

Training in ade4 in R - Module I: Basic methods

Correspondence analysis and Principal Coordinates Analysis

Stéphane Dray

2021-04-20

Data structure



- One table with n rows and m columns
- Data are counts of individuals (sums by row and columns are meaningful)
- For instance
 - sites species

Objectives

- Identify what is the main information contained in the table
 - Identify the principal differences/similarities between row categories
 - Identify the principal differences/similarities between column categories
 - Identify the principal differences/similarities between row and column categories

Data

We consider the **meaudret** data set

```
library(ade4)
data(meaudret)
names(meaudret)
```

```
## [1] "env"      "design"    "spe"      "spe.names"
```

```
dim(meaudret$spe)
```

```
## [1] 20 13
```

```
head(meaudret$spe.names)
```

```
## [1] "Ephemera_danica" "Baetis_sp"      "Baetis_rhodani" "Baetis_niger"
## [6] "Centropilum_sp"
```

The data set contains the abundances of 13 Ephemeroptera species in 20 samples. The measurements have been made in 6 sites at each season along a small French stream (see [?meaudret](#))

```
head(meaudret$design)
```

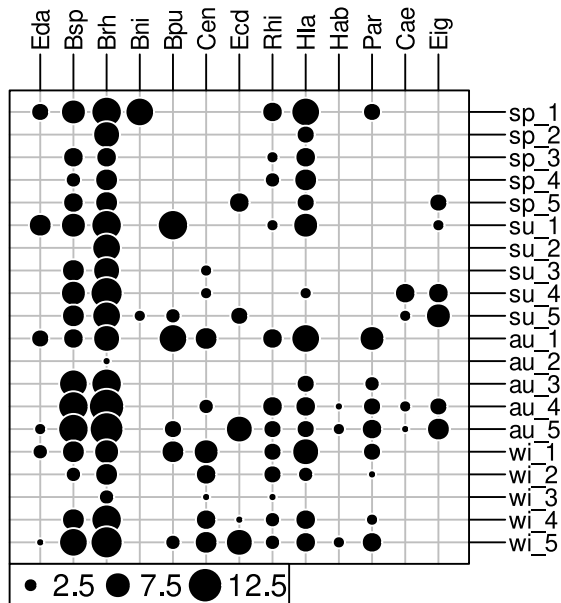
```
##      season site
## sp_1 spring  S1
## sp_2 spring  S2
## sp_3 spring  S3
## sp_4 spring  S4
## sp_5 spring  S5
## su_1 summer  S1
```

We want to know

- which species have similar distributions
- which sites have similar composition
- which species are mainly present in which sites

Contingency table

```
library(adegraphics)
table.value(meaudret$spe, symbol
```



```
head(rowSums(meaudret$spe))
```

```
## sp_1 sp_2 sp_3 sp_4 sp_5 su_1
##    48   12   17   18   24   44
```

```
head(colSums(meaudret$spe))
```

```
## Eda Bsp Brh Bni Bpu Cen
##   20 104 163  11  35  37
```

```
sum(meaudret$spe)
```

```
## [1] 595
```

Chi-square test

The test allows to measure and evaluate the significance of the association between species and sites (the null hypothesis is the random distribution)



```
chisq.test(meaudret$spe)
```

```
## Warning in chisq.test(meaudret$spe): Chi-squared approximation may be inco
```

```
##
```

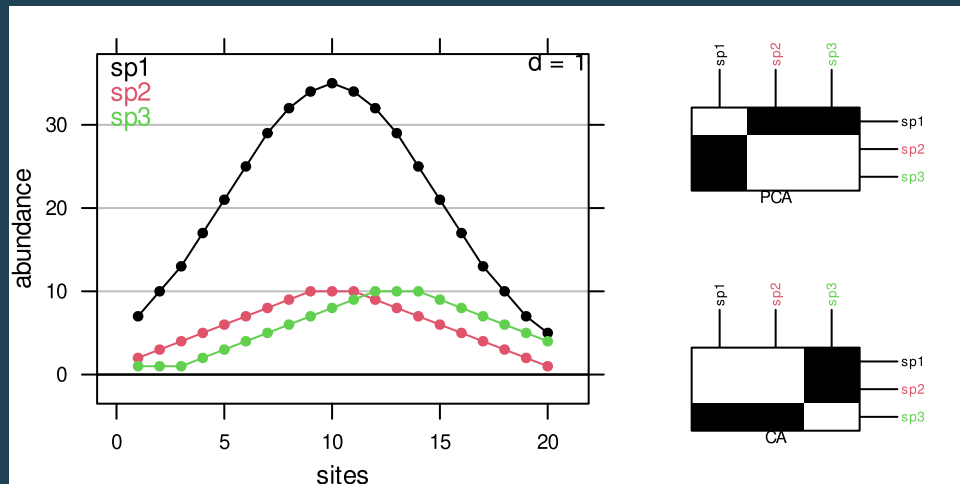
```
##      Pearson's Chi-squared test
```

```
##
```

```
## data:  meaudret$spe
```

```
## X-squared = 534.52, df = 228, p-value < 2.2e-16
```


Absolute and relative frequencies



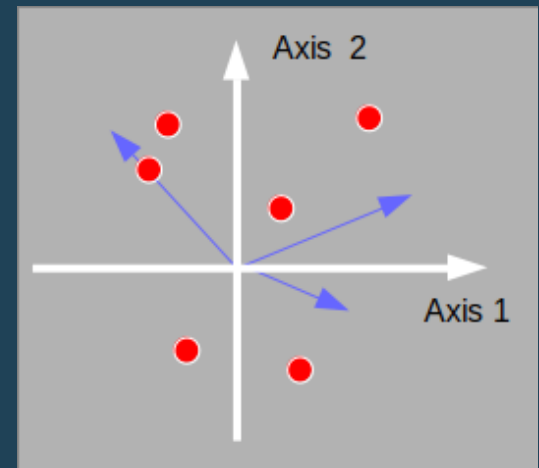
- In PCA (Euclidean distances), species 2 and 3 are closer
- In COA (distances), species 1 and 2 are closer (null distance)

Correspondence analysis

- $X = \frac{D}{PD} \frac{1}{1}$ is the transformed and centred table of relative frequencies with P
- $Q = \frac{D}{D}$ where D contains the column category frequencies
- $D = \frac{D}{D}$ where D contains the row category frequencies



dudi.coa



Maximized criteria

- For rows

$$\mathbf{a} \quad \mathbf{D} \quad \mathbf{P} \quad \mathbf{D} \quad \mathbf{D} \quad \mathbf{a}_D \quad \mathbf{D} \quad \mathbf{P} \quad \mathbf{a}_D$$

In this viewpoint, columns have a unit-variance score \mathbf{a} that maximises the variance between the row barycenters.

- For columns

$$\mathbf{D} \quad \mathbf{P} \quad \mathbf{D} \quad \mathbf{D} \quad \mathbf{b}_D \quad \mathbf{D} \quad \mathbf{P} \quad \mathbf{b}_D$$

In this viewpoint, rows have a unit-variance score \mathbf{b} that maximises the variance between the column barycenters.

The `dudi.coa` function

Arguments

```
args(dudi.coa)
```

```
## function (df, scannf = TRUE, nf = 2)  
## NULL
```

- `df` is a `data.frame` with the positive values (counts)
- `scannf` and `nf` allow to set the number of dimensions to interpret

```
coa.meau <- dudi.coa(meaudret$spe, scannf = FALSE)
```

Returned values

```
names(coa.meau)
```

```
## [1] "tab" "cw" "lw" "eig" "rank" "nf" "c1" "li" "co" "l1"
```

It returns an object of class **dudi** containing:

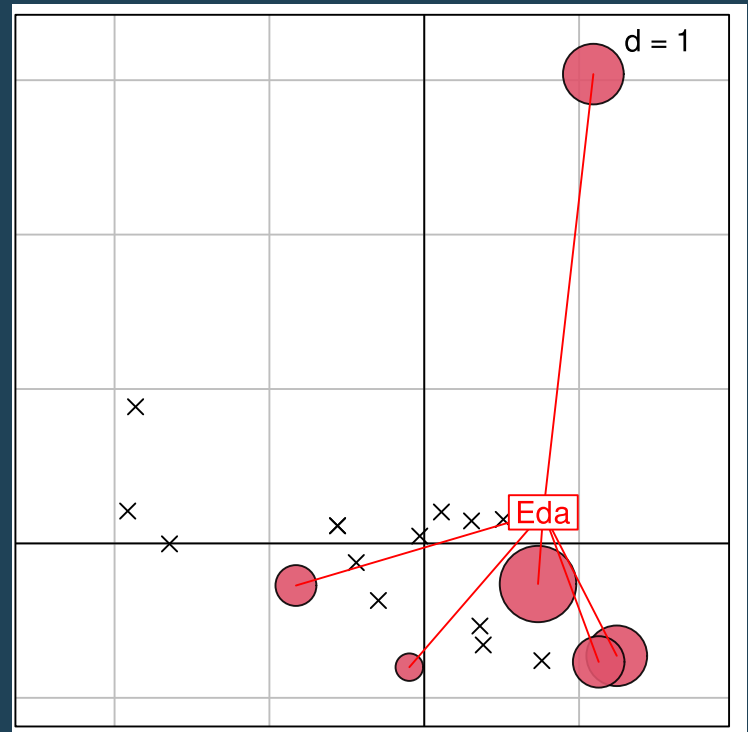
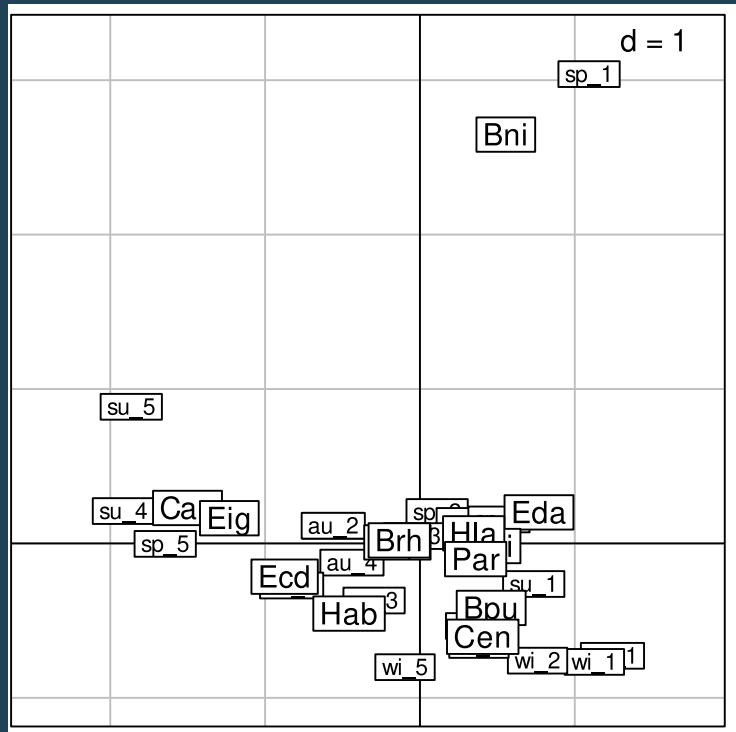
- **\$eig**: eigenvalues (Λ)
- **\$cw**: column weights (\mathbf{D})
- **\$lw**: row weights (\mathbf{D})
- **\$tab**: centred relative frequencies table (\mathbf{D} \mathbf{P} \mathbf{D})
- **\$c1**: unit-variance column scores (\mathbf{A})
- **\$li**: row scores as weighted averages (\mathbf{L} \mathbf{D} \mathbf{P} \mathbf{A})
- **\$l1**: unit-variance row scores (\mathbf{B})
- **\$co**: column scores as weighted averages (\mathbf{C} \mathbf{D} \mathbf{P} \mathbf{B})
- **\$N**: total sum ()

Graphical representations

Biplot can be produced for CA using the `biplot` function. Three types of biplots can be produced using the argument `method`

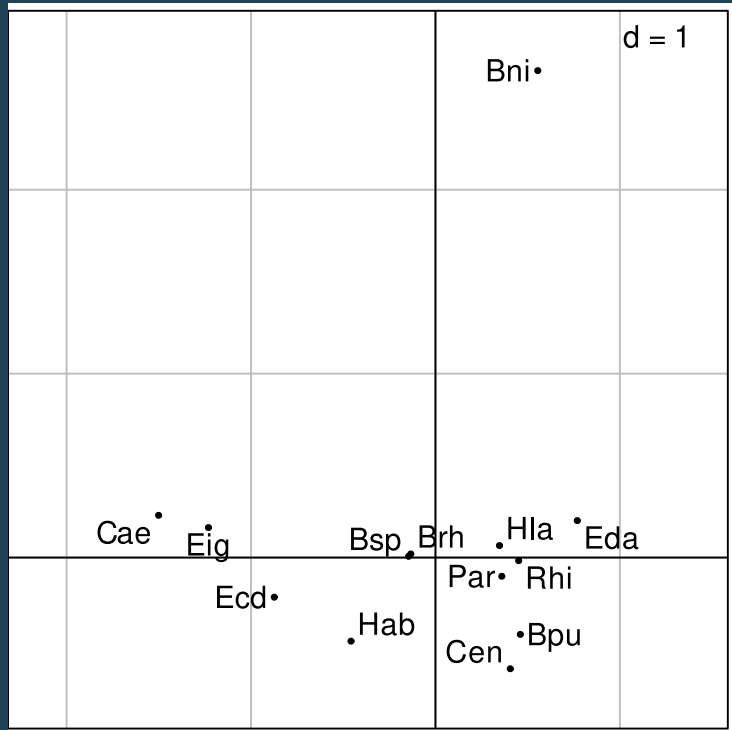
- If `method = 2`, species are positioned by a unit-variance score (`$c1`) and sites by weighted averaging (`li`).
- If `method = 3`, sites are positioned by a unit-variance score (`$l1`) and species by weighted averaging (`$co`).
- By default, `method = 1` corresponds to a compromise between these two representations (`$li` and `$co`).

Weighted averaging

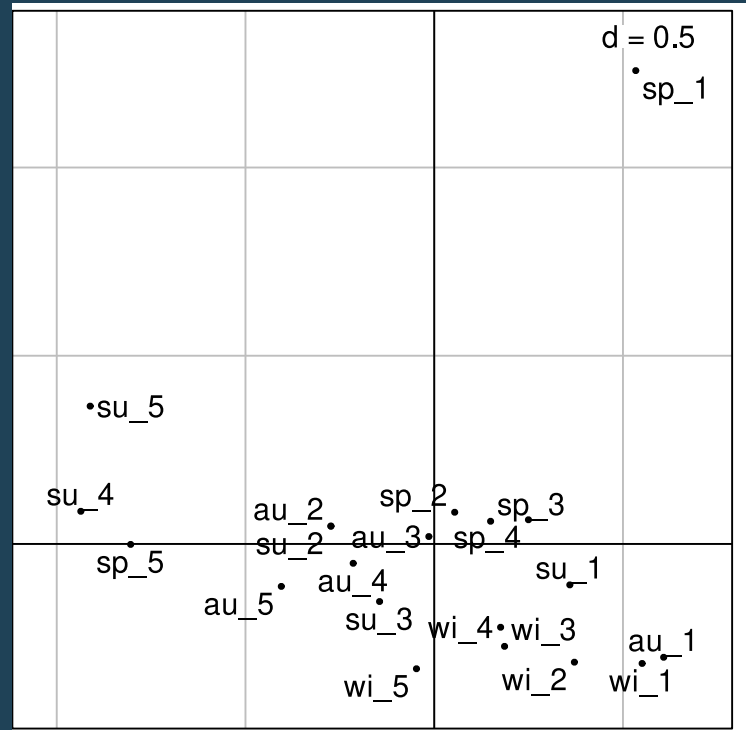


Separate representations

```
s.label(coa.meau$co, plabels.optf
```

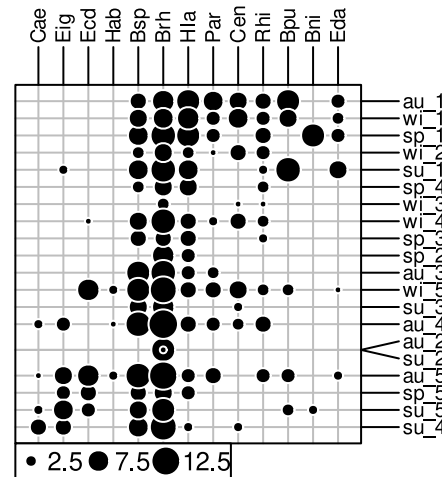
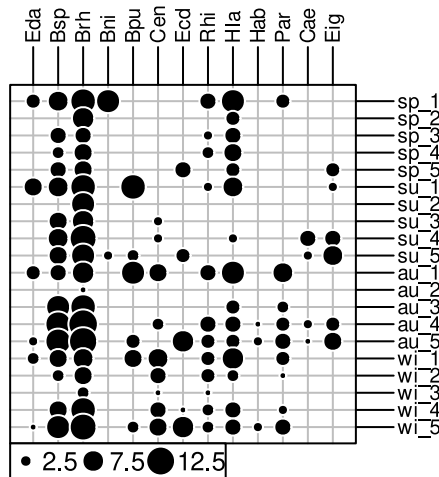


```
s.label(coa.meau$li, plabels.optf
```



Reordering of a table

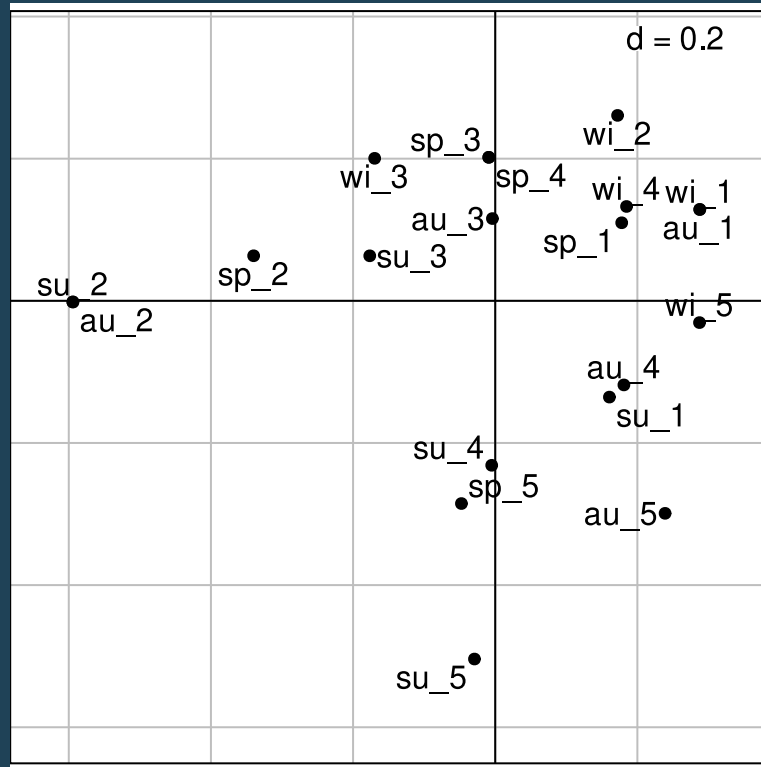
```
g1 <- table.value(meaudret$spe, ppoint.cex = 0.5, symbol = "circle",
  plot = FALSE)
g2 <- table.value(meaudret$spe, ppoint.cex = 0.5, symbol = "circle",
  coordsx = rank(coa.meau$co[, 1]), coordsy = rank(coa.meau$li[,
    1]), plot = FALSE)
cbindADEg(g1, g2, plot = TRUE)
```



Principal coordinates analysis

- PCA, CA methods induce implicitly a way to compute distances
- Several other distances have been proposed (e.g., genetic, presence-absence)
- PCoA takes a distance matrix as input and returns coordinates in a low dimensional space that best preserve the original distances.
- 😊 it allows to choose a particular distance measure between sites (or species).
- ☹ it focuses either on individuals or variables, not both.
- Useful if distances are directly recorded or computed from raw data tables

```
dJ <- dist.binary(meaudret$spe, method = 1) # Jaccard
pcoJ <- dudi.pco(dJ, scannf = FALSE)
s.label(pcoJ$li, plabels.optim = TRUE)
```



COA in practice

[Go to practical 4](#)