

Package: tclust (via r-universe)

September 6, 2024

Type Package

Encoding UTF-8

Title Robust Trimmed Clustering

Version 2.0-4

VersionNote Released 2.0-3 on 2024-04-17 on CRAN

Maintainer Valentin Todorov <valentin.todorov@chello.at>

Description Provides functions for robust trimmed clustering. The methods are described in Garcia-Escudero (2008) <[doi:10.1214/07-AOS515](https://doi.org/10.1214/07-AOS515)>, Fritz et al. (2012) <[doi:10.18637/jss.v047.i12](https://doi.org/10.18637/jss.v047.i12)>, Garcia-Escudero et al. (2011) <[doi:10.1007/s11222-010-9194-z](https://doi.org/10.1007/s11222-010-9194-z)> and others.

Depends R(>= 3.6.2)

Imports Rcpp (>= 1.0.7), doParallel, parallel, foreach, MASS

Suggests mclust, cluster, sn

LazyLoad yes

License GPL-3

LinkingTo Rcpp, RcppArmadillo

NeedsCompilation yes

RoxygenNote 7.2.3

URL <https://github.com/valentint/tclust>

BugReports <https://github.com/valentint/tclust/issues>

Author Valentin Todorov [aut, cre]
(<<https://orcid.org/0000-0003-4215-0245>>), Luis Angel García Escudero [aut], Agustín Mayo Iscar [aut], Javier Crespo Guerrero [aut], Heinrich Fritz [aut]

Repository <https://valentint.r-universe.dev>

RemoteUrl <https://github.com/valentint/tclust>

RemoteRef HEAD

RemoteSha 6d0eb288cb97cd322fb0beb8dd7942e793b81fca

Contents

ctlcurves	2
DiscrFact	4
geyser2	5
LG5data	6
M5data	7
pine	8
plot.ctlcurves	8
plot.DiscrFact	10
plot.rlg	11
plot.tclust	12
rlg	14
simula.rlg	16
simula.tclust	17
summary.DiscrFact	18
swissbank	19
tclust	20
tkmeans	25
wholesale	28
Index	30

 ctlcurves

Classification Trimmed Likelihood Curves

Description

The function applies `tclust` several times on a given dataset while parameters `alpha` and `k` are altered. The resulting object gives an idea of the optimal trimming level and number of clusters considering a particular dataset.

Usage

```
ctlcurves(
  x,
  k = 1:4,
  alpha = seq(0, 0.2, len = 6),
  restr.fact = 50,
  parallel = FALSE,
  trace = 1,
  ...
)
```

Arguments

x	A matrix or data frame of dimension $n \times p$, containing the observations (row-wise).
k	A vector of cluster numbers to be checked. By default cluster numbers from 1 to 5 are examined.
alpha	A vector containing the alpha levels to be checked. By default alpha levels from 0 to 0.2 (continuously increased by 0.01), are checked.
restr.fact	The restriction factor passed to <code>tclust</code> .
parallel	A logical value, to be passed further to <code>tclust()</code> .
trace	Defines the tracing level, which is set to 1 by default. Tracing level 2 gives additional information on the current iteration.
...	Further arguments (as e.g. <code>restr</code>), passed to <code>tclust</code>

Details

These curves show the values of the trimmed classification (log-)likelihoods when altering the trimming proportion alpha and the number of clusters k. The careful examination of these curves provides valuable information for choosing these parameters in a clustering problem. For instance, an appropriate k to be chosen is one that we do not observe a clear increase in the trimmed classification likelihood curve for k with respect to the k+1 curve for almost all the range of alpha values. Moreover, an appropriate choice of parameter alpha may be derived by determining where an initial fast increase of the trimmed classification likelihood curve stops for the final chosen k. A more detailed explanation can be found in García-Escudero et al. (2011).

Value

The function returns an S3 object of type `ctlcurves` containing the following components:

- `par` A list containing all the parameters passed to this function
- `obj` An array containing the objective functions values of each computed cluster-solution
- `min.weights` An array containing the minimum cluster weight of each computed cluster-solution

References

García-Escudero, L.A.; Gordaliza, A.; Matrán, C. and Mayo-Iscar, A. (2011), "Exploring the number of groups in robust model-based clustering." *Statistics and Computing*, **21** pp. 585-599, <doi:10.1007/s11222-010-9194-z>

Examples

```
## Not run:

#--- EXAMPLE 1 -----

sig <- diag (2)
cen <- rep (1, 2)
```

```

x <- rbind(MASS::mvrnorm(108, cen * 0, sig),
           MASS::mvrnorm(162, cen * 5, sig * 6 - 2),
           MASS::mvrnorm(30, cen * 2.5, sig * 50))

ctl <- ctlcurves(x, k = 1:4)
ctl

## ctl-curves
plot(ctl) ## --> selecting k = 2, alpha = 0.08

## the selected model
plot(tclust(x, k = 2, alpha = 0.08, restr.fact = 7))

#--- EXAMPLE 2 -----

data(geyser2)
ctl <- ctlcurves(geyser2, k = 1:5)
ctl

## ctl-curves
plot(ctl) ## --> selecting k = 3, alpha = 0.08

## the selected model
plot(tclust(geyser2, k = 3, alpha = 0.08, restr.fact = 5))

#--- EXAMPLE 3 -----

data(swissbank)
ctl <- ctlcurves(swissbank, k = 1:5, alpha = seq(0, 0.3, by = 0.025))
ctl

## ctl-curves
plot(ctl) ## --> selecting k = 2, alpha = 0.1

## the selected model
plot(tclust(swissbank, k = 2, alpha = 0.1, restr.fact = 50))

## End(Not run)

```

DiscrFact

Discriminant Factor analysis for tclust objects

Description

Analyzes a `tclust`-object by calculating discriminant factors and comparing the quality of the actual cluster assignments to that of the second best possible assignment for each observation.

Cluster assignments of observations with large discriminant factors are considered "doubtful" decisions. Silhouette plots give a graphical overview of the discriminant factors distribution (see [plot.DiscrFact](#)). More details can be found in García-Escudero et al. (2011).

Usage

```
DiscrFact(x, threshold = 1/10)
```

Arguments

x	A tclust object.
threshold	A cluster assignment or a trimming decision for an observation with a discriminant factor larger than $\log(\text{threshold})$ is considered a "doubtful" decision.

Value

The function returns an S3 object of type `DiscrFact` containing the following components:

- x A tclust object.
- ylimin A minimum y-limit calculated for plotting purposes.
- ind The actual cluster assignment.
- ind2 The second most likely cluster assignment for each observation.
- lik The (weighted) likelihood of the actual cluster assignment of each observation.
- lik2 The (weighted) likelihood of the second best cluster assignment of each observation.
- assignfact The factor $\log(\text{disc}/\text{disc2})$.
- threshold The threshold used for deciding whether `assignfact` indicates a "doubtful" assignment.
- mean.DiscrFact A vector of length $k + 1$ containing the mean discriminant factors for each cluster (including the outliers).

References

García-Escudero, L.A.; Gordaliza, A.; Matrán, C. and Mayo-Iscar, A. (2011), "Exploring the number of groups in robust model-based clustering." *Statistics and Computing*, **21** pp. 585-599, <doi:10.1007/s11222-010-9194-z>

Description

A bivariate data set obtained from the Old Faithful Geyser, containing the eruption length and the length of the previous eruption for 271 eruptions of this geyser in minutes.

Usage

```
data(geyser2)
```

Format

A data frame containing 272 observations in 2 variables. The variables are as follows:

- Eruption length The eruption length in minutes.
- Previous eruption length The length of the previous eruption in minutes.

Source

This particular data structure can be obtained by applying the following code to the "Old Faithful Geysers" (faithful data set (Härdle 1991) in the package datasets):

```
f1 <- faithful[, 1]
geyser2 <- cbind(f1[-length(f1)], f1[-1])
colnames(geyser2) <- c("Eruption length",
"Previous eruption length")
```

References

García-Escudero, L.A. and Gordaliza, A. (1999). Robustness properties of k-means and trimmed k-means. *Journal of the American Statistical Assoc.*, Vol.94, No.447, 956–969.

Härdle, W. (1991). *Smoothing Techniques with Implementation in S.*, New York: Springer.

LG5data

LG5data data

Description

A data set in dimension 10 with three clusters around affine subspaces of common intrinsic dimension. A 10% background noise is added uniformly distributed in a rectangle containing the three main clusters.

Usage

```
data(LG5data)
```

Format

The first 10 columns are the variables. The last column is the true classification vector where symbol "0" stands for the contaminating data points.

Examples

```
#--- EXAMPLE 1 -----
data (LG5data)
x <- LG5data[, 1:10]
clus <- rlg(x, d = c(2,2,2), alpha=0.1, trace=TRUE)
plot(x, col=clus$cluster+1)
```

M5data

M5data data

Description

A bivariate data set obtained from three normal bivariate distributions with different scales and proportions 1:2:2. One of the components is very overlapped with another one. A 10% background noise is added uniformly distributed in a rectangle containing the three normal components and not very overlapped with the three mixture components. A precise description of the M5 data set can be found in García-Escudero et al. (2008).

Usage

```
data(M5data)
```

Format

The first two columns are the two variables. The last column is the true classification vector where symbol "0" stands for the contaminating data points.

Source

García-Escudero, L.A.; Gordaliza, A.; Matrán, C. and Mayo-Isacar, A. (2008), "A General Trimming Approach to Robust Cluster Analysis". *Annals of Statistics*, Vol.36, pp. 1324-1345.

Examples

```
#--- EXAMPLE 1 -----
data (M5data)
x <- M5data[, 1:2]
clus <- tclust(x, k=3, alpha=0.1, nstart=200, niter1=3, niter2=17,
  nkeep=10, opt="HARD", equal.weights=FALSE, restr.fact=50, trace=TRUE)
plot (x, col=clus$cluster+1)
```

pine

Pinus nigra dataset

Description

To study the growth of the wood mass in a cultivated forest of *Pinus nigra* located in the north of Palencia (Spain), a sample of 362 trees was studied. The data set is made of measurements of heights (in meters), in variable "HT", and diameters (in millimetres), in variable "Diameter", of these trees. The presence of three linear groups can be guessed apart from a small group of trees forming its own cluster with larger heights and diameters one isolated tree with the largest diameter but small height. More details on the interpretation of this dataset in García-Escudero et al (2010).

Usage

```
data(pine)
```

Format

A data frame containing 362 observations in 2 variables. The variables are as follows:

- Diameter Diameter
- HT Height

References

García-Escudero, L. A., Gordaliza, A., Mayo-Isar, A., and San Martín, R. (2010). Robust clusterwise linear regression through trimming. *Computational Statistics & Data Analysis*, 54(12), 3057–3069.

plot.ctlcurves

The plot method for objects of class ctlcurves

Description

The plot method for class `ctlcurves`: This function implements a series of plots, which display characteristic values of the each model, computed with different values for `k` and `alpha`.

Usage

```
## S3 method for class 'ctlcurves'  
plot(  
  x,  
  what = c("obj", "min.weights", "doubtful"),  
  main,  
  xlab,
```



```

    ylab,
    xlim,
    ylim,
    col,
    lty = 1,
    ...
)

```

Arguments

x	The ctlcurves object to be shown
what	A string indicating which type of plot shall be drawn. See the details section for more information.
main	A character-string containing the title of the plot.
xlab, ylab, xlim, ylim	Arguments passed to plot().
col	A single value or vector of line colors passed to lines .
lty	A single value or vector of line colors passed to lines .
...	Arguments to be passed to or from other methods.

Details

These curves show the values of the trimmed classification (log-)likelihoods when altering the trimming proportion α and the number of clusters k . The careful examination of these curves provides valuable information for choosing these parameters in a clustering problem. For instance, an appropriate k to be chosen is one that we do not observe a clear increase in the trimmed classification likelihood curve for k with respect to the $k+1$ curve for almost all the range of α values. Moreover, an appropriate choice of parameter α may be derived by determining where an initial fast increase of the trimmed classification likelihood curve stops for the final chosen k . A more detailed explanation can be found in García-Escudero et al. (2011).

This function implements a series of plots, which display characteristic values of the each model, computed with different values for k and α .

"obj" Objective function values.

"min.weights" The minimum cluster weight found for each computed model. This plot is intended to spot spurious clusters, which in general yield quite small weights.

"doubtful" The number of "doubtful" decisions identified by [DiscrFact](#).

References

García-Escudero, L.A.; Gordaliza, A.; Matrán, C. and Mayo-Iscar, A. (2011), "Exploring the number of groups in robust model-based clustering." *Statistics and Computing*, **21** pp. 585-599, <doi:10.1007/s11222-010-9194-z>

Examples

```
#--- EXAMPLE 1 -----

sig <- diag (2)
cen <- rep (1, 2)
x <- rbind(MASS::mvrnorm(108, cen * 0, sig),
          MASS::mvrnorm(162, cen * 5, sig * 6 - 2),
          MASS::mvrnorm(30, cen * 2.5, sig * 50))

(ctl <- ctlcurves(x, k = 1:4))

plot(ctl)
```

plot.DiscrFact *The plot method for objects of class DiscrFact*

Description

The plot method for class DiscrFact: Next to a plot of the tclust object which has been used for creating the DiscrFact object, a silhouette plot indicates the presence of groups with a large amount of doubtfully assigned observations. A third plot similar to the standard tclust plot serves to highlight the identified doubtful observations.

Usage

```
## S3 method for class 'DiscrFact'
plot(
  x,
  enum.plots = FALSE,
  xlab = "Discriminant Factor",
  ylab = "Clusters",
  print.DiscrFact = TRUE,
  xlim,
  col.nodoubt = grey(0.8),
  by.cluster = FALSE,
  ...
)
```

Arguments

x	An object of class DiscrFact as returned from DiscrFact()
enum.plots	A logical value indicating whether the plots shall be enumerated in their title ("(a)", "(b)", "(c)").
xlab, ylab, xlim	Arguments passed to function plot.tclust()
print.DiscrFact	A logical value indicating whether each clusters mean discriminant factor shall be plotted

col.nodoubt	Color of all observations not considered as to be assigned doubtfully.
by.cluster	Logical value indicating whether optional parameters pch and col (if present) refer to observations (FALSE) or clusters (TRUE)
...	Arguments to be passed to or from other methods

Details

plot_DiscrFact_p2 displays a silhouette plot based on the discriminant factors of the observations. A solution with many large discriminant factors is not reliable. Such clusters can be identified with this silhouette plot. Thus plot_DiscrFact_p3 displays the dataset, highlighting observations with discriminant factors greater than the given threshold. The function plot.DiscrFact() combines the standard plot of a tclust object, and the two plots introduced here.

References

García-Escudero, L.A.; Gordaliza, A.; Matrán, C. and Mayo-Iscar, A. (2011), "Exploring the number of groups in robust model-based clustering." *Statistics and Computing*, **21** pp. 585-599, <doi:10.1007/s11222-010-9194-z>

Examples

```
sig <- diag (2)
cen <- rep (1, 2)
x <- rbind(MASS::mvrnorm(360, cen * 0, sig),
          MASS::mvrnorm(540, cen * 5, sig * 6 - 2),
          MASS::mvrnorm(100, cen * 2.5, sig * 50))

clus.1 <- tclust(x, k = 2, alpha=0.1, restr.fact=12)
clus.2 <- tclust(x, k = 3, alpha=0.1, restr.fact=1)

dsc.1 <- DiscrFact(clus.1)
plot(dsc.1)

dsc.2 <- DiscrFact(clus.2)
plot(dsc.2)
```

plot.rlg

Plot an 'rlg' object

Description

Different plots for the results of 'rlg' analysis, stored in an rlg object, see Details.

Usage

```
## S3 method for class 'rlg'
plot(
  x,
  which = c("all", "scores", "loadings", "eigenvalues"),
  sort = TRUE,
  ask = (which == "all" && dev.interactive(TRUE)),
  ...
)
```

Arguments

x	An rlg object to plot.
which	Select the required plot.
sort	Whether to sort.
ask	if TRUE, the user is <i>asked</i> before each plot, see par(ask=.). Default is ask = which=="all" && dev.interactive().
...	Other parameters to be passed to the lower level functions.

Examples

```
data (LG5data)
x <- LG5data[, 1:10]
clus <- rlg(x, d = c(2,2,2), alpha=0.1)
plot(clus, which="eigenvalues")
plot(clus, which="scores")
```

plot.tclust

Plot Method for tclust and tkmeans Objects

Description

One and two dimensional structures are treated separately (e.g. tolerance intervals/ellipses are displayed). Higher dimensional structures are displayed by plotting the two first Fisher's canonical coordinates (evaluated by `tclust::discr_coords`) and derived from the final cluster assignments (trimmed observations are not taken into account). `plot.tclust.Nd` can be called with one or two-dimensional `tclust`- or `tkmeans`-objects too. The function fails, if `store.x = FALSE` is specified in the `tclust()` or `tkmeans()` call, because the original data matrix is required here.

Usage

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

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

Arguments

- x The tclust or tkmeans object to be displayed
- ... Further (optional) arguments which specify the details of the resulting plot (see section "Further Arguments").

Details

The plot method for classes tclust and tkmeans.

Further Arguments

- xlab, ylab, xlim, ylim, pch, col Arguments passed to plot().
- main The title of the plot. Use "/p" for displaying the chosen parameters alpha and k or "/r" for plotting the chosen restriction.
- main.pre An optional string which is added to the plot's caption.
- sub A string specifying the subtitle of the plot. Use "/p" (default) for displaying the chosen parameters alpha and k, "/r" for plotting the chosen restriction and "/pr" for both.
- sub1 A secondary (optional) subtitle.
- labels A string specifying the type of labels to be drawn. Either labels="none" (default), labels="cluster" or labels="observation" can be specified. If specified, parameter pch is ignored.
- text A vector of length n (the number of observations) containing strings which are used as labels for each observation. If specified, the parameters labels and pch are ignored.
- by.cluster Logical value indicating whether parameters pch and col refer to observations (FALSE) or clusters (TRUE).
- jitter.y Logical value, specifying whether the drawn values shall be jittered in y-direction for better visibility of structures in 1 dimensional data.
- tol The tolerance interval. 95% tolerance ellipsoids (assuming normality) are plotted by default.
- tol.col, tol.lty, tol.lwd Vectors of length k or 1 containing the col, lty and lwd arguments for the tolerance ellipses/lines.

Examples

```
#--- EXAMPLE 1-----
sig <- diag (2)
cen <- rep (1, 2)
x <- rbind(MASS::mvrnorm(360, cen * 0, sig),
           MASS::mvrnorm(540, cen * 5, sig * 6 - 2),
           MASS::mvrnorm(100, cen * 2.5, sig * 50))
# Two groups and 10% trimming level
a <- tclust(x, k = 2, alpha = 0.1, restr.fact = 12)
plot (a)
plot (a, labels = "observation")
plot (a, labels = "cluster")
plot (a, by.cluster = TRUE)
```

```
#--- EXAMPLE 2-----
sig <- diag (2)
cen <- rep (1, 2)
x <- rbind(MASS::mvrnorm(360, cen * 0, sig),
           MASS::mvrnorm(540, cen * 5, sig),
           MASS::mvrnorm(100, cen * 2.5, sig))
# Two groups and 10% trimming level
a <- tkmeans(x, k = 2, alpha = 0.1)
plot (a)
plot (a, labels = "observation")
plot (a, labels = "cluster")
plot (a, by.cluster = TRUE)
```

rlg

*Robust Linear Grouping***Description**

The function `rlg()` searches for clusters around affine subspaces of dimensions given by vector `d` (the length of that vector is the number of clusters). For instance `d=c(1,2)` means that we are clustering around a line and a plane. For robustifying the estimation, a proportion `alpha` of observations is trimmed. In particular, the trimmed k-means method is represented by the `rlg` method, if `d=c(0,0,...0)` (a vector of length `k` with zeroes).

Usage

```
rlg(
  x,
  d,
  alpha = 0.05,
  nstart = 500,
  niter1 = 3,
  niter2 = 20,
  nkeep = 5,
  scale = FALSE,
  parallel = FALSE,
  n.cores = -1,
  trace = FALSE
)
```

Arguments

- `x` A matrix or data.frame of dimension $n \times p$, containing the observations (row-wise).
- `d` A numeric vector of length equal to the number of clusters to be detected. Each component of vector `d` indicates the intrinsic dimension of the affine subspace where observations on that cluster are going to be clustered. All the elements of vector `d` should be smaller than the problem dimension minus 1.

<code>alpha</code>	The proportion of observations to be trimmed.
<code>nstart</code>	The number of random initializations to be performed.
<code>niter1</code>	The number of concentration steps to be performed for the <code>nstart</code> initializations.
<code>niter2</code>	The maximum number of concentration steps to be performed for the <code>nkeep</code> solutions kept for further iteration. The concentration steps are stopped, whenever two consecutive steps lead to the same data partition.
<code>nkeep</code>	The number of iterated initializations (after <code>niter1</code> concentration steps) with the best values in the target function that are kept for further iterations
<code>scale</code>	A robust centering and scaling (using the median and MAD) is done if TRUE.
<code>parallel</code>	A logical value, specifying whether the <code>nstart</code> initializations should be done in parallel.
<code>n.cores</code>	The number of cores to use when parallellizing, only taken into account if <code>parallel=T</code> .
<code>trace</code>	Defines the tracing level, which is set to 0 by default. Tracing level 1 gives additional information on the stage of the iterative process.

Details

The procedure allows to deal with robust clustering around affine subspaces with an alpha proportion of trimming level by minimizing the trimmed sums of squared orthogonal residuals. Each component of vector `d` indicates the intrinsic dimension of the affine subspace where observations on that cluster are going to be clustered. Therefore a component equal to 0 on that vector implies clustering around centres, equal to 1 around lines, equal to 2 around planes and so on. The procedure so allows simultaneous clustering and dimensionality reduction.

This iterative algorithm performs "concentration steps" to improve the current cluster assignments. For approximately obtaining the global optimum, the procedure is randomly initialized `nstart` times and `niter1` concentration steps are performed for them. The `nkeep` most "promising" iterations, i.e. the `nkeep` iterated solutions with the initial best values for the target function, are then iterated until convergence or until `niter2` concentration steps are done.

Value

Returns an object of class `rlg` which is basically a list with the following elements:

- `centers` - A matrix of size $p \times k$ containing the location vectors (column-wise) of each cluster.
- `U` - A list with k elements where each element is $p \times d_j$ matrix whose d_j columns are unitary and orthogonal vectors generating the affine subspace (after subtracting the corresponding cluster's location parameter in `centers`). d_j is the intrinsic dimension of the affine subspace approximation in the j -th cluster, i.e., the elements of vector `d`.
- `cluster` - A numerical vector of size n containing the cluster assignment for each observation. Cluster names are integer numbers from 1 to k , 0 indicates trimmed observations.
- `obj` - The value of the objective function of the best (returned) solution.
- `cluster.ini` - A matrix with `nstart` rows and number of columns equal to the number of observations and where each row shows the final clustering assignments (0 for trimmed observations) obtained after the `niter1` iteration of the `nstart` random initializations.

- obj.ini -A numerical vector of length nstart containing the values of the target function obtained after the niter1 iteration of the nstart random initializations.
- x - The input data set.
- dimensions - The input d vector with the intrinsic dimensions. The number of clusters is the length of that vector.
- alpha - The input trimming level.

Author(s)

Javier Crespo Guerrero, Jesús Fernández Iglesias, Luis Angel Garcia Escudero, Agustin Mayo Iscar.

References

García-Escudero, L. A., Gordaliza, A., San Martín, R., Van Aelst, S., & Zamar, R. (2009). Robust linear clustering. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71, 301-318.

Examples

```
##--- EXAMPLE 1 -----
data (LG5data)
x <- LG5data[, 1:10]
clus <- rlg(x, d = c(2,2,2), alpha=0.1)
plot(x, col=clus$cluster+1)
plot(clus, which="eigenvalues")
plot(clus, which="scores")

##--- EXAMPLE 2 -----
data (pine)
clus <- rlg(pine, d = c(1,1,1), alpha=0.035)
plot(pine, col=clus$cluster+1)
```

simula.rlg

Simulate contaminated data set for applying rlg

Description

Simulate $\alpha \cdot 100\%$ contaminated data set for applying rlg by generating a $k=3$ components with equal size and # common underlying dimension $q_1=q_2=q_3=q$

Usage

```
simula.rlg(q = 2, p = 10, n = 200, var = 0.01, sep.means = 0, alpha = 0.05)
```


Arguments

q	intrinsic dimension
p	dimension ($p \geq 2$ and $p > q$)
n	number of observations
var	The smaller 'var' the smaller the scatter around the lower dimensional space
sep.means	Parameter controlling the location vectors separation
alpha	contamination level

Value

a list with the following items

- x - The generated dataset
- true - The true classification

Examples

```
res <- simula.rlg(q=5, p=200, n=150, var=0.01, sep.means=0.00)
plot(res$x,col=res$true+1)
```

simula.tclust

Simulate contaminated data set for applying TCLUST

Description

Simulate 10% contaminated data set for applying TCLUST

Usage

```
simula.tclust(n, p = 4, k = 3, type = 2, balanced = 1)
```

Arguments

n	number of observations
p	dimension ($p \geq 2$ and $p > q$)
k	number of clusters (only $k=3$ and $k=6$ are allowed!!!)
type	1 (spherical for $\text{rest.fact}=1$) or 2 (elliptical for $\text{rest.fact}=9^2$)
balanced	1 (all clusters equal size) or 2 [proportions (25,30,35)% if $k=3$ and (12.5,15,17.5,12.5,15,17.5)% if $k=6$]

Value

a list with the following items

- x - The generated dataset
- true - The true classification

Examples

```
res <- simula.tclust(n=400,k=3,p=8,type=2,balanced=1)
plot(res$x,col=res$true+1)
```

summary.DiscrFact *The summary method for objects of class DiscrFact*

Description

The summary method for class DiscrFact.

Usage

```
## S3 method for class 'DiscrFact'
summary(object, hide.empty = TRUE, show.clust, show.alt, ...)
```

Arguments

object	An object of class DiscrFact as returned from DiscrFact().
hide.empty	A logical value specifying whether clusters without doubtful assignment shall be hidden.
show.clust	A logical value specifying whether the number of doubtful assignments per cluster shall be displayed.
show.alt	A logical value specifying whether the alternative cluster assignment shall be displayed.
...	Arguments passed to or from other methods.

References

García-Escudero, L.A.; Gordaliza, A.; Matrán, C. and Mayo-Iscar, A. (2011), "Exploring the number of groups in robust model-based clustering." *Statistics and Computing*, **21** pp. 585-599, <doi:10.1007/s11222-010-9194-z>

Examples

```
sig <- diag(2)
cen <- rep(1, 2)
x <- rbind(MASS::mvrnorm(360, cen * 0, sig),
          MASS::mvrnorm(540, cen * 5, sig * 6 - 2),
          MASS::mvrnorm(100, cen * 2.5, sig * 50)
)

clus.1 <- tclust(x, k = 2, alpha=0.1, restr.fact=12)
clus.2 <- tclust(x, k = 3, alpha=0.1, restr.fact=1)

dsc.1 <- DiscrFact(clus.1)
```

```
summary(dsc.1)

dsc.2 <- DiscrFact(clus.2)
summary(dsc.2)
```

swissbank

Swiss banknotes data

Description

Six variables measured on 100 genuine and 100 counterfeit old Swiss 1000-franc bank notes (Flury and Riedwyl, 1988).

Usage

```
data(swissbank)
```

Format

A data frame containing 200 observations in 6 variables. The variables are as follows:

- Length Length of the bank note
- Ht_Left Height of the bank note, measured on the left
- Ht_Right Height of the bank note, measured on the right
- IF_Lower Distance of inner frame to the lower border
- IF_Upper Distance of inner frame to the upper border
- Diagonal Length of the diagonal

Details

Observations 1–100 are the genuine bank notes and the other 100 observations are the counterfeit bank notes.

Source

Flury, B. and Riedwyl, H. (1988). *Multivariate Statistics, A Practical Approach*, Cambridge University Press.

tclust

*TCLUST method for robust clustering***Description**

This function searches for k (or less) clusters with different covariance structures in a data matrix x . Relative cluster scatter can be restricted when `restr="eigen"` by constraining the ratio between the largest and the smallest of the scatter matrices eigenvalues by a constant value `restr.fact`. Relative cluster scatters can be also restricted with `restr="deter"` by constraining the ratio between the largest and the smallest of the scatter matrices' determinants.

For robustifying the estimation, a proportion `alpha` of observations is trimmed. In particular, the trimmed k-means method is represented by the `tclust()` method, by setting parameters `restr.fact=1`, `opt="HARD"` and `equal.weights=TRUE`.

Usage

```
tclust(
  x,
  k,
  alpha = 0.05,
  nstart = 500,
  niter1 = 3,
  niter2 = 20,
  nkeep = 5,
  iter.max,
  equal.weights = FALSE,
  restr = c("eigen", "deter"),
  restr.fact = 12,
  cshape = 1e+10,
  opt = c("HARD", "MIXT"),
  center = FALSE,
  scale = FALSE,
  store_x = TRUE,
  parallel = FALSE,
  n.cores = -1,
  zero_tol = 1e-16,
  drop.empty.clust = TRUE,
  trace = 0
)
```

Arguments

<code>x</code>	A matrix or data.frame of dimension $n \times p$, containing the observations (row-wise).
<code>k</code>	The number of clusters initially searched for.
<code>alpha</code>	The proportion of observations to be trimmed.

nstart	The number of random initializations to be performed.
niter1	The number of concentration steps to be performed for the nstart initializations.
niter2	The maximum number of concentration steps to be performed for the nkeep solutions kept for further iteration. The concentration steps are stopped, whenever two consecutive steps lead to the same data partition.
nkeep	The number of iterated initializations (after niter1 concentration steps) with the best values in the target function that are kept for further iterations
iter.max	(deprecated, use the combination nkeep, niter1 and niter2) The maximum number of concentration steps to be performed. The concentration steps are stopped, whenever two consecutive steps lead to the same data partition.
equal.weights	A logical value, specifying whether equal cluster weights shall be considered in the concentration and assignment steps.
restr	Restriction type to control relative cluster scatters. The default value is restr="eigen", so that the maximal ratio between the largest and the smallest of the scatter matrices eigenvalues is constrained to be smaller then or equal to restr.fact (Garcia-Escudero, Gordaliza, Matran, and Mayo-Iscar, 2008). Alternatively, restr="deter" imposes that the maximal ratio between the largest and the smallest of the scatter matrices determinants is smaller or equal than restr.fact (see Garcia-Escudero, Mayo-Iscar and Riani, 2020)
restr.fact	The constant restr.fact ≥ 1 constrains the allowed differences among group scatters in terms of eigenvalues ratio (if restr="eigen") or determinant ratios (if restr="deter"). Larger values imply larger differences of group scatters, a value of 1 specifies the strongest restriction.
cshape	constraint to apply to the shape matrices, cshape ≥ 1 , (see Garcia-Escudero, Mayo-Iscar and Riani, 2020)). This options only works if restr=="deter". In this case the default value is cshape=1e10 to ensure the procedure is (virtually) affine equivariant. On the other hand, cshape values close to 1 would force the clusters to be almost spherical (without necessarily the same scatters if restr.fact is strictly greater than 1).
opt	Define the target function to be optimized. A classification likelihood target function is considered if opt="HARD" and a mixture classification likelihood if opt="MIXT".
center	Optional centering of the data: a function or a vector of length p which can optionally be specified for centering x before calculation
scale	Optional scaling of the data: a function or a vector of length p which can optionally be specified for scaling x before calculation
store_x	A logical value, specifying whether the data matrix x shall be included in the result object. By default this value is set to TRUE, because some of the plotting functions depend on this information. However, when big data matrices are handled, the result object's size can be decreased noticeably when setting this parameter to FALSE.
parallel	A logical value, specifying whether the nstart initializations should be done in parallel.
n.cores	The number of cores to use when paralellizing, only taken into account if parallel=T.

<code>zero_tol</code>	The zero tolerance used. By default set to 1e-16.
<code>drop.empty.clust</code>	Logical value specifying, whether empty clusters shall be omitted in the resulting object. (The result structure does not contain center and covariance estimates of empty clusters anymore. Cluster names are reassigned such that the first l clusters ($l \leq k$) always have at least one observation.
<code>trace</code>	Defines the tracing level, which is set to 0 by default. Tracing level 1 gives additional information on the stage of the iterative process.

Details

The procedure allows to deal with robust clustering with an α proportion of trimming level and searching for k clusters. We are considering classification trimmed likelihood when using `opt="HARD"` so that "hard" or "crisp" clustering assignments are done. On the other hand, mixture trimmed likelihood are applied when using `opt="MIXT"` so providing a kind of clusters "posterior" probabilities for the observations. Relative cluster scatter can be restricted when `restr="eigen"` by constraining the ratio between the largest and the smallest of the scatter matrices eigenvalues by a constant value `restr.fact`. Setting `restr.fact=1`, yields the strongest restriction, forcing all clusters to be spherical and equally scattered. Relative cluster scatters can be also restricted with `restr="deter"` by constraining the ratio between the largest and the smallest of the scatter matrices' determinants.

This iterative algorithm performs "concentration steps" to improve the current cluster assignments. For approximately obtaining the global optimum, the procedure is randomly initialized `nstart` times and `niter1` concentration steps are performed for them. The `nkeep` most "promising" iterations, i.e. the `nkeep` iterated solutions with the initial best values for the target function, are then iterated until convergence or until `niter2` concentration steps are done.

The parameter `restr.fact` defines the cluster scatter matrices restrictions, which are applied on all clusters during each concentration step. It restricts the ratio between the maximum and minimum eigenvalue of all clusters' covariance structures to that parameter. Setting `restr.fact=1`, yields the strongest restriction, forcing all clusters to be spherical and equally scattered.

Cluster components with similar sizes are favoured when considering `equal.weights=TRUE` while `equal.weights=FALSE` admits possible different prior probabilities for the components and it can easily return empty clusters when the number of clusters is greater than apparently needed.

Value

The function returns the following values:

- `cluster` - A numerical vector of size n containing the cluster assignment for each observation. Cluster names are integer numbers from 1 to k , 0 indicates trimmed observations. Note that it could be empty clusters with no observations when `equal.weights=FALSE`.
- `obj` - The value of the objective function of the best (returned) solution.
- `size` - An integer vector of size k , returning the number of observations contained by each cluster.
- `weights` - Vector of Cluster weights
- `centers` - A matrix of size $p \times k$ containing the centers (column-wise) of each cluster.

- `cov` - An array of size $p \times p \times k$ containing the covariance matrices of each cluster.
- `code` - A numerical value indicating if the concentration steps have converged for the returned solution (2).
- `posterior` - A matrix with k columns that contains the posterior probabilities of membership of each observation (row-wise) to the k clusters. This posterior probabilities are 0-1 values in the `opt="HARD"` case. Trimmed observations have 0 membership probabilities to all clusters.
- `cluster.ini` - A matrix with `nstart` rows and number of columns equal to the number of observations and where each row shows the final clustering assignments (0 for trimmed observations) obtained after the `niter1` iteration of the `nstart` random initializations.
- `obj.ini` - A numerical vector of length `nstart` containing the values of the target function obtained after the `niter1` iteration of the `nstart` random initializations.
- `x` - The input data set.
- `k` - The input number of clusters.
- `alpha` - The input trimming level.

Author(s)

Javier Crespo Guerrero, Luis Angel Garcia Escudero, Agustin Mayo Iscar.

References

- Fritz, H.; Garcia-Escudero, L.A.; Mayo-Iscar, A. (2012), "tclust: An R Package for a Trimming Approach to Cluster Analysis". *Journal of Statistical Software*, 47(12), 1-26. URL <http://www.jstatsoft.org/v47/i12/>
- Garcia-Escudero, L.A.; Gordaliza, A.; Matran, C. and Mayo-Iscar, A. (2008), "A General Trimming Approach to Robust Cluster Analysis". *Annals of Statistics*, Vol.36, 1324–1345.
- García-Escudero, L. A., Gordaliza, A. and Mayo-Íscar, A. (2014). A constrained robust proposal for mixture modeling avoiding spurious solutions. *Advances in Data Analysis and Classification*, 27–43.
- García-Escudero, L. A., and Mayo-Íscar, A. and Riani, M. (2020). Model-based clustering with determinant-and-shape constraint. *Statistics and Computing*, 30, 1363–1380.]

Examples

```
##--- EXAMPLE 1 -----
sig <- diag(2)
cen <- rep(1,2)
x <- rbind(MASS::mvrnorm(360, cen * 0, sig),
          MASS::mvrnorm(540, cen * 5, sig * 6 - 2),
          MASS::mvrnorm(100, cen * 2.5, sig * 50))

## Two groups and 10% trimming level
clus <- tclust(x, k = 2, alpha = 0.1, restr.fact = 8)

plot(clus)
plot(clus, labels = "observation")
plot(clus, labels = "cluster")
```

```

## Three groups (one of them very scattered) and 0% trimming level
clus <- tclus(x, k = 3, alpha=0.0, restr.fact = 100)

plot(clus)

##--- EXAMPLE 2 -----
data(geyser2)
(clus <- tclus(geyser2, k = 3, alpha = 0.03))

plot(clus)

## Not run:

##--- EXAMPLE 3 -----
data(M5data)
x <- M5data[, 1:2]

clus.a <- tclus(x, k = 3, alpha = 0.1, restr.fact = 1,
               restr = "eigen", equal.weights = TRUE)
clus.b <- tclus(x, k = 3, alpha = 0.1, restr.fact = 50,
               restr = "eigen", equal.weights = FALSE)
clus.c <- tclus(x, k = 3, alpha = 0.1, restr.fact = 1,
               restr = "deter", equal.weights = TRUE)
clus.d <- tclus(x, k = 3, alpha = 0.1, restr.fact = 50,
               restr = "deter", equal.weights = FALSE)

pa <- par(mfrow = c(2, 2))
plot(clus.a, main = "(a)")
plot(clus.b, main = "(b)")
plot(clus.c, main = "(c)")
plot(clus.d, main = "(d)")
par(pa)

##--- EXAMPLE 4 -----

data (swissbank)
## Two clusters and 8\
(clus <- tclus(swissbank, k = 2, alpha = 0.08, restr.fact = 50))

## Pairs plot of the clustering solution
pairs(swissbank, col = clus$cluster + 1)
## Two coordinates
plot(swissbank[, 4], swissbank[, 6], col = clus$cluster + 1,
     xlab = "Distance of the inner frame to lower border",
     ylab = "Length of the diagonal")
plot(clus)

## Three clusters and 0\
clus<- tclus(swissbank, k = 3, alpha = 0.0, restr.fact = 110)

## Pairs plot of the clustering solution
pairs(swissbank, col = clus$cluster + 1)

```



```

## Two coordinates
plot(swissbank[, 4], swissbank[, 6], col = clus$cluster + 1,
      xlab = "Distance of the inner frame to lower border",
      ylab = "Length of the diagonal")

plot(clus)

##--- EXAMPLE 5 -----
data(M5data)
x <- M5data[, 1:2]

## Classification trimmed likelihood approach
clus.a <- tclust(x, k = 3, alpha = 0.1, restr.fact = 50,
                opt="HARD", restr = "eigen", equal.weights = FALSE)
## Mixture trimmed likelihood approach
clus.b <- tclust(x, k = 3, alpha = 0.1, restr.fact = 50,
                opt="MIXT", restr = "eigen", equal.weights = FALSE)

## Hard 0-1 cluster assignment (all 0 if trimmed unit)
head(clus.a$posterior)

## Posterior probabilities cluster assignment for the
## mixture approach (all 0 if trimmed unit)
head(clus.b$posterior)

## End(Not run)

```

tkmeans

TKMEANS method for robust K-means clustering

Description

This function searches for k (or less) spherical clusters in a data matrix x , whereas the ceiling(αn) most outlying observations are trimmed.

Usage

```

tkmeans(
  x,
  k,
  alpha = 0.05,
  nstart = 500,
  niter1 = 3,
  niter2 = 20,
  nkeep = 5,
  iter.max,

```

```

points = NULL,
center = FALSE,
scale = FALSE,
store_x = TRUE,
parallel = FALSE,
n.cores = -1,
zero_tol = 1e-16,
drop.empty.clust = TRUE,
trace = 0
)

```

Arguments

<code>x</code>	A matrix or data.frame of dimension $n \times p$, containing the observations (row-wise).
<code>k</code>	The number of clusters initially searched for.
<code>alpha</code>	The proportion of observations to be trimmed.
<code>nstart</code>	The number of random initializations to be performed.
<code>niter1</code>	The number of concentration steps to be performed for the <code>nstart</code> initializations.
<code>niter2</code>	The maximum number of concentration steps to be performed for the <code>nkeep</code> solutions kept for further iteration. The concentration steps are stopped, whenever two consecutive steps lead to the same data partition.
<code>nkeep</code>	The number of iterated initializations (after <code>niter1</code> concentration steps) with the best values in the target function that are kept for further iterations
<code>iter.max</code>	(deprecated, use the combination <code>nkeep</code> , <code>niter1</code> and <code>niter2</code>) The maximum number of concentration steps to be performed. The concentration steps are stopped, whenever two consecutive steps lead to the same data partition.
<code>points</code>	Optional initial mean vectors, NULL or a matrix with <code>k</code> vectors used as means to initialize the algorithm. If initial mean vectors are specified, <code>nstart</code> should be 1 (otherwise the same initial means are used for all runs).
<code>center</code>	Optional centering of the data: a function or a vector of length <code>p</code> which can optionally be specified for centering <code>x</code> before calculation
<code>scale</code>	Optional scaling of the data: a function or a vector of length <code>p</code> which can optionally be specified for scaling <code>x</code> before calculation
<code>store_x</code>	A logical value, specifying whether the data matrix <code>x</code> shall be included in the result object. By default this value is set to TRUE, because some of the plotting functions depend on this information. However, when big data matrices are handled, the result object's size can be decreased noticeably when setting this parameter to FALSE.
<code>parallel</code>	A logical value, specifying whether the <code>nstart</code> initializations should be done in parallel.
<code>n.cores</code>	The number of cores to use when paralellizing, only taken into account if <code>parallel=TRUE</code> .
<code>zero_tol</code>	The zero tolerance used. By default set to 1e-16.

<code>drop.empty.clust</code>	Logical value specifying, whether empty clusters shall be omitted in the resulting object. (The result structure does not contain center estimates of empty clusters anymore. Cluster names are reassigned such that the first l clusters ($l \leq k$) always have at least one observation.
<code>trace</code>	Defines the tracing level, which is set to 0 by default. Tracing level 1 gives additional information on the stage of the iterative process.

Value

The function returns the following values:

- `cluster` - A numerical vector of size n containing the cluster assignment for each observation. Cluster names are integer numbers from 1 to k , 0 indicates trimmed observations. Note that it could be empty clusters with no observations when `equal.weights=FALSE`.
- `obj` - The value of the objective function of the best (returned) solution.
- `size` - An integer vector of size k , returning the number of observations contained by each cluster.
- `centers` - A matrix of size $p \times k$ containing the centers (column-wise) of each cluster.
- `code` - A numerical value indicating if the concentration steps have converged for the returned solution (2).
- `cluster.ini` - A matrix with `nstart` rows and number of columns equal to the number of observations and where each row shows the final clustering assignments (0 for trimmed observations) obtained after the `niter1` iteration of the `nstart` random initializations.
- `obj.ini` - A numerical vector of length `nstart` containing the values of the target function obtained after the `niter1` iteration of the `nstart` random initializations.
- `x` - The input data set.
- `k` - The input number of clusters.
- `alpha` - The input trimming level.

Author(s)

Valentin Todorov, Luis Angel Garcia Escudero, Agustin Mayo Iscar.

References

Cuesta-Albertos, J. A.; Gordaliza, A. and Matrán, C. (1997), "Trimmed k-means: an attempt to robustify quantizers". *Annals of Statistics*, Vol. 25 (2), 553-576.

Examples

```
##--- EXAMPLE 1 -----
sig <- diag(2)
cen <- rep(1,2)
x <- rbind(MASS::mvrnorm(360, cen * 0, sig),
           MASS::mvrnorm(540, cen * 5, sig),
```

```
      MASS::mvrnorm(100, cen * 2.5, sig))

## Two groups and 10% trimming level
(clus <- tkmeans(x, k = 2, alpha = 0.1))

plot(clus)
plot(clus, labels = "observation")
plot(clus, labels = "cluster")

#--- EXAMPLE 2 -----
data(geyser2)
(clus <- tkmeans(geyser2, k = 3, alpha = 0.03))
plot(clus)
```

wholesale

Wholesale customers dataset

Description

The data set refers to clients of a wholesale distributor. It includes the annual spending in monetary units on diverse product categories.

Usage

```
data(wholesale)
```

Format

A data frame containing 440 observations in 8 variables (6 numerical and two categorical). The variables are as follows:

- Region Customers' Region - Lisbon (coded as 1), Porto (coded as 2) or Other (coded as 3)
- Fresh Annual spending on fresh products
- Milk Annual spending on milk products
- Grocery Annual spending on grocery products
- Frozen Annual spending on frozen products
- Detergents Annual spending on detergents and paper products
- Delicatessen Annual spending on and delicatessen products
- Channel Customers' Channel - Horeca (Hotel/Restaurant/Café) or Retail channel. Horeca is coded as 1 and Retail channel is coded as 2

Source

Abreu, N. (2011). Análise do perfil do cliente Recheio e desenvolvimento de um sistema promocional. Mestrado em Marketing, ISCTE-IUL, Lisbon. url=<https://api.semanticscholar.org/CorpusID:124027622>

Examples

```
#--- EXAMPLE 1 -----  
data (wholesale)  
x <- wholesale[, -c(1, ncol(wholesale))]  
clus <- tclust(x, k=3, alpha=0.1, nstart=200, niter1=3, niter2=17,  
  nkeep=10, opt="HARD", equal.weights=FALSE, restr.fact=50, trace=TRUE)  
plot (x, col=clus$cluster+1)  
plot(clus)
```

Index

* datasets

- geyser2, 5
- LG5data, 6
- M5data, 7
- pine, 8
- swissbank, 19
- wholesale, 28

ctlcurves, 2

DiscrFact, 4, 9

geyser2, 5

LG5data, 6

lines, 9

M5data, 7

pine, 8

plot.ctlcurves, 8

plot.DiscrFact, 5, 10

plot.rlg, 11

plot.tclust, 12

plot.tkmeans (plot.tclust), 12

print.ctlcurves (ctlcurves), 2

print.DiscrFact (DiscrFact), 4

print.tclust (tclust), 20

print.tkmeans (tkmeans), 25

rlg, 14

simula.rlg, 16

simula.tclust, 17

summary.DiscrFact, 18

swissbank, 19

tclust, 2, 3, 20

tkmeans, 25

wholesale, 28