

Package ‘nomclust’

July 12, 2021

Title Hierarchical Cluster Analysis of Nominal Data

Author Zdenek Sulc [aut, cre],
Jana Cibulkova [aut],
Hana Rezankova [aut]

Maintainer Zdenek Sulc <zdenek.sulc@vse.cz>

Version 2.5.0

Date 2021-7-12

Description Similarity measures for hierarchical clustering of objects characterized by nominal (categorical) variables. Evaluation criteria for nominal data clustering.

Depends cluster, methods

License GPL (>= 2)

RoxygenNote 7.1.1

NeedsCompilation yes

Encoding UTF-8

Imports Rcpp (>= 0.11.0)

LinkingTo Rcpp

Repository CRAN

Date/Publication 2021-07-12 09:40:02 UTC

R topics documented:

as.agnes	2
CA.methods	3
data20	4
dend.plot	4
eskin	6
eval.plot	7
evalclust	9
good1	10
good2	11
good3	13

good4	14
iof	15
lin	16
lin1	17
nomclust	19
nomprox	21
of	24
sm	25
ve	26
vm	27

Index	29
--------------	-----------

as.agnes	<i>Convert Objects to Class agnes, twins</i>
----------	--

Description

Converts objects of the class "nomclust" to the class "agnes, twins".

Usage

```
as.agnes(x, ...)
```

Arguments

x	The "nomclust" object containing components "dend" and "prox".
...	Further arguments passed to or from other methods.

Value

The function returns an object of class "agnes, twins".

Author(s)

Zdenek Sulc.
Contact: <zdenek.sulc@vse.cz>

See Also

[agnes](#), [as.hclust](#) and [hclust](#).

Examples

```
# sample data
data(data20)

# creating an object with results of hierarchical clustering of
hca.object <- nomclust(data20, measure = "lin", method = "average",
  clu.high = 5, prox = TRUE)

# nomclust plot
plot(hca.object)

# obtaining the agnes, twins object
hca.object.agnes <- as.agnes(hca.object)

# agnes plot
plot(hca.object.agnes)

# obtaining the hclust object
hca.object.hclust <- as.hclust(hca.object)

# hclust plot
plot(hca.object.hclust)
```

CA.methods

Selected clustering algorithms

Description

The dataset contains five different characteristics of 24 clustering algorithms. The "Type" variable expresses the principle on which the clustering is based. There are five possible categories: density, grid, hierarchical, model-based, and partitioning. The binary variable "OptClu" indicates if the clustering algorithm offers the optimal number of clusters. The variable "Large" indicates if the clustering algorithm was designed to cluster large datasets. The "TypicalType" variable presents the typical data type for which the clustering algorithm was determined. There are three possible categories: categorical, mixed, and quantitative. Since some clustering algorithms support more data types, the binary variable "MoreTypes" indicates this support.

Usage

```
data("CA.methods")
```

Format

A data frame containing 5 variables and 24 cases.

Source

created by the authors of the nomclust package

`data20`*Artificial nominal dataset*

Description

This dataset consists of 5 nominal variables and 20 cases. Its main aim is to demonstrate the desired entry data structure for the nomclust package.

Usage

```
data(data20)
```

Format

A data frame containing 5 variables and 20 cases.

Source

created by the authors of the nomclust package

`dend.plot`*Visualization of Cluster Hierarchy using a Dendrogram*

Description

The function visualizes the hierarchy of clusters using a dendrogram. The function also enables a user to distinguish the individual clusters with colors. The number of displayed clusters can be defined by a user or by one of the six evaluation criteria.

Usage

```
dend.plot(  
  x,  
  clusters = "BIC",  
  style = "greys",  
  colorful = TRUE,  
  clu.col = NA,  
  main = "Dendrogram",  
  ac = TRUE,  
  ...  
)
```

Arguments

x	An output of the <code>nomclust()</code> or <code>nomprox()</code> functions containing the dend component.
clusters	Either a numeric value or a character string with the name of the evaluation criterion expressing the number of displayed clusters in a dendrogram. The following evaluation criteria can be used: "AIC", "BIC", "BK", "PSFE", "PSFM", and "SI".
style	A character string or a vector of colors defines the graphical style of the produced plots. There are two predefined styles in the nomclust package, namely "greys" and "dark", but a custom color scheme can be set by a user as a vector of a length four.
colorful	A logical argument that is specifying if the output will be colorful or black and white.
clu.col	An optional vector of colors that allows a researcher to apply user-defined colors for displayed (marked) clusters in a dendrogram.
main	A character string with the chart title.
ac	A logical argument that indicates if an agglomerative coefficient will be present in the output.
...	Other graphical arguments compatible with the generic <code>plot()</code> function.

Details

The function can be applied to a `nomclust()` or `nomprox()` output containing the dend component.

Value

The function returns a dendrogram describing the hierarchy of clusters that can help to identify the optimal number of clusters.

Author(s)

Jana Cibulkova and Zdenek Sulc.
Contact: <jana.cibulkova@vse.cz>

See Also

[eval.plot](#), [nomclust](#), [nomprox](#).

Examples

```
# sample data
data(data20)

# creating an object with results of hierarchical clustering
hca.object <- nomclust(data20, measure = "iof", eval = TRUE)
```

```
# a basic plot
dend.plot(hca.object)

# a dendrogram with color-coded clusters according to the BIC index
dend.plot(hca.object, clusters = "BIC", colorful = TRUE)

# using a dark style and specifying own colors in a solution with three clusters
dend.plot(hca.object, clusters = 3, style = "dark", clu.col = c("blue", "red", "green"))

# a black and white dendrogram
dend.plot(hca.object, clusters = 3, style = "dark", colorful = FALSE)
```

eskin

Eskin (ES) Measure

Description

The function calculates a dissimilarity matrix based on the ES similarity measure.

Usage

```
eskin(data)
```

Arguments

`data` A data.frame or a matrix with cases in rows and variables in columns.

Details

The Eskin similarity measure was proposed by Eskin et al. (2002) and examined by Boriah et al., (2008). It is constructed to assign higher weights to mismatches on variables with more categories.

Value

The function returns an object of the class "dist".

Author(s)

Zdenek Sulc.
Contact: <zdenek.sulc@vse.cz>

References

Boriah S., Chandola V., Kumar V. (2008). Similarity measures for categorical data: A comparative evaluation. In: Proceedings of the 8th SIAM International Conference on Data Mining, SIAM, p. 243-254.

Eskin E., Arnold A., Prerau M., Portnoy L. and Stolfo S. (2002). A geometric framework for unsupervised anomaly detection. In D. Barbara and S. Jajodia (Eds): Applications of Data Mining in Computer Security, p. 78-100. Norwell: Kluwer Academic Publishers.

See Also

[good1](#), [good2](#), [good3](#), [good4](#), [iof](#), [lin](#), [lin1](#), [of](#), [sm](#), [ve](#), [vm](#).

Examples

```
# sample data
data(data20)

# dissimilarity matrix calculation
prox.eskin <- eskin(data20)
```

eval.plot

Visualization of Evaluation Criteria

Description

The function visualizes the values of up to eight evaluation criteria for the range of cluster solutions defined by the user in the **nomclust**, **evalclust** or **nomprox** functions. It also indicates the optimal number of clusters determined by these criteria. The charts for the evaluation criteria in the **nomclust** package.

Usage

```
eval.plot(
  x,
  criteria = "all",
  style = "greys",
  opt.col = "red",
  main = "Cluster Evaluation",
  ...
)
```

Arguments

x	An output of the "nomclust" object containing the eval and opt components.
criteria	A character string or character vector specifying the criteria that are going to be visualized. It can be selected one particular criterion, a vector of criteria, or all the available criteria by typing "all".

style	A character string or a vector of colors defines the graphical style of the produced plots. There are two predefined styles in the nomclust package, namely "greys" and "dark", but a custom color scheme can be set by a user as a vector of a length four.
opt.col	An argument specifying a color that is used for the optimal number of clusters identification.
main	A character string with the chart title.
...	Other graphical arguments compatible with the generic plot() function.

Details

The function can display up to eight evaluation criteria. Namely, Within-cluster mutability coefficient (WCM), Within-cluster entropy coefficient (WCE), Pseudo F Indices based on the mutability (PSFM) and the entropy (PSFE), Bayesian (BIC), and Akaike (AIC) information criteria for categorical data, the BK index, and the silhouette index (SI).

Value

The function returns a series of up to eight plots with evaluation criteria values and the graphical indication of the optimal numbers of clusters (for AIC, BIC, BK, PSFE, PSFM, SI).

Author(s)

Jana Cibulkova and Zdenek Sulc.
Contact: <jana.cibulkova@vse.cz>

See Also

[dend.plot](#), [nomclust](#), [evalclust](#), [nomprox](#).

Examples

```
# sample data
data(data20)

# creating an object with results of hierarchical clustering
hca.object <- nomclust(data20, measure = "iof", eval = TRUE)

# a default series of plots
eval.plot(hca.object)

# changing the color indicating the optimum number of clusters
eval.plot(hca.object, opt.col= "darkorange")

# selecting only AIC and BIC criteria with the dark style
eval.plot(hca.object, criteria = c("AIC", "BIC"), style = "dark")

# selecting only SI
eval.plot(hca.object, criteria = "SI")
```


Description

The function calculates a set of evaluation criteria if the original dataset and the cluster membership variables are provided. The function calculates up to eight evaluation criteria described in (Sulc et al., 2018) and provides the optimal number of clusters based on these criteria. It is primarily focused on evaluating hierarchical clustering results obtained by similarity measures different from those that occur in the `nomclust` package. Thus, it can serve for the comparison of various similarity measures for categorical data.

Usage

```
evalclust(data, clusters, diss = NULL)
```

Arguments

<code>data</code>	A <code>data.frame</code> or a matrix with cases in rows and variables in columns.
<code>clusters</code>	A <code>data.frame</code> or a list of cluster memberships obtained based on the dataset defined in the parameter <code>data</code> in the form of a sequence from the two-cluster solution to the maximal-cluster solution.
<code>diss</code>	An optional parameter. A matrix or a <code>dist</code> object containing dissimilarities calculated based on the dataset defined in the parameter <code>data</code> .

Value

The function returns a list with three components.

The `eval` component contains up to eight evaluation criteria as vectors in a list. Namely, Within-cluster mutability coefficient (WCM), Within-cluster entropy coefficient (WCE), Pseudo F Indices based on the mutability (PSFM) and the entropy (PSFE), Bayesian (BIC), and Akaike (AIC) information criteria for categorical data, the BK index, and, if the `diss.matrix` argument is present, the silhouette index (SI).

The `opt` component is present in the output together with the `eval` component. It displays the optimal number of clusters for the evaluation criteria from the `eval` component, except for WCM and WCE, where the optimal number of clusters is based on the elbow method.

The `call` component contains the function call.

Author(s)

Zdenek Sulc.
Contact: <zdenek.sulc@vse.cz>

References

Sulc Z., Cibulkova J., Prochazka J., Rezankova H. (2018). Internal Evaluation Criteria for Categorical Data in Hierarchical Clustering: Optimal Number of Clusters Determination, Metodoloski Zveski, 15(2), p. 1-20.

See Also

[nomclust](#), [nomprox](#), [eval.plot](#).

Examples

```
# sample data
data(data20)

# creating an object with results of hierarchical clustering
hca.object <- nomclust(data20, measure = "iof", method = "average", clu.high = 7)

# the cluster memberships
data20.clu <- hca.object$mem

# obtaining evaluation criteria for the provided dataset and cluster memberships
data20.eval <- evalclust(data20, clusters = data20.clu)

# visualization of the evaluation criteria
eval.plot(data20.eval)

# silhouette index can be calculated if the dissimilarity matrix is provided
data20.eval <- evalclust(data20, clusters = data20.clu, diss = hca.object$prox)
eval.plot(data20.eval, criteria = "SI")
```

good1

Goodall 1 (G1) Measure

Description

The function calculates a dissimilarity matrix based on the G1 similarity measure.

Usage

```
good1(data)
```

Arguments

`data` A data.frame or a matrix with cases in rows and variables in columns.

Details

The Goodall 1 similarity measure was presented in (Borjah et al., 2008). It is a simple modification of the original Goodall measure (Goodall, 1966). The measure assigns higher weights to infrequent matches.

Value

The function returns an object of the class "dist".

Author(s)

Zdenek Sulc.
Contact: <zdenek.sulc@vse.cz>

References

Borjah S., Chandola V., Kumar V. (2008). Similarity measures for categorical data: A comparative evaluation. In: Proceedings of the 8th SIAM International Conference on Data Mining, SIAM, p. 243-254.

Goodall V.D. (1966). A new similarity index based on probability. Biometrics, 22(4), p. 882.

See Also

[eskin](#), [good2](#), [good3](#), [good4](#), [iof](#), [lin](#), [lin1](#), [of](#), [sm](#), [ve](#), [vm](#).

Examples

```
# sample data
data(data20)

# dissimilarity matrix calculation
prox.good1 <- good1(data20)
```

good2

Goodall 2 (G2) Measure

Description

The function calculates a dissimilarity matrix based on the G2 similarity measure.

Usage

```
good2(data)
```

Arguments

`data` A data.frame or a matrix with cases in rows and variables in columns.

Details

The Goodall 2 similarity measure was presented in (Boriah et al., 2008). It is a simple modification of the original Goodall measure (Goodall, 1966). The measure assigns weight to infrequent matches under the condition that there are also other categories, which are even less frequent than the examined one.

Value

The function returns an object of the class "dist".

Author(s)

Zdenek Sulc.
Contact: <zdenek.sulc@vse.cz>

References

Boriah S., Chandola V., Kumar V. (2008). Similarity measures for categorical data: A comparative evaluation. In: Proceedings of the 8th SIAM International Conference on Data Mining, SIAM, p. 243-254.

Goodall V.D. (1966). A new similarity index based on probability. *Biometrics*, 22(4), p. 882.

See Also

[eskin](#), [good1](#), [good3](#), [good4](#), [iof](#), [lin](#), [lin1](#), [of](#), [sm](#), [ve](#), [vm](#).

Examples

```
# sample data
data(data20)

# dissimilarity matrix calculation
prox.good2 <- good2(data20)
```

good3	<i>Goodall 3 (G3) Measure</i>
-------	-------------------------------

Description

The function calculates a dissimilarity matrix based on the G3 similarity measure.

Usage

```
good3(data)
```

Arguments

`data` A data.frame or a matrix with cases in rows and variables in columns.

Details

The Goodall 3 similarity measure was presented in (Boriah et al., 2008). It is a simple modification of the original Goodall measure (Goodall, 1966). The measure assigns higher weight if the infrequent categories match regardless on frequencies of other categories.

Value

The function returns an object of the class "dist".

Author(s)

Zdenek Sulc.
Contact: <zdenek.sulc@vse.cz>

References

Boriah S., Chandola V., Kumar V. (2008). Similarity measures for categorical data: A comparative evaluation. In: Proceedings of the 8th SIAM International Conference on Data Mining, SIAM, p. 243-254.

Goodall V.D. (1966). A new similarity index based on probability. *Biometrics*, 22(4), p. 882.

See Also

[eskin](#), [good1](#), [good2](#), [good4](#), [iof](#), [lin](#), [lin1](#), [of](#), [sm](#), [ve](#), [vm](#).

Examples

```
# sample data
data(data20)

# dissimilarity matrix calculation
prox.good3 <- good3(data20)
```

good4	<i>Goodall 4 (G4) Measure</i>
-------	-------------------------------

Description

The function calculates a dissimilarity matrix based on the G4 similarity measure.

Usage

```
good4(data)
```

Arguments

data A data.frame or a matrix with cases in rows and variables in columns.

Details

The Goodall 4 similarity measure was presented in (Borjah et al., 2008). It is a simple modification of the original Goodall measure (Goodall, 1966). It assigns higher weights to the frequent categories matches.

Value

The function returns an object of the class "dist".

Author(s)

Zdenek Sulc.
Contact: <zdenek.sulc@vse.cz>

References

Borjah S., Chandola V., Kumar V. (2008). Similarity measures for categorical data: A comparative evaluation. In: Proceedings of the 8th SIAM International Conference on Data Mining, SIAM, p. 243-254.

Goodall V.D. (1966). A new similarity index based on probability. Biometrics, 22(4), p. 882.

See Also

[eskin](#), [good1](#), [good2](#), [good3](#), [iof](#), [lin](#), [lin1](#), [of](#), [sm](#), [ve](#), [vm](#).

Examples

```
# sample data
data(data20)

# dissimilarity matrix calculation
prox.good4 <- good4(data20)
```

iof

Inverse Occurrence Frequency (IOF) Measure

Description

The function calculates a dissimilarity matrix based on the IOF similarity measure.

Usage

```
iof(data)
```

Arguments

`data` A data.frame or a matrix with cases in rows and variables in columns.

Details

The IOF (Inverse Occurrence Frequency) measure was originally constructed for the text mining tasks, see (Sparck-Jones, 1972), later, it was adjusted for categorical variables, see (Borjiah et al., 2008). The measure assigns higher weight to mismatches on less frequent values and vice versa.

Value

The function returns an object of the class "dist".

Author(s)

Zdenek Sulc.
Contact: <zdenek.sulc@vse.cz>

References

Boriah S., Chandola V., Kumar V. (2008). Similarity measures for categorical data: A comparative evaluation. In: Proceedings of the 8th SIAM International Conference on Data Mining, SIAM, p. 243-254.

Spark-Jones K. (1972). A statistical interpretation of term specificity and its application in retrieval. In Journal of Documentation, 28(1), 11-21. Later: Journal of Documentation, 60(5) (2002), 493-502.

See Also

[eskin](#), [good1](#), [good2](#), [good3](#), [good4](#), [lin](#), [lin1](#), [of](#), [sm](#), [ve](#), [vm](#).

Examples

```
# sample data
data(data20)

# dissimilarity matrix calculation
prox.iof <- iof(data20)
```

lin	<i>Lin (LIN) Measure</i>
-----	--------------------------

Description

The function calculates a dissimilarity matrix based on the LIN similarity measure.

Usage

```
lin(data)
```

Arguments

data A data.frame or a matrix with cases in rows and variables in columns.

Details

The Lin measure was introduced by Lin (1998) and presented in (Boriah et al., 2008). The measure assigns higher weights to more frequent categories in case of matches and lower weights to less frequent categories in case of mismatches.

Value

The function returns an object of the class "dist".

Author(s)

Zdenek Sulc.
Contact: <zdenek.sulc@vse.cz>

References

Boriah S., Chandola V., Kumar V. (2008). Similarity measures for categorical data: A comparative evaluation. In: Proceedings of the 8th SIAM International Conference on Data Mining, SIAM, p. 243-254.

Lin D. (1998). An information-theoretic definition of similarity. In: ICML '98: Proceedings of the 15th International Conference on Machine Learning. San Francisco, p. 296-304.

See Also

[eskin](#), [good1](#), [good2](#), [good3](#), [good4](#), [iof](#), [lin1](#), [of](#), [sm](#), [ve](#), [vm](#).

Examples

```
# sample data
data(data20)

# dissimilarity matrix calculation
prox.lin <- lin(data20)
```

lin1	<i>Lin 1 (LIN1) Measure</i>
------	-----------------------------

Description

The function calculates a dissimilarity matrix based on the LIN1 similarity measure.

Usage

```
lin1(data)
```

Arguments

data A data.frame or a matrix with cases in rows and variables in columns.

Details

The Lin 1 similarity measure was introduced in (Boriah et al., 2008) as a modification of the original Lin measure (Lin, 1998). It has a complex system of weights. In case of mismatch, lower similarity is assigned if either the mismatching values are very frequent or their relative frequency is in between the relative frequencies of mismatching values. Higher similarity is assigned if the mismatched categories are infrequent and there are a few other infrequent categories. In case of match, lower similarity is given for matches on frequent categories or matches on categories that have many other values of the same frequency. Higher similarity is given to matches on infrequent categories.

Value

The function returns an object of the class "dist".

Author(s)

Zdenek Sulc.
Contact: <zdenek.sulc@vse.cz>

References

Boriah S., Chandola V., Kumar V. (2008). Similarity measures for categorical data: A comparative evaluation. In: Proceedings of the 8th SIAM International Conference on Data Mining, SIAM, p. 243-254.

Lin D. (1998). An information-theoretic definition of similarity. In: ICML '98: Proceedings of the 15th International Conference on Machine Learning. San Francisco, p. 296-304.

See Also

[eskin](#), [good1](#), [good2](#), [good3](#), [good4](#), [iof](#), [lin](#), [of](#), [sm](#), [ve](#), [vm](#).

Examples

```
# sample data
data(data20)

# dissimilarity matrix calculation
prox.lin1 <- lin1(data20)
```

Description

The function runs hierarchical cluster analysis (HCA) with objects characterized by nominal variables (without natural order of categories). It completely covers the clustering process, from the dissimilarity matrix calculation to the cluster quality evaluation. The function enables a user to choose from twelve similarity measures for nominal data summarized by (Boriah et al., 2008) and by (Sulc and Rezankova, 2019). Next, it offers to choose from three linkage methods that can be used for categorical data. The obtained clusters can be evaluated by up to eight evaluation criteria (Sulc et al., 2018). The output of the `nomclust()` function may serve as an input for the visualization functions `dend.plot` and `eval.plot` in the `nomclust` package.

Usage

```
nomclust(
  data,
  measure = "lin",
  method = "average",
  clu.high = 6,
  eval = TRUE,
  prox = 100
)
```

Arguments

<code>data</code>	A <code>data.frame</code> or a matrix with cases in rows and variables in columns.
<code>measure</code>	A character string defining the similarity measure used for computation of proximity matrix in HCA: "eskin", "good1", "good2", "good3", "good4", "iof", "lin", "lin1", "of", "sm", "ve", "vm".
<code>method</code>	A character string defining the clustering method. The following methods can be used: "average", "complete", "single".
<code>clu.high</code>	A numeric value expressing the maximal number of cluster for which the cluster memberships variables are produced.
<code>eval</code>	A logical operator; if TRUE, evaluation of the clustering results is performed.
<code>prox</code>	A logical operator or a numeric value. If a logical value TRUE indicates that the proximity matrix is a part of the output. A numeric value (integer) of this argument indicates the maximal number of cases in a dataset for which a proximity matrix will occur in the output.

Value

The function returns a list with up to six components.

The `mem` component contains cluster membership partitions for the selected numbers of clusters

in the form of a list.

The `eval` component contains up to eight evaluation criteria as vectors in a list. Namely, Within-cluster mutability coefficient (WCM), Within-cluster entropy coefficient (WCE), Pseudo F Indices based on the mutability (PSFM) and the entropy (PSFE), Bayesian (BIC), and Akaike (AIC) information criteria for categorical data, the BK index, and, if the `prox` component is present, the silhouette index (SI).

The `opt` component is present in the output together with the `eval` component. It displays the optimal number of clusters for the evaluation criteria from the `eval` component, except for WCM and WCE, where the optimal number of clusters is based on the elbow method.

The `dend` component can be found in the output together with the `prox` component. It contains all the necessary information for dendrogram creation.

The `prox` component contains the dissimilarity matrix in the form of the "dist" object.

The `call` component contains the function call.

Author(s)

Zdenek Sulc.

Contact: <zdenek.sulc@vse.cz>

References

Boriah S., Chandola V. and Kumar, V. (2008). Similarity measures for categorical data: A comparative evaluation. In: Proceedings of the 8th SIAM International Conference on Data Mining, SIAM, p. 243-254.

Sulc Z., Cibulkova J., Prochazka J., Rezankova H. (2018). Internal Evaluation Criteria for Categorical Data in Hierarchical Clustering: Optimal Number of Clusters Determination, Metodoloski Zveski, 15(2), p. 1-20.

Sulc Z. and Rezankova H. (2019). Comparison of Similarity Measures for Categorical Data in Hierarchical Clustering. Journal of Classification. 2019, 35(1), p. 58-72. DOI: 10.1007/s00357-019-09317-5.

See Also

[evalclust](#), [nomprox](#), [eval.plot](#), [dend.plot](#).

Examples

```
# sample data
data(data20)

# creating an object with results of hierarchical clustering of
hca.object <- nomclust(data20, measure = "lin", method = "average",
  clu.high = 5, prox = TRUE)
```

```
# quick clustering summary
summary(hca.object)

# quick cluster quality evaluation
print(hca.object)

# visualization of the evaluation criteria
eval.plot(hca.object)

# a quick dendrogram
plot(hca.object)

# a dendrogram with three designated clusters
dend.plot(hca.object, clusters = 3)

# obtaining values of evaluation indices as a data.frame
data20.eval <- as.data.frame(hca.object$eval)

# getting the optimal numbers of clusters as a data.frame
data20.opt <- as.data.frame(hca.object$opt)

# extracting cluster membership variables as a data.frame
data20.mem <- as.data.frame(hca.object$mem)

# obtaining a proximity matrix
data20.prox <- as.matrix(hca.object$prox)

# setting the maximal number of objects for which a proximity matrix is provided in the output to 30
hca.object <- nomclust(data20, measure = "iof", method = "complete",
  clu.high = 5, prox = 30)

# transforming the nomclust object to the class "hclust"
hca.object.hclust <- as.hclust(hca.object)

# transforming the nomclust object to the class "agnes, twins"
hca.object.agnes <- as.agnes(hca.object)
```

Description

The function performs hierarchical cluster analysis in situations when the proximity (dissimilarity) matrix was calculated externally. For instance, in a different R package, in an own-created function, or in other software. It offers three linkage methods that can be used for categorical data. The obtained clusters can be evaluated by up to eight evaluation indices (Sulc et al., 2018).

Usage

```
nomprox(
  diss,
  data = NULL,
  method = "average",
  clu.high = 6,
  eval = TRUE,
  prox = 100
)
```

Arguments

diss	A proximity matrix or a dist object calculated based on the dataset defined in a parameter data.
data	A data.frame or a matrix with cases in rows and variables in columns.
method	A character string defining the clustering method. The following methods can be used: "average", "complete", "single".
clu.high	A numeric value that expresses the maximal number of clusters for which the cluster membership variables are produced.
eval	A logical operator; if TRUE, evaluation of clustering results is performed.
prox	A logical operator or a numeric value. If a logical value TRUE indicates that the proximity matrix is a part of the output. A numeric value (integer) of this argument indicates the maximal number of cases in a dataset for which a proximity matrix will occur in the output.

Value

The function returns a list with up to six components:

The mem component contains cluster membership partitions for the selected numbers of clusters in the form of a list.

The eval component contains up to eight evaluation criteria as vectors in a list. Namely, Within-cluster mutability coefficient (WCM), Within-cluster entropy coefficient (WCE), Pseudo F Indices based on the mutability (PSFM) and the entropy (PSFE), Bayesian (BIC), and Akaike (AIC) information criteria for categorical data, the BK index, and, if the prox component is present, the silhouette index (SI).

The opt component is present in the output together with the eval component. It displays the optimal number of clusters for the evaluation criteria from the eval component, except for WCM and WCE, where the optimal number of clusters is based on the elbow method.

The dend component can be found in the output only together with the prox component. It contains all the necessary information for dendrogram creation.

The prox component contains the dissimilarity matrix in the form of the "dist" object.

The call component contains the function call.

Author(s)

Zdenek Sulc.
Contact: <zdenek.sulc@vse.cz>

References

Sulc Z., Cibulkova J., Prochazka J., Rezankova H. (2018). Internal Evaluation Criteria for Categorical Data in Hierarchical Clustering: Optimal Number of Clusters Determination, Metodoloski Zveski, 15(2), p. 1-20.

See Also

[nomclust](#), [evalclust](#), [eval.plot](#).

Examples

```
# sample data
data(data20)

# computation of a dissimilarity matrix using the iof similarity measure
diss.matrix <- iof(data20)

# creating an object with results of hierarchical clustering
hca.object <- nomprox(diss = diss.matrix, data = data20, method = "complete",
  clu.high = 5, eval = TRUE, prox = FALSE)

# quick clustering summary
summary(hca.object)

# quick cluster quality evaluation
print(hca.object)

# visualization of the evaluation criteria
eval.plot(hca.object)

# a dendrogram can be displayed if the object contains the prox component
hca.object <- nomprox(diss = diss.matrix, data = data20, method = "complete",
  clu.high = 5, eval = TRUE, prox = TRUE)

# a quick dendrogram
plot(hca.object)

# a dendrogram with three designated clusters
dend.plot(hca.object, clusters = 3)
```

of *Occurrence Frequency (OF) Measure*

Description

The function calculates a dissimilarity matrix based on the OF similarity measure.

Usage

```
of(data)
```

Arguments

`data` A data.frame or a matrix with cases in rows and variables in columns.

Details

The OF (Occurrence Frequency) measure was originally constructed for the text mining tasks, see (Sparck-Jones, 1972), later, it was adjusted for categorical variables, see (Boriah et al., 2008) It assigns higher weight to mismatches on less frequent values and otherwise.

Value

The function returns an object of the class "dist".

Author(s)

Zdenek Sulc.
Contact: <zdenek.sulc@vse.cz>

References

Boriah S., Chandola V., Kumar V. (2008). Similarity measures for categorical data: A comparative evaluation. In: Proceedings of the 8th SIAM International Conference on Data Mining, SIAM, p. 243-254.

Spark-Jones K. (1972). A statistical interpretation of term specificity and its application in retrieval. In Journal of Documentation, 28(1), p. 11-21. Later: Journal of Documentation, 60(5) (2002), p. 493-502.

See Also

[eskin](#), [good1](#), [good2](#), [good3](#), [good4](#), [iof](#), [lin](#), [lin1](#), [sm](#), [ve](#), [vm](#).

Examples

```
# sample data
data(data20)

# dissimilarity matrix calculation
prox.of <- of(data20)
```

sm	<i>Simple Matching Coefficient (SM)</i>
----	---

Description

The function calculates a dissimilarity matrix based on the SM similarity measure.

Usage

```
sm(data)
```

Arguments

data A data.frame or a matrix with cases in rows and variables in columns.

Details

The simple matching coefficient (Sokal, 1958) represents the simplest way of measuring similarity. It does not impose any weights. By a given variable, it assigns the value 1 in case of match and value 0 otherwise.

Value

The function returns an object of the class "dist".

Author(s)

Zdenek Sulc.
Contact: <zdenek.sulc@vse.cz>

References

Boriah S., Chandola V., Kumar V. (2008). Similarity measures for categorical data: A comparative evaluation. In: Proceedings of the 8th SIAM International Conference on Data Mining, SIAM, p. 243-254.

Sokal R., Michener C. (1958). A statistical method for evaluating systematic relationships. In: Science bulletin, 38(22), The University of Kansas.

See Also

[eskin](#), [good1](#), [good2](#), [good3](#), [good4](#), [iof](#), [lin](#), [lin1](#), [of](#), [ve](#), [vm](#).

Examples

```
# sample data
data(data20)

# dissimilarity matrix calculation
prox.sm <- sm(data20)
```

ve

Variable Entropy (VE) Measure

Description

The function calculates a dissimilarity matrix based on the VE similarity measure.

Usage

```
ve(data)
```

Arguments

`data` A data.frame or a matrix with cases in rows and variables in columns.

Details

The Variable Entropy similarity measure was introduced in (Sulc and Rezankova, 2019). It treats the similarity between two categories based on the within-cluster variability expressed by the normalized entropy. The measure assigns higher weights to rare categories.

Value

The function returns an object of the class "dist".

Author(s)

Zdenek Sulc.
Contact: <zdenek.sulc@vse.cz>

References

Boriah S., Chandola V., Kumar V. (2008). Similarity measures for categorical data: A comparative evaluation. In: Proceedings of the 8th SIAM International Conference on Data Mining, SIAM, p. 243-254.

Sulc Z. and Rezankova H. (2019). Comparison of Similarity Measures for Categorical Data in Hierarchical Clustering. *Journal of Classification*. 2019, 35(1), p. 58-72. DOI: 10.1007/s00357-019-09317-5.

See Also

[eskin](#), [good1](#), [good2](#), [good3](#), [good4](#), [iof](#), [lin](#), [lin1](#), [of](#), [sm](#), [vm](#).

Examples

```
# sample data
data(data20)

# dissimilarity matrix calculation
prox.ve <- ve(data20)
```

vm

Variable Mutability (VM) measure

Description

The function calculates a dissimilarity matrix based on the VM similarity measure.

Usage

```
vm(data)
```

Arguments

`data` A data.frame or a matrix with cases in rows and variables in columns.

Details

The Variable Mutability similarity measure was introduced in (Sulc and Rezankova, 2019). It treats the similarity between two categories based on the within-cluster variability expressed by the normalized mutability. The measure assigns higher weights to rarer categories.

Value

The function returns an object of the class "dist".

Author(s)

Zdenek Sulc.
Contact: <zdenek.sulc@vse.cz>

References

Sulc Z. and Režanková H. (2019). Comparison of Similarity Measures for Categorical Data in Hierarchical Clustering. *Journal of Classification*. 2019, 35(1), p. 58-72. DOI: 10.1007/s00357-019-09317-5.

See Also

[eskin](#), [good1](#), [good2](#), [good3](#), [good4](#), [iof](#), [lin](#), [lin1](#), [of](#), [sm](#), [ve](#).

Examples

```
#sample data
data(data20)

# dissimilarity matrix calculation
prox.vm <- vm(data20)
```

Index

* **clustering**

CA.methods, 3

* **datasets**

data20, 4

agnes, 2

as.agnes, 2

as.hclust, 2

CA.methods, 3

data20, 4

dend.plot, 4, 8, 20

eskin, 6, 11–13, 15–18, 24, 26–28

eval.plot, 5, 7, 10, 20, 23

evalclust, 8, 9, 20, 23

good1, 7, 10, 12, 13, 15–18, 24, 26–28

good2, 7, 11, 11, 13, 15–18, 24, 26–28

good3, 7, 11, 12, 13, 15–18, 24, 26–28

good4, 7, 11–13, 14, 16–18, 24, 26–28

hclust, 2

iof, 7, 11–13, 15, 15, 17, 18, 24, 26–28

lin, 7, 11–13, 15, 16, 16, 18, 24, 26–28

lin1, 7, 11–13, 15–17, 17, 24, 26–28

nomclust, 5, 8, 10, 19, 23

nomprox, 5, 8, 10, 20, 21

of, 7, 11–13, 15–18, 24, 26–28

sm, 7, 11–13, 15–18, 24, 25, 27, 28

ve, 7, 11–13, 15–18, 24, 26, 26, 28

vm, 7, 11–13, 15–18, 24, 26, 27, 27