

Tutorial on Exploratory Data Analysis

Julie Josse, François Husson, Sébastien Lê

julie.josse at agrocampus-ouest.fr
francois.husson at agrocampus-ouest.fr

Applied Mathematics Department, Agrocampus Ouest

useR-2008
Dortmund, August 11th 2008

Why a tutorial on Exploratory Data Analysis?

- Our research focus on multi-table data analysis
- We teach Exploratory Data analysis from a long time (but in French ...)
- We made an R package:
 - Possibility to add supplementary information
 - The use of a more geometrical point of view allowing to draw graphs
 - 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)

Outline

- 1 Principal Component Analysis (PCA)
- 2 Correspondence Analysis (CA)
- 3 Multiple Correspondence Analysis (MCA)
- 4 Some extensions

Please ask questions!

Multivariate data analysis

- Principal Component Analysis (PCA) \Rightarrow continuous variables
- Correspondence Analysis (CA) \Rightarrow contingency table
- Multiple Correspondence Analysis (MCA) \Rightarrow categorical variables

- Dimensionality reduction \Rightarrow describe the dataset with smaller number of variables

- Techniques widely used for applications such as: data compression, data reconstruction; preprocessing before clustering, and ...

Exploratory data analysis

- ⇒ Descriptive methods
 - ⇒ Data visualization
 - ⇒ Geometrical approach: importance to graphical outputs
 - ⇒ Identification of clusters, detection of outliers
-
- ⇒ French school (Benzécri)

Principal Component Analysis

PCA in R

- `prcomp`, `princomp` from `stats`
- `dudi.pca` from `ade4` (<http://pbil.univ-lyon1.fr/ADE-4>)
- PCA from `FactoMineR` (<http://factominer.free.fr>)

PCA deals with which kind of data?

	1	k	K
1			
i		x_{ik}	
I			

Figure: Data table in PCA.

Notations:

x_i . the individual i ,

$S = \frac{1}{n} X_c' X_c$ the covariance matrix,

$W = X_c X_c'$ the inner products matrix.

- PCA deals with continuous variables, but categorical variables can also be included in the analysis

Some examples

- Many examples
 - Sensory analysis: products - descriptors
 - Environmental data: plants - measurements; waters - physico-chemical analyses
 - Economy: countries - economic indicators
 - Microbiology: cheeses - microbiological analyses
 - etc.
- Today we illustrate PCA with:
 - data decathlon: athletes performances during two athletics meetings
 - data chicken: genomics data

Decathlon data

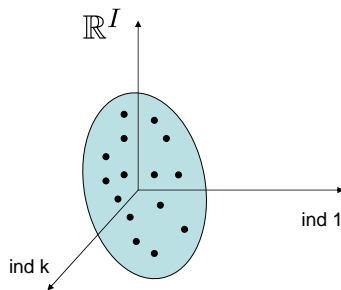
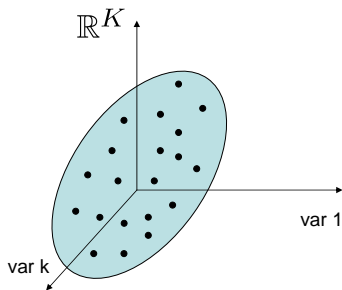
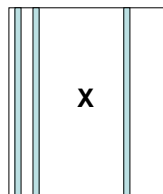
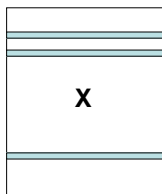
- 41 athletes (rows)
- 13 variables (columns):
 - 10 continuous variables corresponding to the performances
 - 2 continuous variables corresponding to the rank and the points obtained
 - 1 categorical variable corresponding to the athletics meeting: Olympic Game and Decastar (2004)

	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

Problems - objectives

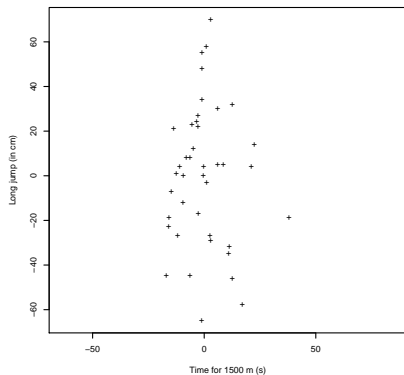
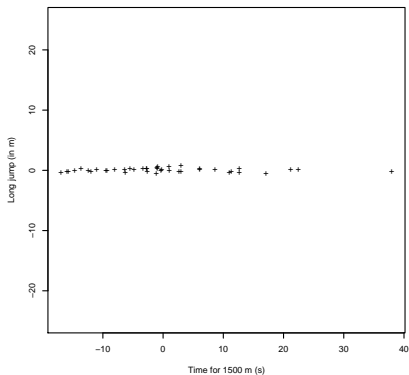
- Individuals study: similarity between individuals for all the variables (Euclidian distance) → partition between individuals
- Variables study: are there linear relationships between variables? Visualization of the correlation matrix; find some synthetic variables
- Link between the two studies: characterization of the groups of individuals by the variables; particular individuals to better understand the links between variables

Two points clouds



Pre-processing: Mean centering, Scaling?

- Mean centering does not modify the shape of the cloud
- Scaling: variables are always scaled when they are not in the same units



⇒ PCA always centered and often scaled

Individuals cloud



- Individuals are in \mathbb{R}^K
- Similarity between individuals: Euclidean distance
- Study the structure, i.e. the shape of the individual cloud

Inertia

- Total inertia \Rightarrow multidimensional variance \Rightarrow distance between the data and the barycenter:

$$\begin{aligned} I_g &= \sum_{i=1}^I p_i \|x_i - g\|^2 = \frac{1}{I} \sum_{i=1}^I (x_i - g)'(x_i - g), \\ &= \text{tr}(S) = \sum_s \lambda_s = K. \end{aligned}$$

Fit the individuals cloud

Find the subspace who better sum up the data: the closest one by projection.

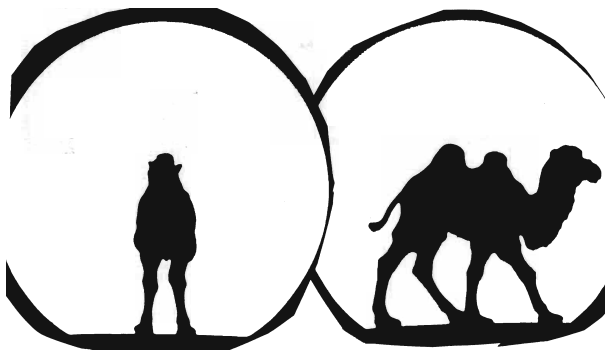
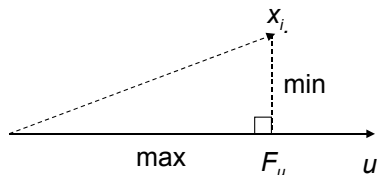


Figure: Camel vs dromedary?

Fit the individuals cloud



$$\begin{aligned}
 P_{u_1}(x_i) &= u_1(u_1' u_1)^{-1} u_1' x_i \\
 &= \langle x, u_1 \rangle u_1 \\
 F_{u_1}(i) &= \langle x, u_1 \rangle
 \end{aligned}$$

- Maximize the variance of the projected data
- Minimize the distance between individuals and their projections

⇒ Best representation of the diversity, variability of the individuals

⇒ Do not distort the distances between individuals

Find the subspace (1)

- Find the first axis u_1 , for which variance of $F_{u_1} = Xu_1$ is maximized:

$$u_1 = \operatorname{argmax}_{u_1} \operatorname{var}(Xu_1) \text{ with } u_1' u_1 = 1$$

$$\operatorname{var}(F_{u_1}) = \frac{1}{J} (Xu_1)' Xu_1 = u_1' \frac{1}{J} X' X u_1$$

It leads to:

$$\max u_1' S u_1 \text{ with } u_1' u_1 = 1$$

$\Rightarrow u_1$ first eigenvector of S (associated with the largest eigenvalue λ_1):

$$S u_1 = \lambda_1 u_1.$$

\Rightarrow This eigenvector is known as the first axis (loadings).

Find the subspace (2)

⇒ Projected inertia on the first axis:

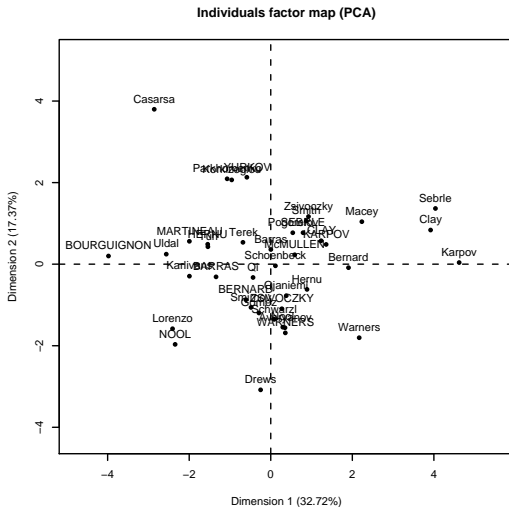
$$\text{var}(F_{u_1}) = \text{var}(Xu_1) = \lambda_1$$

⇒ Percentage of variance explained by the first axis:

$$\frac{\lambda_1}{\sum_k \lambda_k}$$

- Additional axes are defined in an incremental fashion: each new direction is chosen by maximizing the projected variance among all orthogonal directions
- Solution: K eigenvectors u_1, \dots, u_K of the data covariance matrix corresponding to the K largest eigenvalues $\lambda_1, \dots, \lambda_K$

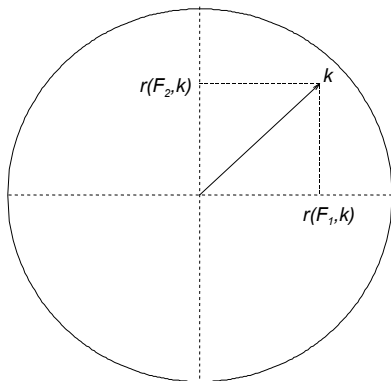
Example: graph of the individuals



⇒ Need variables to interpret the dimension of variability

Can we interpret the individuals' graph with the variables?

- Correlation between variable $x_{.k}$ and F_u (the vector of individuals coordinates in \mathbb{R}^I)



⇒ Correlation circle graph

Can we interpret the individuals' graph with the variables?

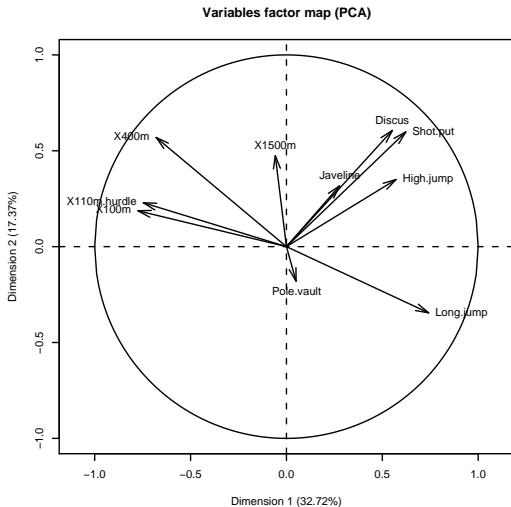
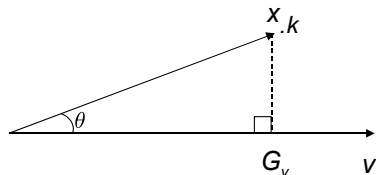


Figure: Graph of the variables.

Find the subspace (1)

- Variables are in \mathbb{R}^I
- Find the first dimension v_1 (in \mathbb{R}^I) which maximizes the variance of the projected data



$$P_{v_1}(x.k) = v_1(v_1'v_1)^{-1}v_1'x.k$$

$$G_{v_1}(k) = \frac{\langle v_1, x.k \rangle}{\|v_1\|}$$

$$\max \sum_{k=1}^K G_{v_1}(k)^2 = \max \sum_{k=1}^K \text{cor}(v_1, x.k)^2 = \max \sum_{k=1}^K \cos^2(\theta)$$

with $v'v = 1$

$\Rightarrow v_1$ is the best synthetic variable

Find the subspace (2)

Solution: v_1 is the first eigenvector of $W = XX'$ the individuals inner product matrix (associated with the largest eigenvalue λ_1):

$$Wv_1 = \lambda_1 v_1$$

⇒ The next dimensions are the other eigenvectors

⇒ Dimensionality reduction: principal components are linear combination of the variables

⇒ A subset of components to sum up the data

Fit the variables cloud

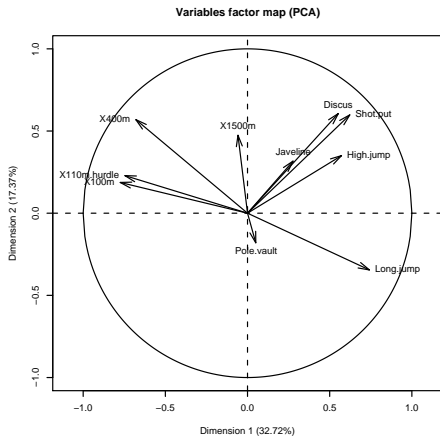
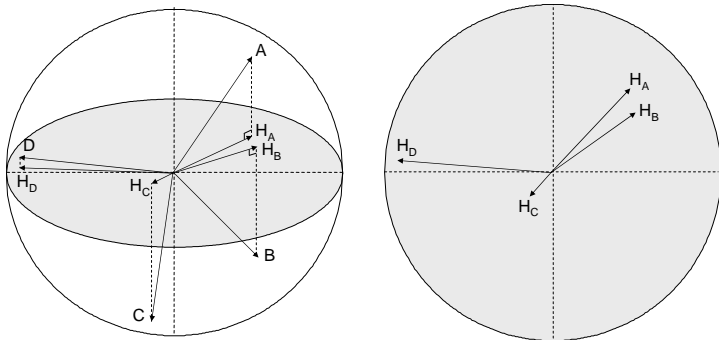


Figure: Graph of the variables.

⇒ Same representation! What a wonderful result!

Projections...

Only well projected variables (high \cos^2 between the variable and its projection) can be interpreted!



Exercices...

- With 200 independent variables and 7 individuals, how does your correlation circle look like?

Exercices...

- With 200 independent variables and 7 individuals, how does your correlation circle look like?

```
mat=matrix(rnorm(7*200,0,1),ncol=200)
PCA(mat)
```

Link between the two representations: transition formulae

- $Su = X'Xu = \lambda u$
- $XX'Xu = X\lambda u \rightarrow W(Xu) = \lambda(Xu)$
- $WF_u = \lambda F_u$ and since $Wv = \lambda v$ then F_u and v are colinear
- Since, $\|F_u\| = \lambda$ and $\|v\| = 1$ we have:

$$\begin{aligned} v &= \frac{1}{\sqrt{\lambda}} F_u &\Rightarrow G_v &= X'v = \frac{1}{\sqrt{\lambda}} X'F_u \\ u &= \frac{1}{\sqrt{\lambda}} G_v &\Rightarrow F_u &= Xu = \frac{1}{\sqrt{\lambda}} XG_v \end{aligned}$$

$$F_s(i) = \frac{1}{\sqrt{\lambda_s}} \sum_{k=1}^K x_{ik} G_s(k)$$

$$G_s(k) = \frac{1}{\sqrt{\lambda_s}} \sum_{i=1}^I x_{ik} F_s(i)$$

Example on decathlon

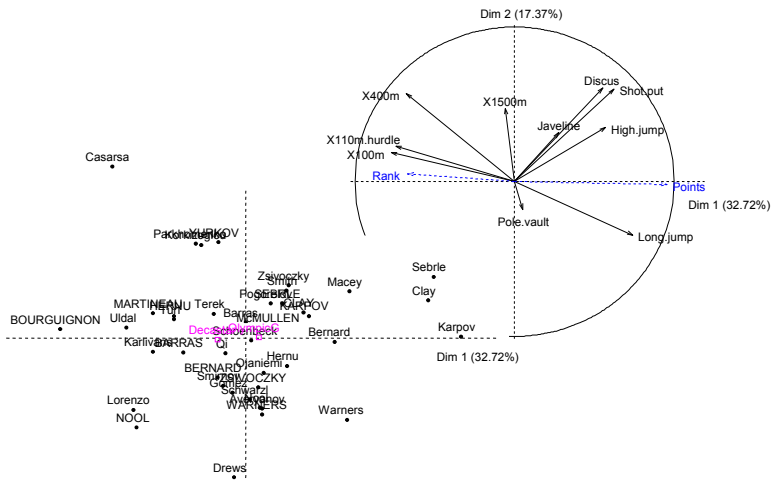
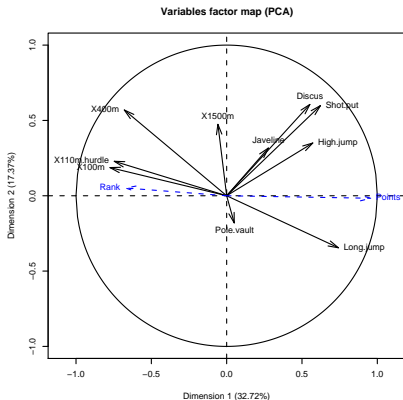


Figure: Individuals and variables representations.

Supplementary information (1)

- For the continuous variables: projection of these supplementary variables on the dimensions



- For the individuals: projection
- ⇒ Supplementary information do not participate to create the axes

Supplementary information (2)

→ How to deal with (supplementary) categorical variables?

	X100m	Long.jump	Shot.put	High.jump	Competition
HERNU	11.37	7.56	14.41	1.86	Decastar
BARRAS	11.33	6.97	14.09	1.95	Decastar
NOOL	11.33	7.27	12.68	1.98	Decastar
BOURGUIGNON	11.36	6.80	13.46	1.86	Decastar
Sebrle	10.85	7.84	16.36	2.12	OlympicG
Clay	10.44	7.96	15.23	2.06	OlympicG



	X100m	Long.jump	Shot.put	High.jump
HERNU	11.37	7.56	14.41	1.86
BARRAS	11.33	6.97	14.09	1.95
NOOL	11.33	7.27	12.68	1.98
BOURGUIGNON	11.36	6.80	13.46	1.86
Sebrle	10.85	7.84	16.36	2.12
Clay	10.44	7.96	15.23	2.06
Decastar	11.18	7.25	14.16	1.98
Olympic.G	10.92	7.27	14.62	1.98

→ The categories are projected at the barycenter of the individuals who take the categories

Number of dimensions?

- Percentage of variance explained by each axis: information brought by the dimension
- Quality of the approximation: $\frac{\sum_s^Q \lambda_s}{\sum_s^K \lambda_s}$
- Dimensionality Reduction implies Information Loss

- Number of components to retain? Retain much of the variability in our data (the other components are noise)
 - Bar plot of the eigenvalues: scree test
 - Test on eigenvalues, confidence interval,...

Percentage of variance obtained under independence

⇒ Is there a structure on my data?

```
nr=50
nc=8
iner=rep(0,1000)
for (i in 1:1000)
{
mat=matrix(rnorm(nr*nc,0,1),ncol=nc)
iner[i]=PCA(mat,graph=F)$eig[2,3]
}
quantile(iner,0.95)
```

Percentage of variance obtained under independence

⇒ Is there a structure on my data?

nbind	Number of variables												
	4	5	6	7	8	9	10	11	12	13	14	15	16
5	96.5	93.1	90.2	87.6	85.5	83.4	81.9	80.7	79.4	78.1	77.4	76.6	75.5
6	93.3	88.6	84.8	81.5	79.1	76.9	75.1	73.2	72.2	70.8	69.8	68.7	68.0
7	90.5	84.9	80.9	77.4	74.4	72.0	70.1	68.3	67.0	65.3	64.3	63.2	62.2
8	88.1	82.3	77.2	73.8	70.7	68.2	66.1	64.0	62.8	61.2	60.0	59.0	58.0
9	86.1	79.5	74.8	70.7	67.4	65.1	62.9	61.1	59.4	57.9	56.5	55.4	54.3
10	84.5	77.5	72.3	68.2	65.0	62.4	60.1	58.3	56.5	55.1	53.7	52.5	51.5
11	82.8	75.7	70.3	66.3	62.9	60.1	58.0	56.0	54.4	52.7	51.3	50.1	49.2
12	81.5	74.0	68.6	64.4	61.2	58.3	55.8	54.0	52.4	50.9	49.3	48.2	47.2
13	80.0	72.5	67.2	62.9	59.4	56.7	54.4	52.2	50.5	48.9	47.7	46.6	45.4
14	79.0	71.5	65.7	61.5	58.1	55.1	52.8	50.8	49.0	47.5	46.2	45.0	44.0
15	78.1	70.3	64.6	60.3	57.0	53.9	51.5	49.4	47.8	46.1	44.9	43.6	42.5
16	77.3	69.4	63.5	59.2	55.6	52.9	50.3	48.3	46.6	45.2	43.6	42.4	41.4
17	76.5	68.4	62.6	58.2	54.7	51.8	49.3	47.1	45.5	44.0	42.6	41.4	40.3
18	75.5	67.6	61.8	57.1	53.7	50.8	48.4	46.3	44.6	43.0	41.6	40.4	39.3
19	75.1	67.0	60.9	56.5	52.8	49.9	47.4	45.5	43.7	42.1	40.7	39.6	38.4
20	74.1	66.1	60.1	55.6	52.1	49.1	46.6	44.7	42.9	41.3	39.8	38.7	37.5
25	72.0	63.3	57.1	52.5	48.9	46.0	43.4	41.4	39.6	38.1	36.7	35.5	34.5
30	69.8	61.1	55.1	50.3	46.7	43.6	41.1	39.1	37.3	35.7	34.4	33.2	32.1
35	68.5	59.6	53.3	48.6	44.9	41.9	39.5	37.4	35.6	34.0	32.7	31.6	30.4
40	67.5	58.3	52.0	47.3	43.4	40.5	38.0	36.0	34.1	32.7	31.3	30.1	29.1
45	66.4	57.1	50.8	46.1	42.4	39.3	36.9	34.8	33.1	31.5	30.2	29.0	27.9
50	65.6	56.3	49.9	45.2	41.4	38.4	35.9	33.9	32.1	30.5	29.2	28.1	27.0
100	60.9	51.4	44.9	40.0	36.3	33.3	31.0	28.9	27.2	25.8	24.5	23.3	22.3

Table: 95 % quantile inertia on the two first dimensions of 10000 PCA on data with independent variables

Percentage of variance obtained under independence

nbind	Number of variables												
	17	18	19	20	25	30	35	40	50	75	100	150	200
5	74.9	74.2	73.5	72.8	70.7	68.8	67.4	66.4	64.7	62.0	60.5	58.5	57.4
6	67.0	66.3	65.6	64.9	62.3	60.4	58.9	57.6	55.8	52.9	51.0	49.0	47.8
7	61.3	60.7	59.7	59.1	56.4	54.3	52.6	51.4	49.5	46.4	44.6	42.4	41.2
8	57.0	56.2	55.4	54.5	51.8	49.7	47.8	46.7	44.6	41.6	39.8	37.6	36.4
9	53.6	52.5	51.8	51.2	48.1	45.9	44.4	42.9	41.0	38.0	36.1	34.0	32.7
10	50.6	49.8	49.0	48.3	45.2	42.9	41.4	40.1	38.0	35.0	33.2	31.0	29.8
11	48.1	47.2	46.5	45.8	42.8	40.6	39.0	37.7	35.6	32.6	30.8	28.7	27.5
12	46.2	45.2	44.4	43.8	40.7	38.5	36.9	35.5	33.5	30.5	28.8	26.7	25.5
13	44.4	43.4	42.8	41.9	39.0	36.8	35.1	33.9	31.8	28.8	27.1	25.0	23.9
14	42.9	42.0	41.3	40.4	37.4	35.2	33.6	32.3	30.4	27.4	25.7	23.6	22.4
15	41.6	40.7	39.8	39.1	36.2	34.0	32.4	31.1	29.0	26.0	24.3	22.4	21.2
16	40.4	39.5	38.7	37.9	35.0	32.8	31.1	29.8	27.9	24.9	23.2	21.2	20.1
17	39.4	38.5	37.6	36.9	33.8	31.7	30.1	28.8	26.8	23.9	22.2	20.3	19.2
18	38.3	37.4	36.7	35.8	32.9	30.7	29.1	27.8	25.9	22.9	21.3	19.4	18.3
19	37.4	36.5	35.8	34.9	32.0	29.9	28.3	27.0	25.1	22.2	20.5	18.6	17.5
20	36.7	35.8	34.9	34.2	31.3	29.1	27.5	26.2	24.3	21.4	19.8	18.0	16.9
25	33.5	32.5	31.8	31.1	28.1	26.0	24.5	23.3	21.4	18.6	17.0	15.2	14.2
30	31.2	30.3	29.5	28.8	26.0	23.9	22.3	21.1	19.3	16.6	15.1	13.4	12.5
35	29.5	28.6	27.9	27.1	24.3	22.2	20.7	19.6	17.8	15.2	13.7	12.1	11.1
40	28.1	27.3	26.5	25.8	23.0	21.0	19.5	18.4	16.6	14.1	12.7	11.1	10.2
45	27.0	26.1	25.4	24.7	21.9	20.0	18.5	17.4	15.7	13.2	11.8	10.3	9.4
50	26.1	25.3	24.6	23.8	21.1	19.1	17.7	16.6	14.9	12.5	11.1	9.6	8.7
100	21.5	20.7	19.9	19.3	16.7	14.9	13.6	12.5	11.0	8.9	7.7	6.4	5.7

Table: 95 % quantile inertia on the two first dimensions of 10000 PCA on data with independent variables

Quality of the representation: \cos^2

- For the variables: only well projected variables (high \cos^2 between the variable and its projection) can be interpreted!
- For the individuals: (same idea) distance between individuals can only be interpreted for well projected individuals

```
res.pca$ind$cos
```

	Dim.1	Dim.2
Sebrle	0.70	0.08
Clay	0.71	0.03
Karpov	0.85	0.00

Contribution

⇒ Contribution to the inertia to create the axis:

- For the individuals: $Ctr_s(i) = \frac{F_s^2(i)}{\sum_{i=1}^I F_s^2(i)} = \frac{F_s^2(i)}{\lambda_s}$

⇒ Individuals with large coordinate contribute the most to the construction of the axis

```
round(res.pca$ind$contrib,2)
```

	Dim.1	Dim.2
Sebrle	12.16	2.62
Clay	11.45	0.98
Karpov	15.91	0.00

- For the variables: $Ctr_s(k) = \frac{G_s^2(k)}{\lambda_s} = \frac{cor(x_k, v_s)^2}{\lambda_s}$

⇒ Variables highly correlated with the principal component contribute the most to the construction of the dimension

Description of the dimensions (1)

By the quantitative variables:

- The correlation between each variable and the coordinate of the individuals (principal components) on the axis s is calculated
- The correlation coefficients are sorted and significant ones are given

\$Dim.1

\$Dim.1\$quanti

Dim.1

Points 0.96

Long.jump 0.74

Shot.put 0.62

Rank -0.67

400m -0.68

110m.hurdle -0.75

100m -0.77

\$Dim.2

\$Dim.2\$quanti

Dim.2

Discus 0.61

Shot.put 0.60

Description of the dimensions (2)

By the categorical variables:

- Perform a one-way analysis of variance with the coordinates of the individuals on the axis explained by the categorical variable
- A F -test by variable
- For each category, a t -test to compare the average of the category with the general mean

```
$Dim.1$quali
```

	P-value
Competition	0.155

```
$Dim.1$category
```

	Estimate	P-value
OlympicG	0.4393	0.155
Decastar	-0.4393	0.155

Practice

```
library(FactoMineR)
data(decathlon)
res <- PCA(decathlon, quanti.sup=11:12, quali.sup=13)
plot(res, habillage=13)
res$eig
x11()
barplot(res$eig[,1], main="Eigenvalues", names.arg=1:nrow(res$eig))
res$ind$coord
res$ind$cos2
res$ind$contrib
dimdesc(res)
aa=cbind.data.frame(decathlon[,13], res$ind$coord)
bb=coord.ellipse(aa, bary=TRUE)
plot.PCA(res, habillage=13, ellipse=bb)
#write.infile(res, file="my_FactoMineR_results.csv") #to export a list
```

Application

Chicken data:

- 43 chickens (individuals)
- 7407 genes (variables)
- One categorical variable: 6 diets corresponding to different stresses
- Do genes differentially expressed from one stress to another?

⇒ Dimensionality reduction: with few principal components, we identify the structure in the data

