

# Package ‘ordinalClust’

November 19, 2020

**Type** Package

**Title** Ordinal Data Clustering, Co-Clustering and Classification

**Version** 1.3.5

**Date** 2020-10-26

**Author** Margot Selosse, Julien Jacques, Christophe Biernacki

**Maintainer** Margot Selosse <margot.selosse@gmail.com>

**Description** Ordinal data classification, clustering and co-clustering using model-based approach with the BOS distribution for ordinal data (Christophe Biernacki and Julien Jacques (2016) <doi:10.1007/s11222-015-9585-2>).

**License** GPL (>= 2)

**Imports** Rcpp (>= 0.12.11), methods

**LinkingTo** Rcpp, RcppProgress, RcppArmadillo, BH

**Suggests** knitr, rmarkdown, caret, ggplot2

**VignetteBuilder** knitr

**LazyData** true

**SystemRequirements** C++11

**Depends** R (>= 3.3)

**NeedsCompilation** yes

**Repository** CRAN

**Date/Publication** 2020-11-19 08:20:07 UTC

## R topics documented:

bosclassif . . . . .	2
bosclust . . . . .	4
boscoclust . . . . .	6
dataqol . . . . .	9
dataqol.classif . . . . .	9
Msimulated . . . . .	10
pejSim . . . . .	10

plot . . . . .	11
predict . . . . .	11
summary . . . . .	12

<b>Index</b>	<b>13</b>
--------------	-----------

---

bosclassif	<i>Function to perform a classification</i>
------------	---

---

### Description

This function performs a classification algorithm on a dataset with ordinal features, and a label variable that belongs to  $(1,2,\dots,kr)$ . The classification function provides two classification models. The first model, (chosen by the argument  $kc=0$ ), is a multivariate BOS model with the assumption that, conditional on the class of the observations, the features are independent. The second model is a parsimonious version of the first model. Parsimony is introduced by grouping the features into clusters (as in co-clustering) and assuming that the features of a cluster have a common distribution.

### Usage

```
bosclassif(x, y, idx_list=c(1), kr, kc=0, init, nbSEM, nbSEMBurn,
           nbindmini, m=0, percentRandomB=0)
```

### Arguments

<code>x</code>	Matrix made of ordinal data of dimension $N \times J_{tot}$ . The features with same numbers of levels must be placed side by side. The missing values should be coded as NA.
<code>y</code>	Vector of length $N$ . It should represent the classes corresponding to each row of <code>x</code> . Must be labeled with numbers $(1,2,\dots,kr)$ .
<code>idx_list</code>	Vector of length $D$ . This argument is useful when variables have different numbers of levels. Element $d$ should indicate where the variables with number of levels $m[d]$ begin in matrix <code>x</code> .
<code>kr</code>	Number of row classes.
<code>kc</code>	Vector of length $D$ . The $d^{\text{th}}$ element indicates the number of column clusters. Set to 0 to choose a classical multivariate BOS model.
<code>m</code>	Vector of length $D$ . The $d^{\text{th}}$ element defines the number of levels of the ordinal data.
<code>nbSEM</code>	Number of SEM-Gibbs iterations realized to estimate parameters.
<code>nbSEMBurn</code>	Number of SEM-Gibbs burn-in iterations for estimating parameters. This parameter must be inferior to <code>nbSEM</code> .
<code>nbindmini</code>	Minimum number of cells belonging to a block.
<code>init</code>	String that indicates the kind of initialisation. Must be one of the following strings: "kmeans", "random" or "randomBurnin".
<code>percentRandomB</code>	Vector of length 1. Indicates the percentage of resampling when <code>init</code> is equal to "randomBurnin".

**Value**

Return an object. The slots are:

@zr	Vector of length N with resulting row partitions.
@zc	List of length D. The d <sup>th</sup> item is a vector of length J[d] representing the column partitions for the group of variables d.
@J	Vector of length D. The d <sup>th</sup> item represents the number of columns for d <sup>th</sup> group of variables.
@W	List of length D. Item d is a matrix of dimension J*kc[d] such that W[j,h]=1 if j belongs to cluster h.
@V	Matrix of dimension N*kr such that V[i,g]=1 if i belongs to cluster g.
@icl	ICL value for co-clustering.
@kr	Number of row classes.
@name	Name of the result.
@number_distrib	Number of groups of variables.
@pi	Vector of length kr. Row mixing proportions.
@rho	List of length D. The d <sup>th</sup> item represents the column mixing proportion for the d <sup>th</sup> group of variables.
@dlist	List of length d. The d <sup>th</sup> item represents the indexes of group of variables d.
@kc	Vector of length D. The d <sup>th</sup> element represents the number of clusters column H for the d <sup>th</sup> group of variables.
@m	Vector of length D. The d <sup>th</sup> element represents the number of levels of the d <sup>th</sup> group of variables.
@nbSEM	Number of SEM-Gibbs algorithm iteration.
@params	List of length D. The d <sup>th</sup> item represents the blocks parameters for a group of variables d.
@xhat	List of length D. The d <sup>th</sup> item represents the dataset of the d <sup>th</sup> group of variables, with missing values completed.

**Author(s)**

Margot Seloisse, Julien Jacques, Christophe Biernacki.

**Examples**

```
# loading the real dataset
data("dataqol.classif")

set.seed(5)

# loading the ordinal data
M <- as.matrix(dataqol.classif[,2:29])
```

```

# creating the classes values
y <- as.vector(dataqol.classif$death)

# sampling datasets for training and to predict
nb.sample <- ceiling(nrow(M)*2/3)
sample.train <- sample(1:nrow(M), nb.sample, replace=FALSE)

M.train <- M[sample.train,]
M.validation <- M[-sample.train,]
nb.missing.validation <- length(which(M.validation==0))
m <- c(4)
M.validation[which(M.validation==0)] <- sample(1:m, nb.missing.validation,replace=TRUE)

y.train <- y[sample.train]
y.validation <- y[-sample.train]

# configuration for SEM algorithm
nbSEM=50
nbSEMBurn=40
nbindmini=1
init="kmeans"

# number of classes to predict
kr <- 2
# different kc to test with cross-validation
kcol <- 1

res <- bosclassif(x=M.train,y=y.train,kr=kr,kc=kcol,m=m,
                 nbSEM=nbSEM,nbSEMBurn=nbSEMBurn,
                 nbindmini=nbindmini,init=init)

predictions <- predict(res, M.validation)

```

---

bosclust

*Function to perform a clustering*


---

### Description

This function performs a clustering algorithm on ordinal data by using the multiple latent block model (see references for further details). It allows the user to define  $D$  groups of variables that have different numbers of levels. The BOS distribution is used, and the parameters inference is obtained using the SEM-Gibbs algorithm.

**Usage**

```
bosclust(x, idx_list=c(1), kr, init, nbSEM, nbSEMBurn,
         nbindmini, m=0, percentRandomB=0)
```

**Arguments**

<code>x</code>	Matrix made of ordinal data of dimension $N \times J_{tot}$ . The features with the same numbers of levels must be placed side by side. The missing values should be coded as NA.
<code>idx_list</code>	Vector of length $D$ . This argument is useful when variables have different numbers of levels. Element $d$ should indicate where the variables with number of levels $m[d]$ begin in matrix $x$ .
<code>kr</code>	Number of row clusters.
<code>m</code>	Vector of length $D$ . The $d^{\text{th}}$ element defines the number of levels of the ordinal data.
<code>nbSEM</code>	Number of SEM-Gibbs iterations realized to estimate the parameters.
<code>nbSEMBurn</code>	Number of SEM-Gibbs burn-in iterations for estimating parameters. This parameter must be inferior to <code>nbSEM</code> .
<code>nbindmini</code>	Minimum number of cells belonging to a block.
<code>init</code>	String that indicates the kind of initialisation. Must be one of the following words : "kmeans", "random" or "randomBurnin".
<code>percentRandomB</code>	Vector of length 1. Indicates the percentage of resampling when <code>init</code> is equal to "randomBurnin".

**Value**

<code>@v</code>	Matrix of dimension $N \times kr$ such that $V[i,g]=1$ if $i$ belongs to cluster $g$ .
<code>@zr</code>	Vector of length $N$ with resulting row partitions.
<code>@pi</code>	Vector of length $kr$ . This corresponds to the row mixing proportions.
<code>@m</code>	Vector of length $D$ . The $d^{\text{th}}$ element represents the number of levels of the $d^{\text{th}}$ group of variables.
<code>@icl</code>	ICL value for clustering.
<code>@name</code>	Name of the result.
<code>@params</code>	List of length $D$ . The $d^{\text{th}}$ item stores the resulting position and precision parameters $\mu$ and $\pi$ .
<code>@paramschain</code>	List of length <code>nbSEMBurn</code> . The parameters of the blocks are stored for each iteration of the SEM-Gibbs algorithm.
<code>@xhat</code>	List of length $D$ . The $d^{\text{th}}$ item represents the dataset of the $d^{\text{th}}$ group of variables, with missing values completed.
<code>@zrchain</code>	Matrix of dimension <code>nbSEM</code> $\times$ $N$ . Row $i$ represents the row cluster partitions at iteration $i$ .
<code>@pichain</code>	List of length <code>nbSEM</code> . Item $i$ is a vector of length $kr$ that contains the row mixing proportions at iteration $i$ .

**Author(s)**

Margot Seloisse, Julien Jacques, Christophe Biernacki.

**Examples**

```
library(ordinalClust)
data("dataqol")
set.seed(5)

# loading the ordinal data
M <- as.matrix(dataqol[,2:29])

m = 4

krow = 4

nbSEM=50
nbSEMBurn=40
nbindmini=2
init = "random"

object <- boscoclust(x=M,kr=krow, m=m, nbSEM=nbSEM,
  nbSEMBurn=nbSEMBurn, nbindmini=nbindmini, init=init)
```

---

boscoclust

*Function to perform a co-clustering*

---

**Description**

This function runs a co-clustering algorithm on ordinal data by using the latent block model (see references for further details). A BOS distribution is used, and the parameters inference is obtained using the SEM-Gbbs algorithm.

**Usage**

```
boscoclust(x=matrix(0,nrow=1,ncol=1), idx_list=c(1), kr, kc, init, nbSEM, nbSEMBurn,
  nbRepeat=1, nbindmini, m=0, percentRandomB=0)
```

**Arguments**

x Matrix made of ordinal data of dimension  $N \times J_{tot}$ . The features with the same numbers of levels must be placed side by side. The missing values should be coded as NA.

<code>idx_list</code>	Vector of length D. This argument is useful when variables have different numbers of levels. Element d should indicate where the variables with number of levels $m[d]$ begin in matrix $x$ .
<code>kr</code>	Number of row classes.
<code>kc</code>	Vector of length D. The $d^{\text{th}}$ element indicates the number of column clusters.
<code>m</code>	Vector of length D. The $d^{\text{th}}$ element defines the number of levels of the ordinal data.
<code>nbSEM</code>	Number of SEM-Gibbs iterations realized to estimate parameters.
<code>nbSEMBurn</code>	Number of SEM-Gibbs burn-in iterations for estimating parameters. This parameter must be inferior to <code>nbSEM</code> .
<code>nbRepeat</code>	Number of times sampling on rows and columns will be done for each SEM-Gibbs iteration.
<code>nbindmini</code>	Minimum number of cells belonging to a block.
<code>init</code>	String that indicates the kind of initialisation. Must be one of the following words : "kmeans", "random" or "randomBurnin".
<code>percentRandomB</code>	Vector of length 2. Indicates the percentage of resampling when <code>init</code> is equal to "randomBurnin".

**Value**

<code>@V</code>	Matrix of dimension $N \times kr$ such that $V[i,g]=1$ if $i$ belongs to cluster $g$ .
<code>@icl</code>	ICL value for co-clustering.
<code>@name</code>	Name of the result.
<code>@paramschain</code>	List of length <code>nbSEMBurn</code> . The parameters of the blocks are stored for each iteration of the SEM-Gibbs algorithm.
<code>@pichain</code>	List of length <code>nbSEM</code> . Item $i$ is a vector of length $kr$ that contains the row mixing proportions at iteration $i$ .
<code>@rhochain</code>	List of length <code>nbSEM</code> . Item $i$ is a list of length D whose $d^{\text{th}}$ element contains the column mixing proportions of the group of variables $d$ , for iteration $i$ .
<code>@zc</code>	List of length D. The $d^{\text{th}}$ item is a vector of length $J[d]$ representing the column partitions for the group of variables $d$ .
<code>@zr</code>	Vector of length N with resulting row partitions.
<code>@W</code>	List of length D. Item $d$ is a matrix of dimension $J \times kc[d]$ such that $W[j,h]=1$ if $j$ belongs to cluster $h$ .
<code>@m</code>	Vector of length D. The $d^{\text{th}}$ element represents the number of levels of $d^{\text{th}}$ group of variables.
<code>@params</code>	List of length D. The $d^{\text{th}}$ item represents the blocks parameters for a group of variables $d$ .
<code>@pi</code>	Vector of length $kr$ . This corresponds to the row mixing proportions.
<code>@rho</code>	List of length D. The $d^{\text{th}}$ item represents the column mixing proportion for the $d^{\text{th}}$ group of variables.

@xhat	List of length $D$ . The $d^{\text{th}}$ item represents the dataset of the $d^{\text{th}}$ group of variables, with missing values completed.
@zrchain	Matrix of dimension $\text{nbSEM} \times N$ . Row $i$ represents the row cluster partitions at iteration $i$ .
@zcchain	List of length $D$ . Item $d$ is a matrix of dimension $\text{nbSEM} \times J[d]$ . Row $i$ represents the column cluster partitions at iteration $i$ .

### Author(s)

Margot Selosse, Julien Jacques, Christophe Biernacki.

### Examples

```

library(ordinalClust)

# loading the real dataset
data("dataqol")
set.seed(5)

# loading the ordinal data
M <- as.matrix(dataqol[,2:29])

# defining different number of categories:
m=4

# defining number of row and column clusters
krow = 4
kcol = 4

# configuration for the inference
nbSEM=50
nbSEMBurn=40
nbindmini=2
init = "randomBurnin"
percentRandomB=c(20,20)

# Co-clustering execution
object <- boscoClust(x=M,kr=krow,kc=kcol,m=m,nbSEM=nbSEM,
                    nbSEMBurn=nbSEMBurn,nbindmini=nbindmini,init=init,percentRandomB=percentRandomB)

```



---

dataqol	<i>Questionnaire Responses Of Patients Affected By Breast Cancer</i>
---------	--

---

**Description**

This dataset contains the responses of 121 patients to 30 questions about their quality of life.

**Usage**

dataqol

**Format**

A dataframe with 121 lines and 31 columns. A line represents a patient and a column contains information about the patients.

**Id** patient Id

**q1-q28** responses to 28 questions with the number of levels equal to 4

**q29-q30** responses to 22 questions with the number of levels equal to 7

**Source**

The table was determined based on data associated with the package available on: <https://cran.r-project.org/package=QoLR>

---

dataqol.classif	<i>Questionnaire Responses Of Patients Affected By Breast Cancer</i>
-----------------	--

---

**Description**

This dataset contains the responses of 40 patients to 30 questions about their quality of life. Furthermore, a variable indicates if the patient survived the disease.

**Usage**

dataqol.classif

**Format**

A dataframe with 40 lines and 32 columns. A line represents a patient and a column contains information about the patients.

**Id** patient Id

**q1-q28** responses to 28 questions with the number of levels equal to 4

**q29-q30** responses to 22 questions with the number of levels equal to 7

**death** if the patient survived (1) or not (2)

**Source**

The table was determined based on data associated with the package available at: <https://cran.r-project.org/package=QoLR>

---

Msimulated	<i>Matrix of simulated ordinal data</i>
------------	---

---

**Description**

This is a toy dataset for running simple examples.

**Usage**

Msimulated

**Format**

An ordinal data matrix with 60 lines and 50 columns. The number of levels is equal to 3. Four blocks are simulated with (mu,pi) parameters equal to (3,0.5), (2,0.7), (1,0.8) and (2,0.6).

---

pejSim	<i>pejSim</i>
--------	---------------

---

**Description**

This function computes the probability for a level  $e_j$  to be sampled from a BOS distribution of parameters (mu,pi), with the number of levels equal to  $m$ . It can be used to generate data with a BOS distribution.

**Usage**

pejSim(ej,m,mu,p)

**Arguments**

ej	Levels to be sampled
m	Number of levels.
mu	mu parameter for BOS distribution.
p	pi parameter for BOS distribution.

**Value**

Returns the probability of  $e_j$  to be sampled from a BOS distribution of parameters (mu,pi), with the number of levels equal to  $m$ .

**Author(s)**

Margot Seloisse, Julien Jacques, Christophe Biernacki.

**Examples**

```
library(ordinalClust)
data("dataqol")
set.seed(5)

m=7
nr=10000
mu=5
pi=0.5

probaBOS=rep(0,m)
for (im in 1:m) probaBOS[im]=pejSim(im,m,mu,pi)
M <- sample(1:m,nr,prob = probaBOS, replace=TRUE)
```

---

plot

*~~ Method for the function plot in the **ordinalClust** package ~~*

---

**Description**

Plots the result of a classification, clustering or co-clustering performed using the following functions: bosclassif,bosclust,boscoclust.

**Methods**

```
signature(object = "ResultClassifOrdinal")
signature(object = "ResultClustOrdinal")
signature(object = "ResultCoclustOrdinal")
```

---

predict

*~~ Method for the function predict in the **ordinalClust** package ~~*

---

**Description**

Method for the function predict in the **ordinalClust** package.

**Methods**

```
signature(object = "ResultClassifOrdinal") Use this method with the result of the function
bosclassif, and a new sample to predict the classes.
```



# Index

## \* **datasets**

- dataqol, [9](#)
- dataqol.classif, [9](#)
- Msimulated, [10](#)

## \* **methods**

- plot, [11](#)
- predict, [11](#)
- summary, [12](#)

bosclassif, [2](#)

bosclust, [4](#)

boscoclust, [6](#)

dataqol, [9](#)

dataqol.classif, [9](#)

Msimulated, [10](#)

pejSim, [10](#)

plot, [11](#)

plot, ResultClassifOrdinal-method  
(plot), [11](#)

plot, ResultClustOrdinal-method (plot),  
[11](#)

plot, ResultCoclustOrdinal-method  
(plot), [11](#)

predict, [11](#)

predict, ANY-method (predict), [11](#)

predict, ResultClassifOrdinal-method  
(predict), [11](#)

summary, [12](#)

summary, ResultClassifOrdinal-method  
(summary), [12](#)

summary, ResultClustOrdinal-method  
(summary), [12](#)

summary, ResultCoclustOrdinal-method  
(summary), [12](#)