

Package ‘ecodive’

September 21, 2025

Type Package

Title Parallel and Memory-Efficient Ecological Diversity Metrics

Version 2.0.0

Date 2025-09-20

Description Computes alpha and beta diversity metrics using concurrent 'C' threads.

Metrics include 'UniFrac', Faith's phylogenetic diversity, Bray-Curtis dissimilarity, Shannon diversity index, and many others.

Also parses newick trees into 'phylo' objects and rarefies feature tables.

URL <https://cmmr.github.io/ecodive/>, <https://github.com/cmmr/ecodive>

BugReports <https://github.com/cmmr/ecodive/issues>

License MIT + file LICENSE

Encoding UTF-8

LazyData true

Depends R (>= 3.6.0)

RoxygenNote 7.3.3

Config/Needs/website rmarkdown

Config/testthat/edition 3

Imports parallel, utils

Suggests knitr, parallelly, rmarkdown, testthat (>= 3.0.0)

VignetteBuilder knitr

NeedsCompilation yes

Author Daniel P. Smith [aut, cre] (ORCID:

[<https://orcid.org/0000-0002-2479-2044>](https://orcid.org/0000-0002-2479-2044)),

Alkek Center for Metagenomics and Microbiome Research [cph, fnd]

Maintainer Daniel P. Smith <dansmith01@gmail.com>

Repository CRAN

Date/Publication 2025-09-21 19:50:02 UTC

Contents

<i>adiv_functions</i>	2
<i>alpha_div</i>	5
<i>bdiv_functions</i>	6
<i>beta_div</i>	11
<i>ex_counts</i>	13
<i>ex_tree</i>	14
<i>list_metrics</i>	14
<i>n_cpus</i>	17
<i>rarefy</i>	17
<i>read_tree</i>	18

Index

20

<i>adiv_functions</i>	<i>Alpha Diversity Metrics</i>
-----------------------	--------------------------------

Description

Alpha Diversity Metrics

Usage

```
ace(counts, cutoff = 10, cpus = n_cpus())
berger(counts, rescale = TRUE, cpus = n_cpus())
brillouin(counts, cpus = n_cpus())
chao1(counts, cpus = n_cpus())
faith(counts, tree = NULL, cpus = n_cpus())
fisher(counts, digits = 3L, cpus = n_cpus())
inv_simpson(counts, rescale = TRUE, cpus = n_cpus())
margalef(counts, cpus = n_cpus())
mcintosh(counts, cpus = n_cpus())
menhinick(counts, cpus = n_cpus())
observed(counts, cpus = n_cpus())
shannon(counts, rescale = TRUE, cpus = n_cpus())
```

```

simpson(counts, rescale = TRUE, cpus = n_cpus())
squares(counts, cpus = n_cpus())

```

Arguments

counts	A numeric matrix of count data where each column is a feature, and each row is a sample. Any object coercible with <code>as.matrix()</code> can be given here, as well as <code>phyloseq</code> , <code>rbiom</code> , <code>SummarizedExperiment</code> , and <code>TreeSummarizedExperiment</code> objects.
cutoff	The maximum number of observations to consider "rare". Default: 10.
cpus	How many parallel processing threads should be used. The default, <code>n_cpus()</code> , will use all logical CPU cores.
rescale	Normalize each sample's counts so they sum to 1. Default: TRUE
tree	A phylo-class object representing the phylogenetic tree for the OTUs in <code>counts</code> . The OTU identifiers given by <code>colnames(counts)</code> must be present in <code>tree</code> . Can be omitted if a tree is embedded with the <code>counts</code> object or as <code>attr(counts, 'tree')</code> .
digits	Precision of the returned values, in number of decimal places. E.g. the default <code>digits=3</code> could return 6.392.

Value

A numeric vector.

Formulas

Prerequisite: all counts are whole numbers.

Given:

- n : The number of features (e.g. species, OTUs, ASVs, etc).
- X_i : Integer count of the i -th feature.
- X_T : Total of all counts (i.e. sequencing depth). $X_T = \sum_{i=1}^n X_i$
- P_i : Proportional abundance of the i -th feature. $P_i = X_i/X_T$
- F_1 : Number of features where $X_i = 1$ (i.e. singletons).
- F_2 : Number of features where $X_i = 2$ (i.e. doubletons).

Abundance-based Coverage Estimator (ACE)	<code>ace()</code>	See below.
Berger-Parker Index	<code>berger()</code>	$\max(P_i)$
Brillouin Index	<code>brillouin()</code>	$\frac{\ln [(\sum_{i=1}^n X_i)!] - \sum_{i=1}^n \ln (X_i!)}{\sum_{i=1}^n X_i}$
Chao1	<code>chao1()</code>	$n + \frac{(F_1)^2}{2F_2}$
Faith's Phylogenetic Diversity	<code>faith()</code>	See below.

Fisher's Alpha (α) <code>fisher()</code>	$\frac{n}{\alpha} = \ln \left(1 + \frac{X_T}{\alpha} \right)$	The value of α must be solved for iteratively.
Gini-Simpson Index <code>simpson()</code>	$1 - \sum_{i=1}^n P_i^2$	
Inverse Simpson Index <code>inv_simpson()</code>	$1 / \sum_{i=1}^n P_i^2$	
Margalef's Richness Index <code>margalef()</code>	$\frac{n-1}{\ln X_T}$	
McIntosh Index <code>mcintosh()</code>	$\frac{X_T - \sqrt{\sum_{i=1}^n (X_i)^2}}{X_T - \sqrt{X_T}}$	
Menhinick's Richness Index <code>menhinick()</code>	$\frac{n}{\sqrt{X_T}}$	
Observed Features <code>observed()</code>	n	
Shannon Diversity Index <code>shannon()</code>	$-\sum_{i=1}^n P_i \times \ln(P_i)$	
Squares Richness Estimator <code>squares()</code>	$n + \frac{(F_1)^2 \sum_{i=1}^n (X_i)^2}{X_T^2 - nF_1}$	

Abundance-based Coverage Estimator (ACE):

Given:

- n : The number of features (e.g. species, OTUs, ASVs, etc).
- r : Rare cutoff. Features with $\leq r$ counts are considered rare.
- X_i : Integer count of the i -th feature.
- F_i : Number of features with exactly i counts.
- F_1 : Number of features where $X_i = 1$ (i.e. singletons).
- F_{rare} : Number of rare features where $X_i \leq r$.
- F_{abund} : Number of abundant features where $X_i > r$.
- X_{rare} : Total counts belonging to rare features.
- C_{ace} : The sample abundance coverage estimator, defined below.
- γ_{ace}^2 : The estimated coefficient of variation, defined below.
- D_{ace} : Estimated number of features in the sample.

$$C_{ace} = 1 - \frac{F_1}{X_{rare}}$$

$$\gamma_{ace}^2 = \max \left[\frac{F_{rare} \sum_{i=1}^r i(i-1)F_i}{C_{ace} X_{rare}(X_{rare} - 1)} - 1, 0 \right]$$

$$D_{ace} = F_{abund} + \frac{F_{rare}}{C_{ace}} + \frac{F_1}{C_{ace}} \gamma_{ace}^2$$

Faith's Phylogenetic Diversity (Faith's PD):

Given n branches with lengths L and a sample's abundances A on each of those branches coded as 1 for present or 0 for absent:

$$\sum_{i=1}^n L_i A_i$$

Examples

```
# Example counts matrix
t(ex_counts)

ace(ex_counts)
```

```
chao1(ex_counts)
squares(ex_counts)
```

alpha_div*Alpha Diversity Wrapper Function***Description**

Alpha Diversity Wrapper Function

Usage

```
alpha_div(counts, metric, ...)
```

Arguments

counts	A numeric matrix of count data where each column is a feature, and each row is a sample. Any object coercible with <code>as.matrix()</code> can be given here, as well as <code>phyloseq</code> , <code>rbiom</code> , <code>SummarizedExperiment</code> , and <code>TreeSummarizedExperiment</code> objects.
metric	The name of an alpha diversity metric. One of <code>c('ace', 'berger', 'brillouin', 'chao1', 'faith', 'fisher', 'inv_simpson', 'margalef', 'mcintosh', 'menhinick', 'observed', 'shannon', 'simpson', 'squares')</code> . Case-insensitive and partial name matching is supported. Programmatic access via <code>list_metrics('alpha')</code> .
...	Additional options to pass through to the called function. I.e. <code>cpus</code> or <code>tree</code> .

Details**Integer Count Requirements:**

A frequent and critical error in alpha diversity analysis is providing the wrong type of data to a metric's formula. Some indices are mathematically defined based on counts of individuals and require raw, integer abundance data. Others are based on proportional abundances and can accept either integer counts (which are then converted to proportions) or pre-normalized proportional data. Using proportional data with a metric that requires integer counts will return an error message.

Requires Integer Counts Only	Can Use Proportional Data
Chao1	Observed Features
ACE	Shannon Index
Squares Richness Estimator	Gini-Simpson Index
Margalef's Index	Inverse Simpson Index
Menhinick's Index	Berger-Parker Index
Fisher's Alpha	McIntosh Index
Brillouin Index	Faith's PD (presence/absence)

Value

A numeric vector.

Examples

```
# Example counts matrix
ex_counts

# Shannon diversity values
alpha_div(ex_counts, 'Shannon')

# Chao1 diversity values
alpha_div(ex_counts, 'c')

# Faith PD values
alpha_div(ex_counts, 'faith', tree = ex_tree)
```

Description

Beta Diversity Metrics

Usage

```
aitchison(counts, pseudocount = NULL, pairs = NULL, cpus = n_cpus())
bhattacharyya(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())
bray(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())
canberra(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())
chebyshev(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())
chord(counts, pairs = NULL, cpus = n_cpus())
clark(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())
divergence(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())
euclidean(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())
```

```
gower(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())  
hellinger(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())  
horn(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())  
jensen(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())  
jsd(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())  
lorentzian(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())  
manhattan(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())  
matusita(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())  
minkowski(counts, rescale = TRUE, power = 1.5, pairs = NULL, cpus = n_cpus())  
morisita(counts, pairs = NULL, cpus = n_cpus())  
motyka(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())  
psym_chisq(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())  
soergel(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())  
squared_chisq(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())  
squared_chord(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())  
squared_euclidean(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())  
topsoe(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())  
wave_hedges(counts, rescale = TRUE, pairs = NULL, cpus = n_cpus())  
hamming(counts, pairs = NULL, cpus = n_cpus())  
jaccard(counts, pairs = NULL, cpus = n_cpus())  
ochiai(counts, pairs = NULL, cpus = n_cpus())  
sorensen(counts, pairs = NULL, cpus = n_cpus())  
unweighted_unifrac(counts, tree = NULL, pairs = NULL, cpus = n_cpus())  
weighted_unifrac(counts, tree = NULL, pairs = NULL, cpus = n_cpus())
```

```

normalized_unifrac(counts, tree = NULL, pairs = NULL, cpus = n_cpus())

generalized_unifrac(
  counts,
  tree = NULL,
  alpha = 0.5,
  pairs = NULL,
  cpus = n_cpus()
)

variance_adjusted_unifrac(counts, tree = NULL, pairs = NULL, cpus = n_cpus())

```

Arguments

counts	A numeric matrix of count data where each column is a feature, and each row is a sample. Any object coercible with <code>as.matrix()</code> can be given here, as well as <code>phyloseq</code> , <code>rbiom</code> , <code>SummarizedExperiment</code> , and <code>TreeSummarizedExperiment</code> objects.
pseudocount	The value to add to all counts in <code>counts</code> to prevent taking $\log(0)$ for unobserved features. The default, <code>NULL</code> , selects the smallest non-zero value in <code>counts</code> .
pairs	Which combinations of samples should distances be calculated for? The default value (<code>NULL</code>) calculates all-vs-all. Provide a numeric or logical vector specifying positions in the distance matrix to calculate. See examples.
cpus	How many parallel processing threads should be used. The default, <code>n_cpus()</code> , will use all logical CPU cores.
rescale	Normalize each sample's counts so they sum to 1. Default: <code>TRUE</code>
power	Scaling factor for the magnitude of differences between communities (p). Default: 1.5
tree	A <code>phylo</code> -class object representing the phylogenetic tree for the OTUs in <code>counts</code> . The OTU identifiers given by <code>colnames(counts)</code> must be present in <code>tree</code> . Can be omitted if a tree is embedded with the <code>counts</code> object or as <code>attr(counts, 'tree')</code> .
alpha	How much weight to give to relative abundances; a value between 0 and 1, inclusive. Setting <code>alpha=1</code> is equivalent to <code>normalized_unifrac()</code> .

Value

A `dist` object.

Formulas

Given:

- n : The number of features.
- X_i, Y_i : Absolute counts for the i -th feature in samples X and Y .
- X_T, Y_T : Total counts in each sample. $X_T = \sum_{i=1}^n X_i$

- P_i, Q_i : Proportional abundances of X_i and Y_i . $P_i = X_i/X_T$
- X_L, Y_L : Mean log of abundances. $X_L = \frac{1}{n} \sum_{i=1}^n \ln X_i$
- R_i : The range of the i -th feature across all samples (max - min).

Aitchison distance aitchison()	$\sqrt{\frac{\sum_{i=1}^n [(\ln X_i - X_L) - (\ln Y_i - Y_L)]^2}{-\ln \sum_{i=1}^n \sqrt{P_i Q_i}}}$
Bhattacharyya distance bhattacharyya()	$\frac{\sum_{i=1}^n P_i - Q_i }{\sum_{i=1}^n (P_i + Q_i)}$
Bray-Curtis dissimilarity bray()	$\frac{\sum_{i=1}^n P_i - Q_i }{\sum_{i=1}^n (P_i + Q_i)}$
Canberra distance canberra()	$\frac{\sum_{i=1}^n P_i - Q_i }{\sum_{i=1}^n (P_i + Q_i)}$
Chebyshev distance chebyshev()	$\max(P_i - Q_i)$
Chord distance chord()	$\sqrt{\sum_{i=1}^n \left(\frac{X_i}{\sqrt{\sum_{j=1}^n X_j^2}} - \frac{Y_i}{\sqrt{\sum_{j=1}^n Y_j^2}} \right)^2}$
Clark's divergence distance clark()	$\sqrt{\sum_{i=1}^n \left(\frac{P_i - Q_i}{P_i + Q_i} \right)^2}$
Divergence divergence()	$2 \sum_{i=1}^n \frac{(P_i - Q_i)^2}{(P_i + Q_i)^2}$
Euclidean distance euclidean()	$\sqrt{\sum_{i=1}^n (P_i - Q_i)^2}$
Gower distance gower()	$\frac{1}{n} \sum_{i=1}^n \frac{ P_i - Q_i }{R_i}$
Hellinger distance hellinger()	$\sqrt{\sum_{i=1}^n (\sqrt{P_i} - \sqrt{Q_i})^2}$
Horn-Morisita dissimilarity horn()	$1 - \frac{2 \sum_{i=1}^n P_i Q_i}{\sum_{i=1}^n P_i^2 + \sum_{i=1}^n Q_i^2}$
Jensen-Shannon distance jensen()	$\sqrt{\frac{1}{2} \left[\sum_{i=1}^n P_i \ln \left(\frac{2P_i}{P_i + Q_i} \right) + \sum_{i=1}^n Q_i \ln \left(\frac{2Q_i}{P_i + Q_i} \right) \right]}$
Jensen-Shannon divergence (JSD) jsd()	$\frac{1}{2} \left[\sum_{i=1}^n P_i \ln \left(\frac{2P_i}{P_i + Q_i} \right) + \sum_{i=1}^n Q_i \ln \left(\frac{2Q_i}{P_i + Q_i} \right) \right]$
Lorentzian distance lorentzian()	$\sum_{i=1}^n \ln(1 + P_i - Q_i)$
Manhattan distance manhattan()	$\sum_{i=1}^n P_i - Q_i $
Matusita distance matusita()	$\sqrt{\sum_{i=1}^n (\sqrt{P_i} - \sqrt{Q_i})^2}$
Minkowski distance minkowski()	$\sqrt[p]{\sum_{i=1}^n (P_i - Q_i)^p}$ Where p is the geometry of the space.
Morisita dissimilarity * Integers Only morisita()	$1 - \frac{2 \sum_{i=1}^n X_i Y_i}{\left(\frac{\sum_{i=1}^n X_i (X_i - 1)}{X_T (X_T - 1)} + \frac{\sum_{i=1}^n Y_i (Y_i - 1)}{Y_T (Y_T - 1)} \right) X_T Y_T}$
Motyka dissimilarity motyka()	$\frac{\sum_{i=1}^n \max(P_i, Q_i)}{\sum_{i=1}^n (P_i + Q_i)}$
Probabilistic Symmetric χ^2 distance psym_chisq()	$2 \sum_{i=1}^n \frac{(P_i - Q_i)^2}{P_i + Q_i}$

Soergel distance soergel()	$\frac{\sum_{i=1}^n P_i - Q_i }{\sum_{i=1}^n \max(P_i, Q_i)}$
Squared χ^2 distance squared_chisq()	$\sum_{i=1}^n \frac{(P_i - Q_i)^2}{P_i + Q_i}$
Squared Chord distance squared_chord()	$\sum_{i=1}^n \frac{(\sqrt{P_i} - \sqrt{Q_i})^2}{(P_i - Q_i)^2}$
Squared Euclidean distance squared_euclidean()	$\sum_{i=1}^n (P_i - Q_i)^2$
Topsoe distance topsoe()	$\sum_{i=1}^n P_i \ln \left(\frac{2P_i}{P_i + Q_i} \right) + \sum_{i=1}^n Q_i \ln \left(\frac{2Q_i}{P_i + Q_i} \right)$
Wave Hedges distance wave_hedges()	$\frac{\sum_{i=1}^n P_i - Q_i }{\sum_{i=1}^n \max(P_i, Q_i)}$

Presence / Absence:

Given:

- A, B : Number of features in each sample.
- J : Number of features in common.

Dice-Sorensen dissimilarity sorensen()	$\frac{2J}{(A + B)}$
Hamming distance hamming()	$(A + B) - 2J$
Jaccard distance jaccard()	$1 - \frac{J}{(A + B - J)}$
Otsuka-Ochiai dissimilarity ochiai()	$1 - \frac{J}{\sqrt{AB}}$

Phylogenetic:

Given n branches with lengths L and a pair of samples' binary (A and B) or proportional abundances (P and Q) on each of those branches.

Unweighted UniFrac unweighted_unifrac()	$\frac{1}{n} \sum_{i=1}^n L_i A_i - B_i $
Weighted UniFrac weighted_unifrac()	$\sum_{i=1}^n L_i P_i - Q_i $
Normalized Weighted UniFrac normalized_unifrac()	$\frac{\sum_{i=1}^n L_i P_i - Q_i }{\sum_{i=1}^n L_i (P_i + Q_i)}$
Generalized UniFrac (GUniFrac) generalized_unifrac()	$\frac{\sum_{i=1}^n L_i (P_i + Q_i)^\alpha \left \frac{P_i - Q_i}{P_i + Q_i} \right }{\sum_{i=1}^n L_i (P_i + Q_i)^\alpha}$ Where α is a scalar
Variance-Adjusted Weighted UniFrac variance_adjusted_unifrac()	$\frac{\sum_{i=1}^n L_i \frac{ P_i - Q_i }{\sqrt{(P_i + Q_i)(2 - P_i - Q_i)}}}{\sum_{i=1}^n L_i \frac{P_i + Q_i}{\sqrt{(P_i + Q_i)(2 - P_i - Q_i)}}}$

See `vignette('unifrac')` for detailed example UniFrac calculations.

References

- Levy, A., Shalom, B. R., & Chalamish, M. (2024). A guide to similarity measures. *arXiv*.
- Cha, S.-H. (2007). Comprehensive survey on distance/similarity measures between probability density functions. *International Journal of Mathematical Models and Methods in Applied Sciences*, 1(4), 300–307.

Examples

```
# Example counts matrix
t(ex_counts)

bray(ex_counts)

jaccard(ex_counts)

generalized_unifrac(ex_counts, tree = ex_tree)

# Only calculate distances for Saliva vs all.
bray(ex_counts, pairs = 1:3)
```

`beta_div`

Beta Diversity Wrapper Function

Description

Beta Diversity Wrapper Function

Usage

```
beta_div(counts, metric, ...)
```

Arguments

<code>counts</code>	A numeric matrix of count data where each column is a feature, and each row is a sample. Any object coercible with <code>as.matrix()</code> can be given here, as well as <code>phyloseq</code> , <code>rbiom</code> , <code>SummarizedExperiment</code> , and <code>TreeSummarizedExperiment</code> objects.
<code>metric</code>	The name of a beta diversity metric. One of c('aitchison', 'bhattacharyya', 'bray', 'canberra', 'chebyshev', 'chord', 'clark', 'divergence', 'euclidean', 'generalized_unifrac', 'gower', 'hamming', 'hellinger', 'horn', 'jaccard', 'jensen', 'jsd', 'lorentzian', 'manhattan', 'matusita', 'minkowski', 'morisita', 'motska', 'normalized_unifrac', 'ochiai', 'psym_chisq',

'soergel', 'sorensen', 'squared_chisq', 'squared_chord', 'squared_euclidean', 'topsoe', 'unweighted_unifrac', 'variance_adjusted_unifrac', 'wave_hedges', 'weighted_unifrac'). Flexible matching is supported (see below). Programmatic access via `list_metrics('beta')`.

... Additional options to pass through to the called function. I.e. `tree`, `pairs`, `alpha`, or `cpus`.

Details

List of Beta Diversity Metrics

Option / Function Name	Metric Name
aitchison	Aitchison distance
bhattacharyya	Bhattacharyya distance
bray	Bray-Curtis dissimilarity
canberra	Canberra distance
chebyshev	Chebyshev distance
chord	Chord distance
clark	Clark's divergence distance
divergence	Divergence
euclidean	Euclidean distance
generalized_unifrac	Generalized UniFrac (GUniFrac)
gower	Gower distance
hamming	Hamming distance
hellinger	Hellinger distance
horn	Horn-Morisita dissimilarity
jaccard	Jaccard distance
jensen	Jensen-Shannon distance
jsd	Jesen-Shannon divergence (JSD)
lorentzian	Lorentzian distance
manhattan	Manhattan distance
matusita	Matusita distance
minkowski	Minkowski distance
morisita	Morisita dissimilarity
motyka	Motyka dissimilarity
normalized_unifrac	Normalized Weighted UniFrac
ochiai	Otsuka-Ochiai dissimilarity
psym_chisq	Probabilistic Symmetric Chi-Squared distance
soergel	Soergel distance
sorensen	Dice-Sorensen dissimilarity
squared_chisq	Squared Chi-Squared distance
squared_chord	Squared Chord distance
squared_euclidean	Squared Euclidean distance
topsoe	Topsoe distance
unweighted_unifrac	Unweighted UniFrac
variance_adjusted_unifrac	Variance-Adjusted Weighted UniFrac (VAW-UniFrac)
wave_hedges	Wave Hedges distance
weighted_unifrac	Weighted UniFrac

Flexible name matching

Case insensitive and partial matching. Any runs of non-alpha characters are converted to underscores. E.g. metric = 'Weighted UniFrac' selects weighted_unifrac.

UniFrac names can be shortened to the first letter plus "unifrac". E.g. uunifrac, w_unifrac, or V UniFrac. These also support partial matching.

Finished code should always use the full primary option name to avoid ambiguity with future additions to the metrics list.

Value

A numeric vector.

Examples

```
# Example counts matrix
ex_counts

# Bray-Curtis distances
beta_div(ex_counts, 'bray')

# Generalized UniFrac distances
beta_div(ex_counts, 'GUniFrac', tree = ex_tree)
```

ex_counts

Example counts matrix

Description

Genera found on four human body sites.

Usage

ex_counts

Format

A matrix of 4 samples (columns) x 6 genera (rows).

Source

Derived from The Human Microbiome Project dataset. <https://commonfund.nih.gov/hmp>

ex_tree	<i>Example phylogenetic tree</i>
---------	----------------------------------

Description

Companion tree for ex_counts.

Usage

ex_tree

Format

A phylo object.

Details

ex_tree encodes this tree structure:

```

+-----44----- Haemophilus
+-2-|
|   +-----68----- Bacteroides
|
|   +---18---- Streptococcus
|   +-12--|
|   |   +-11-- Staphylococcus
+-11--|
|   +----24----- Corynebacterium
+-12--|
|   +-13-- Propionibacterium

```

list_metrics	<i>Find and Browse Available Metrics</i>
--------------	--

Description

Programmatic access to the lists of available metrics, and their associated functions.

Usage

```
list_metrics(
  div = c(NA, "alpha", "beta"),
  val = c("data.frame", "list", "func", "id", "name", "div", "phylo", "weighted",
         "int_only", "true_metric"),
  nm = c(NA, "id", "name"),
  phylo = NULL,
  weighted = NULL,
  int_only = NULL,
  true_metric = NULL
)

match_metric(
  metric,
  div = NULL,
  phylo = NULL,
  weighted = NULL,
  int_only = NULL,
  true_metric = NULL
)
```

Arguments

<code>div, phylo, weighted, int_only, true_metric</code>	Consider only metrics matching specific criteria. For example, <code>div = "alpha"</code> will only return alpha diversity metrics. Default: <code>NULL</code>
<code>val</code>	Sets the return value for this function call. See "Value" section below. Default: <code>"data.frame"</code>
<code>nm</code>	What value to use for the names of the returned object. Default is <code>"id"</code> when <code>val</code> is <code>"list"</code> or <code>"func"</code> , otherwise the default is <code>NA</code> (no name).
<code>metric</code>	The name of an alpha/beta diversity metric to search for. Supports partial matching. All non-alpha characters are ignored.

Value

`match_metric()`

A list with the following elements.

- `name` : Metric name, e.g. "Faith's Phylogenetic Diversity"
- `id` : Metric ID - also the name of the function, e.g. "faith"
- `div` : Either `"alpha"` or `"beta"`.
- `phylo` : TRUE if metric requires a phylogenetic tree; FALSE otherwise.
- `weighted` : TRUE if metric takes relative abundance into account; FALSE if it only uses presence/absence.
- `int_only` : TRUE if metric requires integer counts; FALSE otherwise.

- `true_metric` : TRUE if metric is a true metric and satisfies the triangle inequality; FALSE if it is a non-metric dissimilarity; NA for alpha diversity metrics.
- `func` : The function for this metric, e.g. `ecodive::faith`
- `params` : Formal args for `func`, e.g. `c("counts", "tree", "cpus")`

`list_metrics()`

The returned object's type and values are controlled with the `val` and `nm` arguments.

- `val = "data.frame"` : The data.frame from which the below options are sourced.
- `val = "list"` : A list of objects as returned by `match_metric()` (above).
- `val = "func"` : A list of functions.
- `val = "id"` : A character vector of metric IDs.
- `val = "name"` : A character vector of metric names.
- `val = "div"` : A character vector "alpha" and/or "beta".
- `val = "phylo"` : A logical vector indicating which metrics require a tree.
- `val = "weighted"` : A logical vector indicating which metrics take relative abundance into account (as opposed to just presence/absence).
- `val = "int_only"` : A logical vector indicating which metrics require integer counts.
- `val = "true_metric"` : A logical vector indicating which metrics are true metrics and satisfy the triangle inequality, which work better for ordinations such as PCoA.

If `nm` is set, then the names of the vector or list will be the metric ID (`nm="id"`) or name (`nm="name"`). When `val="data.frame"`, the names will be applied to the `rownames()` property of the `data.table`.

Examples

```
# A data.frame of all available metrics.
head(list_metrics())

# All alpha diversity function names.
list_metrics('alpha', val = 'id')

# Try to find a metric named 'otus'.
m <- match_metric('otus')

# The result is a list that includes the function.
str(m)
```

n_cpus	<i>Number of CPU Cores</i>
--------	----------------------------

Description

A thin wrapper around `parallel::availableCores()`. If the `parallel` package is not installed, then it falls back to `parallel::detectCores(all.tests = TRUE, logical = TRUE)`. Returns 1 if `pthread` support is unavailable or when the number of cpus cannot be determined.

Usage

```
n_cpus()
```

Value

A scalar integer, guaranteed to be at least 1.

Examples

```
n_cpus()
```

rarefy	<i>Rarefy OTU counts.</i>
--------	---------------------------

Description

Sub-sample OTU observations such that all samples have an equal number. If called on data with non-integer abundances, values will be re-scaled to integers between 1 and depth such that they sum to depth.

Usage

```
rarefy(
  counts,
  depth = 0.1,
  n_samples = NULL,
  seed = 0,
  times = NULL,
  cpus = n_cpus()
)
```

Arguments

<code>counts</code>	A numeric matrix of count data where each column is a feature, and each row is a sample. Any object coercible with <code>as.matrix()</code> can be given here, as well as <code>phyloseq</code> , <code>rbiom</code> , <code>SummarizedExperiment</code> , and <code>TreeSummarizedExperiment</code> objects.
<code>depth</code>	How many observations to keep per sample. When $0 < \text{depth} < 1$, it is taken as the minimum percentage of the dataset's observations to keep. Ignored when <code>n_samples</code> is specified. Default: <code>0.1</code>
<code>n_samples</code>	The number of samples to keep. When $0 < \text{n_samples} < 1$, it is taken as the percentage of samples to keep. If negative, that number of samples is dropped. If <code>0</code> , all samples are kept. If <code>NULL</code> , then <code>depth</code> is used instead. Default: <code>NULL</code>
<code>seed</code>	An integer seed for randomizing which observations to keep or drop. If you need to create different random rarefactions of the same data, set the seed to a different number each time.
<code>times</code>	How many independent rarefactions to perform. If set, <code>rarefy()</code> will return a list of matrices. The seeds for each matrix will be sequential, starting from <code>seed</code> .
<code>cpus</code>	How many parallel processing threads should be used. The default, <code>n_cpus()</code> , will use all logical CPU cores.

Value

An integer matrix.

Examples

```
# Create an OTU matrix with 4 samples (A-D) and 5 OTUs.
counts <- matrix(
  data      = c(4,0,3,2,6,0,8,0,0,5,0,9,0,0,7,0,10,0,0,1),
  nrow      = 5,
  dimnames = list(paste0('OTU', 1:5), LETTERS[1:4]) )
counts
colSums(counts)

counts <- rarefy(counts, depth = 14)
counts
colSums(counts)
```

`read_tree`

Read a newick formatted phylogenetic tree.

Description

A phylogenetic tree is required for computing UniFrac distance matrices. You can load a tree from a file or by providing the tree string directly. This tree must be in Newick format, also known as parenthetic format and New Hampshire format.

Usage

```
read_tree(newick, underscores = FALSE)
```

Arguments

- newick Input data as either a file path, URL, or Newick string. Compressed (gzip or bzip2) files are also supported.
- underscores If TRUE, underscores in unquoted names will remain underscores. If FALSE, underscores in unquoted named will be converted to spaces.

Value

A phylo class object representing the tree.

Examples

```
tree <- read_tree("((A:0.99,((B:0.87,C:0.89):0.51,(((D:0.16,(E:0.83,F:0.96):0.94):0.69,(G:0.92,(H:0.62,I:0.85):0.54):0.23):0.74,J:0.12):0.43):0.67);")  
class(tree)
```

Index

* **adiv_functions**
 adiv_functions, 2
* **bdiv_functions**
 bdiv_functions, 6
* **datasets**
 ex_counts, 13
 ex_tree, 14

ace (adiv_functions), 2
adiv_functions, 2
aitchison (bdiv_functions), 6
alpha_div, 5

bdiv_functions, 6
berger (adiv_functions), 2
beta_div, 11
bhattacharyya (bdiv_functions), 6
bray (bdiv_functions), 6
brillouin (adiv_functions), 2

canberra (bdiv_functions), 6
chao1 (adiv_functions), 2
chebyshev (bdiv_functions), 6
chord (bdiv_functions), 6
clark (bdiv_functions), 6

divergence (bdiv_functions), 6

euclidean (bdiv_functions), 6
ex_counts, 13
ex_tree, 14

faith (adiv_functions), 2
fisher (adiv_functions), 2

generalized_unifrac (bdiv_functions), 6
gower (bdiv_functions), 6

hamming (bdiv_functions), 6
hellinger (bdiv_functions), 6
horn (bdiv_functions), 6

inv_simpson (adiv_functions), 2

jaccard (bdiv_functions), 6
jensen (bdiv_functions), 6
jsd (bdiv_functions), 6

list_metrics, 14
lorentzian (bdiv_functions), 6

manhattan (bdiv_functions), 6
margalef (adiv_functions), 2
match_metric (list_metrics), 14
matusita (bdiv_functions), 6
mcintosh (adiv_functions), 2
menhinick (adiv_functions), 2
minkowski (bdiv_functions), 6
morisita (bdiv_functions), 6
motyka (bdiv_functions), 6

n_cpus, 17
normalized_unifrac (bdiv_functions), 6

observed (adiv_functions), 2
ochiai (bdiv_functions), 6

psym_chisq (bdiv_functions), 6

rarefy, 17
read_tree, 18

shannon (adiv_functions), 2
simpson (adiv_functions), 2
soergel (bdiv_functions), 6
sorensen (bdiv_functions), 6
squared_chisq (bdiv_functions), 6
squared_chord (bdiv_functions), 6
squared_euclidean (bdiv_functions), 6
squares (adiv_functions), 2

topsoe (bdiv_functions), 6

unweighted_unifrac (bdiv_functions), 6

`variance_adjusted_unifrac
(bdiv_functions), 6`

`wave_hedges (bdiv_functions), 6
weighted_unifrac (bdiv_functions), 6`