

# Package ‘specificity’

May 17, 2021

**Title** Calculate Environmental or Host Phylogenetic Specificity

**Version** 0.1.6.9000

**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

**Depends** R (>= 3.5.0)

**Encoding** UTF-8

**LazyData** true

**RoxygenNote** 7.1.1

**Imports** ape, parallel, Rcpp, fields, graphics, grDevices, stats

**LinkingTo** Rcpp

**Suggests** testthat

**NeedsCompilation** yes

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

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

## R topics documented:

bl_distance_ns . . . . .	2
calculate_spec_and_pval . . . . .	3
check_pes_inputs . . . . .	5
circularize2dist . . . . .	6
distcalc . . . . .	7
endophyte . . . . .	8
env_spec_sim . . . . .	8
geo_spec_sim . . . . .	10
get_ga_defaults . . . . .	12
make_nested_set . . . . .	13
occ_threshold . . . . .	14
onto2nwk . . . . .	15

pairwise_product . . . . .	16
phy_or_env_spec . . . . .	17
phy_spec_sim . . . . .	20
plot_grid_abunds . . . . .	22
plot_pairwise_spec . . . . .	23
plot_specs_stacks . . . . .	24
plot_specs_violin . . . . .	26
prop_abund . . . . .	27
randomgrid . . . . .	28
random_rep_positions . . . . .	29
raolsp . . . . .	30
raoperms . . . . .	31
rao_genetic_max . . . . .	32
rao_sort_max . . . . .	33
tips_from_node . . . . .	34
tree2mat . . . . .	35
wpd . . . . .	36
wpd_table . . . . .	37
<b>Index</b>	<b>40</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(specificity)
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
```

---

```
calculate_spec_and_pval
```

```
calculate_spec_and_pval
```

---

**Description**

This function is called by `phy_or_env_spec()`. It is made available as a standalone function in the (rare) case a user wishes to calculate Spec using their own null model. `calculate_spec_and_pval()` takes empirical rao values and sim rao values (from a null model) and calculates Spec and P-values. To do that, use your own null model to make species data, and use `rao1sp()` and/or `raoperms()` to get raw rao values. This function expects a vector of empirical values, and a list of vectors of sim values (see below). Most of the inputs for this function are the same as `phy_or_env_spec()`. Think of this function as the final component of "build your own `phy_or_env_spec()`". Note that for this custom approach, the environmental variable must be a dist.

**Usage**

```
calculate_spec_and_pval(
  emp_raos,
  sim_raos,
  abunds_mat,
  env,
  p_adj = "fdr",
  tails = 1,
  n_cores = 2,
  verbose = TRUE,
  p_method = "raw",
  center = "mean",
  denom_type = "index_full",
  diagnostic = FALSE,
  ga_params = get_ga_defaults()
)
```

**Arguments**

<code>emp_raos</code>	vector. Empirical rao values, one per species ("feature").
<code>sim_raos</code>	list of numeric vectors. Sim rao values, generated under null hypothesis. Each item in list corresponds to an entry in <code>emp_raos</code> . As such, <code>length(emp_raos)</code> must equal <code>length(sim_raos)</code> . Each item within <code>sim_raos</code> is a vector of rao values ( <code>length=n_sim</code> in the case of <code>phy_or_env_spec()</code> ).
<code>abunds_mat</code>	site x species matrix. See <code>?phy_or_env_spec</code> .
<code>env</code>	MUST BE A dist OBJECT!!!! VERY IMPORTANT!!!! See <code>?phy_or_env_spec</code> .
<code>p_adj</code>	string. Type of multiple hypothesis testing correction performed on P-values. Can take any valid method argument to <code>p.adjust</code> , including "none", "bonferroni", "holm", "fdr", and others (default: "fdr").
<code>tails</code>	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).
<code>n_cores</code>	integer. Number of CPU cores to use for parallel operations. If set to 1, <code>lapply</code> will be used instead of <code>mclapply</code> (default: 2).
<code>verbose</code>	logical. Should status messages be displayed? (default: TRUE).
<code>p_method</code>	string. "raw" for quantile method, or "gamma_fit" for calculating P by fitting a gamma distribution (default: "raw").
<code>center</code>	string. Type of central tendency to use for simulated RQE values. Options are "mean", "median", and "mode". If mode is chosen, a reversible gamma distribution is fit and mode is calculated using that distribution (default: mean).
<code>denom_type</code>	string. Type of denominator (d) to use (default: "index"). Note that denominator type does NOT affect P-values.

**"ses":** d for species s is calculated as the standard deviation of RQE values calculated from permuted species weights. This makes the output specificity a standardized effect size (SES). Unfortunately, this makes SES counterintuitively sensitive to occupancy, where species with high occupancy have more extreme SES than rare species, due to their more deterministic sim specificities. Included for comparative purposes, 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 not scale-invariant to the species weights in `abunds_mat`, either. Not sensitive to number of samples. Not suggested because units are strange, and isn't comparable between variables.

**"index\_full":** d is the mean of simulated (permuted) RQE values for species that have stronger specificity than expected by chance, resulting in specificity values with range [-1, 0), with 0 as the null hypothesis. In this case, -1 indicates perfect specificity, where a species is associated with zero environmental variability. In the euclidean sense, this could be a species that is always found at the exact same elevation or the exact same pH. For species that have weaker specificity than expected by chance, d is x minus the center (see above) of simulated RQE values, where x is the maximum possible dissimilarity observable given species weights. In "index\_full", x is estimated using a genetic algorithm. This d has other useful properties: scale

invariance to env/hosts\_phylo, insensitivity to the number of samples, insensitivity to occupancy, and strong sensitivity to specificity (default).

**"index\_rough"**: Same as "index\_full", but for species where specificity is weaker than expected by chance, a rough approximation is used to vastly reduce computational load. This approximation does NOT give values where 1 is maximized generality. instead, the values will tightly correlate to the correct values, but the slope of that relationship will not be known a priori. In other words, if you don't care about species that are more general than expected by chance, this is fine.

**"index\_fast"**: Same as "index\_full", but results as per "index\_rough" are used together with a subset of results from "index\_full" to create a model relating the two. Thus, speed is a compromise between the two approaches, with medium-high accuracy.

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)
ga_params	list. Parameters for genetic algorithm that maximizes RQE. Only used with denom_type="index_full/rough/fast". Default is the output of get_ga_defaults(). If different parameters are desired, start with output of get_ga_defaults and modify accordingly.

### 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

### Examples

```
# None yet. Forthcoming examples:
# 1. calculating regular old elevational specificity the hard way
# 2. same thing, but using vazquez null model from bipartite package
```

---

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**

abunds\_mat (required, see phy\_or\_env\_spec)  
 env (required, can be NULL, see phy\_or\_env\_spec)  
 hosts (required, can be NULL, see phy\_or\_env\_spec)  
 hosts\_phylo (required, can be NULL, see phy\_or\_env\_spec)  
 verbose logical. Should status messages be displayed? (default: TRUE).

**Value**

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

**Examples**

```
library(specificity)
attach(endophyte)
m <- occ_threshold(prop_abund(otutable), threshold=10)
check_pes_inputs(m, env=metadata$Elevation, hosts=NULL, hosts_phylo=NULL)
check_pes_inputs(m, env=NULL, hosts=metadata$PlantGenus, hosts_phylo=supertree)
aspect_dis <- circularize2dist(metadata$Aspect, 360)
check_pes_inputs(m, env=aspect_dis, hosts=NULL, hosts_phylo=NULL)
```

---

circularize2dist	<i>circularize2dist</i>
------------------	-------------------------

---

**Description**

Circularizes a vector into a dist object. For example, a vector of days of the year, where the distance between 365 and 2 should be less than the distance between 350 and 365. Another example may be direction, where 0.1 and 2pi radians are close together.

**Usage**

```
circularize2dist(x, maxx)
```

**Arguments**

x a numeric vector. All values should be >0.  
 maxx the maximum theoretical value (also the zero value!) of variable x. In the example of months of the year, maxx would be 12, even if you only had data for months 1-8. For degrees, maxx=360. For radians, maxx=2\*pi. Must be greater than or equal to values of x.

**Value**

a vector of differences, ordered identically to a "dist" object.

**Author(s)**

John L. Darcy

**Examples**

```
library(specificity)
# make some fake data to represent months of the year
months <- c(1, 4, 11)
# run circularize2dist() on the months. Must specify that
# maxx = 12, since december is both 12 and 0 for these data.
circularize2dist(months, 12)
# output is a distance matrix.
# rows and cols of months_circdm are months - it's ordered.
# notice the distance between 11 and 1 is 2, not 10!
```

---

distcalc

*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

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

---

endophyte	<i>Foliar endophytic fungi across the Hawaiian Archipelago</i>
-----------	--

---

Description

A dataset containing an OTU table (species-by-site), environmental metadata, and host plant phylogeny.

Usage

```
endophyte
```

Format

- A list containing 3 objects:
- otutable:** data.frame object where each row is a sample and each column is a fungal OTU (actually ASV from DADA2). Rownames are sample IDs.
  - metadata:** data.frame object containing environmental metadata for samples in otutable. SampleID column of metadata matches rownames of otutable.
  - supertree:** Phylogenetic tree containing all host plant genera in PlantGenus column of metadata.

Source

Darcy et al. (2020) Fungal communities living within leaves of native Hawaiian dicots are structured by landscape-scale variables as well as by host plants. Mol Ecol 29:3102-3115 <https://doi.org/10.1111/mec.15544>

---

env_spec_sim	<i>env_spec_sim</i>
--------------	---------------------

---

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 weaker specificity, and a low SD means stronger 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**

sdev	numeric vector. Standard deviation of the probability distribution $P(\text{species})$ , in the same units of env. Low values mean that the species is found across only a narrow range of env, i.e. specificity. High values mean that the species is found across a wide range of env, i.e. cosmopolitanism. Multiple values can be input in order to simulate a range of specificities simultaneously. Can be length 1 or n.
ideal	numeric vector. Value of env that is ideal for the simulated species. This is the mode of the probability distribution $P(\text{species})$ . Can be length 1 or n.
ideal2	numeric vector. Value of env that is the second ideal for the simulated species. Only used if $n\_ideal \geq 2$ . This is the second mode of the probability distribution $P(\text{species})$ . Can be length 1 or n.
ideal3	numeric vector. Value of env that is the third ideal for the simulated species. Only used if $n\_ideal = 3$ . This is the third mode of the probability distribution $P(\text{species})$ . Can be length 1 or n.
n_ideal	integer vector. Number of ideal values for the simulated species, i.e. modality of that species' distribution across env; 1 for unimodal, etc. Only can use values 1, 2, or 3, which correspond to ideal, ideal2, and ideal3. Can be length 1 or n (default: 1).
env	numeric vector. Real or fake environmental variable.
n_obs	integer vector. Number of positive observations to make, i.e. occupancy of simulated species. Can be length 1 or n (default: 1).
up	numeric vector. up=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 up, and the weight for the normal distribution will be 1-up (default: 0).
oceanp	numeric vector. oceanp=ocean proportion. This is the proportion of samples in env that are "in the ocean", i.e. samples where the species would not expect to

be found even if env 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

*geo\_spec\_sim*

---

### 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).
ideal_y3	numeric vector. y-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 up, and the weight for the normal distribution will be 1-up (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).

**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**

```
library(specificity)
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")
```

---

get\_ga\_defaults

get\_ga\_defaults

---

**Description**

Simply returns default parameters for the genetic algorithm in rao\_genetic\_max(). This function has no parameters.

**Usage**

```
get_ga_defaults()
```

**Value**

named list of parameters.

**Author(s)**

John L. Darcy

**Examples**

```
# get_ga_defaults()
```

---

make_nested_set	<i>make_nested_set</i>
-----------------	------------------------

---

**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

**Examples**

```
library(specificity)
library(ape)
phy <- get(data(endophyte))$supertree
# check if tree is rooted:
ape::is.rooted(phy)
# make nested set table:
phy_ns <- make_nested_set(phy)
# show that nested set table matches up with edges table in phy:
all(phy$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

```
library(specificity)
attach(endophyte)
dim(otutable)
otutable_over25 <- occ_threshold(otutable, 25)
dim(otutable_over25)
```

---

onto2nwk	<i>onto2nwk</i>
----------	-----------------

---

## Description

Converts an ontology (hierarchical categories) into a nwk phylogeny.

## Usage

```
onto2nwk(df)
```

## Arguments

df	a data.frame object where columns represent ontology levels, which are assumed to be nested hierarchically. this function does not check for proper hierarchical nestedness - it is the user's job to check that each node and tip name is monophyletic. Lower levels (e.g. tips) should be the rightmost column of df, and higher levels (e.g. roots) should be leftmost column, with intermediate columns ordered between.
----	--

## Value

A newick (nwk) format string.

## Author(s)

John L. Darcy

## Examples

```
library(specificity)
library(ape)
df <- data.frame(
  l1 = c( "a", "a", "a", "a", "a", "a", "a", "b", "b", "b", "b", "b", "b", "c", "d"),
  l2 = c( "e", "e", "e", "e", "f", "f", "g", "h", "h", "h", "i", "j", "j", "k", "l"),
  l3 = c( "m", "n", "o", "o", "p", "p", "q", "r", "r", "s", "t", "u", "v", "w", "x")
)
```

```
nwk_str <- onto2nwk(df)
a <- ape::read.tree(text=nwk_str)
plot(a, show.node.label=TRUE)
```

---

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  $(l^2-1)/2$ , where  $l=\text{length}(x)$ .

## Author(s)

John L. Darcy

## Examples

```
x <- 1:6
y_cpp <- pairwise_product(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 Rao's Quadratic Entropy to simulated RQE (same but with permuted abundances). By default (*denom\_type* = "index"), an index is calculated from *emp* and *sim* values such that *Spec*=0 indicates 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,
  p_adj = "fdr",
  seed = 1234567,
  tails = 1,
  n_cores = 2,
  verbose = TRUE,
  p_method = "raw",
  center = "mean",
  denom_type = "index_full",
  diagnostic = F,
  chunksize = 1000,
  ga_params = get_ga_defaults()
)
```

## Arguments

<i>abunds_mat</i>	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, specificity cannot be calculated, and NAs will be returned. Negative values in <i>abunds_mat</i> are not allowed (REQUIRED).
<i>env</i>	numeric vector, dist, or square matrix. Environmental variable corresponding to <i>abunds</i> . For example, temperature, or geographic distance. Not required for computing phylogenetic specificity (default: NULL).
<i>hosts</i>	character vector. Host identities corresponding to <i>abunds</i> . Only required if calculating phylogenetic specificity (default: NULL).
<i>hosts_phylo</i>	phylo object. Tree containing all unique hosts as tips. Only required if calculating 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).
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. Same seed will be used for each species in abunds_mat, so all species will experience the same permutations. This can be disabled by setting seed=0, which will make permutation is both non deterministic (not repeatable) AND each species will experience different permutations (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. "raw" for quantile method, or "gamma_fit" for calculating P by fitting a gamma distribution (default: "raw").
center	string. Type of central tendency to use for simulated RQE values. Options are "mean", "median", and "mode". If mode is chosen, a reversible gamma distribution is fit and mode is calculated using that distribution (default: mean).
denom_type	string. Type of denominator (d) to use (default: "index"). Note that denominator type does NOT affect P-values.

**"ses":** d for species s is calculated as the standard deviation of RQE values calculated from permuted species weights. This makes the output specificity a standardized effect size (SES). Unfortunately, this makes SES counterintuitively sensitive to occupancy, where species with high occupancy have more extreme SES than rare species, due to their more deterministic sim specificities. Included for comparative purposes, 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 not scale-invariant to the species weights in aunds\_mat, either. Not sensitive to number of samples. Not suggested because units are strange, and isn't comparable between variables.

**"index\_full":** d is the mean of simulated (permuted) RQE values for species that have stronger specificity than expected by chance, resulting in specificity values with range [-1, 0), with 0 as the null hypothesis. In this case, -1 indicates perfect specificity, where a species is associated with zero environmental variability. In the euclidean sense, this could be a species that is always found at the exact same elevation or the exact same pH. For species that have weaker specificity than expected by chance, d is x minus the center (see above) of simulated RQE values, where x is the maximum possible dissimilarity observable given species weights. In "index\_full", x is estimated using a genetic algorithm. This d has other useful properties: scale

invariance to env/hosts\_phylo, insensitivity to the number of samples, insensitivity to occupancy, and strong sensitivity to specificity (default).

**"index\_rough"**: Same as "index\_full", but for species where specificity is weaker than expected by chance, a rough approximation is used to vastly reduce computational load. This approximation does NOT give values where 1 is maximized generality. Instead, the values will tightly correlate to the correct values, but the slope of that relationship will not be known a priori. In other words, if you don't care about species that are more general than expected by chance, this is fine.

**"index\_fast"**: Same as "index\_full", but results as per "index\_rough" are used together with a subset of results from "index\_full" to create a model relating the two. Thus, speed is a compromise between the two approaches, with medium-high accuracy.

diagnostic	logical. If true, changes output to include different parts of Spec. This includes Pval, Spec, raw, denom, emp, and all sim values with column labels as simN where N is the number of sims (default: FALSE)
chunksize	integer. If greater than zero, computation of sim RAO values will be done using chunked evaluation, which lowers memory use considerably for larger data sets. Can be disabled by setting to 0. Default value is 1000 species per chunk (default: 1000).
ga_params	list. Parameters for genetic algorithm that maximizes RQE. Only used with denom_type="index". Default is the output of get_ga_defaults(). If different parameters are desired, start with output of get_ga_defaults and modify accordingly.

## 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
- Rao CR (2010) Quadratic entropy and analysis of diversity. Sankhya 72:70-80. doi: 10.1007/s13171-010-0016-3
- Rao CR (1982) Diversity and dissimilarity measurements: A unified approach. Theor Popul Biol 21:24-43.

## Examples

```
library(specificity)
attach(endophyte)
# only analyze species with occupancy >= 20
```

```

m <- occ_threshold(prop_abund(otutable), 20)
# create list to hold phy_or_env_spec outputs
specs_list <- list()

# phylogenetic specificity using endophyte data set
specs_list$host <- phy_or_env_spec(
  abunds_mat=m,
  hosts=metadata$PlantGenus,
  hosts_phylo=supertree,
  n_sim=100, p_method="gamma_fit",
  n_cores=4
)

# environmental specificity using elevation from endophyte data set:
specs_list$elev <- phy_or_env_spec(
  abunds_mat=m,
  env=metadata$Elevation,
  n_sim=100, p_method="gamma_fit",
  n_cores=4
)

# geographic specificity using spatial data from endophyte data set:
specs_list$geo <- phy_or_env_spec(
  abunds_mat=m,
  env=distcalc(metadata$Lat, metadata$Lon),
  n_sim=100, p_method="gamma_fit",
  n_cores=4
)

plot_specs_violin(specs_list, cols=c("forestgreen", "red", "black"))

```

---

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

sdev	numeric vector. Standard deviation of the probability distribution $P(s)$ , in units of patristic distance in <code>hosts_phylo</code> . Low values mean that species $s$ is found with a narrow grouping of hosts, i.e. specificity. High values mean that $s$ 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 $n$ .
ideal	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 MUST be in <code>hosts_phylo</code> . Can be length 1 or $n$ .
ideal2	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 MUST be in <code>hosts_phylo</code> . Can be blank (""), if corresponding <code>n_ideal</code> < 2. Can be length 1 or $n$ (default: "").
ideal3	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 MUST be in <code>hosts_phylo</code> . Can be blank (""), if corresponding <code>n_ideal</code> < 3. Can be length 1 or $n$ (default: "").
n_ideal	integer vector. number of ideal hosts to use. Must be 1, 2, or 3 (default: 1).
hosts	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> .
hosts_phylo	phylo object. Tree containing all unique hosts as tips.
n_obs	integer vector. Number of positive observations to make, i.e. occupancy of simulated species. Can be length 1 or $n$ .
up	numeric vector. <code>up=uniform proportion</code> . 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).
oceanp	numeric vector. See <code>?env_spec_sim</code> for help.

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).

**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**

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

**Value**

returns nothing, just makes a plot.

**Author(s)**

John L. Darcy

**Examples**

```

library(specificity)
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_pairwise_spec	<i>plot_pairwise_spec</i>
--------------------	---------------------------

---

**Description**

Plots pairwise correlations between specificity to multiple variables. Specificity results are supplied to this function as a list of specificity tables, i.e. a list where each object within the list is an output of `phy_or_env_spec`, and all were created using the same `abunds_mat` object (see: `?phy_or_env_spec`).

**Usage**

```

plot_pairwise_spec(
  sl,
  label_cex = 1,
  point_cex = 1,
  cor_cex = 2,
  cor_red_lim = 0.7,
  method = "pearson"
)

```

**Arguments**

<code>sl</code>	"specs list" list of outputs from <code>phy_or_env_spec</code> as described above.
<code>label_cex</code>	float. Size of variable labels, which will be displayed along the plot's diagonal. Use cex units; see <code>?par</code> (default: 1).
<code>point_cex</code>	float. Size of points in the plot's lower triangle. Useful to reduce this if you are plotting lots of species. Use cex units; see <code>?par</code> (default: 1).
<code>cor_cex</code>	float. Size of text for correlations displayed in plot's upper triangle. Use cex units; see <code>?par</code> (default: 1).

`cor_red_lim` float. Correlation coefficients will be shown in red if they are equal to or more extreme than this value (default: 0.70).

`method` string. Preferred correlation method. see `?cor` for options (default: "pearson").

### Value

Returns nothing. Plots correlations in a square matrix of subplots, where variable names are shown in the diagonal, pairwise specificities are plotted in the lower triangle, and correlation coefficients are displayed in the upper triangle. For plots in the lower triangle, each point represents a species.

### Author(s)

John L. Darcy

### Examples

```
library(specificity)
attach(endophyte)
# only analyze species with occupancy >= 20
m <- occ_threshold(prop_abund(otutable), 20)
# create list to hold phy_or_env_spec outputs
specs_list <- list()
specs_list$NDVI <- phy_or_env_spec(m, env=metadata$NDVI,
  n_cores=10, n_sim=50, p_method="gamma_fit")
specs_list$Evapotranspiration <- phy_or_env_spec(m,
  env=metadata$Evapotranspiration, n_cores=10,
  n_sim=100, p_method="gamma_fit")
specs_list$Rainfall <- phy_or_env_spec(m, env=metadata$Rainfall,
  n_cores=10, n_sim=50, p_method="gamma_fit")
plot_pairwise_spec(specs_list)
```

---

plot\_specs\_stacks      *plot\_specs\_stacks*

---

### Description

Visualizes results from `phy_or_env_spec` as stacked histograms. Aliased to `plot_specificities()` for backward compatibility.

### Usage

```
plot_specs_stacks(
  specs_list,
  n_bins = 20,
  col_sig = "black",
  col_nsig = "gray",
  col_bord = NA,
```



```

    alpha = 0.05,
    label_cex = 1
  )

```

### Arguments

specs_list	list of data.frames. Each data.frame must be an output from phy_or_env_spec; must have columns "Spec" 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).
label_cex	float. Used to change size of x-axis labels (default: 1).

### Value

returns nothing (a plot is made).

### Author(s)

John L. Darcy

### Examples

```

library(specificity)
attach(endophyte)
# only analyze species with occupancy >= 20
m <- occ_threshold(prop_abund(otutable), 20)
# create list to hold phy_or_env_spec outputs
specs_list <- list()
specs_list$NDVI <- phy_or_env_spec(m, env=metadata$NDVI,
  n_cores=10, n_sim=50, p_method="gamma_fit")
specs_list$Evapotranspiration <- phy_or_env_spec(m,
  env=metadata$Evapotranspiration, n_cores=10,
  n_sim=100, p_method="gamma_fit")
specs_list$Rainfall <- phy_or_env_spec(m, env=metadata$Rainfall,
  n_cores=10, n_sim=50, p_method="gamma_fit")
plot_specs_stacks(specs_list)

```

---

plot_specs_violin	<i>plot_specs_violin</i>
-------------------	--------------------------

---

## Description

Visualizes results from `phy_or_env_spec` as violins. Violin area is proportionally divided such that lighter colors represent density of non-significant features, and darker colors represent statistically significant features.

## Usage

```
plot_specs_violin(
  specs_list,
  cols = "black",
  cols_bord = "white",
  alpha = 0.05,
  label_cex = 1,
  nsig_trans = 0.3,
  minval = -1,
  maxval = 1,
  ylab = "Spec",
  ...
)
```

## Arguments

<code>specs_list</code>	list of data.frames. Each data.frame must be an output from <code>phy_or_env_spec</code> ; must have columns "Spec" and "Pval".
<code>cols</code>	character vector of color names or hex codes. If only one value is given, all violins will be that color. Otherwise, one value may be given per item in <code>specs_list</code> , corresponding to its order (default: "black").
<code>cols_bord</code>	character vector of color names or hex codes. Color name or hex code for borders drawn around and within violins. Length 1 or length n, just like <code>cols</code> . For no borders, use <code>cols_bord=NA</code> (default: "white").
<code>alpha</code>	float. alpha value for determining statistical significance (default: 0.05).
<code>label_cex</code>	float. Used to change size of x-axis labels (default: 1).
<code>nsig_trans</code>	float between 0 and 1 (inclusive). Determines how transparent violin area will be for nonsignificant features, with 0 meaning totally transparent and 1 meaning totally opaque (default: 0.4).
<code>minval</code>	minimum possible value for specificity statistic (default: -1).
<code>maxval</code>	maximum possible value for specificity statistic (default: 1).
<code>ylab</code>	y-axis label for plot (default: "Spec").
<code>...</code>	additional arguments to be passed to <code>polygons()</code> .

**Value**

returns nothing (a plot is made).

**Author(s)**

John L. Darcy

**Examples**

```
library(specificity)
attach(endophyte)
# only analyze species with occupancy >= 20
m <- occ_threshold(prop_abund(otutable), 20)
# create list to hold phy_or_env_spec outputs
specs_list <- list()
specs_list$NDVI <- phy_or_env_spec(m, env=metadata$NDVI,
  n_cores=10, n_sim=50, p_method="gamma_fit")
specs_list$Evapotranspiration <- phy_or_env_spec(m,
  env=metadata$Evapotranspiration, n_cores=10,
  n_sim=100, p_method="gamma_fit")
specs_list$Rainfall <- phy_or_env_spec(m, env=metadata$Rainfall,
  n_cores=10, n_sim=50, p_method="gamma_fit")
# default black
plot_specs_violin(specs_list)
# with colors
plot_specs_violin(specs_list, cols=c("forestgreen", "gold", "blue"))
# with border colors
plot_specs_violin(specs_list, cols=c("forestgreen", "gold", "blue"),
  cols_bord=c("red", "blue", "black"))
# with thicker borders (arg "lwd" is passed to polygon())
plot_specs_violin(specs_list, cols=c("forestgreen", "gold", "blue"),
  cols_bord="black", lwd=3)
# with NO borders
plot_specs_violin(specs_list, cols=c("forestgreen", "gold", "blue"),
  cols_bord=NA)
```

---

prop\_abund

*prop\_abund*


---

**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

<code>m</code>	matrix or data frame of numeric values. Columns represent species, rows are samples.
<code>to_int</code>	logical. Should output matrix be transformed into integers from 0 to <code>max_int</code> ? 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).
<code>max_int</code>	integer. Maximum integer value used for <code>to_int</code> . 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.
<code>speciesRows</code>	logical. Do rows represent species (instead of samples)? (default:FALSE)

### Value

matrix of proportional abundances.

### Author(s)

John L. Darcy

### Examples

```

library(specificity)
attach(endophyte)
m_dbl <- prop_abund(otutable)
m_int <- prop_abund(otutable, 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])

```

---

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**

```
library(specificity)
g <- randomgrid()
plot(g)
g2 <- randomgrid(n_samp=50, xmin=0, ymin=0)
plot(g2)
```

---

random\_rep\_positions    *random\_rep\_positions*

---

**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**

x	vector
nbins	number of bins to use

**Value**

integer vector of positions that were selected

**Author(s)**

John L. Darcy

**Examples**

```
library(specificity)
attach(endophyte)
m_dbl <- prop_abund(otutable)
m_int <- prop_abund(otutable, 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])
```

---

rao1sp

*rao1sp*


---

**Description**

Calculate's Rao's quadratic entropy for one species (rao1sp = "rao one species").

**Usage**

```
rao1sp(p, D, perm, seed)
```

**Arguments**

<code>p</code>	numeric vector of length <code>n</code> . a species weights vector.
<code>D</code>	numeric vector of length $n(n-1)/2$ . i.e. a dist object whose full matrix is $n \times n$ .
<code>perm</code>	bool. Whether or not the permute order of <code>p</code> before calculating Rao (default: FALSE).
<code>seed</code>	integer. a seed to be used if <code>perm=TRUE</code> . setting <code>seed=0</code> will give nondeterministic random results, as if no seed were set (default: 0).

**Value**

A single Rao value.

**Author(s)**

John L. Darcy

**Examples**

```
library(specificity)
p <- 1:5/sum(1:5)
D <- dist(6:10)
rao1sp(p, D)
```

---

raoperms

*raoperms*


---

**Description**

C++ function used by `phy_or_env_spec()`. Not meant for use otherwise.

**Usage**

```
raoperms(p,D,n_sim,seed)
```

**Arguments**

<code>p</code>	numeric vector of species weights.
<code>D</code>	numeric vector (dist) of distances corresponding to a lower triangle of a matrix whose rows and cols correspond to <code>p</code> ; i.e. an $l \times l$ matrix where $l$ is <code>length(p)</code> . R's <code>dist()</code> function does this for you!
<code>n_sim</code>	integer. number of sims to do (default: 1000).
<code>seed</code>	integer. a seed to be used for permutation. setting <code>seed=0</code> will give nondeterministic random results, as if no seed were set (default: 0).

**Value**

Vector of Rao values, `length = n_sim`.

**Author(s)**

John L. Darcy

**Examples**

```
library(specificity)
p <- runif(100)
D <- dist(runif(100))
a <- raoperms(p,D,100,12345)
hist(a)
```

---

rao\_genetic\_max

*rao\_genetic\_max*


---

**Description**

Uses a genetic algorithm to find the optimum permutation of p to maximize Rao(p,D).

**Usage**

```
rao_genetic_max(p, D, term_cycles, maxiters, popsize, keep, prc)
```

**Arguments**

p	numeric vector of length n - a species weights vector.
D	numeric vector of length $n(n-1)/2$ - i.e. a dist object whose full matrix is nxn.
term_cycles	integer, number of cycles with no improvement to trigger termination (default: 10).
maxiters	integer, maximum number of iterations to run algorithm (default: 400).
popsize	integer, population size for genetic algorithm (default: 300).
keep	integer, number of individuals to keep during each iteration (default: 5).
prc	double, precision for calculating termination with term_cycles (default: 0.001).

**Value**

List object containing results of genetic algorithm:

**best\_rao:** Maximum Rao value found.

**iter\_raos:** Max Rao value for each iteration. If termination condition was met, rest of values after final iteration are NA.

**iterations:** Iteration numbers, corresponding to iter\_raos.

**best\_p:** The best permutation of p found (corresponds to best\_rao).

**Author(s)**

John L. Darcy



## Examples

```
library(specificity)
set.seed(12345)
p <- runif(100)
D <- dist(sample(p))
a <- rao_genetic_max(p,D)
plot(a$iter_raos ~ a$iterations)
```

---

rao\_sort\_max

*rao\_sort\_max*

---

## Description

Sorts pairwise\_product(p) and D to approximate the maximum of Rao(p,D) under permutations of p.

## Usage

```
rao_sort_max(p, D)
```

## Arguments

p                      numeric vector of length n - a species weights vector.  
D                      numeric vector of length  $n(n-1)/2$  - i.e. a dist object whose full matrix is nxn.

## Value

A single value, approximating maximum rao under permutations of p.

## Author(s)

John L. Darcy

## Examples

```
data(endophyte)
p <- prop_abund(endophyte$otutable)[,1]
D <- dist(endophyte$metadata$Elevation)
rsm <- rao_sort_max(p,D)
```

---

tips_from_node	<i>tips_from_node</i>
----------------	-----------------------

---

## Description

Determines which tip indices in a phylogeny descend from a given node. Called by `make_nested_set()`, not intended for use otherwise, but some may find it handy. Data should come from a rooted phylogeny, but this function doesn't check that so be careful.

## Usage

```
tips_from_node(nodes, anc, des)
```

## Arguments

nodes	integer vector or scalar. The node index or indices for which tip indices are desired.
anc	integer vector. "ancestor" column vector from an adjacency matrix. For an <code>ape::phylo</code> object <code>phy</code> , <code>anc=phy\$edge[,1]</code> .
des	integer vector. "descendant" column vector from an adjacency matrix. For an <code>ape::phylo</code> object <code>phy</code> , <code>des=phy\$edge[,2]</code> .

## Value

integer vector of tip indices, in no particular order.

## Author(s)

John L. Darcy

## See Also

`ape::phylo`

## Examples

```
library(specificity)
library(ape)
phy <- get(data(endophyte))$supertree
# check if tree is rooted:
ape::is.rooted(phy)
# which tips are in the Cucurbitales?
plot(phy) # need to stretch out the plot to see...
nodelabels(adj=c(0,-1), bg="yellow") # node numbers
nodelabels(phy$node.label, adj=c(0,1), bg="lightblue") # node names
# we can see that Cucurbitales is node 107
cuc_tips <- tips_from_node( nodes=107, anc=phy$edge[,1], des=phy$edge[,2] )
cuc_tips
```

```
phy$tip.label[cuc_tips]
```

---

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 <code>n_cores = 1</code> . 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 <code>tree\$tip.label</code> . 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**

```
library(specificity)
library(ape)
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. 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*.

## Usage

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

## Arguments

<i>s</i>	character vector. One species name per observation. If no species was observed for a given datum, use NA. <i>s</i> can also be provided as a vector of unique species identities, in which case counts of those species can be given as <i>w</i> .
<i>s_phylo</i>	phylo object. Tree containing all unique names in <i>s</i> as tips. Must not contain duplicate tip labels.
<i>w</i>	numeric vector. Optional weights for <i>s</i> , 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 <i>w</i> is not provided but a weighted metric is specified, <i>w</i> will be set to 1 for each value of <i>s</i> . Thus, weights for each unique species in <i>s</i> would be equal to the number of times that species appears in <i>s</i> . <i>w</i> is not used for unweighted metrics (PD). Any NA values in <i>w</i> will be pairwise removed from <i>w</i> and <i>s</i> (default: NULL).
<i>nested_set</i>	matrix. The output of <code>make_nested_set(s_phylo)</code> . 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

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.

### Examples

```
library(specificity)
set.seed(12345)
s_phylo <- get(data(endophyte))$supertree
w <- sample(c(0, 1), replace=TRUE, size=10)
s <- sample(s_phylo$tip.label, replace=TRUE, size=10)
wpa(s, s_phylo, w, metric="Hp")
```

---

wpa\_table

*wpa\_table*

---

### 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,
  s_names = NULL,
  nested_set = NULL,
  metric = "Hp",
  ncores = 4
)
```

**Arguments**

<code>m</code>	numeric matrix or data.frame of weights, where columns are species and rows are samples.
<code>s_phylo</code>	phylo object. Tree containing all unique names in <code>s</code> as tips. Must not contain duplicate tip labels.
<code>s_names</code>	species names for <code>m</code> if not <code>colnames(m)</code> . <code>NULL</code> will use <code>colnames</code> (default: <code>NULL</code> )
<code>nested_set</code>	matrix. The output of <code>make_nested_set(s_phylo)</code> . If not provided, will be calculated on the fly. Precalculation only provides speedup with very large trees (default: <code>NULL</code> ).
<code>metric</code>	character. Abbreviated name of desired tree-based phylogenetic diversity metric. Available metrics are:  <b>Hp:</b> Phylogenetic Entropy. Insensitive to 0 weights, cannot increase with removal of taxa. Allen et al. 2009. <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. <b>PD:</b> Original Faith's Phylogenetic Diversity. Unweighted. Simply a sum of branch-lengths in your tree (but only for taxa in <code>s</code> ). Faith 1992.
<code>ncores</code>	integer. Number of CPU cores to use for parallel operations (default: 4).

**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.

### Examples

```
library(specificity)
set.seed(12345)
s_phylo <- get(data(endophyte))$supertree
w <- sample(c(0, 1), replace=TRUE, size=10)
nspec <- 12
m <- t(as.matrix(data.frame(
  a=runif(nspec, 0, 100),
  b=runif(nspec, 0, 100),
  c=runif(nspec, 0, 100)
)))
colnames(m) <- sample(s_phylo$tip.label, ncol(m))
wpa_table(m, s_phylo)
```

# Index

## \*Topic **datasets**

endophyte, [8](#)

bl\_distance\_ns, [2](#)

calculate\_spec\_and\_pval, [3](#)

check\_pes\_inputs, [5](#)

circularize2dist, [6](#)

distcalc, [7](#)

endophyte, [8](#)

env\_spec\_sim, [8](#)

geo\_spec\_sim, [10](#)

get\_ga\_defaults, [12](#)

make\_nested\_set, [13](#)

occ\_threshold, [14](#)

onto2nwk, [15](#)

pairwise\_product, [16](#)

phy\_or\_env\_spec, [17](#)

phy\_spec\_sim, [20](#)

plot\_grid\_abunds, [22](#)

plot\_pairwise\_spec, [23](#)

plot\_specificities (plot\_specs\_stacks),  
[24](#)

plot\_specs\_stacks, [24](#)

plot\_specs\_violin, [26](#)

prop\_abund, [27](#)

random\_rep\_positions, [29](#)

randomgrid, [28](#)

rao1sp, [30](#)

rao\_genetic\_max, [32](#)

rao\_sort\_max, [33](#)

raoperms, [31](#)

tips\_from\_node, [34](#)

tree2mat, [35](#)

wpd, [36](#)

wpd\_table, [37](#)