



Exploratory multidimensional statistical methods Introductory presentation

Anàlisi de Dades i Explotació de la Informació

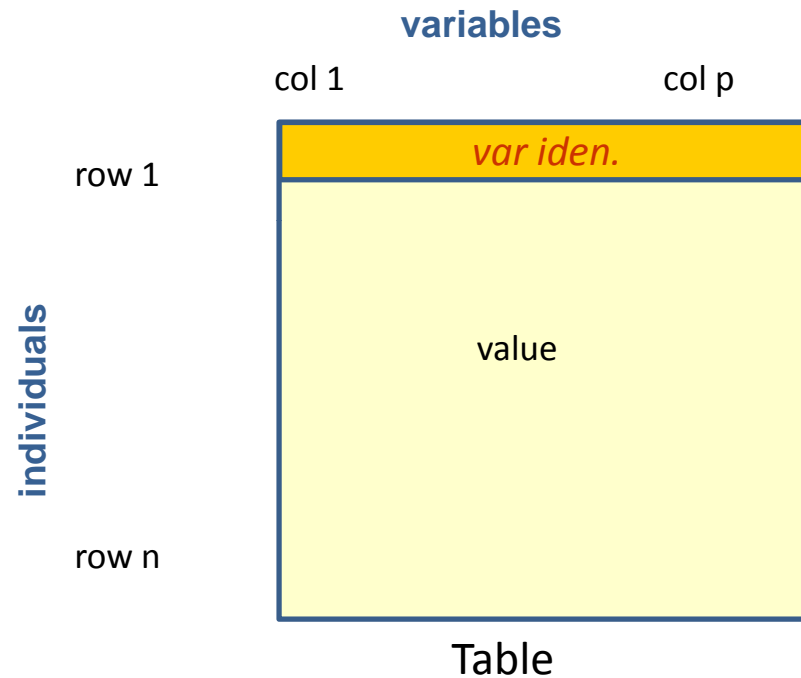
Grau d'Enginyeria Informàtica.

Information System track

Prof. Mónica Bécue Bertaut & Lidia Montero

Monica.becue@upc.edu lidia.montero@upc.edu

Information as stored in rectangular tables



Information is stored in tables

Rows of the table

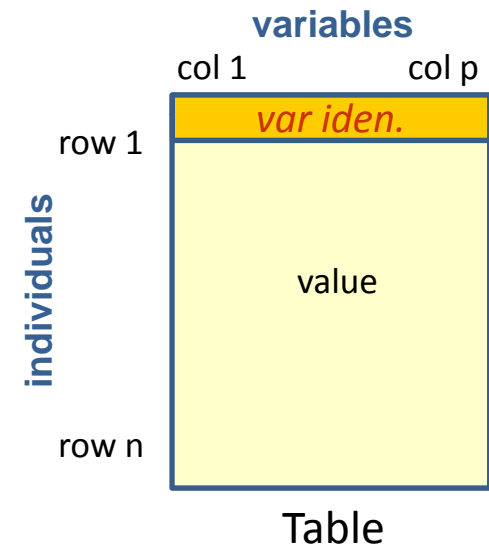
Individuals or instances , sample, example, record, ...
forming the sample under study, extracted from a population

Columns of the table

Variables or “attributes”.

Main attribute types : quantitative, binary, nominal, ordinal,
interval, ratio, textual, ...

Variables/ Attributes observed on the individuals (in the whole
population, the variable follows a probabilistic law...) or
constructed a posteriori



Encoding the variables

- **Continuous or quantitative**

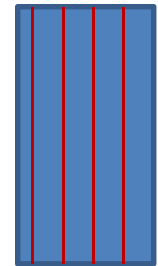
- Interval : temperature,... Laplace-Gauss distribution)
income, ... Exponential distribution)
- Count data (“number of words of a sentence of a same author”,
.... Poisson distribution)



One column

- **Categorical**

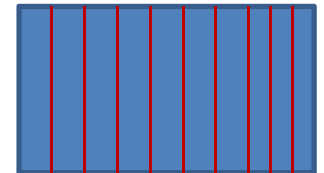
- Binary (yes/no variable, boolean attribute.
Binomial/Hypergeometric distribution)
- Nominal (marital status, region, Multinomial distribution)
- Ordinal (clothe’s size, social class,, Ordered multinomial)



as many columns as categories

- **Frequency data**

- Series of columns with count data that must be treated as a whole
(counts of the occurrences of the different words used to answer
an open-ended question; counts of mortality data by causes;
counts of the occurrences of all the different species present in an ecological site)



as many columns as species

- **Textual data**

Role of variables

- **Response**

Variable to be explained or predicted
either continuous, categorical or frequency

- **Explanatory**

Variables used to explain the behaviour of the response variables
continuous or categorical or frequency

Types of data matrix

With or without response(s) variable

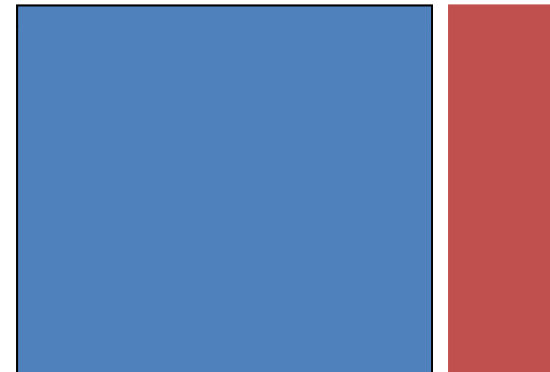
i.e. transactions data



Data to explore, to describe, to find associations (i.e. itemsets), ...

Inputs

Output(s)



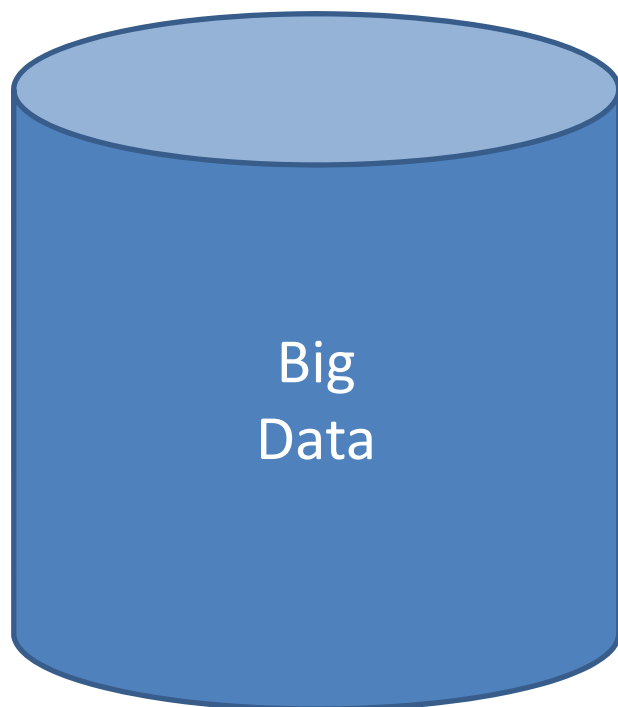
we want to **find a model to predict the response**

R function available in FactoMineR

We will use FactoMineR Package (cran R)

You can also consult (and download this R function from) <http://factominer.free.fr/> where a large documentation is provided, with theoretical background, examples, tutorials and so on.

Some details about the exploratory methods



Big data require
**Exploratory
multidimensional
statistical methods**

JOB: Data Scientist Technology – Business Intelligence | Mind Candy, London

Mind Candy – Posted by [Advertiser](#) – London, England, United Kingdom

Job Description

We are looking for truly talented individuals to become an integral part in driving the (...)

The Role

Due to our continued success Mind Candy is rapidly expanding and we have a truly fantastic opportunity for a Data Scientist to come on board and play a key role in (...)

Minimum Requirements:

Good business and technical skills in data analytics. Technical skills must include:

Highly proficient data mining skills in **small and very large data sets**.

Great ETL skills using a variety of languages (e.g. SQL, **R**, Python, Scala) and **big data tools** (e.g. Hive, Scalding, Pig, Elastic MapReduce).

Great statistical skills and a **passion for data** and **data visualisation**.

Ability to continuously adapt to the data needs in a rapidly changing environment.

This would include quickly and efficiently integrating new data sources using various methods (from internal or external databases, using REST API, etc.).

Experience and confidence in gathering business requirements from the product teams and delivering reports, analysis and innovative, fit for purpose information solutions.

Experienced in managing your own priorities based on business goals.

Strong communication skills.

Experience of working in an Agile environment.

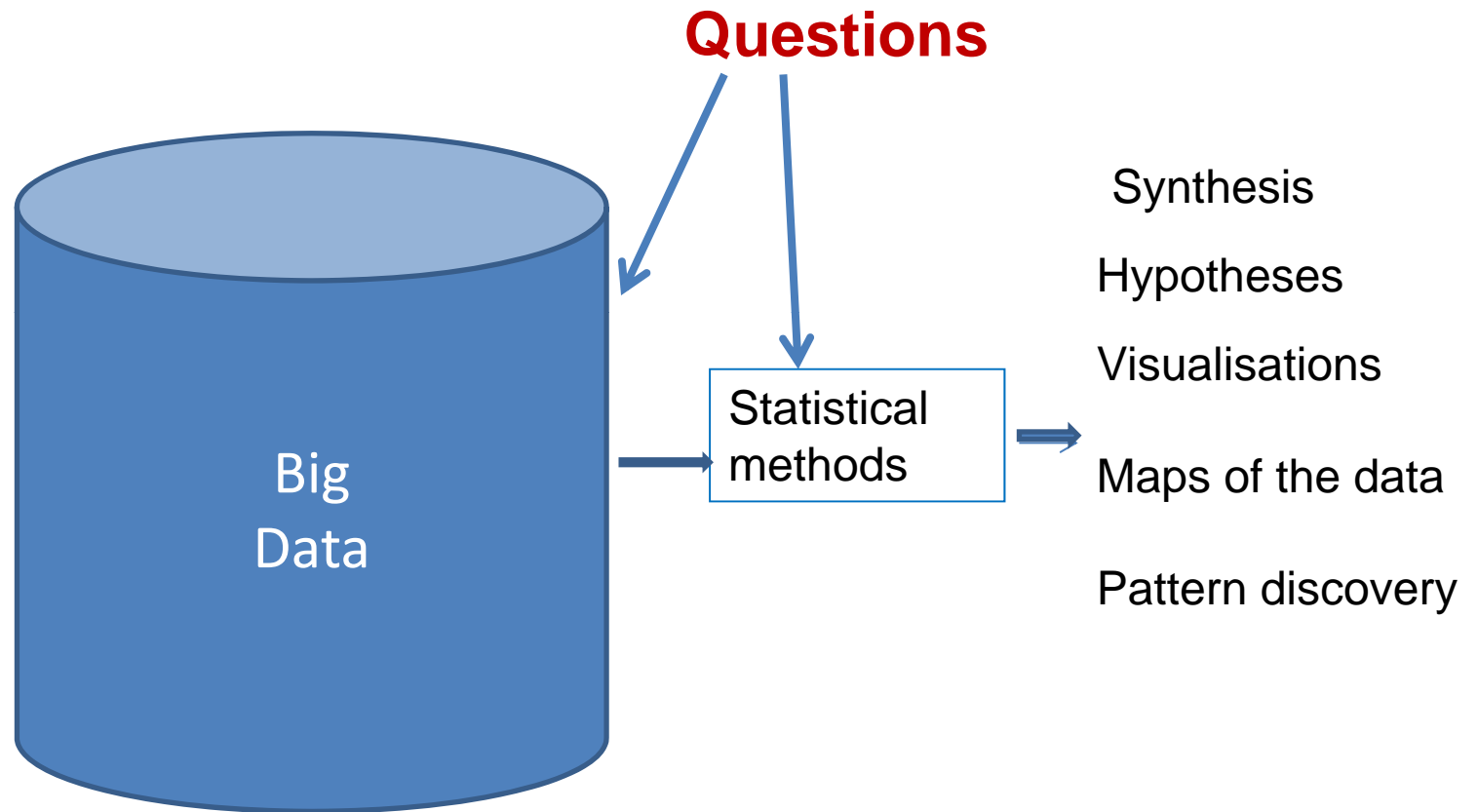
Preferred Requirements: (...)

Resumen de los tipos de descripciones de **Big Data** que Ward y Barker han descubierto de varias organizaciones influyentes:

1. **Gartner.** En 2001, un informe de Meta (hoy día Gartner) tomó nota del aumento del tamaño de los datos, la tasa de aumento a la que se producen y la creciente **variedad de formatos** y representaciones empleadas. Este informe es anterior a la expresión 'big data', pero proponía una definición triple con tres 'V': **volumen, velocidad y variedad**. Desde entonces, esta idea se ha hecho muy popular y, a veces, incluye una cuarta V: **veracidad**, para cubrir la cuestión de la confianza y la incertidumbre.
2. **Oracle.** 'Big data' es la derivación de valor a partir de la toma de decisiones de negocio en función de bases de datos relacionales tradicionales, aumentada con nuevas **fuentes de datos no estructurados**.
3. **Intel.** Las oportunidades de trabajo con grandes volúmenes de datos surgen en organizaciones que generen un **promedio de 300 terabytes de información** a la semana. La clase de datos más común es la de las transacciones comerciales almacenadas en bases de datos relacionales, seguida de **documentos**, correo electrónico, datos de sensores, **blogs y redes sociales**.



4. **Microsoft.** "**Big data**" es un término cada vez más utilizado para describir el proceso de aplicación de una significativa potencia de computación (lo último en el aprendizaje de máquinas e inteligencia artificial) a conjuntos de información de **enorme tamaño** y, a menudo, de **alta complejidad**".
5. El proyecto de código abierto **MIKE** (siglas en inglés de **Method for an Integrated Knowledge Environment**). El proyecto MIKE argumenta que **los grandes volúmenes de datos no tienen que ver con el tamaño sino con la complejidad**. Por consiguiente, lo que define un conjunto de datos como 'big data' es su alto grado de permutaciones e interacciones.
6. El **Instituto Nacional de Estándares y Tecnología** de EEUU. El Instituto afirma que los grandes volúmenes de datos se refieren a **aquellos que "superan la capacidad o la habilidad de los métodos y sistemas actuales o convencionales"**. En otras palabras, la noción de 'grande' está relacionada con el estándar de computación actual.

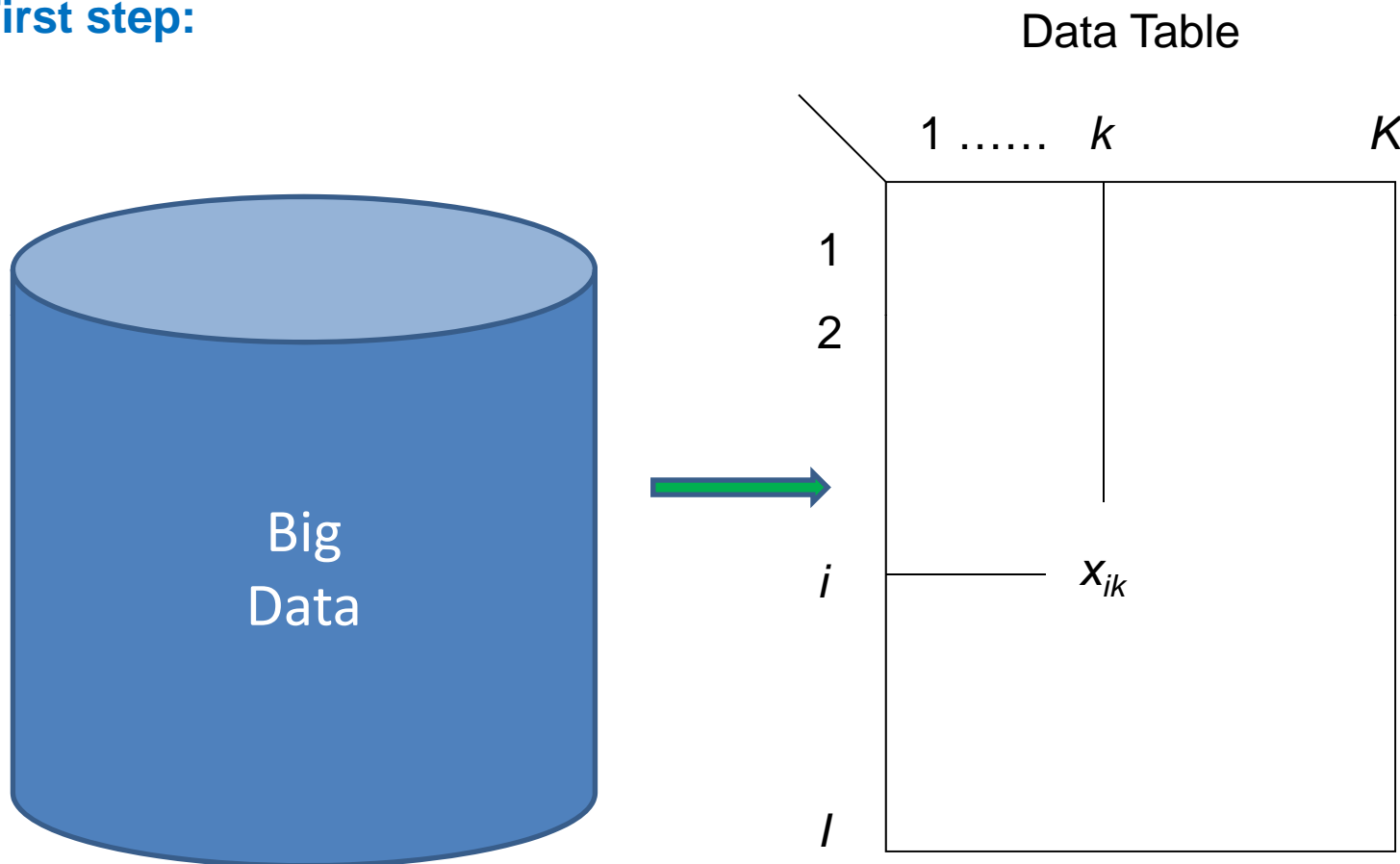




Among the Statistical methods: **multidimensional exploratory statistical methods**

What are they for?

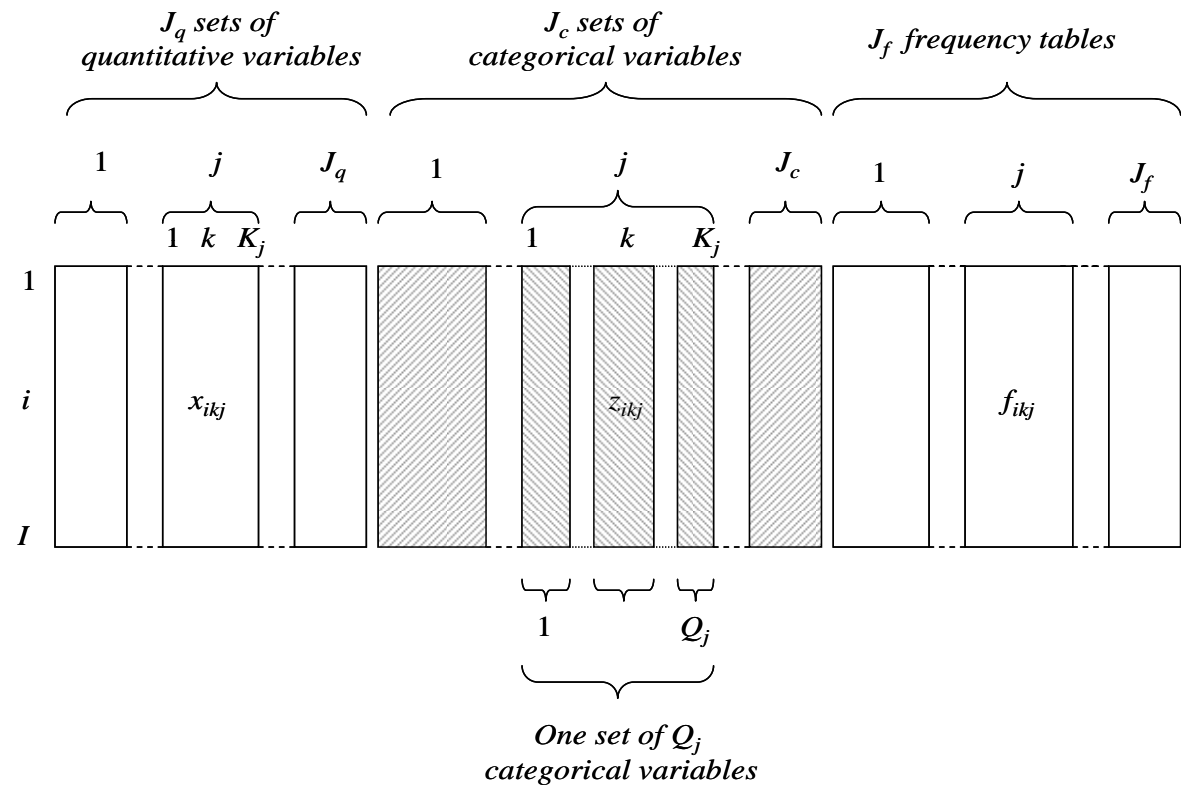
First step:



It tabular data led to data analysis, it can also be pointed out that tabular data led to the computer (Fionn Murtagh *Electronic Journ@l for History of Probability and Statistics*- Vol 4, n°2; Décembre/December 2008 www.jehps.net)

Or First
step:

Multiple Table



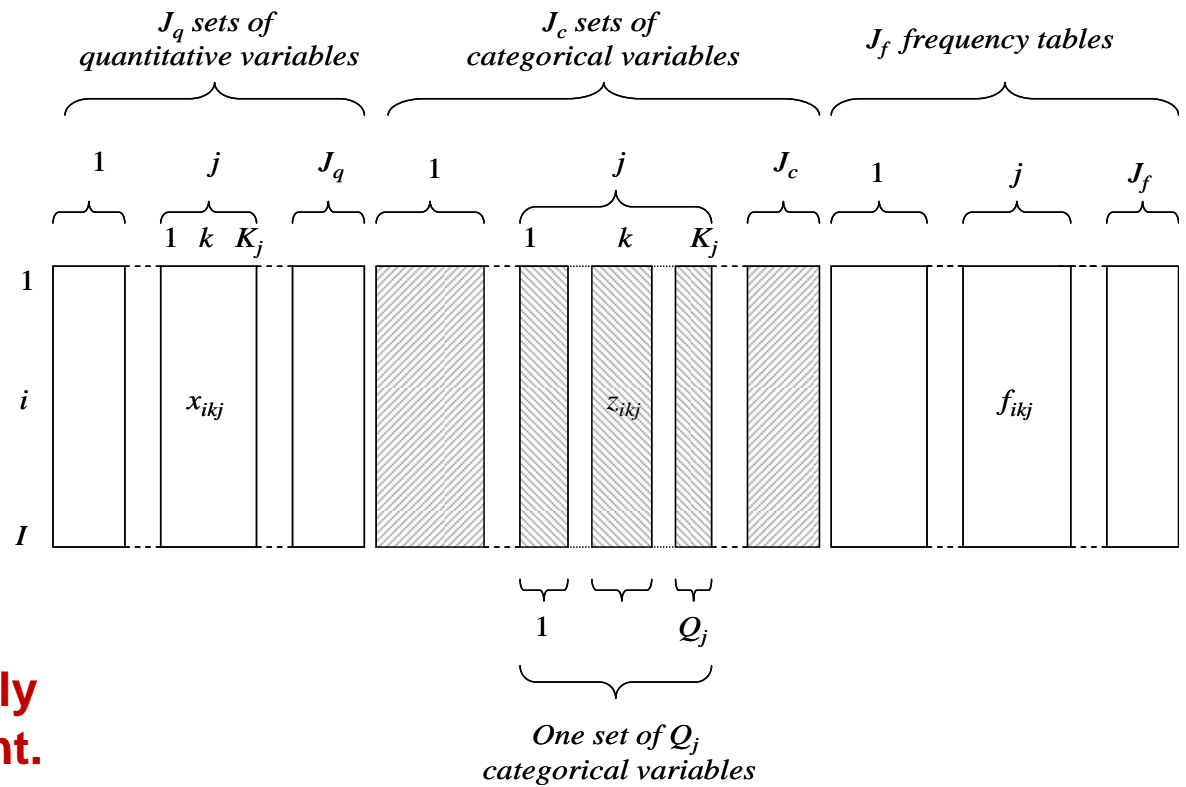
Or First
step:

Multiple Table



**Data
Coding**

**Extremely
important.
This
conditions
the results**



And now exploratory data analysis gets going!

To answer the questions

“In data analysis numerous disciplines have to collaborate.

The role of mathematics, although essential, remains modest in the sense that classical theorems are used almost exclusively, or elementary demonstration techniques.

But it is necessary that certain abstract conceptions penetrate the spirit of the users, who are the specialists collecting the data and having to orientate the analysis in accordance with the problems that are fundamental to their particular science.”

Fionn Murtagh

Electronic Journ @I for History of Probability and Statistics

Vol 4, n°2; Décembre/December 2008

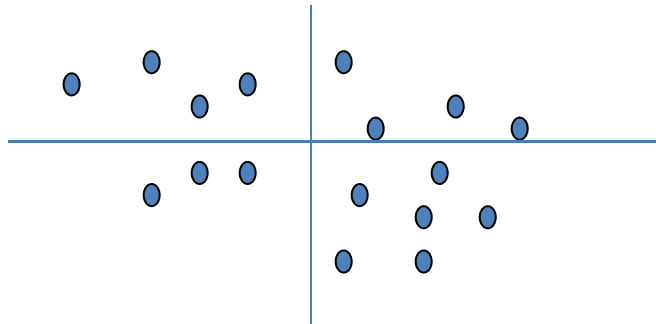
www.jehps.net



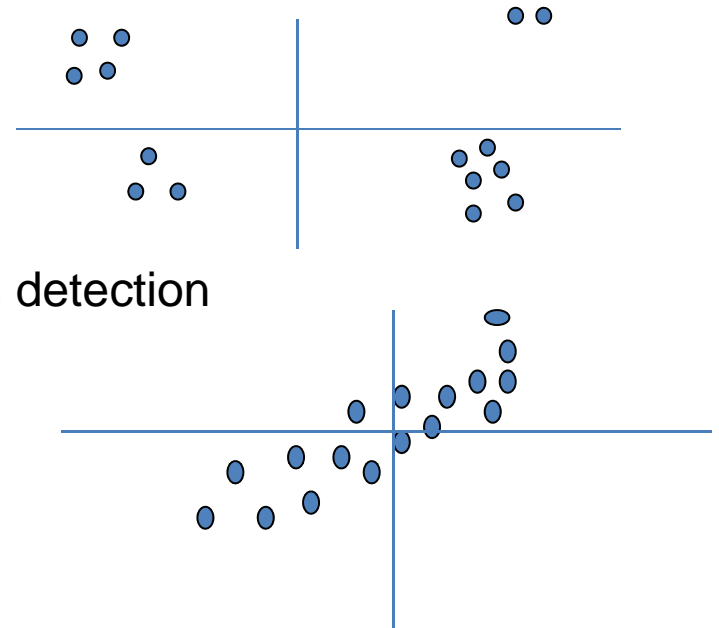
Among the Statistical methods: **multidimensional exploratory statistical methods**

What are they for?

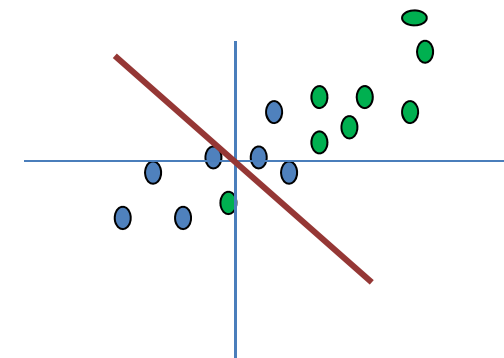
Visualisation of data



pattern detection

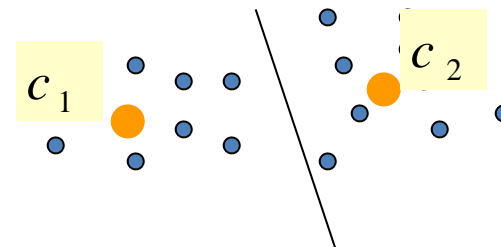


Mapping the individuals from their similarity



For classifying the individuals

or clustering them





Examples

1. Method: PCA

Ejemplo: cata de chocolates

Ejemplo: Chocolates

- 10 chocolates negros
 - 3 marcas: Lindt, Valrhona y Hacendado
 - Porcentaje de cacao entre 55% y 85%
- Método de recogida de datos: QDA (Quantitative Descriptive Analysis)
- 16 panelistas y 2 sesiones
- 14 descriptores
 - Olor: cacao, leche
 - Sabor: azúcar, ácido, amargo, cacao, leche, caramelo, vainilla
 - Características: astringencia, crujiente, fusión en la boca, pegajoso, granuloso
- Notas entre 0 y 10
- **Diseño de experimentos completo balanceado para los rangos y efectos de arrastre de orden 1**

Data capture

sesión	panelista	rango	producto	primero	O.Cacao	O.leche	Azúcar	Ácido	Amargo	Cacao	Leche	Caramelo	Vainilla	Astringencia	Crujiente	Fusión boca	Pegajoso	Granuloso
1	2	1	VALRHONA (64)	1	6	1	4	7		6	1	1	1	3	3	7	1	2
1	2	2	VALRHONA (66)	0	7	1	4	6	6	7	1	1	1	4	3	7	1	2
1	2	3	LINDT-NS (70)	0	4	1	6	5	3	5	1	2	2	2	2	5	1	2
1	2	4	HACENDADO (72)	0	8	1	2	4	7	8	1	1	1	8	4	8	1	3
1	2	5	HACENDADO (55)	0	4	1	7	2	1	4	1	3	1	1	2	6	3	4
1	2	6	LINDT-EF (85)	0	8	1	2	6	7	9	1	1	1	8	6	7	1	1
1	2	7	HACENDADO (85)	0	8	1	2	7	6	7	1	1	1	7	4	5	1	2
1	2	8	VALRHONA (70)	0	6	1	5	7	5	5	1	1	1	2	6	6	1	2
1	2	9	LINDT-EF (70)	0	7	1	2	7	6	6	1	1	1	7	6	6	1	2
1	2	10	VALRHONA (85)	0	9	1	1	6	8	9	1	1	1	9	6	7	2	2
2	20	8	LINDT-EF (70)	0	8	2	6	7	7	6	2	2	1	6	3	7	4	1
2	20	9	HACENDADO (55)	0	7	2	9	2	2	4	2	2	2	3	2	4	5	1
2	20	10	HACENDADO (85)	0	8	2	3	8	10	7	2	1	1	9	6	7	3	1

Building the products x variables table

Diagram illustrating the process of building the products x variables table:

Diagram components:

- productos** (products)
- descriptores** (descriptors)
- medias (ajustadas)** (adjusted means)
- Producto x juez** (Product x judge)

Blue arrow indicates the flow from the diagram components to the data table.

ProductoxJuez-Descriptores

	o_cacao	o_leche	s_azucar	s_acido	s_amargo	a_cacao	a_leche	a_caramelo	a_vainilla	astringencia	crujiente	fusion	pegajoso	granuloso
VAL(66)*2	7	1	4	6	6	7	1	1	1	4	3	7	1	2
LINDT(70)*2	4	1	6	5	3	5	1	2	2	2	2	5	1	2
HAC(72)*2	8	1	2	4	7	8	1	1	1	8	4	8	1	3
HAC(55)*2	4	1	7	2	1	4	1	3	1	1	2	6	3	4
LINDT(70)*20	8	2	6	7	7	6	2	2	1	6	3	7	4	1
HAC(55)*20	7	2	9	2	2	4	2	2	2	3	2	4	5	1
HAC(85)*20	8	2	3	8	10	7	2	1	1	9	6	7	3	1

Objetivos de ACP

El estudio de los individuos (A)

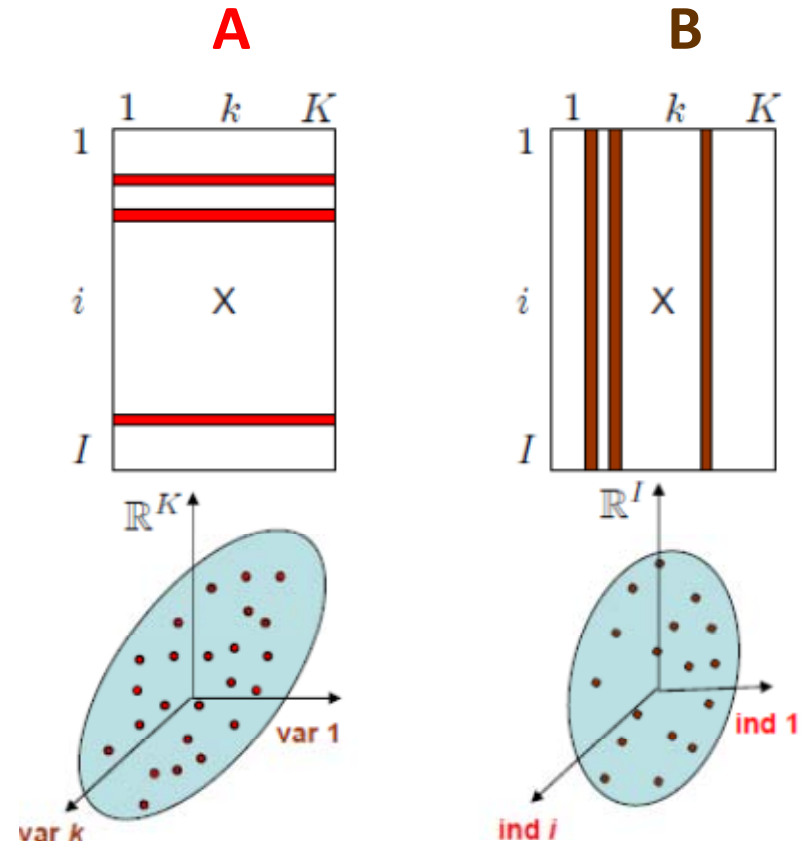
- Similitud entre los individuos respecto a todas las variables
- partición entre los individuos

El estudio de las variables (B)

- Relaciones lineales entre las variables = Visualización de la matriz de correlaciones (**S**)
- Encontrar variables sintéticas

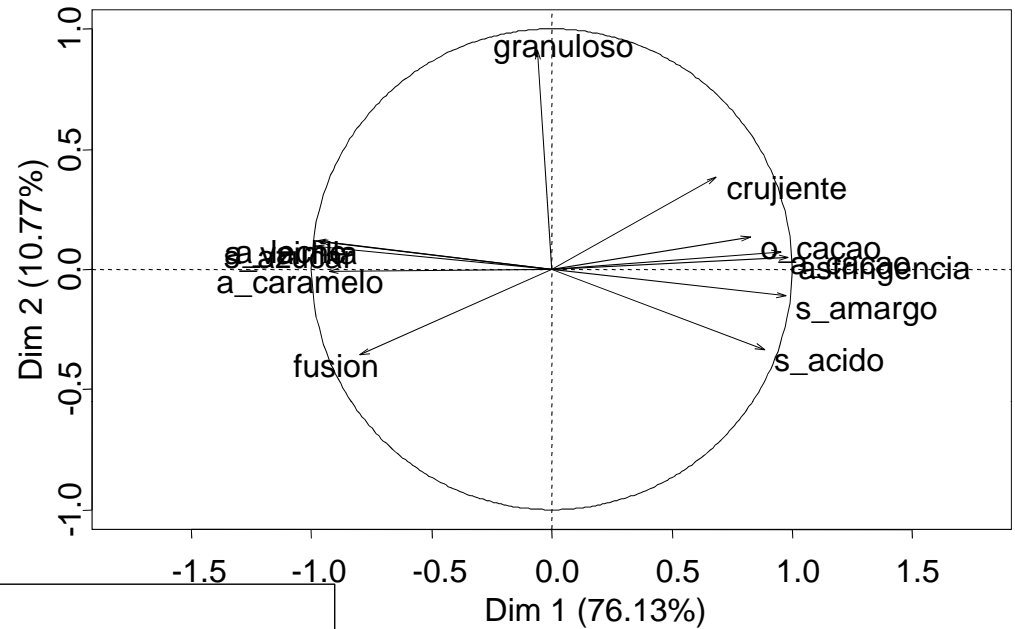
Relación entre los dos estudios

- Caracterización de grupos de individuos por variables; individuos particulares para entender mejor la relación entre las variables

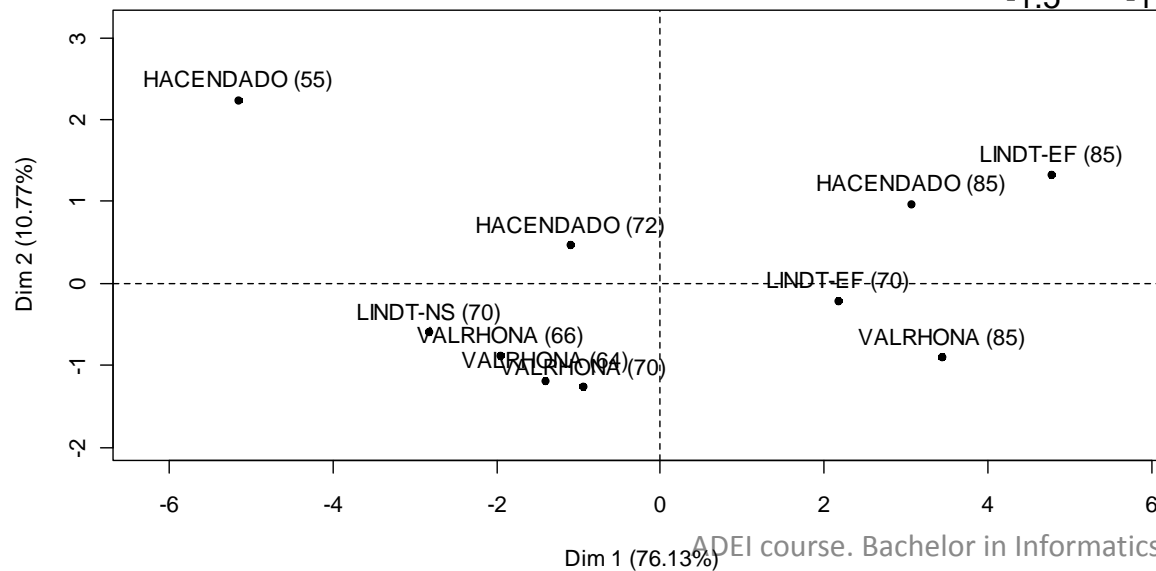


Analysis by PCA

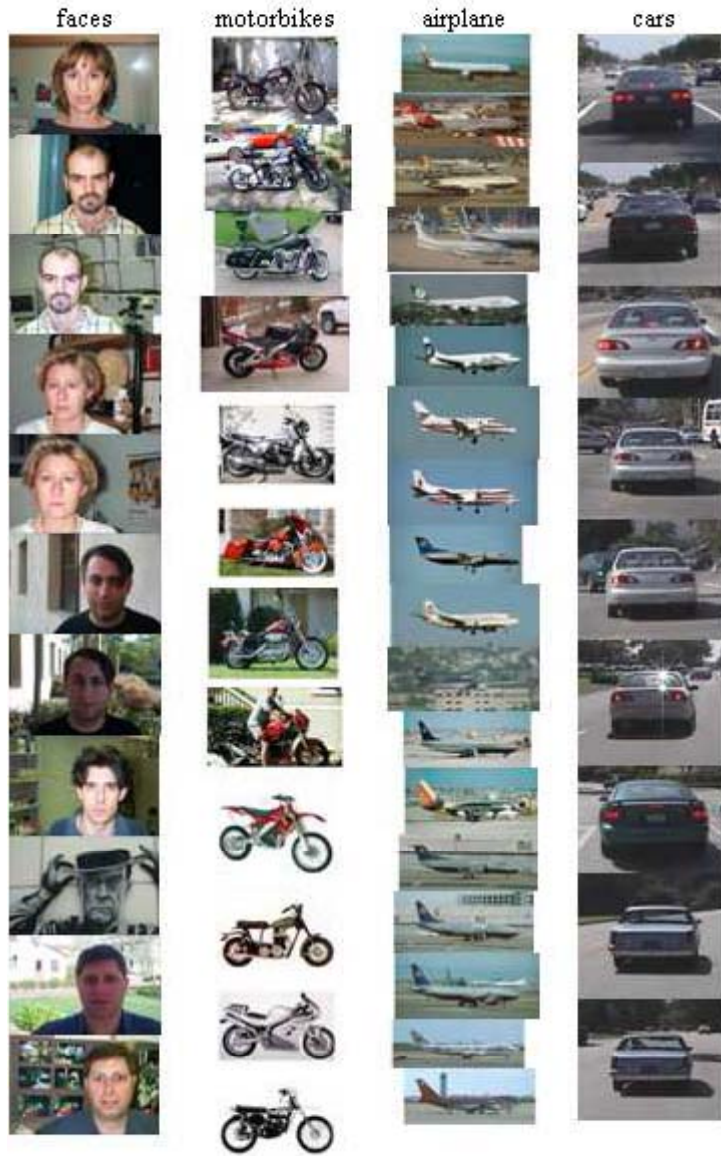
Variables factor map (PCA)



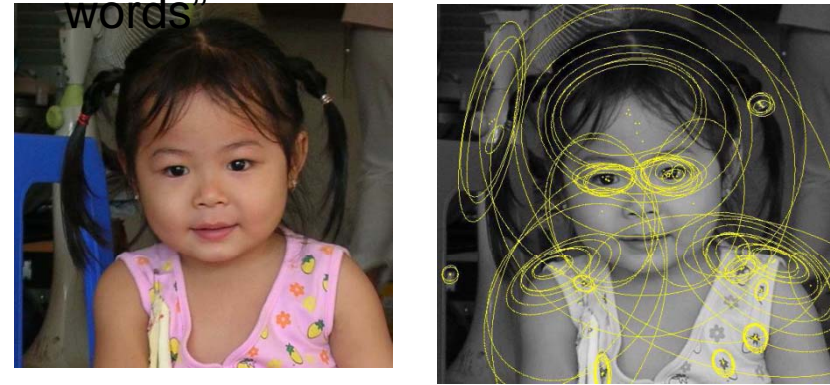
Individuals factor map (PCA)



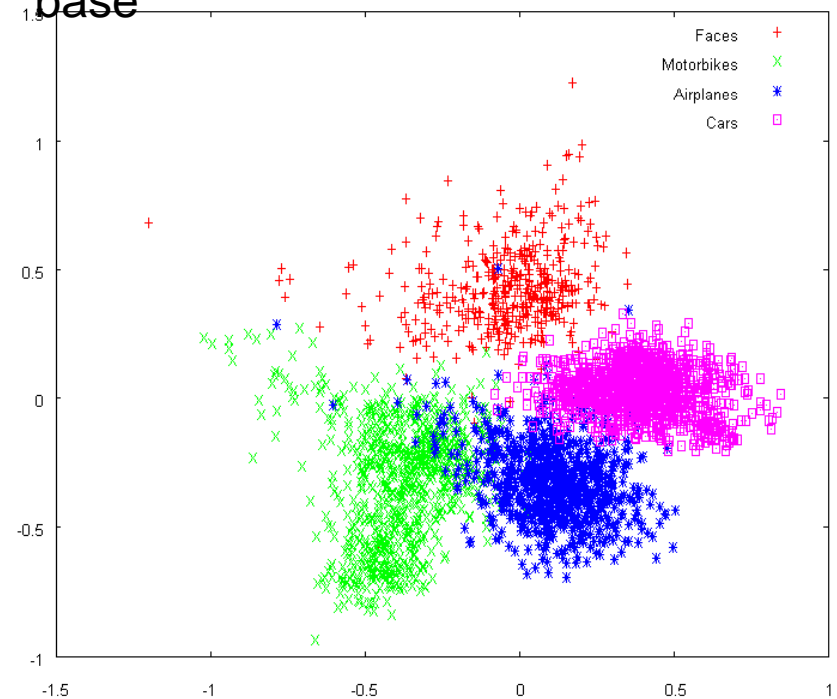
2. Method: CA Image data base



Coding the images into “visual words”



CA is able to “organize” the image data base



3. Multiple Correspondence Analysis (MCA)

Example: Perfumes



Angel



Aromatics
Elixir



Chanel n°5



Cinéma



Coco
Mademoiselle



L'Instant



Lolita
Lempicka



Pleasures



Pure Poison



Shalimar



J'adore
(ET)



J'adore (EP)

The panellists (100)





Reagrupar y describir los grupos (Juez 1)



18 de agosto de 2012

La sensometría estadística aplicada a la
ADEI course. Bachelor in Informatics
percepción sensorial de productos
Engineering. Session 3. Teaching. Tomás
alimentos
Aluja & Lidia Montero



Reagrupar y describir los grupos (Juez 2)

« léger, vanille,
orange »



« fort,
boisé »



« vieux, WC »



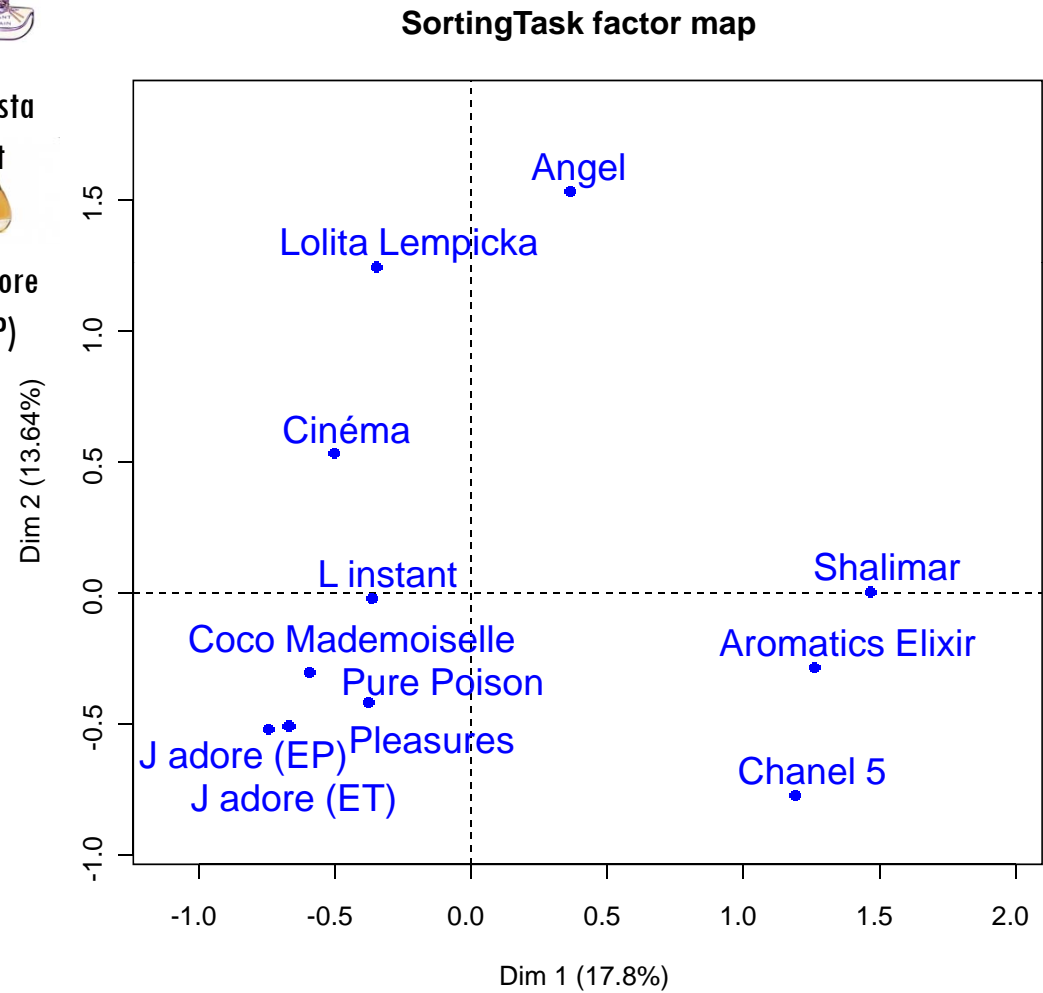
Data coding

- Una tabla con
 - ❖ Los perfumes en línea
 - ❖ Los jueces en columna
 - ❖ Cada juez se considera como una variable cualitativa, cuyas modalidades (categorías) son las palabras utilizadas
- **Tratamiento estadístico:** Análisis de Correspondencias Múltiple (ACM)

produit	juge 12	juge 13	juge 14	juge 15	juge 16
Angel	fleuri doux	fruité fort	vanillé épice esprit des îles	à manger sucré	nourriture épice
Aromatic Elixir	fort homme	capiteux grand-mère	rude fort	le vieux	ménager cire
Chanel n°5	Gr 4	capiteux grand-mère	toilettes	savon	connu classique
Cinéma	fleuri artificiel herbe	fruité moyen	sucré	doux	nourriture épice
Coco Mademoiselle	fleuri doux	fruité moyen	douceur fleuri	doux	connu classique
J'adore (EP)	fleuri doux	sucré faible	douceur fleuri	fleuri	connu classique
J'adore (ET)	fleuri artificiel herbe	sucré faible	douceur fleuri	fleuri	connu classique
L'instant	fleuri doux	fruité fort	sucré	le vieux	fleuri
Lolita Lempicka	fleuri doux	fruité moyen	vanillé épice esprit des îles	à manger sucré	nourriture épice
Pleasures	fort homme	fruité fort	sucré	fleuri	fleuri
Pure Poison	fleuri doux	acidulé désodorisant	douceur fleuri	doux	fleuri
Shalimar	fleuri artificiel herbe	fort lavande eau de cologne	renfermé agressif	le vieux	ménager cire

4. Perfumes

Representation of the rows





Representation of the labels

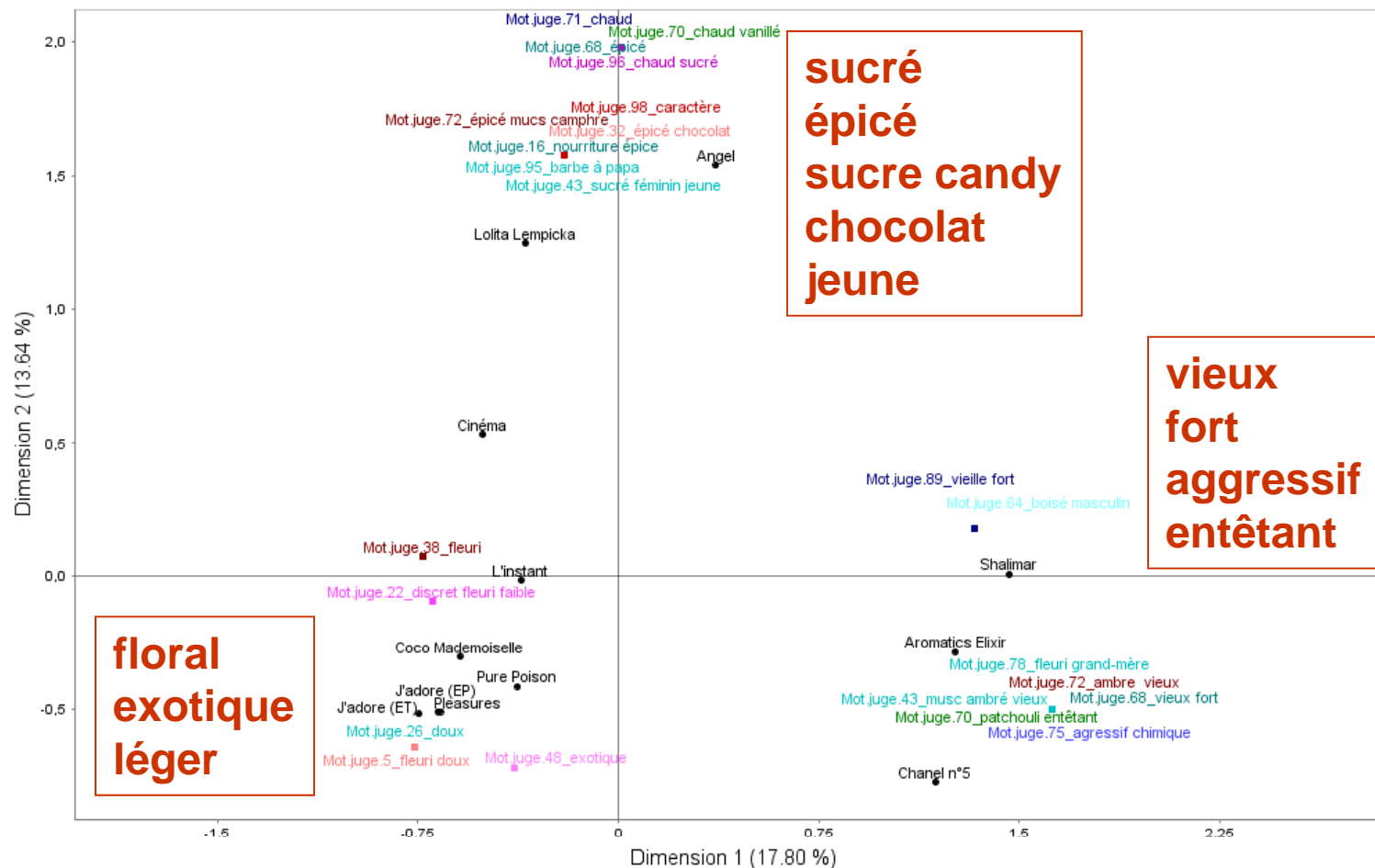


Figura 9.3.1.2. Resumen grupos en función de su nivel de redención.

Clustering

**SENSIBILIDAD
AL PRECIO ALTA**

**SENSIBILIDAD
AL PRECIO
MEDIA**

**SENSIBILIDAD
AL PRECIO BAJA**

SEGMENTO 2 (27.02%)

- Preferencia por productos baratos.
- Segmento con mayor sensibilidad al cupón.
- Gasto medio-bajo.
- Número de visitas medio.
- Gasto por visita muy reducido (14,1€).
- Es el segmento de mayor edad destacando el grupo de mayores de 60 años.
- Segmento con mayor preferencia por productos de marca blanca.

SEGMENTO 3 (44.80%)

- Compran por igual productos baratos, caros y productos de precio medio.
- Gasto mensual alto. Los que más gastan junto con el grupo 4.
- Número de visitas medio-alto.
- Nivel de redención alto
- El grupo con mayor porcentaje de clientes que viven con menores de edad.

SEGMENTO 1 (18.34%)

- Preferencia de productos de gama media.
- Gasto medio-alto.
- Número de visitas medio
- Nivel de sensibilidad al cupón medio.
- Es el segundo segmento de mayor edad después del segmento 2. Predominio del grupo de mayores de 60 años.

SEGMENTO 4 (9.85%)

- Preferencia por productos caros.
- Es el segmento menos sensible al cupón.
- Los que más gastan junto con el segmento 3.
- Son los que más gastan por visita.
- Gasto mensual alto.
- Número de visitas medio.
- Es el grupo que menos productos de marca blanca adquiere.
- Predominio del grupo de 45-60 años.

**Typology of the
customers**

MCA + clustering