

Package ‘fuzzySim’

July 7, 2017

Type Package

Title Fuzzy Similarity in Species Distributions

Version 1.8.0

Date 2017-07-07

Author Barbosa A.M.

Maintainer A. Marcia Barbosa <barbosa@uevora.pt>

Suggests maptools, parallel, PBSmapping, sp, tools

Description Functions to calculate fuzzy versions of species' occurrence patterns based on presence-absence data (including inverse distance interpolation, trend surface analysis and prevalence-independent favourability GLM), and pair-wise fuzzy similarity (based on fuzzy versions of commonly used similarity indices) among those occurrence patterns. Includes also functions for model comparison (overlap and fuzzy similarity, loss or gain), and for data preparation, such as obtaining unique abbreviations of species names, converting species lists (long format) to presence-absence tables (wide format), transposing part of a data frame, assessing the false discovery rate, or analysing and dealing with multicollinearity among variables. Includes also sample datasets for providing practical examples.

License GPL-3

R topics documented:

| | |
|----------------------------|----|
| fuzzySim-package | 1 |
| corSelect | 4 |
| distPres | 6 |
| Fav | 8 |
| FDR | 11 |
| fuzSim | 13 |
| fuzzyOverlay | 15 |
| fuzzyRangeChange | 18 |
| getPreds | 20 |
| integerCols | 22 |
| modelTrim | 23 |
| modOverlap | 25 |
| multConvert | 26 |

| | |
|-----------------------------|----|
| multGLM | 28 |
| multicol | 31 |
| multTSA | 33 |
| pairwiseRangemaps | 35 |
| percentTestData | 37 |
| rangemapSim | 38 |
| rotif.env | 39 |
| rotifers | 41 |
| simFromSetOps | 42 |
| simMat | 44 |
| spCodes | 46 |
| splist2presabs | 48 |
| stepByStep | 49 |
| timer | 51 |
| transpose | 52 |
| triMatInd | 53 |

| | |
|------------------|--|
| fuzzySim-package | <i>Fuzzy Similarity in Species Distributions</i> |
|------------------|--|

Description

Functions to calculate fuzzy versions of species' occurrence patterns based on presence-absence data (including inverse distance interpolation, trend surface analysis and prevalence-independent favourability GLM), and pair-wise fuzzy similarity (based on fuzzy versions of commonly used similarity indices) among those occurrence patterns. Includes also functions for data preparation, such as obtaining unique abbreviations of species names, converting species lists (long format) to presence-absence tables (wide format), transposing part of a data frame, assessing the false discovery rate, or analysing and dealing with multicollinearity among variables. Includes also sample datasets for providing practical examples. A step-by-step illustrated tutorial is available from the package homepage (<http://fuzzysim.r-forge.r-project.org>).

Details

| | |
|----------|------------|
| Package: | fuzzySim |
| Type: | Package |
| Version: | 1.8.0 |
| Date: | 2017-07-07 |
| License: | GPL-3 |

Author(s)

A. Marcia Barbosa

Maintainer: A. Marcia Barbosa <barbosa@uevora.pt>

References

Barbosa A.M. (2015) fuzzySim: applying fuzzy logic to binary similarity indices in ecology. *Methods in Ecology and Evolution*, 6: 853-858.

Examples

```
data(rotifers)

head(rotifers)

# add column with species name abbreviations:

rotifers$spcode <- spCodes(rotifers$species, sep.species = "_",
nchar.gen = 1, nchar.sp = 5, nchar.ssp = 0)

head(rotifers)

# convert species list (long format) to presence-absence table
# (wide format):

rotifers.presabs <- splist2presabs(rotifers, sites.col = "TDWG4",
sp.col = "spcode", keep.n = FALSE)

head(rotifers.presabs)

# get 3rd-degree spatial trend surface for each species' distribution:

data(rotif.env)

names(rotif.env)

rotifers.tsa <- multTSA(rotif.env, sp.cols = 18:47,
coord.cols = c("Longitude", "Latitude"), id.col = 1)

head(rotifers.tsa)

# get inverse squared distance to presence for each species:

rotifers.isqd <- distPres(rotif.env, sp.cols = 18:47,
coord.cols = c("Longitude", "Latitude"), id.col = 1, p = 2, inv = TRUE)

head(rotifers.isqd)

# get prevalence-independent environmental favourability models
# for each species:

data(rotif.env)
```

```

names(rotif.env)

rotifers.fav <- multGLM(data = rotif.env, sp.cols = 18:47,
var.cols = 5:17, id.col = 1, step = FALSE, trim = TRUE,
Favourability = TRUE)

# get matrix of fuzzy similarity between rotifer species distributions:

# either based on inverse squared distance to presence:
rot.fuz.sim.mat <- simMat(rotifers.isqd[ , -1], method = "Jaccard")

# or on environmental favourability for presence:
rot.fuz.sim.mat <- simMat(rotifers.fav$predictions[ , 32:61],
method = "Jaccard")

head(rot.fuz.sim.mat)

# transpose fuzzy rotifer distribution data to compare
# regional species composition rather than species' distributions:

names(rotifers.isqd)

rot.fuz.reg <- transpose(rotifers.fav$predictions, sp.cols = 32:61,
reg.names = 1)

head(rot.fuz.reg)

# get matrix of fuzzy similarity between (some) regions'
# species compositions:

reg.fuz.sim.mat <- simMat(rot.fuz.reg[ , 1:100], method = "Jaccard")

head(reg.fuz.sim.mat)

```

corSelect

Select among correlated variables based on a given criterion

Description

This function calculates pairwise correlations among the variables in a dataset and, among each pair of variables correlated above a given threshold, excludes the variable with either the highest variance inflation factor (VIF), or the least significant or least informative bivariate (individual) relationship with the response variable (if supplied), according to a specified criterion.

Usage

```
corSelect(data, sp.cols = NULL, var.cols, cor.thresh = 0.8,
select = "p.value", ...)
```

Arguments

| | |
|-------------------------|---|
| <code>data</code> | a data frame containing the response and predictor variables. |
| <code>sp.cols</code> | index number of the column of 'data' that contains the response (e.g. species) variable. Currently, only one 'sp.cols' can be used at the same time, so an error message is returned if <code>length(sp.cols) > 1</code> . If <code>sp.cols = NULL</code> , the function returns only the pairs of variables that are correlated over the given threshold. |
| <code>var.cols</code> | index numbers of the columns of 'data' that contain the predictor variables. |
| <code>cor.thresh</code> | threshold value of correlation coefficient above which predictor variables should be excluded. The default is 0.8. |
| <code>select</code> | character value indicating the criterion for excluding variables among those that are correlated. Can be "p.value" (the default), "AIC", "BIC", or "VIF" (see Details). |
| <code>...</code> | additional arguments to pass to <code>cor</code> , namely the <code>method</code> to use (either "pearson", "kendall" or "spearman" correlation coefficient; the first is the default) and the way to deal with missing values (<code>use = "everything"</code> , "all.obs", "complete.obs", "na.or.complete", or "pairwise.complete.obs"). |

Details

Correlations among variables are problematic in multivariate models, as they inflate the variance of coefficients and thus may bias the interpretation of the effects of those variables on the response (Legendre & Legendre 2012). One of the strategies to circumvent this problem is to eliminate one from each pair of correlated variables, but it is not always straightforward to choose the right variable to exclude a priori.

This function selects among correlated variables, based either on their variance inflation factor (VIF: Marquardt 1970; Mansfield & Helms 1982) within the variables dataset (obtained with the `multicol` function and recalculated iteratively after each variable exclusion); or on their relationship with the response, by building a bivariate model of each individual variable against the response and excluding, among each of two correlated variables, the one with the largest (worst) p-value, AIC (Akaike's Information Criterion: Akaike, 1973) or BIC (Bayesian Information Criterion, also known as Schwarz criterion, SBC or SBIC: Schwarz, 1978), which it calculates with the FDR function.

If `sp.cols` is left `NULL` and the 'select' criterion is other than "VIF", the function returns only the pairs of variables that are correlated above the given threshold. If the 'select' criterion requires assessing bivariate relationships and `sp.cols` is provided, the function uses only the rows of the dataset where this column (used as the response variable) contains finite values against which the predictor variables can be modelled; rows with NA or NaN in `sp.cols` are thus excluded from the calculation of correlations among predictor variables.

Value

This function returns a list of 7 elements, unless `sp.cols = NULL`, in which case it returns only the first of these elements:

`high.correlations`
data frame showing the pairs of input variables that are correlated above the given threshold, and their correlation coefficient.

`bivariate.significance`
data frame with the individual p-value, AIC and BIC (if one of these was the 'select' criterion) of each of the highly correlated variables against the response variable.

`excluded.vars`
character vector containing the names of the variables to be excluded (i.e., from each highly correlated pair, the variable with the larger (worse) 'select' score.

`selected.vars`
character vector containing the names of the variables to be selected (i.e., the non-correlated variables and, from each correlated pair, the variable with the smaller (better) 'select' score).

`selected.var.cols`
integer vector containing the column indices of the selected variables in 'data'.

`strongest.remaining.corr`
numerical value indicating the highest (absolute) correlation coefficient among the selected variables.

`remaining.multicollinearity`
data frame showing the multicollinearity among the selected variables.

Author(s)

A. Marcia Barbosa

References

- Akaike, H. (1973) Information theory and an extension of the maximum likelihood principle. In: Petrov B.N. & Csaki F., 2nd International Symposium on Information Theory, Tsahkadsor, Armenia, USSR, September 2-8, 1971, Budapest: Akademiai Kiado, p. 267-281.
- Legendre P. & Legendre L. (2012) Numerical ecology (3rd edition). Elsevier, Amsterdam: 990 pp.
- Marquardt D.W. (1970) Generalized inverses, ridge regression, biased linear estimation, and non-linear estimation. *Technometrics* 12: 591-612.
- Mansfield E.R. & Helms B.P. (1982) Detecting multicollinearity. *The American Statistician* 36: 158-160.
- Schwarz, G.E. (1978) Estimating the dimension of a model. *Annals of Statistics*, 6 (2): 461-464.

See Also

`multicol`, `FDR`, `cor`

Examples

```
data(rotif.env)

corSelect(rotif.env, var.cols = 5:17)

corSelect(rotif.env, sp.cols = 46, var.cols = 5:17)

corSelect(rotif.env, sp.cols = 46, var.cols = 5:17, cor.thresh = 0.7)

corSelect(rotif.env, sp.cols = 46, var.cols = 5:17, method = "spearman")
```

| | |
|----------|---|
| distPres | <i>(Inverse) distance to the nearest presence</i> |
|----------|---|

Description

This function takes a matrix or data frame containing species presence (1) and absence (0) data and their spatial coordinates (optionally also a pre-calculated distance matrix between all localities), and calculates the (inverse) distance from each locality to the nearest presence locality for each species.

Usage

```
distPres(data, sp.cols, coord.cols = NULL, id.col = NULL,
dist.mat = NULL, method = "euclidian", suffix = "_D", p = 1,
inv = TRUE)
```

Arguments

| | |
|-------------------------|---|
| <code>data</code> | a matrix or data frame containing, at least, two columns with spatial coordinates, and one column per species containing their presence (1) and absence (0) data, with localities in rows. |
| <code>sp.cols</code> | names or index numbers of the columns containing the species presences and absences in <code>data</code> . It must contain only zeros (0) for absences and ones (1) for presences. |
| <code>coord.cols</code> | names or index numbers of the columns containing the spatial coordinates in <code>data</code> (in this order, x and y, or longitude and latitude). |
| <code>id.col</code> | optionally, the name or index number of a column (to be included in the output) containing locality identifiers in <code>data</code> . |
| <code>dist.mat</code> | optionally, if you do not want distances calculated with any of the methods available in <code>dist</code> , you may provide a distance matrix calculated elsewhere for the localities in <code>data</code> . |
| <code>method</code> | the method with which to calculate distances between localities. Available options are those of <code>dist</code> . The default is "euclidian". |
| <code>suffix</code> | character indicating the suffix to add to the distance columns in the resulting data frame. The default is "_D". |

| | |
|------------------|--|
| <code>p</code> | the power to which distance should be raised. The default is 1; use 2 or higher if you want more conservative distances. |
| <code>inv</code> | logical value indicating whether distance should be inverted, so that it varies between 0 and 1 and higher values mean closer to presence. The default is TRUE, which is adequate as a fuzzy version of presence-absence (for using e.g. with <code>fuzSim</code> and <code>simMat</code>). In this case, presences maintain the value 1, and inverse distance to presence is calculated only for absence localities. |

Details

This function can be used to calculate a simple spatial interpolation model of a species' distribution (e.g. Barbosa 2015, Areias-Guerreiro et al. 2016).

Value

`distPres` returns a matrix or data frame containing the identifier column (if provided in `id.col`) and one column per species containing the distance (inverse squared by default) from each locality to the nearest presence of that species.

Author(s)

A. Marcia Barbosa

References

Areias-Guerreiro J., Mira A. & Barbosa A.M. (2016) How well can models predict changes in species distributions? A 13-year-old otter model revisited. *Hystrix - Italian Journal of Mammalogy*, in press. DOI: <http://dx.doi.org/10.4404/hystrix-27.1-11867>

Barbosa A.M. (2015) fuzzySim: applying fuzzy logic to binary similarity indices in ecology. *Methods in Ecology and Evolution*, 6: 853-858

See Also

`dist`

Examples

```
data(rotif.env)

head(rotif.env)

names(rotif.env)

# calculate plain distance to presence:

rotifers.dist <- distPres(rotif.env, sp.cols = 18:47,
  coord.cols = c("Longitude", "Latitude"), id.col = 1, p = 1,
  inv = FALSE, suffix = "_D")
```



```

head(rotifers.dist)

# calculate inverse squared distance to presence:

rotifers.invd2 <- distPres(rotif.env, sp.cols = 18:47,
  coord.cols = c("Longitude", "Latitude"), id.col = 1, p = 2,
  inv = TRUE, suffix = "_iDsqr")

head(rotifers.invd2)

```

| | |
|-----|----------------------|
| Fav | <i>Favourability</i> |
|-----|----------------------|

Description

Environmental (prevalence-independent) favourability for a species' presence

Usage

```

Fav(model = NULL, obs = NULL, pred = NULL, nln0 = NULL,
  sample.preval = NULL, method = "RBV", true.preval = NULL)

```

Arguments

| | |
|----------------------------|---|
| <code>model</code> | a model object of class "glm" and binomial family. |
| <code>obs</code> | a vector of the 1/0 values of the modelled binary variable. This argument is ignored if <code>model</code> is provided. |
| <code>pred</code> | a vector of predicted probability values for <code>obs</code> , given e.g. by logistic regression. This argument is ignored if <code>model</code> is provided. |
| <code>nln0</code> | alternatively to <code>obs</code> , an integer vector of length 2 providing the total numbers of modelled ones and zeros, in this order. Ignored if <code>obs</code> or <code>model</code> is provided. |
| <code>sample.preval</code> | alternatively to <code>obs</code> or <code>nln0</code> , the prevalence (proportion of positive cases) of the modelled binary variable in the modelled data. Ignored if <code>model</code> is provided. |
| <code>method</code> | either "RBV" for the original Real, Barbosa & Vargas (2006) procedure, or "AT" for the modification proposed by Albert & Thuiller (2008) (but see Details). |
| <code>true.preval</code> | the true prevalence (as opposed to sample prevalence), necessary if you want to use the AT method. |

Details

Logistic regression (Generalised Linear Model with binomial error distribution and a logit link) is widely used for modelling species' potential distributions using presence/absence data and a set of categorical or continuous predictor variables. However, this GLM incorporates the prevalence (proportion of presences) of the species in the training sample, which affects the probability values produced. Barbosa (2006) and Real, Barbosa & Vargas (2006) proposed an environmental favourability function which is based on logistic regression but cancels out uneven proportions of presences and absences in the modelled data. Favourability thus assesses the extent to which the environmental conditions change the probability of occurrence of a species with respect to its overall prevalence in the study area. Model predictions become, therefore, directly comparable among species with different prevalences. The favourability function is implemented in the **fuzzySim** package and is also in the SAM (Spatial Analysis in Macroecology) software (Rangel et al. 2010).

Using simulated data, Albert & Thuiller (2008) proposed a modification to the favourability function, but it requires knowing the true prevalence of the species (not just the prevalence in the studied sample), which is rarely possible in real-world modelling. Besides, this suggestion was based on the misunderstanding that the favourability function was a way to obtain the probability of occurrence when prevalence differs from 50%, which is incorrect (see Acevedo & Real 2012).

To get environmental favourability with either the Real, Barbosa & Vargas ("RBV") or the Albert & Thuiller ("AT") method, you just need to get a probabilistic model (e.g. logistic regression) from your data and then use the `Fav` function. Input data for this function are either a model object resulting from the `glm` function, or the presences-absences (1-0) of your species and the corresponding presence probability values, obtained e.g. with `predict(mymodel, mydata, type = "response")`. Alternatively to the presences-absences, you can provide either the sample prevalence or the numbers of presences and absences. In case you want to use the "AT" method, you also need to provide the true (absolute) prevalence of your species.

Value

A numeric vector of the favourability values corresponding to the input probability values.

Author(s)

A. Marcia Barbosa

References

- Acevedo P. & Real R. (2012) Favourability: concept, distinctive characteristics and potential usefulness. *Naturwissenschaften* 99: 515-522
- Albert C.H. & Thuiller W. (2008) Favourability functions versus probability of presence: advantages and misuses. *Ecography* 31: 417-422.
- Barbosa A.M.E. (2006) Modelacion de relaciones biogeograficas entre predadores, presas y parasitos: implicaciones para la conservacion de mamiferos en la Peninsula Iberica. PhD Thesis, University of Malaga (Spain).
- Rangel T.F.L.V.B, Diniz-Filho J.A.F & Bini L.M. (2010) SAM: a comprehensive application for Spatial Analysis in Macroecology. *Ecography* 33: 46-50.
- Real R., Barbosa A.M. & Vargas J.M. (2006) Obtaining environmental favourability functions from logistic regression. *Environmental and Ecological Statistics* 13: 237-245.

See Also

`glm, multGLM`

Examples

```
# obtain a probability model and its predictions:

data(rotif.env)

names(rotif.env)

mod <- with(rotif.env, glm(Abrigh ~ Area + Altitude +
  AltitudeRange + HabitatDiversity + HumanPopulation,
  family = binomial))

prob <- predict(mod, data = rotif.env, type = "response")

# obtain predicted favourability in different ways:

Fav(model = mod)

Fav(obs = rotif.env$Abrigh, pred = prob)

Fav(pred = mod$fitted.values, nln0 = c(112, 179))

Fav(pred = mod$fitted.values, sample.preval = 0.38)
```

FDR

False Discovery Rate

Description

Calculate the false discovery rate (type I error) under repeated testing and determine which variables to select and to exclude from multivariate analysis.

Usage

```
FDR(data = NULL, sp.cols = NULL, var.cols = NULL, pvalues = NULL,
  model.type = NULL, family = "auto", correction = "fdr", q = 0.05,
  verbose = TRUE, simplif = FALSE)
```

Arguments

| | |
|----------------------|---|
| <code>data</code> | a data frame containing the response and predictor variables (one in each column). |
| <code>sp.cols</code> | index number of the column containing the response variable (currently implemented for only one response variable at a time). |

| | |
|-------------------------|---|
| <code>var.cols</code> | index numbers of the columns containing the predictor variables. |
| <code>pvalues</code> | optionally, instead of <code>data</code> , <code>sp.cols</code> and <code>var.cols</code> , a data frame with the names of the predictor variables in the first column and their bivariate p-values (obtained elsewhere) in the second column. Example: <code>pvalues <- data.frame(var = letters[1:5], pval = c(0.02, 0.004, 0.07, 0.03, 0.05))</code> . |
| <code>model.type</code> | this argument (previously a character value, either "LM" or "GLM") is now deprecated and ignored with a warning if provided. This information is now included in argument <code>family</code> - e.g., if you want linear models (LM), you can set <code>family = 'gaussian'</code> . |
| <code>family</code> | The error distribution and (optionally) the link function to use (see <code>glm</code> or <code>family</code> for details). The default, <code>'auto'</code> , automatically uses 'binomial' family for response variables containing only values 0 and 1, 'poisson' for positive integer responses (i.e. count data), and 'gaussian' (i.e. linear models) otherwise. |
| <code>correction</code> | the correction procedure to apply to the p-values; see <code>p.adjust.methods</code> for available options and <code>p.adjust</code> for more information. The default is "fdr". |
| <code>q</code> | the threshold value of FDR-corrected significance above which to reject variables. Defaults to 0.05. |
| <code>verbose</code> | logical value indicating whether to display messages. |
| <code>simplif</code> | logical value indicating if simplified results should be provided (see Value). |

Details

It is common in ecology to search for statistical relationships between species' occurrence and a set of predictor variables. However, when a large number of variables is analysed (compared to the number of observations), false findings may arise due to repeated testing. Garcia (2003) recommended controlling the false discovery rate (FDR; Benjamini & Hochberg 1995) in ecological studies. The `p.adjust` R function performs this and other corrections to the significance (p) values of variables under repeated testing. The `FDR` function performs repeated regressions (either linear or binary logistic) or uses already-obtained p values for a set of variables; calculates the FDR with `p.adjust`; and shows which variables should be retained for or excluded from further multivariate analysis according to their corrected p values (see, for example, Barbosa, Real & Vargas 2009).

The `FDR` function uses the Benjamini & Hochberg ("BH", alias "fdr") correction by default, but check the `p.adjust` documentation for other available methods, namely "BY", which allows for non-independent data. Input data may be the response variable (for example, the presence-absence or abundance of a species) and the predictors (a table with one independent variable in each column, with the same number of rows and in the same order as the response); there should be no missing values in the data. Alternatively, you may already have performed the univariate regressions and have a set of variables and corresponding p values which you want to correct with FDR; in this case, get a table with your variables' names in the first column and their p values in the second column, and supply it as the `pvalues` argument (no need to provide response or predictors in this case).

Value

If `simplif = TRUE`, this function returns a data frame with the variables' names as row names and 4 columns containing, respectively, their individual (bivariate) coefficients against the response,

their individual AIC (Akaike's Information Criterion; Akaike, 1973), BIC (Bayesian Information Criterion, also known as Schwarz criterion, SBC, SBIC; Schwarz, 1978), p-value and adjusted p-value according to the applied correction. If `simplif = FALSE` (the default), the result is a list of two such data frames:

```
exclude      with the variables to exclude.
select      with the variables to select (under the given q value).
```

Author(s)

A. Marcia Barbosa

References

- Akaike, H. (1973) Information theory and an extension of the maximum likelihood principle. In: Petrov B.N. & Csaki F., 2nd International Symposium on Information Theory, Tsahkadsor, Armenia, USSR, September 2-8, 1971, Budapest: Akademiai Kiado, p. 267-281.
- Barbosa A.M., Real R. & Vargas J.M (2009) Transferability of environmental favourability models in geographic space: The case of the Iberian desman (*Galemys pyrenaicus*) in Portugal and Spain. *Ecological Modelling* 220: 747-754
- Benjamini Y. & Hochberg Y. (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society, Series B* 57: 289-300
- Garcia L.V. (2003) Controlling the false discovery rate in ecological research. *Trends in Ecology and Evolution* 18: 553-554
- Schwarz, G.E. (1978) Estimating the dimension of a model. *Annals of Statistics*, 6 (2): 461-464.

See Also

```
p.adjust
```

Examples

```
data(rotif.env)

names(rotif.env)

FDR(data = rotif.env, sp.cols = 18, var.cols = 5:17)

FDR(data = rotif.env, sp.cols = 18, var.cols = 5:17, simplif = TRUE)
```

fuzSim

*Fuzzy similarity***Description**

This function calculates fuzzy similarity, based on a fuzzy version of the binary similarity index specified in `method`, between two binary (0 or 1) or fuzzy (between 0 and 1) variables.

Usage

```
fuzSim(x, y, method, na.rm = TRUE)
```

Arguments

| | |
|---------------------|---|
| <code>x</code> | a vector of (optionally fuzzy) presence-absence data, with 1 meaning presence, 0 meaning absence, and values in between meaning fuzzy presence (or the degree to which each locality belongs to the set of species presences, or to which each species belongs to the locality; Zadeh, 1965). Fuzzy presence-absence can be obtained, for example, with functions <code>multGLM</code> , <code>distPres</code> or <code>multTSA</code> in this package. |
| <code>y</code> | a vector similar to <code>x</code> , of the same length and in the same order. |
| <code>method</code> | the similarity index to use. Currently available options are "Jaccard", "Sorensen", "Simpson" and "Baroni" (see Details). |
| <code>na.rm</code> | logical value indicating whether NA values should be ignored. The default is TRUE. |

Details

Similarity between ecological communities, beta diversity patterns, biotic regions, and distributional relationships among species are commonly determined based on pair-wise (dis)similarities in species' occurrence patterns. Some of the most commonly employed similarity indices are those of Jaccard (1901), Sorensen (1948), Simpson (1960) and Baroni-Urbani & Buser (1976), which are here implemented in their fuzzy versions (Barbosa, 2015), able to deal with both binary and fuzzy data. Jaccard's and Baroni's indices have associated tables of significant values (Baroni-Urbani & Buser 1976, Real & Vargas 1996, Real 1999).

Value

The function returns a value between 0 and 1 representing the fuzzy similarity between `x` and `y`. Note, for example, that Jaccard similarity can be converted to dissimilarity (or Jaccard distance) if subtracted from 1, while 1-Sorensen is not a proper distance metric as it lacks the property of triangle inequality (see http://en.wikipedia.org/wiki/S%C3%B8rensen%E2%80%99s_dice_coefficient).

Note

The formulas used in this function may look slightly different from some of their published versions (e.g. Baroni-Urbani & Buser 1976), not only because the letters are switched, but because here the A and B are the numbers of attributes present in each element, whether or not they are also present in the other one. Thus, our 'A+B' is equivalent to 'A+B+C' in formulas where A and B are the numbers of attributes present in one but not the other element, and our A+B-C is equivalent to their A+B+C. The formulas used here (adapted from Olivero et al. 1998) are faster to calculate, visibly for large datasets.

Author(s)

A. Marcia Barbosa

References

- Barbosa A.M. (2015) fuzzySim: applying fuzzy logic to binary similarity indices in ecology. *Methods in Ecology and Evolution*, 6: 853-858.
- Baroni-Urbani C. & Buser M.W. (1976) Similarity of Binary Data. *Systematic Zoology*, 25: 251-259
- Jaccard P. (1901) Etude comparative de la distribution florale dans une portion des Alpes et des Jura. *Memoires de la Societe Vaudoise des Sciences Naturelles*, 37: 547-579
- Olivero J., Real R. & Vargas J.M. (1998) Distribution of breeding, wintering and resident waterbirds in Europe: biotic regions and the macroclimate. *Ornis Fennica*, 75: 153-175
- Real R. (1999) Tables of significant values of Jaccard's index of similarity. *Miscellanea Zoologica* 22: 29:40
- Real R. & Vargas J.M (1996) The probabilistic basis of Jaccard's index of similarity. *Systematic Biology* 45: 380-385
- Simpson, G.G. (1960) Notes on the measurement of faunal resemblance. *Amer. J. Sci.* 258A, 300-311
- Sorensen T. (1948) A method of establishing groups of equal amplitude in plant sociology based on similarity of species and its application to analyses of the vegetation on Danish commons. *Kongelige Danske Videnskabernes Selskab*, 5(4): 1-34
- Zadeh L.A. (1965) Fuzzy sets. *Information and Control*, 8: 338-353

See Also

`simMat`; `modOverlap`

Examples

```
data(rotif.env)

names(rotif.env)

# you can calculate similarity between binary species occurrence patterns:
```

```
fuzSim(rotif.env[, "Abrigh"], rotif.env[, "Afissa"], method = "Jaccard")
fuzSim(rotif.env[, "Abrigh"], rotif.env[, "Afissa"], method = "Sorensen")
fuzSim(rotif.env[, "Abrigh"], rotif.env[, "Afissa"], method = "Simpson")
fuzSim(rotif.env[, "Abrigh"], rotif.env[, "Afissa"], method = "Baroni")

# or you can model environmental favourability for these species
# and calculate fuzzy similarity between their environmental predictions
# which goes beyond the strict coincidence of their occurrence records:

fav <- multGLM(rotif.env, sp.cols = 18:19, var.cols = 5:17, step = TRUE,
FDR = TRUE, trim = TRUE, P = FALSE, Fav = TRUE) $ predictions

fuzSim(fav[, "Abrigh_F"], fav[, "Afissa_F"], method = "Jaccard")
fuzSim(fav[, "Abrigh_F"], fav[, "Afissa_F"], method = "Sorensen")
fuzSim(fav[, "Abrigh_F"], fav[, "Afissa_F"], method = "Simpson")
fuzSim(fav[, "Abrigh_F"], fav[, "Afissa_F"], method = "Baroni")
```

fuzzyOverlay

Row-wise overlay operations based on fuzzy logic

Description

Logical and set operations are useful for comparative distribution modelling, to assess consensus or mismatches between the predictions of different models, and to quantify differences between models obtained for different time periods. Fuzzy set theory (Zadeh 1965, Barbosa & Real 2012) allows performing such operations without converting the predictions from continuous to binary, with the inherent application of arbitrary thresholds and over-simplification of model predictions. The result is a continuous numerical value quantifying the intersection, union, sum, or other operation among model predictions, whether binary or continuous.

Usage

```
fuzzyOverlay(data, overlay.cols = 1:ncol(data), op = "intersection",
na.rm = FALSE, round.digits = 2)
```

Arguments

| | |
|---------------------------|--|
| <code>data</code> | matrix or data frame containing the model predictions to compare. |
| <code>overlay.cols</code> | vector of the names or index numbers of the columns to compare. The default is all columns in data. |
| <code>op</code> | character value indicating the operation to perform between the prediction columns in data. Can be "consensus" for the arithmetic mean of predictions (or the fuzzy equivalent of the proportion of models that agree that the species occurs at each site), "fuzzy_and" or "intersection" for fuzzy intersection; "fuzzy_or" or "union" for fuzzy union; "prob_and" or "prob_or" for probabilistic and/or, respectively (see Details); "maintenance" for the values where all predictions for the same row (rounded to the number of digits specified in the next argument) |

are the same. If `data` has only two columns to compare, you can also calculate "xor" for exclusive 'or', "AnotB" for the occurrence of the species in column 1 in detriment of that in column 2, "expansion" for the prediction increase in rows where column 2 has higher values than column 1, "contraction" for the prediction decrease in rows where column 2 has lower values than column 1, or "change" for a mix of the latter two, with positive values where there has been an increase and negative values where there was decrease in favourability from columns 1 to 2. For expansion, contraction and maintenance, rows where the values do not satisfy the condition (i.e. second column larger, smaller, or roughly equal to the first column) get a value of zero.

`na.rm` logical value indicating if NA values should be ignored. The default is `FALSE`, so rows with NA in any of the prediction columns get NA as a result.

`round.digits` integer value indicating the number of decimal places to be used if `op = "maintenance"`. The default is 2.

Details

If your predictions are probabilities, "prob_and" (probabilistic 'and') gives the probability of all species in `data` occurring simultaneously by multiplying all probabilities; and "prob_or" (probabilistic 'or') gives the probability of any of them occurring at each site. These can be quite restrictive, though; probabilistic "and" can give particularly unrealistically small values.

If you have (or convert your probabilities to) favourability predictions, which can be used directly with fuzzy logic (Real et al. 2006; see `Fav` function), you can use "fuzzy_and" or "intersection" to get the favourability for all species co-occurring at each site, and "fuzzy_or" or "union" to get favourability for any of them to occur at each site (Barbosa & Real 2012).

Value

This function returns a vector, with length equal to the number of rows in `data`, containing the row-wise result of the operation performed.

Author(s)

A. Marcia Barbosa

References

- Barbosa A.M. & Real R. (2012) Applying fuzzy logic to comparative distribution modelling: a case study with two sympatric amphibians. *The Scientific World Journal*, 2012, Article ID 428206
- Real R., Barbosa A.M. & Vargas J.M. (2006) Obtaining environmental favourability functions from logistic regression. *Environmental and Ecological Statistics* 13: 237-245.
- Zadeh, L.A. (1965) Fuzzy sets. *Information and Control*, 8: 338-353

See Also

`fuzSim`, `modOverlap` and `fuzzyRangeChange` for overall (not row-wise) comparisons among model predictions.

Examples

```

data(rotif.env)

names(rotif.env)

# get model predictions for 3 of the species in rotif.env:

mods <- multGLM(rotif.env, sp.cols = 18:20, var.cols = 5:17, id.col = 1,
step = TRUE, FDR = TRUE, trim = TRUE)

preds <- mods$predictions[ , c("Abrigh_F", "Afissa_F", "Aperiod_F")]

# calculate intersection and union among those predictions:

preds$intersect <- fuzzyOverlay(preds, op = "intersection")

preds$union <- fuzzyOverlay(preds, op = "union")

head(preds)

# imagine you have a model prediction for species 'Abrigh' in a future time
# (here we will create one by randomly jittering the current predictions)

preds$Abrigh_imag <- jitter(preds[ , "Abrigh_F"], amount = 0.2)
preds$Abrigh_imag[preds$Abrigh_imag < 0] <- 0
preds$Abrigh_imag[preds$Abrigh_imag > 1] <- 1

# you can calculate row-wise prediction changes from Abrigh to Abrigh_imag:

preds$Abrigh_exp <- fuzzyOverlay(preds, overlay.cols = c("Abrigh_F",
"Abrigh_imag"), op = "expansion")

preds$Abrigh_contr <- fuzzyOverlay(preds, overlay.cols = c("Abrigh_F",
"Abrigh_imag"), op = "contraction")

preds$Abrigh_chg <- fuzzyOverlay(preds, overlay.cols = c("Abrigh_F",
"Abrigh_imag"), op = "change")

preds$Abrigh_maint <- fuzzyOverlay(preds, overlay.cols = c("Abrigh_F",
"Abrigh_imag"), op = "maintenance")

head(preds)

```

Description

This function quantifies overall range change (expansion, contraction, maintenance and balance) based on either presence-absence data or the continuous predictions of two models.

Usage

```
fuzzyRangeChange(pred1, pred2, number = TRUE, prop = TRUE,
na.rm = TRUE, round.digits = 2, measures = c("Gain", "Loss",
"Stable presence", "Stable absence", "Balance"), plot = TRUE, ...)
```

Arguments

| | |
|---------------------------|---|
| <code>pred1</code> | numeric vector containing the predictions (between 0 and 1) of the model that will serve as reference. |
| <code>pred2</code> | numeric vector containing the predictions (between 0 and 1) of the model whose change will be calculated. Must be of the same length and in the same order as <code>pred1</code> . |
| <code>number</code> | logical value indicating if results should include the fuzzy number of cases. The default is TRUE. |
| <code>prop</code> | logical value indicating if results should include the proportion of the total number of cases. The default is TRUE. |
| <code>na.rm</code> | logical value indicating whether NA values should be ignored. The default is TRUE. |
| <code>round.digits</code> | argument to pass to <code>fuzzyOverlay</code> , indicating the number of decimal places to which to round <code>pred</code> for calculating 'maintenance' or 'stability'. The default is 2. |
| <code>measures</code> | character vector listing the range change measures to calculate. The default is all available measures. |
| <code>plot</code> | logical value indicating whether to make a barplot with the results. The default is TRUE. |
| <code>...</code> | additional arguments to be passed to the <code>barplot</code> function (if <code>plot = TRUE</code>). |

Value

This function returns a data frame with the following values in different rows (among those that are included in `measures`):

| | |
|------------------------------|---|
| <code>Gain</code> | sum of the predicted values that have increased from <code>pred1</code> to <code>pred2</code> (fuzzy equivalent of the number of gained presences) |
| <code>Loss</code> | sum of the predicted values that have decreased from <code>pred1</code> to <code>pred2</code> (fuzzy equivalent of the number of lost presences) |
| <code>Stable_presence</code> | fuzzy equivalent of the number of predicted presences that have remained as such (when rounded to <code>round.digits</code>) between <code>pred1</code> and <code>pred2</code> |
| <code>Stable_absence</code> | fuzzy equivalent of the number of predicted absences that have remained as such (when rounded to <code>round.digits</code>) between <code>pred1</code> and <code>pred2</code> |

Balance sum of the change in predicted values from pred1 to pred2 (fuzzy equivalent of the balance of gained and lost presences)

If `prop = TRUE` (the default), there is an additional column named "Proportion" in which these values are divided by the total number of reference values (i.e., the fuzzy range or non-range size). If `plot = TRUE` (the default), a barplot is also produced representing the last column of the result data frame.

Author(s)

A. Marcia Barbosa

See Also

`fuzSim`, `modOverlap` for other ways to compare models; `fuzzyOverlay` for row-wise model comparisons

Examples

```
# get an environmental favourability model for a rotifer species:

data(rotif.env)

names(rotif.env)

fav_current <- multGLM(rotif.env, sp.cols = 18, var.cols = 5:17,
step = TRUE, FDR = TRUE, trim = TRUE, P = FALSE, Fav = TRUE) $
predictions

# imagine you have a model prediction for this species in a future time
# (here we will create one by randomly jittering the current predictions)

fav_imag <- jitter(fav_current, amount = 0.2)
fav_imag[fav_imag < 0] <- 0
fav_imag[fav_imag > 1] <- 1

# calculate range change given by current and imaginary future predictions:

fuzzyRangeChange(fav_current, fav_imag)

fuzzyRangeChange(fav_current, fav_imag, number = FALSE)

fuzzyRangeChange(fav_current, fav_imag, ylim = c(-1, 1),
ylab = "Proportional change")
```

`getPreds`*Get model predictions*

Description

This function allows getting the predictions of multiple models when applied to a given dataset. It can be useful if you have a list of model objects (e.g. resulting from `multGLM`) and want to apply them to a new data set containing the same variables for another region or time period. There are options to include the logit link (`Y`) and/or Favourability (see `Fav`).

Usage

```
getPreds(data, models, id.col = NULL, Y = FALSE, P = TRUE,
  Favourability = TRUE, incl.input = FALSE)
```

Arguments

| | |
|----------------------------|--|
| <code>data</code> | the data frame to which to apply the <code>models</code> to get their predictions; must contain all variables (with the same names, case-sensitive) included in any of the <code>models</code> . |
| <code>models</code> | a list of model objects obtained e.g. with function <code>glm</code> or <code>multGLM</code> . |
| <code>id.col</code> | optionally, the index number of a column of data containing row identifiers, to be included in the result. Ignored if <code>incl.input = TRUE</code> . |
| <code>Y</code> | logical, whether to include the logit link (<code>y</code>) value in the predictions. |
| <code>P</code> | logical, whether to include the probability value in the predictions. |
| <code>Favourability</code> | logical, whether to include Favourability in the predictions (see <code>Fav</code>). |
| <code>incl.input</code> | logical, whether to include input columns in the output. The default is <code>FALSE</code> . |

Value

This function returns a data frame containing the model predictions, next to the `id.col` if provided, and to the input data if `incl.input = TRUE`.

Author(s)

A. Marcia Barbosa

See Also

`multGLM`, `predict`

Examples

```
data(rotif.env)

names(rotif.env)

# identify rotifer data in the Eastern and Western hemispheres:
unique(rotif.env$CONTINENT)

rotif.env$HEMISPHERE <- "Eastern"

rotif.env$HEMISPHERE[rotif.env$CONTINENT %in%
c("NORTHERN_AMERICA", "SOUTHERN_AMERICA")] <- "Western"

head(rotif.env)

# separate the rotifer data into hemispheres

east.hem <- rotif.env[rotif.env$HEMISPHERE == "Eastern", ]
west.hem <- rotif.env[rotif.env$HEMISPHERE == "Western", ]

# make models for 3 of the species in rotif.env based on their distribution
# in the Eastern hemisphere:

mods <- multGLM(east.hem, sp.cols = 18:20, var.cols = 5:17,
id.col = 1, step = FALSE, FDR = FALSE, trim = FALSE)

# get the models' predictions for the Western hemisphere dataset:

preds <- getPreds(west.hem, models = mods$models, P = TRUE,
Favourability = TRUE)

head(preds)
```

*integerCols**Classify integer columns*

Description

This function detects which numeric columns in a data frame contain only whole numbers, and converts those columns to integer class, so that they take up less space.

Usage

```
integerCols(data)
```

Arguments

`data` a data frame containing possibly integer columns classified as `numeric`.

Value

The function returns a data frame with the same columns as `data`, but with those that are numeric and contain only whole numbers (possibly including NA) now classified as `integer`.

Author(s)

A. Marcia Barbosa

See Also

`is.integer`, `as.integer`, `multConvert`

Examples

```
dat <- data.frame(
  var1 = 1:10,
  var2 = as.numeric(1:10),
  var3 = as.numeric(c(1:4, NA, 6:10)),
  var4 = as.numeric(c(1:3, NaN, 5, Inf, 7, -Inf, 9:10)),
  var5 = as.character(1:10),
  var6 = seq(0.1, 1, by = 0.1),
  var7 = letters[1:10]
) # creates a sample data frame

dat

str(dat)
# var2 classified as 'numeric' but contains only whole numbers
# var3 same as var2 but containing also NA values
# var4 same as var2 but containing also NaN and infinite values
# var5 contains only whole numbers but initially classified as factor

dat <- integerCols(dat)

str(dat)
# var2 and var3 now classified as 'integer'
# var4 remains as numeric because contains infinite and NaN
# (not integer) values
# var5 remains as factor
```

modelTrim

Trim off non-significant variables from a model

Description

This function performs a stepwise removal of non-significant variables from a model.

Usage

```
modelTrim(model, method = "summary", alpha = 0.05)
```

Arguments

| | |
|--------|---|
| model | a model object. |
| method | the method for getting the individual p-values. Can be either "summary" for the p-values of the coefficient estimates, or "anova" for the p-values of the variables themselves (see Details). |
| alpha | the p-value above which a variable is removed. |

Details

Stepwise variable selection is a common procedure for simplifying models. It maximizes predictive efficiency in an objective and reproducible way, and is useful when the individual importance of the predictors is not known a priori (Hosmer & Lemeshow, 2000). The `step` R function performs such procedure using an information criterion (AIC) to select the variables, but it often leaves variables that are not significant in the model. Such variables can be subsequently removed with a manual stepwise procedure (e.g. Crawley 2007, p. 442; Barbosa & Real 2010, 2012; Estrada & Arroyo 2012). The `modelTrim` function performs such removal automatically until all remaining variables are significant. It can also be applied to a full model (i.e., without previous use of the `step` function), as it serves as a backward stepwise selection procedure based on the significance of the coefficients (if `method = "summary"`, the default) or on the significance of the variables themselves (if `method = "anova"`, better when there are categorical variables in the model).

Value

The input model object after removal of non-significant variables.

Author(s)

A. Marcia Barbosa

References

Barbosa A.M. & Real R. (2010) Favourable areas for expansion and reintroduction of Iberian lynx accounting for distribution trends and genetic diversity of the European rabbit. *Wildlife Biology in Practice* 6: 34-47

Barbosa A.M. & Real R. (2012) Applying fuzzy logic to comparative distribution modelling: a case study with two sympatric amphibians. The Scientific World Journal, Article ID 428206

Crawley M.J. (2007) The R Book. John Wiley & Sons, Chichester (UK)

Estrada A. & Arroyo B. (2012) Occurrence vs abundance models: Differences between species with varying aggregation patterns. Biological Conservation, 152: 37-45

Hosmer D. W. & Lemeshow S. (2000) Applied Logistic Regression (2nd ed). John Wiley and Sons, New York

See Also

step

Examples

```
# load sample data:

data(rotif.env)

names(rotif.env)

# build a stepwise model of a species' occurrence based on
# some of the variables:

mod <- with(rotif.env, step(glm(Abrigh ~ Area + Altitude + AltitudeRange +
HabitatDiversity + HumanPopulation, family = binomial)))

# examine the model:

summary(mod) # contains non-significant variables

# use modelTrim to get rid of non-significan effects:

mod <- modelTrim(mod)

summary(mod) # only significant variables now
```

modOverlap

Overall overlap between model predictions

Description

This function calculates the degree of overlap between the predictions of two models, using niche comparison metrics such as Schoener's D, Hellinger distance and Warren's I.

Usage

```
modOverlap(pred1, pred2, na.rm = TRUE)
```

Arguments

| | |
|--------------------|---|
| <code>pred1</code> | numeric vector of the predictions of a generalized linear model (values between 0 and 1). |
| <code>pred2</code> | numeric vector of the predictions of another generalized linear model; must be of the same length and in the same order as <code>pred1</code> . |
| <code>na.rm</code> | logical value indicating whether NA values should be removed prior to calculation. |

Details

See Warren et al. (2008).

Value

This function returns a list of 3 metrics:

| | |
|----------------------------|--|
| <code>SchoenerD</code> | Schoener's (1968) D statistic for niche overlap, varying between 0 (no overlap) and 1 (identical niches). |
| <code>WarrenI</code> | the I index of Warren et al. (2008), based on Hellinger distance (below) but re-formulated to also vary between 0 (no overlap) and 1 (identical niches). |
| <code>HellingerDist</code> | Hellinger distance (as in van der Vaart 1998, p. 211) between probability distributions, varying between 0 and 2. |

Note

A function providing similar measures, `niche.overlap`, is available in package **phyloclim**, but it requires complex and software-specific input data formats.

Author(s)

A. Marcia Barbosa

References

Schoener T.W. (1968) Anolis lizards of Bimini: resource partitioning in a complex fauna. *Ecology* 49: 704-726

van der Vaart A.W. (1998) Asymptotic statistics. Cambridge Univ. Press, Cambridge (UK)

Warren D.L., Glor R.E. & Turelli M. (2008) Environmental niche equivalency versus conservatism: quantitative approaches to niche evolution. *Evolution*, 62: 2868-83 (and further ERRATUM)

See Also

`fuzSim`; `niche.overlap` in package **phyloclim**

Examples

```
# get an environmental favourability model for a rotifer species:

data(rotif.env)

names(rotif.env)

fav_current <- multGLM(rotif.env, sp.cols = 18, var.cols = 5:17,
step = TRUE, FDR = TRUE, trim = TRUE, P = FALSE, Fav = TRUE) $
predictions

# imagine you have a model prediction for this species in a future time
# (here we will create one by randomly jittering the current predictions)

fav_imag <- jitter(fav_current, amount = 0.2)
fav_imag[fav_imag < 0] <- 0
fav_imag[fav_imag > 1] <- 1

# calculate niche overlap between current and imaginary future predictions:

modOverlap(fav_current, fav_imag)
```

| | |
|-------------|----------------------------|
| multConvert | <i>Multiple conversion</i> |
|-------------|----------------------------|

Description

This function can simultaneously convert multiple columns of a matrix or data frame.

Usage

```
multConvert(data, conversion, cols = 1:ncol(data))
```

Arguments

| | |
|------------|--|
| data | A matrix or data frame containing columns that need to be converted |
| conversion | the conversion to apply, e.g. <code>as.factor</code> or a custom-made function |
| cols | the columns of data to convert |

Details

Sometimes we need to change the data type (class, mode) of a variable in R. There are various possible conversions, performed by functions like `as.integer`, `as.factor` or `as.character`. If we need to perform the same conversion on a number of variables (columns) in a data frame, we can convert them all simultaneously using this function. By default it converts all the columns in the data frame, but you can specify just a few of them. `multConvert` can also be used to apply other kinds of transformations - for example, if you need to divide some of your columns by 100, just write a function to do this and then use `multConvert` to apply this function to any group of columns.

Value

The input data with the specified columns converted as asked.

Author(s)

A. Marcia Barbosa

Examples

```
data(rotif.env)

str(rotif.env)

# convert the first 4 columns to character:
converted.rotif.env <- multConvert(data = rotif.env,
  conversion = as.character, cols = 1:4)

str(converted.rotif.env)

names(rotif.env)

# divide some columns by 100:

div100 <- function(x) x / 100

rotif.env.cent <- multConvert(data = rotif.env,
  conversion = div100, cols = c(6:10, 12:17))

head(rotif.env.cent)
```

multGLM

GLMs for multiple species with multiple options

Description

This function calculates generalized linear models for a set of (species) presence/absence records in a data frame, with a wide set of options for data partition, variable selection, and output form.

Usage

```
multGLM(data, sp.cols, var.cols, id.col = NULL, family = "binomial",
  test.sample = 0, FDR = FALSE, correction = "fdr", corSelect = FALSE,
  cor.thresh = 0.8, step = TRUE, trace = 0, start = "null.model",
  direction = "both", select = "AIC", Y.prediction = FALSE,
  P.prediction = TRUE, Favourability = TRUE, group.preds = TRUE,
  trim = TRUE, ...)
```

Arguments

| | |
|----------------------------|---|
| <code>data</code> | a data frame in wide format (see <code>splist2presabs</code>) containing, in separate columns, your species' binary (0/1) occurrence data and the predictor variables. |
| <code>sp.cols</code> | index numbers of the columns containing the species data to be modelled. |
| <code>var.cols</code> | index numbers of the columns containing the predictor variables to be used for modelling. |
| <code>id.col</code> | (optional) index number of column containing the row identifiers (if defined, it will be included in the output <code>predictions</code> data frame). |
| <code>family</code> | argument to be passed to the <code>glm</code> function; only 'binomial' is implemented in <code>multGLM</code> so far. |
| <code>test.sample</code> | a subset of data to set aside for subsequent model testing. Can be a value between 0 and 1 for a proportion of the data to choose randomly (e.g. 0.2 for 20%), or an integer number for a particular number of cases to choose randomly among the records in <code>data</code> , or a vector of integers for the index numbers of the particular rows to set aside, or "Huberty" for his rule of thumb based on the number of variables (Huberty 1994, Fielding & Bell 1997). |
| <code>FDR</code> | logical value indicating whether to do a preliminary exclusion of variables based on the false discovery rate (see <code>FDR</code>). The default is <code>FALSE</code> . |
| <code>correction</code> | argument to pass to the <code>FDR</code> function if <code>FDR = TRUE</code> . The default is "fdr", but see <code>p.adjust</code> for more options. |
| <code>corSelect</code> | logical value indicating whether to do a preliminary exclusion of highly correlated variables (see <code>corSelect</code>). The default is <code>FALSE</code> . |
| <code>cor.thresh</code> | numerical value indicating the correlation threshold to pass to <code>corSelect</code> (used only if <code>corSelect = TRUE</code>). |
| <code>step</code> | logical, whether to use the <code>step</code> function to perform a stepwise variable selection (based on AIC or BIC). |
| <code>trace</code> | if positive, information is printed during the running of <code>step</code> . Larger values may give more detailed information. |
| <code>start</code> | character, whether to start with the 'null.model' (so that variable selection starts forward) or with the 'full.model' (so selection starts backward). Used only if <code>step = TRUE</code> . |
| <code>direction</code> | argument to be passed to <code>step</code> specifying the direction of variable selection ('forward', 'backward' or 'both'). Used only if <code>step = TRUE</code> . |
| <code>select</code> | character string specifying the criterion for stepwise selection of variables. Options are "AIC" (Akaike's Information Criterion; Akaike, 1973), the default; or BIC (Bayesian Information Criterion, also known as Schwarz criterion, SBC or SBIC; Schwarz, 1978). Used only if <code>step = TRUE</code> . |
| <code>Y.prediction</code> | logical, whether to include output predictions in the scale of the predictor variables (<code>type = "link"</code> in <code>predict.glm</code>). |
| <code>P.prediction</code> | logical, whether to include output predictions in the scale of the response variable, i.e. probability (<code>type = "response"</code> in <code>predict.glm</code>). |
| <code>Favourability</code> | logical, whether to apply the <code>Favourability</code> function to extract the effect of prevalence on probability (Real et al. 2006) and include its results in the output. |

| | |
|--------------------------|---|
| <code>group.preds</code> | logical, whether to group together predictions of similar type (Y, P or F) in the output <code>predictions</code> table (e.g. if <code>FALSE</code> : <code>sp1_Y</code> , <code>sp1_P</code> , <code>sp1_F</code> , <code>sp2_Y</code> , <code>sp2_P</code> , <code>sp2_F</code> ; if <code>TRUE</code> : <code>sp1_Y</code> , <code>sp2_Y</code> , <code>sp1_P</code> , <code>sp2_P</code> , <code>sp1_F</code> , <code>sp2_F</code>). |
| <code>trim</code> | logical, whether to trim non-significant variables off the models using the <code>modelTrim</code> function; can be used whether or not <code>step</code> is <code>TRUE</code> ; works as a backward variable elimination procedure based on significance. |
| <code>...</code> | additional arguments to be passed to <code>modelTrim</code> . |

Details

This function automatically calculates binomial GLMs for one or more species (or other binary variables) in a data frame. The function can optionally perform `stepwise` variable selection (and it does so by default) instead of forcing all variables into the models, starting from either the null model (the default, so selection starts forward) or from the full model (so selection starts backward) and using Akaike's information criterion (AIC) as a variable selection criterion. Instead or subsequently, it can also perform stepwise removal of non-significant variables from the models using the `modelTrim` function.

There is also an optional preliminary selection of non-correlated variables, and/or of variables with a significant bivariate relationship with the response, based on the false discovery rate (FDR). Note, however, that some variables can be significant in a multivariate model even if they would not have been selected by FDR.

Favourability is also calculated, removing the effect of species prevalence from occurrence probability and thus allowing direct comparisons between models (Real et al. 2006).

By default, all data are used in model training, but you can define an optional `test.sample` to be reserved for model testing afterwards. You may also want to do a previous check for `multicollinearity` among variables, e.g. the variance inflation factor (VIF).

The `multGLM` function will create a list of the resulting models (each with the name of the corresponding species column) and a data frame with their predictions (Y, P and/or F, all of which are optional). If you plan on representing these predictions in a GIS based on `.dbf` tables, remember that `dbf` only allows up to 10 characters in column names; `multGLM` predictions will add 2 characters (`_Y`, `_P` and/or `_F`) to each of your species column names, so use species names/codes with up to 8 characters in the data set that you are modelling. You can create (sub)species name abbreviations with the `spCodes` function.

Value

This function returns a list with the following components:

| | |
|--------------------------|--|
| <code>predictions</code> | a data frame with the model predictions (if either of <code>Y.prediction</code> , <code>P.prediction</code> or <code>Favourability</code> are <code>TRUE</code>). |
| <code>models</code> | a list of the resulting model objects. |

Author(s)

A. Marcia Barbosa

References

- Akaike, H. (1973) Information theory and an extension of the maximum likelihood principle. In: Petrov B.N. & Csaki F., 2nd International Symposium on Information Theory, Tsahkadsor, Armenia, USSR, September 2-8, 1971, Budapest: Akademiai Kiado, p. 267-281.
- Fielding A.H. & Bell J.F. (1997) A review of methods for the assessment of prediction errors in conservation presence/absence models. *Environmental Conservation* 24: 38-49
- Huberty C.J. (1994) *Applied Discriminant Analysis*. Wiley, New York, 466 pp.
- Schaafsma W. & van Vark G.N. (1979) Classification and discrimination problems with applications. Part IIa. *Statistica Neerlandica* 33: 91-126
- Real R., Barbosa A.M. & Vargas J.M. (2006) Obtaining environmental favourability functions from logistic regression. *Environmental and Ecological Statistics* 13: 237-245.
- Schwarz, G.E. (1978) Estimating the dimension of a model. *Annals of Statistics*, 6 (2): 461-464.

See Also

glm, Fav, step, modelTrim, multicol, corSelect

Examples

```
data(rotif.env)

names(rotif.env)

# make models for 3 of the species in rotif.env:

mods <- multGLM(rotif.env, sp.cols = 45:47, var.cols = 5:17, id.col = 1,
step = TRUE, FDR = TRUE, trim = TRUE)

names(mods)

head(mods$predictions)

names(mods$models)

mods$models[[1]]

mods$models[["Ttetra"]]
```

multicol

Analyse multicollinearity in a dataset, including VIF

Description

This function analyses multicollinearity in a set of variables or in a model, including the R-squared, tolerance and variance inflation factor (VIF).

Usage

```
multicol(vars = NULL, model = NULL, reorder = TRUE)
```

Arguments

| | |
|----------------------|---|
| <code>vars</code> | A matrix or data frame containing the numeric variables for which to calculate multicollinearity. Only the 'independent' (predictor, explanatory, right hand side) variables should be entered, as the result obtained for each variable depends on all the other variables present in the analysed data set. |
| <code>model</code> | Alternatively to <code>vars</code> , a <code>glm</code> model object to calculate <code>multicol</code> among the included variables. |
| <code>reorder</code> | logical, whether variables should be output in decreasing order or VIF value rather than in their input order. The default is <code>TRUE</code> . |

Details

Testing collinearity among covariates is a recommended step of data exploration before applying a statistical model (Zuur et al. 2010). However, you can also calculate multicollinearity among the variables already included in a model.

The `multicol` function calculates the degree of multicollinearity in a set of numeric variables, using three closely related measures: R squared (the coefficient of determination of a linear regression of each predictor variable on all other predictor variables, i.e., the amount of variation in each variable that is accounted for by other variables in the dataset); tolerance ($1 - R^2$), i.e. the amount of variation in each variable that is not included in the remaining variables; and the variance inflation factor: $VIF = 1 / (1 - R^2)$, which, in a linear model with these variables as predictors, reflects the degree to which the variance of an estimated regression coefficient is increased due only to the correlations among covariates (Marquardt 1970; Mansfield & Helms 1982).

Value

The function returns a matrix with one row per analysed variable, the names of the variables as row names, and 3 columns: R-squared, Tolerance and VIF.

Author(s)

A. Marcia Barbosa

References

- Marquardt D.W. (1970) Generalized inverses, ridge regression, biased linear estimation, and non-linear estimation. *Technometrics* 12: 591-612.
- Mansfield E.R. & Helms B.P. (1982) Detecting multicollinearity. *The American Statistician* 36: 158-160.
- Zuur A.F., Ieno E.N. & Elphick C.S. (2010) A protocol for data exploration to avoid common statistical problems. *Methods in Ecology and Evolution* 1: 3-14.

See Also

vif in package **HH**, vif in package **usdm**

Examples

```
data(rotif.env)
names(rotif.env)

# calculate multicollinearity among the predictor variables:
multicol(rotif.env[ , 5:17], reorder = FALSE)
multicol(rotif.env[ , 5:17])

# you can also calculate multicol among the variables included in a model:
mod <- step(glm(Abrigh ~ Area + Altitude + AltitudeRange +
HabitatDiversity + HumanPopulation + Latitude + Longitude +
Precipitation + PrecipitationSeasonality + TemperatureAnnualRange
+ Temperature + TemperatureSeasonality + UrbanArea,
data = rotif.env))
multicol(model = mod)

# more examples using R datasets:
multicol(trees)

# you'll get a warning and some NA results if any of the variables
# is not numeric:
multicol(OrchardSprays)

# so define the subset of numeric 'vars' to calculate 'multicol' for:
multicol(OrchardSprays[ , 1:3])
```

multTSA

Trend Surface Analysis for multiple species

Description

This function performs trend surface analysis for multiple species at a time. It converts categorical presence-absence (1-0) data into continuous surfaces denoting the spatial trend in species' occurrence patterns.

Usage

```
multTSA(data, sp.cols, coord.cols, id.col = NULL, degree = 3,
step = TRUE, Favourability = FALSE, suffix = "_TS",
save.models = FALSE)
```

Arguments

| | |
|----------------------------|---|
| <code>data</code> | a matrix or data frame containing, at least, two columns with spatial coordinates, and one column per species containing their presence (1) and absence (0) data, with localities in rows. |
| <code>sp.cols</code> | names or index numbers of the columns containing the species presences and absences in data. Must contain only zeros (0) for absences and ones (1) for presences. |
| <code>coord.cols</code> | names or index numbers of the columns containing the spatial coordinates in data (x and y, or longitude and latitude, in this order!). |
| <code>id.col</code> | optionally, the name or index number of a column (to be included in the output) containing locality identifiers in data. |
| <code>degree</code> | the degree of the spatial polynomial to use (see Details). The default is 3. |
| <code>step</code> | logical value indicating whether the regression of presence-absence on the spatial polynomial should do a stepwise inclusion of the polynomial terms (using the <code>step</code> function with default settings, namely backward AIC selection), rather than forcing all terms into the equation. The default is <code>TRUE</code> . |
| <code>Favourability</code> | logical value indicating whether the probability values obtained from the regression should be converted to favourability, so that they are more directly comparable among species with different prevalence (see Real et al. 2006, Acevedo & Real 2012). The default is <code>FALSE</code> . |
| <code>suffix</code> | character indicating the suffix to add to the trend surface columns in the resulting data frame. The default is <code>"_TS"</code> . |
| <code>save.models</code> | logical value indicating whether the models obtained from the regressions should be saved in the results. The default is <code>FALSE</code> . |

Details

Trend Surface Analysis is a way to model the spatial structure in species' distributions by regressing occurrence data on the spatial coordinates x and y, for a linear trend, or on polynomial terms of these coordinates (x^2 , y^2 , $x*y$, etc.), for curvilinear trends (Legendre & Legendre, 1998; Borcard et al., 2011). Second- and third-degree polynomials are often used. `multTSA` allows specifying the degree of the spatial polynomial to use. By default, it uses a 3rd-degree polynomial and performs stepwise AIC selection of the polynomial terms to include.

Value

This function returns a matrix or data frame containing the identifier column (if provided in `id.col`) and one column per species containing the value predicted by the trend surface analysis.

Author(s)

A. Marcia Barbosa

References

- Acevedo P. & Real R. (2012) Favourability: concept, distinctive characteristics and potential usefulness. *Naturwissenschaften* 99: 515-522
- Borcard D., Gillet F. & Legendre P. (2011) *Numerical Ecology with R*. Springer, New York.
- Legendre P. & Legendre L. (1998) *Numerical Ecology*. Elsevier, Amsterdam.
- Real R., Barbosa A.M. & Vargas J.M. (2006) Obtaining environmental favourability functions from logistic regression. *Environmental and Ecological Statistics* 13: 237-245

See Also

`distPres`, `poly`

Examples

```
data(rotif.env)

head(rotif.env)

names(rotif.env)

tsa <- multTSA(rotif.env, sp.cols = 18:20,
coord.cols = c("Longitude", "Latitude"), id.col = 1)

head(tsa)
```

`pairwiseRangemaps` *Pairwise intersection (and union) of range maps*

Description

This function takes a set of rangemaps and returns a matrix containing the areas of their pairwise intersections; optionally, also their individual areas and pairwise unions.

Usage

```
pairwiseRangemaps(rangemaps, projection, diag = TRUE, unions = TRUE,
verbosity = 2, Ncpu = 1, nchunks = 1, subchunks = NULL,
filename = "rangemap_matrix.csv")
```

Arguments

- | | |
|-------------------------|--|
| <code>rangemaps</code> | a character vector of rangemap filenames, including folder paths if not in the working directory. ESRI shapefile (.shp) is currently the only accepted format. Specifying the extension is optional. |
| <code>projection</code> | argument to be passed to function <code>importShapefile</code> of package PBSmapping |

| | |
|------------------------|--|
| <code>diag</code> | logical, whether to fill the diagonal of the resulting matrix with the area of each rangemap. The default is <code>TRUE</code> , and it is also automatically set to <code>TRUE</code> (as it is necessary) if <code>unions = TRUE</code> . |
| <code>unions</code> | logical, whether to fill the upper triangle of the resulting matrix with the area of union of each pair of rangemaps. The default is <code>TRUE</code> . It is not as computationally intensive as the intersection, as it is calculated not with spatial but with algebraic operations within the matrix ($\text{union} = \text{area1} + \text{area2} - \text{intersection}$). |
| <code>verbosity</code> | integer number indicating the amount of messages to display. |
| <code>Ncpu</code> | integer indicating the number of CPUs (central processing units) to employ if parallel computing is to be used. The default is 1 CPU, which implies no parallel computing, but you might want to increase this if you have many and/or large rangemaps and your machine has more cores that can be used simultaneously. You can find out the total number of cores in you machine with the <code>detectCores</code> function of package parallel ; a usually wise option is to use all cores except one (i.e., <code>Ncpu = parallel::detectCores() - 1</code>). |
| <code>nchunks</code> | either an integer indicating the number of chunks of rows in which to divide the results matrix for calculations, or character value "decreasing" indicating that the matrix should be divided into chunks of decreasing number of rows (as intersections are calculated in the lower triangle, rows further down the matrix have an increasing number of intersections to compute; note, however, that rangemap size, not rangemap number, is the main determinant of computation time). The default is 1 (no division of the matrix) but, if you have many rangemaps, the memory can get clogged. With chunks, each set of rows of the matrix is calculated and saved to disk, and the memory is cleaned before the next chunk begins. |
| <code>subchunks</code> | optional integer vector specifying which chunks to actually calculate. This is useful if a previous, time-consuming run of <code>pairwiseRangemaps</code> was interrupted (e.g. by a power outage) and you want to calculate only the remaining chunks. |
| <code>filename</code> | optional character vector indicating the name of the file to save the resulting matrix to. |

Details

This calculation can be very intensive and slow, especially if you have many and/or large rangemaps, due to the time needed for spatial operations between maps. You can set `nchunks = "decreasing"` for the matrix to be calculated in parts and the memory cleaned between one part and the next; and, if your computer has more than one core that you can use, you can increase `Ncpu` to get parallel computing.

Value

This function returns a square matrix containing, in the lower triangle, the area of the pair-wise intersections among the input `rangemaps`; in the diagonal (if `diag = TRUE` or `union = TRUE`), the area of each rangemap; and in the upper triangle (if `union = TRUE`), the area of the pair-wise unions among the rangemaps.

Note

This function uses the **PBSmapping** package to import and intersect the polygon rangemaps and to calculate areas. Remember to use projected rangemaps, preferably with an equal-area reference system, so that calculations are correct.

Author(s)

A. Marcia Barbosa

References

Barbosa A.M. & Estrada A. (in press) Calcular corotipos sin dividir el territorio en OGUs: una adaptacion de los indices de similitud para su utilizacion directa sobre areas de distribucion. In: Areas de distribucion: entre puentes y barreras. Universidad de Granada, Spain.

See Also

rangemapSim

| | |
|-----------------|--------------------------|
| percentTestData | <i>Percent test data</i> |
|-----------------|--------------------------|

Description

Based on the work of Schaafsma & van Vark (1979), Huberty (1994) provided a heuristic ("rule of thumb") for determining an adequate proportion of data to set aside for testing species presence/absence models, based on the number of predictor variables that are used (Fielding & Bell 1997). The `percentTestData` function calculates this proportion as a percentage.

Usage

```
percentTestData(nvar)
```

Arguments

| | |
|------|---------------------------------------|
| nvar | the number of variables in the model. |
|------|---------------------------------------|

Value

A numeric value of the percentage of data to leave out of the model for further model testing.

Author(s)

A. Marcia Barbosa

References

- Huberty C.J. (1994) Applied Discriminant Analysis. Wiley, New York, 466 pp.
- Schaafsma W. & van Vark G.N. (1979) Classification and discrimination problems with applications. Part IIa. Statistica Neerlandica 33: 91-126
- Fielding A.H. & Bell J.F. (1997) A review of methods for the assessment of prediction errors in conservation presence/absence models. Environmental Conservation 24: 38-49

See Also

multGLM

Examples

```
# say you're building a model with 15 variables:

percentTestData(15)

# the result tells you that 21% is an appropriate percentage of data
# to set aside for testing your model, so train it with 79% of the data
```

| | |
|-------------|--|
| rangemapSim | <i>Pairwise similarity between rangemaps</i> |
|-------------|--|

Description

Calculate pairwise similarity among rangemaps from a matrix of their areas of intersection and union

Usage

```
rangemapSim(rangemap.matrix, total.area,
method = c("Jaccard", "Sorensen", "Simpson", "Baroni"),
diag = FALSE, upper = FALSE, verbosity = 2)
```

Arguments

| | |
|-----------------|---|
| rangemap.matrix | a matrix like the one produced by function pairwiseRangemaps, containing the areas of pairwise intersection among rangemaps in the lower triangle, individual rangemap areas in the diagonal, and pairwise union areas in the upper diagonal. |
| total.area | numeric value indicating the total size of the study area, in the same units as the areas in the rangemap.matrix. Used only if the method uses shared absences (as is the case of "Baroni") |
| method | character value indicating the similarity index to use. Currently implemented indices are "Jaccard", "Sorensen", "Simpson" and "Baroni". The default is the first one. |

| | |
|------------------------|---|
| <code>diag</code> | logical value indicating if the diagonal of the resulting matrix should be filled |
| <code>upper</code> | logical value indicating if the upper triangle of the resulting matrix should be filled (symmetrical to the lower triangle) |
| <code>verbosity</code> | integer number indicating the amount of messages to display. |

Details

Distributional relationships among species are commonly determined based on pair-wise (dis)similarities in species' occurrence patterns. Some of the most commonly employed similarity indices are those of Jaccard (1901), Sorensen (1948), Simpson (1960) and Baroni-Urbani & Buser (1976), which are here implemented for comparing rangemaps based on their areas of intersection and union (Barbosa & Estrada, in press).

Value

This function returns a square matrix of pairwise similarities between the rangemaps in 'rangemap.matrix', calculated with the (first) similarity index specified in 'method'.

Author(s)

A. Marcia Barbosa

References

- Barbosa A.M. & Estrada A. (in press) Calcular corotipos sin dividir el territorio en OGUs: una adaptacion de los indices de similitud para su utilizacion directa sobre areas de distribucion. In: Areas de distribucion: entre puentes y barreras. Universidad de Granada, Spain.
- Baroni-Urbani C. & Buser M.W. (1976) Similarity of Binary Data. *Systematic Zoology*, 25: 251-259
- Jaccard P. (1901) Etude comparative de la distribution florale dans une portion des Alpes et des Jura. *Memoires de la Societe Vaudoise des Sciences Naturelles*, 37: 547-579
- Simpson G.G. (1960) Notes on the measurement of faunal resemblance. *Amer. J. Sci.* 258A, 300-311
- Sorensen T. (1948) A method of establishing groups of equal amplitude in plant sociology based on similarity of species and its application to analyses of the vegetation on Danish commons. *Kongelige Danske Videnskabernes Selskab*, 5(4): 1-34

See Also

`pairwiseRangemaps`; `simFromSetOps`; `simMat`

rotif.env

*Rotifers and environmental variables on TDWG level 4 regions of the world***Description**

These data were extracted from a database of monogonont rotifer species presence records on the geographical units used by the Biodiversity Information Standards (formerly Taxonomic Database Working Group, TDWG; base maps available from www.kew.org/science-conservation/research-data/resources/gis-unit/tdwg-world) and a few environmental (including human and spatial) variables on the same spatial units. The original data were compiled and published by Fontaneto et al. (2012) in long (narrow, stacked) format. Here they are presented in wide or unstacked format (presence-absence table, obtained with the `splist2presabs` function), reduced to the species recorded in at least 100 (roughly one third) different TDWG level 4 units, and with abbreviations of the species' names (obtained with the `spCodes` function). Mind that this is not a complete picture of these species' distributions, due to insufficient sampling in many regions.

Usage

```
data(rotif.env)
```

Format

A data frame with 291 observations on the following 47 variables.

TDWG4 a factor with 291 levels indicating the abbreviation code of each TDWG4 region

LEVEL_NAME a factor with 291 levels indicating the name of each TDWG4 region

REGION_NAME a factor with 47 levels indicating the name of the main geographical region to which each TDWG4 level belongs

CONTINENT a factor with 9 levels indicating the continent to which each TDWG4 level belongs

Area a numeric vector

Altitude a numeric vector

AltitudeRange a numeric vector

HabitatDiversity a numeric vector

HumanPopulation a numeric vector

Latitude a numeric vector

Longitude a numeric vector

Precipitation a numeric vector

PrecipitationSeasonality a numeric vector

TemperatureAnnualRange a numeric vector

Temperature a numeric vector

TemperatureSeasonality a numeric vector

UrbanArea a numeric vector
Abrigh a numeric vector
Afissa a numeric vector
Apriod a numeric vector
Bangul a numeric vector
Bcalyc a numeric vector
Bplica a numeric vector
Bquadr a numeric vector
Burceo a numeric vector
Cgibba a numeric vector
Edilat a numeric vector
Flongi a numeric vector
Kcochl a numeric vector
Kquadr a numeric vector
Ktropi a numeric vector
Lbulla a numeric vector
Lclost a numeric vector
Lhamat a numeric vector
Lluna a numeric vector
Llunar a numeric vector
Lovali a numeric vector
Lpatel a numeric vector
Lquadr a numeric vector
Mventr a numeric vector
Ppatul a numeric vector
Pquadr a numeric vector
Pvulga a numeric vector
Specti a numeric vector
Tpatin a numeric vector
Tsimil a numeric vector
Ttetra a numeric vector

Source

Fontaneto D., Barbosa A.M., Segers H. & Pautasso M. (2012) The 'rotiferologist' effect and other global correlates of species richness in monogonont rotifers. *Ecography*, 35: 174-182.

Examples

```
data(rotif.env)  
  
head(rotif.env)
```

`rotifers`*Rotifer species on TDWG level 4 regions of the world*

Description

These data were extracted from a database of monogonont rotifer species records on the geographical units used by the Biodiversity Information Standards (formerly Taxonomic Database Working Group, TDWG; base maps available at www.kew.org/science-conservation/research-data/resources/gis-unit/tdwg-world). The original data were compiled and published by Fontaneto et al. (2012) for all TDWG levels. Here they are reduced to the TDWG - level 4 units and to the species recorded in at least 100 (roughly one third) of these units. Mind that this is not a complete picture of these species' distributions, due to insufficient sampling in many regions.

Usage

```
data("rotifers")
```

Format

A data frame with 3865 observations on the following 2 variables.

`TDWG4` a factor with 274 levels corresponding to the code names of the TDWG level 4 regions in which the records were taken

`species` a factor with 30 levels corresponding to the names of the (sub)species recorded in at least 100 different TDWG level 4 regions

Source

Fontaneto D., Barbosa A.M., Segers H. & Pautasso M. (2012) The 'rotiferologist' effect and other global correlates of species richness in monogonont rotifers. *Ecography*, 35: 174-182.

Examples

```
data(rotifers)
```

```
head(rotifers, 10)
```

simFromSetOps*Calculate similarity from set operations*

Description

This function calculates pair-wise similarity based on the results of set operations (intersection, union) among the subjects.

Usage

```
simFromSetOps(size1, size2, intersection, union, total.size = NULL,  
method = c("Jaccard", "Sorensen", "Simpson", "Baroni"),  
verbosity = 1)
```

Arguments

| | |
|--------------|--|
| size1 | size of subject 1 (e.g., area of the distribution range of a species, or its number of presences within a grid). Not needed if <code>method = "Jaccard"</code> . |
| size2 | the same for subject 2. |
| intersection | size of the intersection among subjects 1 and 2 (area of the intersection among their distribution ranges, or number of grid cells in which they co-occur). |
| union | size of the union of subjects 1 and 2. |
| total.size | total size of the study area. Needed only when calculating a similarity index that takes shared absences into account (i.e., <code>method = "Baroni"</code>). |
| method | the similarity index to use. Currently implemented options are 'Jaccard', 'Sorensen', 'Simpson' and 'Baroni'. |
| verbosity | integer indicating whether to display messages. |

Details

Similarities among ecological communities, beta diversity patterns, biotic regions, and distributional relationships among species are commonly determined based on pair-wise (dis)similarities in species' occurrence patterns. This function implements some of the most commonly employed similarity indices, namely those of Jaccard (1901), Sorensen (1948), Simpson (1960) and Baroni-Urbani & Buser (1976), based on the amount of occupied and overlap area between two species.

Value

The numeric value of similarity among subjects 1 and 2.

Author(s)

A. Marcia Barbosa

References

- Baroni-Urbani C. & Buser M.W. (1976) Similarity of Binary Data. *Systematic Zoology*, 25: 251-259
- Jaccard P. (1901) Etude comparative de la distribution florale dans une portion des Alpes et des Jura. *Memoires de la Societe Vaudoise des Sciences Naturelles*, 37: 547-579
- Simpson, G.G. (1960) Notes on the measurement of faunal resemblance. *Amer. J. Sci.* 258A, 300-311
- Sorensen T. (1948) A method of establishing groups of equal amplitude in plant sociology based on similarity of species and its application to analyses of the vegetation on Danish commons. *Kongelige Danske Videnskabernes Selskab*, 5(4): 1-34

See Also

fuzSim, simMat

Examples

```
# take two species which occur in 22 and 35 area units, respectively
# and which overlap in 8 of those units:

sp1 <- 22
sp2 <- 35
int <- 8
uni <- sp1 + sp2 - int

# calculate similarity between their distributions based on
# different indices:

simFromSetOps(intersection = int, union = uni, method = "Jaccard")

simFromSetOps(sp1, sp2, int, uni, method = "Sorensen")

simFromSetOps(sp1, sp2, int, uni, method = "Simpson")

# if you want Baroni-Urbani & Buser's index
# you need to provide also the total size of your study area:

simFromSetOps(sp1, sp2, int, uni, total = 100, method = "Baroni")
```

simMat

Pair-wise (fuzzy) similarity matrix

Description

`simMat` takes a matrix or data frame containing species occurrence data or regional species composition, either categorical (0 or 1) or fuzzy (between 0 and 1), and uses the `fuzSim` function to calculate a square matrix of pair-wise similarities between them, using a fuzzy logic version (Barbosa, 2015) of the specified similarity index.

Usage

```
simMat(data, method, diag = TRUE, upper = TRUE)
```

Arguments

| | |
|---------------------|---|
| <code>data</code> | a matrix or data frame containing (optionally fuzzy) species presence-absence data (in wide format, i.e. one column per species), with 1 meaning presence, 0 meaning absence, and values in between for fuzzy presence (or the degree to which each locality belongs to the set of species presences; see Zadeh, 1965). Fuzzy presence-absence can be obtained, for example, with <code>multGLM</code> , <code>distPres</code> or <code>multTSA</code> . These data can also be transposed for comparing regional species compositions. |
| <code>method</code> | the similarity index whose fuzzy version to use. See <code>fuzSim</code> for available options. |
| <code>diag</code> | logical value indicating whether the diagonal of the matrix should be filled (with ones). Defaults to <code>TRUE</code> . |
| <code>upper</code> | logical value indicating whether the upper triangle of the matrix (symmetric to the lower triangle) should be filled. Defaults to <code>TRUE</code> . |

Details

The fuzzy versions of species occurrence data and of binary similarity indices introduce tolerance for small spatial differences in species' occurrence localities, allow for uncertainty about species occurrence, and may compensate for under-sampling and geo-referencing errors (Barbosa, 2015).

Value

This function returns a square matrix of pair-wise similarities among the species distributions (columns) in `data`. Similarity is calculated with the fuzzy version of the index specified in `method`, which yields traditional binary similarity if the data are binary (0 or 1), or fuzzy similarity if the data are fuzzy (between 0 and 1) (Barbosa, 2015).

Author(s)

A. Marcia Barbosa

References

Barbosa A.M. (2015) fuzzySim: applying fuzzy logic to binary similarity indices in ecology. *Methods in Ecology and Evolution*, 6: 853-858.

See Also

fuzSim

Examples

```
# load and look at the rotif.env presence-absence data:

data(rotif.env)

head(rotif.env)

names(rotif.env)

# build a matrix of similarity among these binary data
# using e.g. Jaccard's index:

bin.sim.mat <- simMat(rotif.env[ , 18:47], method = "Jaccard")

head(bin.sim.mat)

# calculate a fuzzy version of the presence-absence data
# based on inverse distance to presences:

rotifers.invd <- distPres(rotif.env, sp.cols = 18:47,
  coord.cols = c("Longitude", "Latitude"), id.col = 1, suffix = ".d",
  p = 1, inv = TRUE)

head(rotifers.invd)

# build a matrix of fuzzy similarity among these fuzzy
# distribution data, using the fuzzy version of Jaccard's index:

fuz.sim.mat <- simMat(rotifers.invd[ , -1], method = "Jaccard")

head(fuz.sim.mat)

# plot the similarity matrices as colours:

image(x = 1:ncol(bin.sim.mat), y = 1:nrow(bin.sim.mat),
  z = bin.sim.mat, col = rev(heat.colors(256)), xlab = "", ylab = "",
  axes = FALSE, main = "Binary similarity")
axis(side = 1, at = 1:ncol(bin.sim.mat), tick = FALSE,
  labels = colnames(bin.sim.mat), las = 2)
axis(side = 2, at = 1:nrow(bin.sim.mat), tick = FALSE,
  labels = rownames(bin.sim.mat), las = 2)

image(x = 1:ncol(fuz.sim.mat), y = 1:nrow(fuz.sim.mat),
```

```

z = fuz.sim.mat, col = rev(heat.colors(256)), xlab = "", ylab = "",
axes = FALSE, main = "Fuzzy similarity")
axis(side = 1, at = 1:ncol(fuz.sim.mat), tick = FALSE,
labels = colnames(fuz.sim.mat), las = 2, cex = 0.5)
axis(side = 2, at = 1:nrow(fuz.sim.mat), tick = FALSE,
labels = rownames(fuz.sim.mat), las = 2)

# plot a UPGMA dendrogram from each similarity matrix:

plot(hclust(as.dist(1 - bin.sim.mat), method = "average"),
main = "Binary cluster dendrogram")

plot(hclust(as.dist(1 - fuz.sim.mat), method = "average"),
main = "Fuzzy cluster dendrogram")

# you can get fuzzy chorotypes from these similarity matrices
# (or fuzzy biotic regions if you \code{\link{transpose}} 'data'),
# so that localities are in columns and species in rows)
# using the \pkg{RMACOQUI} package (Olivero et al. 2011)

```

spCodes

Obtain unique abbreviations of species names

Description

This function takes a vector of species names and converts them to abbreviated species codes containing the specified numbers of characters from the genus, the specific and optionally also the subspecific name. Separators can be specified by the user. The function checks that the resulting codes are unique.

Usage

```

spCodes(species, nchar.gen = 3, nchar.sp = 3, nchar.ssp = 0,
sep.species = " ", sep.spcode = "")

```

Arguments

| | |
|------------------------|--|
| <code>species</code> | a character vector containig the species names to be abbreviated. |
| <code>nchar.gen</code> | the number of characters from the genus name to be included in the resulting species code. |
| <code>nchar.sp</code> | the number of characters from the specific name to be included in the resulting species code. |
| <code>nchar.ssp</code> | optionally, the number of characters from the subspecific name to be included in the resulting species code. Set it to 0 if you have subspecific names in <code>species</code> but do not want them included in the resulting species codes. |

- `sep.species` the character that separates genus, specific and subspecific names in `species`. The default is a white space.
- `sep.spcode` the character you want separating genus and species abbreviations in the resulting species codes. The default is an empty character (no separator).

Value

This function returns a character vector containing the species codes resulting from the abbreviation. If the numbers of characters specified do not make for unique codes, an error message is displayed showing which `species` names caused it, so that you can try again with different `nchar.gen`, `nchar.sp` and/or `nchar.ssp`.

Author(s)

A. Marcia Barbosa

See Also

`substr`, `strsplit`

Examples

```
data(rotifers)

head(rotifers)

## add a column to 'rotifers' with shorter versions of the species names:

## Not run:
rotifers$spcode <- spCodes(rotifers$species, sep.species = "_",
nchar.gen = 1, nchar.sp = 4, nchar.ssp = 0, sep.spcode = ".")

# this produces an error due to resulting species codes not being unique

## End(Not run)

rotifers$spcode <- spCodes(rotifers$species, sep.species = "_",
nchar.gen = 1, nchar.sp = 5, nchar.ssp = 0, sep.spcode = ".")

# with a larger number of characters from the specific name,
# resulting codes are now unique

## check out the result:
head(rotifers)
```


splist2presabs

*Convert a species list to a presence-absence table***Description**

This function takes a locality+species dataset in long (stacked) format, i.e., a matrix or data frame containing localities in one column and their recorded species in another column, and converts them to a presence-absence table (wide format) suitable for mapping and for computing distributional similarities (see e.g. `simMat`). Try out the Examples below for an illustration).

Usage

```
splist2presabs(data, sites.col, sp.col, keep.n = FALSE)
```

Arguments

| | |
|------------------------|--|
| <code>data</code> | a matrix or data frame with localities in one column and species in another column. Type <code>data(rotifers)</code> ; <code>head(rotifers)</code> for an example. |
| <code>sites.col</code> | the name or index number of the column containing the localities in <code>data</code> . |
| <code>sp.col</code> | the name or index number of the column containing the species names or codes in <code>data</code> . |
| <code>keep.n</code> | logical value indicating whether to get in the resulting table the number of times each species appears in each locality; if <code>FALSE</code> (the default), only presence (1) or absence (0) is recorded. |

Value

A data frame containing the localities in the first column and then one column per species indicating their presence (or their number of records if `keep.n = TRUE`) and absence. Type `data(rotif.env)`; `head(rotif.env[,18:47])` for an example.

Author(s)

A. Marcia Barbosa

See Also

`table`

Examples

```
data(rotifers)

head(rotifers)

rotifers.presabs <- splist2presabs(rotifers, sites.col = "TDWG4",
sp.col = "species", keep.n = FALSE)

head(rotifers.presabs)
```

Description

This function builds a generalized linear model with forward stepwise inclusion of variables, using AIC as the selection criterion, and provides the values predicted at each step, as well as their correlation with the final model predictions.

Usage

```
stepByStep(data, sp.col, var.cols, family = binomial(link = "logit"),
  Favourability = FALSE, trace = 0, cor.method = "pearson")
```

Arguments

| | |
|----------------------------|---|
| <code>data</code> | a data frame containing your target and predictor variables. |
| <code>sp.col</code> | index number of the column of <code>data</code> that contains the target variable. |
| <code>var.cols</code> | index numbers of the columns of <code>data</code> that contain the predictor variables. |
| <code>family</code> | argument to be passed to the <code>glm</code> function indicating the family (and error distribution) to use in modelling. The default is binomial distribution with logit link (for binary target variables). |
| <code>Favourability</code> | logical, whether to apply the <code>Favourability</code> function to remove the effect of prevalence from predicted probability (Real et al. 2006). Applicable only to binomial GLMs. Defaults to <code>FALSE</code> . |
| <code>trace</code> | argument to pass to the <code>step</code> function. If positive, information is printed during the stepwise procedure. Larger values may give more detailed information. The default is 0 (silent). |
| <code>cor.method</code> | character string to pass to function <code>cor</code> indicating which coefficient should be used for correlating predictions at each step with those of the final model. Can be "pearson" (the default), "kendall", or "spearman". |

Details

Stepwise variable selection often includes more variables than would a model selected after examining all possible combinations of the variables (e.g. with packages **MuMIn** and **glmulti**). The `stepByStep` function can be useful to assess if a stepwise model with just the first few variables could already provide predictions very close to the final ones (see e.g. Fig. 3 in Munoz et al., 2005). It can also be useful to see which variables determine the more general trends in the model predictions, and which just add (local) additional nuances.

Value

This function returns a list of the following components:

| | |
|---------------------------|---|
| <code>predictions</code> | a data frame with the model's fitted values at each step of the variable selection. |
| <code>correlations</code> | a numeric vector of the correlation between the predictions at each step and those of the final model. |
| <code>variables</code> | a character vector of the variables in the final model, named with the step at which each was included. |
| <code>model</code> | the resulting model object. |

Author(s)

A. Marcia Barbosa

References

Munoz, A.R., Real R., BARBOSA A.M. & Vargas J.M. (2005) Modelling the distribution of Bonelli's Eagle in Spain: Implications for conservation planning. *Diversity and Distributions* 11: 477-486

Real R., Barbosa A.M. & Vargas J.M. (2006) Obtaining environmental favourability functions from logistic regression. *Environmental and Ecological Statistics* 13: 237-245.

See Also

`step`, `glm`, `modelTrim`

Examples

```
data(rotif.env)

stepByStep(data = rotif.env, sp.col = 18, var.cols = 5:17,
cor.method = "spearman")

stepByStep(data = rotif.env, sp.col = 18, var.cols = 5:17,
cor.method = "spearman", Favourability = TRUE)

stepByStep(data = rotif.env, sp.col = 9, var.cols = c(5:8, 10:17),
family = poisson)
```

timer

Timer

Description

Reporting of time elapsed since a given start time. This function is used internally by other functions in the package.

Usage

```
timer(start.time)
```

Arguments

`start.time` A date-time object of class `POSIXct`, e.g. as given by `Sys.time`.

Value

The function returns a message informing of the time elapsed since the input `start.time`.

Author(s)

A. Marcia Barbosa

See Also

`Sys.time`, `proc.time`, `difftime`

Examples

```
# get starting time:
start <- Sys.time()

# do some random analysis:
sapply(rnorm(50000), function(x) x*5)

# see how long it took:
timer(start)
```

transpose

Transpose (part of) a matrix or dataframe

Description

This function transposes (a specified part of) a matrix or data frame, optionally using one of its columns as column names for the transposed result. It can be useful for turning a species presence-absence table into a regional species composition table.

Usage

```
transpose(data, sp.cols = 1:ncol(data), reg.names = NULL)
```

Arguments

| | |
|------------------------|---|
| <code>data</code> | a matrix or data frame containing the species occurrence data to transpose. |
| <code>sp.cols</code> | names or index numbers of the columns containing the species occurrences in <code>data</code> which are meant to be transposed. |
| <code>reg.names</code> | name or index number of the column in <code>data</code> containing the region names, to be used as column names in the transposed result. |

Value

This function returns the transposed `sp.cols` of `data`, with the column specified in `reg.names` as column names.

Author(s)

A. Marcia Barbosa

See Also

`t`

Examples

```
data(rotif.env)

head(rotif.env)

names(rotif.env)

rotif.reg <- transpose(rotif.env, sp.cols = 18:47, reg.names = 1)

head(rotif.reg)
```

| | |
|-----------|----------------------------------|
| triMatInd | <i>Triangular matrix indices</i> |
|-----------|----------------------------------|

Description

This function outputs the indices of one triangle (the lower one by default) of an input square matrix. It is used by `simMat` and, for large matrices, makes it faster than e.g. with `lower.tri` or `upper.tri`.

Usage

```
triMatInd(mat, lower = TRUE, list = FALSE)
```

Arguments

| | |
|--------------------|---|
| <code>mat</code> | a square matrix. |
| <code>lower</code> | logical indicating whether the indices should correspond to the lower triangle. The default is TRUE; FALSE produces the upper triangle indices. |
| <code>list</code> | logical indicating whether the results should be output as a list instead of a matrix. The default is FALSE. |

Value

The indices (row, column) of the elements of the matrix that belong to the requested triangle.

Author(s)

A. Marcia Barbosa

References

<http://stackoverflow.com/questions/20898684/how-to-efficiently-generate-lower-triangle-indices-of-a-symmetric-matrix>

See Also

`lower.tri`, `upper.tri`

Examples

```
mat <- matrix(nrow = 4, ncol = 4)
mat
triMatInd(mat)
triMatInd(mat, list = TRUE)
```