


# FACTOMINER<sup>R</sup>

An R package for exploratory data analysis for teaching and research

François Husson, Julie Josse & Sébastien Lê



# Why **FACTOMINER** ?

- To make exploratory multivariate data analysis with a free software 
- The possibility to propose new methods (taking into account different structure on the data)
- To have a package user friendly and oriented to practitioner (a very easy GUI)

# 1 – The classical methods

- Methods implemented are similar in their main objective: to sum up and simplify the data by reducing the dimensionality of the dataset
  - ❑ Continuous variables: Principal Components Analysis
  - ❑ Contingency table: Correspondence Analysis
  - ❑ Categorical variables: Multiple Correspondence Analysis
  - ❑ Continuous and categorical variables: Mixed Data Analysis

# PCA Example

Data : performances of 41 athletes during two meetings of decathlon

	100m	Long.jump	Shot.put	High.jump	400m	110m.hurdle	Discus	Pole.vault	Javeline	1500m	Rank	Points	Competition
SEBRLE	11.04	7.58	14.83	2.07	49.81	14.69	43.75	5.02	63.19	291.70	1	8217	Decastar
CLAY	10.76	7.40	14.26	1.86	49.37	14.05	50.72	4.92	60.15	301.50	2	8122	Decastar
KARPOV	11.02	7.30	14.77	2.04	48.37	14.09	48.95	4.92	50.31	300.20	3	8099	Decastar
BERNARD	11.02	7.23	14.25	1.92	48.93	14.99	40.87	5.32	62.77	280.10	4	8067	Decastar
YURKOV	11.34	7.09	15.19	2.10	50.42	15.31	46.26	4.72	63.44	276.40	5	8036	Decastar
Sebrle	10.85	7.84	16.36	2.12	48.36	14.05	48.72	5.00	70.52	280.01	1	8893	OlympicG
Clay	10.44	7.96	15.23	2.06	49.19	14.13	50.11	4.90	69.71	282.00	2	8820	OlympicG
Karpov	10.50	7.81	15.93	2.09	46.81	13.97	51.65	4.60	55.54	278.11	3	8725	OlympicG
Macey	10.89	7.47	15.73	2.15	48.97	14.56	48.34	4.40	58.46	265.42	4	8414	OlympicG
Warners	10.62	7.74	14.48	1.97	47.97	14.01	43.73	4.90	55.39	278.05	5	8343	OlympicG

# PCA example

➤ Introduction of supplementary information:

- supplementary continuous variables

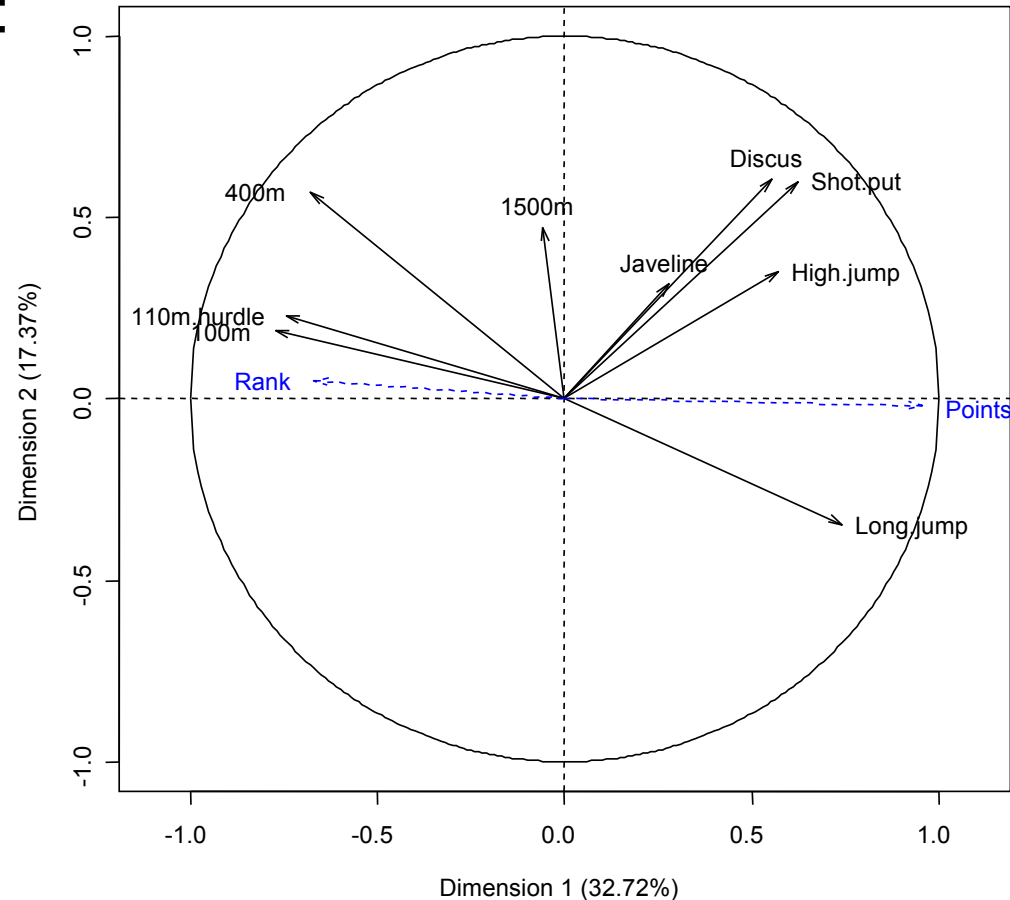
➤ Graphs enriched by :

- representing the variables according to their quality of representation

➤ Indicators:

- contribution
- quality of representation

Variables factor map (PCA)



# PCA example

## ➤ Introduction of supplementary information:

- supplementary individuals
- supplementary categorical variables

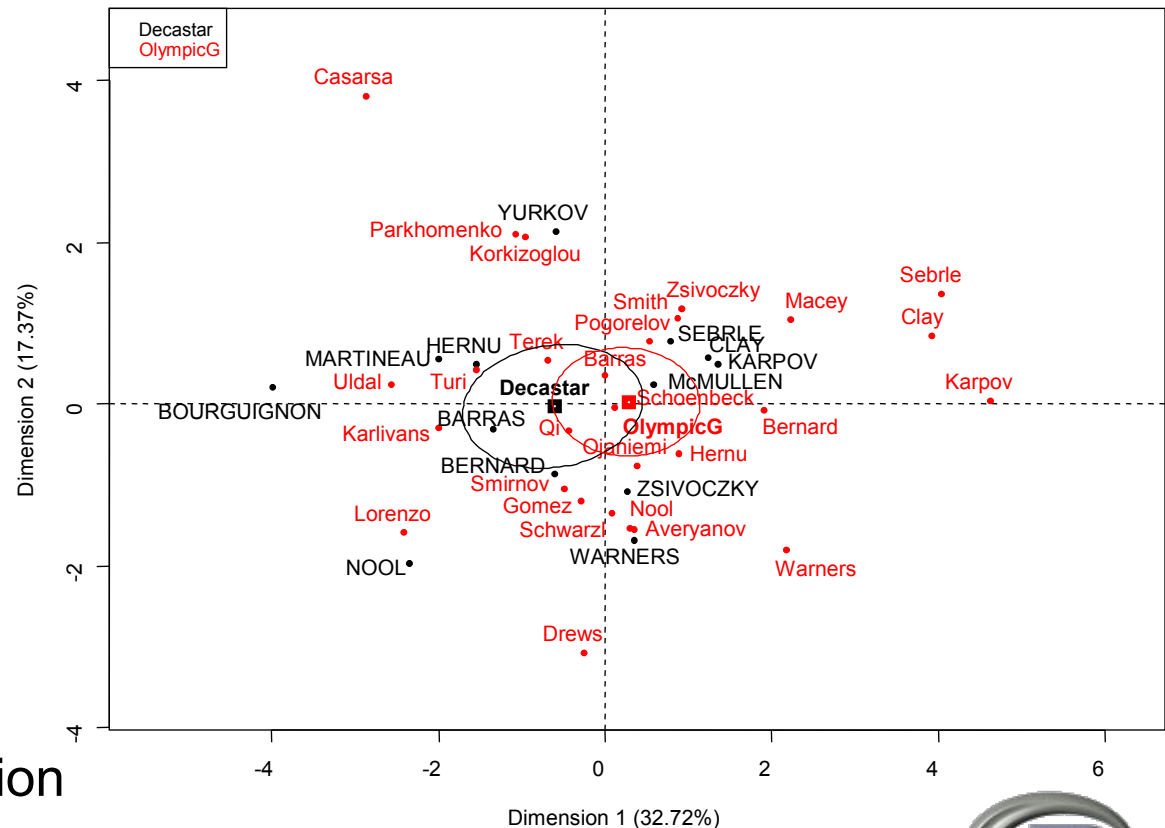
## ➤ Graphs enriched by:

- coloring according to supplementary information
- confidence ellipses around the categories

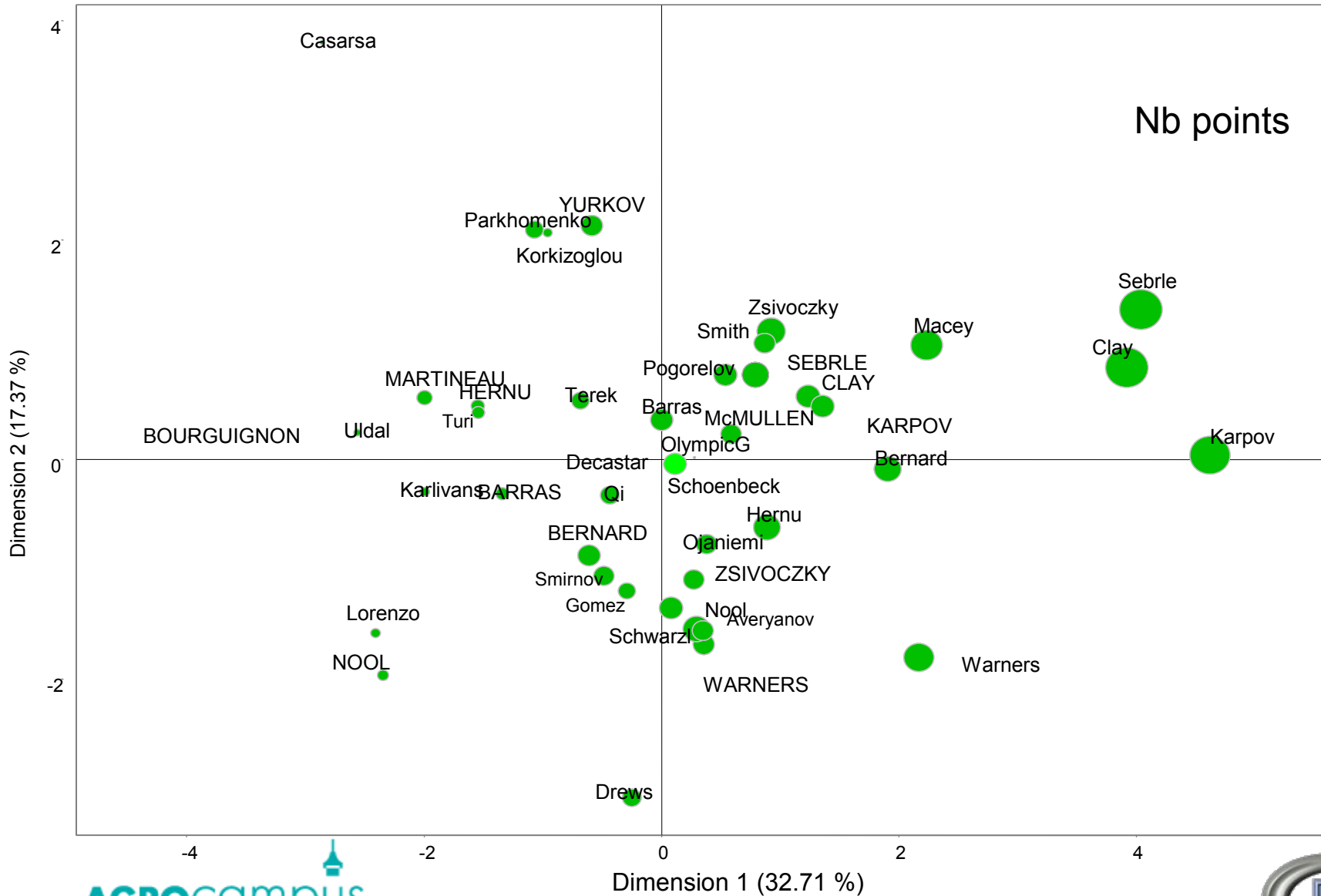
## ➤ Indicators:

- contribution
- quality of representation

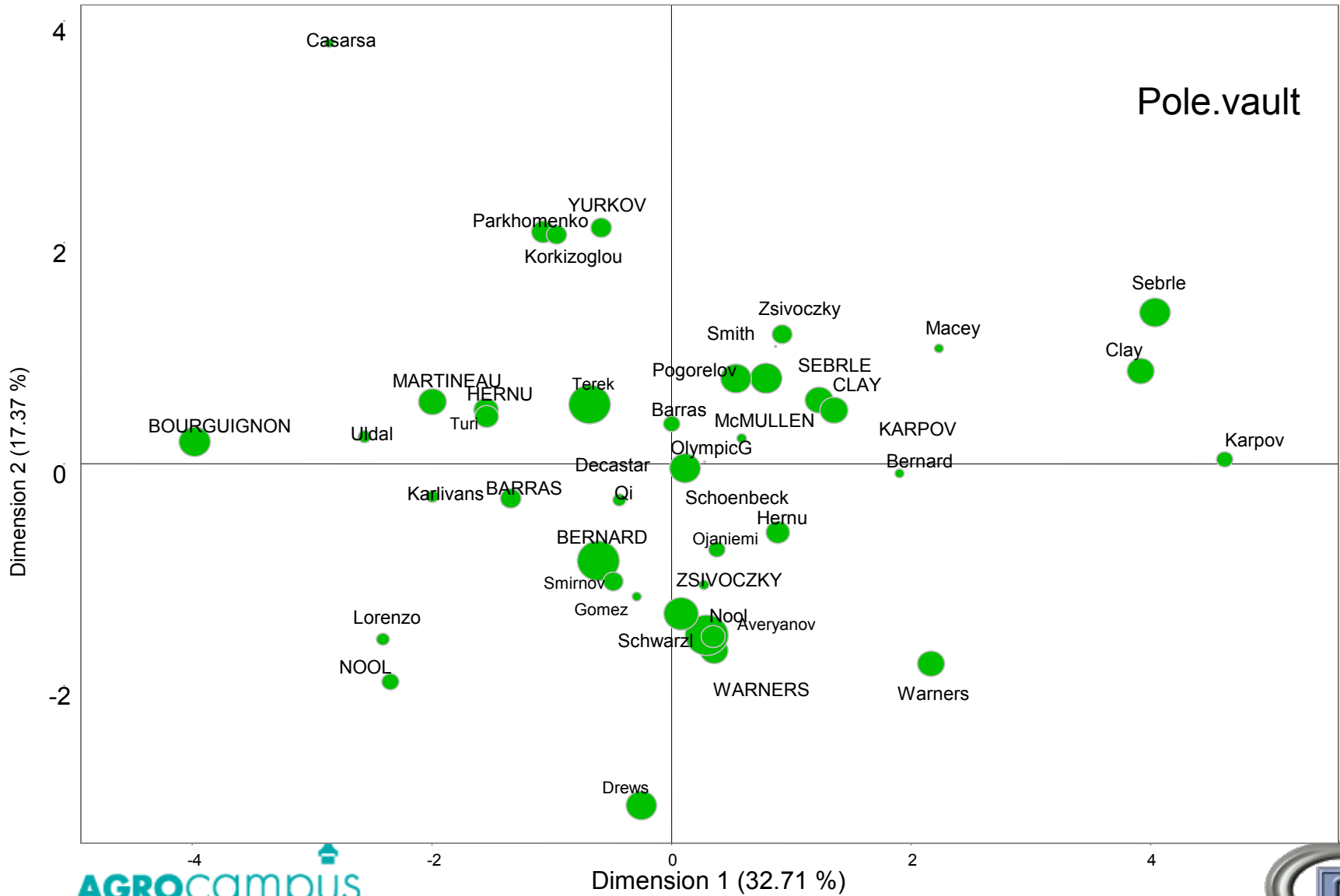
Individuals factor map (PCA)



# PCA example



# PCA example





# Description of the dimensions

## ➤ By the quantitative variables:

- The correlation between each variable and the coordinate of the individuals on the axis  $s$  is calculated
- The correlation coefficients are sorted
- Only the significant correlations are given

	\$Dim.1		\$Dim.2
	\$Dim.1\$quanti		\$Dim.2\$quanti
Best variable to describe the 1 <sup>st</sup> dimension		Dim.1	Dim.2
→ <b>Points</b>		<b>0.96</b>	<b>Discus 0.61</b>
	Long.jump	0.74	Shot.put 0.60
	Shot.put	0.62	
	<b>Rank</b>	<b>-0.67</b>	
	400m	-0.68	
	110m.hurdle	-0.75	
	100m	-0.77	

Significant level = 0.05

# Description of the dimensions

## ➤ By the qualitative variables:

- Perform a one-way analysis of variance with the coordinates of the individuals on the axis explained by the qualitative variable

\$Dim.1\$quali

	P-value
Competition	0.155

- A *F*-test by variable

\$Dim.1\$category

	Estimate	P-value
OlympicG	0.4393	0.155
Decastar	-0.4393	0.155

- For each category, a student *T*-test to compare the average of the category with the general mean

Significant level = 0.2

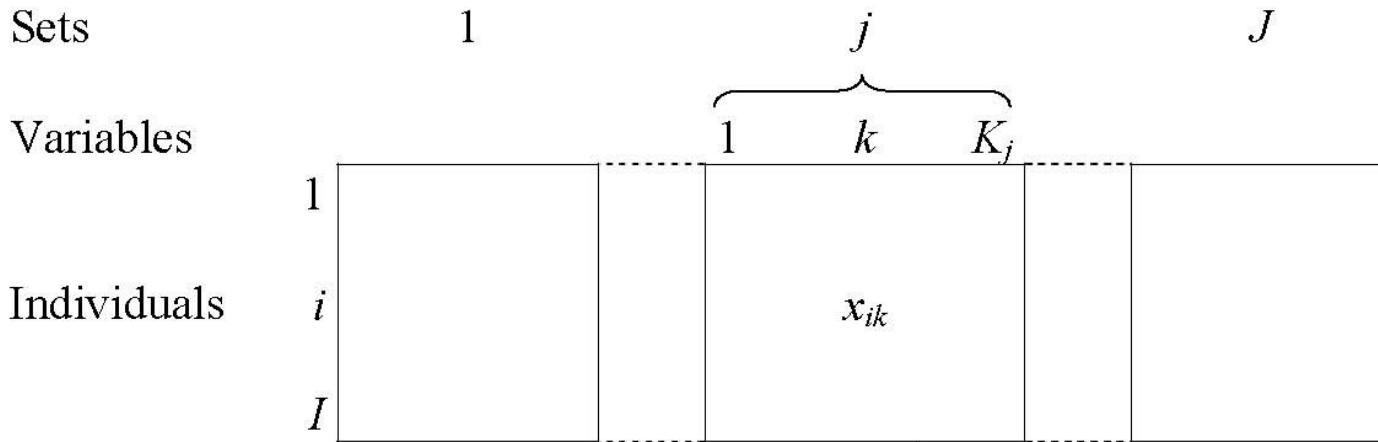
## 2 – Structure on the data

Different structure on the data are proposed:

- ❑ a partition on the variables: several sets of variables are simultaneously studied: **Multiple Factor Analysis, Generalized Procrustes Analysis**
- ❑ a hierarchy on the variables: variables are grouped and subgrouped (like in questionnaires structured in topics and subtopics): **Hierarchical Multiple Factor Analysis**
- ❑ a partition on the individuals: several sets of individuals described by the same variables: **Dual Multiple Factor**

Analysis

# Groups of variables (MFA)

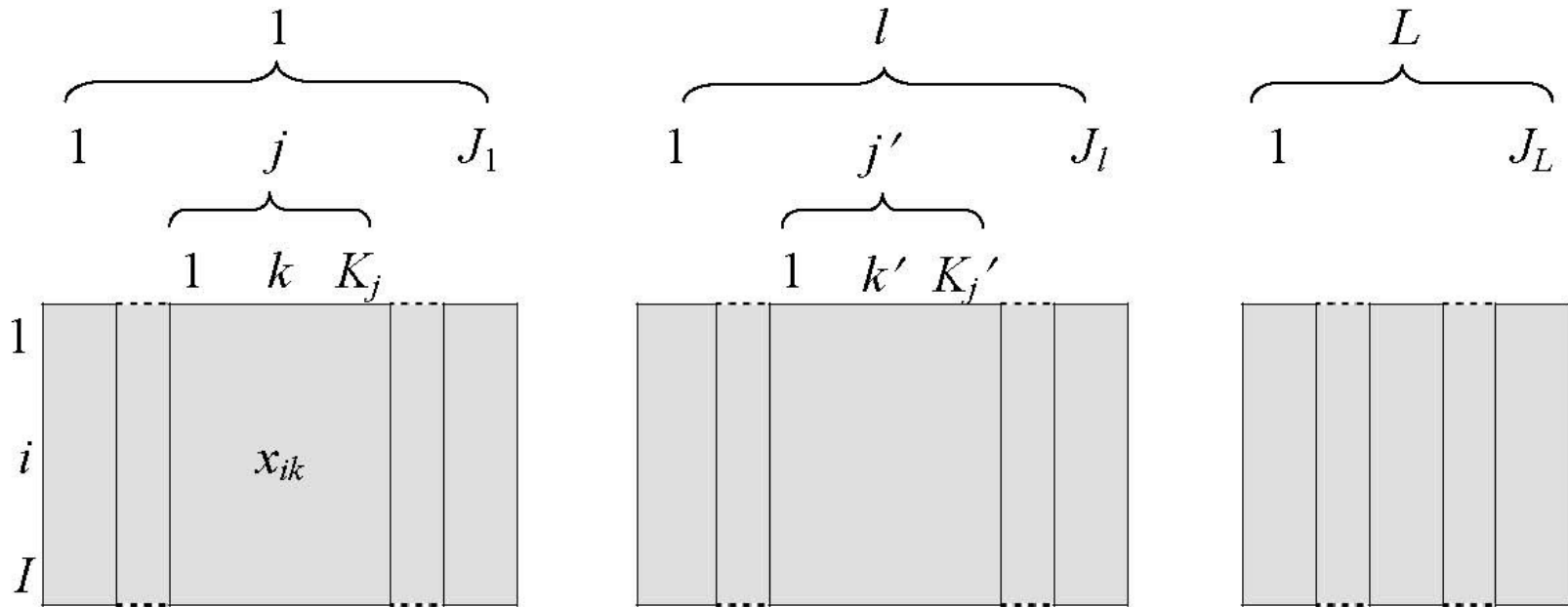


Groups of variables are quantitative and/or qualitative

- Objectives :
- study the link between the sets of variables
  - balance the influence of each group of variables
  - give the classical graphs but also specific graphs: groups of variables - partial representation

- Examples :
- Genomic: DNA, protein
  - Sensory analysis: sensorial, physico-chemical
  - Comparison of coding (quantitative / qualitative)

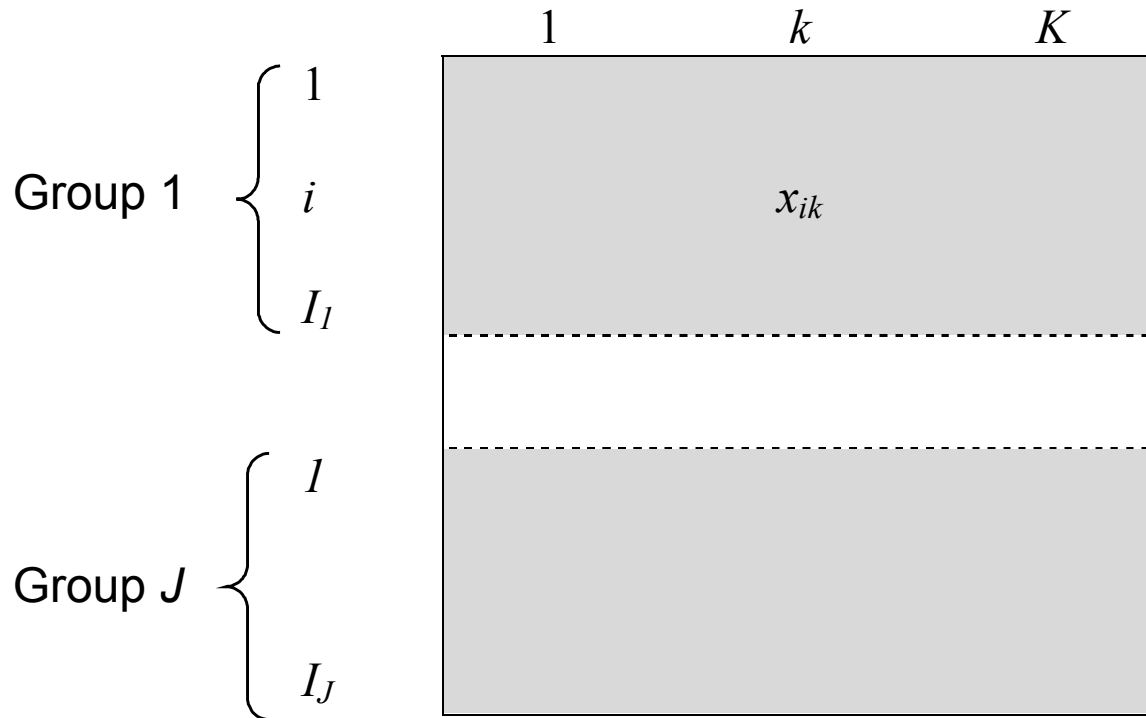
# Hierarchy on the variables (HMFA)



Two levels for the hierarchy: the first one contains  $L$  groups, each  $l$  group contains  $J_l$  subgroups, and each subgroup have  $K_j$  variables

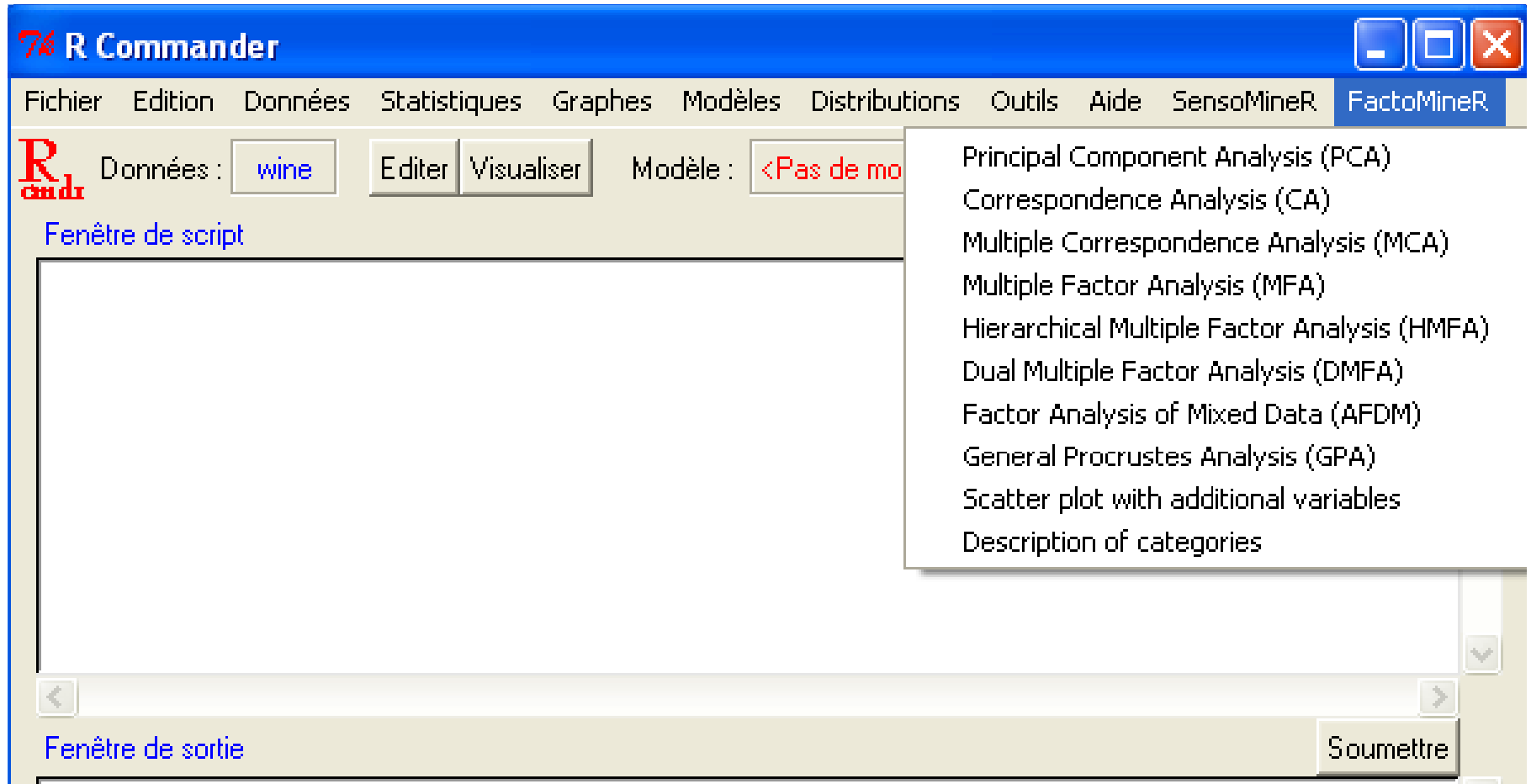
Objective: to balance the groups and the subgroups of variables

# Partition on the individuals (DMFA)



Objective: to compare the covariance matrices

# 3 – Graphical User Interface



Menu of the FactoMineR GUI

# 3 – Graphical User Interface

Main window  
of the PCA

**PCA**

## Principal Components Analysis (PCA)

Select active variables (by default all the variables are active)

- X100m
- Long.jump
- Shot.put
- High.jump
- X400m
- X110m.hurdle
- Discus
- Pole.vault
- Javeline
- X1500m

Modify supplementary factors    Modify supplementary variables    Select supplementary individuals

Graphical options    Outputs    Restart

**Main options**

Name of the result object:

Number of dimensions:

Scaled the variables of the group:

Graphical output : select the dimensions :

Apply

OK    Annuler    Aide



# 3 – Graphical User Interface

Graphical options

**74 Graphical options**

**Plot individuals graph**

Title of the graph

Hide some elements:

ind  ind sup  quali

Label for the active individuals

Label for the supplementary factor

Color of the active individuals  Change Color

Color for factors  Change Color

Coloring for individuals

by.individual  
Competition

x limits of the graph:

y limits of the graph:

**Plot variables graph**

Title of the graph

Draw variables with a  $\cos^2 >$ :

Labels for the active variables

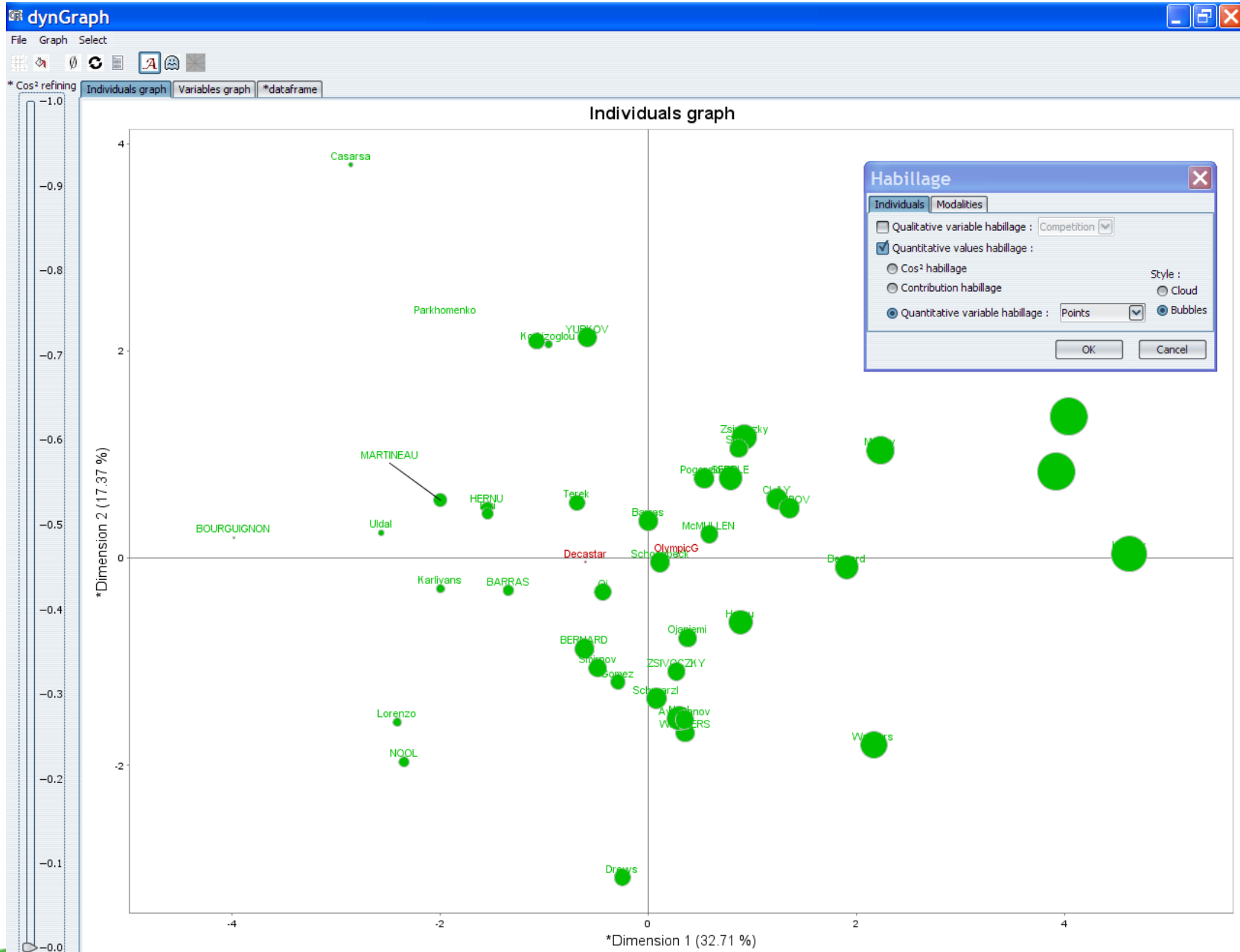
Labels for the supplementary variables

Color for active variables  Change Color

Color for supplementary variables  Change Color

OK Annuler Aide

# 3 – Graphical User Interface



# 4 – Conclusion

For researchers, practitioners and students: with classical and advanced methods

The FactoMineR package is available on the CRAN

The GUI can be simply loaded:

```
source("http://factominer.free.fr/install-facto.r")
```

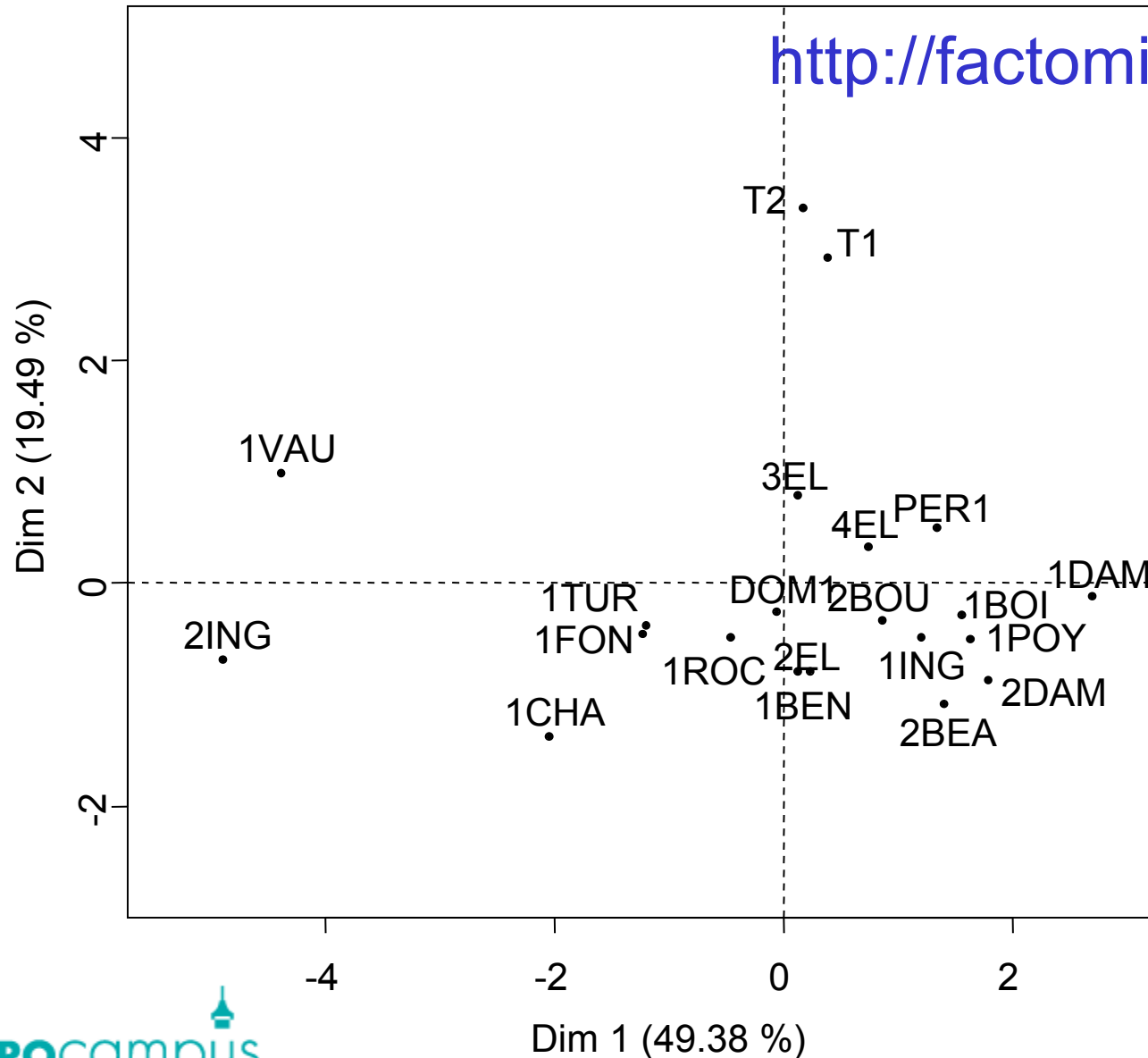
A website is dedicated to this package: <http://factominer.free.fr>

Future: dynamical graphs

Perspective: UseR!2008 (2 tutorials), UseR!2009 at Rennes

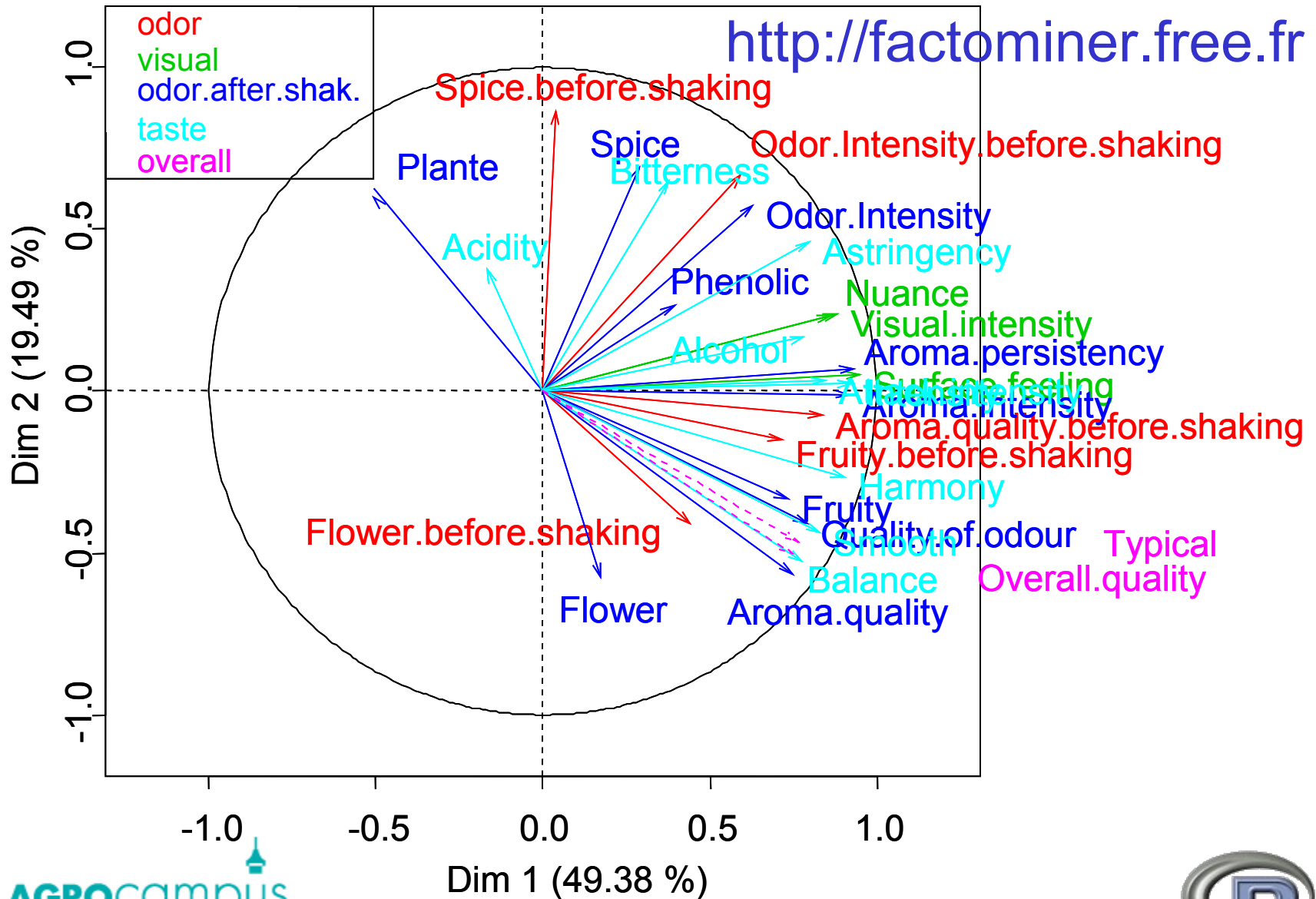
# MFA example: representation of the individuals

<http://factominer.free.fr>



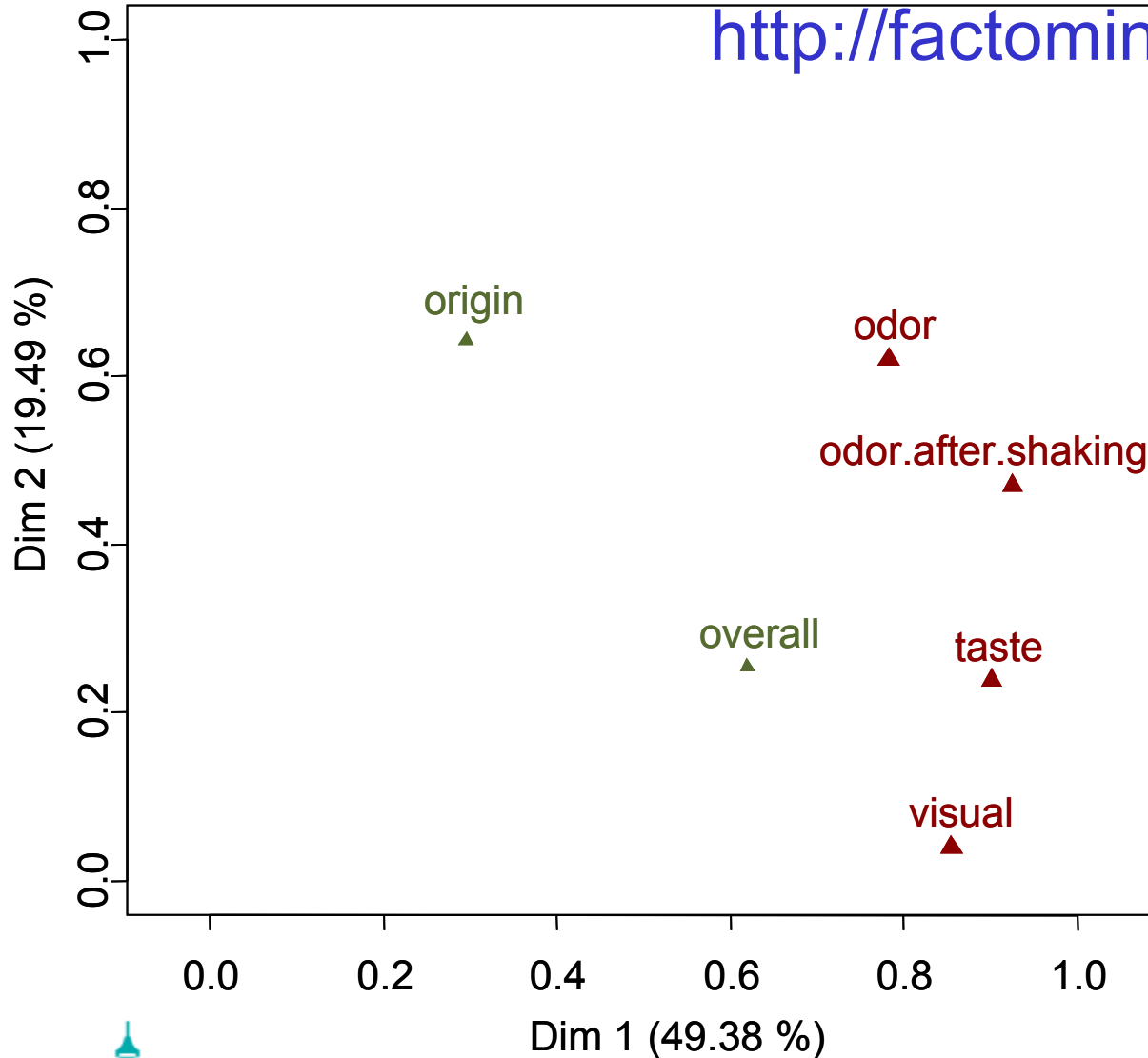
# MFA example: representation of the variables

<http://factominer.free.fr>

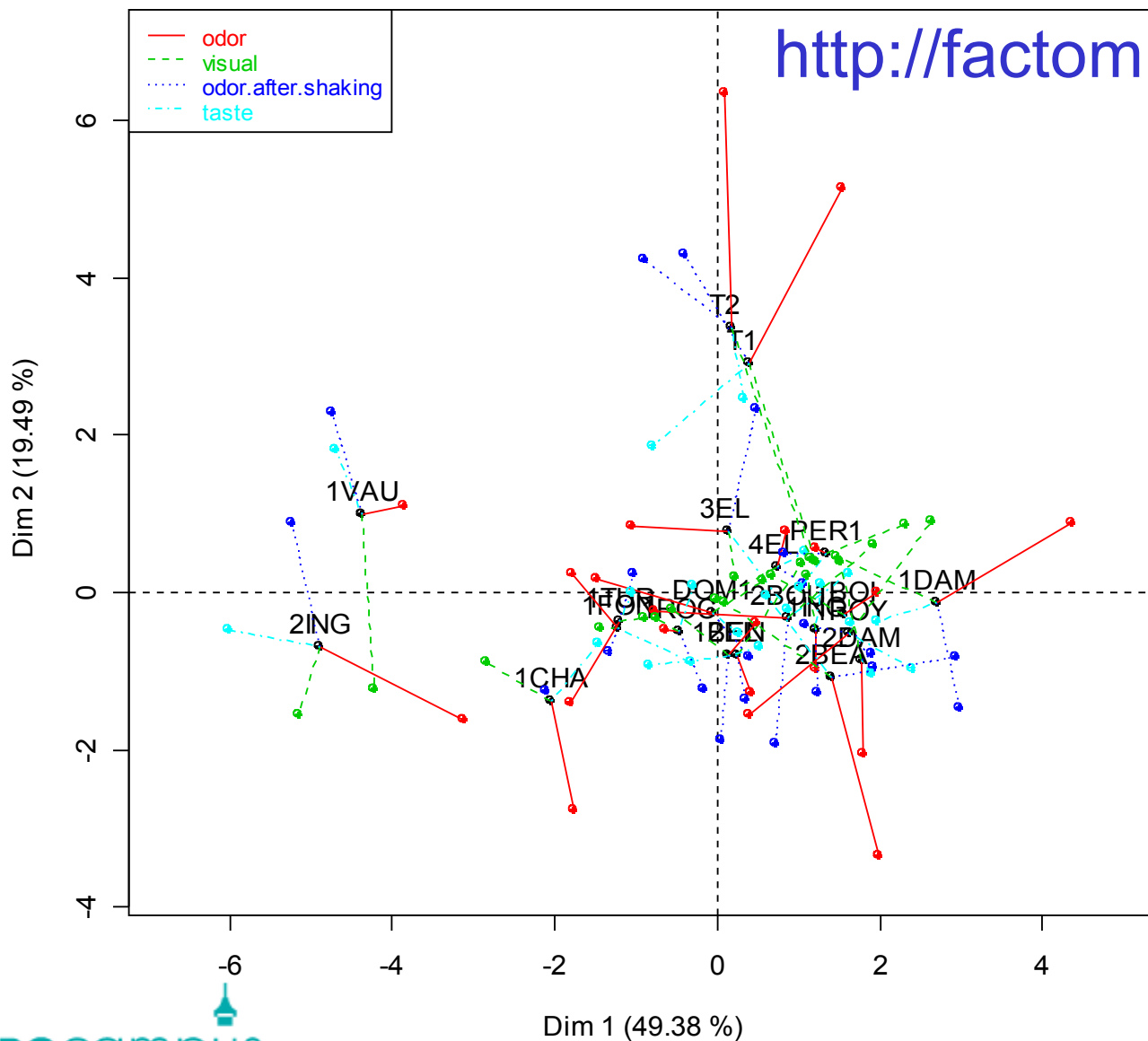


# MFA example: representation of the groups

<http://factominer.free.fr>



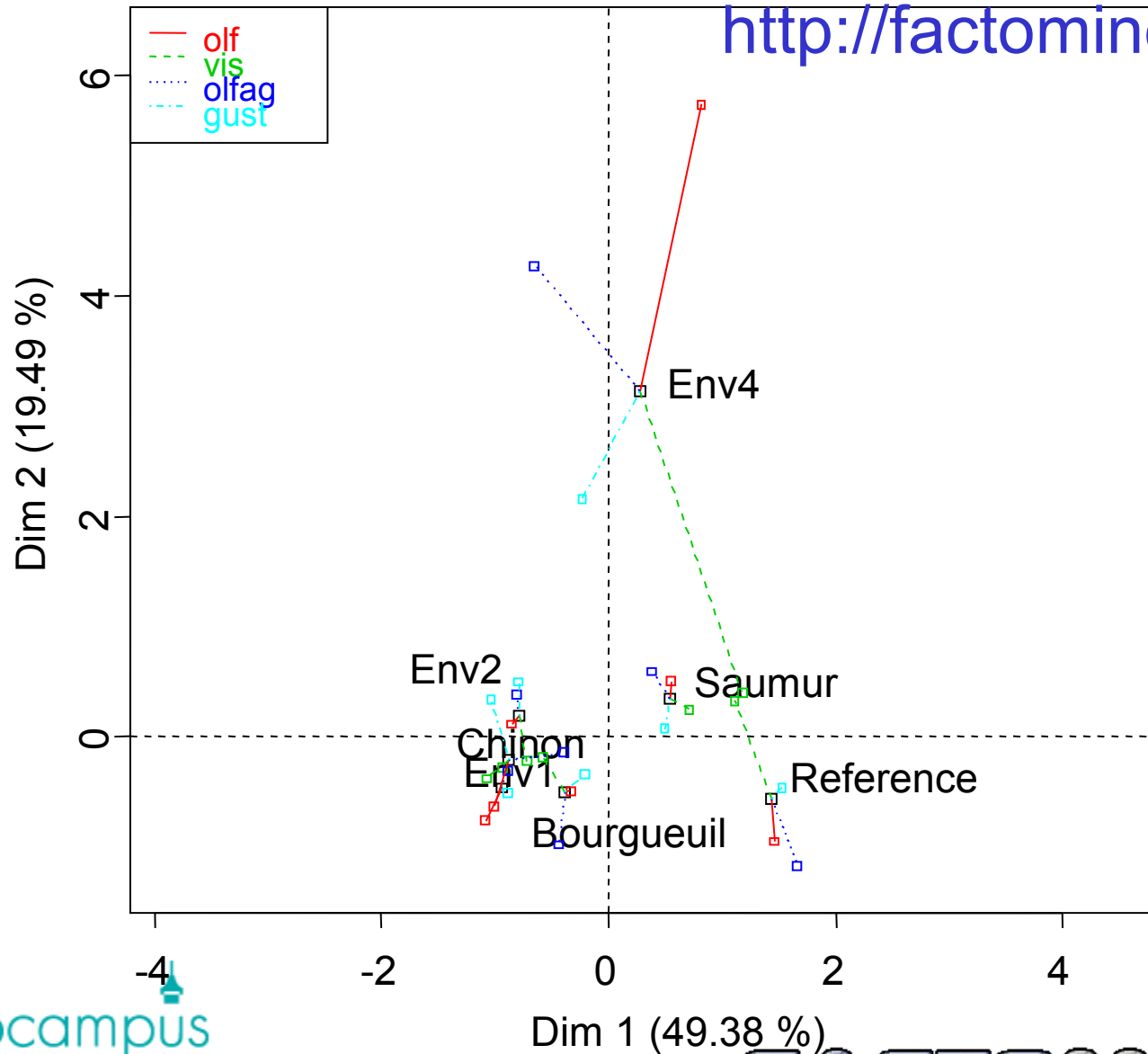
# MFA example: representation of the partial points



<http://factominer.free.fr>

# MFA example: representation of the partial points

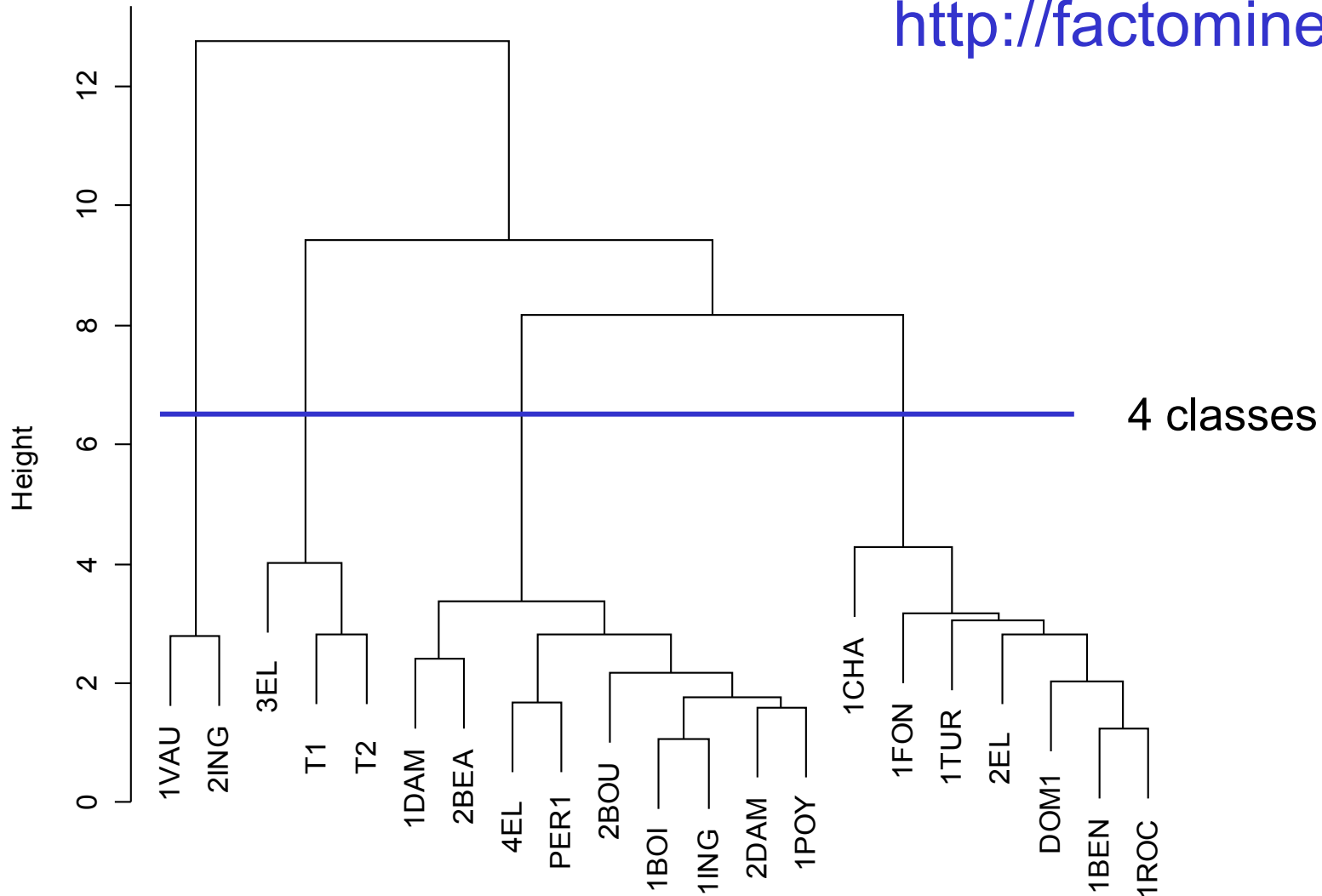
<http://factominer.free.fr>





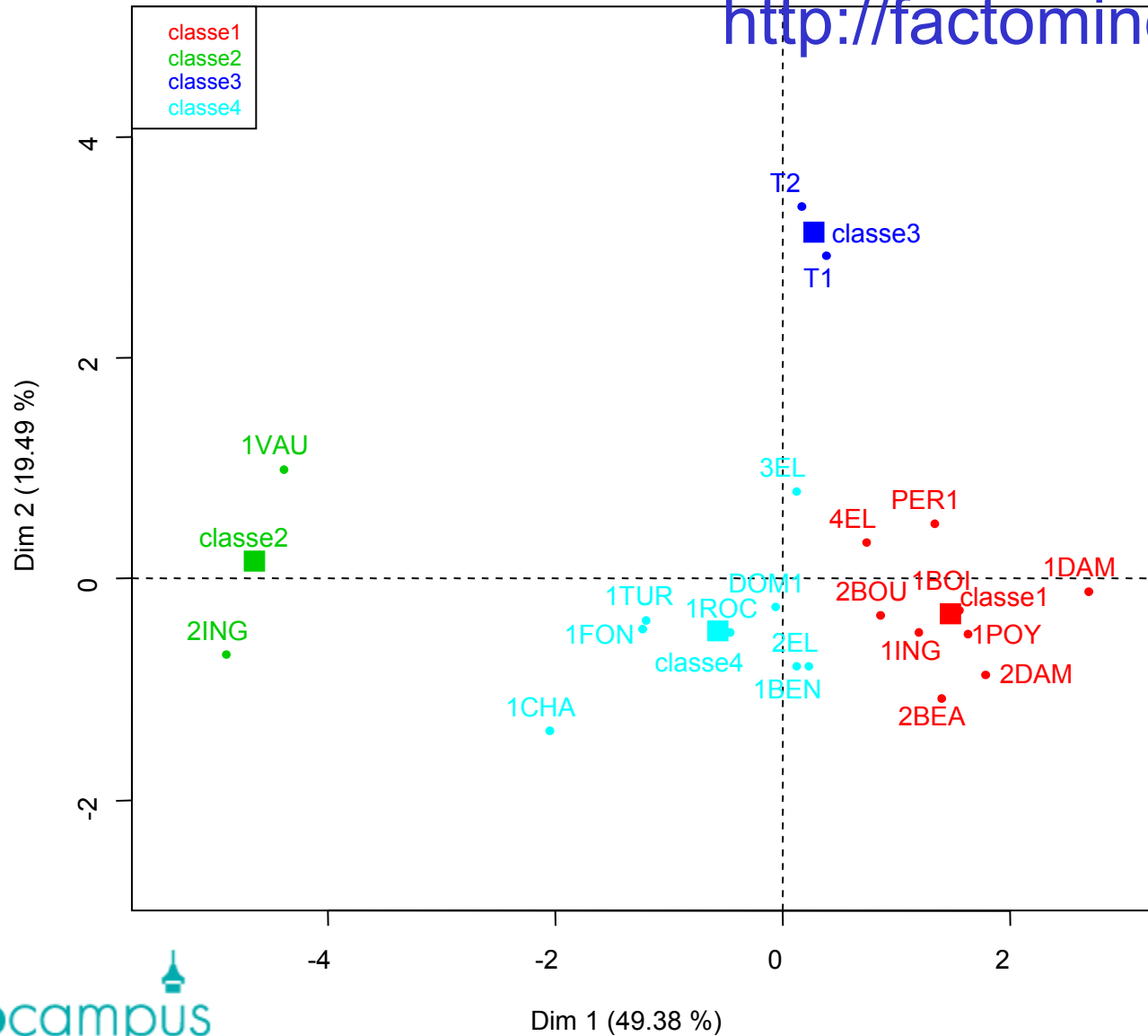
# Unsupervised classification

<http://factominer.free.fr>



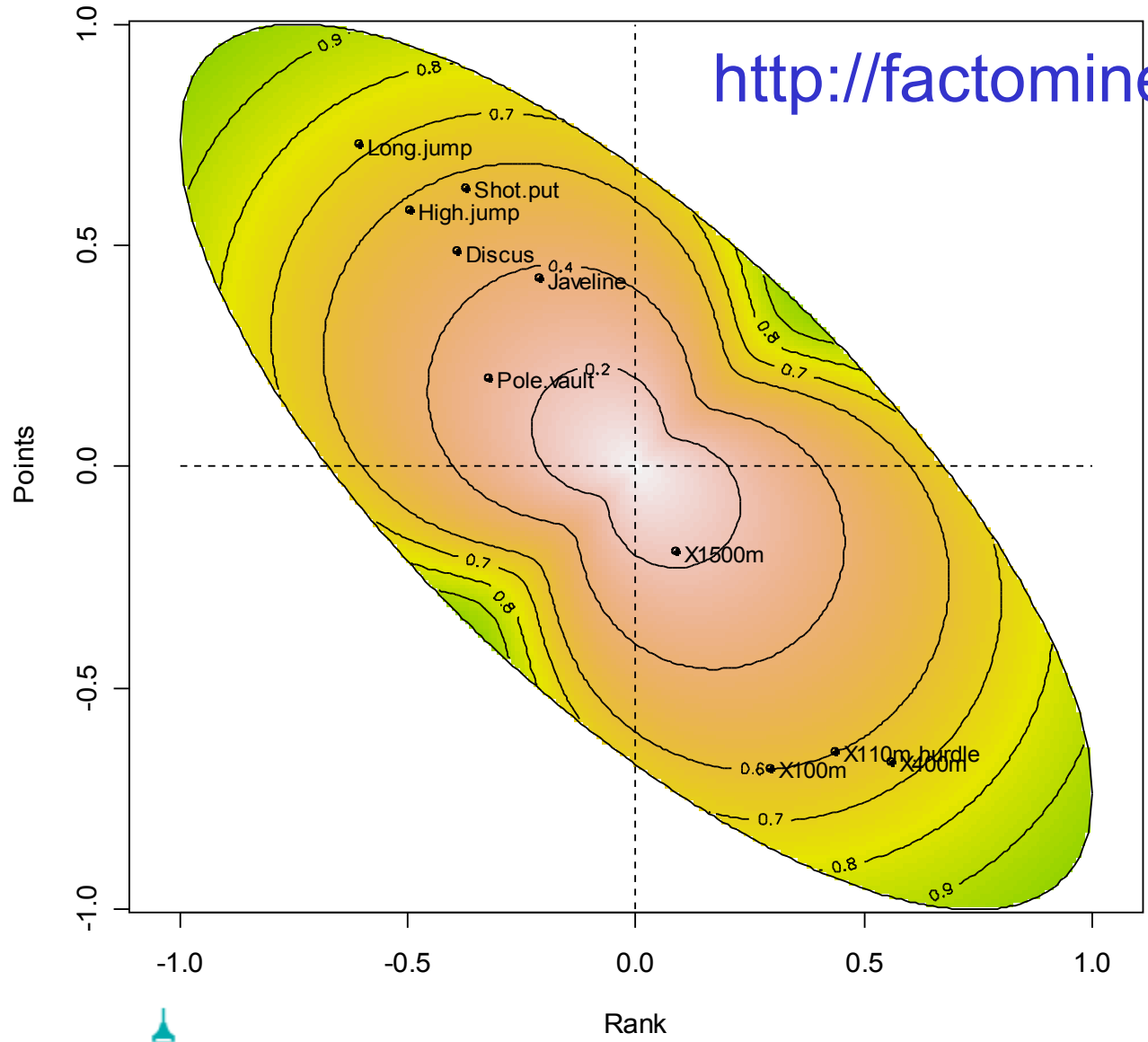
# MFA example: representation of the individuals

<http://factominer.free.fr>



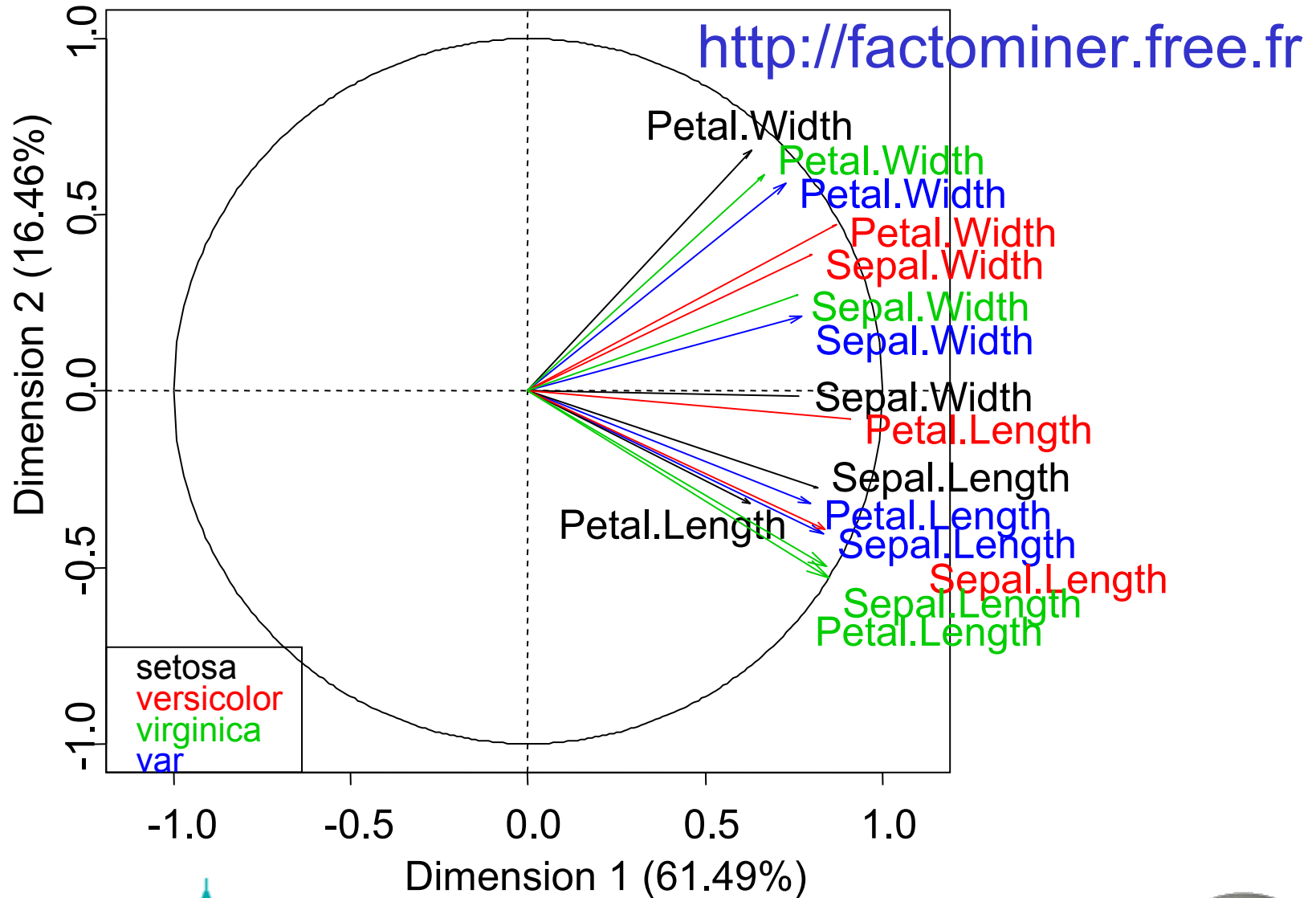
# Prefmap-PLS graph between Rank and Points

<http://factominer.free.fr>



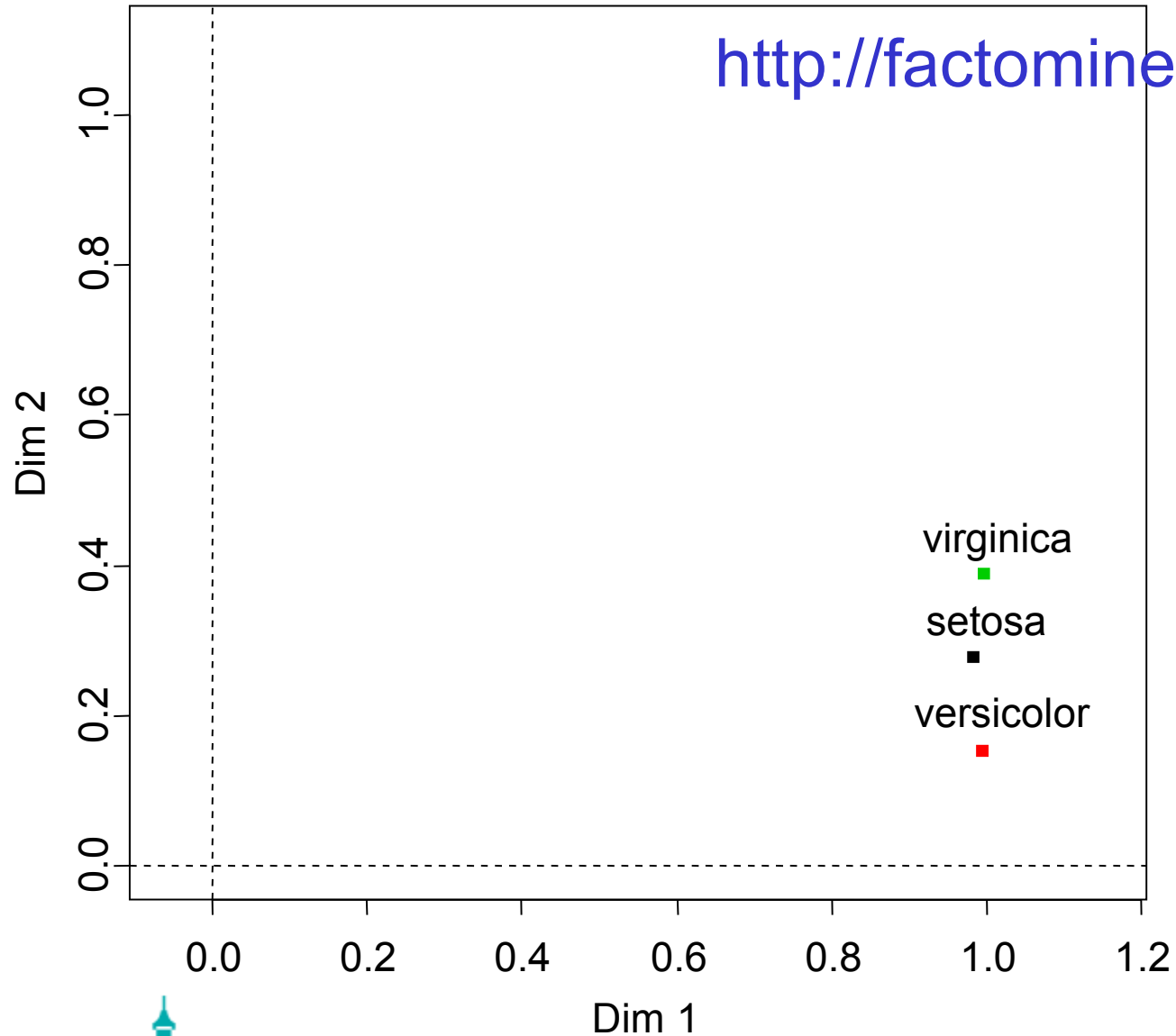
Correlation between Rank and Points : -0.7392

# MFA example: representation of the variables



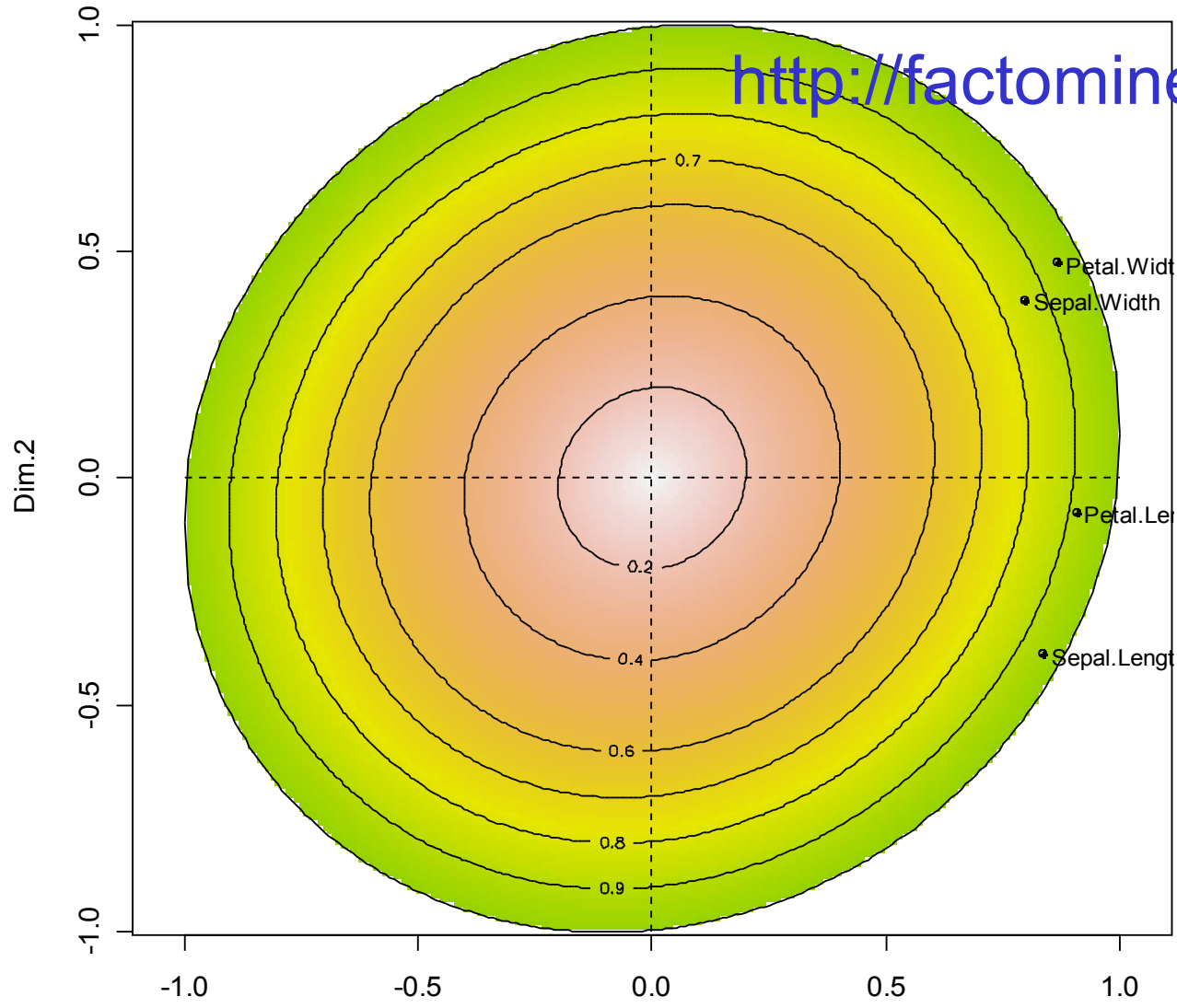
# Projection of the groups

<http://factominer.free.fr>



# Biplot between axes 1 and 2 for group versicolor

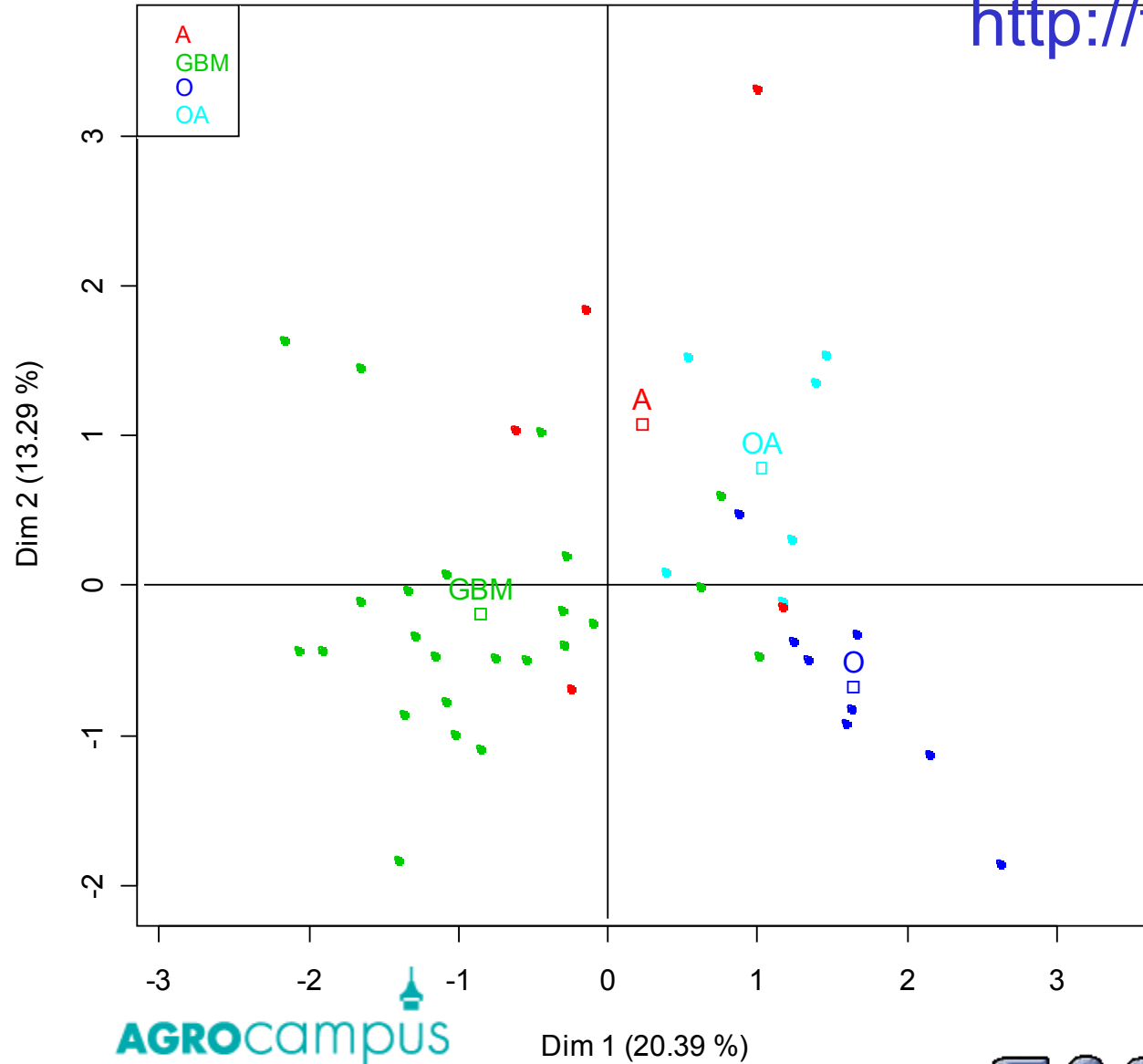
<http://factominer.free.fr>



Dim.1  
Correlation between Dim.1 and Dim.2 : 0.09613

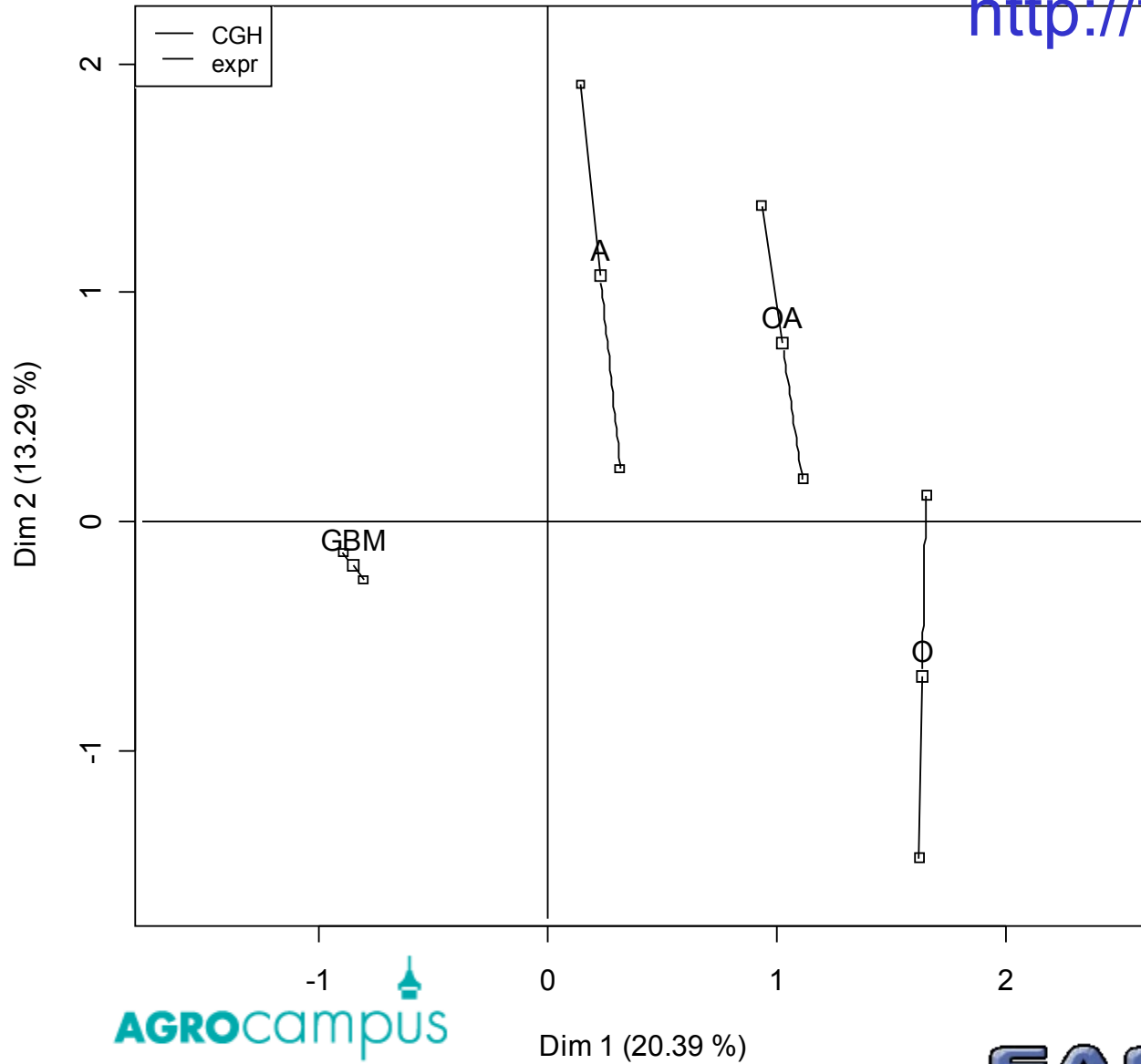
# Individual factor map

<http://factominer.free.fr>



# Individual factor map

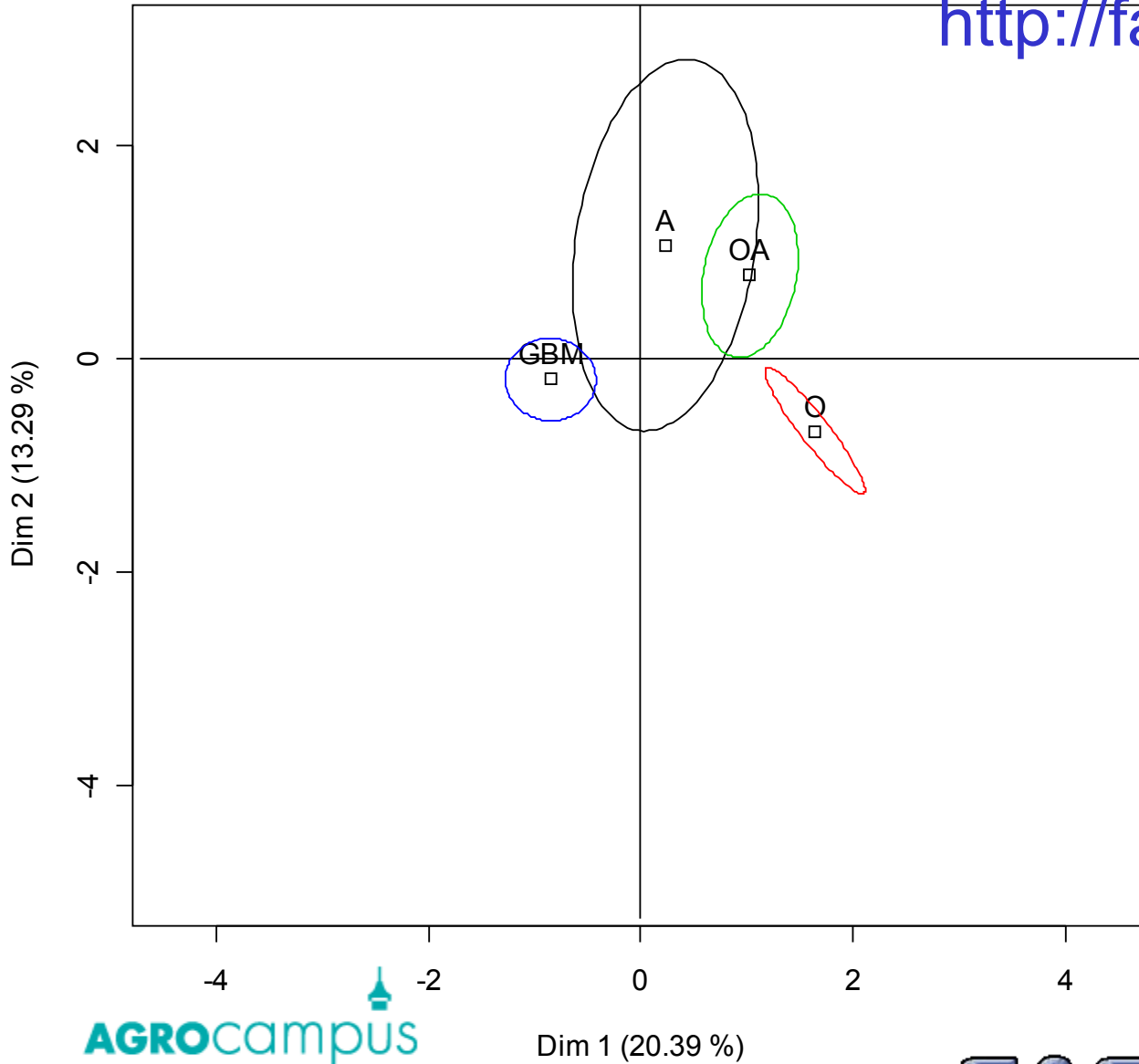
<http://factominer.free.fr>





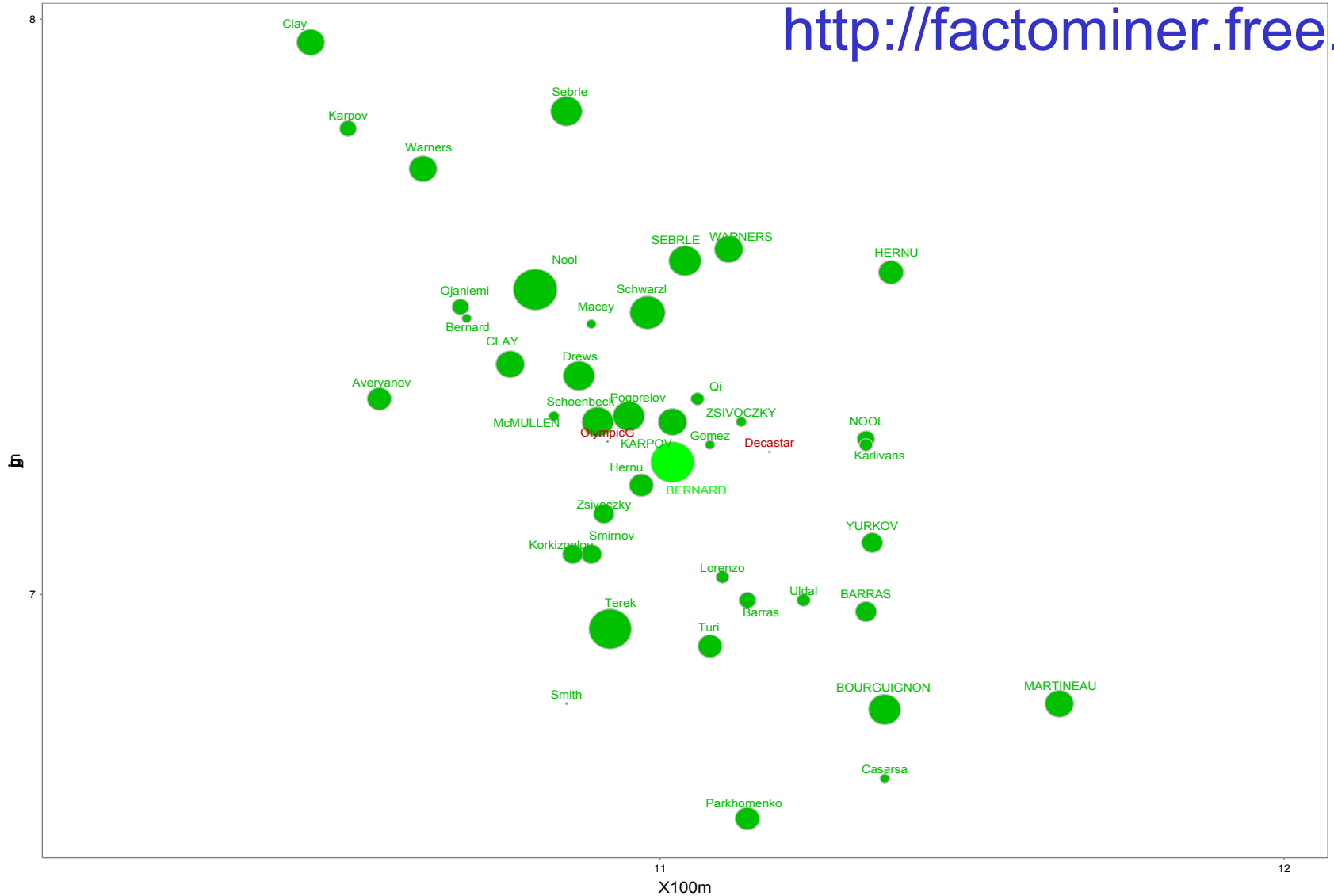
# Individual factor map

<http://factominer.free.fr>



\*Dataframe

<http://factominer.free.fr>



The FactoMineR team is nearly all the time ready to improve the package

