

Master in Data Science (MDS)

Multivariate Analysis (MVA) Multiple Correspondence Analysis

Prof. Arturo Palomino

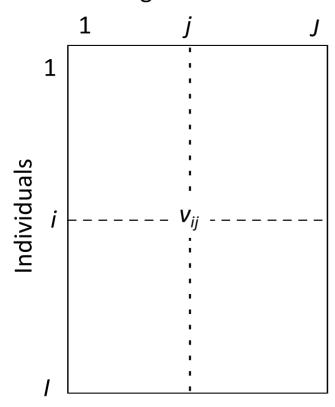
Slides modified from François Husson, with permission (Belchin Kostov belchin.adriyanov.kostov@upc.edu)



- 1 Data issues
- Studying the individuals
- Studying the categories
- 4 Interpretation aids



Categorical variables



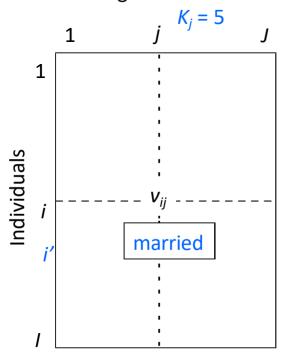
I individualsJ qualitative variables

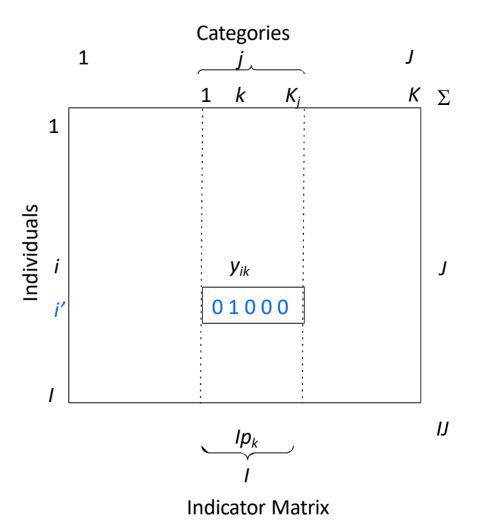
 V_{ij} : category of the j-th variable possessed by the i-th individual

Example: survey where I people reply to J multiple-choice questions



Categorical variables







Studying the individuals

One individual = one row of the CDT = set of categories Similarity of individuals – Inter-individual variability Principal axes of the inter-individual variability (in relation to the categories)

Studying the variables

Links between qualitative variables

(in relation to the categories)

Visualization of the set of associations between categories

Synthetic variables

(quantitative indicators based on the qualitative variables)

⇒ Similar problem to PCA



- Extract from 2003 INSEE survey on identity construction, called the "history of life" survey
- 8403 individuals
- 2 sorts of variables :
 - Which of the following leisure activities do you practice regularly: Reading, Listening to music, Cinema, Shows, Exhibitions, Computer, Sport, Walking, Travel, Playing a musical instrument, Collecting, Voluntary work, Home improvement, Gardening, Knitting, Cooking, Fishing, Number of hours of TV per day on average
 - Ssupplementary variables (4 questions) : sex, gender, profession, marital status



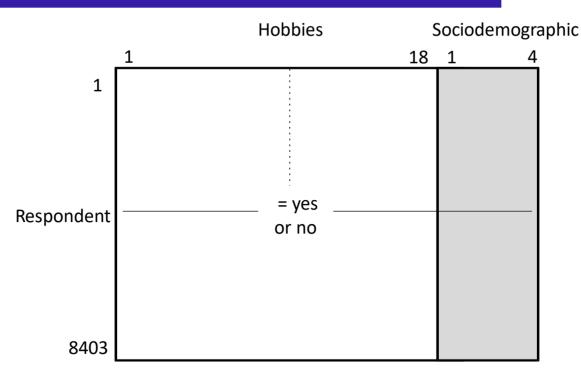
Hobbies

Hobbies		Number
Listening music		5947
Reading		5646
Walking		4175
Cooking		3686
Mechanic		3539
Travelling		3363
Cinema		3359
Gardening		3356
Computer		3158
Sport		3095
Exhibition		2595
Show		2425
Playing music		1460
Knitting		1413
Volunteering		1285
Fishing		945
Collecting		862
Number of hours watching TV	0	1017
	1	1223
	2	2156
	3	1775
	4	2232

Sociodemographic variables

Sex	Female	4616
	Male	3787
Age	[15,25]	857
	(25,35]	1302
	(35,45]	1646
	(45,55]	1837
	(55,65]	1257
	(65,75]	937
	(75,85]	482
	(85,100]	85
Marital	Divorcee	792
status	Married	4333
	Remarried	404
	Single	2140
	Widower	734
Profession	employee	2552
	foreman	735
	management	1052
	manual labourer	1161
	technician	401
	unskilled worker	792
	other	212
	No answer	1498





MCA 1: active = leisure activity, then use supplementary data for interpretation

- 1 individual = vector of leisure activities
- Principal axes of variability of leisure vectors
- Links between these axes and the supplementary variables

MCA 2 : active = supplementary variables, leisure activities as supplementary information

MCA 3 : active = BOTH

Transforming the complete disjunctive table



An individual's weight is $\frac{1}{I}$

- $y_{ik} = 1$ if the i-th individual is in k-th category of the j-th variable
 - (for each pk)
 - = 0 otherwise

Idea:
$$x_{ik} = y_{ik}/p_k$$

$$\frac{\sum_{i=1}^{I} x_{ik}}{I} = \frac{1}{I} \frac{\sum_{i=1}^{I} y_{ik}}{p_k} = \frac{1}{I} \frac{I \times p_k}{p_k} = 1$$

Centering :
$$x_{ik} = y_{ik}/p_k - 1$$

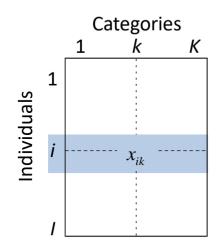


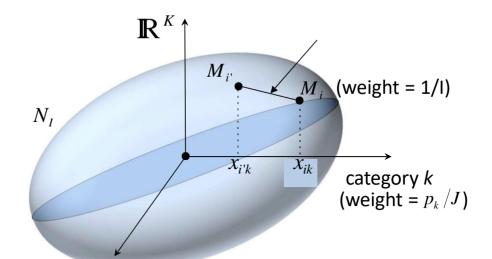
- Data issues
- Studying the individuals
- Studying the categories
- 4 Interpretation aids

Point cloud of individuals



Indicator matrix





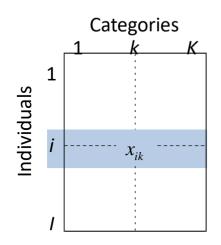
$$d_{i,i'}^2 = \sum_{k=1}^K \frac{p_k}{J} (x_{ik} - x_{i'k})^2 = \sum_{k=1}^K \frac{p_k}{J} \left(\frac{y_{ik}}{p_k} - \frac{y_{i'k}}{p_k} \right)^2 = \frac{1}{J} \sum_{k=1}^K \frac{1}{p_k} (y_{ik} - y_{i'k})^2$$

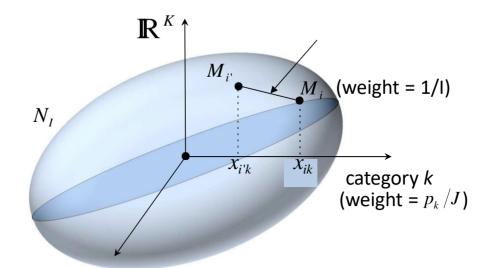
- 2 individuals with same categories : distance = 0
- 2 individuals with many shared categories: small distance
- 2 individuals, only 1 with a rare category: large distance to indicate this
- 2 individuals share rare category: small distance to indicate this shared specificity

Point cloud of individuals



Indicator matrix





$$d_{i,i'}^2 = \sum_{k=1}^K \frac{p_k}{J} (x_{ik} - x_{i'k})^2 = \sum_{k=1}^K \frac{p_k}{J} \left(\frac{y_{ik}}{p_k} - \frac{y_{i'k}}{p_k} \right)^2 = \frac{1}{J} \sum_{k=1}^K \frac{1}{p_k} (y_{ik} - y_{i'k})^2$$

$$d(i, G_I)^2 = \sum_{k=1}^K \frac{p_k}{J} (x_{ik})^2 = \sum_{k=1}^K \frac{p_k}{J} \left(\frac{y_{ik}}{p_k} - 1 \right)^2 = \frac{1}{J} \sum_{k=1}^K \frac{y_{ij}}{p_k} - 1$$

Inertia
$$(N_I) = \sum_{i=1}^{I} \underbrace{\frac{1}{I} d^2(i, O)}_{\text{inertia of } i} = \sum_{i=1}^{I} \left(\frac{1}{IJ} \sum_{k=1}^{K} \frac{y_{ik}}{p_k} - \frac{1}{I} \right) = \frac{K}{J} - 1$$

Building the point cloud of individuals



Getting factor axes, as usual, like for all factor analysis methods

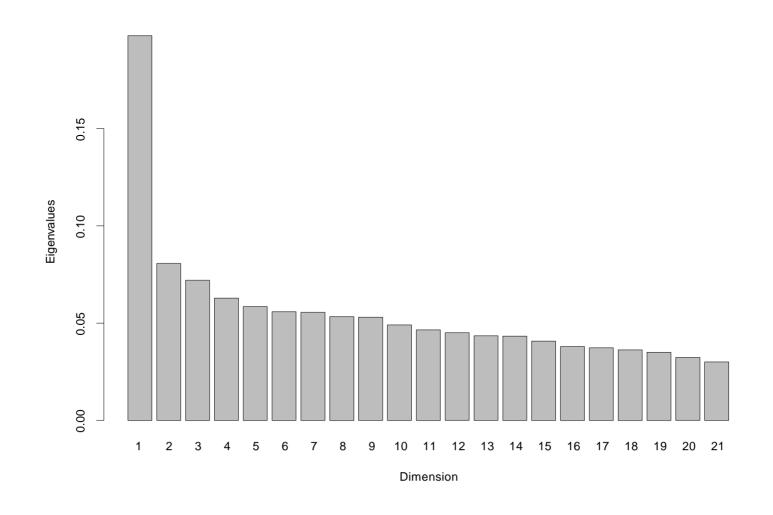
Sequential construction: look for the axis maximizing the inertia and orthogonal to previous axes



- Extract from 2003 INSEE survey on identity construction, called the "history of life" survey
- 8403 individuals
- 2 sorts of variables :
 - Which of the following leisure activities do you practice regularly: Reading, Listening to music, Cinema, Shows, Exhibitions, Computer, Sport, Walking, Travel, Playing a musical instrument, Collecting, Voluntary work, Home improvement, Gardening, Knitting, Cooking, Fishing, Number of hours of TV per day on average
 - Supplementary variables (4 questions) : sex, gender, profession, marital status

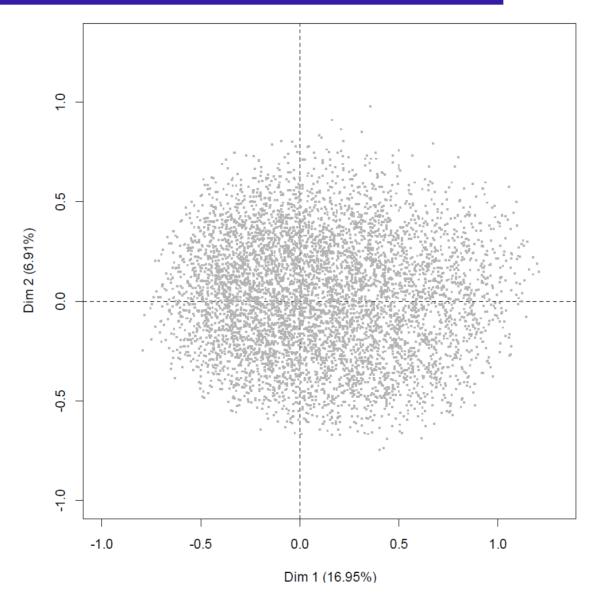
Diagram showing the inertia





Representation of the point cloud of individuals





Representation of the point cloud of individuals



What kind of pattern might we see?

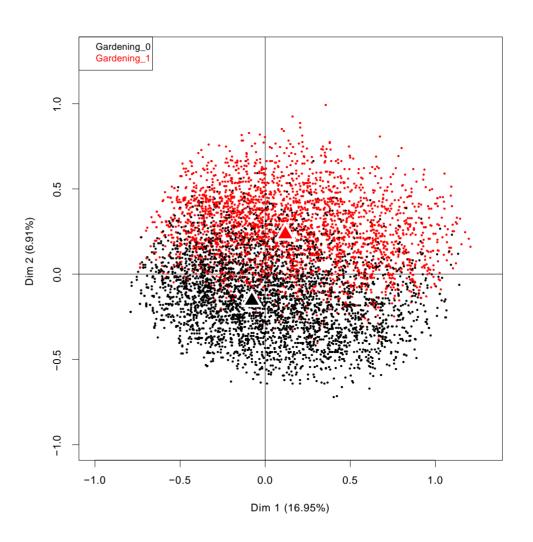


Individuals shown in terms of the gardening variable



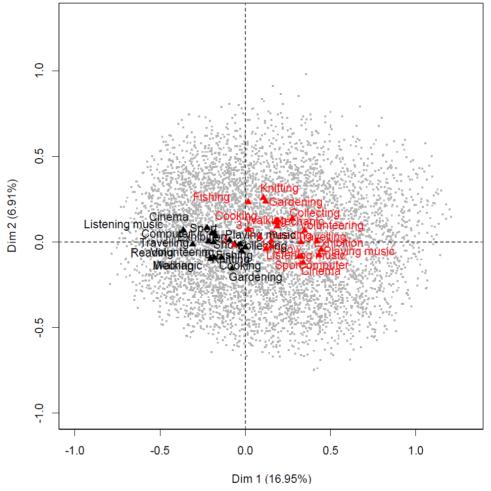
Idea: use the categories and variables to interpret the plot of the individuals

Put a category at the barycenter of the individuals in it



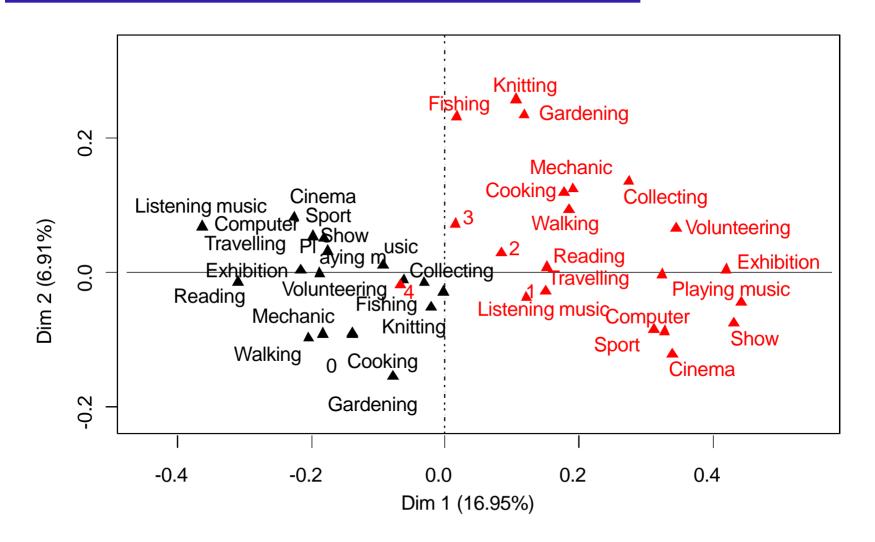


Each category is at the barycenter of the individuals in it



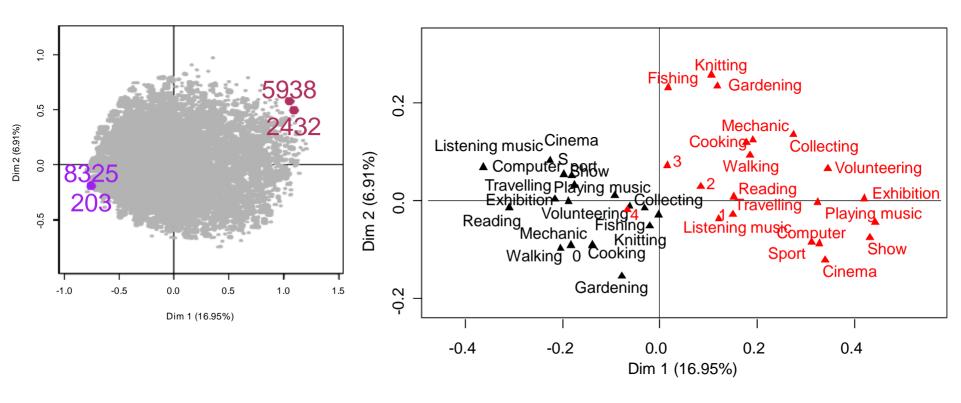
Activity not performed – activity performed





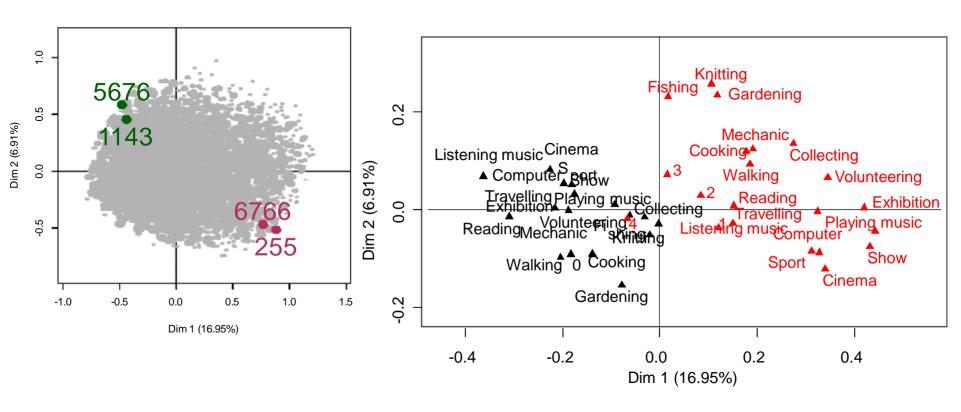
Activity not performed – activity performed





		Listen								Play								
	Read	music	Cinema	Show	Exhib	Comput	Sport	Walk	Travel	music	Collec	Volunteering	${\tt Mechanic}$	Garden	Knitt	Cook	Fish	TV
5938	У	У	n	Y	Y	У	У	Y	У	Y	У	У	У	У	У	Y	n	3
2432	Y	У	У	Y	У	У	n	Y	У	У	У	У	У	У	У	Y	n	2
8325	n	n	n	n	n	n	n	n	n	n	n	n	n	n	n	n	n	4
203	n	n	n	n	n	n	n	n	n	n	n	n	n	n	n	n	n	4





		Listen								Play								
	Read	music	Cinema	Show	Exhib	Comput	Sport	Walk	Travel	music	Collec	Volunteering	${\tt Mechanic}$	${\tt Garden}$	Knitt	Cook	Fish	TV
255	У	У	У	Y	Y	Y	Y	Y	У	Y	n	У	n	n	n	n	n	1
6766	У	У	У	Y	Y	У	Y	Y	У	У	n	n	n	n	n	У	n	0
5676	n	n	n	n	n	n	n	n	n	n	n	n	У	У	У	Y	n	4
1143	Y	n	n	n	n	n	n	n	n	n	n	n	У	У	Y	n	n	4

Showing the variables to help interpret the axes

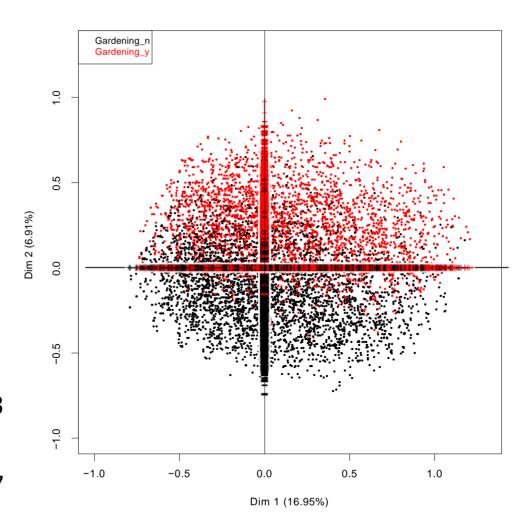


Idea: look at coordinates of projected individuals on each axis, and calculate a value for the connection between these coordinates and each qualitative variable

Correlation ratio between the j -th variable and s-th component : $\eta(v_{.j}, F_s)$

$$\eta^2(F_2, Gardening) = 0.453$$

$$\eta^2(F_1, Gardening) = 0.047$$



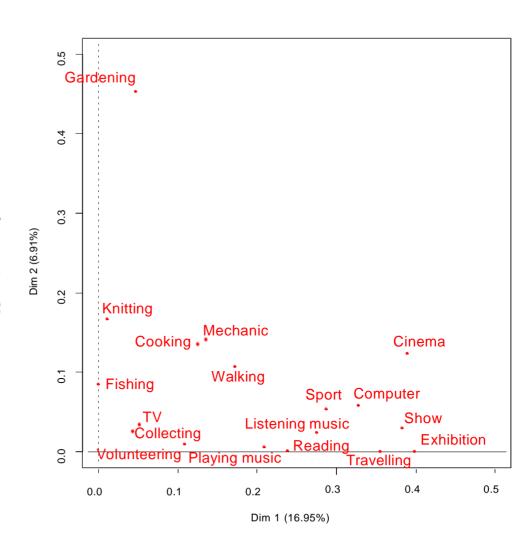
Showing the variables to help interpret the axes



Using the squared correlation ratios

The s-th axis is orthogonal to the t-th for all t < s, and the most related to the qualitative variables in the η^2 sense :

$$F_s = \max_F \sum_{j=1}^J \eta^2(F, v_{.j})$$

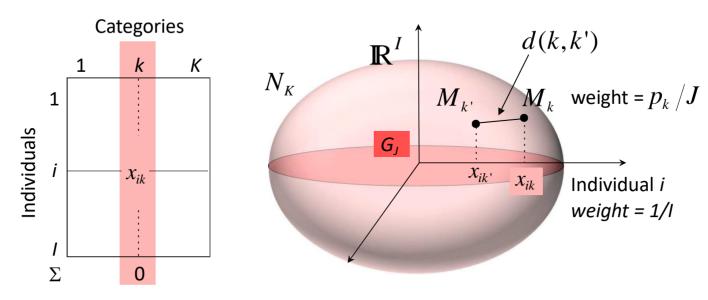




- Data issues
- Studying the individuals
- 3 Studying the categories
- 4 Interpretation aids

Point cloud of categories





Inertia of categories or variables



$$Inertia(k) = rac{1-p_k}{J}$$
 $Inertia(j) = rac{1}{J}\sum_{k=1}^{K_j}(1-p_k) = rac{K_j-1}{J}$

Variable	No. of categories	Inertia	No. dim. of subspace
sex	2	1/ <i>J</i>	1
region	21	20/ <i>J</i>	20
district	96	95/ J	95

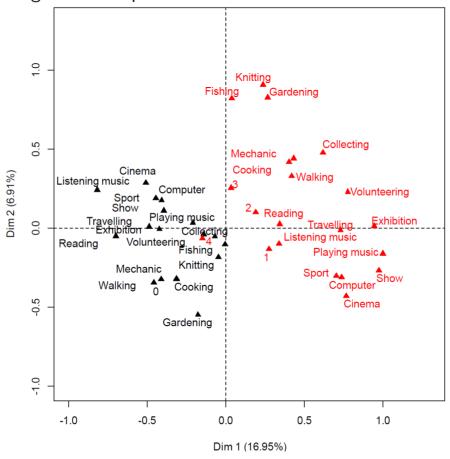
BUT : the inertia $\frac{K_j-1}{J}$ is spread across a K_j-1 dim. subspace

Total inertia
$$=\sum_{j=1}^{J} \frac{K_j-1}{J} = \frac{K}{J}-1$$

Representing the point cloud of categories



Sequential search for axes – as usual in factor analysis : each axis must maximize the inertia and be orthogonal to all previous ones

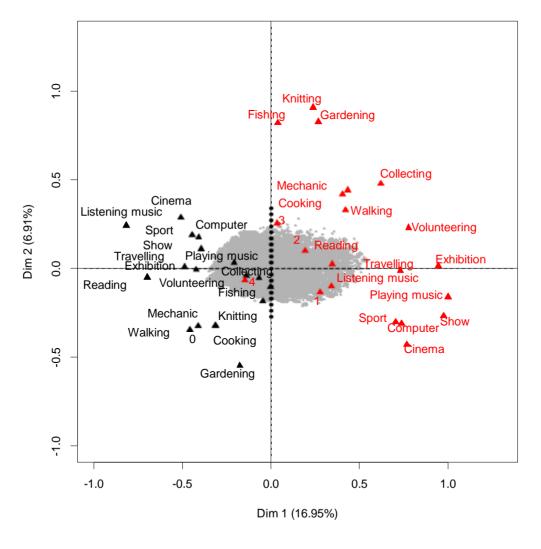


Activity not performed - activity performed

Projections of the individuals



Each individual put at barycenter of the categories they possess



Barycentric representation – simultaneous

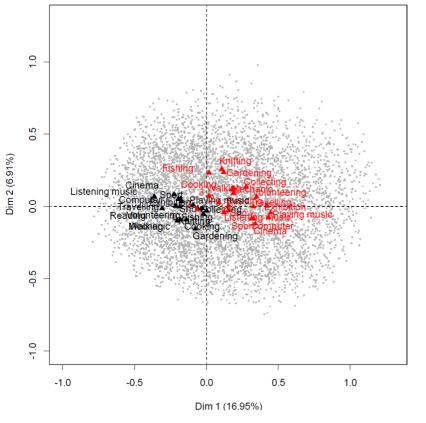


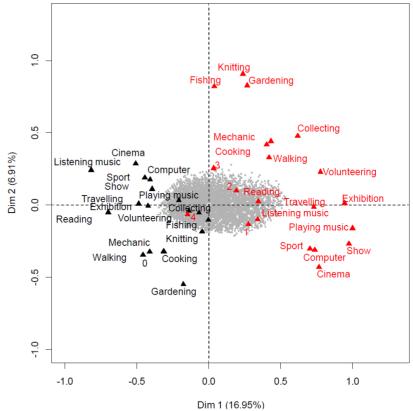
Optimal representation of individuals Categories at the barycenter:

$$G_s(k) = \sum_{i=1}^{l} \frac{y_{ik}}{I_k} F_s(i)$$

Optimal representation of categories Individuals at the barycenter:

$$F_s(i) = \sum_{i=1}^J \frac{y_{ik}}{J} G_s(k)$$





Barycentric representation – simultaneous

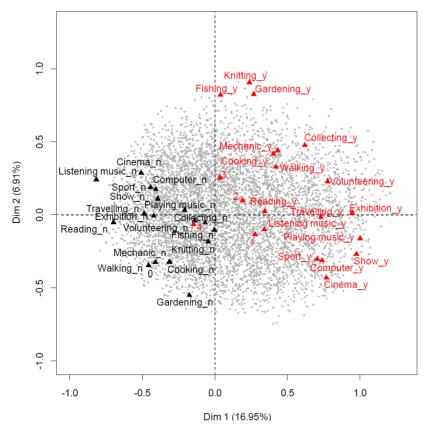


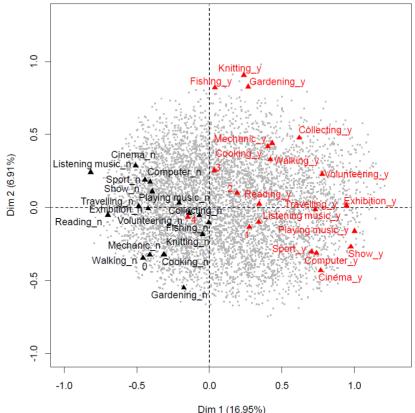
Optimal representation of individuals Categories at the pseudo-barycenter:

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

Optimal representation of categories Individuals at the pseudo-barycenter:

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







- Data issues
- Studying the individuals
- Studying the categories
- 4 Interpretation aids

Inertia and percentage of inertia in MCA



$$\lambda_s = \frac{1}{J} \sum_{j=1}^J \eta^2(F_s, v_{.j})$$

- $\Rightarrow \lambda_s$ is the mean of the squared correlation ratios
 - Individuals live in $\mathbb{R}^{K-J} \Rightarrow$ low percentages of inertia
 - Maximal percentage for given axis s :

$$\frac{\lambda_s}{\sum_{t=1}^{K-J} \lambda_t} \times 100 \leq \frac{1}{\frac{K-J}{J}} \times 100$$

$$\leq \frac{J}{K-J} \times 100$$

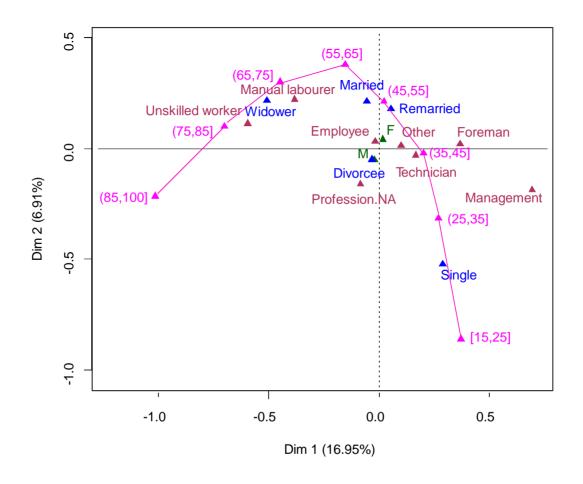
With K = 100, $J = 10 : \lambda_s \le 11.1 \%$

• Mean of non-zero eigenvalues : $\frac{1}{K-J} \times \sum_t \lambda_t = \frac{1}{K-J} \times \left(\frac{K}{J} - 1\right) = \frac{1}{J}$ \Rightarrow interpret the axes of inertia above 1/J

Representing supplementary elements



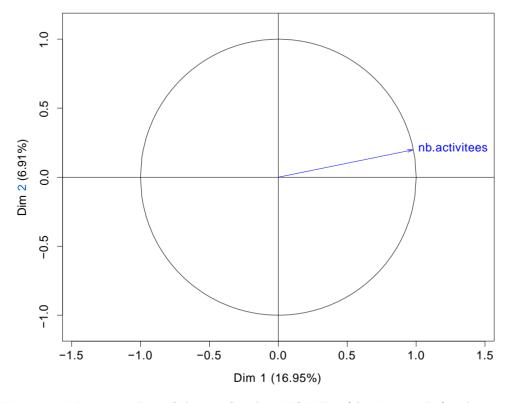
Use transition formulas to represent supplementary elements (individuals, variables, etc.)



Quantitative supplementary variables



- ⇒ What can we do with quantitative variables?
 - Supplementary information : project onto the axes, calculate correlation coefficients with each axis
 - break up quantitative variable into categories/classes



Describing the axes



Using qualitative variables (Fisher test), using categories (Student test), using quantitative variables (correlations)

Quantitative variables correlation p.value nb.activitees 0.9753459 0

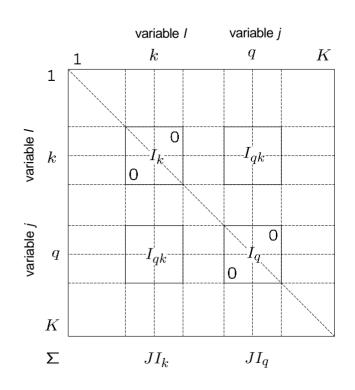
Cate	egorica	al variable	es	Cate	egories
	R2	p.value		Estimate	p.value
Reading	0.239	0.00e+00	Playing music_Y	0.268	0
Listening music	0.275	0.00e+00	Travelling_Y	0.270	0
Cinema	0.389	0.00e+00	${\tt Walking_Y}$	0.184	0
Show	0.383	0.00e+00	Sport_Y	0.247	0
Exhibition	0.399	0.00e+00	Computer_Y	0.263	0
Computer	0.327	0.00e+00	Exhibition_Y	0.304	0
Sport	0.287	0.00e+00	Show_Y	0.304	0
Walking	0.172	0.00e+00	Sport_N	-0.247	0
Travelling	0.355	0.00e+00	Computer_N	-0.263	0
Playing music	0.209	0.00e+00	Exhibition_N	-0.304	0
Mechanic	0.135	8.82e-267	${\tt Show_N}$	-0.304	0
Cooking	0.125	9.42e-247	Cinema_N	-0.283	0
Profession	0.128	7.20e-245	Listening music_N	-0.257	0
Volunteering	0.109	2.25e-212	Reading_N	-0.231	0

Different MCA strategy: Burt table



Burt table:

- Pairwise links between variables (like a correlation matrix between quantitative variables)
- Correspondence analysis on Burt table
- Gives results uniquely for categories : same representation but different Eigenvalues: $\lambda_s^{Burt} = (\lambda_s^{TDC})^2$
- λ_s^{TDC} mean of squared correlation ratios



⇒ The MCA only depends on pairwise links between variables (just like PCA only depends on the correlation matrix)

Conclusion



- MCA is the best factor analysis method for tables of individuals with qualitative variables
- Eigenvalues represent the means of squared correlation ratios
- The values of these squared links are particularly important when there are lots of variables
- Return to the data by analyzing the contingency table with CA
- Convergence of CDT analysis and Burt table analysis is a strong argument in favor of the general method
- MCA can be use to pre-treat data before doing classification



For more information about MCA, have a look at these videos:

https://www.youtube.com/watch?v=gZ_7WWEVITg

https://www.youtube.com/watch?v=b4kRAt4mkB8

https://www.youtube.com/watch?v=OZW1KrKO5Ac

https://www.youtube.com/watch?v=wTCuThGiy4w

https://www.youtube.com/watch?v=SI1-UYD6iac

Husson F., Lê S. & Pagès J. (2017) Exploratory Multivariate Analysis by Example Using R 2nd edition, 230 p., CRC/Press.