbf{x}_i^R | Z_{ik} = 1)$$

Variable selection is difficult because they are many models can be considered!

Objectives

- Deal with many features
- Deal with mixed-type data (with missing values)

Solution: use simpler models for a better search

- Only two types of variables: relevant and irrelevant.
- Assumption of independence of the relevant variables given the cluster.
- Non-classifying variables are independent.
- Optimization of an information criterion (BIC or MICL).

Objectives

- Deal with many features
- Deal with mixed-type data (with missing values)

Solution: use simpler models for a better search

- Only two types of variables: relevant and irrelevant.
- Assumption of independence of the relevant variables given the cluster.
- Non-classifying variables are independent.
- Optimization of an information criterion (BIC or MICL).

Mixture with conditional independence

Conditional independence is not a big deal when the number of features is large.

Pmf of mixture model with conditional independence (Gaussian case)

$$f(\mathbf{x}_i | g, \theta) = \sum_{k=1}^g \pi_k \prod_{j=1}^d \phi(x_{ij}, \mu_{kj}, \sigma_{kj}^2).$$

d	3	4	5	10	20
True model	0.155	0.175	0.198	0.231	0.321
Locally independence	0.111	0.114	0.118	0.101	0.125

Table: Misclassification rate obtained with $n = 50$ (best error 10%)

d	3	4	5	10	20
True model	0.101	0.101	0.104	0.106	0.110
Locally independence	0.103	0.105	0.107	0.108	0.104

Table: Misclassification rate obtained with $n = 2000$ (best error 10%)

Variable selection for model-based clustering

Main idea

Only a subset of variables explains the partition (but the partition is not observed).

Twin objective

Clustering + Find the relevant variables

A variable is irrelevant for clustering if its distribution is the same among components.

Modeling the role of the variables

Model is defined by $\omega = (\omega_j; j = 1, \dots, d)$ with

$$\omega_j = \begin{cases} 0 & \text{if } X_j \text{ is irrelevant} \\ 1 & \text{if } X_j \text{ is relevant} \end{cases} \quad \text{thus} \quad \omega_j = \begin{cases} 0 & \text{if } \theta_{1j} = \dots = \theta_{gj} \\ 1 & \text{otherwise} \end{cases} .$$

Pdf

$$f(\mathbf{x}_i; \omega, \theta) = \prod_{j \in \Omega^c} f_{1j}(x_{ij}; \theta_{1j}) \sum_{k=1}^g \pi_k \prod_{j \in \Omega} f_{kj}(x_{ij}; \theta_{kj}),$$

with $\Omega = \{j : \omega_j = 1\}$ et $\Omega^c = \{1, \dots, d\} \setminus \Omega$.

Type of variable	$X_j \mathbf{Z}_j = \mathbf{z}_j$	Priors	Jeffrey
categorical	\mathcal{M}	\mathcal{D}	yes
continuous	\mathcal{N}	\mathcal{N} et \mathcal{IG}	no
count	\mathcal{P}	\mathcal{Ga}	no

How to select ω (g is assumed to be known)?

Variable selection for model-based clustering

Main idea

Only a subset of variables explains the partition (but the partition is not observed).

Twin objective

Clustering + Find the relevant variables

A variable is irrelevant for clustering if its distribution is the same among components.

Modeling the role of the variables

Model is defined by $\omega = (\omega_j; j = 1, \dots, d)$ with

$$\omega_j = \begin{cases} 0 & \text{if } X_j \text{ is irrelevant} \\ 1 & \text{if } X_j \text{ is relevant} \end{cases} \quad \text{thus} \quad \omega_j = \begin{cases} 0 & \text{if } \theta_{1j} = \dots = \theta_{gj} \\ 1 & \text{otherwise} \end{cases} .$$

Pdf

$$f(\mathbf{x}_i; \omega, \theta) = \prod_{j \in \Omega^c} f_{1j}(x_{ij}; \theta_{1j}) \sum_{k=1}^g \pi_k \prod_{j \in \Omega} f_{kj}(x_{ij}; \theta_{kj}),$$

with $\Omega = \{j : \omega_j = 1\}$ et $\Omega^c = \{1, \dots, d\} \setminus \Omega$.

Type of variable	$X_j \mathbf{Z}_j = \mathbf{z}_j$	Priors	Jeffrey
categorical	\mathcal{M}	\mathcal{D}	yes
continuous	\mathcal{N}	\mathcal{N} et \mathcal{IG}	no
count	\mathcal{P}	\mathcal{Ga}	no

How to select ω (g is assumed to be known)?

Variable selection with BIC

Models in competition

$$\mathcal{M} = \{\omega; \omega \in \{0, 1\}^d\}.$$

g is supposed to be known (for ease of explanations).

If g is unknown, the procedure is repeated for each possible values of g .

Objective

From the sample \mathbf{x} , find ω^* maximizing the BIC

$$\omega^* = \operatorname{argmax}_{\omega \in \mathcal{M}} \text{BIC}(\omega) \text{ with } \text{BIC}(\omega) = \ell(\hat{\theta}_\omega; \omega, \mathbf{x}) - \frac{\nu_\omega}{2} \ln n,$$

where $\hat{\theta}_\omega$ is the MLE and ν_ω is the number of parameters, for model ω .

Classical (but sub-optimal) approaches for maximizing BIC

- Exhaustive: no doable if d large $\text{card}(\mathcal{M}) = 2^d$.
- Backward, forward: Raftery, A. and Dean, D. (2006).

Variable selection with BIC

Models in competition

$$\mathcal{M} = \{\omega; \omega \in \{0, 1\}^d\}.$$

g is supposed to be known (for ease of explanations).

If g is unknown, the procedure is repeated for each possible values of g .

Objective

From the sample \mathbf{x} , find ω^* maximizing the BIC

$$\omega^* = \operatorname{argmax}_{\omega \in \mathcal{M}} \operatorname{BIC}(\omega) \text{ with } \operatorname{BIC}(\omega) = \ell(\hat{\theta}_\omega; \omega, \mathbf{x}) - \frac{\nu_\omega}{2} \ln n,$$

where $\hat{\theta}_\omega$ is the MLE and ν_ω is the number of parameters, for model ω .

Classical (but sub-optimal) approaches for maximizing BIC

- Exhaustive: no doable if d large $\operatorname{card}(\mathcal{M}) = 2^d$.
- Backward, forward: Raftery, A. and Dean, D. (2006).

Variable selection with BIC

Models in competition

$$\mathcal{M} = \{\omega; \omega \in \{0, 1\}^d\}.$$

g is supposed to be known (for ease of explanations).

If g is unknown, the procedure is repeated for each possible values of g .

Objective

From the sample \mathbf{x} , find ω^* maximizing the BIC

$$\omega^* = \operatorname{argmax}_{\omega \in \mathcal{M}} \operatorname{BIC}(\omega) \text{ with } \operatorname{BIC}(\omega) = \ell(\hat{\theta}_\omega; \omega, \mathbf{x}) - \frac{\nu_\omega}{2} \ln n,$$

where $\hat{\theta}_\omega$ is the MLE and ν_ω is the number of parameters, for model ω .

Classical (but sub-optimal) approaches for maximizing BIC

- Exhaustive: no doable if d large $\operatorname{card}(\mathcal{M}) = 2^d$.
- Backward, forward: Raftery, A. and Dean, D. (2006).

Variable selection with BIC

Objective

$$\arg \max_{(\omega, \theta)} \ell(\theta; \omega, \mathbf{x}) - \frac{\nu_{\omega}}{2} \ln n,$$

Penalized complete-data likelihood

$$\ell_p(\theta | \omega, \mathbf{x}, \mathbf{z}) = \sum_{j=1}^d v_j(\theta_{1j}, \dots, \theta_{gj} | \omega_j, \mathbf{x}, \mathbf{z}) + \sum_{k=1}^g \sum_{i=1}^n z_{ik} \ln \pi_k - (g-1) \frac{\ln n}{2},$$

with $v_j(\theta_{1j}, \dots, \theta_{gj} | \omega_j, \mathbf{x}, \mathbf{z})$ has closed form for $\omega_j = 1$ and $\omega_j = 0$ (and is easily maximized).

If X_j is binary, then

$$v_j(\theta_{1j}, \dots, \theta_{gj} | \omega_j, \mathbf{x}, \mathbf{z}) = \begin{cases} \sum_{k=1}^g \sum_{i=1}^n z_{ik} \ln [(\alpha_{kj})^{x_{ij}} (1 - \alpha_{kj})^{1-x_{ij}}] & \text{if } \omega_j = 1 \\ \sum_{i=1}^n \ln [(\alpha_{1j})^{x_{ij}} (1 - \alpha_{1j})^{1-x_{ij}}] & \text{if } \omega_j = 0. \end{cases}$$

The penalized complete-data likelihood is easily maximized on (ω, θ) .

Variable selection with BIC

Objective

$$\arg \max_{(\omega, \theta)} \ell(\theta; \omega, \mathbf{x}) - \frac{\nu_{\omega}}{2} \ln n,$$

Penalized complete-data likelihood

$$\ell_p(\theta | \omega, \mathbf{x}, \mathbf{z}) = \sum_{j=1}^d v_j(\theta_{1j}, \dots, \theta_{gj} | \omega_j, \mathbf{x}, \mathbf{z}) + \sum_{k=1}^g \sum_{i=1}^n z_{ik} \ln \pi_k - (g-1) \frac{\ln n}{2},$$

with $v_j(\theta_{1j}, \dots, \theta_{gj} | \omega_j, \mathbf{x}, \mathbf{z})$ has closed form for $\omega_j = 1$ and $\omega_j = 0$ (and is easily maximized).

If X_j is binary, then

$$v_j(\theta_{1j}, \dots, \theta_{gj} | \omega_j, \mathbf{x}, \mathbf{z}) = \begin{cases} \sum_{k=1}^g \sum_{i=1}^n z_{ik} \ln [(\alpha_{kj})^{x_{ij}} (1 - \alpha_{kj})^{1-x_{ij}}] & \text{if } \omega_j = 1 \\ \sum_{i=1}^n \ln [(\alpha_{1j})^{x_{ij}} (1 - \alpha_{1j})^{1-x_{ij}}] & \text{if } \omega_j = 0. \end{cases}$$

The penalized complete-data likelihood is easily maximized on (ω, θ) .

Variable selection with BIC

Objective

$$\arg \max_{(\boldsymbol{\omega}, \boldsymbol{\theta})} \ell(\boldsymbol{\theta}; \boldsymbol{\omega}, \mathbf{x}) - \frac{\nu_{\boldsymbol{\omega}}}{2} \ln n,$$

Penalized complete-data likelihood

$$\ell_p(\boldsymbol{\theta} | \boldsymbol{\omega}, \mathbf{x}, \mathbf{z}) = \sum_{j=1}^d v_j(\boldsymbol{\theta}_{1j}, \dots, \boldsymbol{\theta}_{gj} | \omega_j, \mathbf{x}, \mathbf{z}) + \sum_{k=1}^g \sum_{i=1}^n z_{ik} \ln \pi_k - (g-1) \frac{\ln n}{2},$$

with $v_j(\boldsymbol{\theta}_{1j}, \dots, \boldsymbol{\theta}_{gj} | \omega_j, \mathbf{x}, \mathbf{z})$ has closed form for $\omega_j = 1$ and $\omega_j = 0$ (and is easily maximized).

If X_j is binary, then

$$v_j(\boldsymbol{\theta}_{1j}, \dots, \boldsymbol{\theta}_{gj} | \omega_j, \mathbf{x}, \mathbf{z}) = \begin{cases} \sum_{k=1}^g \sum_{i=1}^n z_{ik} \ln [(\alpha_{kj})^{x_{ij}} (1 - \alpha_{kj})^{1-x_{ij}}] & \text{if } \omega_j = 1 \\ \sum_{i=1}^n \ln [(\alpha_{1j})^{x_{ij}} (1 - \alpha_{1j})^{1-x_{ij}}] & \text{if } \omega_j = 0. \end{cases}$$

The penalized complete-data likelihood is easily maximized on $(\boldsymbol{\omega}, \boldsymbol{\theta})$.

Variable selection with BIC

The penalized complete-data likelihood is a sum over the variables.

EM Algorithm for assessing the model and its parameters

- EM algorithm for maximizing $\ell(\boldsymbol{\theta}; \boldsymbol{\omega}, \mathbf{x}) - \frac{\nu_{\boldsymbol{\omega}}}{2} \ln n$ on $(\boldsymbol{\omega}, \boldsymbol{\theta})$.
- E-step: unchanged.
- M-step: update ω_j and $\boldsymbol{\theta}_{1j}, \dots, \boldsymbol{\theta}_{gj}$ for each j independently.

This algorithm can be used to perform variable selection according to any IC defined by

$$\ell_p(\boldsymbol{\theta}|\boldsymbol{\omega}, \mathbf{x}) = \ell(\boldsymbol{\theta}|\boldsymbol{\omega}, \mathbf{x}) - \nu_{\boldsymbol{\omega}}c,$$

where c is a constant.

Some experiences

	Data		VarSelLCM			clustvarsel (forward/backward)		
	n	d	NRV	ARI	Time	NRV	ARI	Time
banknote	200	6	5	0.96	1	4/4	0.98/0.98	5/4
coffee	43	12	5	1.00	1	5/6	1.00/1.00	5/13
cancer	569	30	28	0.68	16	17/19	0.66/0.64	2160/14435
golub	83	3051	769	0.70	14	8/●	0.11/●	3414/●

Table: Results obtained on four benchmark datasets: number of relevant variables (NRV), Adjusted Rand Index computed on the selected model (ARI) and computing time in seconds (Time)

Asymptotic and exact criteria

Bayesian model selection

The aim is to find the model maximizing the integrated likelihood:

$$p(\mathbf{x}|\boldsymbol{\omega}) = \int p(\mathbf{x}|\boldsymbol{\omega}, \boldsymbol{\theta})p(\boldsymbol{\theta}|\boldsymbol{\omega})d\boldsymbol{\theta}.$$

Asymptotic criterion

BIC is asymptotic because

$$\ln p(\mathbf{x}|\boldsymbol{\omega}) = \ln p(\mathbf{x}|\boldsymbol{\omega}, \hat{\boldsymbol{\theta}}_{\boldsymbol{\omega}}) - \frac{\nu_{\boldsymbol{\omega}}}{2} \ln n + \mathcal{O}(1),$$

but its approximation can leads to poor results if $n \ll d$.

⇒ An exact criterion is more relevant.

Clustering purpose

In clustering, we want to obtain well-separated classes.

⇒ A criterion considering the entropy between components should be used.

Variable selection with MICL

Integrated complete-data likelihood

The complete data likelihood has a closed form if conjugated priors are used

$$p(\mathbf{x}, \mathbf{z}|\boldsymbol{\omega}) = \int p(\mathbf{x}, \mathbf{z}|\boldsymbol{\omega}, \boldsymbol{\theta})p(\boldsymbol{\theta}|\boldsymbol{\omega})d\boldsymbol{\theta}.$$

For the following Jeffrey's priors are used (Dirichlet with hyper-parameters 1/2).

$$p(\mathbf{x}, \mathbf{z}|\boldsymbol{\omega}) = p(\mathbf{z}) \prod_{j=1}^d p(\mathbf{x}_{\bullet j}|\omega_j, \mathbf{z}),$$

where $p(\mathbf{x}_{\bullet j}|\omega_j, \mathbf{z})$ has closed form for $\omega_j = 1$ and $\omega_j = 0$

If X_j is binary,

$$p(\mathbf{x}_{\bullet j}|g, \omega_j, \mathbf{z}) = \begin{cases} \left(\frac{1}{\Gamma(\frac{1}{2})^2} \right)^g \prod_{k=1}^g \frac{\Gamma(n_k - n_{k1} + \frac{1}{2})\Gamma(n_{k1} + \frac{1}{2})}{\Gamma(n_k + 1)} & \text{si } \omega_j = 1 \\ \frac{1}{\Gamma(\frac{1}{2})^2} \frac{\Gamma(n - \sum_{i=1}^n x_{ij} + \frac{1}{2})\Gamma(\sum_{i=1}^n x_{ij} + \frac{1}{2})}{\Gamma(n + 1)} & \text{si } \omega_j = 0 \end{cases}$$

où $n_{k1} = \sum_{i=1}^n z_{ik}x_{ij}$.

Variable selection with MICL

Integrated complete-data likelihood

The complete data likelihood has a closed form if conjugated priors are used

$$p(\mathbf{x}, \mathbf{z} | \boldsymbol{\omega}) = \int p(\mathbf{x}, \mathbf{z} | \boldsymbol{\omega}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | \boldsymbol{\omega}) d\boldsymbol{\theta}.$$

For the following Jeffrey's priors are used (Dirichlet with hyper-parameters 1/2).

$$p(\mathbf{x}, \mathbf{z} | \boldsymbol{\omega}) = p(\mathbf{z}) \prod_{j=1}^d p(\mathbf{x}_{\bullet j} | \omega_j, \mathbf{z}),$$

where $p(\mathbf{x}_{\bullet j} | \omega_j, \mathbf{z})$ has closed form for $\omega_j = 1$ and $\omega_j = 0$

If X_j is binary,

$$p(\mathbf{x}_{\bullet j} | g, \omega_j, \mathbf{z}) = \begin{cases} \left(\frac{1}{\Gamma(\frac{1}{2})^2} \right)^g \prod_{k=1}^g \frac{\Gamma(n_k - n_{k1} + \frac{1}{2}) \Gamma(n_{k1} + \frac{1}{2})}{\Gamma(n_k + 1)} & \text{si } \omega_j = 1 \\ \frac{1}{\Gamma(\frac{1}{2})^2} \frac{\Gamma(n - \sum_{i=1}^n x_{ij} + \frac{1}{2}) \Gamma(\sum_{i=1}^n x_{ij} + \frac{1}{2})}{\Gamma(n + 1)} & \text{si } \omega_j = 0 \end{cases}$$

où $n_{k1} = \sum_{i=1}^n z_{ik} x_{ij}$.

Variable selection with MICL

The complete-data likelihood is

$$p(\mathbf{x}, \mathbf{z}|\boldsymbol{\omega}) = p(\mathbf{z}) \prod_{j=1}^d p(\mathbf{x}_{\bullet j}|\omega_j, \mathbf{z}),$$

Criterion requiring the computation of the MLE

$$\text{ICL}(\boldsymbol{\omega}) = \ln p(\mathbf{x}, \hat{\mathbf{z}}_{\boldsymbol{\omega}}|\boldsymbol{\omega}) \text{ with } \hat{\mathbf{z}}_{\boldsymbol{\omega}} = \arg \max_{\mathbf{z}} \ln p(\mathbf{z}|\boldsymbol{\omega}, \mathbf{x}, \hat{\boldsymbol{\theta}}_{\boldsymbol{\omega}}).$$

Criterion avoiding the computation of the MLE

$$\text{MICL}(\boldsymbol{\omega}) = \ln p(\mathbf{x}, \mathbf{z}_{\boldsymbol{\omega}}^*|\boldsymbol{\omega}) \text{ with } \mathbf{z}_{\boldsymbol{\omega}}^* = \arg \max_{\mathbf{z}} \ln p(\mathbf{z}|\boldsymbol{\omega}, \mathbf{x}).$$

Variable selection with MICL

Objective

Maximization of MICL implies the maximization on (ω, \mathbf{z}) of

$$p(\mathbf{x}, \mathbf{z}|\omega) = p(\mathbf{z}) \prod_{j=1}^d p(\mathbf{x}_{\bullet j}|\omega_j, \mathbf{z}).$$

⇒ Model selection without parameter estimation!

Algorithm for assessing the model

An iterative algorithm of optimization is used, where iteration $[r]$ is:

- **Partition step (iterative):**

Find $\mathbf{z}^{[r]}$ with

$$\ln p(\mathbf{x}, \mathbf{z}^{[r]}|\omega^{[r]}) \geq \ln p(\mathbf{x}, \mathbf{z}^{[r-1]}|\omega^{[r]}).$$

- **Model step (explicit):**

Find $\omega^{[r+1]} = \operatorname{argmax}_{\omega \in \mathcal{M}} \ln p(\mathbf{x}, \mathbf{z}^{[r]}|\omega)$, where

$$\omega_j^{[r+1]} = \operatorname{argmax}_{\omega_j \in \{0,1\}} p(\mathbf{x}_{\bullet j}|\omega_j, \mathbf{z}^{[r]}).$$

Properties

- The objective function increases at each iteration.
- Convergence to local optimum (like for EM algorithm).

Variable selection with MICL

Objective

Maximization of MICL implies the maximization on (ω, \mathbf{z}) of

$$p(\mathbf{x}, \mathbf{z}|\omega) = p(\mathbf{z}) \prod_{j=1}^d p(\mathbf{x}_{\bullet j}|\omega_j, \mathbf{z}).$$

⇒ Model selection without parameter estimation!

Algorithm for assessing the model

An iterative algorithm of optimization is used, where iteration $[r]$ is:

- **Partition step (iterative):**

Find $\mathbf{z}^{[r]}$ with

$$\ln p(\mathbf{x}, \mathbf{z}^{[r]}|\omega^{[r]}) \geq \ln p(\mathbf{x}, \mathbf{z}^{[r-1]}|\omega^{[r]}).$$

- **Model step (explicit):**

Find $\omega^{[r+1]} = \operatorname{argmax}_{\omega \in \mathcal{M}} \ln p(\mathbf{x}, \mathbf{z}^{[r]}|\omega)$, where

$$\omega_j^{[r+1]} = \operatorname{argmax}_{\omega_j \in \{0,1\}} p(\mathbf{x}_{\bullet j}|\omega_j, \mathbf{z}^{[r]}).$$

Properties

- The objective function increases at each iteration.
- Convergence to local optimum (like for EM algorithm).

Variable selection with MICL

Objective

Maximization of MICL implies the maximization on (ω, \mathbf{z}) of

$$p(\mathbf{x}, \mathbf{z}|\omega) = p(\mathbf{z}) \prod_{j=1}^d p(\mathbf{x}_{\bullet j}|\omega_j, \mathbf{z}).$$

⇒ Model selection without parameter estimation!

Algorithm for assessing the model

An iterative algorithm of optimization is used, where iteration $[r]$ is:

- **Partition step (iterative):**

Find $\mathbf{z}^{[r]}$ with

$$\ln p(\mathbf{x}, \mathbf{z}^{[r]}|\omega^{[r]}) \geq \ln p(\mathbf{x}, \mathbf{z}^{[r-1]}|\omega^{[r]}).$$

- **Model step (explicit):**

Find $\omega^{[r+1]} = \operatorname{argmax}_{\omega \in \mathcal{M}} \ln p(\mathbf{x}, \mathbf{z}^{[r]}|\omega)$, where

$$\omega_j^{[r+1]} = \operatorname{argmax}_{\omega_j \in \{0,1\}} p(\mathbf{x}_{\bullet j}|\omega_j, \mathbf{z}^{[r]}).$$

Properties

- The objective function increases at each iteration.
- Convergence to local optimum (like for EM algorithm).

Details on the partition step

At iteration $[r]$: $\mathbf{z}^{[r]}$ is initialized $\mathbf{z}^{(0)} = \mathbf{z}^{[r-1]}$. S iterations are made where iteration (s) is composed of two steps:

- **Sampling:** $i^{(s)} \sim \mathcal{U}\{1, \dots, n\}$.
- **Optimization:** hold $\mathcal{Z}^{(s)} = \{\mathbf{z} : z_i = z_i^{(s-1)}, \forall i \neq i^{(s)}\}$ and find $\mathbf{z}^{(s)}$ such that

$$\mathbf{z}^{(s)} = \operatorname{argmax}_{\mathbf{z} \in \mathcal{Z}^{(s)}} \ln p(\mathbf{x}, \mathbf{z} | \boldsymbol{\omega}^{[r]}).$$

Stopping criterion : $\ln p(\mathbf{x}, \mathbf{z} | \boldsymbol{\omega}^{[r]})$ does not increase after S iterations

Summary

- Two fast and simple steps.
- $\ln p(\mathbf{x}, \mathbf{z} | \boldsymbol{\omega})$ increases at each iteration.
- Many random starting points are needed.

Details on the partition step

At iteration $[r]$: $\mathbf{z}^{[r]}$ is initialized $\mathbf{z}^{(0)} = \mathbf{z}^{[r-1]}$. S iterations are made where iteration (s) is composed of two steps:

- **Sampling:** $i^{(s)} \sim \mathcal{U}\{1, \dots, n\}$.
- **Optimization:** hold $\mathcal{Z}^{(s)} = \{\mathbf{z} : z_i = z_i^{(s-1)}, \forall i \neq i^{(s)}\}$ and find $\mathbf{z}^{(s)}$ such that

$$\mathbf{z}^{(s)} = \operatorname{argmax}_{\mathbf{z} \in \mathcal{Z}^{(s)}} \ln p(\mathbf{x}, \mathbf{z} | \boldsymbol{\omega}^{[r]}).$$

Stopping criterion : $\ln p(\mathbf{x}, \mathbf{z} | \boldsymbol{\omega}^{[r]})$ does not increase after S iterations

Summary

- Two fast and simple steps.
- $\ln p(\mathbf{x}, \mathbf{z} | \boldsymbol{\omega})$ increases at each iteration.
- Many random starting points are needed.

Some experiences

Data	n	d	VarSelLCM (BIC / MICL)			clustvarsel (forward/backward)		
			NRV	ARI	Time	NRV	ARI	Time
banknote	200	6	5/5	0.96/0.96	1/1	4/4	0.98/0.98	5/4
coffee	43	12	5/5	1.00/1.00	1/1	5/6	1.00/1.00	5/13
cancer	569	30	28/15	0.68/0.75	18/24	17/19	0.66/0.64	2160/14435
golub	83	3051	769/553	0.70/0.79	16/15	8/●	0.11/●	3414/●

Table: Results obtained on four benchmark datasets: number of relevant variables (NRV), Adjusted Rand Index computed on the selected model (ARI) and computing time in seconds (Time)

EDEN

Data

- Obesity prevention
- Data collected on 1436 individuals between 2 and 5 years
- 44 features (continuous and categorical, with missing values) describing:
 - Dietary intake and eating habits
 - Physical activity measures
 - Sedentary behavior measures
 - Sleep measures

Results

- 17/44 features are detected as discriminative with two clusters.
- The discriminative features mainly describe dietary intake and eating habits.
- Whats about the other families of features?

EDEN

Data

- Obesity prevention
- Data collected on 1436 individuals between 2 and 5 years
- 44 features (continuous and categorical, with missing values) describing:
 - Dietary intake and eating habits
 - Physical activity measures
 - Sedentary behavior measures
 - Sleep measures

Results

- 17/44 features are detected as discriminative with two clusters.
- The discriminative features mainly describe dietary intake and eating habits.
- Whats about the other families of features?

Multiple partitions in model-based clustering

Several clustering variables

- The variables in the data can convey several clustering view points with respect to different groups of variables
- Allow to find some clustering which could be hidden by other variables

Main assumptions

- The variables can be decomposed in B independent blocks.
- The block b follows a mixture with g_b components (for $b = 1, \dots, B$), with the assumption of class conditional independence of the variables.

Notations

- $\omega = (\omega_j; j = 1, \dots, d)$ the repartition of the variables in blocks; $\omega_j = b$ if variable j belongs to block b .
- $m = (g_1, \dots, g_B, \omega)$ defines the model
- $\Omega_b = \{j : \omega_j = b\}$ the subset of variables belonging to block b

Probability distribution function of x_i

$$p(x_i; m, \theta) = \prod_{b=1}^B p(x_{i\{\omega_b\}}; m, \theta) \text{ with } p(x_{i\{\omega_b\}}; m, \theta) = \sum_{k=1}^{g_b} \pi_{bk} \prod_{j \in \Omega_b} p(x_{ij} | \theta_{jk}),$$

Several clustering variables

- The variables in the data can convey several clustering view points with respect to different groups of variables
- Allow to find some clustering which could be hidden by other variables

Main assumptions

- The variables can be decomposed in B independent blocks.
- The block b follows a mixture with g_b components (for $b = 1, \dots, B$), with the assumption of class conditional independence of the variables.

Notations

- $\omega = (\omega_j; j = 1, \dots, d)$ the repartition of the variables in blocks; $\omega_j = b$ if variable j belongs to block b .
- $m = (g_1, \dots, g_B, \omega)$ defines the model
- $\Omega_b = \{j : \omega_j = b\}$ the subset of variables belonging to block b

Probability distribution function of x_i

$$p(x_i; m, \theta) = \prod_{b=1}^B p(x_{i\{\omega_b\}}; m, \theta) \text{ with } p(x_{i\{\omega_b\}}; m, \theta) = \sum_{k=1}^{g_b} \pi_{bk} \prod_{j \in \Omega_b} p(x_{ij} | \theta_{jk}),$$

Several clustering variables

- The variables in the data can convey several clustering view points with respect to different groups of variables
- Allow to find some clustering which could be hidden by other variables

Main assumptions

- The variables can be decomposed in B independent blocks.
- The block b follows a mixture with g_b components (for $b = 1, \dots, B$), with the assumption of class conditional independence of the variables.

Notations

- $\omega = (\omega_j; j = 1, \dots, d)$ the repartition of the variables in blocks; $\omega_j = b$ if variable j belongs to block b .
- $\mathbf{m} = (g_1, \dots, g_B, \omega)$ defines the model
- $\Omega_b = \{j : \omega_j = b\}$ the subset of variables belonging to block b

Probability distribution function of x_i

$$p(x_i; \mathbf{m}, \theta) = \prod_{b=1}^B p(x_{i\{\omega_b\}}; \mathbf{m}, \theta) \text{ with } p(x_{i\{b\}}; \mathbf{m}, \theta) = \sum_{k=1}^{g_b} \pi_{bk} \prod_{j \in \Omega_b} p(x_{ij} | \theta_{jk}),$$

Probability distribution function of \mathbf{x}_i

$$p(\mathbf{x}_i; \mathbf{m}, \boldsymbol{\theta}) = \prod_{b=1}^B \left(\sum_{k=1}^{g_b} \pi_{bk} \prod_{j \in \Omega_b} p(x_{ij} | \boldsymbol{\theta}_{jk}) \right)$$

Remarks

- Different partitions explained by subsets of variables.
- Generalizes approaches used for variable selection in model-based clustering (if $B = 2$ and $g_1 = 1$ then variables belonging to block 1 are not relevant for the clustering, while variables belonging to block 2 are relevant).
- Natural extension to the heterogeneous data setting.
- Model identifiability is directly obtained from the identifiability of Gaussian mixture with local independence (Teicher, 1963, 1967).
- Model selection (*i.e.*, find the blocks of variables) can be carried out via modified EM algorithm (BIC) or via the optimization of the integrated complete-data likelihood (MICL).

EDEN (continued)

Data

- Obesity prevention
- Data collected on 1436 individuals between 2 and 5 years
- 44 features (continuous and categorical, with missing values) describing:
 - Dietary intake and eating habits
 - Physical activity measures
 - Sedentary behavior measures
 - Sleep measures

Results

- Best model considers 3 blocks of variables.
- 17/44 features are grouped into block 1 (mainly dietary intake and eating habits) with two components.
- 7/44 features are grouped into block 2 (mainly physical activity measures and sedentary behavior measures) with two components.
- 20/44 features are grouped into block 3 which contains non-discriminative variables.

EDEN (continued)

Data

- Obesity prevention
- Data collected on 1436 individuals between 2 and 5 years
- 44 features (continuous and categorical, with missing values) describing:
 - Dietary intake and eating habits
 - Physical activity measures
 - Sedentary behavior measures
 - Sleep measures

Results

- Best model considers 3 blocks of variables.
- 17/44 features are grouped into block 1 (mainly dietary intake and eating habits) with two components.
- 7/44 features are grouped into block 2 (mainly physical activity measures and sedentary behavior measures) with two components.
- 20/44 features are grouped into block 3 which contains non-discriminative variables.

Model-based clustering

Take home message

- Conditional independence facilitates variable selection...
- ... but if d is small and n large enough, use more complex models!

VarSelLCM

- Available on CRAN: <https://cran.r-project.org/web/packages/VarSelLCM/>
- Tutorial: <http://varsellcm.r-forge.r-project.org/>
 - Clustering of mixed data (continuous, categorical and integer) with missing values
 - Variable selection (BIC or MICL)

References

- Marbac, M. and Sedki, M. Variable selection for model-based clustering using the integrated complete-data likelihood. *Statistics and Computing*, 27 (4), 2017
- Marbac, M. and Sedki, M. VarSelLCM: an R/C++ package for variable selection in model-based clustering of mixed-data with missing values. *Bioinformatics*, 2018
- Marbac, M. and Patin, E. and Sedki, M. Variable selection for mixed data clustering: Application in human population genomics. *Journal of Classification*, (forthcoming)
- Marbac, M., and Vandewalle, V. A tractable Multi-Partitions Clustering. *Computational Statistics and Data Analysis*, (forthcoming)