

Package ‘MixtureMissing’

April 14, 2023

Type Package

Title Robust Model-Based Clustering for Data Sets with Missing Values at Random

Version 2.0.0

Description Implementation of robust model-based cluster analysis for data sets with missing values at random. The models used are: Multivariate Contaminated Normal Mixture (MCNM, Tong and Tortora, 2022, <[doi:10.1007/s11634-021-00476-1](https://doi.org/10.1007/s11634-021-00476-1)>), Multivariate Generalized Hyperbolic Mixture (MGHM, Wei et al., 2019, <[doi:10.1016/j.csda.2018.08.016](https://doi.org/10.1016/j.csda.2018.08.016)>), Multivariate Skew's t Mixture (MStM, Wei et al., 2019, <[doi:10.1016/j.csda.2018.08.016](https://doi.org/10.1016/j.csda.2018.08.016)>), Multivariate t Mixture (MtM, Wang et al., 2004, <[doi:10.1016/j.patrec.2004.01.010](https://doi.org/10.1016/j.patrec.2004.01.010)>), and Multivariate Normal Mixture (MNM, Ghahramani and Jordan, 1994, <[doi:10.21236/ADA295618](https://doi.org/10.21236/ADA295618)>).

Imports mvtnorm (>= 1.1-2), mnormt (>= 2.0.2), cluster (>= 2.1.2), MASS (>= 7.3), numDeriv (>= 8.1.1), Bessel (>= 0.6.0)

Suggests mice (>= 3.10.0)

License GPL (>= 2)

Encoding UTF-8

LazyData true

Repository CRAN

RoxygenNote 7.2.3

Depends R (>= 3.5.0)

NeedsCompilation no

Author Hung Tong [aut, cre],
Cristina Tortora [aut, ths, dgs]

Maintainer Hung Tong <hungtongmx@gmail.com>

Date/Publication 2023-04-13 22:40:05 UTC

R topics documented:

auto	2
bankruptcy	4
cluster_impute	4
cnm_close_100	5
cnm_close_500	6
cnm_far_100	6
cnm_far_500	7
evaluation_metrics	8
generate_patterns	9
hide_values	10
initialize_clusters	11
MCNM	12
mean_impute	16
MGHM	17
MNM	20
MStM	23
MtM	26
nm_1_noise_close_100	29
nm_1_noise_close_500	30
nm_1_noise_far_100	31
nm_1_noise_far_500	32
nm_5_noise_close_100	32
nm_5_noise_close_500	33
nm_5_noise_far_100	34
nm_5_noise_far_500	34
plot.MixtureMissing	35
summary.MixtureMissing	36
tm_close_100	37
tm_close_500	38
tm_far_100	38
tm_far_500	39
Index	40

auto

*Automobile Data Set***Description**

This data set consists of three types of entities: (a) the specification of an auto in terms of various characteristics, (b) its assigned insurance risk rating, (c) its normalized losses in use as compared to other cars. The second rating corresponds to the degree to which the auto is more risky than its price indicates. Cars are initially assigned a risk factor symbol associated with its price. Then, if it is more risky (or less), this symbol is adjusted by moving it up (or down) the scale. Actuarians call this process "symboling". A value of +3 indicates that the auto is risky, -3 that it is probably pretty safe.

Usage

auto

Format

A data frame with 205 rows and 26 variables. The first 15 variables are continuous, while the last 11 variables are categorical. There are 45 rows with missing values.

normalized_losses continuous from 65 to 256.

wheel_base continuous from 86.6 to 120.9.

length continuous from 141.1 to 208.1.

width continuous from 60.3 to 72.3.

height continuous from 47.8 to 59.8.

curb_weight continuous from 1488 to 4066.

engine_size continuous from 61 to 326.

bore continuous from 2.54 to 3.94.

stroke continuous from 2.07 to 4.17.

compression_ratio continuous from 7 to 23.

horsepower continuous from 48 to 288.

peak_rpm continuous from 4150 to 6600.

city_mpg continuous from 13 to 49.

highway_mpg continuous from 16 to 54.

price continuous from 5118 to 45400.

symboling -3, -2, -1, 0, 1, 2, 3.

make alfa-romero, audi, bmw, chevrolet, dodge, honda, isuzu, jaguar, mazda, mercedes-benz, mercury, mitsubishi, nissan, peugot, plymouth, porsche, renault, saab, subaru, toyota, volkswagen, volvo

fuel_type diesel, gas.

aspiration std, turbo.

num_doors four, two.

body_style hardtop, wagon, sedan, hatchback, convertible.

drive_wheels 4wd, fwd, rwd.

engine_location front, rear.

engine_type dohc, dohcvt, l, ohc, ohcf, ohcv, rotor.

num_cylinders eight, five, four, six, three, twelve, two.

fuel_system 1bbl, 2bbl, 4bbl, idi, mfi, mpfi, spdi, spfi.

Source

Kibler, D., Aha, D.W., & Albert, M. (1989). Instance-based prediction of real-valued attributes. *Computational Intelligence*, Vol 5, 51–57. <https://archive.ics.uci.edu/ml/datasets/automobile>

bankruptcy	<i>Bankruptcy Data Set</i>
------------	----------------------------

Description

The data set contain the ratio of retained earnings (RE) to total assets, and the ratio of earnings before interests and taxes (EBIT) to total assets of 66 American firms recorded in the form of ratios. Half of the selected firms had filed for bankruptcy.

Usage

```
bankruptcy
```

Format

A data frame with 66 rows and 3 variables:

Y Status of the firm: 0 for bankruptcy and 1 for financially sound.

RE Ratio of retained earnings.

EBIT ratio of earnings before interests and taxes.

Source

Altman E.I. (1968). Financial ratios, discriminant analysis and the prediction of corporate bankruptcy. *J Finance* 23(4): 589-609 <https://www.jstor.org/stable/2978933>

cluster_impute	<i>Imputation using Cluster Means</i>
----------------	---------------------------------------

Description

Replace missing values within each cluster with the corresponding cluster mean obtained by other observed values. In other words, a separate mean imputation is applied for every cluster.

Usage

```
cluster_impute(X, clusters)
```

Arguments

X	An $n \times d$ matrix or data frame where n is the number of observations and d is the number of columns or variables. Alternately, X can be a vector of n observations.
clusters	A numeric vector containing cluster memberships. Every integer from 1 to G must be present.

Value

A complete data matrix with missing values imputed accordingly.

Examples

```
X <- matrix(nrow = 6, ncol = 3, byrow = TRUE, c(
  NA, 2, 2,
  3, NA, 5,
  4, 3, 2,
  NA, NA, 3,
  7, 2, NA,
  NA, 4, 2
))

cluster_impute(X, clusters = c(1, 1, 1, 2, 2, 2))
```

cnm_close_100

A Mixture of Two Close Contaminated Normal Distributions - 100 Observations

Description

A simulated mixture of two close contaminated normal distributions. Refer to Punzo and McNicholas (2016) for more information about the underlying distribution that generates this data set.

Usage

```
cnm_close_100
```

Format

A matrix with 100 rows and 3 variables. The first two variables make the bivariate data, while the last variable refers to cluster memberships. The first 30 rows belong to cluster 1, and the last 70 rows belong to cluster 2

d1 variable 1.

d2 variable 2.

cluster cluster memberships

Source

Punzo, A. and McNicholas, P.D., 2016. Parsimonious mixtures of multivariate contaminated normal distributions. *Biometrical Journal*, 58(6), pp.1506-1537.

cnm_close_500	<i>A Mixture of Two Close Contaminated Normal Distributions - 500 Observations</i>
---------------	--

Description

A simulated mixture of two close contaminated normal distributions. Refer to Punzo and McNicholas (2016) for more information about the underlying distribution that generates this data set.

Usage

cnm_close_500

Format

A matrix with 500 rows and 3 variables. The first two variables make the bivariate data, while the last variable refers to cluster memberships. The first 150 rows belong to cluster 1, and the last 350 rows belong to cluster 2

d1 variable 1.

d2 variable 2.

cluster cluster memberships

Source

Punzo, A. and McNicholas, P.D., 2016. Parsimonious mixtures of multivariate contaminated normal distributions. *Biometrical Journal*, 58(6), pp.1506-1537.

cnm_far_100	<i>A Mixture of Two Far Contaminated Normal Distributions - 100 Observations</i>
-------------	--

Description

A simulated mixture of two far contaminated normal distributions. Refer to Punzo and McNicholas (2016) for more information about the underlying distribution that generates this data set.

Usage

cnm_far_100

Format

A matrix with 100 rows and 3 variables. The first two variables make the bivariate data, while the last variable refers to cluster memberships. The first 30 rows belong to cluster 1, and the last 70 rows belong to cluster 2

d1 variable 1.

d2 variable 2.

cluster cluster memberships

Source

Punzo, A. and McNicholas, P.D., 2016. Parsimonious mixtures of multivariate contaminated normal distributions. *Biometrical Journal*, 58(6), pp.1506-1537.

cnm_far_500

A Mixture of Two Far Contaminated Normal Distributions - 500 Observations

Description

A simulated mixture of two far contaminated normal distributions. Refer to Punzo and McNicholas (2016) for more information about the underlying distribution that generates this data set.

Usage

cnm_far_500

Format

A matrix with 500 rows and 3 variables. The first two variables make the bivariate data, while the last variable refers to cluster memberships. The first 150 rows belong to cluster 1, and the last 350 rows belong to cluster 2

d1 variable 1.

d2 variable 2.

cluster cluster memberships

Source

Punzo, A. and McNicholas, P.D., 2016. Parsimonious mixtures of multivariate contaminated normal distributions. *Biometrical Journal*, 58(6), pp.1506-1537.

evaluation_metrics *Binary Classification Evaluation*

Description

Evaluate the performance of a classification model by comparing its predicted labels to the true labels. Various metrics are returned to give an insight on how well the model classifies the observations. This function is added to aid outlier detection evaluation of MCNM and MtM in case that true outliers are known in advance.

Usage

```
evaluation_metrics(true_labels, pred_labels)
```

Arguments

true_labels	An 0-1 or logical vector denoting the true labels. The meaning of 0 and 1 (or TRUE and FALSE) is up to the user.
pred_labels	An 0-1 or logical vector denoting the true labels. The meaning of 0 and 1 (or TRUE and FALSE) is up to the user.

Value

A list with the following slots:

matr	The confusion matrix built upon true labels and predicted labels.
TN	True negative.
FP	False positive (type I error).
FN	False negative (type II error).
TP	True positive.
TPR	True positive rate (sensitivity).
FPR	False positive rate.
TNR	True negative rate (specificity).
FNR	False negative rate.
precision	Precision or positive predictive value (PPV).
accuracy	Accuracy.
error_rate	Error rate.
FDR	False discovery rate.

Examples

```
#++++ Inputs are 0-1 vectors +++++#

evaluation_metrics(
  true_labels = c(1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 1),
  pred_labels = c(1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 1, 0, 1, 1, 1)
)

#++++ Inputs are logical vectors +++++#

evaluation_metrics(
  true_labels = c(TRUE, FALSE, FALSE, FALSE, TRUE, TRUE, TRUE, TRUE, FALSE, FALSE),
  pred_labels = c(FALSE, FALSE, TRUE, FALSE, TRUE, FALSE, FALSE, TRUE, FALSE, FALSE)
)
```

generate_patterns *Missing-Data Pattern Generation*

Description

Generate all possible missing patterns in a multivariate data set. The function can be used to complement the function `ampute()` from package `mice` in which a matrix of patterns is needed to allow for general missing-data patterns with missing-data mechanism missing at random (MAR). Using this function, each observation can have more than one missing value.

Usage

```
generate_patterns(d)
```

Arguments

`d` The number of variables or columns of the data set. `d` must be an integer greater than 1.

Details

An observation cannot have all values missing values. A complete observation is not qualified for missing-data pattern. Note that a large value of `d` may result in memory allocation error.

Value

A matrix where 0 indicates that a variable should have missing values and 1 indicates that a variable should remain complete. This matrix has `d` columns and $2^d - 2$ rows.

Examples

```
generate_patterns(4)

##### To use with the function ampute() from package mice #####
library(mice)

patterns_matr <- generate_patterns(4)
data_missing <- ampute(iris[1:4], prop = 0.5, patterns = patterns_matr)$amp
```

hide_values

Missing Values Generation

Description

A convenient function that randomly introduces missing values to an at-least-bivariate data set. The user can specify either the proportion of observations that contain some missing values or the exact number of observations that contain some missing values. Note that the function does not guarantee that underlying missing-data mechanism to be missing at random (MAR).

Usage

```
hide_values(X, prop_cases = 0.1, n_cases = NULL)
```

Arguments

<code>X</code>	An n by d matrix or data frame where n is the number of observations and d is the number of columns or variables. X must have at least 2 rows and 2 columns.
<code>prop_cases</code>	(optional) Proportion of observations that contain some missing values. <code>prop_cases</code> must be a number in $(0, 1)$. <code>prop_cases = 0.1</code> by default, but will be ignored if <code>n_cases</code> is specified.
<code>n_cases</code>	(optional) Number of observations that contain some missing values. <code>n_cases</code> must be an integer ranging from 1 to $\text{nrow}(X) - 1$.

Details

If subject to missingness, an observation can have at least 1 and at most $\text{ncol}(X) - 1$ missing values. Depending on the data set, it is not guaranteed that the resulting matrix will have the number of rows with missing values matches the specified proportion.

Value

The original n by d matrix or data frame with missing values.

Examples

```
set.seed(1234)

hide_values(iris[1:4])
hide_values(iris[1:4], prop_cases = 0.5)
hide_values(iris[1:4], n_cases = 80)
```

initialize_clusters *Cluster Initialization using a Heuristic Method*

Description

Initialize cluster memberships and component parameters to start the EM algorithm using a heuristic clustering method or user-defined labels.

Usage

```
initialize_clusters(
  X,
  G,
  init_method = c("kmedoids", "kmeans", "hierarchical", "manual"),
  clusters = NULL
)
```

Arguments

<code>X</code>	An $n \times d$ matrix or data frame where n is the number of observations and d is the number of columns or variables. Alternately, X can be a vector of n observations.
<code>G</code>	The number of clusters, which must be at least 1. If $G = 1$, then user-defined clusters is ignored.
<code>init_method</code>	(optional) A string specifying the method to initialize the EM algorithm. "kmedoids" clustering is used by default. Alternative methods include "kmeans", "hierarchical", "manual". When "manual" is chosen, a vector <code>clusters</code> of length n must be specified. When $G = 1$ and "kmedoids" clustering is used, the medoid will be returned, not the sample mean.
<code>clusters</code>	A numeric vector of length n that specifies the initial cluster memberships of the user when <code>init_method</code> is set to "manual". This argument is NULL by default, so that it is ignored whenever other given initialization methods are chosen.

Details

Available heuristic methods include k-medoids clustering, k-means clustering, and hierarchical clustering. Alternately, the user can also enter pre-specified cluster memberships, making other initialization methods possible. If the given data set contains missing values, only observations with complete records will be used to initialize clusters. However, in this case, except when $G = 1$, the resulting cluster memberships will be set to NULL since they represent those complete records rather than the original data set as a whole.

Value

A list with the following slots:

pi	Component mixing proportions.
mu	A G by d matrix where each row is the component mean vector.
Sigma	A G -dimensional array where each d by d matrix is the component covariance matrix.
clusters	An numeric vector with values from 1 to G indicating initial cluster memberships if X is a complete data set; NULL otherwise.

References

Everitt, B., Landau, S., Leese, M., and Stahl, D. (2011). *Cluster Analysis*. John Wiley & Sons.

Kaufman, L. and Rousseeuw, P. J. (2009). *Finding groups in data: an introduction to cluster analysis*, volume 344. John Wiley & Sons.

Hartigan, J. A. and Wong, M. A. (1979). Algorithm AS 136: A K-means clustering algorithm. *Applied Statistics*, **28**, 100-108. doi: 10.2307/2346830.

Examples

```
##### Initialization using a heuristic method #####

##### Initialization using user-defined labels #####

init <- initialize_clusters(iris[1:4], G = 3, init_method = 'manual',
                           clusters = as.numeric(iris$Species))

##### Initial parameters and pairwise scatterplot showing the mapping #####

init$pi
init$mu
init$Sigma
init$clusters

pairs(iris[1:4], col = init$clusters, pch = 16)
```

Description

Carries out model-based clustering using a multivariate contaminated normal mixture (MCNM). The function will determine itself if the data set is complete or incomplete and fit the appropriate model accordingly. In the incomplete case, the data set must be at least bivariate, and missing values are assumed to be missing at random (MAR).

Usage

```
MCNM(
  X,
  G,
  max_iter = 20,
  epsilon = 0.01,
  init_method = c("kmedoids", "kmeans", "hierarchical", "manual", "emEM", "RndEM"),
  clusters = NULL,
  impute = FALSE,
  equal_prop = FALSE,
  identity_cov = FALSE,
  eta_min = 1.001,
  progress = TRUE,
  n_run = 100,
  n_short = NULL,
  short_eps = 0.1
)
```

Arguments

<code>X</code>	An $n \times d$ matrix or data frame where n is the number of observations and d is the number of variables.
<code>G</code>	The number of clusters, which must be at least 1. If $G = 1$, then both <code>init_method</code> and <code>clusters</code> are ignored.
<code>max_iter</code>	(optional) A numeric value giving the maximum number of iterations each EM algorithm is allowed to use; 20 by default.
<code>epsilon</code>	(optional) A number specifying the epsilon value for the Aitken-based stopping criterion used in the EM algorithm: 0.01 by default.
<code>init_method</code>	(optional) A string specifying the method to initialize the EM algorithm. "kmedoids" clustering is used by default. Alternative methods include "kmeans", "hierarchical", "manual", "emEM", and "RndEM". When "manual" is chosen, a vector <code>clusters</code> of length n must be specified.
<code>clusters</code>	(optional) A numeric vector of length n that specifies the initial cluster memberships of the user when <code>init_method</code> is set to "manual". This argument is NULL by default, so that it is ignored whenever other given initialization methods are chosen.
<code>impute</code>	(optional) A logical value indicating whether missing values should be imputed for initialization. It is FALSE by default, in which only complete observations are used for obtaining initial parameters. When it is TRUE, imputation varies depending on the initialization method selected. For "emEM" and "RndEM",

after observations are randomly assigned cluster memberships, missing values are replaced by the corresponding cluster means. For other heuristic methods, mean imputation is applied on the whole data set as a pre-processing step.

equal_prop	(optional) A logical value indicating whether mixing proportions should be equal when initialized with emEM or RndEM; FALSE by default.
identity_cov	(optional) A logical value indicating whether covariance matrices should be set to identity matrices when initialized with emEM or RndEM; FALSE by default.
eta_min	(optional) A numeric value close to 1 to the right specifying the minimum value of eta; 1.001 by default.
progress	(optional) A logical value indicating whether the fitting progress should be displayed; TRUE by default.
n_run	(optional) Number of random sets to consider for initialization if <code>init_method = "emEM"</code> or <code>init_method = "RndEM"</code> ; 100 by default.
n_short	(optional) Number of iterations in each run of the short EM phase if <code>init_method = "emEM"</code> . It is ignored when another initialization method is used. When <code>init_method = "emEM"</code> , emEM reduces to RndEM. It is NULL by default.
short_eps	(optional) The epsilon value for the Aitken-based stopping criterion used the short EM phase. The value is ignored if <code>n_short</code> is specified (not NULL). By default, it is 0.1.

Value

An object of class `MixtureMissing` with:

<code>model</code>	The model used to fit the data set
<code>pi</code>	Mixing proportions.
<code>mu</code>	Component mean vectors.
<code>Sigma</code>	Component covariance matrices.
<code>alpha</code>	Component proportions of good observations.
<code>eta</code>	Component degrees of contamination.
<code>z_tilde</code>	An n by G matrix where each row indicates the expected probabilities that the corresponding observation belongs to each cluster.
<code>v_tilde</code>	An n by G matrix where each row indicates the expected probabilities that the corresponding observation is good with respect to each cluster.
<code>clusters</code>	A numeric vector of length n indicating cluster memberships determined by the model.
<code>outliers</code>	A logical vector of length n indicating observations that are outliers.
<code>data</code>	The original data set if it is complete; otherwise, this is the data set with missing values imputed by appropriate expectations.
<code>complete</code>	A logical vector of length n indicating which observation(s) have no missing values.
<code>npar</code>	The breakdown of the number of parameters to estimate.
<code>max_iter</code>	Maximum number of iterations allowed in the EM algorithm.

<code>iter_stop</code>	The actual number of iterations needed when fitting the data set.
<code>final_loglik</code>	The final value of log-likelihood.
<code>loglik</code>	All the values of log-likelihood.
<code>AIC</code>	Akaike information criterion.
<code>BIC</code>	Bayesian information criterion.
<code>KIC</code>	Kullback information criterion.
<code>KICc</code>	Corrected Kullback information criterion.
<code>AIC3</code>	Modified AIC.
<code>CAIC</code>	Bozdogan's consistent AIC.
<code>AICc</code>	Small-sample version of AIC.
<code>ent</code>	Entropy
<code>ICL</code>	Integrated Completed Likelihood criterion.
<code>AWE</code>	Approximate weight of evidence.
<code>CLC</code>	Classification likelihood criterion.
<code>init_method</code>	The initialization method used in model fitting.
<code>n_run</code>	Number of random sets considered for initialization if emEM or RndEM is used.
<code>n_short</code>	Number of iterations used in each run of the short EM phase.
<code>short_eps</code>	The epsilon value for the Aitken-based stopping criterion used the short EM phase.

References

Punzo, A. and McNicholas, P.D., 2016. Parsimonious mixtures of multivariate contaminated normal distributions. *Biometrical Journal*, 58(6), pp.1506-1537.

Tong, H. and, Tortora, C., 2022. Model-based clustering and outlier detection with missing data. *Advances in Data Analysis and Classification*.

Examples

```
data('nm_5_noise_close_100')

##### With no missing values #####

X <- nm_5_noise_close_100[, 1:2]
mod <- MCNM(X, G = 2, init_method = 'kmedoids', max_iter = 10)

summary(mod)
plot(mod)

##### With missing values #####

set.seed(1234)
```

```
X <- hide_values(nm_5_noise_close_100[, 1:2], prop_cases = 0.1)
mod <- MCM(X, G = 2, init_method = 'kmedoids', max_iter = 10)

summary(mod)
plot(mod)
```

mean_impute

Mean Imputation

Description

Replace missing values of data set by the mean of other observed values.

Usage

```
mean_impute(X)
```

Arguments

X An $n \times d$ matrix or data frame where n is the number of observations and d is the number of columns or variables. Alternately, X can be a vector of n observations.

Value

A complete data matrix with missing values imputed accordingly.

References

Schafer, J. L. and Graham, J. W. (2002). Missing data: our view of the state of the art. *Psychological Methods*, 7(2):147–177.

Little, R. J. A. and Rubin, D. B. (2020). *Statistical analysis with missing data*. Wiley Series in Probability and Statistics. Wiley, Hoboken, NJ, 3rd edition

Examples

```
X <- matrix(nrow = 6, ncol = 3, byrow = TRUE, c(
  NA, 2, 2,
  3, NA, 5,
  4, 3, 2,
  NA, NA, 3,
  7, 2, NA,
  NA, 4, 2
))

mean_impute(X)
```


MGHM

*Multivariate Generalized Hyperbolic Mixture (MCNM)***Description**

Carries out model-based clustering using a multivariate generalized hyperbolic mixture (MGHM). The function will determine itself if the data set is complete or incomplete and fit the appropriate model accordingly. In the incomplete case, the data set must be at least bivariate, and missing values are assumed to be missing at random (MAR).

Usage

```
MGHM(
  X,
  G,
  max_iter = 20,
  epsilon = 0.01,
  init_method = c("kmedoids", "kmeans", "hierarchical", "manual", "emEM", "RndEM"),
  clusters = NULL,
  impute = FALSE,
  equal_prop = FALSE,
  identity_cov = FALSE,
  deriv_ctrl = list(eps = 1e-08, d = 1e-04, zero.tol = sqrt(.Machine$double.eps/7e-07), r
    = 6, v = 2, show.details = FALSE),
  progress = TRUE,
  n_run = 100,
  n_short = NULL,
  short_eps = 0.1
)
```

Arguments

<code>X</code>	An $n \times d$ matrix or data frame where n is the number of observations and d is the number of variables.
<code>G</code>	The number of clusters, which must be at least 1. If $G = 1$, then both <code>init_method</code> and <code>clusters</code> are ignored.
<code>max_iter</code>	(optional) A numeric value giving the maximum number of iterations each EM algorithm is allowed to use; 20 by default.
<code>epsilon</code>	(optional) A number specifying the epsilon value for the Aitken-based stopping criterion used in the EM algorithm: 0.01 by default.
<code>init_method</code>	(optional) A string specifying the method to initialize the EM algorithm. "kmedoids" clustering is used by default. Alternative methods include "kmeans", "hierarchical", "manual", "emEM", and "RndEM". When "manual" is chosen, a vector <code>clusters</code> of length n must be specified.

clusters	(optional) A numeric vector of length n that specifies the initial cluster memberships of the user when <code>init_method</code> is set to "manual". This argument is NULL by default, so that it is ignored whenever other given initialization methods are chosen.
impute	(optional) A logical value indicating whether missing values should be imputed for initialization. It is FALSE by default, in which only complete observations are used for obtaining initial parameters. When it is TRUE, imputation varies depending on the initialization method selected. For "emEM" and "RndEM", after observations are randomly assigned cluster memberships, missing values are replaced by the corresponding cluster means. For other heuristic methods, mean imputation is applied on the whole data set as a pre-processing step.
equal_prop	(optional) A logical value indicating whether mixing proportions should be equal when initialized with emEM or RndEM; FALSE by default.
identity_cov	(optional) A logical value indicating whether covariance matrices should be set to identity matrices when initialized with emEM or RndEM; FALSE by default.
deriv_ctrl	(optional) A list containing arguments to control the numerical procedures for calculating the first and second derivatives. Some values are suggested by default. Refer to functions <code>grad</code> and <code>hessian</code> under the package <code>numDeriv</code> for more information.
progress	(optional) A logical value indicating whether the fitting progress should be displayed; TRUE by default.
n_run	(optional) Number of random sets to consider for initialization if <code>init_method = "emEM"</code> or <code>init_method = "RndEM"</code> ; 100 by default.
n_short	(optional) Number of iterations in each run of the short EM phase if <code>init_method = "emEM"</code> . It is ignored when another initialization method is used. When <code>init_method = "emEM"</code> , emEM reduces to RndEM. It is NULL by default.
short_eps	(optional) The epsilon value for the Aitken-based stopping criterion used the short EM phase. The value is ignored if <code>n_short</code> is specified (not NULL). By default, it is 0.1.

Value

An object of class `MixtureMissing` with:

model	The model used to fit the data set
pi	Mixing proportions.
mu	Component mean vectors (location).
Sigma	Component covariance matrices (dispersion).
alpha	Component skewness vectors.
lambda	Component index parameters.
omega	Component concentration parameters.
z_tilde	An n by G matrix where each row indicates the expected probabilities that the corresponding observation belongs to each cluster.
clusters	A numeric vector of length n indicating cluster memberships determined by the model.

data	The original data set if it is complete; otherwise, this is the data set with missing values imputed by appropriate expectations.
complete	A logical vector of length n indicating which observation(s) have no missing values.
npar	The breakdown of the number of parameters to estimate.
max_iter	Maximum number of iterations allowed in the EM algorithm.
iter_stop	The actual number of iterations needed when fitting the data set.
final_loglik	The final value of log-likelihood.
loglik	All the values of log-likelihood.
AIC	Akaike information criterion.
BIC	Bayesian information criterion.
KIC	Kullback information criterion.
KICc	Corrected Kullback information criterion.
AIC3	Modified AIC.
CAIC	Bozdogan's consistent AIC.
AICc	Small-sample version of AIC.
ent	Entropy
ICL	Integrated Completed Likelihood criterion.
AWE	Approximate weight of evidence.
CLC	Classification likelihood criterion.
init_method	The initialization method used in model fitting.
n_run	Number of random sets considered for initialization if emEM or RndEM is used.
n_short	Number of iterations used in each run of the short EM phase.
short_eps	The epsilon value for the Aitken-based stopping criterion used the short EM phase.

References

Browne, R. P. and McNicholas, P. D. (2015). A mixture of generalized hyperbolic distributions. *Canadian Journal of Statistics*, 43(2):176–198.

Wei, Y., Tang, Y., and McNicholas, P. D. (2019). Mixtures of generalized hyperbolic distributions and mixtures of skew- t distributions for model-based clustering with incomplete data. *Computational Statistics & Data Analysis*, 130:18–41.

Examples

```
data('bankruptcy')

##### With no missing values #####

X <- bankruptcy[, 2:3]
```

```

mod <- MGHM(X, G = 2, init_method = 'kmedoids', max_iter = 10)

summary(mod)
plot(mod)

##### With missing values #####

set.seed(1234)

X <- hide_values(bankruptcy[, 2:3], prop_cases = 0.1)
mod <- MGHM(X, G = 2, init_method = 'kmedoids', max_iter = 10)

summary(mod)
plot(mod)

```

MNM

Multivariate Normal Mixture (MNM)

Description

Carries out model-based clustering using a multivariate normal mixture (MNM). The function will determine itself if the data set is complete or incomplete and fit the appropriate model accordingly. In the incomplete case, the data set must be at least bivariate, and missing values are assumed to be missing at random (MAR).

Usage

```

MNM(
  X,
  G,
  max_iter = 20,
  epsilon = 0.01,
  init_method = c("kmedoids", "kmeans", "hierarchical", "manual", "emEM", "RndEM"),
  clusters = NULL,
  impute = FALSE,
  equal_prop = FALSE,
  identity_cov = FALSE,
  progress = TRUE,
  n_run = 100,
  n_short = NULL,
  short_eps = 0.1
)

```

Arguments

X An n by d matrix or data frame where n is the number of observations and d is the number of variables.

<code>G</code>	The number of clusters, which must be at least 1. If <code>G = 1</code> , then both <code>init_method</code> and <code>clusters</code> are ignored.
<code>max_iter</code>	(optional) A numeric value giving the maximum number of iterations each EM algorithm is allowed to use; 20 by default.
<code>epsilon</code>	(optional) A number specifying the epsilon value for the Aitken-based stopping criterion used in the EM algorithm; 0.01 by default.
<code>init_method</code>	(optional) A string specifying the method to initialize the EM algorithm. "kmeans" clustering is used by default. Alternative methods include "kmeans", "hierarchical", "manual", "emEM", and "RndEM". When "manual" is chosen, a vector <code>clusters</code> of length n must be specified.
<code>clusters</code>	(optional) A vector of length n that specifies the initial cluster memberships of the user when <code>init_method</code> is set to "manual". This argument is NULL by default, so that it is ignored whenever other given initialization methods are chosen.
<code>impute</code>	(optional) A logical value indicating whether missing values should be imputed for initialization. It is FALSE by default, in which only complete observations are used for obtaining initial parameters. When it is TRUE, imputation varies depending on the initialization method selected. For "emEM" and "RndEM", after observations are randomly assigned cluster memberships, missing values are replaced by the corresponding cluster means. For other heuristic methods, mean imputation is applied on the whole data set as a pre-processing step.
<code>equal_prop</code>	(optional) A logical value indicating whether mixing proportions should be equal when initialized with emEM or RndEM; FALSE by default.
<code>identity_cov</code>	(optional) A logical value indicating whether covariance matrices should be set to identity matrices when initialized with emEM or RndEM; FALSE by default.
<code>progress</code>	(optional) A logical value indicating whether the fitting progress should be displayed; TRUE by default.
<code>n_run</code>	(optional) Number of random sets to consider for initialization if <code>init_method = "emEM"</code> or <code>init_method = "RndEM"</code> ; 100 by default.
<code>n_short</code>	(optional) Number of iterations in each run of the short EM phase if <code>init_method = "emEM"</code> . It is ignored when another initialization method is used. When <code>init_method = "emEM"</code> , emEM reduces to RndEM. It is NULL by default.
<code>short_eps</code>	(optional) The epsilon value for the Aitken-based stopping criterion used the short EM phase. The value is ignored if <code>n_short</code> is specified (not NULL). By default, it is 0.1.

Value

An object of class `MixtureMissing` with:

<code>model</code>	The model used to fit the data set
<code>pi</code>	Mixing proportions.
<code>mu</code>	Component mean vectors.
<code>Sigma</code>	Component covariance matrices.

<code>z_tilde</code>	An n by G matrix where each row indicates the expected probabilities that the corresponding observation belongs to each cluster.
<code>clusters</code>	A numeric vector of length n indicating cluster memberships determined by the model.
<code>data</code>	The original data set if it is complete; otherwise, this is the data set with missing values imputed by appropriate expectations.
<code>complete</code>	A logical vector of length n indicating which observation(s) have no missing values.
<code>npar</code>	The breakdown of the number of parameters to estimate.
<code>max_iter</code>	Maximum number of iterations allowed in the EM algorithm.
<code>iter_stop</code>	The actual number of iterations needed when fitting the data set.
<code>final_loglik</code>	The final value of log-likelihood.
<code>loglik</code>	All the values of log-likelihood.
<code>AIC</code>	Akaike information criterion.
<code>BIC</code>	Bayesian information criterion.
<code>KIC</code>	Kullback information criterion.
<code>KICc</code>	Corrected Kullback information criterion.
<code>AIC3</code>	Modified AIC.
<code>CAIC</code>	Bozdogan's consistent AIC.
<code>AICc</code>	Small-sample version of AIC.
<code>ent</code>	Entropy
<code>ICL</code>	Integrated Completed Likelihood criterion.
<code>AWE</code>	Approximate weight of evidence.
<code>CLC</code>	Classification likelihood criterion.
<code>init_method</code>	The initialization method used in model fitting.
<code>n_run</code>	Number of random sets considered for initialization if emEM or RndEM is used.
<code>n_short</code>	Number of iterations used in each run of the short EM phase.
<code>short_eps</code>	The epsilon value for the Aitken-based stopping criterion used the short EM phase.

References

- Wolfe, J. H. (1965). A computer program for the maximum likelihood analysis of types. Technical report, NAVAL PERSONNEL RESEARCH ACTIVITY SAN DIEGO United States.
- Ghahramani, Z. and Jordan, M. I. (1995). Learning from incomplete data.

Examples

```

data('nm_5_noise_close_100')

##### With no missing values #####

X <- nm_5_noise_close_100[, 1:2]
mod <- MNM(X, G = 2, init_method = 'kmedoids', max_iter = 10)

summary(mod)
plot(mod)

##### With missing values #####

set.seed(1234)

X <- hide_values(nm_5_noise_close_100[, 1:2], prop_cases = 0.1)
mod <- MNM(X, G = 2, init_method = 'kmedoids', max_iter = 10)

summary(mod)
plot(mod)

```

MStM

Multivariate Skew-t Mixture (MStM)

Description

Carries out model-based clustering using a multivariate Skew- t mixture (MStM). The function will determine itself if the data set is complete or incomplete and fit the appropriate model accordingly. In the incomplete case, the data set must be at least bivariate, and missing values are assumed to be missing at random (MAR).

Usage

```

MStM(
  X,
  G,
  max_iter = 20,
  epsilon = 0.01,
  init_method = c("kmedoids", "kmeans", "hierarchical", "manual", "emEM", "RndEM"),
  clusters = NULL,
  impute = FALSE,
  equal_prop = FALSE,
  identity_cov = FALSE,
  df0 = rep(10, G),
  deriv_ctrl = list(eps = 1e-08, d = 1e-04, zero.tol = sqrt(.Machine$double.eps/7e-07), r
    = 6, v = 2, show.details = FALSE),
  progress = TRUE,
  n_run = 100,

```

```

    n_short = NULL,
    short_eps = 0.1
)

```

Arguments

<code>X</code>	An $n \times d$ matrix or data frame where n is the number of observations and d is the number of variables.
<code>G</code>	The number of clusters, which must be at least 1. If $G = 1$, then both <code>init_method</code> and <code>clusters</code> are ignored.
<code>max_iter</code>	(optional) A numeric value giving the maximum number of iterations each EM algorithm is allowed to use; 20 by default.
<code>epsilon</code>	(optional) A number specifying the epsilon value for the Aitken-based stopping criterion used in the EM algorithm; 0.01 by default.
<code>init_method</code>	(optional) A string specifying the method to initialize the EM algorithm. "kmeans" clustering is used by default. Alternative methods include "kmeans", "hierarchical", "manual", "emEM", and "RndEM". When "manual" is chosen, a vector <code>clusters</code> of length n must be specified.
<code>clusters</code>	(optional) A numeric vector of length n that specifies the initial cluster memberships of the user when <code>init_method</code> is set to "manual". This argument is NULL by default, so that it is ignored whenever other given initialization methods are chosen.
<code>impute</code>	(optional) A logical value indicating whether missing values should be imputed for initialization. It is FALSE by default, in which only complete observations are used for obtaining initial parameters. When it is TRUE, imputation varies depending on the initialization method selected. For "emEM" and "RndEM", after observations are randomly assigned cluster memberships, missing values are replaced by the corresponding cluster means. For other heuristic methods, mean imputation is applied on the whole data set as a pre-processing step.
<code>equal_prop</code>	(optional) A logical value indicating whether mixing proportions should be equal when initialized with emEM or RndEM; FALSE by default.
<code>identity_cov</code>	(optional) A logical value indicating whether covariance matrices should be set to identity matrices when initialized with emEM or RndEM; FALSE by default.
<code>df0</code>	(optional) Starting value of component degrees of freedom; 10 by default.
<code>deriv_ctrl</code>	(optional) A list containing arguments to control the numerical procedures for calculating the first and second derivatives. Some values are suggested by default. Refer to functions <code>grad</code> and <code>hessian</code> under the package <code>numDeriv</code> for more information.
<code>progress</code>	(optional) A logical value indicating whether the fitting progress should be displayed; TRUE by default.
<code>n_run</code>	(optional) Number of random sets to consider for initialization if <code>init_method = "emEM"</code> or <code>init_method = "RndEM"</code> ; 100 by default.
<code>n_short</code>	(optional) Number of iterations in each run of the short EM phase if <code>init_method = "emEM"</code> . It is ignored when another initialization method is used. When <code>init_method = "emEM"</code> , emEM reduces to RndEM. It is NULL by default.

`short_eps` (optional) The epsilon value for the Aitken-based stopping criterion used the short EM phase. The value is ignored if `n_short` is specified (not NULL). By default, it is 0.1.

Value

An object of class `MixtureMissing` with:

<code>model</code>	The model used to fit the data set
<code>pi</code>	Mixing proportions.
<code>mu</code>	Component mean vectors (location).
<code>Sigma</code>	Component covariance matrices (dispersion).
<code>alpha</code>	Component skewness vectors.
<code>df</code>	Component degrees of freedom.
<code>z_tilde</code>	An n by G matrix where each row indicates the expected probabilities that the corresponding observation belongs to each cluster.
<code>clusters</code>	A numeric vector of length n indicating cluster memberships determined by the model.
<code>data</code>	The original data set if it is complete; otherwise, this is the data set with missing values imputed by appropriate expectations.
<code>complete</code>	A logical vector of length n indicating which observation(s) have no missing values.
<code>npar</code>	The breakdown of the number of parameters to estimate.
<code>max_iter</code>	Maximum number of iterations allowed in the EM algorithm.
<code>iter_stop</code>	The actual number of iterations needed when fitting the data set.
<code>final_loglik</code>	The final value of log-likelihood.
<code>loglik</code>	All the values of log-likelihood.
<code>AIC</code>	Akaike information criterion.
<code>BIC</code>	Bayesian information criterion.
<code>KIC</code>	Kullback information criterion.
<code>KICc</code>	Corrected Kullback information criterion.
<code>AIC3</code>	Modified AIC.
<code>CAIC</code>	Bozdogan's consistent AIC.
<code>AICc</code>	Small-sample version of AIC.
<code>ent</code>	Entropy
<code>ICL</code>	Integrated Completed Likelihood criterion.
<code>AWE</code>	Approximate weight of evidence.
<code>CLC</code>	Classification likelihood criterion.
<code>init_method</code>	The initialization method used in model fitting.
<code>n_run</code>	Number of random sets considered for initialization if <code>emEM</code> or <code>RndEM</code> is used.
<code>n_short</code>	Number of iterations used in each run of the short EM phase.
<code>short_eps</code>	The epsilon value for the Aitken-based stopping criterion used the short EM phase.

References

Lin, T.-I. (2010). Robust mixture modeling using multivariate skew t distributions. *Statistics and Computing*, 20(3):343–356. Wang, W.-L. and Lin, T.-I. (2015). Robust model-based clustering via mixtures of skew- t distributions with missing information. *Advances in Data Analysis and Classification*, 9(4):423–445. Wei, Y., Tang, Y., and McNicholas, P. D. (2019). Mixtures of generalized hyperbolic distributions and mixtures of skew- t distributions for model-based clustering with incomplete data. *Computational Statistics & Data Analysis*, 130:18–41.

Examples

```
data('bankruptcy')

##### With no missing values #####

X <- bankruptcy[, 2:3]
mod <- MStM(X, G = 2, init_method = 'kmedoids', max_iter = 10)

summary(mod)
plot(mod)

##### With missing values #####

set.seed(1234)

X <- hide_values(bankruptcy[, 2:3], prop_cases = 0.1)
mod <- MStM(X, G = 2, init_method = 'kmedoids', max_iter = 10)

summary(mod)
plot(mod)
```

MtM

Multivariate t Mixture (MtM)

Description

Carries out model-based clustering using a multivariate t mixture (MtM). The function will determine itself if the data set is complete or incomplete and fit the appropriate model accordingly. In the incomplete case, the data set must be at least bivariate, and missing values are assumed to be missing at random (MAR).

Usage

```
MtM(
  X,
  G,
  max_iter = 20,
```

```

epsilon = 0.01,
init_method = c("kmedoids", "kmeans", "hierarchical", "manual", "emEM", "RndEM"),
clusters = NULL,
impute = FALSE,
equal_prop = FALSE,
identity_cov = FALSE,
df0 = rep(10, G),
outlier_cutoff = 0.95,
progress = TRUE,
n_run = 100,
n_short = NULL,
short_eps = 0.1
)

```

Arguments

<code>X</code>	An n by d matrix or data frame where n is the number of observations and d is the number of variables.
<code>G</code>	The number of clusters, which must be at least 1. If $G = 1$, then both <code>init_method</code> and <code>clusters</code> are ignored.
<code>max_iter</code>	(optional) A numeric value giving the maximum number of iterations each EM algorithm is allowed to use; 20 by default.
<code>epsilon</code>	(optional) A number specifying the epsilon value for the Aitken-based stopping criterion used in the EM algorithm; 0.01 by default.
<code>init_method</code>	(optional) A string specifying the method to initialize the EM algorithm. "kmedoids" clustering is used by default. Alternative methods include "kmeans", "hierarchical", "manual", "emEM", and "RndEM". When "manual" is chosen, a vector <code>clusters</code> of length n must be specified.
<code>clusters</code>	(optional) A vector of length n that specifies the initial cluster memberships of the user when <code>init_method</code> is set to "manual". Both numeric and character vectors are acceptable. This argument is NULL by default, so that it is ignored whenever other given initialization methods are chosen.
<code>impute</code>	(optional) A logical value indicating whether missing values should be imputed for initialization. It is FALSE by default, in which only complete observations are used for obtaining initial parameters. When it is TRUE, imputation varies depending on the initialization method selected. For "emEM" and "RndEM", after observations are randomly assigned cluster memberships, missing values are replaced by the corresponding cluster means. For other heuristic methods, mean imputation is applied on the whole data set as a pre-processing step.
<code>equal_prop</code>	(optional) A logical value indicating whether mixing proportions should be equal when initialized with emEM or RndEM; FALSE by default.
<code>identity_cov</code>	(optional) A logical value indicating whether covariance matrices should be set to identity matrices when initialized with emEM or RndEM; FALSE by default.
<code>df0</code>	(optional) Starting value of component degrees of freedom; 10 by default.
<code>outlier_cutoff</code>	(optional) A number between 0 and 1 indicating the percentile cutoff used for outlier detection.

progress	(optional) A logical value indicating whether the fitting progress should be displayed; TRUE by default.
n_run	(optional) Number of random sets to consider for initialization if <code>init_method = "emEM"</code> or <code>init_method = "RndEM"</code> ; 100 by default.
n_short	(optional) Number of iterations in each run of the short EM phase if <code>init_method = "emEM"</code> . It is ignored when another initialization method is used. When <code>init_method = "emEM"</code> , emEM reduces to RndEM. It is NULL by default.
short_eps	(optional) The epsilon value for the Aitken-based stopping criterion used the short EM phase. The value is ignored if <code>n_short</code> is specified (not NULL). By default, it is 0.1.

Value

An object of class `MixtureMissing` with:

model	The model used to fit the data set
pi	Mixing proportions.
mu	Component mean vectors.
Sigma	Component covariance matrices.
df	Component degrees of freedom.
z_tilde	An n by G matrix where each row indicates the expected probabilities that the corresponding observation belongs to each cluster.
clusters	A numeric vector of length n indicating cluster memberships determined by the model.
outliers	A logical vector of length n indicating observations that are outliers.
data	The original data set if it is complete; otherwise, this is the data set with missing values imputed by appropriate expectations.
complete	A logical vector of length n indicating which observation(s) have no missing values.
npar	The breakdown of the number of parameters to estimate.
max_iter	Maximum number of iterations allowed in the EM algorithm.
iter_stop	The actual number of iterations needed when fitting the data set.
final_loglik	The final value of log-likelihood.
loglik	All the values of log-likelihood.
AIC	Akaike information criterion.
BIC	Bayesian information criterion.
KIC	Kullback information criterion.
KICc	Corrected Kullback information criterion.
AIC3	Modified AIC.
CAIC	Bozdogan's consistent AIC.
AICc	Small-sample version of AIC.

ent	Entropy
ICL	Integrated Completed Likelihood criterion.
AWE	Approximate weight of evidence.
CLC	Classification likelihood criterion.
init_method	The initialization method used in model fitting.
n_run	Number of random sets considered for initialization if emEM or RndEM is used.
n_short	Number of iterations used in each run of the short EM phase.
short_eps	The epsilon value for the Aitken-based stopping criterion used the short EM phase.

References

Peel, D. and McLachlan, G. J. (2000). Robust mixture modelling using the t distribution. *Statistics and computing*, 10(4):339-348.

Wang, H., Zhang, Q., Luo, B., and Wei, S. (2004). Robust mixture modelling using multivariate t -distribution with missing information. *Pattern Recognition Letters*, 25(6):701-710.

Examples

```
data('nm_5_noise_close_100')

##### With no missing values #####

X <- nm_5_noise_close_100[, 1:2]
mod <- MtM(X, G = 2, init_method = 'kmedoids', max_iter = 10)

summary(mod)
plot(mod)

##### With missing values #####

set.seed(1234)

X <- hide_values(nm_5_noise_close_100[, 1:2], prop_cases = 0.1)
mod <- MtM(X, G = 2, init_method = 'kmedoids', max_iter = 10)

summary(mod)
plot(mod)
```

Description

A simulated mixture of two close normal distributions with 1 substituted by high atypical points. Refer to Punzo and McNicholas (2016) for more information about the underlying distribution that generates this data set.

Usage

```
nm_1_noise_close_100
```

Format

A matrix with 100 rows and 4 variables. The first two variables make the bivariate data; the last variable refers to cluster memberships; and the last variable refers to outlier indication (1 means outlier, 0 otherwise). The first 30 rows belong to cluster 1, and the last 70 rows belong to cluster 2

d1 variable 1.

d2 variable 2.

cluster cluster memberships

outlier outlier indication

Source

Punzo, A. and McNicholas, P.D., 2016. Parsimonious mixtures of multivariate contaminated normal distributions. *Biometrical Journal*, 58(6), pp.1506-1537.

nm_1_noise_close_500 *A Mixture of Two Close Normal Distributions with 1 by High Atypical Points - 500 Observations*

Description

A simulated mixture of two close normal distributions with 1 substituted by high atypical points. Refer to Punzo and McNicholas (2016) for more information about the underlying distribution that generates this data set.

Usage

```
nm_1_noise_close_500
```

Format

A matrix with 500 rows and 4 variables. The first two variables make the bivariate data; the last variable refers to cluster memberships; and the last variable refers to outlier indication (1 means outlier, 0 otherwise). The first 150 rows belong to cluster 1, and the last 350 rows belong to cluster 2

d1 variable 1.

d2 variable 2.

cluster cluster memberships

outlier outlier indication

Source

Punzo, A. and McNicholas, P.D., 2016. Parsimonious mixtures of multivariate contaminated normal distributions. *Biometrical Journal*, 58(6), pp.1506-1537.

nm_1_noise_far_100	<i>A Mixture of Two Far Normal Distributions with 1 by High Atypical Points - 100 Observations</i>
--------------------	--

Description

A simulated mixture of two far normal distributions with 1 substituted by high atypical points. Refer to Punzo and McNicholas (2016) for more information about the underlying distribution that generates this data set.

Usage

nm_1_noise_far_100

Format

A matrix with 100 rows and 3 variables. The first two variables make the bivariate data, while the last variable refers to cluster memberships. The first 30 rows belong to cluster 1, and the last 70 rows belong to cluster 2

d1 variable 1.

d2 variable 2.

cluster cluster memberships

Source

Punzo, A. and McNicholas, P.D., 2016. Parsimonious mixtures of multivariate contaminated normal distributions. *Biometrical Journal*, 58(6), pp.1506-1537.

nm_1_noise_far_500	<i>A Mixture of Two Far Normal Distributions with 1 by High Atypical Points - 500 Observations</i>
--------------------	--

Description

A simulated mixture of two far normal distributions with 1 substituted by high atypical points. Refer to Punzo and McNicholas (2016) for more information about the underlying distribution that generates this data set.

Usage

```
nm_1_noise_far_500
```

Format

A matrix with 500 rows and 3 variables. The first two variables make the bivariate data, while the last variable refers to cluster memberships. The first 150 rows belong to cluster 1, and the last 350 rows belong to cluster 2

d1 variable 1.

d2 variable 2.

cluster cluster memberships

Source

Punzo, A. and McNicholas, P.D., 2016. Parsimonious mixtures of multivariate contaminated normal distributions. *Biometrical Journal*, 58(6), pp.1506-1537.

nm_5_noise_close_100	<i>A Mixture of Two Close Normal Distributions with 5 by Noise - 100 Observations</i>
----------------------	---

Description

A simulated mixture of two close normal distributions with 1 substituted by noise. Refer to Punzo and McNicholas (2016) for more information about the underlying distribution that generates this data set.

Usage

```
nm_5_noise_close_100
```


Format

A matrix with 100 rows and 4 variables. The first two variables make the bivariate data; the last variable refers to cluster memberships; and the last variable refers to outlier indication (1 means outlier, 0 otherwise). The first 30 rows belong to cluster 1, and the last 70 rows belong to cluster 2

d1 variable 1.

d2 variable 2.

cluster cluster memberships

outlier outlier indication

Source

Punzo, A. and McNicholas, P.D., 2016. Parsimonious mixtures of multivariate contaminated normal distributions. *Biometrical Journal*, 58(6), pp.1506-1537.

nm_5_noise_close_500 *A Mixture of Two Close Normal Distributions with 5 by Noise - 500 Observations*

Description

A simulated mixture of two close normal distributions with 1 substituted by noise. Refer to Punzo and McNicholas (2016) for more information about the underlying distribution that generates this data set.

Usage

nm_5_noise_close_500

Format

A matrix with 500 rows and 4 variables. The first two variables make the bivariate data; the last variable refers to cluster memberships; and the last variable refers to outlier indication (1 means outlier, 0 otherwise). The first 150 rows belong to cluster 1, and the last 350 rows belong to cluster 2

d1 variable 1.

d2 variable 2.

cluster cluster memberships

outlier outlier indication

Source

Punzo, A. and McNicholas, P.D., 2016. Parsimonious mixtures of multivariate contaminated normal distributions. *Biometrical Journal*, 58(6), pp.1506-1537.

nm_5_noise_far_100	<i>A Mixture of Two Far Normal Distributions with 5 by Noise - 100 Observations</i>
--------------------	---

Description

A simulated mixture of two far normal distributions with 1 substituted by noise. Refer to Punzo and McNicholas (2016) for more information about the underlying distribution that generates this data set.

Usage

nm_5_noise_far_100

Format

A matrix with 100 rows and 4 variables. The first two variables make the bivariate data; the last variable refers to cluster memberships; and the last variable refers to outlier indication (1 means outlier, 0 otherwise). The first 30 rows belong to cluster 1, and the last 70 rows belong to cluster 2

d1 variable 1.

d2 variable 2.

cluster cluster memberships

outlier outlier indication

Source

Punzo, A. and McNicholas, P.D., 2016. Parsimonious mixtures of multivariate contaminated normal distributions. *Biometrical Journal*, 58(6), pp.1506-1537.

nm_5_noise_far_500	<i>A Mixture of Two Far Normal Distributions with 5 by Noise - 500 Observations</i>
--------------------	---

Description

A simulated mixture of two far normal distributions with 1 substituted by noise. Refer to Punzo and McNicholas (2016) for more information about the underlying distribution that generates this data set.

Usage

nm_5_noise_far_500

Format

A matrix with 500 rows and 4 variables. The first two variables make the bivariate data; the last variable refers to cluster memberships; and the last variable refers to outlier indication (1 means outlier, 0 otherwise). The first 150 rows belong to cluster 1, and the last 350 rows belong to cluster 2

d1 variable 1.

d2 variable 2.

cluster cluster memberships

outlier outlier indication

Source

Punzo, A. and McNicholas, P.D., 2016. Parsimonious mixtures of multivariate contaminated normal distributions. *Biometrical Journal*, 58(6), pp.1506-1537.

plot.MixtureMissing *Mixture Missing Plotting*

Description

Provide a parallel plot of up to the first 10 variables of a multivariate data sets, and a line plot showing log-likelihood values at every iteration during the EM algorithm. When applicable, pairwise scatter plots highlighting outliers and/or observations whose values are missing but are replaced by expectations obtained in the EM algorithm will be included.

Usage

```
## S3 method for class 'MixtureMissing'  
plot(x, ...)
```

Arguments

x A MixtureMissing object.
... Arguments to be passed to methods, such as graphical parameters.

Value

No return value, called to visualize the fitted model's results

Examples

```
data('nm_5_noise_close_100')

##### With no missing values #####

X <- nm_5_noise_close_100[, 1:2]
mod <- MCNM(X, G = 2, init_method = 'kmedoids', max_iter = 10)
plot(mod)

##### With missing values #####

set.seed(1234)

X <- hide_values(nm_5_noise_close_100[, 1:2], prop_cases = 0.1)
mod <- MCNM(X, G = 2, init_method = 'kmedoids', max_iter = 10)
plot(mod)
```

summary.MixtureMissing

Summary for Mixture Missing

Description

Summarizes main information regarding a MixtureMissing object.

Usage

```
## S3 method for class 'MixtureMissing'
summary(object, ...)
```

Arguments

`object` A MixtureMissing object.
`...` Arguments to be passed to methods, such as graphical parameters.

Details

Information includes the model used to fit the data set, initialization method, clustering table, total outliers, outliers per cluster, mixing proportions, component means and variances, final log-likelihood value, information criteria.

Value

No return value, called to summarize the fitted model's results

Examples

```
data('nm_5_noise_close_100')

##### With no missing values #####

# X <- nm_5_noise_close_100[, 1:2]
# mod <- MCNM(X, G = 2, init_method = 'kmedoids', max_iter = 10)
# summary(mod)

##### With missing values #####

set.seed(1234)

X <- hide_values(nm_5_noise_close_100[, 1:2], prop_cases = 0.1)
mod <- MCNM(X, G = 2, init_method = 'kmedoids', max_iter = 10)
summary(mod)
```

tm_close_100

A Mixture of Two Close Student's t Distributions - 100 Observations

Description

A simulated mixture of two close Student's t distributions. Refer to Punzo and McNicholas (2016) for more information about the underlying distribution that generates this data set.

Usage

```
tm_close_100
```

Format

A matrix with 100 rows and 3 variables. The first two variables make the bivariate data, while the last variable refers to cluster memberships. The first 30 rows belong to cluster 1, and the last 70 rows belong to cluster 2

d1 variable 1.

d2 variable 2.

cluster cluster memberships

Source

Punzo, A. and McNicholas, P.D., 2016. Parsimonious mixtures of multivariate contaminated normal distributions. *Biometrical Journal*, 58(6), pp.1506-1537.

`tm_close_500`*A Mixture of Two Close Student's t Distributions - 500 Observations*

Description

A simulated mixture of two close Student's t distributions. Refer to Punzo and McNicholas (2016) for more information about the underlying distribution that generates this data set.

Usage`tm_close_500`**Format**

A matrix with 500 rows and 3 variables. The first two variables make the bivariate data, while the last variable refers to cluster memberships. The first 150 rows belong to cluster 1, and the last 350 rows belong to cluster 2

d1 variable 1.

d2 variable 2.

cluster cluster memberships

Source

Punzo, A. and McNicholas, P.D., 2016. Parsimonious mixtures of multivariate contaminated normal distributions. *Biometrical Journal*, 58(6), pp.1506-1537.

`tm_far_100`*A Mixture of Two Far Student's t Distributions - 100 Observations*

Description

A simulated mixture of two far Student's t distributions. Refer to Punzo and McNicholas (2016) for more information about the underlying distribution that generates this data set.

Usage`tm_far_100`

Format

A matrix with 500 rows and 4 variables. The first two variables make the bivariate data; the last variable refers to cluster memberships; and the last variable refers to outlier indication (1 means outlier, 0 otherwise). The first 30 rows belong to cluster 1, and the last 70 rows belong to cluster 2

d1 variable 1.

d2 variable 2.

cluster cluster memberships

outlier outlier indication

Source

Punzo, A. and McNicholas, P.D., 2016. Parsimonious mixtures of multivariate contaminated normal distributions. *Biometrical Journal*, 58(6), pp.1506-1537.

tm_far_500

A Mixture of Two Far Student's t Distributions - 500 Observations

Description

A simulated mixture of two far Student's t distributions. Refer to Punzo and McNicholas (2016) for more information about the underlying distribution that generates this data set.

Usage

tm_far_500

Format

A matrix with 500 rows and 3 variables. The first two variables make the bivariate data, while the last variable refers to cluster memberships. The first 150 rows belong to cluster 1, and the last 350 rows belong to cluster 2

d1 variable 1.

d2 variable 2.

cluster cluster memberships

Source

Punzo, A. and McNicholas, P.D., 2016. Parsimonious mixtures of multivariate contaminated normal distributions. *Biometrical Journal*, 58(6), pp.1506-1537.

Index

* datasets

- auto, [2](#)
- bankruptcy, [4](#)
- cnm_close_100, [5](#)
- cnm_close_500, [6](#)
- cnm_far_100, [6](#)
- cnm_far_500, [7](#)
- nm_1_noise_close_100, [29](#)
- nm_1_noise_close_500, [30](#)
- nm_1_noise_far_100, [31](#)
- nm_1_noise_far_500, [32](#)
- nm_5_noise_close_100, [32](#)
- nm_5_noise_close_500, [33](#)
- nm_5_noise_far_100, [34](#)
- nm_5_noise_far_500, [34](#)
- tm_close_100, [37](#)
- tm_close_500, [38](#)
- tm_far_100, [38](#)
- tm_far_500, [39](#)

auto, [2](#)

bankruptcy, [4](#)

cluster_impute, [4](#)
cnm_close_100, [5](#)
cnm_close_500, [6](#)
cnm_far_100, [6](#)
cnm_far_500, [7](#)

evaluation_metrics, [8](#)

generate_patterns, [9](#)

hide_values, [10](#)

initialize_clusters, [11](#)

MCNM, [12](#)
mean_impute, [16](#)
MGHM, [17](#)

MNM, [20](#)

MStM, [23](#)

MtM, [26](#)

nm_1_noise_close_100, [29](#)
nm_1_noise_close_500, [30](#)
nm_1_noise_far_100, [31](#)
nm_1_noise_far_500, [32](#)
nm_5_noise_close_100, [32](#)
nm_5_noise_close_500, [33](#)
nm_5_noise_far_100, [34](#)
nm_5_noise_far_500, [34](#)

plot.MixtureMissing, [35](#)

summary.MixtureMissing, [36](#)

tm_close_100, [37](#)
tm_close_500, [38](#)
tm_far_100, [38](#)
tm_far_500, [39](#)