

Package ‘stepmixr’

January 9, 2024

Note -*- Encoding: utf-8 -*-

Type Package

Title Interface to 'Python' Package 'StepMix'

Version 0.1.2

Date 2024-01-03

Author Éric Lacourse [aut],
Roxane de la Sablonnière [aut],
Charles-Édouard Giguère [aut, cre],
Sacha Morin [aut],
Robin Legault [aut],
Félix Laliberté [aut],
Zsusza Bakk [ctb]

Maintainer Charles-Édouard Giguère <ce.giguere@gmail.com>

Depends R (>= 4.0.0)

Imports reticulate (>= 1.8)

Description This is an interface for the 'Python' package 'StepMix'. It is a 'Python' package following the scikit-learn API for model-based clustering and generalized mixture modeling (latent class/profile analysis) of continuous and categorical data. 'StepMix' handles missing values through Full Information Maximum Likelihood (FIML) and provides multiple stepwise Expectation-Maximization (EM) estimation methods based on pseudolikelihood theory. Additional features include support for covariates and distal outcomes, various simulation utilities, and non-parametric bootstrapping, which allows inference in semi-supervised and unsupervised settings.

License GPL-2

Encoding UTF-8

LazyLoad TRUE

URL <https://github.com/Labo-Lacourse/StepMixr>

NeedsCompilation no

Repository CRAN

Date/Publication 2024-01-09 22:20:02 UTC

R topics documented:

bootstrap	2
bootstrap_stats	3
Datasets	4
fit	5
install.stepmix	7
mixed_descriptor	8
predict.stepmix.stepmix.StepMix	9
savefit	10
stepmix	11
Index	15

bootstrap	<i>Non-parametric bootstrap of StepMix estimator.</i>
-----------	---

Description

Non-parametric bootstrap of StepMix estimator. Fit the estimator on X,Y then fit n_repetitions on resampled datasets. Repetition parameters are aligned with the class order of the main estimator.

Usage

```
## S3 method for class 'stepmix.stepmix.StepMix'
bootstrap(x, X = NULL, y = NULL, n_repetitions = 10, ...)
bootstrap(x, ...)
```

Arguments

x	An object created with the fit function
X	The X matrix or data.frame for the measurement part of the model
y	The Y matrix or data.frame for the structural part of the model
n_repetitions	The number of bootstrap sample
...	For future options. This option is actually unused.

Details

This methods returns a list with bootstrap samples (samples) and the log-likelihood (rep_stats).

Value

A list containing bootstrap samples of the parameters.

Author(s)

Éric Lacourse, Roxane de la Sablonnière, Charles-Édouard Giguère, Sacha Morin, Robin Legault, Félix Laliberté, Zsuzsa Bakk

References

Bolck, A., Croon, M., and Hagenaars, J. Estimating latent structure models with categorical variables: One-step versus three-step estimators. *Political analysis*, 12(1): 3-27, 2004.

Vermunt, J. K. Latent class modeling with covariates: Two improved three-step approaches. *Political analysis*, 18 (4):450-469, 2010.

Bakk, Z., Tekle, F. B., and Vermunt, J. K. Estimating the association between latent class membership and external variables using bias-adjusted three-step approaches. *Sociological Methodology*, 43(1):272-311, 2013.

Bakk, Z. and Kuha, J. Two-step estimation of models between latent classes and external variables. *Psychometrika*, 83(4):871-892, 2018

Examples

```
## Not run:
if (reticulate::py_module_available("stepmix")) {
  require(stepmixr)
  model1 <- stepmix(n_components = 3, n_steps = 2, measurement = "continuous", progress_bar = 0)
  X <- iris[c(1:10, 51:60, 101:110), 1:4]
  fit1 <- fit(model1, X)
  fit1_bs <- bootstrap(fit1, X, n_repetitions = 5, progress_bar = FALSE)
}

## End(Not run)
```

bootstrap_stats	<i>Non-parametric bootstrap of StepMix estimator.</i>
-----------------	---

Description

Non-parametric bootstrap of StepMix estimator. Obtain bootstrapped parameters and some statistics (mean and standard deviation). If a covariate model is used in the structural model, the output keys "cw_mean" and "cw_std" are omitted.

Usage

```
## S3 method for class 'stepmix.stepmix.StepMix'
bootstrap_stats(x, X = NULL, y = NULL, n_repetitions = 10, ...)
bootstrap_stats(x, ...)
```

Arguments

x	An object created with the fit function
X	The X matrix or data.frame for the measurement part of the model
y	The y matrix or data.frame for the structural part of the model
n_repetitions	The number of bootstrap sample
...	for future options. Currently not used

Details

This methods returns a list with bootstrap samples (`samples`) and the log-likelihood (`rep_stats`). Mean and standard deviation are added to the results.

Value

A list containing bootstrap samples of the parameters. The mean and standard of class weights (`cw_mean`, `cw_std`), measurement model parameters (`mm_mean`, `mm_std`), structural model parameters (`sm_mean`, `sm_std`) are also added. If a covariate model is used in the structural model, the output keys `cw_mean` and `cw_std` are omitted.

Author(s)

Éric Lacourse, Roxane de la Sablonnière, Charles-Édouard Giguère, Sacha Morin, Robin Legault, Félix Laliberté, Zsuzsa Bakk

References

- Bolck, A., Croon, M., and Hageaars, J. Estimating latent structure models with categorical variables: One-step versus three-step estimators. *Political analysis*, 12(1): 3-27, 2004.
- Vermunt, J. K. Latent class modeling with covariates: Two improved three-step approaches. *Political analysis*, 18 (4):450-469, 2010.
- Bakk, Z., Tekle, F. B., and Vermunt, J. K. Estimating the association between latent class membership and external variables using bias-adjusted three-step approaches. *Sociological Methodology*, 43(1):272-311, 2013.
- Bakk, Z. and Kuha, J. Two-step estimation of models between latent classes and external variables. *Psychometrika*, 83(4):871-892, 2018

Datasets

Series of function to simulate data.

Description

These functions generates data with multiple groups using different distributions and optionnaly adding a level of missing value.

Usage

```
random_nan(X, Y, nan_ratio, random_state=NULL)
bakk_measurements(n_classes, n_mm, sep_level)
data_bakk_response(n_samples, sep_level, n_classes = 3, n_mm = 6, random_state = NULL)
data_bakk_covariate(n_samples, sep_level, n_mm = 6, random_state = NULL)
data_bakk_complete(n_samples, sep_level, n_mm=6, random_state=NULL, nan_ratio=0.0)
data_generation_gaussian(n_samples, sep_level, n_mm=6, random_state=NULL)
data_gaussian_diag(n_samples, sep_level, n_mm = 6, random_state = NULL, nan_ratio = 0.0)
```

Arguments

<code>X</code>	The X matrix or data.frame for the measurement part of the model
<code>Y</code>	The Y matrix or data.frame for the structural part of the model
<code>nan_ratio</code>	The ratio of missing values. A value between 0 and 1.
<code>random_state</code>	An integer initializing the seed of the random generator.
<code>n_classes</code>	Number of latent classes required.
<code>n_mm</code>	Number of features in the measurement model.
<code>sep_level</code>	Separation level in the measurement data.
<code>n_samples</code>	Number of samples.

Details

These function returns simulated data used to test the package.

Value

list of data.frame simulated according to the function parameters.

Author(s)

Éric Lacourse, Roxane de la Sablonnière, Charles-Édouard Giguère, Sacha Morin, Robin Legault, Félix Laliberté, Zsuzsa Bakk

References

Bakk, Z. and Kuha, J. Two-step estimation of models between latent classes and external variables. *Psychometrika*, 83(4):871-892, 2018

fit

Fit a mixture using the stepmix python package.

Description

This function initializes the stepmix object in python and fit X and optionnally Y to the object.

Usage

```
fit(smx, X = NULL, Y = NULL, ...)  
## S3 method for class 'stepmix.stepmix.StepMix'  
print(x, x_names = NULL, y_names = NULL, ...)  
identify_coef(coef)
```

Arguments

smx	An object created with the stepmix function.
X	The X matrix or data.frame for the measurement part of the model
Y	The Y matrix or data.frame for the structural part of the model
x	An object fitted with the fit method
coef	Matrix of coefficients to be modified
x_names	Optional name of x variables
y_names	Optional name of y variables
...	unused but included to be inline with requirement of generic function

Details

This methods returns a pointer to a python object of type StepMix. It can be used within reticulate but not within R. To save this type of object, you need to use the savefit function. The print method, uses the same print methods used when verbose = TRUE, it takes the last X and Y arguments used with the fit method. identify_coef find a reference configuration of the coefficients.

Value

A pointer to a python object of type StepMix.

Author(s)

Éric Lacourse, Roxane de la Sablonnière, Charles-Édouard Giguère, Sacha Morin, Robin Legault, Félix Laliberté, Zsuzsa Bakk

References

- Bolck, A., Croon, M., and Hageaars, J. Estimating latent structure models with categorical variables: One-step versus three-step estimators. *Political analysis*, 12(1): 3-27, 2004.
- Vermunt, J. K. Latent class modeling with covariates: Two improved three-step approaches. *Political analysis*, 18 (4):450-469, 2010.
- Bakk, Z., Tekle, F. B., and Vermunt, J. K. Estimating the association between latent class membership and external variables using bias-adjusted three-step approaches. *Sociological Methodology*, 43(1):272-311, 2013.
- Bakk, Z. and Kuha, J. Two-step estimation of models between latent classes and external variables. *Psychometrika*, 83(4):871-892, 2018

Examples

```
## Not run:
if (reticulate::py_module_available("stepmix")) {
  model1 <- stepmix(n_components = 3, n_steps = 2, measurement = "continuous", progress_bar = 0)
  X <- iris[c(1:10, 51:60, 101:110), 1:4]
  fit1 <- fit(model1, X)
}

## End(Not run)
```

install.stepmix *Install stepmix python package into python via reticulate.*

Description

Install the stepmix python package in the python instance used by reticulate.

Usage

```
install.stepmix(envname, method, conda, pip, ...)  
check_pystepmix_version()
```

Arguments

envname	Name of the python environment. "r-reticulate" by default.
method	installation method. See doc in reticulate
conda	Path to a conda install. See doc in reticulate
pip	Logical value to choose pip as the install method
...	Not used in function

Details

This methods installs stepmix in the python instance or environment used by reticulate. It uses `reticulate::py_install`.

Value

It doesn't return anything.

Author(s)

Éric Lacourse, Roxane de la Sablonnière, Charles-Édouard Giguère, Sacha Morin, Robin Legault, Félix Laliberté, Zsuzsa Bakk

References

- Bolck, A., Croon, M., and Hagenars, J. Estimating latent structure models with categorical variables: One-step versus three-step estimators. *Political analysis*, 12(1): 3-27, 2004.
- Vermunt, J. K. Latent class modeling with covariates: Two improved three-step approaches. *Political analysis*, 18 (4):450-469, 2010.
- Bakk, Z., Tekle, F. B., and Vermunt, J. K. Estimating the association between latent class membership and external variables using bias-adjusted three-step approaches. *Sociological Methodology*, 43(1):272-311, 2013.
- Bakk, Z. and Kuha, J. Two-step estimation of models between latent classes and external variables. *Psychometrika*, 83(4):871-892, 2018

mixed_descriptor	<i>Utility function for mixture using mixed description.</i>
------------------	--

Description

This function creates a data.frame ordered by continuous, binary and categorical columns. It also creates a list used if the model uses mixed column types.

Usage

```
mixed_descriptor(data, continuous = NULL, binary = NULL,  
                categorical = NULL, covariate = NULL)
```

Arguments

data	Data.frame with the mixed data
continuous	index or name of continuous column
binary	index or name of binary column
categorical	index or name of categorical column
covariate	index or name of covariate column

Details

This methods returns a list of a data.frame sorted by continuous, binary and categorical columns. It contains also a descriptor that can be used in the measurement section.

Value

A list containing data and a descriptor.

Author(s)

Éric Lacourse, Roxane de la Sablonnière, Charles-Édouard Giguère, Sacha Morin, Robin Legault, Félix Laliberté, Zsuzsa Bakk

References

- Bolck, A., Croon, M., and Hagenars, J. Estimating latent structure models with categorical variables: One-step versus three-step estimators. *Political analysis*, 12(1): 3-27, 2004.
- Vermunt, J. K. Latent class modeling with covariates: Two improved three-step approaches. *Political analysis*, 18 (4):450-469, 2010.
- Bakk, Z., Tekle, F. B., and Vermunt, J. K. Estimating the association between latent class membership and external variables using bias-adjusted three-step approaches. *Sociological Methodology*, 43(1):272-311, 2013.
- Bakk, Z. and Kuha, J. Two-step estimation of models between latent classes and external variables. *Psychometrika*, 83(4):871-892, 2018

Examples

```
md <- mixed_descriptor(iris, continuous = 1:4, categorical = 5)
```

```
predict.stepmix.stepmix.StepMix
```

Predict the membership (probabilities) using the fit of the stepmix python package.

Description

Predict the membership (probabilities) of a mixture using a stepmix object in python using X and optionally Y to the object.

Usage

```
## S3 method for class 'stepmix.stepmix.StepMix'
predict(object, X = NULL, Y = NULL, ...)
## S3 method for class 'stepmix.stepmix.StepMix'
predict_proba(object, X = NULL, Y = NULL, ...)
```

Arguments

object	An object created with the fit function.
X	The X matrix or data.frame for the measurement part of the model
Y	The Y matrix or data.frame for the structural part of the model
...	not used in this function

Value

A vector containing the membership (probabilities) of the mixture.

Author(s)

Éric Lacourse, Roxane de la Sablonnière, Charles-Édouard Giguère, Sacha Morin, Robin Legault, Zsusza Bakk

References

- Bolck, A., Croon, M., and Hagenaaars, J. Estimating latent structure models with categorical variables: One-step versus three-step estimators. *Political analysis*, 12(1): 3-27, 2004.
- Vermunt, J. K. Latent class modeling with covariates: Two improved three-step approaches. *Political analysis*, 18 (4):450-469, 2010.
- Bakk, Z., Tekle, F. B., and Vermunt, J. K. Estimating the association between latent class membership and external variables using bias-adjusted three-step approaches. *Sociological Methodology*, 43(1):272-311, 2013.
- Bakk, Z. and Kuha, J. Two-step estimation of models between latent classes and external variables. *Psychometrika*, 83(4):871-892, 2018

Examples

```
## Not run:
if (reticulate::py_module_available("stepmix")) {
  require(stepmixr)
  model1 <- stepmix(n_components = 3, n_steps = 2, measurement = "continuous", progress_bar = 0)
  X <- iris[c(1:10, 51:60, 101:110), 1:4]
  fit1 <- fit(model1, X)
  pr1 <- predict(fit1, X)
}

## End(Not run)
```

savefit

Save the fit of a mixture using the stepmix python package.

Description

This function saves the stepmix fitted object in python using the pickle package.

Usage

```
savefit(fitx, f)
loadfit(f)
```

Arguments

fitx An object created with the stepmix function.
f String indicating the name of the file

Details

This methods allows to save/load the stepmix object in a binary file using the pickle package.

Value

A pointer to a python object of type StepMix.

Author(s)

Éric Lacourse, Roxane de la Sablonnière, Charles-Édouard Giguère, Sacha Morin, Robin Legault, Félix Laliberté, Zsuzsa Bakk

References

- Bolck, A., Croon, M., and Hageaars, J. Estimating latent structure models with categorical variables: One-step versus three-step estimators. *Political analysis*, 12(1): 3-27, 2004.
- Vermunt, J. K. Latent class modeling with covariates: Two improved three-step approaches. *Political analysis*, 18 (4):450-469, 2010.
- Bakk, Z., Tekle, F. B., and Vermunt, J. K. Estimating the association between latent class membership and external variables using bias-adjusted three-step approaches. *Sociological Methodology*, 43(1):272-311, 2013.
- Bakk, Z. and Kuha, J. Two-step estimation of models between latent classes and external variables. *Psychometrika*, 83(4):871-892, 2018

Examples

```
## Not run:
if (reticulate::py_module_available("stepmix")) {
  model1 <- stepmix(n_components = 2, n_steps = 3, progress_bar = 0)
  X <- data.frame(x1 = c(0,1,1,1,1,0,0,0,0,0,1,1,0),
                 x2 = c(0,1,1,0,0,1,1,0,0,0,1,0,1))
  fit1 <- fit(model1, X)
  savefit(fit1, "fit1.pickle")

  ### clean the directory.
  file.remove("fit1.pickle")
}

## End(Not run)
```

stepmix

R interface to stepmix in StepMix python.

Description

This function creates a basic R list that will be used to initialize the stepmix object in python, in order to use the fit and predict function.

Usage

```
stepmix(n_components = 2, n_steps = 1,
        measurement = "bernoulli", structural = "gaussian_unit",
        assignment = "modal", correction = NULL,
        abs_tol = 1e-10, rel_tol = 0, max_iter = 1000,
        n_init = 1, init_params = "random", random_state = NULL,
        verbose = 0, progress_bar = 1, measurement_params = NULL,
        structural_params = NULL)
```

Arguments

n_components	The number of latent class. 2 by default.
n_steps	1, 2, or 3, 1 by default. Number of steps in the estimation. Must be one of : 1: run EM on both the measurement and structural models. 2: first run EM on the measurement model, then on the complete model, but keep the measurement parameters fixed for the second step. See Bakk, 2018. 3: first run EM on the measurement model, assign class probabilities, then fit the structural model via maximum likelihood. See the correction parameter for bias correction. See Bakk & Kuha (2018) for more details.
measurement	String describing the measurement model. See details for the different available model. The default model is "bernouilli"
structural	String describing the structural model. See details for the different available model. The default model is "bernouilli"
assignment	String indicating the type of class assignments for 3-step estimation, "modal" by default. Must be one of: soft: keep class responsibilities (posterior probabilities) as is. modal: assign 1 to the class with max probability, 0 otherwise (one-hot encoding).
correction	Bias correction for 3-step estimation. Must be one of : None: No correction. Run Naive 3-step. BCH: Apply the empirical BCH correction from Vermunt, 2004. ML: Apply the ML correction from Vermunt, 2010, Bakk et al., 2013.
abs_tol	The convergence threshold. EM iterations will stop when the lower bound average gain is below this threshold. The default value is 1e-3.
rel_tol	The convergence threshold. EM iterations will stop when the relative lower bound average gain is below this threshold.
max_iter	The number of EM iterations to perform.
n_init	The number of initializations to perform. The best results are kept.
init_params	"kmeans", or "random", default="random". The method used to initialize the weights, the means and the precisions. Must be one of: kmeans : responsibilities are initialized using kmeans. random : responsibilities are initialized randomly.
random_state	State instance or NULL, default=NULL. Controls the random seed given to the method chosen to initialize the parameters. Pass an int for reproducible output across multiple function calls.
verbose	Default=0. Enable verbose output. If 1, will print detailed report of the model and the performance metrics after fitting.
progress_bar	Display a tqdm progress bar during fitting
measurement_params	Default=NULL, Additional params passed to the measurement model class. Particularly useful to specify optimization parameters for stepmix.emission.covariate.Covariate. Ignored if the measurement descriptor is a nested object (see stepmix.emission.nested.Nested).

structural_params

Default=NULL, Additional params passed to the structural model class. Particularly useful to specify optimization parameters for `stepmix.emission.covariate.Covariate`. Ignored if the structural descriptor is a nested object (see `stepmix.emission.nested.Nested`).

Details

The options for both the measurement and structural part are describe here:

`bernoulli`: The observed data consists of `n_features` bernoulli (binary) random variables.

`bernoulli_nan`: the observed data consists of `n_features` bernoulli (binary) random variables. Supports missing values.

`binary`: alias for `bernoulli`.

`binary_nan`: alias for `bernoulli_nan`.

`categorical`: alias for `multinoulli`.

`categorical_nan`: alias for `multinoulli_nan`.

`continuous`: alias for `gaussian diag`.

`continuous_nan`: alias for `gaussian_diag_nan`. supports missing values.

`covariate`: covariate model where class probabilities are a multinomial logistic model of the features.

`gaussian`: alias for `gaussian_unit`.

`gaussian_nan`: alias for `gaussian_unit`. Supports missing values.

`gaussian_unit`: each gaussian component has unit variance. Only fit the mean.

`gaussian_unit_nan`: each gaussian component has unit variance. Only fit the mean. Supports missing values.

`gaussian_spherical`: each gaussian component has its own single variance.

`gaussian_spherical_nan`: each gaussian component has its own single variance. Supports missing values.

`gaussian_tied`: all gaussian components share the same general covariance matrix.

`gaussian_diag`: each gaussian component has its own diagonal covariance matrix.

`gaussian_diag_nan`: each gaussian component has its own diagonal covariance matrix. Supports missing values.

`gaussian_full`: each gaussian component has its own general covariance matrix.

`multinoulli`: the observed data consists of `n_features` multinoulli (categorical) random variables.

`multinoulli_nan`: the observed data consists of `n_features` multinoulli (categorical) random variables. Supports missing values.

Value

It returns a list of type `stepmixr` that contains the arguments of the object.

Author(s)

Éric Lacourse, Roxane de la Sablonnière, Charles-Édouard Giguère, Sacha Morin, Robin Legault, Félix Laliberté, Zsuzsa Bakk

References

Bolck, A., Croon, M., and Hagenaars, J. Estimating latent structure models with categorical variables: One-step versus three-step estimators. *Political analysis*, 12(1): 3-27, 2004.

Vermunt, J. K. Latent class modeling with covariates: Two improved three-step approaches. *Political analysis*, 18 (4):450-469, 2010.

Bakk, Z., Tekle, F. B., and Vermunt, J. K. Estimating the association between latent class membership and external variables using bias-adjusted three-step approaches. *Sociological Methodology*, 43(1):272-311, 2013.

Bakk, Z. and Kuha, J. Two-step estimation of models between latent classes and external variables. *Psychometrika*, 83(4):871-892, 2018

See Also

[fit](#)

Examples

```
model1 <- stepmix(n_components = 2, n_steps = 3)
```

Index

bakk_measurements (Datasets), 4
bootstrap, 2
bootstrap_stats, 3

check_pystepmix_version
 (install.stepmix), 7

data_bakk_complete (Datasets), 4
data_bakk_covariate (Datasets), 4
data_bakk_response (Datasets), 4
data_gaussian_diag (Datasets), 4
data_generation_gaussian (Datasets), 4
Datasets, 4

fit, 5, 14

identify_coef (fit), 5
install.stepmix, 7

loadfit (savefit), 10

mixed_descriptor, 8

predict.stepmix.stepmix.StepMix, 9
predict_proba
 (predict.stepmix.stepmix.StepMix),
 9
print.stepmix.stepmix.StepMix (fit), 5

random_nan (Datasets), 4

savefit, 10
stepmix, 11