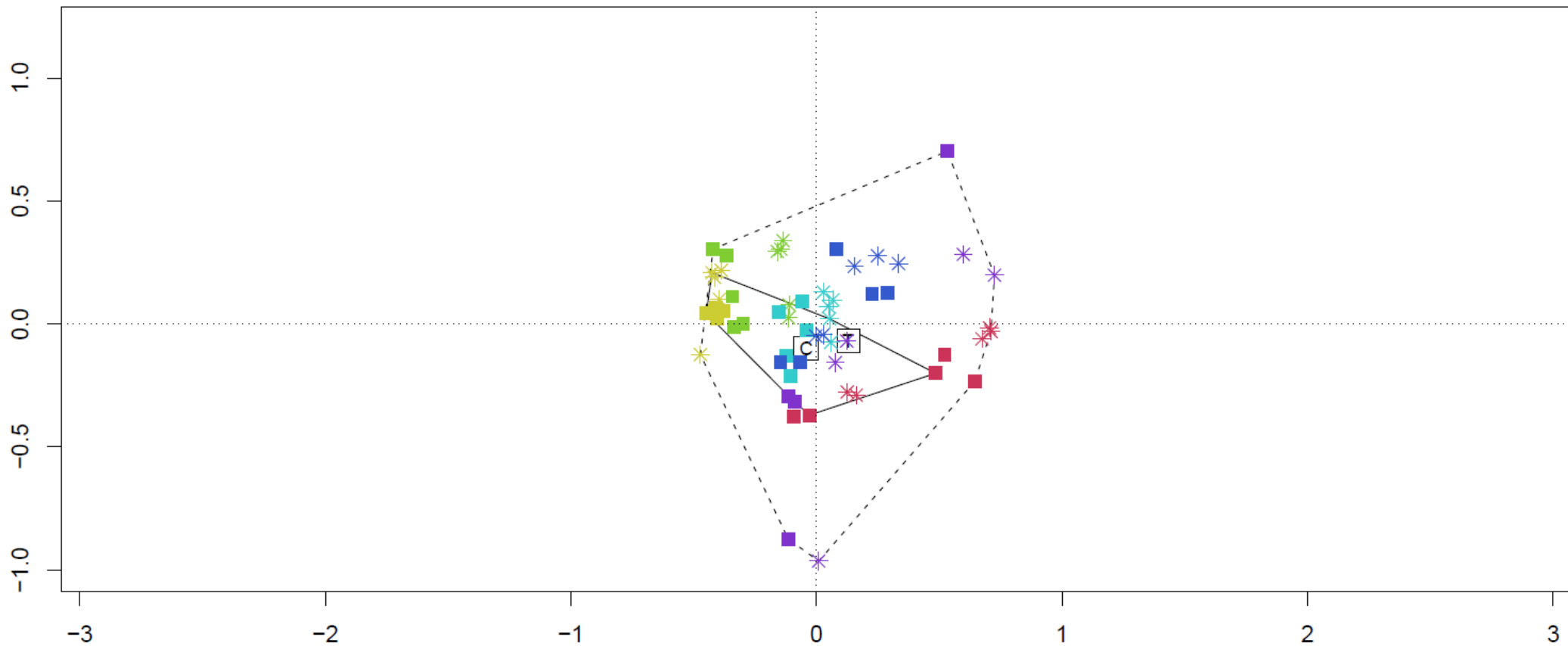
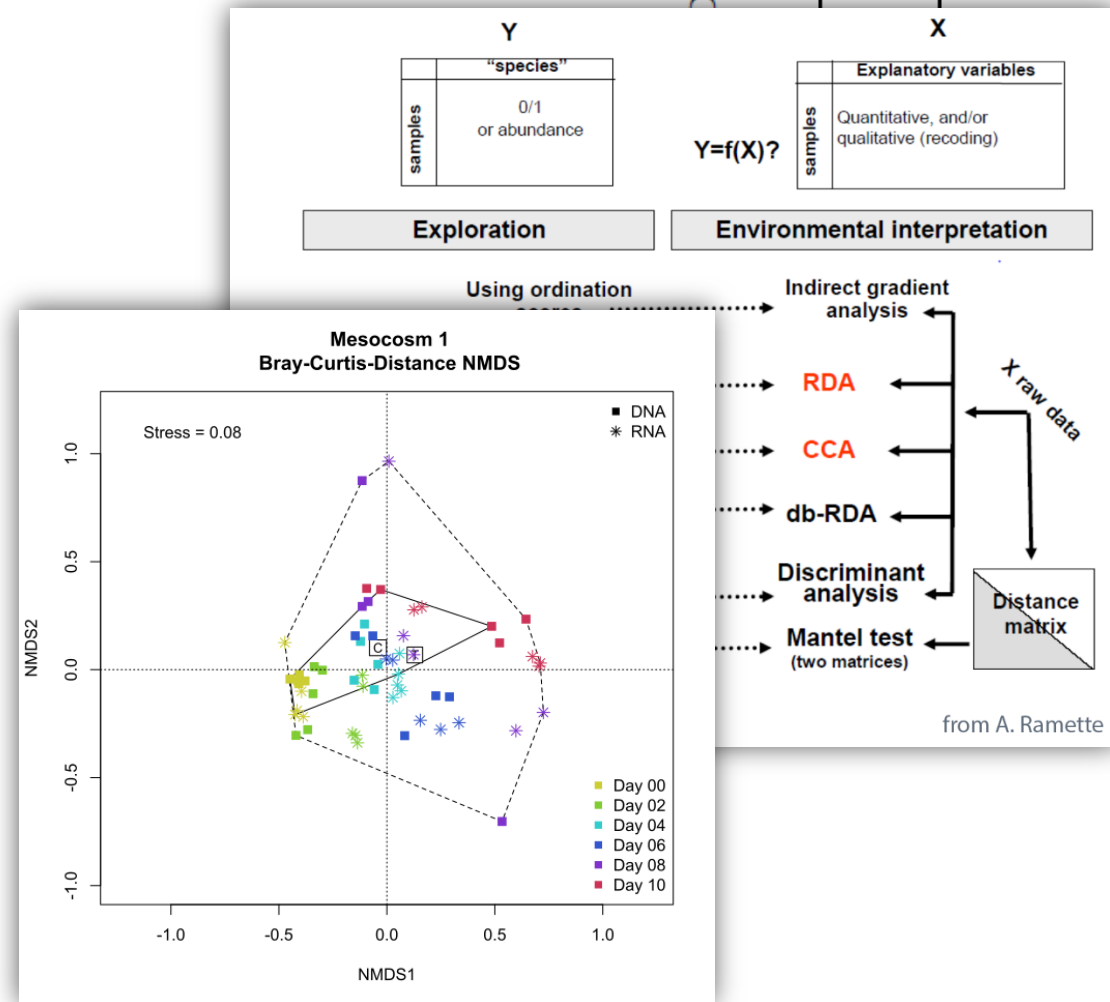
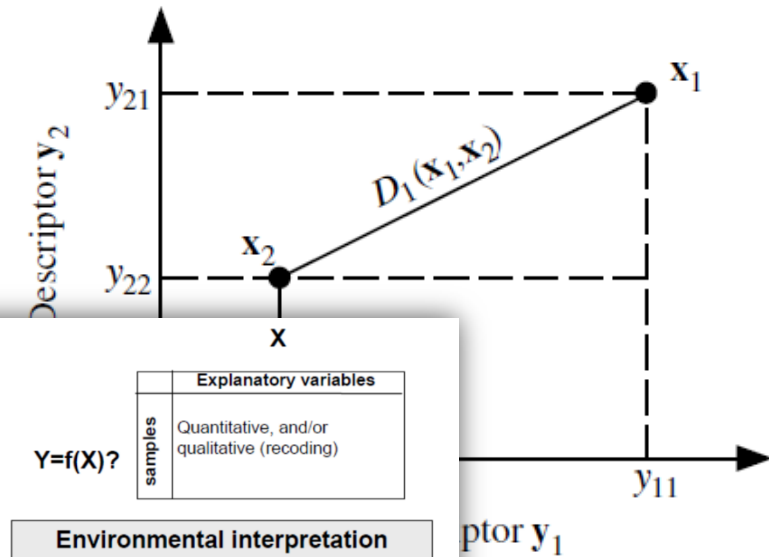


# Distance measures

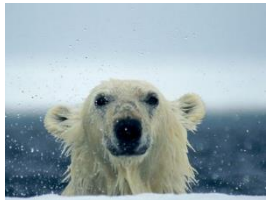


# Distance/Similarity in general

- Measure of relatedness between samples/sites
- Distance matrices used for a variety of multivariate ordination techniques (NMDS, PCoA, Cluster-analysis)



# The double-zero-problem



Abundance Table

Arctic	Hawaii	Deep Sea	Antarctica	Down town Oldenburg
1	0	0	0	0

Shared absence is no measure for similar environments

- Distribution Barriers
- Replacement Species

-> in environmental ecology: asymmetric measures

**symetric distances:** euclidean, Mahalanobis, mean character, coefficient of racial likeness



# Presence/Absence Measures

- Comparing communities based on the absence/presence of species in a sample
- Species-count is unimportant

	Sample_1540_DNA	Sample_1540_RNA		Sample_1540_DNA	Sample_1540_RNA
Bacteria;Actinobacteria	12	8	Bacteria;Actinobacteria	1	1
Bacteria;Bacteroidetes	1415	1326	Bacteria;Bacteroidetes	1	1
Bacteria;Chlamydiae	0	0	Bacteria;Chlamydiae	0	0
Bacteria;Cyanobacteria	9	0	Bacteria;Cyanobacteria	1	0
Bacteria;Fibrobacteres	0	0	Bacteria;Fibrobacteres	0	0
Bacteria;Firmicutes	0	0	Bacteria;Firmicutes	0	0
Bacteria;Gracilibacteria	0	0	Bacteria;Gracilibacteria	0	0
Bacteria;Lentisphaerae	0	0	Bacteria;Lentisphaerae	0	0
Bacteria;Marinimicrobia (SAR406 clade)	0	0	Bacteria;Marinimicrobia (SAR406 clade)	0	0
Bacteria;Peregrinibacteria	0	0	Bacteria;Peregrinibacteria	0	0

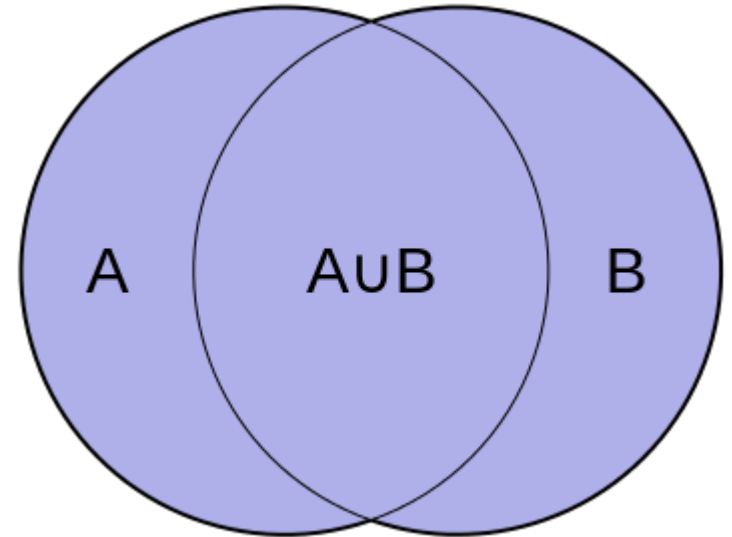
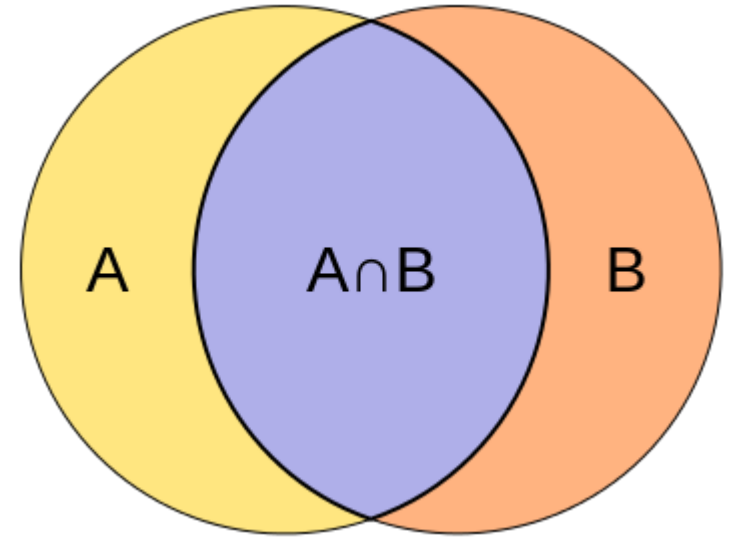
for a comprehensive list of available distance measures: *?vegdist*

# Jaccard-Distance

- Simple coefficient to measure the similarity between two (or more) communities
- Ranges between 0 and 1:
  - 1 = equal samples
  - 0 = no similarities between samples

R-code: `vegan::vegdist(x, method="jaccard")`

$$1 - J(A,B) = \frac{|A \cap B|}{|A \cup B|}$$



# Jaccard-Distance

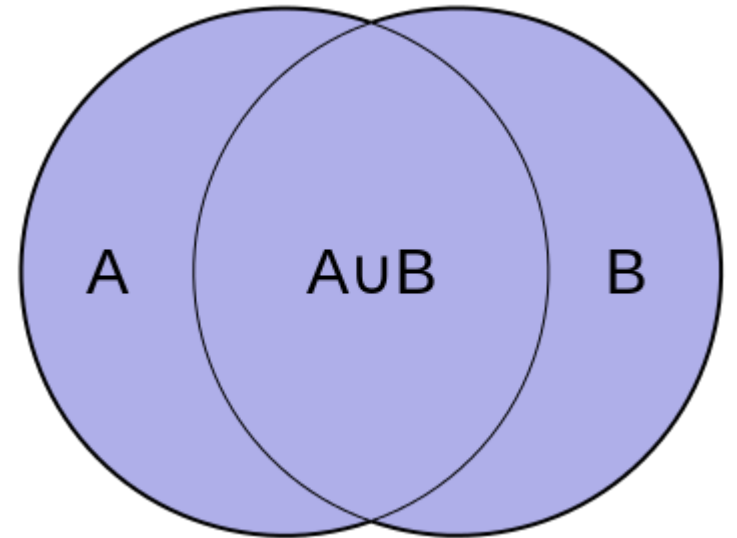
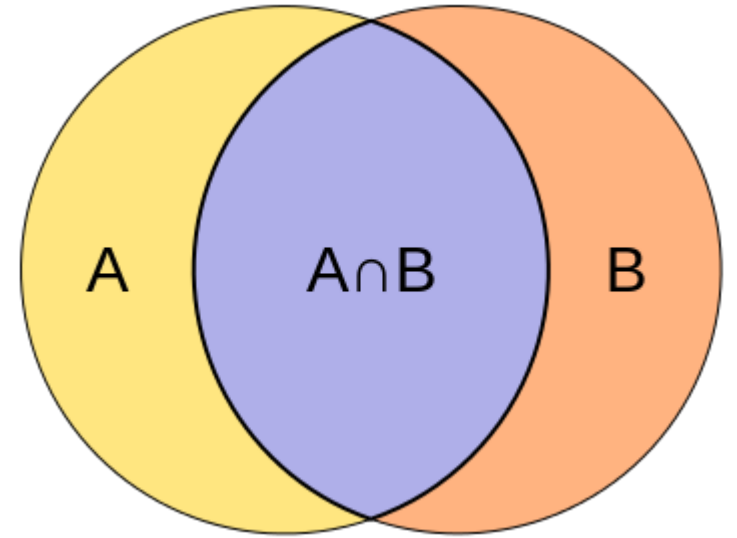
A



B



$$1 - J(A,B) = \frac{|A \cap B|}{|A \cup B|} = \frac{\text{grey starburst, blue circle, orange triangle}}{\text{blue circle, yellow star, green plus, orange triangle, grey starburst, red star, black X}}$$



# Bray-Curtis-(Dis)Similarity

- Compares Sites in terms of minimum abundance of species
- Widely used to compare raw (log transformed) species Data

$$D_{14}(x_1, x_2) = \frac{\sum_{j=1}^p |y_{1j} - y_{2j}|}{\sum_{j=1}^p (y_{1j} + y_{2j})}$$

Quadrats	Species				
	y <sub>1</sub>	y <sub>2</sub>	y <sub>3</sub>	y <sub>4</sub>	y <sub>5</sub>
x <sub>1</sub>	2	5	2	5	3
x <sub>2</sub>	3	5	2	4	3
x <sub>3</sub>	9	1	1	1	1

$$D_{14}(x_1, x_2) = \frac{1 + 0 + 0 + 1 + 0}{17 + 17} = 0.059$$

$$D_{14}(x_1, x_3) = \frac{7 + 4 + 1 + 4 + 2}{17 + 13} = 0.600$$

$$D_{14}(x_2, x_3) = \frac{6 + 4 + 1 + 3 + 2}{17 + 13} = 0.533$$

If the sum of species greatly varies among sites, the Bray-Curtis index may become negative!

-> sqrt-transformation

R-code: `vegan::vegdist(x, method="bray")`

# Sörensen, Sorensen, Sørensen or Dice index

Simply put: **Sörensen index = Bray-Crurtis** distance on **Presence/Absence** data

## **Difference to the Jaccard-Index:**

- more robust against outliers
- more sensitive in heterogeneous Datasets

R-code: `ecodist::distance(x, method=„sorensen“)`



# Euclidean Distance

- metric distance between two points (or vectors) in geometrical „euclidean space“
- symmetrical distance
- Range: zero to positive inf.

R-code: `stats::dist(x, method=„euclidean“)`

$$d(\mathbf{p}, \mathbf{q}) = d(\mathbf{q}, \mathbf{p}) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2 + \dots + (q_n - p_n)^2}$$
$$= \sqrt{\sum_{i=1}^n (q_i - p_i)^2}.$$

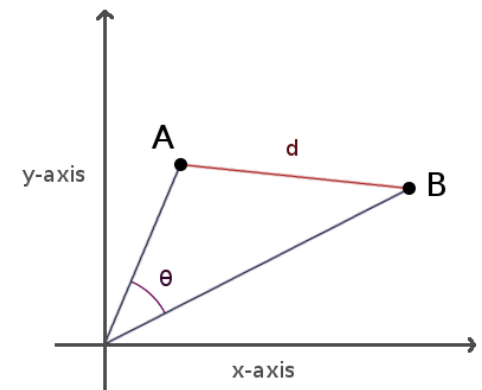
- Hellinger transformation for ecological data  
(used in PCA, PCoA, RDA, CCA)

$$\sqrt{\frac{\text{species count } x}{\text{site's total abundance}}}$$

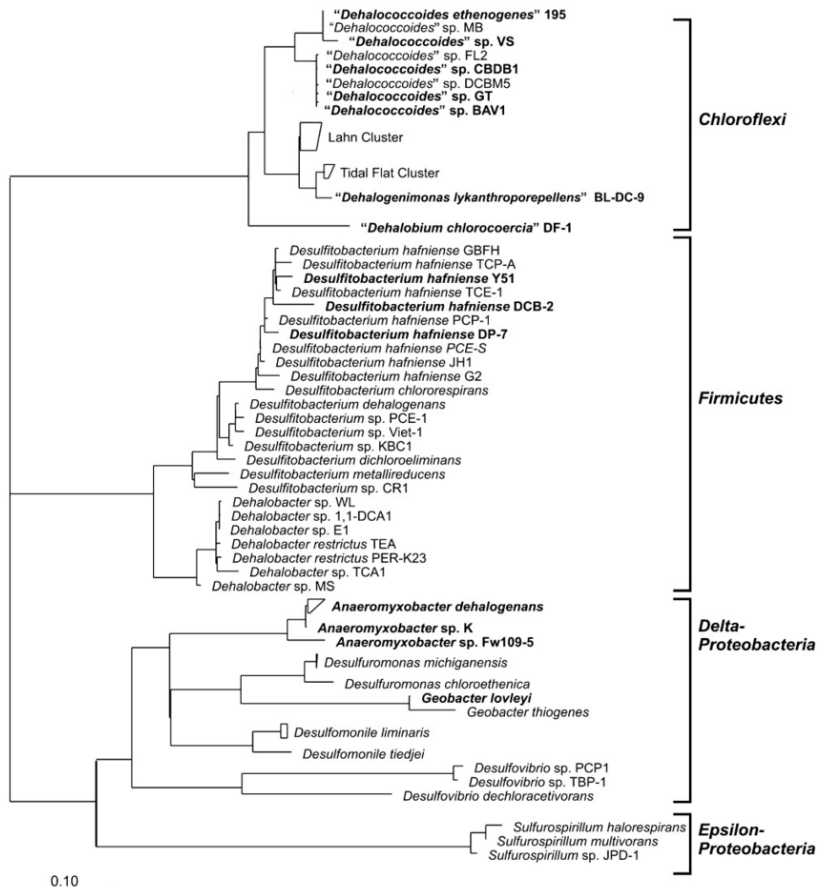
R-code: `vegan::decostand(x, „hel“), stats::dist(x.hel)`

- Chord transformation: euclidean distance normalized to 1

R-code: `vegan::decostand(x, „nor“); stats::dist(x.nor)`



# UniFrac (Lozupone & Knight 2005)



- Phylogenetic distance measure (especially for sequence based community assessments)
- requires a rooted tree auf the OTU-table

## R-code:

```
system("muscle3.8.31_i86win32.exe -in OTU_table.fasta -out OTU_table_red.afa")
```

```
system("muscle3.8.31_i86win32.exe -maketree -in OTU_table_red.afa -out  
OTU_table_red.phy -cluster neighborjoining")
```

(<http://www.drive5.com/muscle/downloads.htm>)

Available unifrac-measures: weighted, unweighted, variance adj.  
weighted, alpha = 0, alpha = .5

R-package: GUniFrac

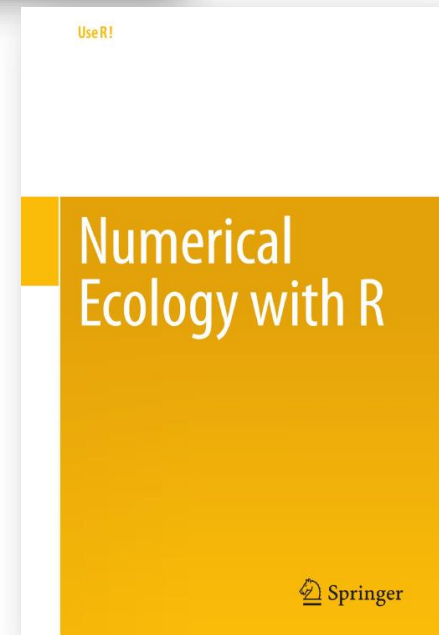
# PERMAVOVA (Permutational Multivariate Analysis of Variance Using Distance Matrices)

- **R-code:** `vegan::adonis(distance.table ~ Env1 + Env2, method = „bray“)`
- Permutational MANOVA explains variance within the dataset using environmental parameters

```
> adonis(mlog ~ Metal$Type + Metal$Day + Metal$Nuc, permutations = 1000)
> Call: adonis(formula = mlog ~ Metal$Type + Metal$Day + Metal$Nuc, permutations = 1000)
> Permutation: free
> Number of permutations: 1000
> Terms added sequentially (first to last)
>
  Df SumsOfSqs MeanSqs F.Model    R2      Pr(>F)
> Metal$Type 1  0.2214  0.22137   7.0298 0.06499 0.000999 ***
> Metal$Day   5  1.1336  0.22672   7.1998 0.33282 0.000999 ***
> Metal$Nuc   1  0.4451  0.44512  14.1351 0.13068 0.000999 ***
> Residuals 51  1.6060  0.03149
> Total      58  3.4061
>
> --- Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

# Further Information

- Legendre & Legendre 1998 - Numerical ecology
  - Mathematical background
- Borcard, Gillet & Legendre 2011 - Numerical Ecology with R
  - Application in R
- GUSTA ME: <https://sites.google.com/site/mb3gustame/reference/dissimilarity>



# Idea... Little Helper Scripts!

```
10 ▾ ##### Anderson-Darling Normality Test #####
11 library(nortest)
12
13 cat("Data\tp-value\n")
14 for(i in 1:12)
15 ▾ {
16   x = species[,i]
17   j = ad.test(x)
18   cat(Species_names[i], "\t", j$p.value, "\n")
19
20 }
21
22 # Testing for normal distribution of the Data; p>0.05 = normal distributed Data. If data are
23 # not normally distributed, Pearsons correlation is not an option
24
25 ▾ ##### Correlation Test #####
26 setwd("C:/Users/icbmadmin/Desktop/R/Transects HE425/Correlations/")
27
28 ##### Full Cor. Table
29 sink("Correlation_table.txt") # File in which all the console output will be saved
30 cat("Group\tParameter\tPearson Correlation\tp-value\n") # Columnnames for the file "\t" = tabstop
31 for(i in 1:12) # Species/whatever from a to b (in species table) to be tested against ...
32 ▾ {
33   x = Species[,i]
34
35   for(j in 3:9) # ...parameters in columns from x to y (in env-data-table)
36   ▾ {
37     y = ENV[,j]
38
39     cor = cor.test(x = x, y = y, method = "pearson")
40
41     cat(Species_names[i], "\t", ENV_names[j], "\t", cor$estimate, "\t", cor$p.value, "\n")
42
43   }
44 }
45 sink()
```