



**HAL**  
open science

## blockcluster: An R Package for Model Based Co-Clustering

Parmeet Bhatia, Serge Iovleff, Gérard Govaert

► **To cite this version:**

Parmeet Bhatia, Serge Iovleff, Gérard Govaert. blockcluster: An R Package for Model Based Co-Clustering. Journal of Statistical Software, 2017, 76 (9), pp.24. 10.18637/jss.v076.i09 . hal-01093554

**HAL Id: hal-01093554**

**<https://inria.hal.science/hal-01093554>**

Submitted on 10 Dec 2014

**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.



## blockcluster: An R Package for Model Based Co-Clustering

Parmeet Bhatia  
INRIA-Lille

Serge Iovleff  
Univ. Lille 1, CNRS, Inria

G erard Govaert  
UTC, CNRS

---

### Abstract

Simultaneous clustering of rows and columns, usually designated by bi-clustering, co-clustering or block clustering, is an important technique in two way data analysis. A new standard and efficient approach have been recently proposed based on latent block model [Govaert and Nadif (2003)] which takes into account the block clustering problem on both the individual and variables sets. This article presents our R package for co-clustering of binary, contingency and continuous data **blockcluster** based on these very models. In this document, we will give a brief review of the model-based block clustering methods, and we will show how the R package **blockcluster** can be used for co-clustering.

*Keywords:* model based clustering, block mixture model, EM and CEM algorithms, simultaneous clustering, co-clustering, R, **blockcluster**.

---

## 1. Introduction

Cluster analysis is an important tool in a variety of scientific areas such as pattern recognition, information retrieval, micro-array, data mining, and so forth. Although many clustering procedures such as hierarchical clustering,  $k$ -means or self-organizing maps, aim to construct an optimal partition of objects or, sometimes, of variables, there are other methods, called block clustering methods, which consider simultaneously the two sets and organize the data into homogeneous blocks. Let  $\mathbf{x}$  denotes a  $n \times d$  data matrix defined by  $\mathbf{x} = \{(x_{ij}); i \in I \text{ and } j \in J\}$ , where  $I$  is a set of  $n$  objects (rows, observations, cases etc) and  $J$  is a set of  $d$  variables (columns, attributes etc). The basic idea of these methods consists in making permutations of objects and variables in order to draw a correspondence structure on  $I \times J$ . For illustration, consider Figure 1 where a binary data set defined on set of  $n = 10$  individuals  $I = A, B, C, D, E, F, G, H, I, J$  and set of  $d = 7$  binary variables  $J = 1, 2, 3, 4, 5, 6, 7$  is re-organized into a set of  $3 \times 3$  clusters by permuting the rows and columns.

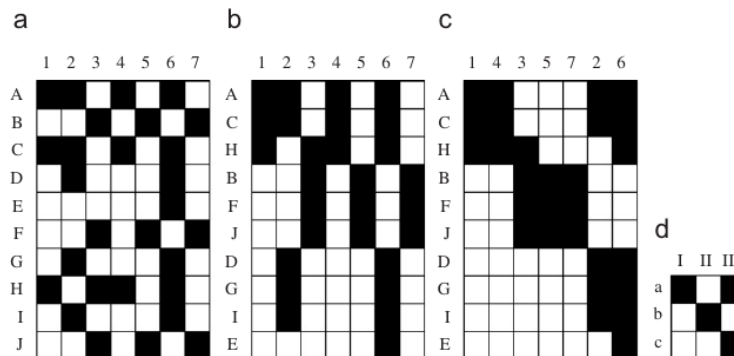


Figure 1: Binary data set (a), data reorganized by a partition on  $I$  (b), by partitions on  $I$  and  $J$  simultaneously (c) and summary matrix (d).

In recent years, co-clustering have found numerous applications in the fields ranging from data mining, information retrieval, biology, computer vision and so forth. [Dhillon (2001)] published an article on text data mining by simultaneously clustering the documents and content (words) using bipartite spectral graph partitioning. This is quite useful technique for instance to manage huge corpus of unlabeled documents. [Xu, Zong, Dolog, and Zhang (2010)] presents another co-clustering application (again using bipartite spectral graph) to understand subset aggregates of web users by simultaneously clustering the users (sessions) and the page view information. [Giannakidou, Koutsonikola, Vakali, and Kompatsiaris (2008)] employed a similarity metric based co-clustering technique for social tagging system. In field of Bio-informatics, co-clustering is mainly used to find structures in gene expression data. This is useful for instance to find set of genes correspond to particular kind of disease. Some of the pioneer material in this context can be found in [Kluger, Basri, Chang, and Gerstein (2003)]. Recently many co-clustering based algorithms have also been developed to target computer vision applications. For instance, [Qiu (2004)] demonstrated the utility of co-clustering in Image grouping by simultaneously clustering images with their low-level visual features. [Guan, Qiu, and Xue (2005)] extended this work and presented opportunity to develop a novel content based image retrieval system. Similarly, [Rasiwasia and Vasconcelos (2009)] used co-clustering to model scenes. We have presented here with a very brief survey demonstrating the diversity of co-clustering applications. Unfortunately, many of the algorithms developed above for co-clustering are suitable in that particular context only and cannot be considered as generic framework. Here, we take the opportunity to emphasize on the fact that the algorithms used in our software can be more or less adopted in most or all the above context.

In practice, several softwares exist for block clustering such as Bicat [Barkow, Bleuler, Prelic, Zimmermann, and Zitzler (2006)], biclust [Kaiser and Leisch (2008)], BiGGEsTS [Goncalves, Madeira, and Oliveira (2009)], or refer to [Madeira and Oliveira (2004)] for a review. In this paper, we introduce our R [R Development Core Team (2008)] package **blockcluster**, which provides a bridge between a newly developed C++ library and the R statistical computing environment. The R package **blockcluster** allows to estimate the parameters of the co-clustering models [Govaert and Nadif (2003)] for binary, contingency and continuous data. This package is unique from the point of view of generative models it implements (latent blocks), the used algorithms (BEM, BCEM) and, apart from that, special attention has been given to design the

library for handling very huge data sets in reasonable time. The R package is already available on CRAN at <http://cran.r-project.org/web/packages/blockcluster/index.html>.

The rest of the article is organized as follows. In Section 2, we give the mathematical foundation (without going into detailed proofs) of latent block models for binary, contingency and continuous data-sets. In Section 3, we elaborate various details of **blockcluster** package and illustrate its usage with the help of various examples. This part can also act as first-hand comprehensive tutorial for the audience who are interested in running the **blockcluster** package directly. Finally we illustrate our package on real data-sets in Section 4 and terminate with some concluding remarks in Section 5.

## 2. Block mixture models

### 2.1. Mixture models: Classic formulation

In model-based clustering, it is assumed that the data arises from a mixture of underlying probability distributions, where each component  $k$  of the mixture represents a cluster. Hence the matrix data  $\mathbf{x}=(\mathbf{x}_1, \dots, \mathbf{x}_n)$  is assumed to be i.i.d and generated from a probability distribution with the density

$$f(\mathbf{x}_i; \boldsymbol{\theta}) = \sum_k \pi_k f_k(\mathbf{x}_i; \boldsymbol{\alpha}),$$

where  $f_k$  is the density function for  $k$ th component. Generally, these densities belong to the same parametric family (say for example Gaussian). The parameters  $(\pi_k)_{k=1}^g$  give the probabilities that an observation belongs to the  $k$ th component and the number of components in the mixture is  $g$  which is assumed to be known. The parameter  $\boldsymbol{\theta}$  of this model is the vector  $(\pi_1, \dots, \pi_g, \boldsymbol{\alpha})$  containing the parameter  $\boldsymbol{\alpha}$  and the mixing proportions  $(\pi_1, \dots, \pi_g)$ . Using this, the density of the observed data  $\mathbf{x}$  can be expressed as

$$f(\mathbf{x}; \boldsymbol{\theta}) = \prod_i \sum_k \pi_k f_k(\mathbf{x}_i; \boldsymbol{\alpha}). \quad (1)$$

An important property of the pdf  $f(\mathbf{x}; \boldsymbol{\theta})$  shown by Govaert and Nadif (2003) is that the density (1) can be rewritten as

$$f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{\mathbf{z} \in \mathcal{Z}} p(\mathbf{z}) f(\mathbf{x}|\mathbf{z}; \boldsymbol{\alpha}) \quad (2)$$

where  $\mathcal{Z}$  denotes the set of all possible partitions of  $I$  in  $g$  clusters,  $f(\mathbf{x}|\mathbf{z}; \boldsymbol{\alpha}) = \prod_i f_{z_i}(\mathbf{x}_i; \boldsymbol{\alpha})$  and  $p(\mathbf{z}) = \prod_i \pi_{z_i}$ . With this formulation, the data matrix is assumed to be a sample of size 1 from a random  $(n, d)$  matrix.

### 2.2. Latent block model: General formulation

Let the data-set  $\mathbf{x}$  be defined on a set  $I$  with  $n$  elements (individuals) and a set  $J$  with  $d$  elements (variables). We represent a partition of  $I$  into  $g$  clusters by  $\mathbf{z} = (z_{11}, \dots, z_{ng})$  with  $z_{ik} = 1$  if  $i$  belongs to cluster  $k$  and  $z_{ik} = 0$  otherwise,  $z_i = k$  if  $z_{ik} = 1$  and we denote by  $z_{\cdot k} = \sum_i z_{ik}$  the cardinality of row cluster  $k$ . Similarly, we represent a partition of  $J$

into  $m$  clusters by  $\mathbf{w} = (w_{11}, \dots, w_{dm})$  with  $w_{j\ell} = 1$  if  $j$  belongs to cluster  $\ell$  and  $w_{j\ell} = 0$  otherwise,  $w_j = \ell$  if  $w_{j\ell} = 1$  and we denote  $w_{\cdot\ell} = \sum_j w_{j\ell}$  the cardinality of column cluster  $\ell$ . Furthermore, note that for any variable  $V_{ij}$ , we denote  $V_{\cdot k}$  to represent  $\sum_i V_{ik}$  and similarly  $V_{i\cdot}$  shall represent  $\sum_k V_{ik}$ .

To study the block clustering problem, the previous formulation (2) is extended to propose block mixture model defined by the following probability density function

$$f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{\mathbf{u} \in \mathcal{U}} p(\mathbf{u}; \boldsymbol{\theta}) f(\mathbf{x}|\mathbf{u}; \boldsymbol{\theta})$$

where  $\mathcal{U}$  denotes the set of all possible labellings of  $I \times J$  and  $\boldsymbol{\theta}$  contains all the unknown parameters of this model. By restricting this model to a set of labellings of  $I \times J$  defined by a product of labellings of  $I$  and  $J$ , and further assuming that the labellings of  $I$  and  $J$  are independent of each other, we obtain the decomposition

$$f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{(\mathbf{z}, \mathbf{w}) \in \mathcal{Z} \times \mathcal{W}} p(\mathbf{z}; \boldsymbol{\theta}) p(\mathbf{w}; \boldsymbol{\theta}) f(\mathbf{x}|\mathbf{z}, \mathbf{w}; \boldsymbol{\theta}) \quad (3)$$

where  $\mathcal{Z}$  and  $\mathcal{W}$  denote the sets of all possible labellings  $\mathbf{z}$  of  $I$  and  $\mathbf{w}$  of  $J$ .

By extending the latent class principle of local independence to our block model, each data point  $x_{ij}$  will be independent once  $z_i$  and  $w_j$  are fixed. Hence we have

$$f(\mathbf{x}|\mathbf{z}, \mathbf{w}; \boldsymbol{\theta}) = \prod_{i,j} f_{z_i w_j}(x_{ij}; \boldsymbol{\alpha})$$

where  $f_{z_i w_j}(x, \boldsymbol{\alpha})$  is either a probability density function (pdf) defined on the real set  $\mathbb{R}$ , or a discrete probability function. Denoting  $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{\rho}, \boldsymbol{\alpha})$ , where  $\boldsymbol{\pi} = (\pi_1, \dots, \pi_g)$  and  $\boldsymbol{\rho} = (\rho_1, \dots, \rho_m)$  are the vectors of probabilities  $\pi_k$  and  $\rho_\ell$  that a row and a column belong to the  $k$ th row component and to the  $\ell$ th column component respectively, we obtain the latent block mixture model with pdf

$$f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{(\mathbf{z}, \mathbf{w}) \in \mathcal{Z} \times \mathcal{W}} \prod_{i,j} \pi_{z_i} \rho_{w_j} f_{z_i w_j}(x_{ij}; \boldsymbol{\alpha}). \quad (4)$$

Using above formulation, the randomized data generation process can be described by the three steps row labellings (R), column labellings(C), and data generation (X) as follows:

- (R) Generate the labellings  $\mathbf{z} = (z_1, \dots, z_n)$  according to the distribution  $\boldsymbol{\pi} = (\pi_1, \dots, \pi_g)$ .
- (C) Generate the labellings  $\mathbf{w} = (w_1, \dots, w_d)$  according to the distribution  $\boldsymbol{\rho} = (\rho_1, \dots, \rho_m)$ .
- (X) Generate for  $i = 1, \dots, n$  and  $j = 1, \dots, d$  a value  $x_{ij}$  according to the (density) distribution  $f_{z_i w_j}(\cdot; \boldsymbol{\alpha})$ .

### 2.3. Model parameters estimation

EM based algorithms are the standard approach to estimate model parameters by maximizing observed data log-likelihood. In our case, the complete data is represented as a vector  $(\mathbf{x}, \mathbf{z}, \mathbf{w})$

where unobservable vectors  $\mathbf{z}$  and  $\mathbf{w}$  are the labels. The *complete data log-likelihood* can then be written

$$L_C(\mathbf{z}, \mathbf{w}, \boldsymbol{\theta}) = \sum_k z_{.k} \log \pi_k + \sum_\ell w_{.\ell} \log \rho_\ell + \sum_{i,j,k,\ell} z_{ik} w_{j\ell} \log f_{k\ell}(x_{ij}; \boldsymbol{\alpha}). \quad (5)$$

The EM algorithm maximizes the log-likelihood  $L(\boldsymbol{\theta})$  iteratively by maximizing the conditional expectation  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(c)})$  of the complete data log-likelihood given a previous current estimate  $\boldsymbol{\theta}^{(c)}$  and  $\mathbf{x}$ :

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(c)}) = \sum_{i,k} t_{ik}^{(c)} \log \pi_k + \sum_{j,\ell} r_{j\ell}^{(c)} \log \rho_\ell + \sum_{i,j,k,\ell} e_{ikj\ell}^{(c)} \log f_{k\ell}(x_{ij}; \boldsymbol{\alpha})$$

where

$$t_{ik}^{(c)} = P(z_{ik} = 1 | \mathbf{x}, \boldsymbol{\theta}^{(c)}), \quad r_{j\ell}^{(c)} = P(w_{j\ell} = 1 | \mathbf{x}, \boldsymbol{\theta}^{(c)}), \quad e_{ikj\ell}^{(c)} = P(z_{ik} w_{j\ell} = 1 | \mathbf{x}, \boldsymbol{\theta}^{(c)})$$

Unfortunately, difficulties arise owing to the dependence structure in the model, and more precisely, to determinate  $e_{ikj\ell}^{(c)}$ . To solve this problem a new approximate solution is proposed [Govaert and Nadif (2003)] using the [Hathaway (1986)] and [Neal and Hinton (1998)] interpretation of the EM algorithm. By constraining the probability  $p$  to verify the relation  $p(z_i, w_j) = p(z_i)p(w_j) \quad \forall i, j$ , it can be shown that the fuzzy clustering criterion proposed by [Neal and Hinton (1998)] takes the form

$$\tilde{F}_C(\mathbf{t}, \mathbf{r}; \boldsymbol{\theta}) = L_C(\mathbf{t}, \mathbf{r}, \boldsymbol{\theta}) + H(\mathbf{t}) + H(\mathbf{r}) \quad (6)$$

in case of latent block model. where

$$H(\mathbf{t}) = - \sum_{ik} t_{ik} \log t_{ik} \quad \text{and} \quad H(\mathbf{r}) = - \sum_{j\ell} r_{j\ell} \log r_{j\ell}$$

and  $L_c$  is the fuzzy complete-data log-likelihood associated to the block latent model given by

$$L_C(\mathbf{t}, \mathbf{r}; \boldsymbol{\theta}) = \sum_k t_{.k} \log \pi_k + \sum_{.\ell} r_{.\ell} \log \rho_\ell + \sum_{i,j,k,\ell} t_{ik} r_{j\ell} \log f_{k\ell}(x_{ij}; \boldsymbol{\alpha}). \quad (7)$$

This is a variational approach and its objective is to maximize the function  $\tilde{F}_C$ . Doing that, it has replace the maximization of the likelihood by the maximization of an approximation of this likelihood defined by

$$\tilde{L}(\boldsymbol{\theta}) = \operatorname{argmax}_{\mathbf{t}, \mathbf{r}} \tilde{F}_C(\mathbf{t}, \mathbf{r}, \boldsymbol{\theta}).$$

To simplify some notations, we will also use the *restricted complete data log-likelihood*

$$L_{CR}(\mathbf{z}, \mathbf{w}, \boldsymbol{\theta}) = \sum_{i,j,k,\ell} z_{ik} w_{j\ell} \log f_{k\ell}(x_{ij}; \boldsymbol{\alpha}) \quad (8)$$

and the *restricted fuzzy complete-data log-likelihood*

$$L_{CR}(\mathbf{t}, \mathbf{r}; \boldsymbol{\theta}) = \sum_{i,j,k,\ell} t_{ik} r_{j\ell} \log f_{k\ell}(x_{ij}; \boldsymbol{\alpha}). \quad (9)$$

## 2.4. Algorithms

Here we give brief outline of two algorithms used in our package.

### *Block expectation maximization (BEM) algorithm*

The fuzzy clustering criterion given in (6) can be maximized using EM algorithm. We here only outline the various expressions evaluated during E and M steps.

**E-Step:** Here we compute conditional row and column class probabilities which are respectively given by:

$$t_{ik} = \log \pi_k + \sum_{j\ell} r_{j\ell} \log f_{k\ell}(x_{ij}; \boldsymbol{\alpha}). \quad (10)$$

and

$$r_{j\ell} = \log \rho_\ell + \sum_{i,k} t_{ik} \log f_{k\ell}(x_{ij}; \boldsymbol{\alpha}) \quad (11)$$

**M-Step:** Here we calculate row proportions  $\boldsymbol{\pi}$  and column proportions  $\boldsymbol{\rho}$ : The maximization of  $\tilde{F}_C$  w.r.t.  $\boldsymbol{\pi}$ , and w.r.t  $\boldsymbol{\rho}$ , is obtained by maximizing  $\sum_k t_{.k} \log \pi_k$ , and  $\sum_\ell r_{.\ell} \log \rho_\ell$  respectively, which leads to  $\pi_k = \frac{t_{.k}}{n}$  and  $\rho_\ell = \frac{r_{.\ell}}{d}$ . Also the estimate of model parameters  $\boldsymbol{\alpha}$  will be obtained by maximizing

$$\sum_{i,j,k,\ell} t_{ik} r_{j\ell} \log f_{k\ell}(x_{ij}; \boldsymbol{\alpha}) \quad (12)$$

which will depend on the pdf's  $f_{k\ell}$  and shall be treated later separately in subsections discussing block mixture models for binary, contingency and continuous case.

Using the **E** and **M** steps defined above, **BEM** algorithm can be enumerated as follows:

1. Initialize  $\mathbf{t}^{(0)}, \mathbf{r}^{(0)}$  and  $\boldsymbol{\theta}^{(0)} = (\boldsymbol{\pi}^{(0)}, \boldsymbol{\rho}^{(0)}, \boldsymbol{\alpha}^{(0)})$ .
2. Compute  $\mathbf{t}^{(c+1)}, \boldsymbol{\pi}^{(c+1)}, \boldsymbol{\alpha}^{(c+(1/2))}$  by using EM algorithm for the data matrix  $u_{i\ell} = \sum_j r_{j\ell} x_{ij}$  and starting from  $\mathbf{t}^{(c)}, \boldsymbol{\pi}^{(c)}, \boldsymbol{\alpha}^{(c)}$ .
3. Compute  $\mathbf{r}^{(c+1)}, \boldsymbol{\rho}^{(c+1)}, \boldsymbol{\alpha}^{(c+1)}$  by using EM algorithm for the data matrix  $v_{jk} = \sum_i t_{ik} x_{ij}$  and starting from  $\mathbf{r}^{(c)}, \boldsymbol{\rho}^{(c)}, \boldsymbol{\alpha}^{(c+(1/2))}$ .
4. Iterate steps (2) and (3) until convergence.

### *Block classification expectation maximization (BCEM) Algorithm*

To apply the complete maximum likelihood (CML) approach for parameter estimation of the latent block model (4), the partitions  $\mathbf{z}$  and  $\mathbf{w}$  are added to the parameters to be estimated, and the objective is to maximize complete data log-likelihood associated with the latent block model (4). Unlike the BEM situation, this maximization does not require an approximation and can be done, for instance, by maximizing with fixed  $\mathbf{z}$  and  $\boldsymbol{\rho}$  and then with fixed  $\mathbf{w}$  and  $\boldsymbol{\pi}$ . This algorithm, can be seen as a variant of **BEM** algorithm described above and to perform this algorithm, it is sufficient to add a C-step in each of the EM algorithms, which converts the  $t_{ik}$ 's and  $r_{j\ell}$ 's to a discrete classification before performing the M-step by assigning each individual and variable to the cluster having the highest posterior probability.

## 2.5. Block mixture models for binary, contingency and continuous datasets

### Binary data

The parameters  $\boldsymbol{\alpha}$  of the underlying distribution of a binary data set is given by the matrix  $\mathbf{p} = (p_{k\ell})$  where  $p_{k\ell} \in ]0, 1[ \forall k = 1, \dots, g$  and  $\ell = 1, \dots, m$  and the probability distribution  $f_{k\ell}(x_{ij}; \mathbf{p}) = f(x_{ij}; p_{k\ell})$  is the Bernoulli distribution

$$f(x_{ij}; p_{k\ell}) = (p_{k\ell})^{x_{ij}}(1 - p_{k\ell})^{1-x_{ij}}.$$

we re-parameterize the model density as follows:

$$f_{k\ell}(x_{ij}; \boldsymbol{\alpha}) = (\varepsilon_{k\ell})^{|x_{ij}-a_{k\ell}|}(1 - \varepsilon_{k\ell})^{1-|x_{ij}-a_{k\ell}|}$$

where

$$\begin{cases} a_{k\ell} = 0, \varepsilon_{k\ell} = p_{k\ell} & \text{if } p_{k\ell} < 0.5 \\ a_{k\ell} = 1, \varepsilon_{k\ell} = 1 - p_{k\ell} & \text{if } p_{k\ell} > 0.5. \end{cases}$$

Hence the parameters  $p_{k\ell}$  of the Bernoulli mixture model are replaced by the following parameters:

- The binary value  $a_{k\ell}$ , which acts as the center of the block  $k, \ell$  and which gives, for each block, the most frequent binary value,
- The value  $\varepsilon_{k\ell}$  belonging to the set  $]0, 1/2[$  that characterizes the dispersion of the block  $k, \ell$  and which is, for each block, represents the probability of having a different value than the center.

For this model, using Equation 10 and 11 the row and column class conditional probabilities are respectively given by:

$$t_{ik} = \log \pi_k + \sum_{\ell} |u_{i\ell} - r_{.\ell} a_{k\ell}| \log \frac{\varepsilon_{k\ell}}{1 - \varepsilon_{k\ell}} + \sum_{\ell} r_{.\ell} \log (1 - \varepsilon_{k\ell})$$

and

$$r_{j\ell} = \log \rho_{\ell} + \sum_k |v_{jk} - t_{.k} a_{k\ell}| \log \frac{\varepsilon_{k\ell}}{1 - \varepsilon_{k\ell}} + \sum_k t_{.k} \log (1 - \varepsilon_{k\ell})$$

where  $u_{i\ell} = \sum_j r_{j\ell} x_{ij}$  and  $v_{jk} = \sum_i t_{ik} x_{ij}$ .

Furthermore, we can write the restrained data log-likelihood for this model in terms of previously defined explicit parameters as

$$L_{CR}(\mathbf{t}, \mathbf{r}, \boldsymbol{\alpha}) = \sum_{k\ell} \left( |y_{k\ell} - t_{.k} r_{.\ell} a_{k\ell}| \log \frac{\varepsilon_{k\ell}}{1 - \varepsilon_{k\ell}} + t_{.k} r_{.\ell} \log (1 - \varepsilon_{k\ell}) \right) \quad (13)$$

where  $y_{k\ell} = \sum_{i,j} t_{ik} r_{j\ell} x_{ij}$ . It can be easily shown that the values of model parameters maximizing Equation 13 is given by

$$\begin{cases} a_{k\ell} = 0, \varepsilon_{k\ell} = \frac{y_{k\ell}}{t_{.k} r_{.\ell}} & \text{if } \frac{y_{k\ell}}{t_{.k} r_{.\ell}} < 0.5 \\ a_{k\ell} = 1, \varepsilon_{k\ell} = 1 - \frac{y_{k\ell}}{t_{.k} r_{.\ell}} & \text{otherwise.} \end{cases}$$



For detailed discussion on Bernoulli latent block models we would kindly refer our readers to [Govaert and Nadif (2008)].

### Contingency data

To apply the block latent mixture model on contingency data, it is assumed that for each block  $k, \ell$ , the values  $x_{ij}$  are distributed according to Poisson distribution  $\mathcal{P}(\mu_i \nu_j \gamma_{k\ell})$  where the Poisson parameter is split into  $\mu_i$  and  $\nu_j$  the effects of the row  $i$  and the column  $j$  respectively and  $\gamma_{k\ell}$  the effect of the block  $k\ell$  [Govaert and Nadif (2010)]. Then, we have

$$f_{k\ell}(x_{ij}; \boldsymbol{\alpha}) = \frac{e^{-\mu_i \nu_j \gamma_{k\ell}} (\mu_i \nu_j \gamma_{k\ell})^{x_{ij}}}{x_{ij}!}$$

where  $\boldsymbol{\alpha} = (\boldsymbol{\mu}, \boldsymbol{\nu}, \boldsymbol{\gamma})$  with  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)$ ,  $\boldsymbol{\nu} = (\nu_1, \dots, \nu_d)$  and  $\boldsymbol{\gamma} = (\gamma_{11}, \dots, \gamma_{gm})$ .

Unfortunately, this parameterization is not identifiable. It is therefore not possible to estimate simultaneously  $\mu_i$ , the  $\nu_j$  and  $\gamma_{k\ell}$ . The following two solutions are considered to tackle this problem:

- $\mu_i$  and  $\nu_j$  are known: In this case, it only remains to estimate the matrix  $\gamma_{k\ell}$ .
- Constraints are imposed on the parameter  $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{\rho}, \boldsymbol{\mu}, \boldsymbol{\nu}, \boldsymbol{\gamma})$ : In this situation, following Constraints

$$\sum_k \pi_k \gamma_{kl} = \sum_\ell \rho_\ell \gamma_{kl} = 1, \quad \sum_i \mu_i = 1, \quad \sum_j \nu_j = 1 \quad (14)$$

appear to make the model identifiable, but it is yet remains to prove it.

The transformation of the parameter  $\boldsymbol{\theta}$  on a parameter  $\boldsymbol{\theta}'$  defined as follows

$$\begin{aligned} \pi'_k &= \pi_k & \text{and} & & \rho'_\ell &= \rho_\ell \\ \mu'_i &= \mu_i N & \text{and} & & \nu'_j &= \nu_j N \\ \gamma'_{k\ell} &= \frac{\gamma_{k\ell}}{N}, \end{aligned}$$

shows that it is possible to parameterize the model with the constraints

$$\sum_k \pi'_k \gamma'_{kl} = \sum_\ell \rho'_\ell \gamma'_{kl} = 1/N, \quad \sum_i \mu'_i = \sum_j \nu'_j = N \quad (15)$$

where  $N$  is any positive real.

### Case-1 : The row and column effects are known

In this situation,  $\boldsymbol{\alpha} = \boldsymbol{\gamma} = (\gamma_{11}, \dots, \gamma_{gm})$  and  $f_{k\ell}(x_{ij}; \boldsymbol{\gamma})$  can be written as

$$f_{k\ell}(x_{ij}; \boldsymbol{\gamma}) = K_{ij} \cdot g_{k\ell}(x_{ij}; \boldsymbol{\gamma})$$

where

$$g_{k\ell}(x_{ij}; \boldsymbol{\gamma}) = e^{-\mu_i \nu_j \gamma_{k\ell}} (\gamma_{k\ell})^{x_{ij}} \quad \text{and} \quad K_{ij} = \frac{(\mu_i \nu_j)^{x_{ij}}}{x_{ij}!}.$$

Then, we can replace the function  $f_{k\ell}$  by the function  $g_{k\ell}$  in the log-likelihoods and in the expressions  $t_{ik}$  and  $r_{j\ell}$ . Hence, using Equation 10 and 11 the row and column class conditional probabilities are respectively given by:

$$t_{ik} = \log \pi_k + \sum_{j\ell} r_{j\ell} x_{ij} \log \gamma_{k\ell} - \mu_i \sum_{\ell} \gamma_{k\ell} \nu_{\ell}$$

$$r_{j\ell} = \log \rho_{\ell} + \sum_{i,k} t_{ik} x_{ij} \log \gamma_{k\ell} - \nu_j \sum_k \gamma_{k\ell} \mu_k.$$

Also the restrained complete data log-likelihood is given by:

$$L_{CR}(\mathbf{t}, \mathbf{r}, \boldsymbol{\alpha}) = \sum_{k,\ell} y_{k\ell} \log \gamma_{k\ell} - \mu_k \nu_{\ell} \gamma_{k\ell}. \quad (16)$$

where  $y_{k\ell} = \sum_{ij} t_{ik} r_{j\ell} x_{ij}$ . The maximization of Equation 16 gives:

$$\gamma_{k\ell} = \frac{y_{k\ell}}{\mu_k \nu_{\ell}} \quad \forall k, \ell$$

**Case-2** : *The row and column effects are unknown*

Using the constraints given by Equation 15, we obtain  $E(x_{i.}) = \mu_i$  and  $E(x_{.j}) = \nu_j$ . We can propose as an estimator of  $\mu_i$  and  $\nu_j$  margins  $x_{i.}$  and  $x_{.j}$  and take the results of **case-1** to calculate row and column class conditional probabilities by replacing  $\mu_i$  and  $\nu_j$  by  $x_{i.}$  and  $x_{.j}$  respectively. Hence the row and column class conditional probabilities are respectively given by:

$$t_{ik} = \log \pi_k + \sum_{j\ell} r_{j\ell} x_{ij} \log \gamma_{k\ell} - \mathbf{x}_{i.} \sum_{\ell} \gamma_{k\ell} y_{\ell}$$

$$r_{j\ell} = \log \rho_{\ell} + \sum_{i,k} t_{ik} x_{ij} \log \gamma_{k\ell} - \mathbf{x}_{.j} \sum_k \gamma_{k\ell} y_k.$$

Furthermore the derivation of  $\gamma_{k\ell}$  is exactly same as in previous case. Using this particular form of  $\gamma_{k\ell}$ , the calculation of  $t_{ik}$  and  $r_{j\ell}$  is simplified and becomes

$$t_{ik} = \log \pi_k + \sum_{j\ell} r_{j\ell} x_{ij} \log \gamma_{k\ell} - \mathbf{x}_{i.}$$

$$r_{j\ell} = \log \rho_{\ell} + \sum_{i,k} t_{ik} x_{ij} \log \gamma_{k\ell} - \mathbf{x}_{.j}$$

*Continuous data*

In this case, the continuous data is modeled using unidimensional normal distribution. Hence the density for each block is given by:

$$f_{k\ell}(x_{ij}; \boldsymbol{\alpha}) = \frac{1}{\sqrt{2\pi\sigma_{k\ell}^2}} \exp -\left\{ \frac{1}{2\sigma_{k\ell}^2} (x_{ij} - \mu_{k\ell})^2 \right\}$$

The parameters of the model is  $\boldsymbol{\alpha} = (\boldsymbol{\alpha}_{11}, \dots, \boldsymbol{\alpha}_{gm})$  where  $\boldsymbol{\alpha}_{k\ell} = (\mu_{k\ell}, \sigma_{k\ell}^2)$ . For this model, using Equation 10 and 11 the row and column class conditional probabilities are respectively given by:

$$t_{ik} = \log \pi_k - \frac{1}{2} \sum_{\ell} (r_{.\ell} \log \sigma_{k\ell}^2 + \frac{1}{\sigma_{k\ell}^2} a_{ik\ell})$$

where

$$a_{ik\ell} = \sum_j r_{j\ell} (x_{ij} - u_{i\ell}/r_{.\ell})^2 + r_{.\ell} \sum_{\ell} \left( \frac{u_{i\ell}}{r_{.\ell}} - \mu_{k\ell} \right)^2$$

and

$$r_{j\ell} = \log \rho_{\ell} - \frac{1}{2} \sum_k (t_{.k} \log \sigma_{k\ell}^2 + \frac{1}{\sigma_{k\ell}^2} b_{jk\ell})$$

where

$$b_{jk\ell} = \sum_i t_{ik} (x_{ij} - v_{jk}/t_{.k})^2 + t_{.k} \sum_k (v_{jk} - \mu_{k\ell})^2$$

The computation of parameters  $\boldsymbol{\alpha}$  is given by maximization of restrained log-likelihood which in this case is given by:

$$L_{CR}(\mathbf{t}, \mathbf{r}, \boldsymbol{\alpha}) = -\frac{1}{2} \sum_{k,\ell} (t_{.k} r_{.\ell} \log \sigma_{k\ell}^2 + \frac{1}{\sigma_{k\ell}^2} \sum_{i,j} t_{ik} r_{j\ell} (x_{ij} - \mu_{k\ell})^2). \quad (17)$$

It can be deduced easily that the value of  $\mu_{k\ell}$  and  $\sigma_{k\ell}^2$  is given by minimization of  $\sum_{i,j} t_{ik} r_{j\ell} (x_{ij} - \mu_{k\ell})^2$  and  $t_{.k} r_{.\ell} \log \sigma_{k\ell}^2 + \frac{1}{\sigma_{k\ell}^2} \sum_{i,j} t_{ik} r_{j\ell} (x_{ij} - \mu_{k\ell})^2 \quad \forall k, \ell$  respectively.

Hence we obtain:

$$\mu_{k\ell} = \frac{\sum_j r_{j\ell} v_{jk}}{t_{.k} r_{.\ell}}, \quad \sigma_{k\ell}^2 = \frac{\sum_{i,j} t_{ik} r_{j\ell} x_{ij}^2}{t_{.k} r_{.\ell}} - \mu_{k\ell}^2$$

### 3. Block clustering with package blockcluster

**blockcluster** is a newly developed R package for co-clustering of binary, contingency and continuous data. The core library is written in C++ and **blockcluster** API acts as a bridge between C++ core library and R statistical computing environment. A brief introduction of the C++ library is given in Appendix A. In the subsequent sections, we shall provide detailed explanation of the various functionalities and utilities of the package with help of examples wherever possible.

#### 3.1. Strategy for parameters estimation

Correct parameters estimation in EM based algorithms has always been a challenging problem due to various local optima. In case of co-clustering algorithms, this problem could be even more apparent and hence the final results shall be directly affected by the choice of a strategy. As being explained in previous sections, we have two block clustering algorithms namely **BEM** and **BCEM** and the way we run these algorithms in our library is named respectively as '**XEMStrategy**' and '**XCEMStrategy**'. Both the strategies are equivalent except that they run different algorithms. The various steps involved in '**XEMStrategy**' (equivalently in '**XCEMStrategy**') are:

1. Set the pseudo log-likelihood  $\mathbf{L}_{XEM}$  to minimum possible value,
2. Run the **BEM** algorithm for some pre-defined number of times (we denote this number by '**nbxem**' in our package) as follows:
  - (a) Set the pseudo log-likelihood  $\mathbf{L}_{xem}$  to minimum possible value,
  - (b) Initialize the model parameters,
  - (c) Run the algorithm with a high tolerance (epsilon) and few iterations using the initialized parameters, and calculate the current pseudo log-likelihood  $\mathbf{L}_{current}$ ,
  - (d) If  $\mathbf{L}_{current}$  is greater than  $\mathbf{L}_{xem}$ , copy the currently estimated parameters to  $\alpha_{start}$  and set  $\mathbf{L}_{xem}$  to  $\mathbf{L}_{current}$ .
3. Starting with the parameters  $\alpha_{start}$ , run the algorithm with a low value of epsilon (low tolerance) and a high number of iterations and calculate the current pseudo log-likelihood  $\mathbf{L}_{current}$ .
4. If  $\mathbf{L}_{current}$  is greater than  $\mathbf{L}_{XEM}$ , copy the currently estimated parameters to  $\alpha_{final}$  and set  $\mathbf{L}_{XEM}$  to  $\mathbf{L}_{current}$ .
5. Repeat step 2 through 3 for some pre-set times (we denote this number by '**nbtry**' in our package) in order to obtain a good estimation of the parameters.

Although the above strategy does not explore all the possible solutions and hence cannot ensure to obtain global optimum, it is still a very promising approach in practice. The tuning of the values of '**nbxem**' and '**nbtry**' need to be done intuitively, and could have a substantial effect on final results. A good way to set these values is to run co-clustering few number of times and check if final log-likelihood is stable. If not, one may need to increase these values. In practice, it is better to increment '**nbxem**' as it could lead to better (stable) results without compromising too much the running time.

In the package **blockcluster**, we have a function called **cocluststrategy**<sup>1</sup> which allow to set the values of various input parameters including '**nbxem**' and '**nbtry**'. In the following example, we get an object **defaultstrategy** of S4 class **strategy** by calling the function **cocluststrategy** without any arguments and then we called the overloaded function **summary** to see various default values.

```
R > defaultstrategy <- cocluststrategy()
R > summary(defaultstrategy)
```

```
*****
```

```
Algorithm: XEMStrategy
```

```
Initialization method(There is no default value):
```

```
Stopping Criteria: Parameter
```

#### Various Iterations

---

<sup>1</sup>This function is also use to set all other input parameters like number of iterations and tolerance values. They are not described here explicitly for the sake of brevity. We refer our readers to the package documentation for details.

```

*****
Number of global iterations while running initialization: 10
Number of iterations for internal E-step: 5
Number of EM iterations used during xem: 50
Number of EM iterations used during XEM: 500
Number of xem iterations: 5
Number of tries: 2

Various epsilons
*****
Tolerance value used while initialization: 0.01
Tolerance value for internal E-step: 0.01
Tolerance value used during xem: 1e-04
Tolerance value used during XEM: 1e-10
*****

```

One thing which is worth noting in the summary output (above) is that there is no default value for initialization method. It will be set automatically depending on the type of input model as explained in Section 3.2. Any of these input parameters can be set by passing appropriate arguments to function `cocluststrategy` as shown in example below where we set `nbtry`, `nbxem` and `algo` parameters.

```
R > newstrategy <- cocluststrategy(nbtry=5, nbxem=10, algo='XCEMstrategy')
```

### 3.2. Model initialization

Initializing model parameters is arguably the most important step in getting the good final estimate of model parameters. The three most important initialization strategies for EM algorithms are:

- **Random:** In this method, the model parameters are initialized randomly using the input data. One may go with several random initializations and chose the one which gives the best log-likelihood.
- **CEM:** The model parameters are first initialized using random values by making use of input data, and then CEM algorithm [Celeux and Govaert (1992)] is run for small number of iterations with high tolerance. The **CEM** algorithm optimizes directly the completed log-likelihood by replacing the M step by a classification step based on MAP. In general, the **CEM** algorithm converges very quickly (not necessarily to good optimum) which is an advantage as we spend few time in initialization compared to the overall run time.
- **FuzzyCEM (Stochastic EM):** This method is equivalent to **CEM** except for the fact that instead of classifying samples into various clusters based on MAP, we stochastically select the class for each sample using the corresponding class probability distribution.

Unfortunately, in **blockcluster**, we do not have the freedom to initialize any model with all of the methods above. Instead, for every model we only have one possible initialization method

Model	Datatype	Proportions	Dispersion/Variance	Initialization
<b>pik_rhol_epsilonkl</b>	binary	unequal	unequal	CEM
<b>pik_rhol_epsilon</b>	binary	unequal	equal	CEM
<b>pi_rho_epsilonkl</b>	binary	equal	unequal	CEM
<b>pi_rho_epsilon</b>	binary	equal	equal	CEM
<b>pik_rhol_sigma2kl</b>	continuous	unequal	unequal	CEM
<b>pik_rhol_sigma</b>	continuous	unequal	equal	CEM
<b>pi_rho_sigma2kl</b>	continuous	equal	unequal	CEM
<b>pi_rho_sigma2</b>	continuous	equal	equal	CEM
<b>pik_rhol_unknown</b>	contingency	unequal	N.A	CEM
<b>pi_rho_unknown</b>	contingency	equal	N.A	CEM
<b>pik_rhol_known</b>	contingency	unequal	N.A	Random
<b>pi_rho_known</b>	contingency	equal	N.A	Random

Table 1: Various models available in package **blockcluster**.

at present (mainly **CEM**) as shown in Table 1. Hence the initialization strategy will be selected automatically by the library based on the input model. This is the reason why we have no default initialization strategy. This is indeed on the priority list of our future works to provide as many initialization strategies as possible for various kinds of models.

### 3.3. Examples with simulated datasets

#### *Binary Data*

We have simulated binary data using the data generation process explained in Section 2 and with parameters given in Table 2. The data consist of 1000 rows (samples) and 100 columns (variables) with two clusters on rows and three clusters on columns.

<b>a, <math>\epsilon</math></b>	0, 0.1	0, 0.3	1, 0.1	<b><math>\pi</math></b>	.6	.4
	1, 0.3	1, 0.2	0, 0.1	<b><math>\rho</math></b>	.3	.3 .4

Table 2: Parameters for simulation of binary data.

The following R commands shows how to load the library, process the data and visualize/-summarize results using **blockcluster**.

```
R > library("blockcluster")
R > data("binarydata")
R > out<-cocluster(binarydata, datatype = "binary", nbcocluster=c(2,3))
R > summary(out)
```

```
*****
Model Family : Bernoulli Latent block model
Model Name : pik_rhol_epsilonkl
```

```
Model Parameters..
```

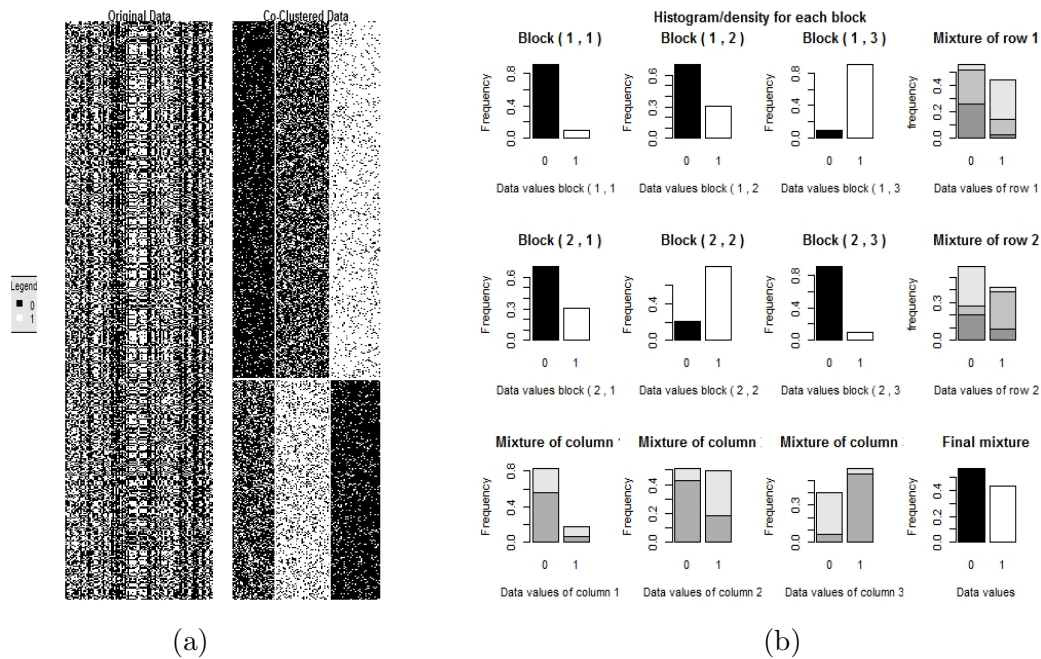


Figure 2: Original and co-clustered binary data (a), and distributions for each block along with various mixture densities (b).

Class Mean:

	[,1]	[,2]	[,3]
[1,]	0	0	1
[2,]	0	1	0

Class Dispersion:

	[,1]	[,2]	[,3]
[1,]	0.09798013	0.3022391	0.1011803
[2,]	0.30176929	0.2003679	0.1006313

Row proportions: 0.618 0.382

Column proportions: 0.29 0.37 0.34

Likelihood: -0.4552042

\*\*\*\*\*

The following R command is used to plot the original and coclustered data (Figure 2(a)) with aspect ratio set to false (it is true by default).

```
R > plot(out, asp = 0)
```

To Plot the various block distributions (Figure 2(b)), the following R command is used with **type** set to 'distribution' (**type** is 'cocluster' by default).

```
R > plot(out, type = 'distribution')
```

## Contingency data

In this case, we simulated contingency data as explained earlier in Section 2 and using the parameters given in Table 3. Again the data dimensions are  $1000 \times 100$  with two clusters on rows and three clusters on columns.

$\gamma(*1e-6)$	1.290	.450	1.260	$\pi$	.6	.4
	.740	1.550	0.710	$\rho$	.4	.3

Table 3: Parameters for simulation of contingency data.

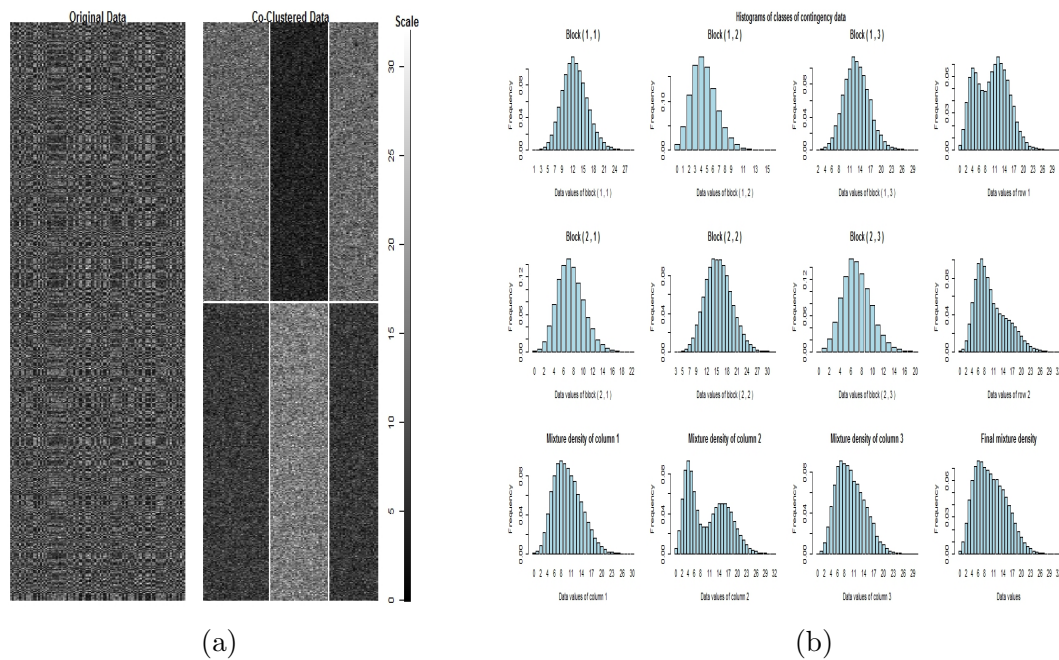


Figure 3: Original and co-clustered contingency data (a), and distributions for each block along with various mixture densities (b).

The following R commands load the simulated contingency data into R environment, co-clusters the contingency data, summarize the results and shows the graphics respectively.

```
R > data("contingencydataunknown")
R > out<-cocluster(contingencydataunknown, datatype = "contingency",
+ nbcocluster=c(2,3))
R > summary(out)
R > plot(out)
R > plot(out, type = 'distribution')
```

```
*****
Model Family : Poisson Latent block model
Model Name : pik_rho1_unknown
```



Model Parameters..

Normalized Class Gamma:

```

          [,1]      [,2]      [,3]
[1,] 1.276872e-06 4.429059e-07 1.311424e-06
[2,] 7.404569e-07 1.512883e-06 7.084545e-07
    
```

Row proportions: 0.484 0.516

Column proportions: 0.3686402 0.34 0.2913598

Likelihood: -14762080

\*\*\*\*\*

*Continuous data*

In this final example, we generated continuous data using parameters given in Table 4. The data-set again contains 1000 rows and 100 columns with two clusters on rows and three clusters on columns.

$\mu, \sigma^2$	-10, 20	0, 10	10, 20	$\pi$	.6	.4
	10, 10	0, 20	-10, 10	$\rho$	.3	.3 .4

Table 4: Parameters for simulation of continuous data.

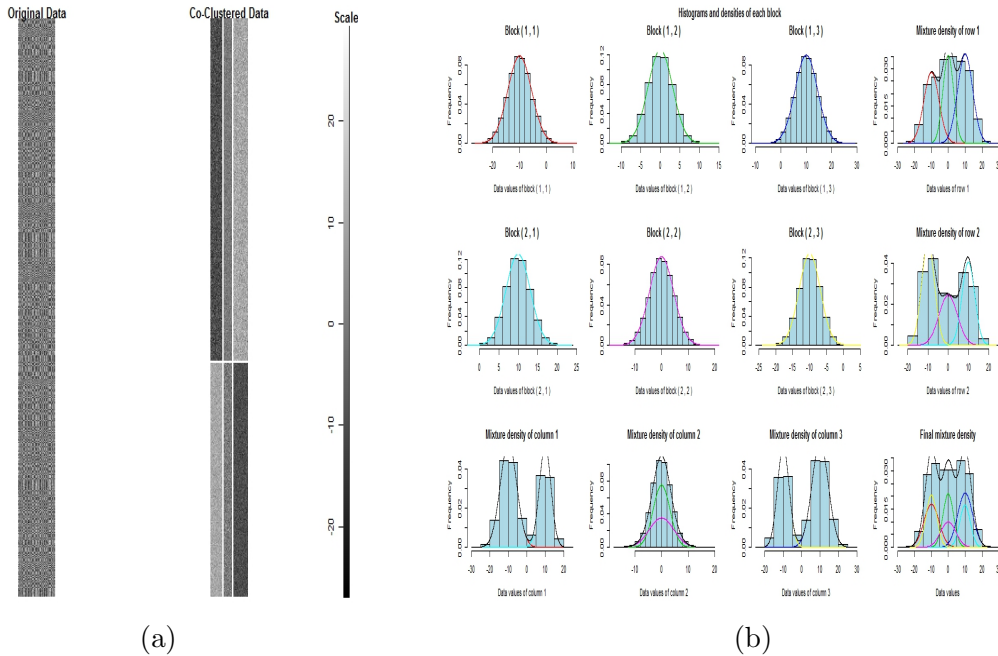


Figure 4: Original and co-clustered continuous data (a), and distributions for each block along with various mixture densities (b).

The following R commands performs co-clustering on simulated continuous data-set and shows the results.

```

R > data("gaussiandata")
R > out<-cocluster(gaussiandata,datatype="continuous",nbcoccluster=c(2,3))
R > summary(out)
R > plot(out, asp = 0)
R > plot(out, type = 'distribution')

*****
Model Family : Gaussian Latent block model
Model Name : pik_rho1_sigma2k1
Model Parameters:
Class Mean:
      [,1]      [,2]      [,3]
[1,] -9.966978  0.0008994504 10.031539
[2,]  9.933649 -0.0113171870 -9.975353

Class Variance:
      [,1]      [,2]      [,3]
[1,] 19.958878 10.00552 19.79938
[2,]  9.854607 20.79533 10.02493

Row proportions:  0.594 0.406
Column proportions:  0.32 0.28 0.4
Likelihood:  -1.846875
*****

```

## 4. Examples with real datasets

In what follows, we shall give the outcome of running package **blockcluster** on real datasets.

### 4.1. Image segmentation

Automatic image segmentation is an important technique and have numerous application especially in fields of Medical imaging. Here we present an interesting application of co-clustering (as pre-processing step) for segmenting object(s) in image. We assume that the object pixels follows Gaussian distribution. Hence we ran the **blockcluster** package with Gaussian Family model 'pik\_rho1\_sigma2k1' on image shown in Figure 5. It can be clearly seen that we are nicely able to segment the snake and insect in two different blocks. Co-clustering have not been exploited in this context so far to the best of our knowledge. Hence it seems they offer a new opportunity for co-clustering.

### 4.2. Document clustering

Document clustering is yet another data mining technique where co-clustering seems to be very useful as demonstrated in [Dhillon (2001)]. Here we run our package on one of the datasets being used in [Dhillon (2001)] which is publicly available at <ftp://ftp.cs.cornell.edu/pub/smart>. We mix Medline (1033 medical abstracts) and Cranfield (1398 aeronautical

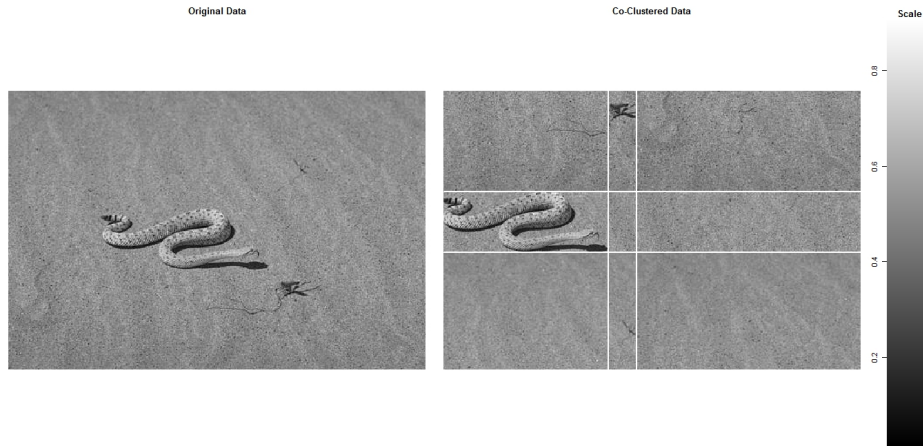


Figure 5: Original and co-clustered (segmented) image.

abstracts) making a total of 2431 documents. Furthermore, we used all the words (excluding stop words) as features making a total of 9275 unique words. The data matrix consist of words on the rows and documents on the columns with each entry giving the term frequency, that is the number of occurrences of corresponding word in corresponding document. To run our package, we assume that the term frequency follows Poisson distribution. Hence we can apply the model 'pik\_rho1\_unknown' available in our package for contingency (Poisson Family) datasets with unknown row and column effects. Table 5 shows the confusion matrix and compare our results with classical bipartite spectral graph partitioning algorithm of [Dhillon (2001)] where we have obtained 100 percent correct classification.

Figure 6 depicts the  $2 \times 2$  checkerboard pattern in the data matrix, hence confirming the more frequent occurrence of particular set of words in one document and vice-versa. Please note that the data matrix images are extremely sparse (data points almost invisible) and have been processed using simple image processing tools for visualization purpose only.

	Medline	Cranfield
Medline	1026	0
Cranfield	7	1400

(a)

	Medline	Cranfield
Medline	1033	0
Cranfield	0	1398

(b)

Table 5: Confusion Matrix: Results reported in [Dhillon (2001)] (a), and Results using **blockcluster** (b). The difference in number of Cranfield documents is because we made use of the readymade data extracted from the documents and there are two less documents data in the same.

## 5. Conclusion and future directions

This article presents our newly developed R package for co-clustering algorithms based on latent block models. In the first part of article, we have presented brief mathematical foundation of co-clustering algorithms based on latent block models. In the second part, we have given details on the strategy adopted to run these algorithms. Furthermore, we have completed

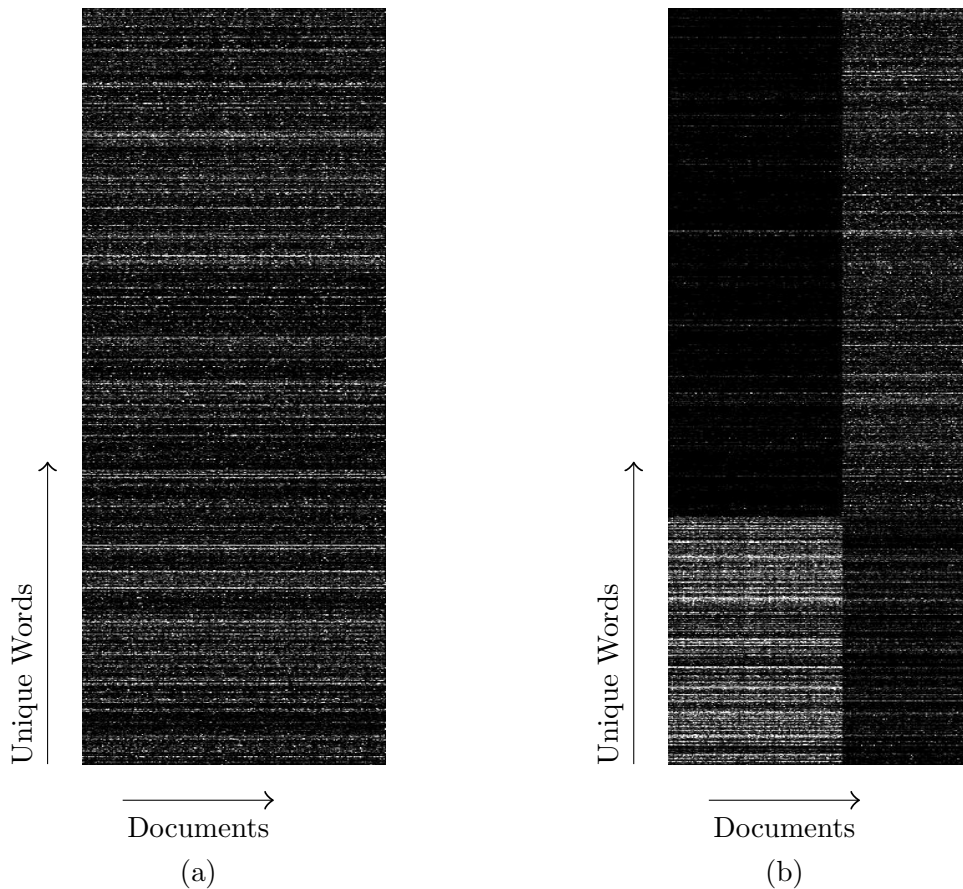


Figure 6: Original data matrix with words on rows and documents on columns (a), and checkerboard pattern in words by documents matrix obtained after performing co-clustering (b).

the analysis of **blockcluster** package by providing the examples on simulated datasets for each data-type. Finally we have shown applicability of our software on real datasets. As the package facilitates co-clustering for all types of data, it seems to provide a generic framework and a one stop solution for wide variety of co-clustering applications.

In the future release(s) of **blockcluster** package, we have planned to provide parallel implementations of the software by exploiting various data and task parallelism in latent block model estimation. Furthermore, we are also working on new initialization strategies as well as model selection criterion that should be part of upcoming releases of the package. Finally we have plans to extend the library by providing Stochastic EM algorithms for parameter estimation as well as to implement semi-supervised co-clustering in which some row and/or column classes are provided by the user.

## 6. Acknowledement

The authors are grateful to Christophe Biernacki for his leading support and constant motivation throughout the project. We are also thankful to Gilles Celeux for his contributions

and providing his support in proposing the project Plan to INRIA-Lille. Furthermore we are thankful to Alhousseiny Pam who worked partly on the graphical part of the R package. Finally we are grateful to INRIA-Lille for funding this project.

## Appendices

### A. The core library: Structure and usage

In this paper we have given detailed analysis on the usage of R package **blockcluster** to estimate parameters of latent block models, but behind the scenes the kernel (computational) part lies in the C++ core library<sup>2</sup>. Our purpose to put this section is only to briefly introduce the kernel part and partly for developers who might be interested in integrating the library with their own codes, but we highly recommend to use only the R package **blockcluster** for research purposes as lot of utilities are absent in the C++ codes for instance the graphics part.

#### A.1. Library structure

The C++ library is based on famous Strategy design pattern [Freeman, Robson, Bates, and Sierra (2004)] and (not so complete) architecture is revealed in Figure 7. The choice of the architecture is motivated from the fact that it allows easy extension to the library (accommodating any future research) with as little changes as possible in the existing codes. Hence it follows famous design principle "Open for extension but closed for modification" [Freeman *et al.* (2004)].

#### A.2. Usage

In our library, we have provided a simple main function (to illustrate its usage) that can be run from console as follows.

```
$ coclustmain -f path/to/optionfile.txt
```

**coclustmain** is binary executable file which take path to a text file as command line argument. A sample input text file is shown below.

```
[InputOptions]
DataType = Binary
nbrowclust = 2
nbcoclust = 3
DataFileName = Data/BinaryData/testbinarypikrholepsk1.dat
nbnititerations = 10
initedepsilon = .01
nbtry = 1
nbxem = 1
epsilon_xemstart = .00001
```

---

<sup>2</sup>The source of the C++ library can be checkout anonymously from the repository [svn://scm.gforge.inria.fr/svnroot/cocluster](http://svn.inria.fr/svnroot/cocluster)

```

epsilon_xem = .00000001
nbiterations_xemstart = 50
nbiterations_xem = 50
epsilon_int = 0.01
nbiterations_int = 5
Algorithm = XEMStrategy
StopCriteria = Parameter
Initialization = CEMInit
ModelName = pik_rhol_epsilonkl

```

This file contains all the input options needed to run the library. The only mandatory inputs are the data-type, number of row/column clusters and the path to data file. As the reader may have already guessed, this file is more or less similar to **blockcluster** package function `cocluststrategy()` where the optional parameters in the sample text file are equivalent to arguments of the function. By looking at the **coclustmain** source file, one can easily understand how to integrate the core library with their own C++ codes.

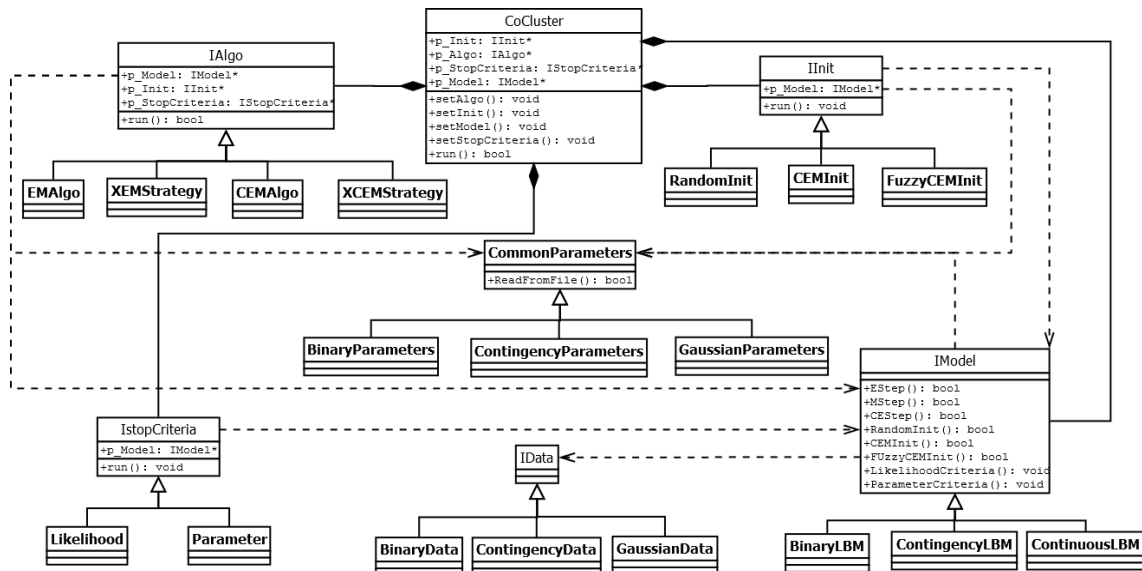


Figure 7: C++ core library software architecture based on strategy design pattern

## References

- Barkow S, Bleuler S, Prelic A, Zimmermann P, Zitzler E (2006). “BicAT: A Biclustering Analysis Toolbox.” *Bioinformatics*, **22**, 1282–1283.
- Celex G, Govaert G (1992). “A Classification EM Algorithm for Clustering and Two Stochastic Versions.” *Computational Statistics and Data Analysis*, **14**(3), 315–332.
- Dhillon IS (2001). “Co-Clustering Documents and Words Using Bipartite Spectral Graph Partitioning.” In *Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining*, KDD '01, pp. 269–274. ACM, New York, NY,

- USA. ISBN 1-58113-391-X. doi:10.1145/502512.502550. URL <http://doi.acm.org/10.1145/502512.502550>.
- Freeman E, Robson E, Bates B, Sierra K (2004). *Head First Design Patterns*. O'Reilly Media, Incorporated.
- Giannakidou E, Koutsonikola V, Vakali A, Kompatsiaris Y (2008). "Co-Clustering Tags and Social Data Sources." In *Web-Age Information Management, 2008. WAIM'08. The Ninth International Conference on*, pp. 317–324. IEEE.
- Goncalves J, Madeira S, Oliveira A (2009). "BiGGES TS: Integrated Environment for Biclustering Analysis of Time Series Gene Expression Data." *BMC Research Notes*, **2**(1), 124. ISSN 1756-0500. doi:10.1186/1756-0500-2-124. URL <http://www.biomedcentral.com/1756-0500/2/124>.
- Govaert G, Nadif M (2003). "Clustering with Block Mixture Models." *Pattern Recognition*, **36**(2), 463 – 473. ISSN 0031-3203. URL <http://www.sciencedirect.com/science/article/pii/S0031320302000742>.
- Govaert G, Nadif M (2008). "Block Clustering with Bernoulli Mixture Models: Comparison of Different Approaches." *Computational Statistics & Data Analysis*, **52**(6), 3233–3245.
- Govaert G, Nadif M (2010). "Latent Block Model for Contingency Table." *Communications in Statistics - Theory and Methods*, **39**(3), 416–425.
- Guan J, Qiu G, Xue X (2005). "Spectral Images and Features Co-Clustering with Application to Content-Based Image Retrieval." In *Multimedia Signal Processing, 2005 IEEE 7th Workshop on*, pp. 1–4. IEEE.
- Hathaway R (1986). "Another Interpretation of the EM Algorithm for Mixture Distributions." *Statistics & Probability Letters*, **4**(2), 53–56.
- Kaiser S, Leisch F (2008). "A Toolbox for Bicluster Analysis in R." *Technical Report 28*, Department of Statistics: Technical Reports.
- Kluger Y, Basri R, Chang J, Gerstein M (2003). "Spectral Biclustering of Microarray Data: Coclustering Genes and Conditions." *Genome research*, **13**(4), 703–716.
- Madeira S, Oliveira A (2004). "Biclustering Algorithms for Biological Data Analysis: A Survey." *IEEE/ACM Trans Comput Biol Bioinform*, **1**(1), 24–45. doi:10.1109/TCBB.2004.2.
- Neal R, Hinton G (1998). "A View of the EM Algorithm that Justifies Incremental, Sparse, and Other Variants." *NATO ASI SERIES D BEHAVIOURAL AND SOCIAL SCIENCES*, **89**, 355–370.
- Qiu G (2004). "Image and Feature Co-Clustering." In *Pattern Recognition, 2004. ICPR 2004. Proceedings of the 17th International Conference on*, volume 4, pp. 991–994. IEEE.
- Rasiwasia N, Vasconcelos N (2009). "Holistic Context Modeling Using Semantic Co-Occurrences." In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pp. 1889–1895. IEEE.

R Development Core Team (2008). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, URL <http://www.R-project.org>.

Xu G, Zong Y, Dolog P, Zhang Y (2010). “Co-Clustering Analysis of Weblogs Using Bipartite Spectral Projection Approach.” *Knowledge-Based and Intelligent Information and Engineering Systems*, pp. 398–407.

### **Affiliation:**

Parmeet Bhatia  
 MODAL Team – Inria Lille - Nord Europe  
 59655 Villeneuve d’Ascq Cedex – France  
 E-mail: [parmeet.bhatia@inria.fr](mailto:parmeet.bhatia@inria.fr)

Serge Iovleff  
 Laboratoire Paul Painlevé – Université Lille 1 & CNRS  
 Modal Team – Inria Lille - Nord Europe  
 59655 Villeneuve d’Ascq Cedex – France  
 E-mail: [serge.iovleff@inria.fr](mailto:serge.iovleff@inria.fr)

G erard Govaert  
 Laboratoire Heudiasyc – Universit e de Technologie de Compi egne & CNRS  
 60205 Compi egne Cedex – France  
 E-mail: [gerard.govaert@utc.fr](mailto:gerard.govaert@utc.fr)