

```
## Warning in fun(libname, pkgname): Package 'ccrepe' is deprecated and
will be removed from Bioconductor
## version 3.24
```

CCREPE: Compositionality Corrected by PERmutation and RENormalization

Emma Schwager, George Weingart, Craig Bielski, Curtis Huttenhower

May 29, 2026

Contents

1	Introduction	1
2	ccrepe	1
2.1	General functionality	1
2.2	Arguments	2
2.3	Output	3
2.4	Usage	4
2.5	Example 1	4
2.6	Example 2	6
2.7	Example 3	8
2.8	Example 4	10
2.9	Example 5	11
3	nc.score	12
3.1	General Functionality	12
3.2	Arguments	12
3.3	Output	13
3.4	Usage	13
3.5	Example 1	13
3.6	Example 2	14
3.7	Example 3	15
4	References	16

1 Introduction

ccrepe is a package for analysis of sparse compositional data. Specifically, it determines the significance of association between features in a composition, using any similarity measure (e.g. Pearson correlation, Spearman correlation, etc.) The CCREPE methodology stands for Compositionality Corrected by Renormalization and Permutation, as detailed below. The package also provides a novel similarity measure, the N-dimensional checkerboard score (NC

CCREPE: Compositionality Corrected by PERmutation and RENormalization

over sample subsets in order to assess confidence in the "true" similarity measure. Finally, the two resulting distributions are compared using a pooled-variance Z-test to give a compositionality-corrected p-value. False discovery rate q-values are additionally calculated using the Benjamin-Hochberg-Yekutieli procedure. For greater detail, see [Faust et al. \[2012\]](#) and [Schwager and Colleagues](#).

CCREPE employs several filtering steps before the data are processed. It removes any missing subjects using `na.omit`: in the two dataset case, any subjects missing in *either* dataset will be removed. Any subjects or features which are all zero are removed as well: an all-zero subject cannot be normalized (its sum is 0) and an all-zero feature has standard deviation 0 (in addition to being uninteresting biologically).

2.2 Arguments

x First *dataframe* or *matrix* containing relative abundances. Columns are features, rows are samples. Rows should therefore sum to a constant. Row names are used for identification if present.

y Second *dataframe* or *matrix* (optional) containing relative abundances. Columns are features, rows are samples. Rows should therefore sum to a constant. If both **x** and **y** are specified, they will be merged by row names. If no row names are specified for either or both datasets, the default is to merge by row number.

sim.score Similarity measure, such as `cor` or `nc.score`. This can be either an existing R function or user-defined. If the latter, certain properties should be satisfied as detailed below (also see examples). The default similarity measure is Spearman correlation.

A user-defined similarity measure should mimic the interface of `cor`:

1. Take either two *vector* inputs one *matrix* or *dataframe* input.
2. In the case of two inputs, return a single number.
3. In the case of one input, return a matrix in which the (i,j) th entry is the similarity score for column **i** and column **j** in the original matrix.
4. The resulting matrix (in the case of one input) must be symmetric.
5. The inputs must be named **x** and **y**.

sim.score.args An optional list of arguments for the measurement function. When given, they are passed to the `sim.score` function directly. For example, in the case of `cor`, the following would be acceptable:

```
sim.score.args = list(method="spearman", use="complete.obs")
```

min.subj Minimum number (count) of samples that must be non-missing in order to apply the similarity measure. This is to ensure that there are sufficient samples to perform a bootstrap (default: 20).

iterations The number of iterations for both bootstrap and permutation calculations (default: 1000).

CCREPE: Compositionality Corrected by PERmutation and RENormalization

subset.cols.x A vector of column indices from *x* to indicate which features to compare

subset.cols.y A vector of column indices from *y* to indicate which features to compare

errthresh If feature has number of zeros greater than $errthresh^{1/n}$, that feature is excluded

verbose If TRUE, print periodic progress of the algorithm through the dataset(s), as well as including more detailed debugging output. (default: FALSE).

iterations.gap If **verbose=TRUE**, the number of iterations between issuing status messages (default: 100).

distributions Optional output file for detailed log (if given) of all intermediate permutation and renormalization distributions.

compare.within.x A boolean value indicating whether to do comparisons given by taking all subsets of size 2 from **subset.cols.x** or to do comparisons given by taking all possible combinations of **subset.cols.x** and **subset.cols.y**. If TRUE but **subset.cols.y=NA**, returns all comparisons involving any features in **subset.cols.x**. This argument is only used when **y=NA**.

concurrent.output Optional output file to which each comparison will be written as it is calculated.

make.output.table A boolean value indicating whether to include table-formatted output.

2.3 Output

`ccrepe` returns a *list* containing both the calculation results and the parameters used:

sim.score *matrix* of similarity scores for all requested comparisons. The (i,j) th element corresponds to the similarity score of column *i* (or the *i*th column of **subset.cols.1**) and column *j* (or the *j*th column of **subset.cols.1**) in one dataset, or to the similarity score of column *i* (or the *i*th column of **subset.cols.1**) in dataset *x* and column *j* (or the *j*th column of **subset.cols.2**) in dataset *y* in the case of two datasets.

p.values *matrix* of the corrected p-values for all requested comparisons. The (i,j) th element corresponds to the p-value of the (i,j) th element of **sim.score**.

q.values *matrix* of the Benjamini-Hochberg-Yekutieli corrected p-values. The (i,j) th element corresponds to the p-value of the (i,j) th element of **sim.score**.

z.stat *matrix* of the z-statistics used in generating the p-values for all requested comparisons. The (i,j) th element corresponds to the z-statistic generating the (i,j) th element of **p.values**.

2.4 Usage

```
ccrepe(  
  x = NA,  
  y = NA,  
  sim.score = cor,  
  sim.score.args = list(),  
  min.subj = 20,  
  iterations = 1000,  
  subset.cols.x = NULL,  
  subset.cols.y = NULL,  
  errthresh = 1e-04,  
  verbose = FALSE,  
  iterations.gap = 100,  
  distributions = NA,  
  compare.within.x = TRUE,  
  concurrent.output = NA,  
  make.output.table = FALSE)
```

2.5 Example 1

An example of how to use `ccrepe` with one dataset.

```
data <- matrix(rlnorm(40,meanlog=0,sdlog=1),nrow=10,ncol=4)  
data[,1] = 2*data[,2] + rnorm(10,0,0.01)  
data.rowsum <- apply(data,1,sum)  
data.norm <- data/data.rowsum  
apply(data.norm,1,sum) # The rows sum to 1, so the data are normalized  
## [1] 1 1 1 1 1 1 1 1 1 1  
test.input <- data.norm  
  
dimnames(test.input) <- list(c(  
  "Sample 1", "Sample 2", "Sample 3", "Sample 4", "Sample 5",  
  "Sample 6", "Sample 7", "Sample 8", "Sample 9", "Sample 10"),  
  c("Feature 1", "Feature 2", "Feature 3", "Feature 4"))  
  
test.output <- ccrepe(x=test.input, iterations=20, min.subj=10)
```

```
par(mfrow=c(1,2))  
plot(data[,1],data[,2],xlab="Feature 1",ylab="Feature 2",main="Non-normalized")  
plot(data.norm[,1],data.norm[,2],xlab="Feature 1",ylab="Feature 2",  
  main="Normalized")
```

CCREPE: Compositionality Corrected by PERmutation and RENormalization

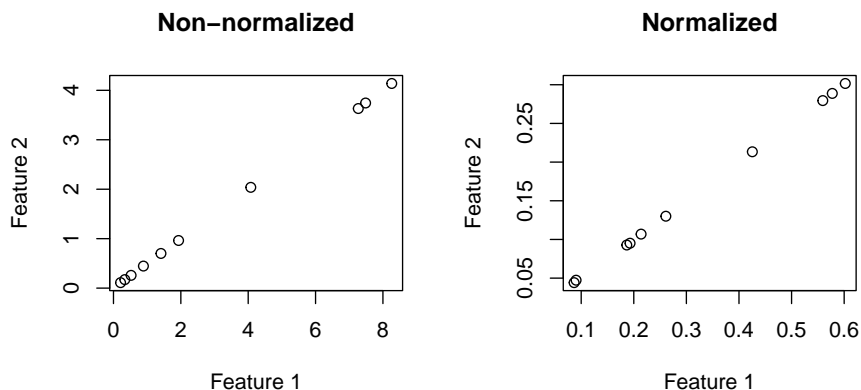


Figure 1: Non-normalized and normalized associations between feature 1 and feature 2. In this case we would expect feature 1 and feature 2 to be associated. In the output we see this by the positive sim.score value in the [1,2] element of test.output\$sim.score and the small q-value in the [1,2] element of test.output\$q.values.

```
test.output
## $p.values
##           Feature 1   Feature 2   Feature 3   Feature 4
## Feature 1           NA 3.577884e-09 0.09491487 0.5358208
## Feature 2 3.577884e-09           NA 0.01862487 0.1171381
## Feature 3 9.491487e-02 1.862487e-02           NA 0.2869634
## Feature 4 5.358208e-01 1.171381e-01 0.28696345           NA
##
## $z.stat
##           Feature 1   Feature 2   Feature 3   Feature 4
## Feature 1           NA 5.902613 -1.670023 -0.6191451
## Feature 2 5.9026132           NA -2.352955 -1.5668998
## Feature 3 -1.6700227 -2.352955           NA 1.0648066
## Feature 4 -0.6191451 -1.566900 1.064807           NA
##
## $sim.score
##           Feature 1   Feature 2   Feature 3   Feature 4
## Feature 1           NA 0.9999560 -0.7968531 -0.4822551
## Feature 2 0.9999560           NA -0.7946476 -0.4854218
## Feature 3 -0.7968531 -0.7946476           NA -0.1449802
## Feature 4 -0.4822551 -0.4854218 -0.1449802           NA
##
## $q.values
##           Feature 1   Feature 2   Feature 3   Feature 4
## Feature 1           NA 5.085551e-08 0.4497019 1.2693461
## Feature 2 5.085551e-08           NA 0.1323655 0.4162459
## Feature 3 4.497019e-01 1.323655e-01           NA 0.8157711
## Feature 4 1.269346e+00 4.162459e-01 0.8157711           NA
```

2.6 Example 2

An example of how to use `ccrepe` with two datasets.

```
data <- matrix(rlnorm(40,meanlog=0,sdlog=1),nrow=10,ncol=4)
data[,1] = 2*data[,2] + rnorm(10,0,0.01)
data.rowsum <- apply(data,1,sum)
data.norm <- data/data.rowsum
apply(data.norm,1,sum) # The rows sum to 1, so the data are normalized
## [1] 1 1 1 1 1 1 1 1 1 1

test.input <- data.norm

data2 <- matrix(rlnorm(105,meanlog=0,sdlog=1),nrow=15,ncol=7)
aligned.rows <- c(seq(1,4),seq(6,9),11,12) # The datasets dont need
# to have subjects line up exactly
data2[aligned.rows,1] <- 2*data[,3] + rnorm(10,0,0.01)
data2.rowsum <- apply(data2,1,sum)
data2.norm <- data2/data2.rowsum
apply(data2.norm,1,sum) # The rows sum to 1, so the data are normalized
## [1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1

test.input.2 <- data2.norm

dimnames(test.input) <- list(paste("Sample",seq(1,10)),paste("Feature",seq(1,4)))
dimnames(test.input.2) <- list(paste("Sample",c(seq(1,4),11,seq(5,8),12,9,10,13,14,15)),paste("Feature",seq(1,7)))

test.output.two.datasets <- ccrepe(x=test.input, y=test.input.2, iterations=20, min.subj=10)
## Warning in preprocess_data(CA): Removing subjects Sample 11, Sample 12, Sample
13, Sample 14, Sample 15 from dataset y because they are not in dataset x.
```

Please note that we receive a warning because the subjects don't match - only paired observations.

```
par(mfrow=c(1,2))
plot(data2[aligned.rows,1],data[,3],xlab="dataset 2: Feature 1",ylab="dataset 1: Feature 3",main="Non-normalized")
plot(data2.norm[aligned.rows,1],data.norm[,3],xlab="dataset 2: Feature 1",ylab="dataset 1: Feature 3",
      main="Normalized")
```

CCREPE: Compositionality Corrected by PERmutation and REnormalization

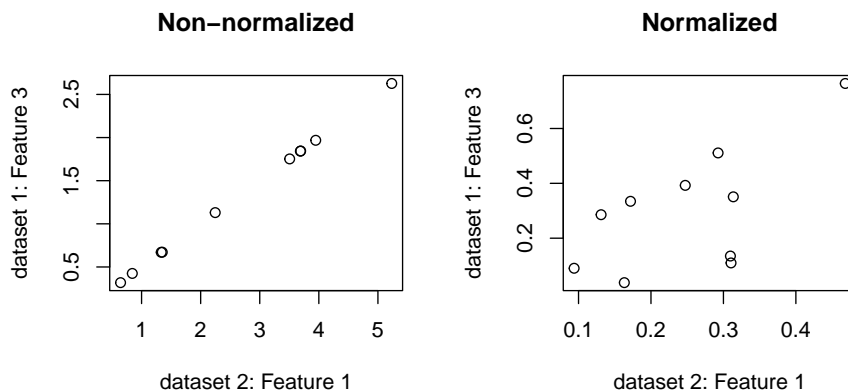


Figure 2: Non-normalized and normalized associations between feature 1 and feature 2. In this case we would expect feature 1 and feature 2 to be associated. In the output we see this by the positive `sim.score` value in the [1,2] element of `test.output$sim.score` and the small q-value in the [1,2] element of `test.output$q.values`.

```
test.output.two.datasets
## $p.values
##      Feature 1 Feature 2 Feature 3 Feature 4 Feature 5 Feature 6
## Feature 1 0.1843066 0.8729028 0.4518731 0.5389576 0.09670803 0.6046174
## Feature 2 0.2264195 0.8419358 0.6675011 0.3322070 0.37402040 0.8306622
## Feature 3 0.2705223 0.9340202 0.6059841 0.3982781 0.14813813 0.9174762
## Feature 4 0.5873301 0.7625890 0.6999057 0.5364062 0.44806839 0.7650211
##      Feature 7
## Feature 1 0.8683891
## Feature 2 0.7941651
## Feature 3 0.5183697
## Feature 4 0.6653469
##
## $z.stat
##      Feature 1 Feature 2 Feature 3 Feature 4 Feature 5 Feature 6
## Feature 1 -1.3276111 0.15997236 -0.7522959 0.6143900 1.6610270 -0.5177719
## Feature 2 -1.2096336 -0.19941799 -0.4295802 0.9696780 0.8889678 -0.2138523
## Feature 3 1.1018605 0.08278787 0.5158143 -0.8447005 -1.4461393 0.1036133
## Feature 4 0.5427091 0.30208278 0.3854477 -0.6182566 -0.7586392 0.2988937
##      Feature 7
## Feature 1 0.1657050
## Feature 2 0.2609059
## Feature 3 -0.6458605
## Feature 4 0.4325429
##
## $sim.score
##      Feature 1 Feature 2 Feature 3 Feature 4 Feature 5
## Feature 1 -0.58625840 0.06098857 -0.059367067 -0.1901324 0.4525316
## Feature 2 -0.58599413 0.06445441 -0.062365664 -0.1903745 0.4506643
## Feature 3 0.65063910 -0.09671751 0.006139096 0.3850556 -0.3954591
## Feature 4 0.09356566 0.03832872 0.115534038 -0.2722102 -0.2568041
```

CCREPE: Compositionality Corrected by PERmutation and RENormalization

```
##           Feature 6 Feature 7
## Feature 1  0.10032905  0.1019982
## Feature 2  0.10075711  0.1014961
## Feature 3 -0.17763715 -0.2119829
## Feature 4  0.09898218  0.1558946
##
## $q.values
##           Feature 1 Feature 2 Feature 3 Feature 4 Feature 5 Feature 6
## Feature 1  6.724966  3.675047  4.946373  4.538179 10.586025  4.412247
## Feature 2  6.196183  3.840061  4.059288  6.060771  5.848812  3.953366
## Feature 3  5.922478  3.651477  4.145831  5.449628  8.107879  3.719644
## Feature 4  4.592240  4.173793  4.032332  4.893087  5.449695  3.987719
##           Feature 7
## Feature 1  3.802286
## Feature 2  3.951468
## Feature 3  5.158427
## Feature 4  4.284199
```

2.7 Example 3

An example of how to use `ccrepe` with `nc.score` as the similarity score.

```
data <- matrix(rlnorm(40,meanlog=0,sdlog=1),nrow=10,ncol=4)
data[,1] = 2*data[,2] + rnorm(10,0,0.01)
data.rowsum <- apply(data,1,sum)
data.norm <- data/data.rowsum
apply(data.norm,1,sum) # The rows sum to 1, so the data are normalized
## [1] 1 1 1 1 1 1 1 1 1 1
test.input <- data.norm
dimnames(test.input) <- list(paste("Sample",seq(1,10)),paste("Feature",seq(1,4)))
test.output.nc.score <- ccrepe(x=test.input, sim.score=nc.score, iterations=20, min.subj=10)

par(mfrow=c(1,2))
plot(data[,1],data[,2],xlab="Feature 1",ylab="Feature 2",main="Non-normalized")
plot(data.norm[,1],data.norm[,2],xlab="Feature 1",ylab="Feature 2",
      main="Normalized")
```

CCREPE: Compositionality Corrected by PERmutation and RENormalization

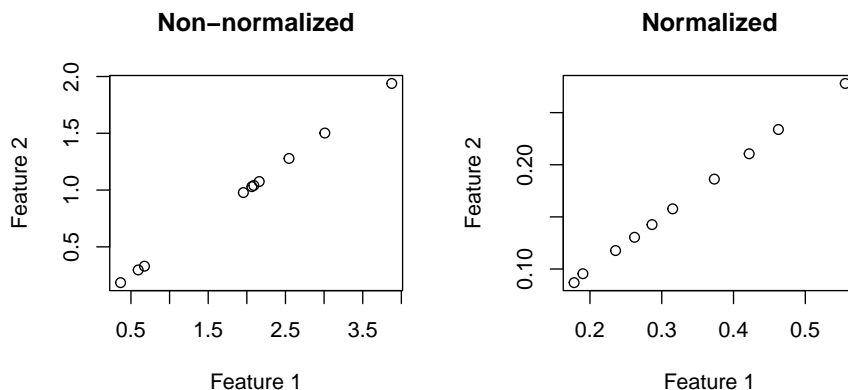


Figure 3: Non-normalized and normalized associations between feature 1 and feature 2. In this case we would expect feature 1 and feature 2 to be associated. In the output we see this by the positive sim.score value in the [1,2] element of test.output\$sim.score and the small q-value in the [1,2] element of test.output\$q.values. In this case, however, the sim.score represents the NC-Score between two features rather than the Spearman correlation.

```
test.output.nc.score

## $p.values
##           Feature 1   Feature 2 Feature 3 Feature 4
## Feature 1           NA 0.0001158502 0.7601491 0.7520893
## Feature 2 0.0001158502           NA 0.7156233 0.4569500
## Feature 3 0.7601490575 0.7156232964           NA 0.6558845
## Feature 4 0.7520892654 0.4569499883 0.6558845           NA
##
## $z.stat
##           Feature 1   Feature 2   Feature 3   Feature 4
## Feature 1           NA 3.8547453 -0.3052851 -0.3158857
## Feature 2 3.8547453           NA -0.3643143 -0.7438785
## Feature 3 -0.3052851 -0.3643143           NA 0.4456023
## Feature 4 -0.3158857 -0.7438785 0.4456023           NA
##
## $sim.score
##           Feature 1   Feature 2   Feature 3   Feature 4
## Feature 1           NA 1.00000 -0.34375 -0.4375
## Feature 2 1.00000           NA -0.34375 -0.4375
## Feature 3 -0.34375 -0.34375           NA -0.1250
## Feature 4 -0.43750 -0.43750 -0.12500           NA
##
## $q.values
##           Feature 1   Feature 2   Feature 3   Feature 4
## Feature 1           NA 0.001646677 1.800774 2.138017
## Feature 2 0.001646677           NA 2.542941 3.247509
## Feature 3 1.800774215 2.542940692           NA 3.107548
## Feature 4 2.138016922 3.247509479 3.107548           NA
```

2.8 Example 4

An example of how to use `ccrepe` with a user-defined `sim.score` function.

```

data <- matrix(rlnorm(40,meanlog=0,sdlog=1),nrow=10,ncol=4)
data[,1] = 2*data[,2] + rnorm(10,0,0.01)
data.rowsum <- apply(data,1,sum)
data.norm <- data/data.rowsum
apply(data.norm,1,sum) # The rows sum to 1, so the data are normalized

## [1] 1 1 1 1 1 1 1 1 1 1

test.input <- data.norm

dimnames(test.input) <- list(paste("Sample",seq(1,10)),paste("Feature",seq(1,4)))

my.test.sim.score <- function(x,y=NA,constant=0.5){
  if(is.vector(x) && is.vector(y)) return(constant)
  if(is.matrix(x) && is.na(y)) return(matrix(rep(constant,ncol(x)^2),ncol=ncol(x)))
  if(is.data.frame(x) && is.na(y)) return(matrix(rep(constant,ncol(x)^2),ncol=ncol(x)))
  else stop('ERROR')
}

test.output.sim.score <- ccrepe(x=test.input, sim.score=my.test.sim.score, iterations=20, min.subj=10, sim

par(mfrow=c(1,2))
plot(data[,1],data[,2],xlab="Feature 1",ylab="Feature 2",main="Non-normalized")
plot(data.norm[,1],data.norm[,2],xlab="Feature 1",ylab="Feature 2",
      main="Normalized")

```

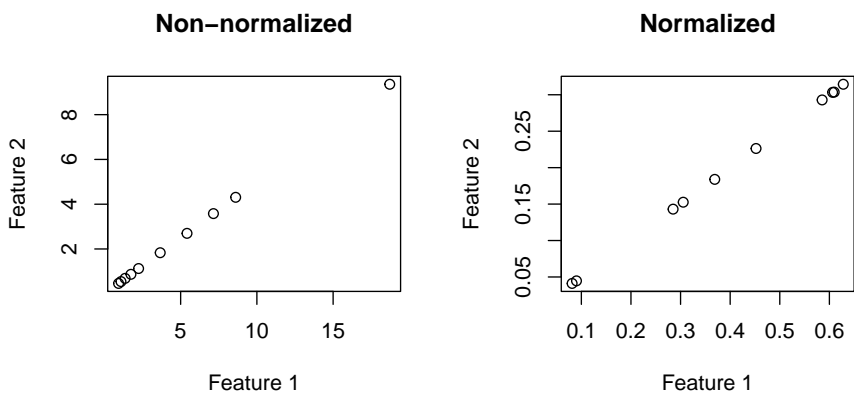


Figure 4: Non-normalized and normalized associations between feature 1 and feature 2. In this case we would expect feature 1 and feature 2 to be associated. Note that the values of `sim.score` are all 0.6 and none of the `p-values` are very small because of the arbitrary definition of the similarity score.

```

test.output.sim.score

## $p.values
##      Feature 1 Feature 2 Feature 3 Feature 4

```

CCREPE: Compositionality Corrected by PERmutation and RENormalization

```
## Feature 1      NA      NaN      NaN      NaN
## Feature 2     NaN      NA      NaN      NaN
## Feature 3     NaN      NaN      NA      NaN
## Feature 4     NaN      NaN      NaN      NA
##
## $z.stat
##           Feature 1 Feature 2 Feature 3 Feature 4
## Feature 1      NA      NaN      NaN      NaN
## Feature 2     NaN      NA      NaN      NaN
## Feature 3     NaN      NaN      NA      NaN
## Feature 4     NaN      NaN      NaN      NA
##
## $sim.score
##           Feature 1 Feature 2 Feature 3 Feature 4
## Feature 1      NA      0.6      0.6      0.6
## Feature 2     0.6      NA      0.6      0.6
## Feature 3     0.6      0.6      NA      0.6
## Feature 4     0.6      0.6      0.6      NA
##
## $q.values
##           Feature 1 Feature 2 Feature 3 Feature 4
## Feature 1      NA      NaN      NaN      NaN
## Feature 2     NaN      NA      NaN      NaN
## Feature 3     NaN      NaN      NA      NaN
## Feature 4     NaN      NaN      NaN      NA
```

2.9 Example 5

An example of how to use `ccrepe` when specifying column subsets.

```
data <- matrix(rlnorm(40,meanlog=0,sdlog=1),nrow=10,ncol=4)
data.rowsum <- apply(data,1,sum)
data.norm <- data/data.rowsum
apply(data.norm,1,sum) # The rows sum to 1, so the data are normalized
## [1] 1 1 1 1 1 1 1 1 1 1

test.input <- data.norm

dimnames(test.input) <- list(paste("Sample",seq(1,10)),paste("Feature",seq(1,4)))

test.output.1.3 <- ccrepe(x=test.input, iterations=20, min.subj=10, subset.cols.x=c(1,3))
test.output.1 <- ccrepe(x=test.input, iterations=20, min.subj=10, subset.cols.x=c(1), compare.within.x=
test.output.1.2.3 <- ccrepe(x=test.input, iterations=20, min.subj=10, subset.cols.x=c(1,2),subset.cols.y=c
test.output.1.3$sim.score

##           Feature 1 Feature 3
## Feature 1      NA 0.2824729
## Feature 3 0.2824729      NA

test.output.1$sim.score

##           Feature 1 Feature 2 Feature 3 Feature 4
```

```
## Feature 1      NA -0.6761607 0.2824729 0.01238038
## Feature 2 -0.67616066      NA      NA      NA
## Feature 3 0.28247292      NA      NA      NA
## Feature 4 0.01238038      NA      NA      NA

test.output.12.3$sim.score

##      Feature 1 Feature 2 Feature 3 Feature 4
## Feature 1      NA      NA 0.2824729      NA
## Feature 2      NA      NA -0.6577615      NA
## Feature 3 0.2824729 -0.6577615      NA      NA
## Feature 4      NA      NA      NA      NA
```

3 nc.score

The `nc.score` similarity measure is an N-dimensional extension of the checkerboard score particularly suited to similarity score calculations between compositions derived from ecological relative abundance measurements. In such cases, features typically represent species abundances, and the NC-score discretizes these continuous values into one of N bins before computing a normalized similarity of co-occurrence or co-exclusion. This can be used as a standalone function or with `ccrepe` as above to obtain compositionality-corrected p-values.

3.1 General Functionality

The NC-score is an extension to Diamond's checkerboard score (see [Cody and Diamond \[1975\]](#)) to ordinal data, and simplifies to a calculation of Kendall's τ on binned data instead of ranked data. Let two features in a dataset with n subjects be denoted by

$$\begin{bmatrix} x_1 & x_2 & \dots & x_n \\ y_1 & y_2 & \dots & y_n \end{bmatrix}.$$

The binning function maps from the original data to b numbered bins in $\{1, \dots, b\}$. Let the binning function be denoted by $B(\cdot)$. The co-variation and co-exclusion patterns are the same as concordant and discordant pairs in Kendall's τ . Considering a 2×2 submatrix of the form

$$\begin{bmatrix} B(x_i) & B(x_j) \\ B(y_i) & B(y_j) \end{bmatrix},$$

a co-variation pattern is counted when $(B(x_i) - B(x_j))(B(y_i) - B(y_j)) > 0$ and a co-exclusion pattern, conversely, when $(B(x_i) - B(x_j))(B(y_i) - B(y_j)) < 0$. The NC-score statistic for features x and y is then defined as

$$(\text{number of co-variation patterns}) - (\text{number of co-exclusion patterns}),$$

normalized by the Kendall's τ normalization factor accounting for ties described in [Kendall \[1970\]](#).

3.2 Arguments

- ✕ First numerical *vector*, or single *dataframe* or *matrix*, containing relative abundances. If the latter, columns are features, rows are samples. Rows should therefore sum to a constant.

CCREPE: Compositionality Corrected by PERmutation and RENormalization

y If provided, second numerical *vector* containing relative abundances. If given, **x** must be a *vector* as well.

nbins A non-negative integer of the number of bins to generate (cutoffs will be generated by the `discretize` function from the `infotheo` package).

bin.cutoffs A list of values demarcating the bin cutoffs. The binning is performed using the `findInterval` function.

use An optional character string giving a method for computing covariances in the presence of missing values. This must be (an abbreviation of) one of the strings "everything", "all.obs", "complete.obs", "na.or.complete", or "pairwise.complete.obs".

3.3 Output

`nc.score` returns either a single number (if called with two vectors) or a *matrix* of all pairwise scores (if called with a *matrix*) of normalized scores. This behaviour is precisely analogous to the `cor` function in R

3.4 Usage

```
nc.score(  
  x,  
  y = NULL,  
  use = "everything",  
  nbins = NULL,  
  bin.cutoffs=NULL)
```

3.5 Example 1

An example of using `nc.score` to get a single similarity score or a matrix.

```
data <- matrix(rlnorm(40,meanlog=0,sdlog=1),nrow=10,ncol=4)  
data.rowsum <- apply(data,1,sum)  
data[,1] = 2*data[,2] + rnorm(10,0,0.01)  
data.norm <- data/data.rowsum  
apply(data.norm,1,sum) # The rows sum to 1, so the data are normalized  
## [1] 1.215751 1.836954 2.474784 1.674961 1.056729 1.398172 1.027829 1.475163  
## [9] 1.733673 1.736387  
  
test.input <- data.norm  
  
dimnames(test.input) <- list(paste("Sample",seq(1,10)),paste("Feature",seq(1,4)))  
  
test.output.matrix <- nc.score(x=test.input)  
test.output.num <- nc.score(x=test.input[,1],y=test.input[,2])
```

```
par(mfrow=c(1, 2))  
plot(data[,1],data[,2],xlab="Feature 1",ylab="Feature 2",main="Non-normalized")
```

CCREPE: Compositionality Corrected by PERmutation and REnormalization

```
plot(data.norm[,1],data.norm[,2],xlab="Feature 1",ylab="Feature 2",  
      main="Normalized")
```

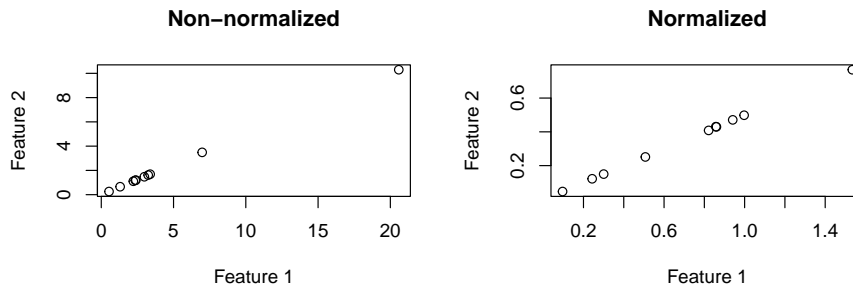


Figure 5: Non-normalized and normalized associations between feature 1 and feature 2 of the second example. Again, we expect to observe a positive association between feature 1 and feature 2. In terms of generalized checkerboard scores, we would expect to see more co-variation patterns than co-exclusion patterns. This is shown by the positive and relatively high value of the [1,2] element of test.output.matrix (which is identical to test.output.num)

```
test.output.matrix  
##           Feature 1 Feature 2 Feature 3 Feature 4  
## Feature 1    1.0000    0.81250  -0.1250  -0.43750  
## Feature 2    0.8125    1.00000  -0.1250  -0.59375  
## Feature 3   -0.1250   -0.12500    1.0000  -0.18750  
## Feature 4   -0.4375   -0.59375  -0.1875    1.00000  
  
test.output.num  
## [1] 0.8125
```

3.6 Example 2

An example of using `nc.score` with an arbitrary bin number.

```
data <- matrix(rlnorm(40,meanlog=0,sdlog=1),nrow=10,ncol=4)  
data.rowsum <- apply(data,1,sum)  
data[,1] = 2*data[,2] + rnorm(10,0,0.01)  
data.norm <- data/data.rowsum  
apply(data.norm,1,sum) # The rows sum to 1, so the data are normalized  
## [1] 1.0650119 1.2419817 2.0909150 1.5049011 1.1518230 0.8355778 1.0987219  
## [8] 0.8906571 0.9735985 1.7738150  
  
test.input <- data.norm  
  
dimnames(test.input) <- list(paste("Sample",seq(1,10)),paste("Feature",seq(1,4)))  
  
test.output <- nc.score(x=test.input,nbins=4)
```

CCREPE: Compositionality Corrected by PERmutation and RENormalization

```
par(mfrow=c(1, 2))
plot(data[,1],data[,2],xlab="Feature 1",ylab="Feature 2",main="Non-normalized")
plot(data.norm[,1],data.norm[,2],xlab="Feature 1",ylab="Feature 2",
      main="Normalized")
```

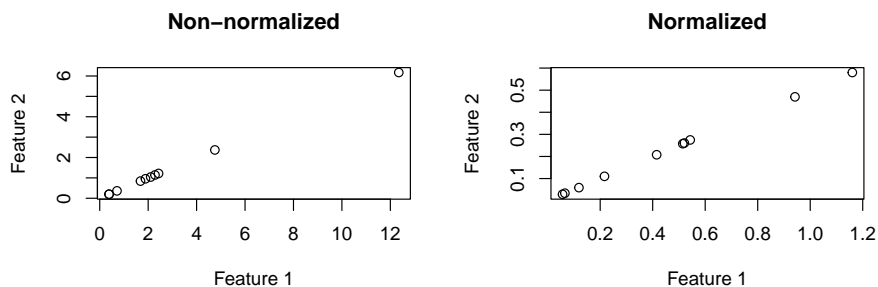


Figure 6: Non-normalized and normalized associations between feature 1 and feature 2 of the second example. Again, we expect to observe a positive association between feature 1 and feature 2. In terms of generalized checkerboard scores, we would expect to see more co-variation patterns than co-exclusion patterns. This is shown by the positive and relatively high value in the [1,2] element of test.output. In this case, the smaller bin number yields a smaller NC-score because of the coarser partitioning of the data.

```
test.output
##           Feature 1  Feature 2  Feature 3  Feature 4
## Feature 1  1.0000000  1.0000000 -0.25714286 -0.54285714
## Feature 2  1.0000000  1.0000000 -0.25714286 -0.54285714
## Feature 3 -0.2571429 -0.2571429  1.00000000 -0.05714286
## Feature 4 -0.5428571 -0.5428571 -0.05714286  1.00000000
```

3.7 Example 3

An example of using `nc.score` with user-defined bin edges.

```
data <- matrix(rlnorm(40,meanlog=0,sdlog=1),nrow=10,ncol=4)
data.rowsum <- apply(data,1,sum)
data[,1] = 2*data[,2] + rnorm(10,0,0.01)
data.norm <- data/data.rowsum
apply(data.norm,1,sum) # The rows sum to 1, so the data are normalized
## [1] 1.0356367 1.0619945 1.1517240 0.9623243 0.9874042 0.9661466 1.4627768
## [8] 1.1143498 1.0734688 0.8209136

test.input <- data.norm

dimnames(test.input) <- list(paste("Sample",seq(1,10)),paste("Feature",seq(1,4)))

test.output <- nc.score(x=test.input,bin.cutoffs=c(0.1,0.2,0.3))
```

```
par(mfrow=c(1, 2))
plot(data[,1],data[,2],xlab="Feature 1",ylab="Feature 2",main="Non-normalized")
```

CCREPE: Compositionality Corrected by PERmutation and RENormalization

```
plot(data.norm[,1],data.norm[,2],xlab="Feature 1",ylab="Feature 2",  
      main="Normalized")
```

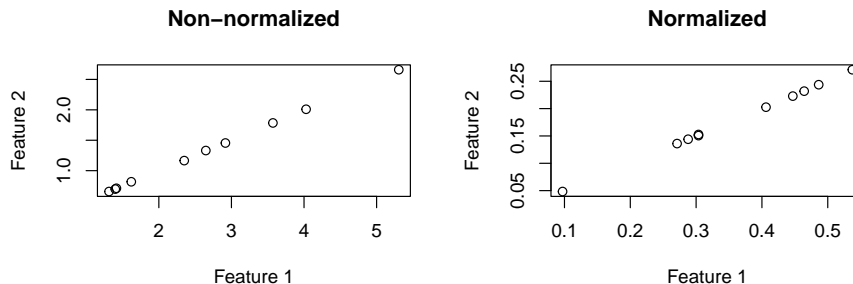


Figure 7: Non-normalized and normalized associations between feature 1 and feature 2 of the second example. Again, we expect to observe a positive association between feature 1 and feature 2. In terms of generalized checkerboard scores, we would expect to see more co-variation patterns than co-exclusion patterns. This is shown by the positive and relatively high value in the [1,2] element of test.output. The bin edges specified here represent almost absent ([0,0.001]), low abundance ([0.001,0.1]), medium abundance ([0.1,0.25]), and high abundance ([0.6,1]).

```
test.output
```

```
##           Feature 1  Feature 2  Feature 3  Feature 4  
## Feature 1  1.0000000  0.7356829  0.2467176 -0.1390096  
## Feature 2  0.7356829  1.0000000  0.3138824 -0.4951876  
## Feature 3  0.2467176  0.3138824  1.0000000 -0.6197798  
## Feature 4 -0.1390096 -0.4951876 -0.6197798  1.0000000
```

4 References

References

- Martin Leonard Cody and Jared Mason Diamond. *Ecology and evolution of communities*. Harvard University Press, 1975.
- Karoline Faust, J Fah Sathirapongsasuti, Jacques Izard, Nicola Segata, Dirk Gevers, Jeroen Raes, and Curtis Huttenhower. Microbial co-occurrence relationships in the human microbiome. *PLoS computational biology*, 8(7):e1002606, 2012.
- M.G. Kendall. *Rank correlation methods*. Charles Griffin & Co., 1970.
- Emma Schwager and Colleagues. Detecting statistically significant associations between sparse and high dimensional compositional data. In Progress.