

# III — Introduction to Multiple Correspondence Analysis (MCA)

Adapted from “What is MCA?”, paper presented by Brigitte Le Roux at  
*Research Methods Festival*, Oxford (July 2006).

Brigitte.LeRoux@math-info.univ-paris5.fr      [www.math-info.univ-paris5.fr/~lerb/](http://www.math-info.univ-paris5.fr/~lerb/)  
rouanet@math-info.univ-paris5.fr      [www.math-info.univ-paris5.fr/~rouanet/](http://www.math-info.univ-paris5.fr/~rouanet/)

# 1 Introduction

Language of questionnaire

Basic data set: Individuals×Questions table

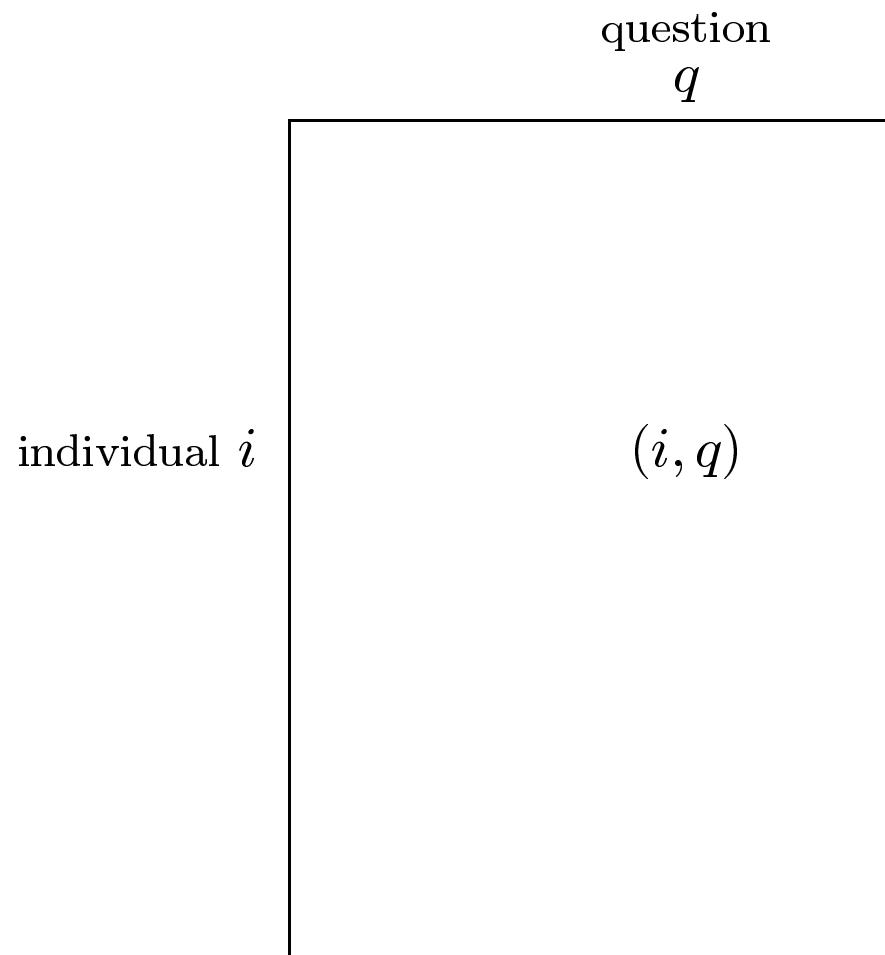
Questions are categorized variables, that is, variables with a finite number of response categories, or *modalities*.

Questionnaire in “standard format”: for each question, each individual chooses *one and only one* response modality.

$I$ : set of  $n$  individuals

$Q$ : set of questions

Basic data table analyzed by MCA:



## 1.1 Historical landmarks

Guttman (1941)

Burt (1950)

Benzécri (1972-1977)

Lebart (1975)

Bourdieu & Saint-Martin 1978 (*Le patronat*).

## 2 Principles of MCA

$$\text{MCA} \longrightarrow \begin{cases} \text{cloud of individuals} \\ \text{cloud of modalities} \end{cases}$$

*Distance between two individuals  $i$  and  $i'$  for question  $q$*   
 *$i$  chooses modality  $k$ ;  $i'$  chooses modality  $k' \neq k$ :*

$$d_q^2(i, i') = \frac{1}{n_k/n} + \frac{1}{n_{k'}/n}$$

$$\text{Overall distance } d^2(i, i') = \frac{1}{Q} \sum_{q \in Q} d_q^2(i, i')$$

*Distance between two modalities  $k$  and  $k'$*

$$d^2(k, k') = \frac{n_k + n_{k'} - 2n_{kk'}}{n_k n_{k'}/n}$$

$n_k$  = number of individuals who have chosen  $k$  n (resp.  $n_{k'}$ );

$n_{kk'}$  = number of individuals who have chosen both  $k$  and  $k'$ .

# Principal axes, eigenvalues and contributions

*Fundamental properties:*

- the principal axes of the cloud of individuals are in a one-one correspondence with those of the cloud of modalities,
- the two clouds have the same eigenvalues.

## Aids to interpretation: Contributions

*Overall contribution of question  $q$ :*  $\frac{K_q - 1}{K - Q}$

( $K_q$ : number of modalities of question  $q$ ,  $K$  overall number of modalities)

*Contribution of point to axis:*  $\frac{py^2}{\lambda}$

( $y$ : coordinate of point on axis;  $p$ : relative weight;  $\lambda$ : eigenvalue)

Contributions add up by grouping → contribution of a question.

### 3 Steps of analysis

1. Choose active individuals, active questions and encode modalities;
2. Decide how many axes to be interpreted;
3. Visualize the two clouds;
4. Interpret the axes retained in the cloud of modalities;
5. Investigate the cloud of individuals (landmark patterns, structuring factors, ellipses).

## The Culture Example

Data from a 1997 survey on the cultural practices of French people conducted by O. Donnat<sup>a</sup> & Col at the Department of Studies and Prospective of the Ministry of Culture and Communication.

**Sample:** 3002 individuals aged 15 or more, representative of the French population; 125 questions

**Our dataset<sup>b</sup>:**  $Q = 6$  questions pertaining to *leisure activities*, and 3 *identification variables*; a set  $I$  of 2720 individuals aged 18 or more who answered fully the six questions.

**Research Questions:** Are there differences about leisure practice between genders, age classes and educational levels ?

---

<sup>a</sup>See O. Donnat (1998), *Les pratiques culturelles des français: enquête 1997*, Paris, La documentation française.

<sup>b</sup>See B. Le Roux & H. Rouanet (2004), p.221-241.

## Data Set

(q1). As a general rule, do you prefer *Leisure activities* that you can do

<i>q1</i>	<i>Leisure activity</i>	freq.	%
<i>q1r1</i> alone	434	16.0	
<i>q1r2</i> with <i>partner</i>	507	18.6	
<i>q1r3</i> with <i>friends</i>	1060	39.0	
<i>q1r4</i> with <i>family</i>	719	26.4	
Total	2720	100.0	

(q2). As a general rule, would you say that during your *Free time*

<i>q2</i>	<i>Free time</i>	freq.	%
<i>q2r1</i> you <i>lack time</i> to do all that you would like to do	1121	41.2	
<i>q2r2</i> you don't lack time but you have <i>always something to do</i>	1157	42.5	
<i>q2r3</i> sometimes you have <i>nothing particular to do</i>	241	8.9	
<i>q2r4</i> often you <i>do nothing in particular</i>	201	7.4	
Total	2720	100.0	

(q3). If you had *More time*, your first choice of activity would be

<i>q3 More time</i>	freq.	%
<i>q3r1</i> to <i>rest</i> , not to do anything in particular	304	11.2
<i>q3r2</i> to <i>take courses</i> to improve your work situation	262	9.6
<i>q3r3</i> to discover or practice more <i>physical activities</i>	573	21.1
<i>q3r4</i> to discover or practice more <i>artistic activities</i>		
<i>q3r5</i> to develop your <i>general knowledge</i>	449	16.5
<i>q3r6</i> to <i>take care</i> of your <i>family</i>	316	11.6
<i>q3r7</i> to do some <i>home DIY</i> (gardening, etc.)	422	15.5
Total	2720	100.0

(q4). When you *Go out* in the evening, do you usually go

<i>q4 Going out</i>	freq.	%
<i>q4r1</i> <i>alone</i>	202	7.4
<i>q4r2</i> with your <i>partner</i>	911	33.5
<i>q4r3</i> with <i>family</i> , children, parents, etc.	546	20.1
<i>q4r4a</i> with <i>friends</i>	538	19.8
<i>q4r4b</i> with a <i>group</i> (workers' council, club, etc.)	54	1.4
<i>q4r5</i> you <i>don't go out</i> in the evening.	469	17.2
Total	2720	100.0

We have constructed two further questions  $q5$  and  $q6$ :

Question  $q5$  was built from questions about the time of watching *TV* and grouped in 5 categories.

Question  $q6$  was built from questions about the number of *Books* and comic strips read during the last 12 months and grouped in 5 categories.

$q5$	<i>TV</i>	freq.	%
$q5r1$	never	257	9.4
$q5r2$	less than 10h	435	16.0
$q5r3$	10-19h	794	29.2
$q5r4$	19-30h	705	25.9
$q5r5$	over 30h	529	19.4
Total		2720	100.0

$q6$	<i>Books</i>	freq.	%
$q6r1$	no books	603	22.2
$q6r2$	1-4 books	482	17.7
$q6r3$	5-12 books	641	23.6
$q6r4$	13-39 books	563	20.7
$q6r5$	40 books or more	431	15.8
Total		2720	100.0

## Choosing active questions and encoding modalities

- *Active questions:* 6 questions about leisure
- *Encoding modalities:* *TV* and *Books* recoded from the distribution of hours by week and # of books read the last 12 months.
- *Rare modalities, non-responses, “junk” modalities*

Rare modalities (say, of frequencies less than 5%) need to be pooled with others whenever feasible, or alternatively be put as “passive” ones (Specific MCA<sup>a</sup>).

For question *Go out*: we pool the rare modality “group” (1.4%) with the modality “friends”.

- *Technique of Supplementary Elements:*

*Supplementary variables; Supplementary individuals*

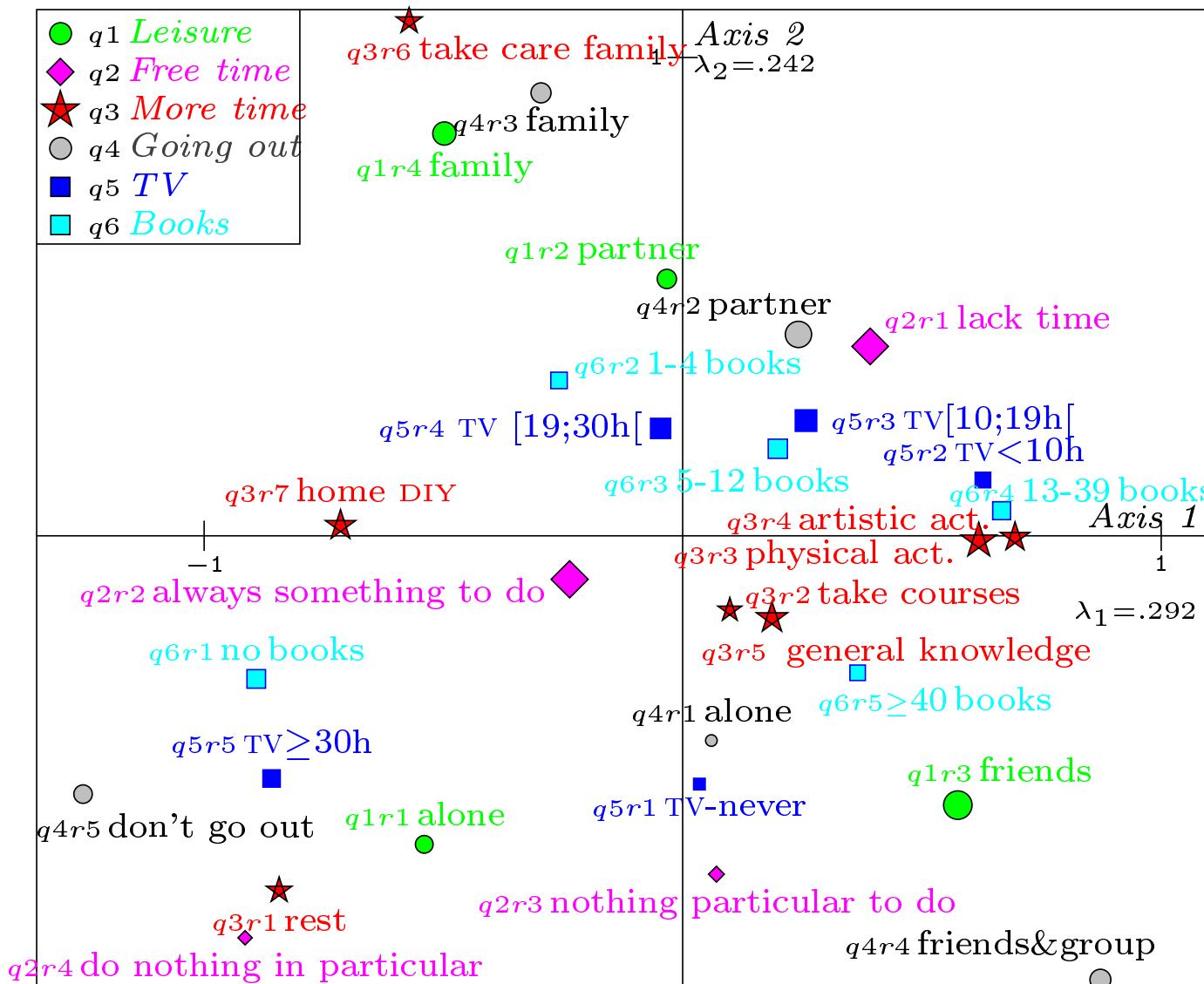
---

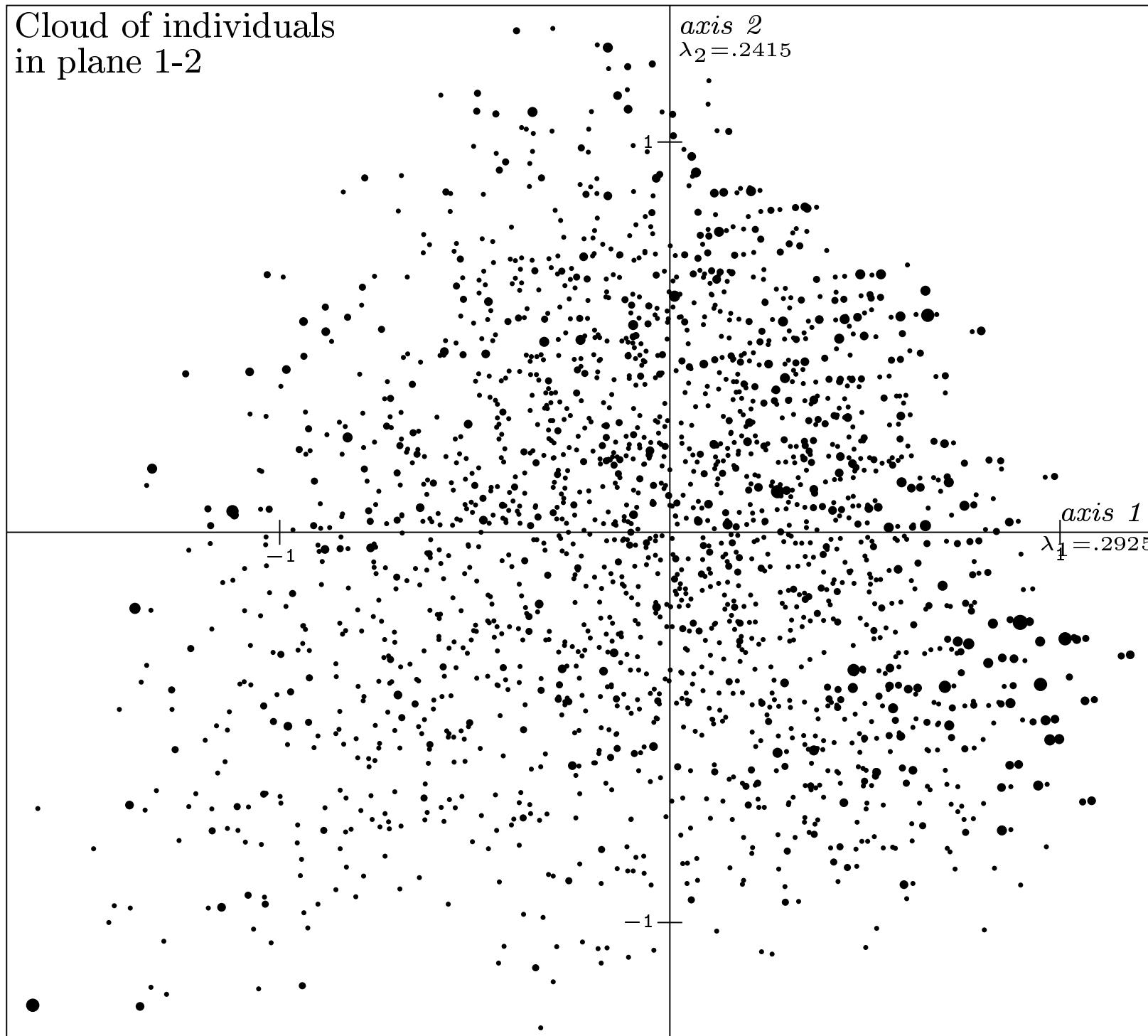
<sup>a</sup>See B. Le Roux & H. Rouanet (2004), p.203-210.

## Basic Results

- (i) the variances of axes (eigenvalues);
- (ii) the principal coordinates of modalities (categories) and of individuals;
- (iii) the contributions of modalities to axes;
- (iv) the geometric representation of the two clouds (cloud of modalities and cloud of individuals).

## Cloud of modalities in plane 1-2





## How many axes should be interpreted?

Eigenvalues: 10 eigenvalues exceed  $\bar{\lambda} = 1/6 = .1667$ :

	Eigen-values	modified rates
$\lambda_1$	.2925	.569
$\lambda_2$	.2415	.201
$\lambda_3$	.2248	.122
$\lambda_4$	.2073	.059
$\lambda_5$	.1950	
$\lambda_6$	.1832	
$\lambda_7$	.1790	
$\lambda_8$	.1758	
$\lambda_9$	.1733	
$\lambda_{10}$	.1688	

**Modified rates:** To assess importance of axes

Benzécri's formula, cf. Benzécri (1992), p.412, Le Roux & Rouanet (2004), p.200 &225

*3 axes will be interpreted* (cumulated modified rate 89%)

## Interpreting axes

It is done in the cloud of modalities and based on the modalities whose contributions to axis exceed a specified threshold (e.g. average contribution) according to the *Method of contributions of points and deviations*<sup>a</sup>, following the principle stated by Benzécri (1992, p. 405):

“Interpreting an axis amounts to finding out what is similar, on the one hand, between all the elements figuring on the right of the origin and, on the other hand between all that is written on the left; and expressing with conciseness and precision, the contrast (or opposition) between the two extremes.”

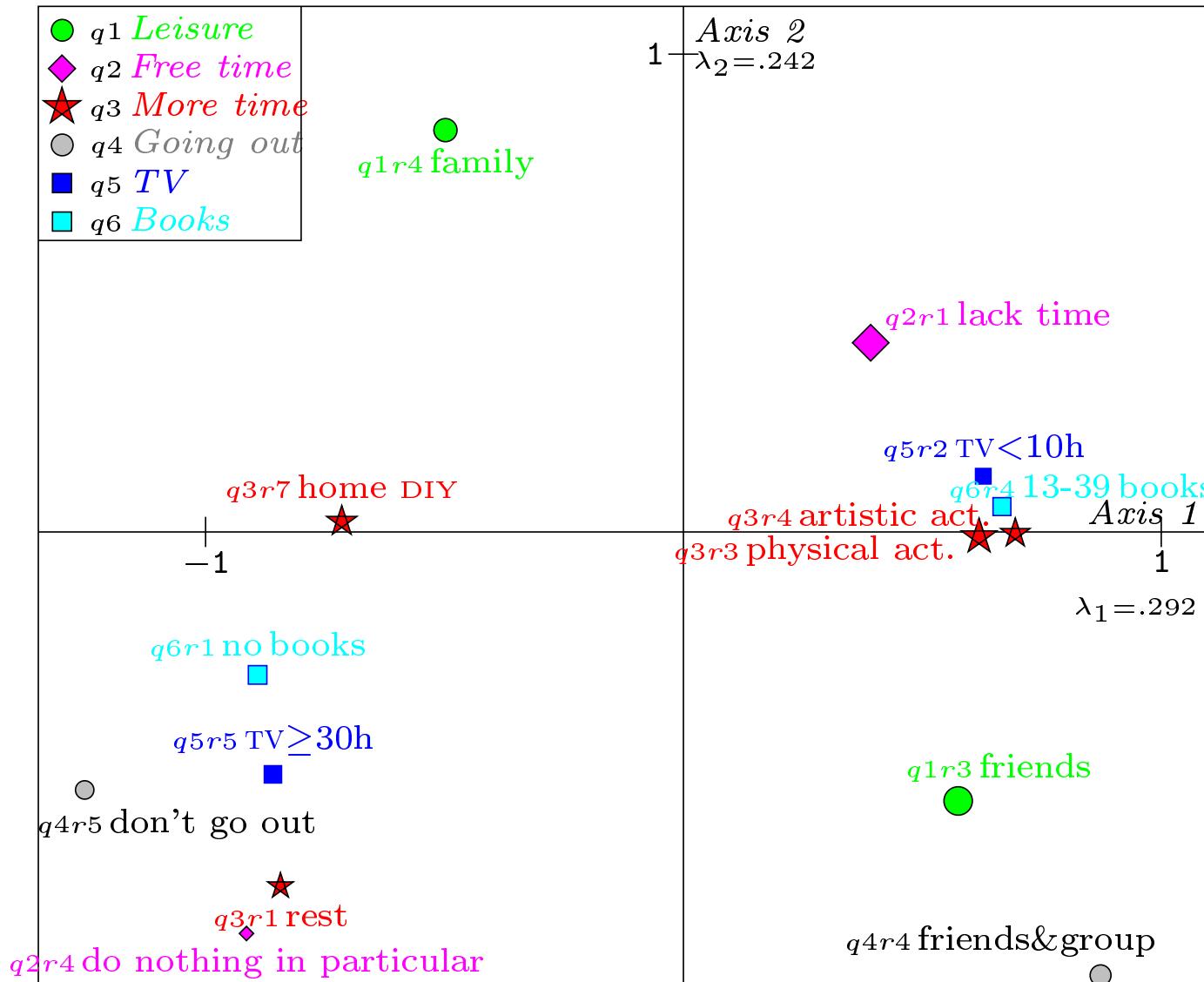
---

<sup>a</sup>Le Roux B. & Rouanet H. (1998). Interpreting axes in Multiple Correspondence Analysis: Method of contributions of points and deviations, in *Visualization of Categorical Data*, (eds Blasius & Greenacre), 197-220, Academic Press.

### 3.1 Interpretation of Axis 1

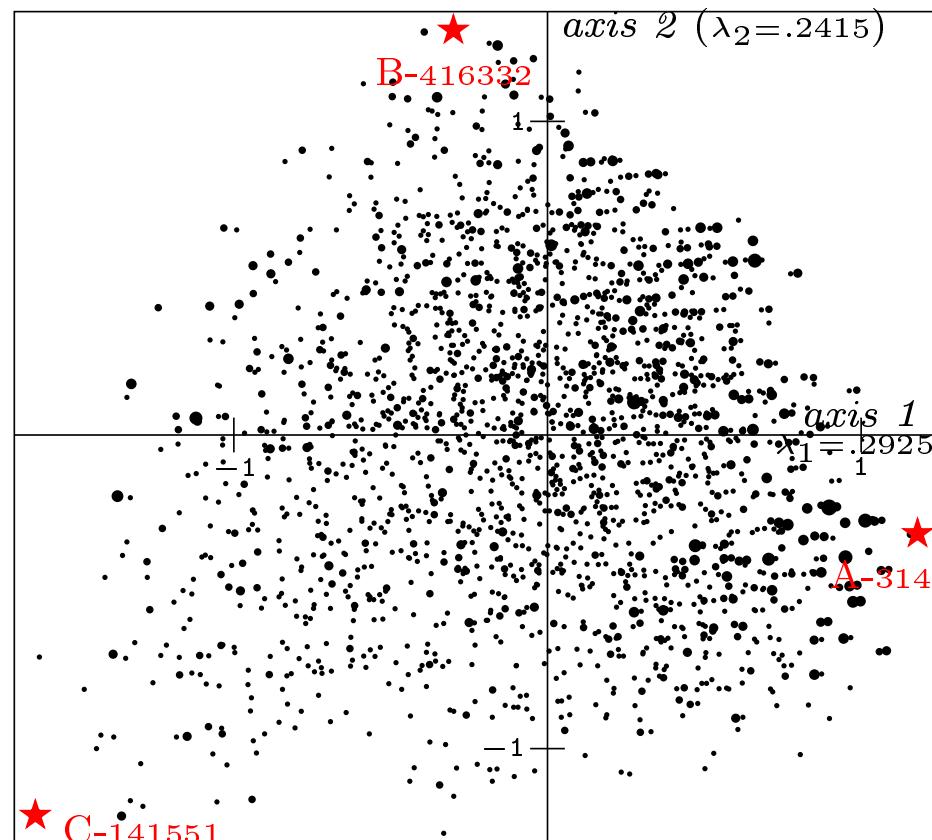
			left	right	deviation
<i>q4 (27%)</i> ● <i>Going out</i>	<i>r5 don't go out</i> <i>r4 friends&amp;group</i>		15.4	9.5	91.8
<i>q3 (20%)</i> ★ <i>More time</i>	<i>r3 physical activities</i> <i>r7 home DIY</i> <i>r1 rest</i> <i>r4 artistic activities</i>		4.5 4.5	4.6 4.0	
<i>q6 (18%)</i> ■ <i>Books</i>	<i>r1 no books</i> <i>r4 13-39 books</i>		10.0	5.3	83.7
<i>q1 (14%)</i> ● <i>Leisure</i>	<i>r3 friends</i> <i>r4 family</i>		3.7	7.3	75.2
<i>q5 (13%)</i> ■ <i>TV</i>	<i>r5 ≥ 30h</i> <i>r2 &lt; 10h</i>		8.2	3.6	85.6
<i>q2 (9%)</i> ◇ <i>Free time</i>	<i>r1 lack time</i> <i>r4 do nothing in particular</i>		3.5	3.6	71.6
			49.8	37.9	

Relative contribution (percentage) of each modality to axis, written either in column “left” or in column “right” according to its position on graph; relative contribution of deviation to question.

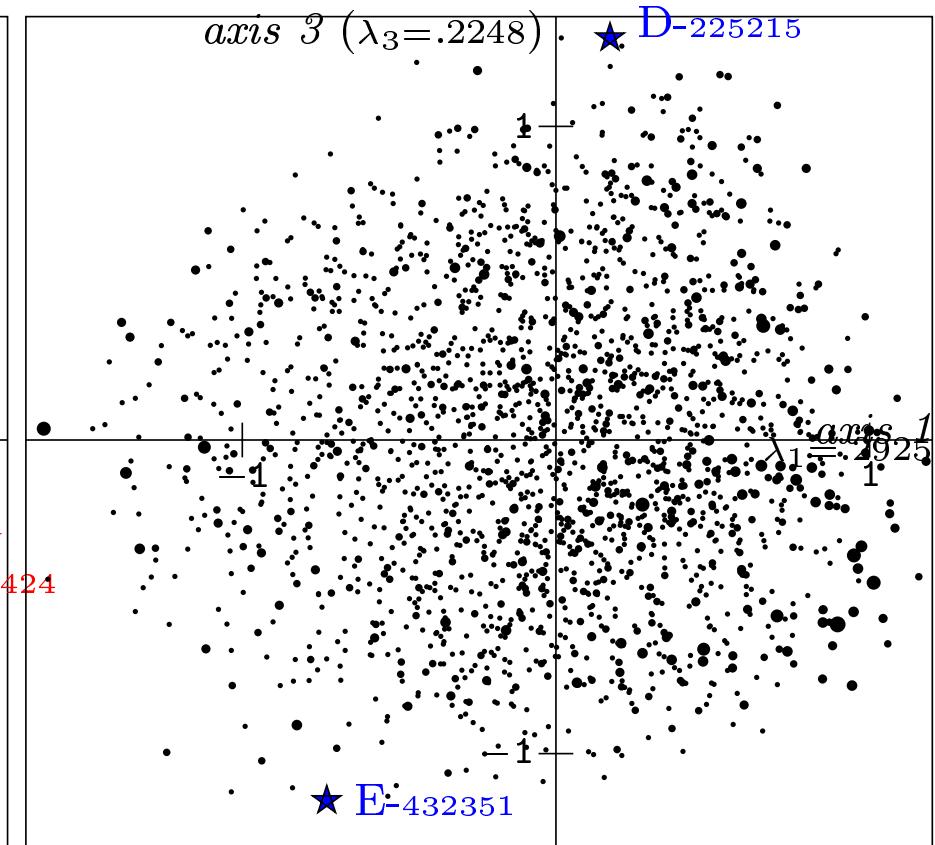


Interpretation of Axis 1: 14 modalities most contributing to axis.

## Investigation of the cloud of individuals



Plane 1-2



Plane 1-3

## Subclouds and modality mean points

Subcloud of the individuals having chosen one modality → its mean point is called a *modality mean-point*.

*Fundamental Property:*

Coordinate of the modality mean-point =  $\sqrt{\lambda} y$

( $y$  is the coordinate of the modality in the space of modalities)

## Concentration ellipses

The concentration ellipse<sup>a</sup> of a subcloud is such that the half-axis of the ellipse is along the principal direction of the subcloud projected in the plane under study and its length is equal to  $2\sqrt{\lambda'}$ .

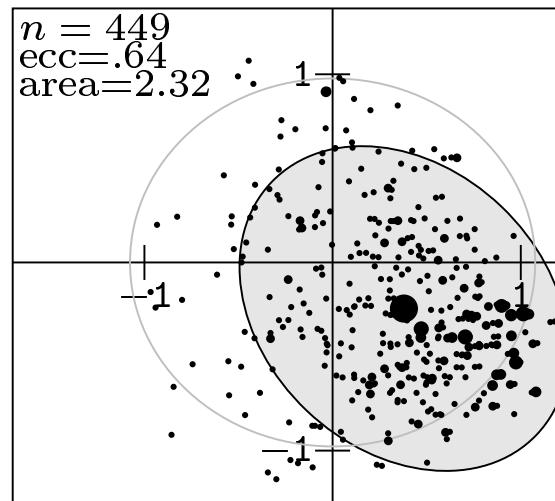
A uniform distribution over the interior of the ellipse has the *same variance* as the subcloud.

For a normally-shaped cloud, the concentration ellipse contains about *86% of the points* of the cloud.

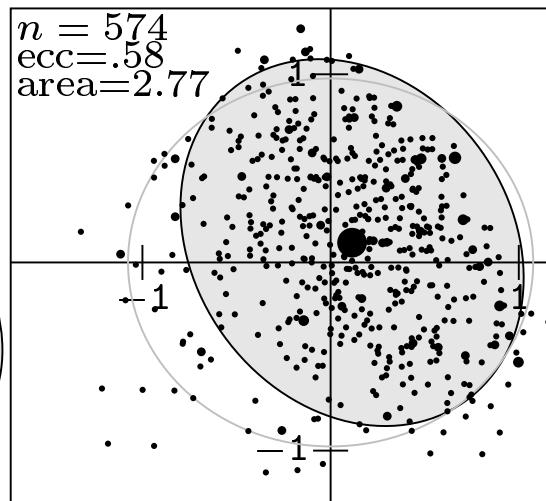
Concentration ellipses are especially useful for studying families of subclouds induced by a structuring factor or a clustering procedure: see e.g. Age in the Culture Example.

---

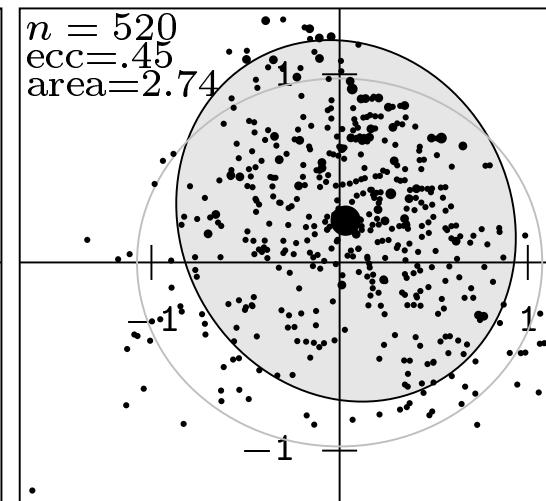
<sup>a</sup>see Cramér, 1946, p. 284; Le Roux & Rouanet (2004), p.95-100



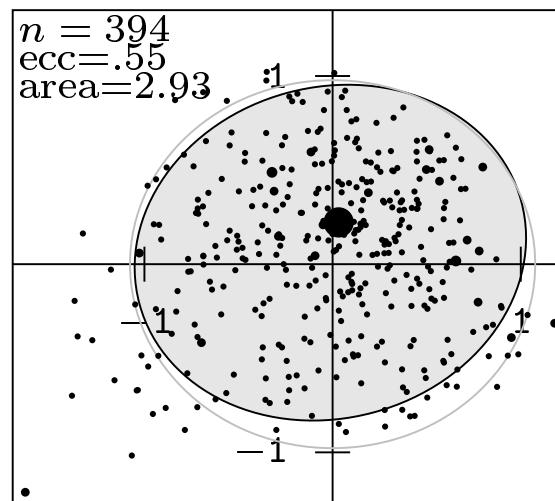
18-25 years class



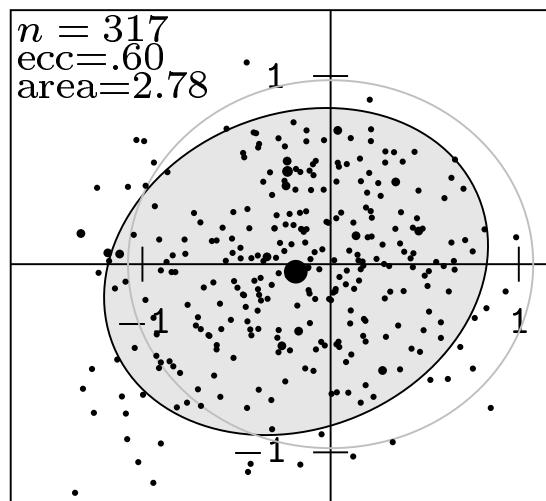
26-35 years class



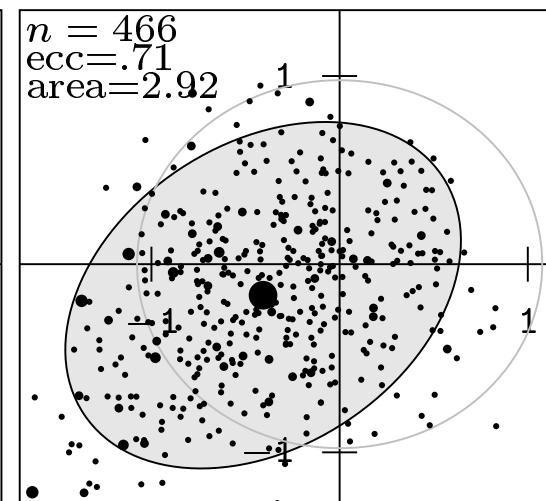
36-45 years class



46-55 years class



56-65 years class



&gt;65 years class

## 4 Final message about MCA

MCA is a method for the *Geometric Analysis* of questionnaires

Investigate the *cloud of individuals*: individuals carry all the information!

## References

- BENZÉCRI J-P. (1977). Sur l'analyse des tableaux binaires associés à une correspondance multiple [On the analysis of binary tables associated with a multiple correspondence], *Les Cahiers de l'Analyse des Données*, 2, 55-71 (from a mimeographed note of 1972).
- BENZÉCRI J-P. (1992). *Correspondence Analysis Handbook*, New York, Dekker.
- ESCOFIER B. & PAGÈS J. (1988), *Analyses factorielles simples et multiples* [Simple and Multiple factor Analyses] (chapter 3 on MCA). Paris: Dunod.
- GREENACRE M. (1984). *Theory and Applications of Correspondence Analysis* (chapter 5 on MCA). London: Academic Press.
- LE ROUX B. (2006). What is MCA?, *Research Methods Festival*, Oxford, [www.ccsr.ac.uk/methods/festival/programme/Wita/leroux.pdf](http://www.ccsr.ac.uk/methods/festival/programme/Wita/leroux.pdf)
- LE ROUX B., & ROUANET H. (2004). *Geometric Data Analysis; From Correspondence Analysis to Structured Analysis*. Dordrecht:

Kluwer (the chapter 5 contains a detailed presentation of MCA and an extensive illustration: “the Culture Example”, p.221-251)

— LEBART L., MORINEAU A. & WARWICK K.M. (1984). *Multivariate Descriptive Statistical Analysis: Correspondence Analysis and Related Techniques for Large Matrices* (chapter 4 on MCA). New York: Wiley.

### **Recent substantive studies**

- BONNET P., LE ROUX B. & LEMAINE G. (1996). Analyse géométrique des données: une enquête sur le racisme [Geometric Data Analysis: a survey on racism], *Mathématiques et Sciences Humaines*, 136, 5-24.
- BOURDIEU P. (1999). Une révolution conservatrice dans l'édition [A conservative revolution in publishing], *Actes de la Recherche en Sciences Sociales*, Vol. 126-127, 3-28.
- CHICHE J., LE ROUX B., PERRINEAU P. & ROUANET H. (2000). L'espace politique des électeurs français à la fin des années 1990 [The political space of French electors in the late 1990s]. *Revue française de sciences politiques*, 50, 463-487.

- HJELLBREKKE J., LE ROUX B., KORSNES O., LEBARON F., ROSENlund L. & ROUANET H. (to appear). The Norwegian field of Power Anno 2000. *European Societies*.
- LE ROUX B. (2006). Que pense votre député de la mondialisation (What does your deputy think of the globalization).  
[www.telos-eu.com/2006/05/que\\_pense\\_votre\\_depute\\_de\\_la\\_m.php](http://www.telos-eu.com/2006/05/que_pense_votre_depute_de_la_m.php)
- LE ROUX B. & ROUANET H. (2003). Geometric Analysis of Individual Differences in Mathematical Performance for EPGY Students in the Third Grade. [www-epgy.stanford.edu/research/](http://www-epgy.stanford.edu/research/).
- See chapter “Case Studies” in Le Roux & Rouanet (2004) (op. cit.)

## About software

These results have been obtained using ADDAD, ellipse and EyelID software freely available from my website

[math-info.univ-paris5.fr/~lerb/](http://math-info.univ-paris5.fr/~lerb/) under the “Logiciels” heading.

All results presented here are obtainable from the September 2006 version of SPAD software distributed by SPAD company ([www.spad.eu](http://www.spad.eu)).