## Package 'PUGMM'

October 23, 2025

```
Version 0.1.2
Title Parsimonious Ultrametric Gaussian Mixture Models
Description Parsimonious Ultrametric Gaussian Mixture Models via grouped coordinate as-
      cent (equivalent to EM) algorithm characterized by the inspection of hierarchical relation-
      ships among variables via parsimonious extended ultrametric covariance structures. The method-
      ologies are described in Cavicchia, Vichi, Zaccaria (2024) <doi:10.1007/s11222-024-10405-
      9>, (2022) <doi:10.1007/s11634-021-00488-x> and (2020) <doi:10.1007/s11634-020-00400-z>.
Depends R (>= 4.0)
Imports ClusterR, doParallel, foreach, igraph, ManlyMix, MASS, Matrix,
      mclust, mcompanion, ppclust, Rcpp
License MIT + file LICENSE
URL https://github.com/giorgiazaccaria/PUGMM
BugReports https://github.com/giorgiazaccaria/PUGMM/issues
LazyData true
Encoding UTF-8
RoxygenNote 7.3.3
Maintainer Giorgia Zaccaria < giorgia.zaccaria@unimib.it>
NeedsCompilation yes
LinkingTo Rcpp
Author Giorgia Zaccaria [aut, cre] (ORCID:
       <https://orcid.org/0000-0001-9119-9104>),
      Carlo Cavicchia [aut] (ORCID: <a href="https://orcid.org/0000-0003-1816-3521">https://orcid.org/0000-0003-1816-3521</a>),
      Lorenzo Balzotti [aut] (ORCID: <a href="https://orcid.org/0000-0001-6191-9801">https://orcid.org/0000-0001-6191-9801</a>),
      Alexa A. Sochaniwsky [aut] (ORCID:
       <a href="https://orcid.org/0009-0005-8043-5091">https://orcid.org/0009-0005-8043-5091</a>),
      Paul D. McNicholas [aut] (ORCID:
       <https://orcid.org/0000-0002-2482-523X>)
Repository CRAN
Date/Publication 2025-10-23 15:20:02 UTC
```

2 Harbour\_metals

## **Contents**

Harbour_metals	. 2
penguins	. 3
plot.pugmm	
pugmm	. 4
pugmm_available_models	. 8
puMmm	
rand.member	
UCM	. 13
	14

Harbour\_metals

Harbour\_metals

#### **Description**

Index

The harbour metals data set contains several metal concentration measurements on 60 seaweed samples. Each sample is either of the Padina or Sargassum species and the samples were collected across five cites in Port Jackson (Australia).

## Usage

```
data(Harbour_metals)
```

#### **Format**

A data frame with 60 observations and 7 variables, which are described as follows.

Location Sample sites (Balls Head, Berrys Bay, Felix Bay, Hermit Bay, Neutral Bay)

Species Seaweed species (Padina, Sargassum)

Rep Unique labels for replicate samples

**Cd** Cadmium ( $\mu$ g/g)

**Cr** Chromium (µg/g)

**Cu** Copper  $(\mu g/g)$ 

**Mn** Manganese ( $\mu$ g/g)

Ni Nickel ( $\mu$ g/g)

**Pb** Lead  $(\mu g/g)$ 

**Zn** Zinc  $(\mu g/g)$ 

## Source

Dataset downloaded from Environmental Computing https://environmentalcomputing.net/datasets/Harbour\_metals.csv.

penguins 3

#### References

Roberts, D.A., Johnston, E.L., Poore, A.G. (2008). Biomonitors and the assessment of ecological impacts: distribution of herbivorous epifauna in contaminated macroalgal beds. *Environmental Pollution*, 156(2), 489-503.

#### **Examples**

data(Harbour\_metals)

penguins

Penguins

## **Description**

The data set contains five measurements made on 342 penguins which are classified into three species.

#### Usage

data(penguins)

#### **Format**

A data frame with 342 observations and 5 variables, which are described as follows.

```
species Penguin species (Chinstrap, Adélie, or Gentoo)
culmen_length_mm Culmen length (mm)
culmen_depth_mm Culmen depth (mm)
flipper_length_mm Flipper length (mm)
body_mass_g Body mass (g)
```

#### **Details**

Data were collected and made available by Dr. Kristen Gorman and the Palmer Station, Antarctica LTER, a member of the Long Term Ecological Research Network. The categorical variables 'island' and 'sex' have been removed from the original dataset, as well as the incomplete observations on the five variables reported herein.

## Source

Dataset downloaded from Kaggle https://www.kaggle.com/code/parulpandey/penguin-dataset-the-new-iris.

plot.pugmm

#### References

Gorman, K.B., Williams T.D., Fraser W.R. (2014). Ecological sexual dimorphism and environmental variability within a community of Antarctic penguins (genus Pygoscelis). *PLoS ONE*, 9(3), e90081.

## **Examples**

```
data(penguins)
```

plot.pugmm

Plotting method for pugmm object

## **Description**

Plots for Parsimonious Ultrametric Gaussian Mixture Models results, such as BIC and path diagrams.

## Usage

```
## S3 method for class 'pugmm'
plot(x, what = NULL, nrow = NULL, ncol = NULL, cluster_names = NULL, ...)
```

## **Arguments**

Χ	Output from pugmm.
what	A string specifying the type of graph requested. Available choices are:
	"BIC" Plot of BIC values for the fitted models. For each $G$ , the best BIC among the ones corresponding to different $m$ is displayed.  "Path Diagram" Path diagram representation of the extended ultrametric co-
	variance matrix per component for the best model.
nrow	Number of rows in the graphical window. A new graphical window is opened every 6 plots, i.e., components of pugmm.
ncol	Number of columns in the graphical window. A new graphical window is opened every 6 plots, i.e., components of pugmm.
cluster_names	String of dimension $G$ with the clusters/components' name.

#### Value

No return value since this is a plot method.

Other graphics parameters.

## See Also

pugmm()

pugmm 5

## **Examples**

```
data(penguins)
x <- scale(penguins[, 2:5])
pugmm.penguins <- pugmm(x, 3, 1)
plot.pugmm(pugmm.penguins, what = c("BIC", "Path Diagram"))</pre>
```

pugmm

Parsimonious Ultrametric Gaussian Mixture Models

## **Description**

Model-based clustering via Parsimonious Ultrametric Gaussian Mixture Models. Hierarchical relationships among variables within and between clusters are inspected. The grouped coordinate ascent algorithm is used for the parameter estimation. The optimal model is selected according to BIC.

#### Usage

```
pugmm(
   X,
   G = NULL,
   m = NULL,
   normalization = NULL,
   model = NULL,
   maxiter = 500,
   tol = 1e-06,
   stop = "aitken",
   rndstart = 1,
   initG = "kmeans",
   initm = "ucms",
   gaussian = "mclust",
   parallel = FALSE
)
```

#### **Arguments**

X  $(n \times p)$  numeric matrix or data frame, where n and p represent the number of units and variables, respectively. Categorical variables are not allowed.

G Integer (vector) specifying the number of mixture components (default: G = 1:5).

Integer (vector) specifying the number of variable groups (default: m = 1:5).

normalization

Character string specifying the data transformation. If NULL, no transformation is applied to the data matrix (default). Other options are: "standard" for the standardization; "center" for centering the data; "range" for the MinMax transformation; "SVD" for the Singular Value Decomposition transformation.

6 pugmm

model	Vector of character strings indicating the model names to be fitted. If NULL, all the possible models are fitted (default). See the possible models using available_models().
maxiter	Integer value specifying the maximum number of iterations of the algorithm (default: maxiter = 500).
tol	Numeric value specifying the tolerance for the convergence criterion (default: $tol = 1e-6$ ).
stop	Character string specifying the convergence criteria. If "aitken", the Aitken acceleration-based stopping rule is used (default); if "relative", the relative log-likelihood in two sequential iterations is evaluated.
rndstart	Integer value specifying the number of random starts (default: rndstart = 1).
initG	Character string specifying the method for the initialization of the unit-component membership. If "kmeans", k-means via RcppArmadillo is used (default). Other options are: "random" for random assignment; "kmeansf" for fuzzy c-means (via the function fcm of the package ppclust).
initm	Character string specifying the method for the initialization of the variable-group membership. If "ucms", the multivariate model to be used for obtaining the variable-group membership estimated is the same model.name used for estimating the Parsimonious Ultrametric Gaussian Mixture Model (default); if "random", a random assignment is performed.
gaussian	Character string specifying the way to compute the log-likelihood. If "mclust", dmvnorm of mclust is used (default); if "canonical", the log-likelihood computation is based upon the canonical representation of an extended ultrametric covariance matrix.
parallel	A logical value, specifying whether the models should be run in parallel.

#### **Details**

The grouped coordinate ascent algorithm used for the estimation of PUGMMs parameters was demonstrated to be equivalent to an Expectation-Maximization algorithm in the GMM framework (Hathaway, 1986).

## Value

An object of class pugmm containing the results of the optimal - according to BIC - Parsimonious Ultrametric Gaussian Mixture Model estimation.

call Matched call.

X Input data matrix.

G Number of components of the best model.

m Number of variable groups of the best model.

label Integer vector of dimension n, taking values in  $\{1,\ldots,G\}$ . It identifies the unit classification according to the maximum a posteriori of the best model.

pp Numeric vector of dimension G containing the prior probabilities for the best model.

 $mu(G \times p)$  numeric matrix containing the component mean vectors (by row) for the best model.

pugmm 7

sigma List of dimension G containing the  $(p \times p)$  numeric component extended ultrametric covariance matrices for the best model.

V List of dimension G containing the  $(p \times m)$  binary variable-group membership matrices for the best model.

Sv List of dimension G containing the  $(m \times m)$  numeric diagonal matrices of the group variances for the best model.

Sw List of dimension G containing the  $(m \times m)$  numeric diagonal matrices of the within-group covariances for the best model.

Sb List of dimension G containing the  $(m \times m)$  numeric hallow matrices of the between-group covariances for the best model.

post  $(n \times G)$  numeric matrix containing the posterior probabilities for the best model.

pm Number of parameters of the best model.

pm. cov Number of covariance parameters of the best model.

pm. free Number of free parameters of the best model (pm-(constraints on V + count.constr.SwSb + count.constr.SvSw)).

count.constr.SwSb Number of times the constraint between Sw and Sb has been turned on for the best model.

count.constr.SvSw Number of times the constraint between Sv and Sw has been turned on for the best model.

BIC BIC values for all the fitted models. If BIC is NA, the model has not been computed since its structure is equal to another model, while if BIC is -Inf the solution has a number of clusters < G.

bic BIC value of the best model.

loglik Log-likelihood of the best model.

loop Random start corresponding to the selected solution of the best model.

iter Number of iterations needed to estimate the best model.

model.name Character string denoting the PUGMM model name of the best model among the ones fitted.

messages Messages.

#### References

Cavicchia, C., Vichi, M., Zaccaria, G. (2024) Parsimonious ultrametric Gaussian mixture models. *Statistics and Computing*, 34, 108.

Cavicchia, C., Vichi, M., Zaccaria, G. (2022) Gaussian mixture model with an extended ultrametric covariance structure. *Advances in Data Analysis and Classification*, 16(2), 399-427.

Hathaway, R. (1986) Another interpretation of the EM algorithm for mixture distributions. *Statistics and Probability Letters*, 4(2), 53-56.

#### See Also

pugmm\_available\_models(), plot.pugmm()

#### **Examples**

```
data(penguins)
x <- scale(penguins[, 2:5])
pugmm.penguins <- pugmm(x, 3, 1)
table(penguins$species, pugmm.penguins$label)
pugmm.penguins <- pugmm(x)
pugmm.penguins$G
pugmm.penguins$m
pugmm.penguins$m</pre>
```

pugmm\_available\_models

**PUGMM Model Names** 

## **Description**

Description of the model names used in the *PUGMM* package.

## Usage

```
pugmm_available_models()
```

#### **Details**

The PUGMM model names in the *PUGMM* package are characterized by four letters:

- First letter: it refers to the variable-group membership matrix V, which can be equal (E) or free to vary (F) across components.
- Second, third, fourth letters: they refer to the matrices of the group variances  $\Sigma_V$ , the withingroup covariances  $\Sigma_W$  and the between-group covariances  $\Sigma_B$ , respectively, by indicating if they are unique (U, i.e., equal within and across components), isotropic (I, i.e., equal within components), equal (E, i.e., equal across components) or free to vary across components (F).

#### Value

Available models in PUGMM, i.e., the thirteen extended ultrametric covariance structures of PUGMM.

#### References

Cavicchia, C., Vichi, M., Zaccaria, G. (2024) Parsimonious ultrametric Gaussian mixture models. *Statistics and Computing*, 34, 108.

#### See Also

pugmm()

puMmm 9

#### **Examples**

```
pugmm_available_models()
```

puMmm

Parsimonious Ultrametric Manly Mixture Models

## **Description**

Model-based clustering via Parsimonious Ultrametric Manly Mixture Models. Hierarchical relationships among variables within and between clusters are inspected. The grouped coordinate ascent algorithm is used for the parameter estimation. The optimal model is selected according to BIC.

## Usage

```
puMmm(
 Χ,
 G = NULL,
 m = NULL,
 lambda = NULL,
 normalization = NULL,
 model = NULL,
 modelselect = "BIC",
 maxiter = 500,
  tol = 1e-06,
  stop = "aitken",
  rndstart = 1,
  initG = "ManlyMix",
  initm = "ucms",
  seed = 123,
  parallel = FALSE
)
```

## **Arguments**

X	$(n \times p)$ numeric matrix or data frame, where $n$ and $p$ represent the number of units and variables, respectively. Categorical variables are not allowed.
G	Integer (vector) specifying the number of mixture components (default: $G = 1:5$ ).
m	Integer (vector) specifying the number of variable groups (default: m = 1:5).
lambda	$(G \times p)$ numeric matrix containing the initial transformation parameters for a single specified G.
normalization	Character string specifying the data transformation. If NULL, no transformation is applied to the data matrix (default). Other options are: "standard" for the standardization; "center" for centering the data; "range" for the MinMax trans-

formation; "SVD" for the Singular Value Decomposition transformation.

10 puMmm

model Vector of character strings indicating the model names to be fitted. If NULL, all

the possible models are fitted (default). See the possible models using pugmm\_available\_models().

modelselect Character string indicating the model selection method to be used. If "BIC",

the best model is selected according to the BIC (default); if "two-step", the best

model is selected according to the two-step model selection method.

maxiter Integer value specifying the maximum number of iterations of the EM algorithm

(default: maxiter = 500).

Numeric value specifying the tolerance for the convergence criteria used in the

EM algorithm (default: tol = 1e-6).

stop Character string specifying the convergence criteria. If "aitken", the Aitken

acceleration-based stopping rule is used (default); if "relative", the relative log-

likelihood in two sequential iterations is evaluated.

rndstart Integer value specifying the number of random starts (default: rndstart = 1).

initG Character string specifying the method for the initialization of the unit-component

membership. If "ManlyMix", the Manly.model() function via the ManlyMix package is used (default). Other options are: "kmeans" for k-means (via Rcp-pArmadillo); "random" for random assignment; "kmeansf" for fuzzy c-means

(via the function fcm of the package ppclust).

initm Character string specifying the method for the initialization of the variable-

group membership. If "ucms", the multivariate model to be used for obtaining the variable-group membership estimated is the same model.name used for estimating the Parsimonious Ultrametric Manly Mixture Model (default); if "ran-

dom", a random assignment is performed.

seed Numeric value specifying the seed (default: seed = 123).

parallel A logical value, specifying whether the models should be run in parallel.

#### **Details**

The grouped coordinate ascent algorithm used for the estimation of PUMMMs parameters was demonstrated to be equivalent to an Expectation-Maximization (EM) algorithm in the GMM framework (Hathaway, 1986).

#### Value

An object of class puMmm containing the results of the optimal -according to the model selection criteria- Parsimonious Ultrametric Manly Mixture Model estimation.

call Matched call.

X Input data matrix.

G Number of components of the best model.

m Number of variable groups of the best model.

label Integer vector of dimension n, taking values in  $\{1, \ldots, G\}$ . It identifies the unit classification according to the maximum a posteriori of the best model.

pp Numeric vector of dimension G containing the prior probabilities for the best model.

puMmm 11

lambda  $(G \times p)$  numeric matrix containing the component transformation vectors (by row) for the best model.

 $mu(G \times p)$  numeric matrix containing the component mean vectors (by row) for the best model.

sigma List of dimension G containing the  $(p \times p)$  numeric component extended ultrametric covariance matrices for the best model.

V List of dimension G containing the  $(p \times m)$  binary variable-group membership matrices for the best model.

Sv List of dimension G containing the  $(m \times m)$  numeric diagonal matrices of the group variances for the best model.

Sw List of dimension G containing the  $(m \times m)$  numeric diagonal matrices of the within-group covariances for the best model.

Sb List of dimension G containing the  $(m \times m)$  numeric hallow matrices of the between-group covariances for the best model.

post  $(n \times G)$  numeric matrix containing the posterior probabilities for the best model.

pm Number of parameters of the best model.

pm. cov Number of covariance parameters of the best model.

pm. free Number of free parameters of the best model (pm - (constraints on V + count.constr.SwSb + count.constr.SvSw)).

count.constr.SwSb Number of times the constraint between Sw and Sb has been turned on for the best model.

count.constr.SvSw Number of times the constraint between Sv and Sw has been turned on for the best model.

BIC BIC values for all the fitted models. If BIC is NA, the model has not been computed since its structure is equal to another model, while if BIC is -Inf the solution has a number of clusters < G.

bic BIC value of the best model.

loglik Log-likelihood of the best model.

loop Random start corresponding to the selected solution of the best model.

iter Number of iterations needed to estimate the best model.

model.name Character string denoting the PUGMM model name of the best model among the ones fitted.

messages Messages.

#### References

Cavicchia, C., Vichi, M., Zaccaria, G. (2024) Parsimonious ultrametric Gaussian mixture models. *Statistics and Computing*, 34, 108.

Cavicchia, C., Vichi, M., Zaccaria, G. (2022) Gaussian mixture model with an extended ultrametric covariance structure. *Advances in Data Analysis and Classification*, 16(2), 399-427.

Hathaway, R. (1986) Another interpretation of the EM algorithm for mixture distributions. *Statistics and Probability Letters*, 4(2), 53-56.

12 rand.member

#### See Also

```
pugmm(), pugmm_available_models(), plot.pugmm()
```

## **Examples**

```
data(Harbour_metals)
x <- scale(Harbour_metals[,4:10])
## Not run:
results <- puMmm(x, G = 1:4, m = 1:7, model = NULL, modelselect = "two-step")
results$G
results$m
results$model.name
table(Harbour_metals$Species, results$label)
plot.pugmm(results, what = c("BIC", "Path Diagram"))
## End(Not run)</pre>
```

rand.member

Random partition of objects into classes

## Description

Performs a random partition of objects into classes.

## Usage

```
rand.member(n.obs, G)
```

## Arguments

n. obsNumber of objectsGNumber of classes

## **Details**

No empty classes can occur.

#### Value

A binary and row-stochastic matrix with n.obs rows and G columns.

## **Examples**

```
rand.member(10, 3)
```

UCM 13

UCM	Ultrametric Correlation Matrix	

#### **Description**

Fit an ultrametric correlation matrix on a nonnegative correlation one.

## Usage

```
UCM(R, m, rndstart, maxiter = 100, eps = 1e-06)
```

## **Arguments**

R  $(p \times p)$  nonnegative correlation matrix.

m Integer specifying the number of variable groups.

rndstart Integer value specifying the number of random starts.

maxiter Integer value specifying the maximum number of iterations of the EM algorithm

(default: maxiter = 100).

eps Numeric value specifying the tolerance for the convergence criterion used in the

coordinate descent algorithm (default: eps = 1e-6).

#### Value

A list with the following elements:

call Matched call.

V Optimal binary and row-stochastic  $(p \times m)$  variable-group membership matrix.

Rt Optimal  $(p \times p)$  ultrametric correlation matrix.

Rw Optimal  $(m \times m)$  within-concept consistency (diagonal) matrix.

Rb Optimal  $(m \times m)$  between-concept correlation matrix.

of Objective function corresponding to the optimal solution.

loop Random start corresponding to the optimal solution.

iter Number of iterations needed to obtain the optimal solution.

## References

Cavicchia, C., Vichi, M., Zaccaria, G. (2020) The ultrametric correlation matrix for modelling hierarchical latent concepts. *Advances in Data Analysis and Classification*, 14(4), 837-853.

## **Examples**

```
data(penguins)
R <- cor(penguins[, 2:5])
UCM(R, 4, 1)</pre>
```

# **Index**

```
* Datasets
    penguins, 3
* datasets
    Harbour_metals, 2

Harbour_metals, 2

penguins, 3
plot.pugmm, 4
plot.pugmm(), 7, 12
pugmm, 5
pugmm(), 4, 8, 12
pugmm_available_models, 8
pugmm_available_models(), 7, 12
puMmm, 9

rand.member, 12

UCM, 13
```