

Package ‘mccca’

January 24, 2024

Type Package

Title Visualizing Class Specific Heterogeneous Tendencies in Categorical Data

Version 1.1.0.1

Author Mariko Takagishi [aut, cre]

Maintainer Mariko Takagishi <m.takagishi0728@gmail.com>

Description Performing multiple-class cluster correspondence analysis(MCCCA). The main functions are create.MCCCAdata() to create a list to be applied to MCCCA, MCCCA() to apply MCCCA, and plot.mccca() for visualizing MCCCA result. Methods used in the package refers to Mariko Takagishi and Michel van de Velden (2022)<[doi:10.1080/10618600.2022.2035737](https://doi.org/10.1080/10618600.2022.2035737)>.

License GPL (>= 2)

Depends R (>= 4.1.0)

Imports magic, stringr, ggplot2, wordcloud, RColorBrewer, stats, utils, grDevices, Rcpp (>= 1.0.9)

NeedsCompilation yes

Encoding UTF-8

RoxygenNote 7.2.1

LinkingTo Rcpp, RcppArmadillo

Repository CRAN

Date/Publication 2024-01-24 14:39:05 UTC

R topics documented:

create.MCCCAdata	2
create.prop	3
generate.cate.list	4
generate.catecls	5
generate.ext	5
generate.onedata	6
MCCCA	7
plot.mccca	10

Index**13**

create.MCCCAdata	<i>this function creates a list (class: mcccadata) to be applied to MCCCA.</i>
------------------	--

Description

Creates a list (named `mcccadata.list`) applied to MCCCA.

Usage

```
create.MCCCAdata(dat, ext.mat=ext.mat, clstr0.vec=NULL)
```

Arguments

<code>dat</code>	An (N×J) matrix of categorical data (N:the number of observations, J:the number of variables). If <code>rownames(dat)</code> is NULL, <code>c(obj1, . . . , objN)</code> are defined as <code>rownames(dat)</code> .
<code>ext.mat</code>	An (N×H) external variable matrix (H:the number of external variable).
<code>clstr0.vec</code>	An integer vector of length N giving each observation's true cluster.

Value

Returns a list with the following elements.

<code>data.mat</code>	data matrix same as <code>dat</code> .
<code>data.list</code>	A list of C (N×J) categorical data matrices for each class (C:the number of classes).
<code>clstr0.list</code>	A list of C vectors where each vector indicates the true cluster (given in <code>clstr0.vec</code>) to which each class of observations belongs (NULL if <code>clstr0.vec</code> is NULL).
<code>N.vec</code>	A vector of length C giving the number of observations in each class.
<code>Ktrue.vec</code>	A vector of length C giving the true number of clusters in each class (NULL if <code>clstr0.vec</code> is NULL).
<code>q.vec</code>	A vector of length J giving the number of categories in each of J categorical variables.
<code>class.n.vec</code>	An integer (from 1:C) vector of length N giving the class index of each observation. <code>names(class.n.vec)=rownames(dat)</code> .
<code>classname.n.vec</code>	A characteristic vector of length N giving the class label each observation belongs to. <code>names(classname.n.vec)=rownames(dat)</code> .
<code>classlabel</code>	A characteristic vector of length C giving the classlabel for each class.
<code>classlab.mat</code>	(C×(H+1)) table, showing which combinations of categories of external variables each class index and class name corresponds to. The first H columns indicate the categories for each of the H external variables, and the last H+1th column indicates the corresponding class label (same as <code>classlabel</code>).

`oriindex.list` A list of length `C`, where each list element corresponds to a row (observation) in `data.list`, indicating which row of observations (in `data.mat`) each observation (in `oriindex.list`) corresponds to.

References

Takagishi & Michel van de Velden (2022): Visualizing Class Specific Heterogeneous Tendencies in Categorical Data, *Journal of Computational and Graphical Statistics*, DOI: 10.1080/10618600.2022.2035737

Examples

```
#setting
N <- 100 ; J <- 5 ; Ktrue <- 2 ; q.vec <- rep(5,J) ; noise.prop <- 0.2
extcate.vec=c(2,3)#the number of categories for each external variable

#generate categorical variable data
catedata.list <- generate.onedata(N=N,J=J,Ktrue=Ktrue,q.vec=q.vec,noise.prop = noise.prop)
data.cate=catedata.list$data.mat
clstr0.vec=catedata.list$clstr0.vec

#generate external variable data
data.ext=generate.ext(N,extcate.vec=extcate.vec)

#create mcca.list to be applied to MCCA function
mcca.data=create.MCCAdata(data.cate,ext.mat=data.ext,clstr0.vec =clstr0.vec)

#check which class each observation belongs to. (given by class name)
mcca.data$classname.n.vec

#A table showing that which combinations of categories of external variables
# each class index and class name corresponds to.
mcca.data$classlab.mat
```

`create.prop`

Creates a list length J of category proportion for each cluster.

Description

Creates a list length J of category proportion for each cluster.

Usage

```
create.prop(
  J = J,
  q.vec = q.vec,
  Ktrue = Ktrue,
  strongprop = 0.8,
  which.noise = NULL
)
```

Arguments

J	The number of active variable.!!!
q.vec	A vector of length J giving the number of categories for each active variable.
Ktrue	The number of clusters in J active variables.
strongprop	A numeric value giving the strongest proportion of categories (common for all J active variables).
which.noise	A vector of length ($\leq J$) giving the index of noise variables in J active variables. NULL indicating all variable is non-noise.

Value

Returns a list length J, each of which is a ($K_{true} \times q_j$) matrix giving the proportion for each q_j category in each K_{true} cluster.

`generate.cate.list` *Generate (NxJ) categorical data matrix.*

Description

Generate an ($N \times J$) categorical data matrix given by `prop.list` and true cluster allocation.

Usage

```
generate.cate.list(N = N, prop.list = prop.list)
```

Arguments

N	The number of observations.
prop.list	a list length J, each of which is a vector of length q_j giving the proportion for each categories.

Value

an ($N \times J$) categorical data matrix.

generate.catecls *Generate (NxJ) clustered categorical data matrix.*

Description

Generate an (NxJ) clustered categorical data matrix given by prop.J.list and true cluster allocation.

Usage

```
generate.catecls(
  N = N,
  J = J,
  q.vec = q.vec,
  Ktrue = Ktrue,
  prop.J.list = prop.J.list,
  clstr.vec = clstr.vec
)
```

Arguments

N	The number of observations.
J	The number of active variables.
q.vec	A vector of length J giving the number of categories for each active variable.
Ktrue	An integer indicating the number of content-based clusters used for CCRS estimation.
prop.J.list	a list of length J, where each list is a (Ktrue x qj) matrix giving the proportion for each qj category in each of the Ktrue cluster.
clstr.vec	A vector of length N giving true clusters for each observations.

Value

an (NxJ) clustered categorical data matrix.

generate.ext *generates an artificial (NxH) external variable matrix.*

Description

Generates an artificial (NxH) external variable matrix.

Usage

```
generate.ext(N,extcate.vec=extcate.vec,unbala.cate=FALSE)
```

Arguments

N	The number of observation.
extcate.vec	A vector of length H, each element indicates the number of category for each H external variables.
unbala.cate	logical value. If TRUE, the proportion of categories in the external variable is unbalanced. The default is FALSE.

Value

An (NxH) external variable matrix.

See Also

[generate.catecls](#)

Examples

```
###data setting
N <- 30 ; extcate.vec=c(2,3)
ext.mat=generate.ext(N,extcate.vec=extcate.vec)
```

generate.onedata *Generate (NxJ) categorical data matrix.*

Description

Generate (NxJ) categorical data matrix.

Usage

```
generate.onedata(N=100, J=5, Ktrue=3, q.vec=rep(3,5), noise.prop=0.3)
```

Arguments

N	The number of observations. Default is 100.
J	The number of active variables. Default is 5.
Ktrue	The number of true clusters. Default is 3.
q.vec	A vector of length J giving the number of categories for each active variable. Default is rep(3,5).
noise.prop	A numeric value between 0 and 1 indicating the proportion of noise variables among J variables. Default is 0.3.

Value

Returns a list with the following elements.

`data.mat` A (NxJ) data frame of categorical data.
`clstr0.vec` A vector of integers (from 1:Ktrue) length N giving the cluster to which each observation is allocated.

See Also

[create.prop](#), [generate.catecls](#)

Examples

```
###data setting
N <- 30 ; J <- 10 ; Ktrue <- 2 ; q.vec <- rep(5,J) ; noise.prop <- 0.3
datagene <- generate.onedata(N=N,J=J,Ktrue=Ktrue,q.vec=q.vec,noise.prop = noise.prop)
```

MCCCA *apply MCCCA for dataset.*

Description

Applies MCCCA to `mcccadata.list`.

Usage

```
MCCCA(
  mcca.data,
  K.vec = K.vec,
  known.vec = NULL,
  knowncluster.list = NULL,
  nstart = 3,
  maxit = 50,
  p = 2,
  tol = 1e-08,
  verbose = TRUE,
  remove.miss = TRUE,
  kmeans.initial = TRUE
)
```

Arguments

`mcca.data` A list created in [create.MCCCAdata](#).
`K.vec` An integer vector of length C (the number of classes). Each element corresponds to the number of clusters in each class specified for estimation.
`known.vec` A vector of length C giving logical values indicating whether a cluster allocation in each class is known or not. The default is all FALSE.

<code>knowncluster.list</code>	A vector of length <code>C</code> giving logical values indicating whether a cluster allocation in each class is known or not. The default is all FALSE.
<code>nstart</code>	An integer indicating the number of random initial values.
<code>maxit</code>	An integer indicating the maximum number of iterations.
<code>p</code>	An integer indicating the dimension of quantification. The default is 2.
<code>tol</code>	A numeric value indicating the absolute convergence tolerance.
<code>verbose</code>	A logical value indicating. If TRUE, tracing information on the progress of the optimization is produced.
<code>remove.miss</code>	A logical value indicating whether categories nobody choose are removed or not. The default is TRUE.
<code>kmeans.initial</code>	A logical value indicating whether the 1st initial value for indicator matrix is generated by <code>kmeans</code> or not. The default is TRUE.

Details

`Bg`, `Gg` and `Qg` are scaled `B`, `G` and `Q` respectively, such that the average squared deviation from the origin of the row and column points is the same (See section 2.3 in the paper).

If you want to specify the cluster allocation for some or all classes, prepare the following two.

-`knowncluster.list`: A list of `C` vectors. The length of each vector in the list should be the same as the number of rows in each matrix in the `data.list` (ex. `length(knowncluster.list[[c]])=nrow(data.list[[c]])`, (`c=1,...,C`)). For example, suppose that `data.list` is a list of 4 matrices (meaning `C=4`), and the cluster assignment is known only for the second class, and the assignments in other classes are estimated. In this case, the second vector of `knowncluster.list` should be specified as the vector of cluster indexes to which the observations in each row of `data.list[[2]]` belong, with length `nrow(data.list[[2]])`, and the other vectors (1, 3, and 4) in the list can be specified as NA. For each vector in the `knowncluster.list`, the specified cluster index should start from 1, and there should not be any skipping numbers.

-`known.vec`: A vector of logical values of length `C`. For example, if `C=4` and you want to know the cluster assignment of only the second class, it should be `known.vec=c(FALSE, TRUE, FALSE, FALSE)`.

Value

Returns a list with the following elements.

<code>G</code>	A (<code>K</code> × <code>p</code>) quantification matrix for all clusters (<code>K=sum(K.vec)</code>).
<code>Gg</code>	Scaled <code>G</code> . See details.
<code>B</code>	A (<code>Q</code> × <code>p</code>) quantification matrix for all categories (<code>Q=sum(q.vec)</code>), and <code>q.vec</code> is given in <code>create.MCCCAdata</code> .
<code>Bg</code>	Scaled <code>B</code> .
<code>Q</code>	A (<code>N</code> × <code>p</code>) quantification matrix for all observations.
<code>Qg</code>	Scaled <code>Q</code> .
<code>clses.list</code>	A list of <code>C</code> vectors, giving the estimated cluster index for each observation in each class.

<code>clses.vec</code>	A vector of length N, where each element represents the cluster index to which the observations in the rows of <code>data.mat</code> (given in <code>mccca.data</code>) belong.
<code>optval</code>	A numeric value giving the optimized value of the objective function that is the smallest among all initial values.
<code>optval.vec</code>	A numeric vector of length <code>nstart</code> giving the optimized values of the objective function for each initial value.
<code>stepconv</code>	An integer giving the number of iterations until convergence at the initial value where the objective function was the smallest.
<code>stepconv.vec</code>	An integer vector of length <code>nstart</code> giving the number of iterations until convergence for each initial value.
<code>catename.vec</code>	A characteristic vector of length Q that combines the category names of each categorical variable into a single vector.
<code>catename.vari.vec</code>	A characteristic vector of length Q with <code>catename.vec</code> plus the name of categorical variable (by default, this is used as the column name of B and Bg).
<code>cate.removed</code>	If there is a category that no one chooses and <code>remove.miss=TRUE</code> , <code>cate.removed</code> gives which category was removed (given by the index of column in dummy matrix). Otherwise, return NULL.
<code>cluster.vec</code>	An integer vector of length K, where each index in the <code>clses.list</code> and <code>clses.vec</code> indicates which class it corresponds to.
<code>q.vec</code>	A vector of length J, same as the one given in <code>mccca.data</code> .
<code>K.vec</code>	A vector of length C, which is used as an input in this MCCCA function.
<code>classlabel</code>	A characteristic vector of length C, same as the one given in <code>mccca.data</code> .

References

Takagishi & Michel van de Velden (2022): Visualizing Class Specific Heterogeneous Tendencies in Categorical Data, *Journal of Computational and Graphical Statistics*, DOI: 10.1080/10618600.2022.2035737

See Also

[create.MCCCAdata](#)

Examples

```
#setting
N <- 100 ; J <- 5 ; Ktrue <- 2 ; q.vec <- rep(5,J) ; noise.prop <- 0.2
extcate.vec=c(2,3)#the number of categories for each external variable

#generate categorical variable data
catedata.list <- generate.onedata(N=N,J=J,Ktrue=Ktrue,q.vec=q.vec,noise.prop = noise.prop)
data.cate=catedata.list$data.mat
clstr0.vec=catedata.list$clstr0.vec

#generate external variable data
data.ext=generate.ext(N,extcate.vec=extcate.vec)
```

```

#create mccca.list to be applied to MCCA function
mccca.data=create.MCCCAdata(data.cate,ext.mat=data.ext,clstr0.vec =clstr0.vec)

#specify the number of cluster for each of C classes
C=length(mccca.data$data.list)
K.vec=rep(2,C)

#apply MCCA
mccca.res=MCCA(mccca.data,K.vec=K.vec)

#plot MCCA result
plot(mccca.res)

#if you want to specify cluster allocation in the 2nd class:
knowncluster.list=rep(list(NA),C)
#specify cluster index for the 2nd class
N2=nrow(mccca.data$data.list[[2]])
knowncluster.list[[2]]=rep(c(1,2),times=c(2,N2-2))
known.vec=c(FALSE,TRUE,FALSE,FALSE,FALSE,FALSE)
mccca.res=MCCA(mccca.data,K.vec=K.vec,known.vec=known.vec,knowncluster.list = knowncluster.list)

```

plot.mccca

plot mccca object.

Description

plot mccca object.

Usage

```

## S3 method for class 'mccca'
plot(
  x,
  main = "MCCA result",
  catelabel = NULL,
  classlabel = NULL,
  classlabel.legend = NULL,
  xlim = NULL,
  ylim = NULL,
  sort.clssize = TRUE,
  break.size = NULL,
  output.coord = FALSE,
  connect.cord = TRUE,
  include.variname = TRUE,
  scale.gamma = TRUE,
  scatter.level = 2,
  plot.setting = list(alp.point = 0.3, alp.seg = 0.8, txtsize = 3, txtsize.legend = 10),
  ...
)

```

Arguments

<code>x</code>	An object of class <code>mccca</code> , a list of MCCA outputs.
<code>main</code>	A character giving the title of biplot.
<code>catelabel</code>	A characteristic vector of length <code>Q</code> giving labels for all categories to be displayed on the biplot ($Q = \text{sum}(q.\text{vec})$). If <code>NULL</code> , <code>rownames(B)</code> are used.
<code>classlabel</code>	A characteristic vector of length <code>C</code> (<code>C</code> :the number of class) giving labels for all classes to be displayed on the biplot. If <code>NULL</code> , labels specified in <code>create.MCCCAdata</code> are used.
<code>classlabel.legend</code>	A characteristic vector of length <code>C</code> giving labels for all classes to be used on the legend (this can be longer). If <code>NULL</code> , <code>classlabel</code> is used.
<code>xlim</code>	A numeric vector of length 2 giving the range of plot on the x (horizontal) axis. If <code>NULL</code> , the range is automatically determined.
<code>ylim</code>	A numeric vector of length 2 for the y (vertical) axis (same role as <code>xlim</code>).
<code>sort.classsize</code>	If <code>TRUE</code> , the class-specific cluster numbers are sorted in the order of cluster size. The default is <code>TRUE</code> .
<code>break.size</code>	An integer vector that adjusts the size of bubble displayed on the legend.
<code>output.coord</code>	If <code>TRUE</code> , the output will be <code>Cocls.mat</code> and <code>Cocate.mat</code> . See value.
<code>connect.cord</code>	If <code>TRUE</code> , lines are drawn between original (estimated by MCCA) coordinates and coordinates moved to avoid overlap.
<code>include.variname</code>	If <code>TRUE</code> , variable name is included in category labels in the biplot (ex.a point of category "male" in "v1"(the name of 1st variable) is displayed as "v1:male" on the biplot).
<code>scale.gamma</code>	If <code>TRUE</code> , quantifications are scaled such that the average squared deviation from the origin of the row and column points is the same (See section 2.3 in the paper).
<code>scatter.level</code>	A numeric value that adjusts the scatter of points in the biplot. The higher the value, the more scattered the points are. The default is 2.
<code>plot.setting</code>	A list of biplot settings. See details.
<code>...</code>	Additional arguments passed to <code>print</code> .

Details

Parameters in `plot.setting` are as follows:

-`alp.point`:A numeric value from 0 to 1 which adjusts the transparency of the bubble point. The default is 0.3.

-`alp.seg`:A numeric value from 0 to 1 which adjusts the transparency of the segments between texts and points. The default is 0.8.

-`txtsize`:A numeric value which adjusts the textsize on the biplot. The default is 3.

-`txtsize.legend`:A numeric value which adjusts the textsize of the legend on the biplot. The default is 10.

Value

If `output.coord` is TRUE, returns a list with the following elements.

<code>Cocls.mat</code>	A (Kx4) coordinate matrix of clusters, where the last two columns are the coordinates estimated by MCCA, and the first two columns are the coordinates moved from the estimated coordinates to prevent overlap.
<code>Cocate.mat</code>	A (Kx4) coordinate matrix of categories (each column plays the same role as <code>Cocls.mat</code>)

References

Takagishi & Michel van de Velden (2022): Visualizing Class Specific Heterogeneous Tendencies in Categorical Data, *Journal of Computational and Graphical Statistics*, DOI: 10.1080/10618600.2022.2035737

See Also

[MCCA](#)

Examples

```
#setting
N <- 100 ; J <- 5 ; Ktrue <- 2 ; q.vec <- rep(5,J) ; noise.prop <- 0.2
extcate.vec=c(2,3)#the number of categories for each external variable

#generate categorical variable data
catedata.list <- generate.onedata(N=N,J=J,Ktrue=Ktrue,q.vec=q.vec,noise.prop = noise.prop)
data.cate=catedata.list$data.mat
clstr0.vec=catedata.list$clstr0.vec
#generate external variable data
data.ext=generate.ext(N,extcate.vec=extcate.vec)

#create mccca.list to be applied to MCCA function
mccca.data=create.MCCCAdata(data.cate,ext.mat=data.ext,clstr0.vec =clstr0.vec)

#specify the number of cluster for each of C classes
C=length(mccca.data$data.list)
K.vec=rep(2,C)
#apply MCCA
mccca.res=MCCA(mccca.data,K.vec=K.vec)

#plot MCCA result
plot(mccca.res)
```

Index

`create.MCCCAdata`, [2](#), [7](#), [9](#)
`create.prop`, [3](#), [7](#)

`generate.cate.list`, [4](#)
`generate.catecls`, [5](#), [6](#), [7](#)
`generate.ext`, [5](#)
`generate.onedata`, [6](#)

MCCCA, [7](#), [12](#)

`plot.mccca`, [10](#)
`print`, [11](#)