

# CCREPE: Compositionality Corrected by PErmutation and REnormalization

Emma Schwager, George Weingart, Craig Bielski, Curtis Huttenhower

May 15, 2016

## Contents

---

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>ccrepe</b>	<b>2</b>
2.1	General functionality . . . . .	2
2.2	Arguments . . . . .	2
2.3	Output . . . . .	3
2.4	Usage . . . . .	3
2.5	Example 1 . . . . .	4
2.6	Example 2 . . . . .	6
2.7	Example 3 . . . . .	8
2.8	Example 4 . . . . .	10
2.9	Example 5 . . . . .	12
<b>3</b>	<b>nc.score</b>	<b>12</b>
3.1	General Functionality . . . . .	12
3.2	Arguments . . . . .	13
3.3	Output . . . . .	13
3.4	Usage . . . . .	13
3.5	Example 1 . . . . .	13
3.6	Example 2 . . . . .	15
3.7	Example 3 . . . . .	16
<b>4</b>	<b>References</b>	<b>17</b>

## 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-score), particularly appropriate to compositions derived from microbial community sequencing data. This results in p-values and false discovery rate q-values corrected for the effects of compositionality. The package contains two functions `ccrepe` and `nc.score` and is maintained by the Huttenhower lab ([ccrepe-users@googlegroups.com](mailto:ccrepe-users@googlegroups.com)).

## 2 ccrepe

---

`ccrepe` is the main package function. It calculates compositionality-corrected p-values and q-values for a user-selected similarity measure, operating on either one or two input matrices. If given one matrix, all features (columns) in the matrix are compared to each other using the selected similarity measure. If given two matrices, each feature in the first are compared against all features in the second.

### 2.1 General functionality

Compositional data induces spurious correlations between features due to the nonindependence of values that must sum to a fixed total. CCREPE abrogates this when determining the significance of a similarity measure for each feature pair using two main steps, permutation/renormalization and bootstrapping. First, given two features to compare, CCREPE generates a null distribution of the similarity expected just due to compositionality by iteratively permuting one feature, renormalizing all samples in the composition to their previous sum, and computing the resulting similarity measures. Second, CCREPE bootstraps 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).

`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")
```

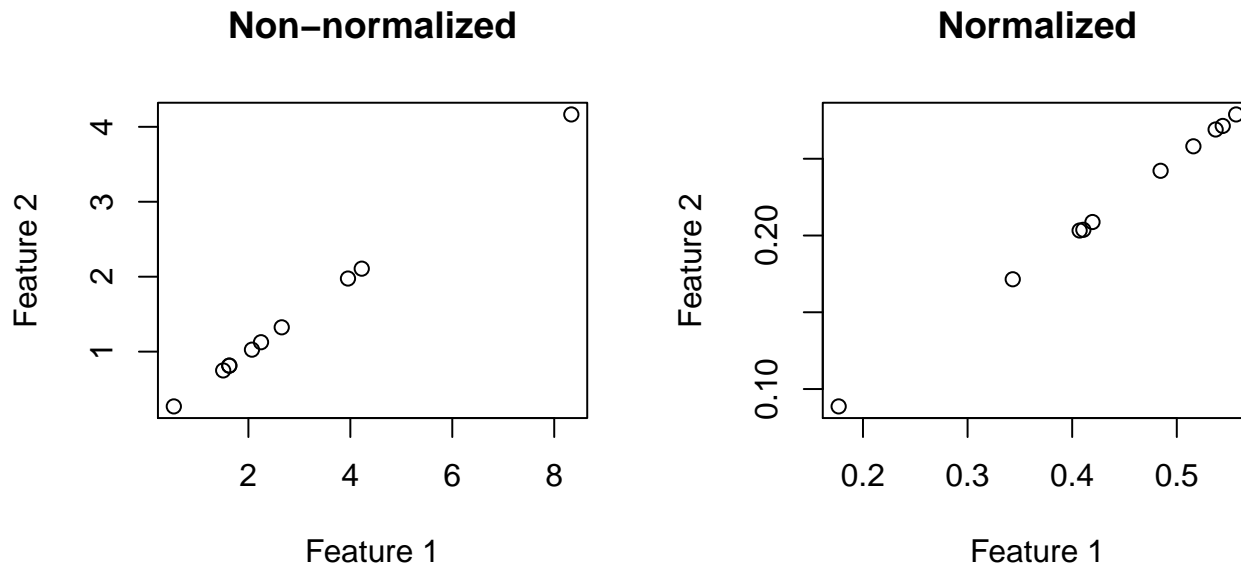


Figure 1: Non-normalized and normalized associations between feature 1 and feature 2. In this case we would expect feature 1 and feature 2 to be associated. In the output we see this by the positive `sim.score` value in the [1,2] element of `test.output$sim.score` and the small `q.value` in the [1,2]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

```
test.output
## $p.values
##           Feature 1   Feature 2   Feature 3 Feature 4
## Feature 1           NA 0.0001485624 0.055542690 0.6945296
## Feature 2 0.0001485624           NA 0.007420849 0.2523613
## Feature 3 0.0555426901 0.0074208493           NA 0.2692857
## Feature 4 0.6945295938 0.2523612974 0.269285732           NA
##
## $z.stat
##           Feature 1 Feature 2 Feature 3 Feature 4
## Feature 1           NA  3.793460 -1.914607 -0.3927155
## Feature 2  3.7934602           NA -2.677344 -1.1446328
## Feature 3 -1.9146066 -2.677344           NA  1.1047089
## Feature 4 -0.3927155 -1.144633  1.104709           NA
##
## $sim.score
##           Feature 1 Feature 2 Feature 3 Feature 4
## Feature 1           NA  0.9999338 -0.9653469 -0.3614493
## Feature 2  0.9999338           NA -0.9652598 -0.3616763
## Feature 3 -0.9653469 -0.9652598           NA  0.1056274
## Feature 4 -0.3614493 -0.3616763  0.1056274           NA
##
## $q.values
##           Feature 1   Feature 2   Feature 3 Feature 4
## Feature 1           NA 0.002111644 0.26315850 1.6453233
## Feature 2 0.002111644           NA 0.05273942 0.8967565
## Feature 3 0.263158503 0.052739422           NA 0.7655174
```

```
## Feature 4 1.645323338 0.896756457 0.76551744 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")
```

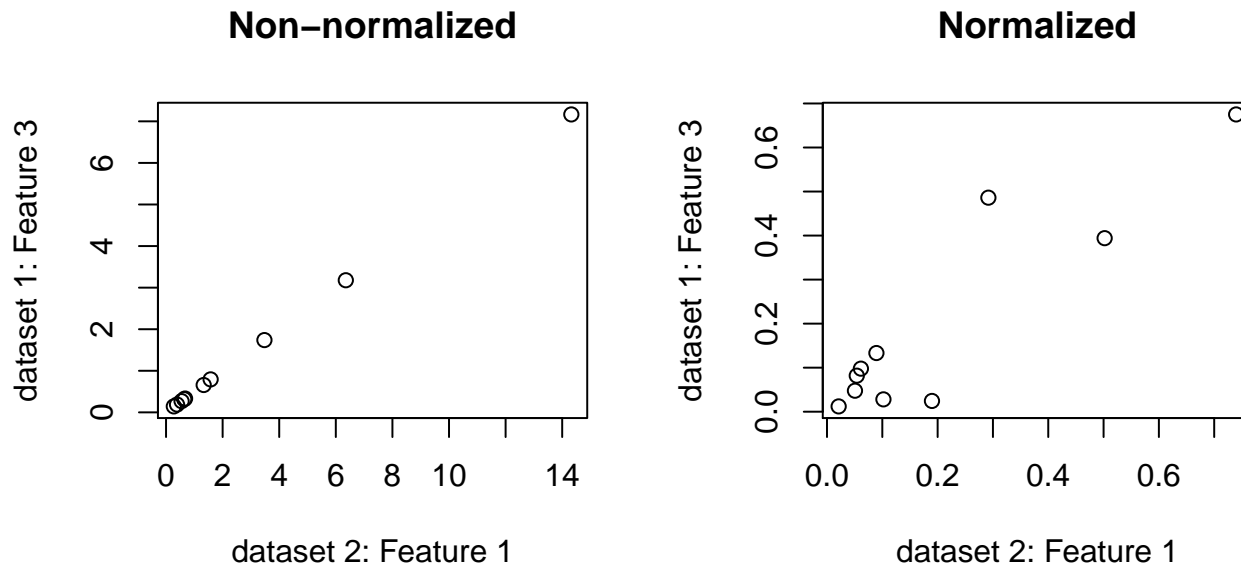


Figure 2: Non-normalized and normalized associations between feature 1 and feature 2. In this case we would expect feature 1 and feature 2 to be associated. In the output we see this by the positive `sim.score` value in the [1,2] element of `test.output$sim.score` and the small `q`-value in the [1,2]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

```
test.output.two.datasets
## $p.values
##           Feature 1 Feature 2 Feature 3 Feature 4 Feature 5 Feature 6
## Feature 1 0.034958088 0.9703375 0.64120682 0.17222247 0.8474826 0.6415652
## Feature 2 0.029774216 0.9433432 0.40847979 0.12310638 0.7782208 0.4919041
## Feature 3 0.002675776 0.5236336 0.10584242 0.01226415 0.6641736 0.8931339
## Feature 4 0.009086378 0.1804207 0.05309024 0.26462478 0.3971822 0.2006745
##           Feature 7
## Feature 1 0.0006413544
## Feature 2 0.0032006537
## Feature 3 0.0047879350
## Feature 4 0.9716468675
##
## $z.stat
##           Feature 1 Feature 2 Feature 3 Feature 4 Feature 5 Feature 6
## Feature 1 -2.108844 0.03718494 0.4660121 1.365097 0.1923314 -0.4655115
## Feature 2 -2.173081 0.07106856 0.8265718 1.541865 -0.2816384 -0.6872835
## Feature 3 3.002721 -0.63775429 -1.6171662 -2.504451 0.4341582 -0.1343398
## Feature 4 -2.608787 1.33946076 1.9341862 1.115527 -0.8466645 1.2796322
##           Feature 7
## Feature 1 3.41349457
## Feature 2 2.94777941
## Feature 3 -2.82096551
## Feature 4 0.03554286
##
## $sim.score
##           Feature 1 Feature 2 Feature 3 Feature 4 Feature 5 Feature 6
```

```
## Feature 1 -0.5990507 -0.01197985 0.2336724 0.4229193 -0.1258659 -0.1670415
## Feature 2 -0.6015587 -0.00745300 0.2351363 0.4276379 -0.1288281 -0.1686927
## Feature 3 0.9110280 -0.17669619 -0.5592206 -0.4742695 0.1920677 -0.1481001
## Feature 4 -0.6407536 0.33927167 0.6206051 0.1433585 -0.1344389 0.5536034
##           Feature 7
## Feature 1 0.7366686767
## Feature 2 0.7352165228
## Feature 3 -0.6862950852
## Feature 4 0.0007602632
##
## $q.values
##           Feature 1 Feature 2 Feature 3 Feature 4 Feature 5 Feature 6
## Feature 1 0.4783305 3.933956 3.5094456 1.5710100 3.865360 3.344197
## Feature 2 0.4655997 3.971611 2.6302197 1.2250626 3.703782 2.991426
## Feature 3 0.1464503 3.016785 1.1585910 0.2237466 3.304679 3.910632
## Feature 4 0.1989258 1.519195 0.6457175 1.9311216 2.717316 1.569042
##           Feature 7
## Feature 1 0.07020506
## Feature 2 0.11678520
## Feature 3 0.13102635
## Feature 4 3.79857587
```

## 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")
```



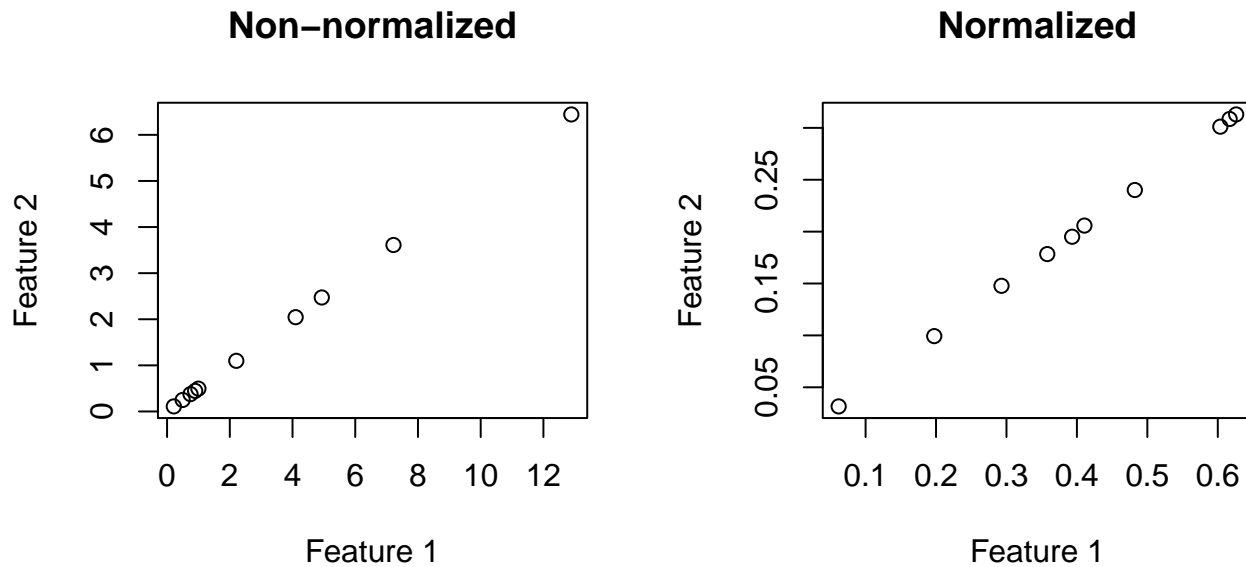


Figure 3: Non-normalized and normalized associations between feature 1 and feature 2. In this case we would expect feature 1 and feature 2 to be associated. In the output we see this by the positive `sim.score` value in the [1,2] element of `test.output$sim.score` and the small `q`-value in the [1,2] element of `test.output$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 1.074433e-06 0.2456618 0.2720205
## Feature 2 1.074433e-06           NA 0.1200195 0.1270534
## Feature 3 2.456618e-01 1.200195e-01           NA 0.1012892
## Feature 4 2.720205e-01 1.270534e-01 0.1012892           NA
##
## $z.stat
##           Feature 1 Feature 2 Feature 3 Feature 4
## Feature 1           NA 4.877492 -1.160951 -1.098421
## Feature 2 4.877492           NA -1.554692 -1.525825
## Feature 3 -1.160951 -1.554692           NA 1.638635
## Feature 4 -1.098421 -1.525825 1.638635           NA
##
## $sim.score
##           Feature 1 Feature 2 Feature 3 Feature 4
## Feature 1           NA 1.00000 -0.6250 -0.59375
## Feature 2 1.00000           NA -0.6250 -0.59375
## Feature 3 -0.62500 -0.62500           NA 0.18750
## Feature 4 -0.59375 -0.59375 0.1875           NA
##
## $q.values
##           Feature 1   Feature 2 Feature 3 Feature 4
## Feature 1           NA 1.527183e-05 0.6983599 0.6444098
## Feature 2 1.527183e-05           NA 0.5686465 0.4514794
```

```
## Feature 3 6.983599e-01 5.686465e-01      NA 0.7198549
## Feature 4 6.444098e-01 4.514794e-01 0.7198549      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,

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")
```

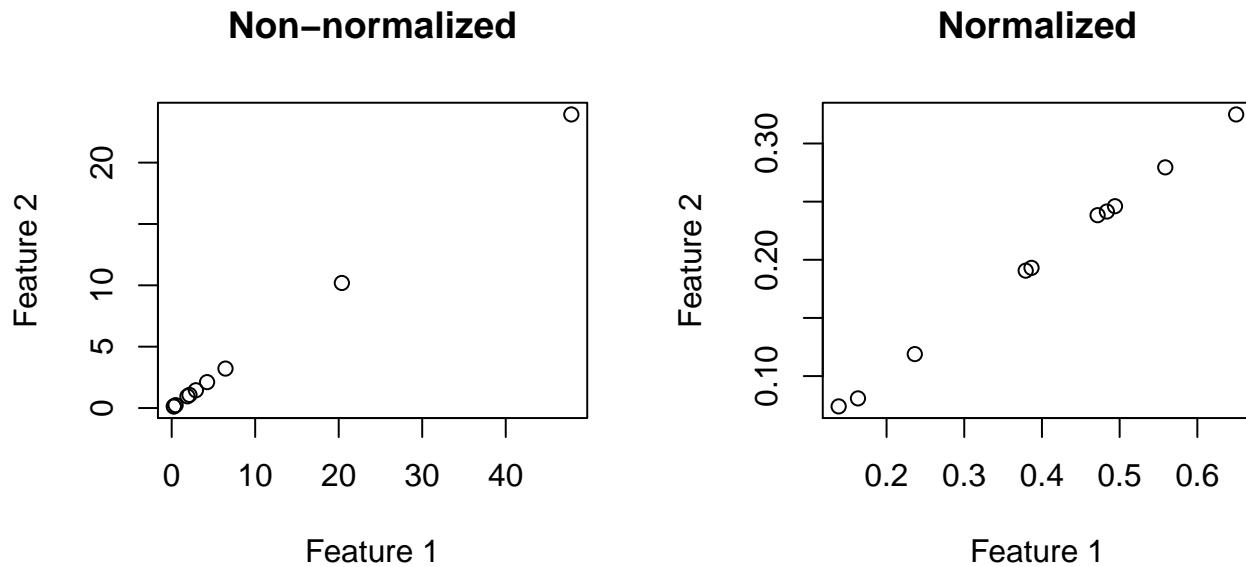


Figure 4: Non-normalized and normalized associations between feature 1 and feature 2. In this case we would expect feature 1 and feature 2 to be associated. Note that the values of sim.score are all 0.6 and none of the p-values are very small because of the arbitrary definition of the similarity score.

```
test.output.sim.score
## $p.values
##      Feature 1 Feature 2 Feature 3 Feature 4
## 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=1)
test.output.12.3 <- ccrepe(x=test.input, iterations=20, min.subj=10, subset.cols.x=c(1,2),subset.cols.y=c(2,3))
test.output.1.3$sim.score
##           Feature 1  Feature 3
## Feature 1          NA -0.05266883
## Feature 3 -0.05266883          NA

test.output.1$sim.score
##           Feature 1  Feature 2  Feature 3  Feature 4
## Feature 1          NA -0.5558889 -0.05266883 -0.2381933
## Feature 2 -0.5558885          NA          NA          NA
## Feature 3 -0.05266883          NA          NA          NA
## Feature 4 -0.23819325          NA          NA          NA

test.output.12.3$sim.score
##           Feature 1  Feature 2  Feature 3  Feature 4
## Feature 1          NA          NA -0.05266883          NA
## Feature 2          NA          NA -0.49337848          NA
## Feature 3 -0.05266883 -0.4933785          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  $\tau$  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  $\tau$ . Considering a  $2 \times 2$  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  $\tau$  normalization factor accounting for ties described in [Kendall \[1970\]](#).

## 3.2 Arguments

**x** 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.

**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] 1.5107606 0.9609102 1.9881809 1.4992654 0.8981335 1.1928673 0.5869541
## [8] 1.0600679 1.7780739 2.0831505

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

```

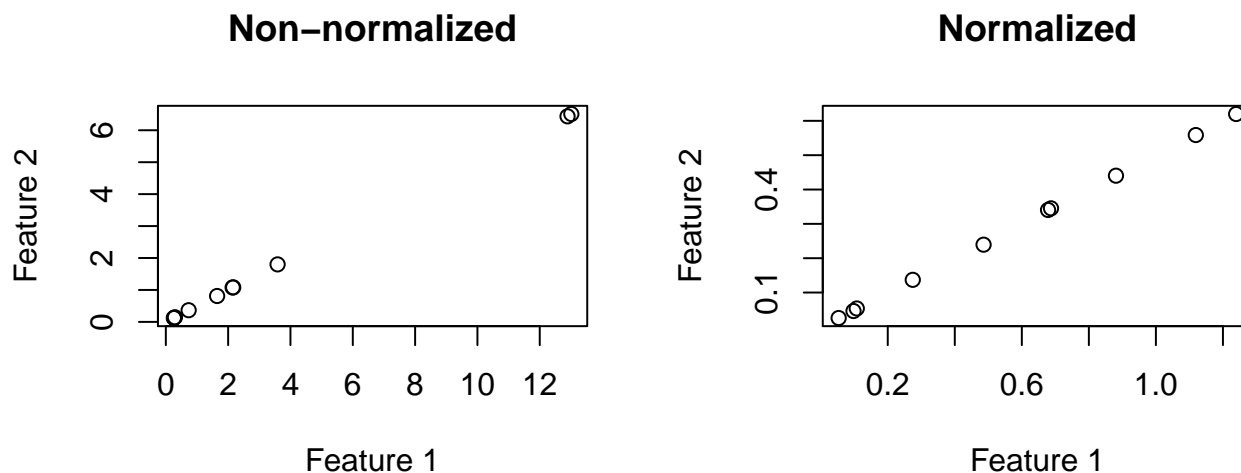


Figure 5: Non-normalized and normalized associations between feature 1 and feature 2 of the second example. Again, we expect to observe a positive association between feature 1 and feature 2. In terms of generalized checkerboard scores, we would expect to see more co-variation patterns than co-exclusion 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  1.00000  1.00000 -0.59375 -0.21875
## Feature 2  1.00000  1.00000 -0.59375 -0.21875
## Feature 3 -0.59375 -0.59375  1.00000 -0.21875
## Feature 4 -0.21875 -0.21875 -0.21875  1.00000

test.output.num
## [1] 1

```

### 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] 1.7974522 1.4345805 1.8447689 1.1590105 1.3031820 1.8836376 1.9292757
## [8] 0.6441619 0.7598349 1.0077756

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)

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")
```

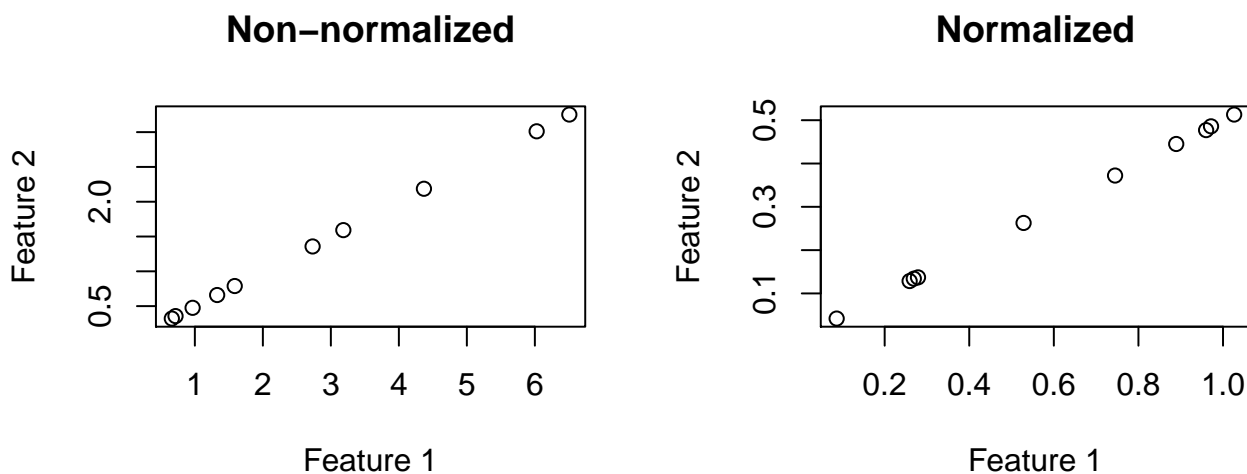


Figure 6: Non-normalized and normalized associations between feature 1 and feature 2 of the second example. Again, we expect to observe a positive association between feature 1 and feature 2. In terms of generalized checkerboard scores, we would expect to see more co-variation patterns than co-exclusion 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  0.3142857 -0.3142857
## Feature 2  1.0000000  1.0000000  0.3142857 -0.3142857
## Feature 3  0.3142857  0.3142857  1.0000000 -0.6285714
## Feature 4 -0.3142857 -0.3142857 -0.6285714  1.0000000
```

### 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] 1.4847395 1.3918864 1.5769226 1.2896420 1.8828316 1.2817689 1.2791274
## [8] 1.8449508 0.8097623 1.8417035

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

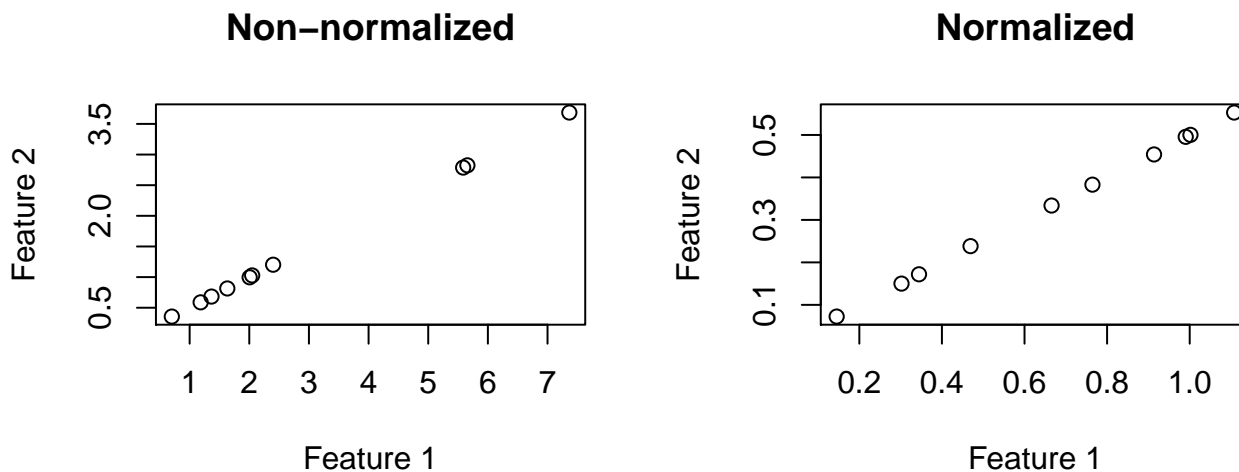


Figure 7: Non-normalized and normalized associations between feature 1 and feature 2 of the second example. Again, we expect to observe a positive association between feature 1 and feature 2. In terms of generalized checkerboard scores, we would expect to see more co-variation patterns than co-exclusion 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 ( $[0.25,1)$ ).

```
test.output
##           Feature 1  Feature 2  Feature 3  Feature 4
## Feature 1  1.0000000  0.5570860 -0.3481553 -0.05479966
## Feature 2  0.55708601  1.0000000 -0.6465082 -0.27475313
## Feature 3 -0.34815531 -0.6465082  1.0000000  0.00000000
## Feature 4 -0.05479966 -0.2747531  0.0000000  1.00000000
```



## 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.