

# Package ‘specificity’

January 8, 2020

**Title** Calculate Environmental or Host Phylogenetic Specificity

**Version** 0.0.0.9000

**Description** The purpose of this package is to calculate phylogenetic and environmental specificity of species. I wrote this software to analyze specificity of microbes to hosts or to environment, but there is no reason that this software wouldn't work with macroorganisms as well.

**License** GPL

**Encoding** UTF-8

**LazyData** true

**RoxygenNote** 6.1.1

**Depends** ape, geiger, parallel, MASS

**Imports** Rcpp, fields, testthat

**LinkingTo** Rcpp

**NeedsCompilation** yes

**Author** John L Darcy [aut, cre]

**Maintainer** John L Darcy <darcyj@colorado.edu>

## R topics documented:

bl_distance_ns . . . . .	2
check_pes_inputs . . . . .	3
daughter_seeds . . . . .	3
distcalc . . . . .	4
env_spec_sim . . . . .	5
geo_spec_sim . . . . .	6
make_nested_set . . . . .	8
occ_threshold . . . . .	9
pairwise_geo_mean . . . . .	10
pairwise_product . . . . .	11
phy_or_env_spec . . . . .	11
phy_spec_sim . . . . .	14
plot_grid_abunds . . . . .	16
plot_specificities . . . . .	17

prop_abund . . . . .	18
pval_from_perms . . . . .	19
randomgrid . . . . .	20
random_rep_positions . . . . .	21
rao_quad_ent . . . . .	21
tree2mat . . . . .	22
wpd . . . . .	23
wpd_table . . . . .	25

<b>Index</b>	<b>27</b>
--------------	-----------

---

<i>bl_distance_ns</i>	<i>bl_distance_ns</i>
-----------------------	-----------------------

---

**Description**

Calculates branch-length distance between tipa and tipb in a phylogenetic tree using nested-set optomization. Requires a pre-calculated nested-set.

**Usage**

```
bl_distance_ns(tipa, tipb, tree, ns)
```

**Arguments**

tipa	string. Name of a tip in tree.
tipb	string. Name of another tip in tree.
tree	phylo object. Tree containing all unique species in x as tips. May contain tips that are not in x.
ns	matrix. Nested-set matrix for tree; use make_nested_set(tree).

**Value**

Distance between tipa and tipb.

**Author(s)**

John L. Darcy

**Examples**

```
library(ape)
example_tree <- ape::read.tree(text=" (((a:1,b:1):1,c:2):1,d:3):1,(e:1,f:1):3);")
plot(example_tree); axis(side=1)
example_ns <- make_nested_set(example_tree)
bl_distance_ns("a", "c", example_tree, example_ns) # should be 4
bl_distance_ns("a", "f", example_tree, example_ns) # should be 8
bl_distance_ns("d", "c", example_tree, example_ns) # should be 6
```

---

check_pes_inputs	<i>check_pes_inputs</i>
------------------	-------------------------

---

### Description

Function used by `phy_or_env_spec`. checks `abunds_mat`, `env`, `hosts`, and `hosts_phylo` inputs to `phy_or_env_spec` to make sure there are no problems. This could include missing species in trees, incompatible dimensions, non-numeric inputs, etc. Returns an input type, which is just a string that can be "mat", "dist", "vec", "phy", or "error".

### Usage

```
check_pes_inputs(abunds_mat, env, hosts, hosts_phylo, verbose = TRUE)
```

### Arguments

<code>abunds_mat</code>	(required, see <code>phy_or_env_spec</code> )
<code>env</code>	(required, can be NULL, see <code>phy_or_env_spec</code> )
<code>hosts</code>	(required, can be NULL, see <code>phy_or_env_spec</code> )
<code>hosts_phylo</code>	(required, can be NULL, see <code>phy_or_env_spec</code> )
<code>verbose</code>	logical. Should status messages be displayed? (DEFAULT: TRUE).

### Value

string. either "mat", "dist", "vec", "phy", or "error".

---

daughter_seeds	<i>daughter_seeds</i>
----------------	-----------------------

---

### Description

Makes `n` daughter seeds from seed `s`. This is useful for processes one wishes to be deterministic, but may not be executed in the same order every time.

### Usage

```
daughter_seeds(n, s = 12345)
```

### Arguments

<code>n</code>	integer. Number of daughter seeds to make.
<code>s</code>	integer. A seed (DEFAULT: 12345).

**Value**

vector of length n containing integer seeds.

**Author(s)**

John L. Darcy

---

distcalc

*geo\_distcalc*

---

**Description**

Calculates pairwise geographic distance between locations on earth. Just a convenient wrapper for `fields::rdist.earth`.

**Usage**

```
distcalc(lat, lng, sampIDs = NULL)
```

**Arguments**

lat	Numeric vector. Latitudes in decimal degree format.
lng	Numeric vector. Longitudes in decimal degree format.
sampIDs	Character vector. Sample identifiers. Only required if output dist should have names associated.

**Value**

matrix containing all pairwise geographic distances in km.

**Author(s)**

John L. Darcy

**Examples**

```
data(endophyte)
geo_dists <- distcalc(metadata$Lat, metadata$Lon, metadata$SampleID)
all(rownames(geo_dists) == metadata$SampleID)
```

env\_spec\_sim

*env\_spec\_sim***Description**

Simulates inputs for `phy_or_env_spec`, by creating a species distribution over an artificial (or real) environmental variable. That distribution has a mean at the "ideal" environmental value for the simulated species, and the standard deviation of that distribution controls the extent to which the species is specific to the variable. A high SD means less specificity, and a low SD means more specificity.

**Usage**

```
env_spec_sim(sdev, ideal, ideal2 = 0, ideal3 = 0, n_ideal = 1, env,
             n_obs, up = 0, oceanp = 0, n_cores = 2, seed = 1234567)
```

**Arguments**

<code>sdev</code>	numeric vector. Standard deviation of the probability distribution $P(\text{species})$ , in the same units of <code>env</code> . Low values mean that the species is found across only a narrow range of <code>env</code> , i.e. specificity. High values mean that the species is found across a wide range of <code>env</code> , i.e. cosmopolitanism. Multiple values can be input in order to simulate a range of specificities simultaneously. Can be length 1 or <code>n</code> .
<code>ideal</code>	numeric vector. Value of <code>env</code> that is ideal for the simulated species. This is the mode of the probability distribution $P(\text{species})$ . Can be length 1 or <code>n</code> .
<code>ideal2</code>	numeric vector. Value of <code>env</code> that is the second ideal for the simulated species. Only used if <code>n_ideal</code> $\geq 2$ . This is the second mode of the probability distribution $P(\text{species})$ . Can be length 1 or <code>n</code> .
<code>ideal3</code>	numeric vector. Value of <code>env</code> that is the third ideal for the simulated species. Only used if <code>n_ideal</code> = 3. This is the third mode of the probability distribution $P(\text{species})$ . Can be length 1 or <code>n</code> .
<code>env</code>	numeric vector. Real or fake environmental variable.
<code>n_obs</code>	integer vector. Number of positive observations to make, i.e. occupancy of simulated species. Can be length 1 or <code>n</code> (default: 1).
<code>up</code>	numeric vector. <code>up</code> =uniform proportion. This is the proportion of the probability distribution $P(\text{species})$ that is composed of a uniform distribution, if desired. If set to a value above zero (and below 1), $P(\text{species})$ will be a weighted sum of the normal distribution described above, and a uniform distribution. The weight for the uniform distribution will be <code>up</code> , and the weight for the normal distribution will be <code>1-up</code> (default: 0).
<code>oceanp</code>	numeric vector. <code>oceanp</code> =ocean proportion. This is the proportion of samples in <code>env</code> that are "in the ocean", i.e. samples where the species would not expect to be found even if <code>env</code> is permissive. If aliens were calculating specificity of cows to temperature, they might look in the ocean at sites where the temperature is

17C (great for cows). But cows are not found in the ocean. This proportion is used to randomly select ocean sites within env, and then  $p(\text{slenv}|\text{ocean}) = \text{up}$ . Can be length 1 or n (default: 0).

n\_cores integer. Number of CPU cores for parallel computation (DEFAULT: 2).

seed integer. Seed for randomization. Daughter seeds will be generated for parallel computations, each with the same number of digits as seed (DEFAULT: 1234567).

Details

Since this process can result in failures (if a species is requested that's highly specific to a region of env that isn't samples), some output species will be failures. Default operation is to remove those failures from output matrix and output params data frame, but this can be changed.

Value

List object containing "matrix" and "params" objects:

**matrix:** matrix where each column is a vector of simulated observation frequencies (counts) corresponding to a value of env; each row represents a simulated species.

**params:** data.frame of parameters (columns) used to simulate each species (rows).

Author(s)

John L. Darcy

Examples

none yet written.

---

geo_spec_sim	<i>geo_spec_sim</i>
--------------	---------------------

---

Description

Simulates inputs for phy\_or\_env\_spec, by creating a species distribution over artificial (or real) geographic space. That distribution has a bivariate mean at the "ideal" location inspace for the simulated species, and the standard deviation of that (normal) distribution controls the extent to which the species specific to geographic space. A high SD means less specificity, and a low SD means more specificity.

Usage

```
geo_spec_sim(sdev, n_obs, grid, ideal_x = 0, ideal_y = 0,
  ideal_x2 = 0, ideal_y2 = 0, ideal_x3 = 0, ideal_y3 = 0,
  n_ideal = 1, up = 0, seed = 123456, n_cores = 2)
```

**Arguments**

sdev	numeric vector. Standard deviation of the probability distribution $P(\text{species})$ , in the same units as grid. $P(\text{species})$ is a function of the distance between a sample site and its closest ideal location (specified with <code>ideal_x/2/3</code> and <code>ideal_y/2/3</code> ). Low values mean that the species is found in abundance within only short distances of ideal locations, high values mean the species is found across a wider area. Multiple values can be input in order to simulate a range of specificities simultaneously. Can be length 1 or n.
n_obs	integer vector. Number of observations to make, i.e. number of times species is observed. Will be the sum of the species' output column. Can be length 1 or n.
grid	data frame with columns x and y, representing cartesian coordinates of sample locations. Can be artificial (generate with <code>randomgrid()</code> ) or real.
ideal_x	numeric vector. x-coordinate of the ideal spatial location for species (DEFAULT=0).
ideal_y	numeric vector. y-coordinate of the ideal spatial location for species (DEFAULT=0).
ideal_x2	numeric vector. x-coordinate for secondary ideal location. Only used if <code>n_ideal&lt;1</code> (DEFAULT=0).
ideal_y2	numeric vector. y-coordinate for secondary ideal location. Only used if <code>n_ideal&lt;1</code> (DEFAULT=0).
ideal_x3	numeric vector. x-coordinate for secondary ideal location. Only used if <code>n_ideal&lt;2</code> (DEFAULT=0).
n_ideal	integer vector. number of ideal locations to use. Must be 1, 2, or 3 (DEFAULT=1).
up	numeric vector. <code>up=uniform</code> proportion. This is the proportion of the probability distribution $P(\text{species})$ that is composed of a uniform distribution, if desired. If set to a value above zero (and below 1), $P(\text{species})$ will be a weighted sum of the normal distribution described above, and a uniform distribution. The weight for the uniform distribution will be <code>up</code> , and the weight for the normal distribution will be <code>1-up</code> (default: 0).
seed	integer. Seed for randomization. Daughter seeds will be generated for parallel computations, each with the same number of digits as seed (DEFAULT: 1234567).
n_cores	integer. Number of CPU cores for parallel computation (DEFAULT: 2).
ideal_x3	numeric vector. x-coordinate for secondary ideal location. Only used if <code>n_ideal&lt;2</code> (DEFAULT=0).

**Value**

List object containing "matrix" and "params" objects:

**matrix** matrix where each column is a vector of simulated observations for each row in grid; each column of matrix represents a simulated species.

**params** data.frame of parameters (columns) used to simulate each species (rows).

**Author(s)**

John L. Darcy

## Examples

```
g1 <- randomgrid()
plot(g1)
a1 <- geo_spec_sim(sdev=c(30, 30, 30, 30), n_obs=1000, grid=g1, up=c(0, 0.20, 0.40, 0.60))
par(mfrow=c(2,2))
plot_grid_abunds(g1, a1$matrix[,1])
plot_grid_abunds(g1, a1$matrix[,2])
plot_grid_abunds(g1, a1$matrix[,3])
plot_grid_abunds(g1, a1$matrix[,4])
a2 <- geo_spec_sim(sdev=c(10, 20, 30, 40), n_obs=1000, grid=g1, ideal_x=-50, ideal_x2=50, n_ideal=2)
par(mfrow=c(2,2))
plot_grid_abunds(g1, a2$matrix[,1], main="sd=10")
plot_grid_abunds(g1, a2$matrix[,2], main="sd=20")
plot_grid_abunds(g1, a2$matrix[,3], main="sd=30")
plot_grid_abunds(g1, a2$matrix[,4], main="sd=40")
```

---

make\_nested\_set

*make\_nested\_set*

---

## Description

Makes a nested set table for a phylo object. Phylo objects made by the ape package store phylogenies as an "adjacency list", which in R is a table within which any given edge is represented by the two node numbers it connects. With this data structure, it is very computationally expensive to figure out which tips are the descendents of a given node. Instead, using a "nested set" data structure, this operation is trivial. A nested set stores the minimum and maximum tip index for each node, such that the descendents of that node are given by the inclusive range between those values.

## Usage

```
make_nested_set(phy, n_cores = 2)
```

## Arguments

phy	phylo object. Must be rooted, and sorted such that tip indices are ordered. This is the default for rooted trees read in using ape's read.tree function.
n_cores	integer. Number of CPU cores to use (DEFAULT: 2). lapply will be used instead of mclapply if ncores is 1.

## Value

Matrix object representing a nested set of nodes. Each row matches rows of the "edges" object within phy. Object has the following columns:

- 1 (node)** Node value in the original phylo object.
- 2 (min)** minimum tip index subtended by node.
- 3 (max)** maximum tip index subtended by node.
- 4 (contig)** Is min:max contiguous? 1 (true) or 0 (false).



**Author(s)**

John L. Darcy

**References**

[https://en.wikipedia.org/wiki/Nested\\_set\\_model](https://en.wikipedia.org/wiki/Nested_set_model) [https://en.wikipedia.org/wiki/Adjacency\\_list](https://en.wikipedia.org/wiki/Adjacency_list)

**See Also**

ape::phylo geiger::tips

**Examples**

```
library(geiger)
library(ape)
library(parallel)
phy <- get(data(geospiza))$phy
# check if tree is rooted:
is.rooted(phy)
# make nested set table:
phy_ns <- make_nested_set(phy)
# show that nested set table matches up with edges table in phy:
all(phy$edge[,2] == phy_ns[,1])
```

---

occ\_threshold

*occ\_threshold*


---

**Description**

removes species (columns) from a matrix that don't meet a minimum occupancy, defined as the number of samples in which that species was observed.

**Usage**

```
occ_threshold(m, threshold, max_absent = 0)
```

**Arguments**

m	matrix or data frame of numeric values. Columns represent species, rows are samples.
threshold	integer. Minimum number of samples a species can occupy without being removed.
max_absent	float. Maximum abundance value at which a species will be considered absent (DEFAULT: 0).

**Value**

matrix with low-occupancy species removed.

**Author(s)**

John L. Darcy

**Examples**

```
attach(endophyte)
dim(zotutable)
zotutable_over25 <- occ_threshold(zotutable, 25)
dim(zotutable_over25)
```

---

pairwise_geo_mean	<i>pairwise_geo_mean</i>
-------------------	--------------------------

---

**Description**

Calculates pairwise geometric means from unique 2-element combinations of vector x. Written in C++ because R slow. The output vector is the same length and same order as a lower triangle of matrix with rows and columns x.

**Usage**

```
pairwise_geo_mean(x)
```

**Arguments**

x	numeric vector.
---	-----------------

**Value**

vector of pairwise geometric means, of length  $(l * l - 1) / 2$ , where  $l = \text{length}(x)$ .

**Author(s)**

John L. Darcy

**Examples**

```
x <- 1:6
y_cpp <- pairwise_geo_mean(x)
y_r <- as.dist(outer(x, x, function(x,y){sqrt(x*y)}))
print("Calculated with R's outer() function:")
y_r
print("As a vector:")
as.vector(y_r)
print("Calculated with pairwise_geo_mean (C++):")
y_cpp
```

---

pairwise_product	<i>pairwise_product</i>
------------------	-------------------------

---

**Description**

Calculates pairwise\_products from unique 2-element combinations of vector x. The output vector is the same length and same order as a lower triangle of matrix with rows and columns x.

**Usage**

```
pairwise_product(x)
```

**Arguments**

x                      numeric vector.

**Value**

vector of pairwise\_products, of length  $(l^2-l)/2$ , where  $l=\text{length}(x)$ .

**Author(s)**

John L. Darcy

**Examples**

```
x <- 1:6
y_cpp <- pairwise_geo_mean(x)
y_r <- as.dist(outer(x, x, function(x,y){x*y}))
print("Calculated with R's outer() function:")
y_r
print("As a vector:")
as.vector(y_r)
print("Calculated with pairwise_product (C++):")
y_cpp
```

---

phy_or_env_spec	<i>phy_or_env_spec</i>
-----------------	------------------------

---

**Description**

Calculates species' specificities to either a 1-dimensional variable (vector), 2-dimensional variable (matrix), or to a phylogeny. Transforms all variable input types into a matrix D, and calculates specificity by comparing empirical RQE\* (weighted mean of D's lower triangle and unique pairwise products of species abundances) to simulated RQE\* (same but with permuted abundances). This "raw" specificity is then divided by some denominator d in order to standardize it such that specificities can be compared between different species and different variables. Values closer to 0 indicate random assortment (null hypothesis), and more negative values indicate stronger specificity.

**Usage**

```
phy_or_env_spec(abunds_mat, env = NULL, hosts = NULL,
  hosts_phylo = NULL, n_sim = 1000, sim_fun = function(m) {
  m[sample(1:nrow(m)), ] }, p_adj = "fdr", seed = 1234567, tails = 1,
  n_cores = 2, verbose = TRUE, lowmem = FALSE, p_method = "raw",
  denom_type = "flat", diagnostic = F)
```

**Arguments**

abunds_mat	matrix or data frame of numeric values. Columns represent species, rows are samples. For columns where the value is nonzero for two or fewer data points, environmental SES cannot be calculated, and NAs will be returned. Negative values in abunds_mat are not allowed (REQUIRED).
env	numeric vector, dist, or square matrix. Environmental variable corresponding to abunds. For example, temperature, or geographic distance. Not required for computing phylogenetic specificity (DEFAULT: NULL).
hosts	character vector. Host identities corresponding to abunds. Only required if calculating SES for phylogenetic specificity (DEFAULT: NULL).
hosts_phylo	phylo object. Tree containing all unique hosts as tips. Only required if calculating SES for phylogenetic specificity (DEFAULT: NULL).
n_sim	integer. Number of simulations of abunds_mat to do under the null hypothesis that host or environmental association is random. P-values will not be calculated if n_sim < 100 (DEFAULT: 500).
sim_fun	function. A function f where f(abunds_mat) returns a matrix object with the same number of rows and columns as abunds_mat. Default is f=function(m) m[sample(1:nrow(m)),], which just permutes the order of rows in abunds_mat. Users may wish to use a null model that is able to preserve row and column totals such as the function permatswap() from the vegan package or the function vaznull() from the bipartite package. Either of these can be easily adapted to return only a single matrix (see examples). However, neither can accomodate non-integer matrices.
p_adj	string. Type of multiple hypothesis testing correction performed on P-values. Can take any valid method argument to p.adjust, including "none", "bonferroni", "holm", "fdr", and others (DEFAULT: "fdr").
seed	integer. Seed to use so that this is repeatable (DEFAULT: 1234557).
tails	integer. 1 = 1-tailed, test for specificity only. 2 = 2-tailed. 3 = 1-tailed, test for cosmopolitanism only. 0 = no test, P=1.0 (DEFAULT: 1).
n_cores	integer. Number of CPU cores to use for parallel operations. If set to 1, lapply will be used instead of mclapply (DEFAULT: 2).
verbose	logical. Should status messages be displayed? (DEFAULT: TRUE).
p_method	string. method argument to pval_from_perms (DEFAULT: "raw").
denom_type	string. Type of denominator (d) to use (DEFAULT: "flat"). Note that denominator type does NOT affect P-values.

- species\_sim:** d for species s is calculated as the standard deviation of specificities calculated from permuted abundances of s. This makes the output specificity a standardized effect size (SES). Unfortunately, this makes SES of specificity counterintuitively sensitive to occupancy, where species with high occupancy have more extreme SES of specificity than rare species, due to their more deterministic sim specificities. Not suggested.
- global\_unif:** d is same for all species. Calculated as variability in RQE\* under random uniform distribution (beta 1,1) of species abunds. This d is comparable between different abundance matrices and between different variables. Specificity is an SES using this denom\_type, but is NOT comparable with results from any other. Fairly sensitive to sample size (number of data points per species), so this is a better option than species\_sim if you really want units of SDs, but is still not suggested.
- raw:** d is 1 for all species, so output specificity has units of distance, i.e. the raw difference between empirical and simulated RQE\*. This means that results from different variables are not comparable, since it is not scale-invariant to env or hosts\_phylo. It IS still scale-invariant to the species weights in abunds\_mat. Not sensitive to number of samples. Not suggested because units are strange, and isn't comparable between variables.
- flat:** d is RQE\* calculated for a flat abundance distribution, i.e. all abundances = 1. Since RQE\* is a weighted mean, this simplifies d to the mean of unique pairwise distances in env or in the matrix decomposition of hosts\_phylo. While that species distribution may seem (and is) arbitrary, this d has useful properties: scale invariance to abunds\_mat, scale invariance to env/hosts\_phylo, insensitivity to number of samples, insensitivity to occupancy, and strong sensitivity to specificity. Interpretation of results is similar to common SES approaches in that specificity is negative, null is 0, and positive results indicate overdispersion. Default.
- diagnostic** logical. If true, changes output to include different parts of SES. This includes Pval, SES, raw, denom, emp, and all sim values with column labels as simN where N is the number of sims (DEFAULT: FALSE)

## Value

data.frame where each row is an input species. First column is P-value (\$Pval), second column is specificity (\$Spec).

## Author(s)

John L. Darcy

## References

Poulin et al. (2011) Host specificity in phylogenetic and geographic space. Trends Parasitol 8:355-361. doi: 10.1016/j.pt.2011.05.003

## Examples

```
# phylogenetic specificity using endophyte data set
```

```

attach(endophyte)
# only analyze species with occupancy >= 20
m <- occ_threshold(prop_abund(zotutable), 20)
ses_host <- phy_or_env_spec(
  abunds_mat=m,
  hosts=metadata$PlantGenus,
  hosts_phylo=supertree,
  n_cores=12
)

# using vazquez null model from bipartite package as an alternate permutation:
# note that the "creating permuted matrices" step will be slow.
library(bipartite)
ses_host_vaz <- phy_or_env_spec(
  abunds_mat=m,
  hosts=metadata$PlantGenus,
  hosts_phylo=supertree,
  n_cores=12,
  sim_fun=function(m){bipartite::vaznull(1, m)[[1]]},
)

# compare naive permutation vs. vazquez:
plot(ses_host$Spec, ses_host_vaz$Spec, ylab="bipartite::vaznull", xlab="naive")
abline(h=0);abline(v=0)
hist(ses_host_vaz$Spec)
hist(ses_host$Spec)

# environmental specificity using elevation from endophyte data set:
ses_elev <- phy_or_env_spec(
  abunds_mat=m,
  env=metadata$Elevation,
  n_cores=12
)

# geographic specificity using spatial data from endophyte data set:
ses_geo <- phy_or_env_spec(
  abunds_mat=m,
  env=distcalc(metadata$Lat, metadata$Lon),
  n_cores=12
)

```

---

phy\_spec\_sim

*phy\_spec\_sim*


---

## Description

Simulates inputs for `phy_or_env_spec`, by creating a species distribution over an artificial (or real) host phylogenetic tree. For a phylogeny, the species probability distribution  $P(s)$  is based on patristic distances within the tree, such that  $P(s)$  is maximized at zero patristic distance between a tip in the

tree and the ideal host species for *s*. This distribution is given by a truncated normal distribution centered on zero, using only positive values. A uniform proportion (*up*) to that distribution may be added as well, to add a baseline probability to *P(s)*. The standard deviation of *P(s)* can be raised or lowered to simulate cosmopolitanism or specificity.

## Usage

```
phy_spec_sim(sdev, ideal, ideal2 = "", ideal3 = "", n_ideal = 1,
  hosts, hosts_phylo, n_obs, up = 0, oceanp = 0, n_cores = 2,
  seed = 1234567)
```

## Arguments

<code>sdev</code>	numeric vector. Standard deviation of the probability distribution <i>P(s)</i> , in units of patristic distance in <code>hosts_phylo</code> . Low values mean that species <i>s</i> is found with a narrow grouping of hosts, i.e. specificity. High values mean that <i>s</i> is found across a wider group of hosts, i.e. cosmopolitanism. Multiple values can be input in order to simulate a range of specificities, simultaneously. To get a handle on this somewhat opaque variable, consider plotting a histogram of patristic distances within <code>hosts_phylo</code> (see: <code>ape::cophenetic.phylo</code> ). Can be length 1 or <i>n</i> .
<code>ideal</code>	character vector. Tip label of <code>hosts_phylo</code> that is ideal (or closest to ideal) for the simulated species. Does not have to be in <code>hosts</code> , but <b>MUST</b> be in <code>hosts_phylo</code> . Can be length 1 or <i>n</i> .
<code>ideal2</code>	character vector. Tip label of <code>hosts_phylo</code> that is secondary ideal host for the simulated species. Does not have to be in <code>hosts</code> , but <b>MUST</b> be in <code>hosts_phylo</code> . Can be blank (""), if corresponding <code>n_ideal</code> < 2. Can be length 1 or <i>n</i> (default: "").
<code>ideal3</code>	character vector. Tip label of <code>hosts_phylo</code> that is tertiary ideal host for the simulated species. Does not have to be in <code>hosts</code> , but <b>MUST</b> be in <code>hosts_phylo</code> . Can be blank (""), if corresponding <code>n_ideal</code> < 3. Can be length 1 or <i>n</i> (default: "").
<code>n_ideal</code>	integer vector. number of ideal hosts to use. Must be 1, 2, or 3 (DEFAULT=1).
<code>hosts</code>	character vector. Real of fake host identities. All must be tips within <code>hosts_phylo</code> . Analogous to <code>env</code> argument to <code>env_spec_sim</code> .
<code>hosts_phylo</code>	phylo object. Tree containing all unique hosts as tips.
<code>n_obs</code>	integer vector. Number of positive observations to make, i.e. occupancy of simulated species. Can be length 1 or <i>n</i> .
<code>up</code>	numeric vector. <code>up</code> =uniform proportion. This is the proportion of the probability distribution <i>P(species)</i> that is composed of a uniform distribution, if desired. If set to a value above zero (and below 1), <i>P(species)</i> will be a weighted sum of the normal distribution described above, and a uniform distribution. The weight for the uniform distribution will be <code>up</code> , and the weight for the normal distribution will be 1- <code>up</code> (default: 0).
<code>oceanp</code>	numeric vector. See <code>?env_spec_sim</code> for help.
<code>n_cores</code>	integer. Number of CPU cores for parallel computation (DEFAULT: 2).

**seed** integer. Seed for randomization. Daughter seeds will be generated for parallel computations, each with the same number of digits as seed (DEFAULT: 1234567).

### Value

List object containing "matrix" and "params" objects:

**matrix:** matrix where each column is a vector of simulated observations corresponding to a value of hosts; each row represents a simulated species.

**params:** data.frame of parameters (columns) used to simulate each species (rows). A column called "index" is included so that simulated species can be mapped back onto original data structures when some species are omitted due to simulation failure (see fail\_rm).

### Author(s)

John L. Darcy

### Examples

none yet written.

---

plot_grid_abunds	<i>plot_grid_abunds</i>
------------------	-------------------------

---

### Description

plots species abundances across spatial sampling locations

### Usage

```
plot_grid_abunds(grid, abunds, pch = "", ...)
```

### Arguments

<b>grid</b>	data frame with columns x and y, representing cartesian coordinates of sample locations. Can be artificial (generate with randomgrid()) or real.
<b>abunds</b>	abundances of a species, corresponding to rows in grid.
<b>pch</b>	pch character code to use for bottom of each abundance line (DEFAULT="")
<b>...</b>	arguments to be passed to plot.

### Value

returns nothing, just makes a plot.



**Author(s)**

John L. Darcy

**Examples**

```

g1 <- randomgrid()
plot(g1)
a1 <- geo_spec_sim(sdev=c(30, 30, 30, 30), n_obs=1000, grid=g1, up=c(0, 0.20, 0.40, 0.60))
par(mfrow=c(2,2))
plot_grid_abunds(g1, a1$matrix[,1])
plot_grid_abunds(g1, a1$matrix[,2])
plot_grid_abunds(g1, a1$matrix[,3])
plot_grid_abunds(g1, a1$matrix[,4])

```

---

plot_specificities	<i>plot_specificities</i>
--------------------	---------------------------

---

**Description**

Visualizes results from phy\_or\_env\_spec

**Usage**

```

plot_specificities(specs_list, n_bins = 20, col_sig = "black",
  col_nsig = "gray", col_bord = NA, alpha = 0.05, label_cex = 0.6)

```

**Arguments**

specs_list	list of data.frames. Each data.frame must be an output from phy_or_env_spec; must have columns "SES" and "Pval".
n_bins	integer. Number of bins for stacked violins (DEFAULT: 20).
col_sig	string. Color name or hex code for species where Pval <= alpha (DEFAULT = "black").
col_nsig	string. Color name or hex code for species where Pval > alpha (DEFAULT = "gray").
col_bord	string. Color name or hex code for border color. Use NA for no border (DEFAULT = NA).
alpha	float. alpha value for determining statistical significance; see col_sig and col_nsig above (DEFAULT = 0.05).

**Value**

returns nothing (a plot is made).

**Author(s)**

John L. Darcy

**Examples**

none yet written.

---

prop_abund	<i>prop_abund</i>
------------	-------------------

---

**Description**

Calculates proportional abundance of each species (columns) across samples (rows) in community data matrix m. Row sums of output matrix will all be 1.

**Usage**

```
prop_abund(m, to_int = FALSE,  
           max_int = floor(sqrt(.Machine$integer.max)), speciesRows = FALSE)
```

**Arguments**

m	matrix or data frame of numeric values. Columns represent species, rows are samples.
to_int	logical. Should output matrix be transformed into integers from 0 to max_int? Integers take up half as much space as doubles, and as weights are equivalent for calculating specificity. The tradeoff is a little bit of precision (DEFAULT: FALSE).
max_int	integer. Maximum integer value used for to_int. If pairwise geometric means will be calculated with these data, it is nice to keep this value as the square root of the maximum integer size, which is the default.
speciesRows	logical. Do rows represent species (instead of samples)? (DEFAULT:FALSE)

**Value**

matrix of proportional abundances.

**Author(s)**

John L. Darcy

**Examples**

```

attach(endophyte)
m_dbl <- prop_abund(zotutable)
m_int <- prop_abund(zotutable, to_int=TRUE)
head(rowSums(m_dbl))
head(rowSums(m_int))
# note that they are off by a little bit. This small loss in precision is OK.
object.size(m_dbl)
object.size(m_int)
random_positions <- random_rep_positions(m_dbl, 100)
plot(m_int[random_positions] ~ m_dbl[random_positions])

```

---

pval_from_perms	<i>pval_from_perms</i>
-----------------	------------------------

---

**Description**

Calculates P-value for permutation tests.

**Usage**

```
pval_from_perms(emp, perm, tails, method = "MASS_fit", threshold = 30)
```

**Arguments**

<b>emp</b>	Numeric scalar. An empirical test statistic value.
<b>perm</b>	Numeric vector. Test statistic values similar to emp, but calculated from permuted data.
<b>tails</b>	integer. <b>1:</b> Left tail only. <b>2:</b> 2-tailed test. <b>3:</b> Right tail only. <b>0:</b> No test, P=1.
<b>method</b>	string. Method by which P should be calculated from perms: \itemraw: P is calculated as the sum of sim values more extreme than the empirical value plus one, divided by the number of sim values. \itemMASS_fit: P is calculated by fitting a normal distribution to sim values, using the MASS package, and calculating area under the curve from (-inf,emp] or [emp,inf) depending on tailedness. \itemdumb_fit: Same as MASS_fit, except a quick-and-dirty fit is used, which is just using mean and sd of sim values. Actually slower than MASS_fit somehow. Only to be used if MASS_fit isn't working or if you can't install MASS.
<b>threshold</b>	integer. Minimum number n of non-NA values in perm that are acceptable. If n < threshold, P=NA (DEFAULT: 50).

**Value**

a P-value.

**Author(s)**

John L. Darcy

---

randomgrid

*randomgrid*

---

**Description**

Generates a random spatial sampling using a bivariate random uniform distribution.

**Usage**

```
randomgrid(n_samp = 1000, xmin = -100, xmax = 100, ymin = -100,  
           ymax = 100, seed = 123456)
```

**Arguments**

n_samp	number of sampling locations to output (DEFAULT=1000).
xmin	minimum x-axis coordinate (DEFAULT=-100).
xmax	maximum x-axis coordinate (DEFAULT=100).
ymin	minimum y-axis coordinate (DEFAULT=-100).
ymax	maximum y-axis coordinate (DEFAULT=100).
seed	integer, seed for randomization.

**Value**

data.frame object with x and y columns, with n\_samp rows.

**Author(s)**

John L. Darcy

**Examples**

```
g <- randomgrid()  
plot(g)  
g2 <- randomgrid(nsamp=50, xmin=0, ymin=0)  
plot(g2)
```

---

random_rep_positions	<i>random_rep_positions</i>
----------------------	-----------------------------

---

**Description**

Finds positions in a vector (or matrix) that are randomly located within `n_bins` evenly sized bins. This is useful for 1:1 comparisons of large vectors where plotting or comparing all points is prohibitive. Only used in an example for the `prop_abund()` function.

**Usage**

```
random_rep_positions(x, nbins = 50)
```

**Arguments**

<code>x</code>	vector
<code>nbins</code>	number of bins to use

**Value**

integer vector of positions that were selected

**Author(s)**

John L. Darcy

---

rao_quad_ent	<i>rao_quad_ent</i>
--------------	---------------------

---

**Description**

Calculates Rao's (1982) quadratic entropy (FDq) from a distance matrix and a vector of weights (e.g. relative abundance data). In simple terms, FDq is sum product of distances and pairwise products of weights. Default operation in this function is to then divide FDq by the sum of pairwise weights, to give a weighted mean of distances (FDq\*). This gives units of distance, and also makes the metric insensitive to the sum of weights, i.e. weights do not need to be normalized before calculation. This behavior can be disabled by setting `raw=TRUE`, which will give FDq.

**Usage**

```
rao_quad_ent(d, w, raw = FALSE)
```

Arguments

d	numeric dist. Distances, as a dist object. Note that dist objects can easily be made from square matrices using as.dist(), or euclidean distances can be calculated from numeric vectors using dist().
w	numeric vector. Per-sample weights, as a vector. For example, abundances of a species across samples. w MUST be sorted such that it corresponds to rows of as.matrix(d). Thus, w must have length l such that $(l^2-1)/2 = \text{length}(d)$ .
raw	logical. If true, FDq will be returned. If false, FDq* will be returned (DEFAULT: FALSE).

Value

A single value.

Author(s)

John L. Darcy

References

Rao R (1982) Diversity: its measurement, decomposition, apportionment and analysis. Sankhyā: The Indian Journal of Statistics 44(1).

Examples

none yet written

---

tree2mat	<i>tree2mat</i>
----------	-----------------

---

Description

Transforms a phylogenetic tree into a dist object containing patristic distances between tips. Dists are just lower triangles of matrices, and the rows and columns of that matrix are defined by a user-supplied vector of tip labels, which can include duplicate values. Contrast with ape::cophenetic.phylo, which produces a distance matrix containing only unique pairwise patristic distances within the phylogeny.

Usage

```
tree2mat(tree, x, n_cores = 1, delim = ";")
```

**Arguments**

tree	phylo object. Tree containing all unique species in x as tips. May contain tips that are not in x.
x	character vector. Vector of species identities, each of which must be in tree as a tip label. May contain any given species identity more than once.
n_cores	integer. Number of cores to use for parallel computation. No parallelization will be done if n_cores = 1. Multithreading should only be used for large trees where x has low redundancy (DEFAULT = 1).
delim	string. Delimiter character or string for internal use. Must not be present in tree\$tip.label. This is checked by the function and will return an error otherwise (DEFAULT: ";").

**Value**

dist object, of vector length equal to  $(l^2-1)/2$  where  $l$  is length(x); i.e. values are the lower triangle of a patristic distance matrix with rows=x and cols=x.

**Author(s)**

John L. Darcy

**Examples**

```
example_tree <- ape::read.tree(text="(((a:1,b:1):1,c:2):1,d:3):1,(e:1,f:1):3);")
example_x <- c("a", "a", "a", "b", "c", "d", "c", "a", "f")
# unique patristic distance matrix:
ape::cophenetic.phylo(example_tree)
# dist object for example_x:
tree2mat(tree=example_tree, x=example_x)

# examples with other delimiters
tree2mat(tree=example_tree, x=example_x, delim="@")
tree2mat(tree=example_tree, x=example_x, delim="i love cats")
# should fail since "a" is in a tip name:
tree2mat(tree=example_tree, x=example_x, delim="a")
```

---

wpd

---

wpd

---

**Description**

Calculates weighted Phylogenetic Diversity for a vector  $s$  of species observations, weighted by the frequency of each species within  $s$ . For example, if  $S=a, a, b, a, b, c, a$ , then species  $a$  will have weight 4, species  $b$  will have weight 2, and species  $c$  will have weight 1. Unobserved species have weight zero.

**Usage**

```
wpd(s, s_phylo, w = NULL, nested_set = NULL, metric = "Hp")
```

**Arguments**

- |            |   |
|------------|---|
| s          | character vector. One species name per observation. If no species was observed for a given datum, use NA. s can also be provided as a vector of unique species identities, in which case counts of those species can be given as w.   |
| s_phylo    | phylo object. Tree containing all unique names in s as tips. Must not contain duplicate tip labels.   |
| w          | numeric vector. Optional weights for s, e.g. number of parasites observed in each sample, or boolean weights corresponding to presence or absence of parasite species, or confidence species was observed, etc. If w is not provided but a weighted metric is specified, w will be set to 1 for each value of s. Thus, weights for each unique species in s would be equal to the number of times that species appears in s. w is not used for unweighted metrics (PD). Any NA values in w will be pairwise removed from w and s (DEFAULT: NULL).   |
| nested_set | matrix. The output of make_nested_set(s_phylo). If not provided, will be calculated on the fly. Precalculation only provides speedup with very large trees (DEFAULT: NULL).   |
| metric     | <p>character. Abbreviated name of desired tree-based phylogenetic diversity metric. Available metrics are:</p> <p><b>Hp:</b> Phylogenetic Entropy. Insensitive to 0 weights, cannot increase with removal of taxa. Allen et al. 2009.</p> <p><b>WF:</b> Weighted Faith's PD. Sensitive to 0 weights, i.e. a clade that was heavily sampled but has lots of zeroes will cause its sister clades to be underrepresented. Swenson 2014.</p> <p><b>PD:</b> Original Faith's Phylogenetic Diversity. Unweighted. Simply a sum of branch-lengths in your tree (but only for taxa in s). Faith 1992.</p> |

**Details**

However, one may wish to exclude observations that do not meet some criterion, such as co-observation of a symbiote or parasite. For this reason, a second set of weights w can be provided as a vector of numeric values that are paired with s. These weights are then implicitly combined with the weights discussed above depending on which weighted metric is chosen. In the case of Phylogenetic Entropy (Hw), per-tip weights are calculated as the sums of w. In the case of Weighted Faith (WF), per-tip weights are averages of w.

**Value**

Single WPD or PD value.

**Author(s)**

John L. Darcy



References

Allen B, Kon M, Bar-Yam Y (2009) A new phylogenetic diversity measure generalizing the Shannon index and its application to Phyllostomid bats. *American Naturalist* 174(2). Swenson NG (2014) *Functional and Phylogenetic Ecology in R*. Springer UseR! Series, Springer, New York, New York, U.S.A. Faith DP (1992) Conservation evaluation and phylogenetic diversity. *Biological Conservation* 61.

See Also

rao\_quad\_ent, a phylogenetic diversity measure that uses a distance matrix instead of a phylogenetic tree.

Examples

```
library(geiger)
set.seed(12345)
s_phylo <- get(data(geospiza))$phy
w <- sample(c(0, 1), replace=T, size=10)
s <- sample(s_phylo$tip.label, replace=T, size=10)
wpd(s, s_phylo, w, metric="Hp")
```

---

wpd_table	<i>wpd_table</i>
-----------	------------------

---

Description

Calculates phylogenetic entropy (Hp) for each column vector s of species observations within matrix m, weighted by the frequency of each species within s. Can also calculate Faith's PD.

Usage

```
wpd_table(m, s_phylo, nested_set, metric = "Hp", ncores = 4)
```

Arguments

- m matrix of species observation vectors (s). See s argument of wpd().
- s\_phylo phylo object. Tree containing all unique names in s as tips. Must not contain duplicate tip labels.
- nested\_set matrix. The output of make\_nested\_set(s\_phylo). If not provided, will be calculated on the fly. Precalculation only provides speedup with very large trees (DEFAULT: NULL).
- metric character. Abbreviated name of desired tree-based phylogenetic diversity metric. Available metrics are:  
**Hp:** Phylogenetic Entropy. Insensitive to 0 weights, cannot increase with removal of taxa. Allen et al. 2009.

**WF:** Weighted Faith's PD. Sensitive to 0 weights, i.e. a clade that was heavily sampled but has lots of zeroes will cause its sister clades to be underrepresented. Swenson 2014.

**PD:** Original Faith's Phylogenetic Diversity. Unweighted. Simply a sum of branch- lengths in your tree (but only for taxa in s). Faith 1992.

### Value

multiple WPD or PD values, one for each column of m.

### Author(s)

John L. Darcy

### References

Allen B, Kon M, Bar-Yam Y (2009) A new phylogenetic diversity measure generalizing the Shannon index and its application to Phyllostomid bats. *American Naturalist* 174(2). Swenson NG (2014) *Functional and Phylogenetic Ecology in R*. Springer UseR! Series, Springer, New York, New York, U.S.A. Faith DP (1992) Conservation evaluation and phylogenetic diversity. *Biological Conservation* 61. Rao R (1982) Diversity: its measurement, decomposition, apportionment and analysis. *Sankhyā: The Indian Journal of Statistics* 44(1).

### Examples

none yet written.

# Index

bl\_distance\_ns, [2](#)  
check\_pes\_inputs, [3](#)  
daughter\_seeds, [3](#)  
distcalc, [4](#)  
env\_spec\_sim, [5](#)  
geo\_spec\_sim, [6](#)  
make\_nested\_set, [8](#)  
occ\_threshold, [9](#)  
pairwise\_geo\_mean, [10](#)  
pairwise\_product, [11](#)  
phy\_or\_env\_spec, [11](#)  
phy\_spec\_sim, [14](#)  
plot\_grid\_abunds, [16](#)  
plot\_specificities, [17](#)  
prop\_abund, [18](#)  
pval\_from\_perms, [19](#)  
random\_rep\_positions, [21](#)  
randomgrid, [20](#)  
rao\_quad\_ent, [21](#)  
tree2mat, [22](#)  
wpd, [23](#)  
wpd\_table, [25](#)