

# Session

## Correspondence analysis

**Anàlisi de Dades i Explotació de la Informació**

**Grau d'Enginyeria Informàtica.**

*Information System tracking*

**Prof. Mónica Bécue Bertaut & Lidia Montero**

[Monica.becue@upc.edu](mailto:Monica.becue@upc.edu) [lidia.montero@upc.edu](mailto:lidia.montero@upc.edu)

Ronald Aylmer **Fisher**,  
1890 –1962



Brigitte Escofier (1941-1994)



Chikio  
Hayashi,  
(1918 - 2002)



Jean Paul Benzécri (1932)

1. Data and notation
2. Relationship between two categorical variables
3. CA: description of the deviation to independence
4. Geometrical approach: cloud of rows and cloud of columns.  
Superposition of both graphics of rows and columns
5. Helps to interpretation
6. Transition relationships: barycentric properties
7. Complements: supplementary elements; intensity of the relationship

## 1. Data

| Indiv    | $V_1$ | $V_2$ |
|----------|-------|-------|
|          |       |       |
| 1        |       |       |
| $\vdots$ |       |       |
| $l$      | $i$   | $j$   |
| $\vdots$ |       |       |
| $n$      |       |       |

Two categorical variables observed  
on the same individuals

| Indiv |       |       |
|-------|-------|-------|
|       | $V_1$ | $V_2$ |
| 1     |       |       |
| $l$   | $i$   | $j$   |
| $n$   |       |       |

Exemple  
Health survey in Croatia

*Edad en clase (7 categorías)*  
and  
*Estado de salud (5 categorías)*

on  $n=5037$  individuals

```
> summary(base$Edad_classe)
18-25 años      26-35 años      36-45 años      46-55 años      56-65 años      66-75 años      76 y más
      639          833          766          794          798          818          389

> summary(base$B1)
health-excellent health-very good      health-good      health-fair      health-poor
          472          833          1367          1322          1043
```

| Indiv    |       |       |
|----------|-------|-------|
|          | $V_1$ | $V_2$ |
| 1        |       |       |
| $\vdots$ |       |       |
| $i$      |       |       |
| $\vdots$ |       |       |
| $l$      |       |       |
| $\vdots$ |       |       |
| $n$      |       |       |

Contingency table

|          | 1 ..... $j$ | $J$ |
|----------|-------------|-----|
| 1        |             |     |
| $\vdots$ |             |     |
| $i$      |             |     |
| $\vdots$ |             |     |
| $I$      |             |     |

$x_{ij}$  = count of individuals who present at the same time category  $i$  of  $V_1$  and category  $j$  of  $V_2$

## Crossed table and margins

|               | health-excellent | health-very good | health-good | health-fair | health-poor |      |
|---------------|------------------|------------------|-------------|-------------|-------------|------|
| 18-25 años    | 181              | 216              | 161         | 69          | 12          | 639  |
| 26-35 años    | 144              | 263              | 259         | 129         | 38          | 833  |
| 36-45 años    | 62               | 150              | 266         | 201         | 87          | 766  |
| 46-55 años    | 35               | 105              | 260         | 239         | 155         | 794  |
| 56-65 años    | 26               | 43               | 190         | 281         | 258         | 798  |
| 66-75 años    | 17               | 38               | 166         | 283         | 314         | 818  |
| 76 y más años | 7                | 18               | 65          | 120         | 179         | 389  |
|               | 472              | 833              | 1367        | 1322        | 1043        | 5037 |



# From the count table to the proportion table

Table **F**

|        |         |          |     |          |
|--------|---------|----------|-----|----------|
|        | 1 ..... | $j$      | $J$ | margin   |
| 1      |         |          |     |          |
| 2      |         |          |     |          |
| ...    |         |          |     |          |
| $i$    |         | $f_{ij}$ |     | $f_{i.}$ |
| $I$    |         |          |     |          |
| margin |         | $f_{.j}$ |     | 1        |

$$f_{ij} = \frac{x_{ij}}{n}$$

$$f_{i.} = \sum_j f_{ij}$$

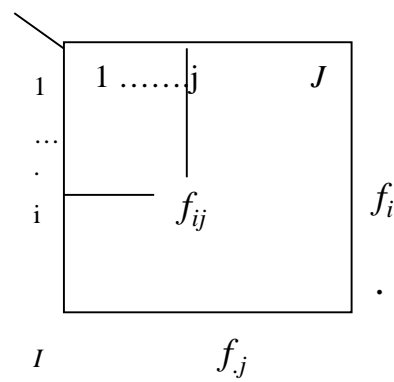
$$f_{.j} = \sum_i f_{ij}$$

Relationship between  $V_1$  and  $V_2$ : desviation between data and independence model

## Proportion table and margins

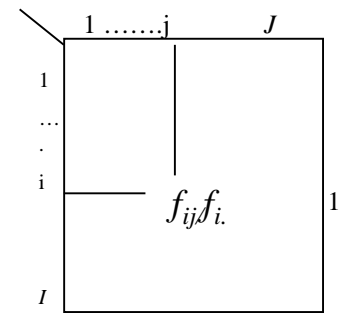
|               | health-excellent | health-very good | health-good  | health-fair  | health-poor  |              |
|---------------|------------------|------------------|--------------|--------------|--------------|--------------|
| 18-25 años    | 0.036            | 0.043            | 0.032        | 0.014        | 0.002        | <b>0.127</b> |
| 26-35 años    | 0.029            | 0.052            | 0.051        | 0.026        | 0.008        | <b>0.166</b> |
| 36-45 años    | 0.012            | 0.030            | 0.053        | 0.040        | 0.017        | <b>0.152</b> |
| 46-55 años    | 0.007            | 0.021            | 0.052        | 0.047        | 0.031        | <b>0.158</b> |
| 56-65 años    | 0.005            | 0.009            | 0.038        | 0.056        | 0.051        | <b>0.159</b> |
| 66-75 años    | 0.003            | 0.008            | 0.033        | 0.056        | 0.062        | <b>0.162</b> |
| 76 y más años | 0.001            | 0.004            | 0.013        | 0.024        | 0.036        | <b>0.078</b> |
|               | <b>0.093</b>     | <b>0.167</b>     | <b>0.272</b> | <b>0.263</b> | <b>0.207</b> | <b>1.000</b> |

## 2. Relationship/independence between two categorical variables



It there were independence

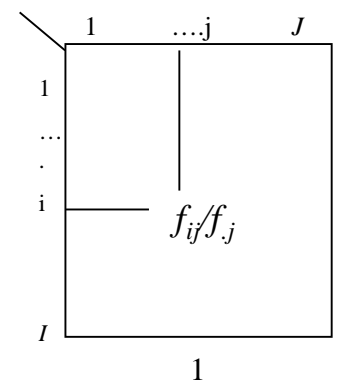
$$f_{ij} = f_{i.} \cdot f_{.j}$$



Row-profiles table

Table  $\mathbf{D_I^{-1}F}$

$$\frac{f_{ij}}{f_{i.}} = f_{.j}$$

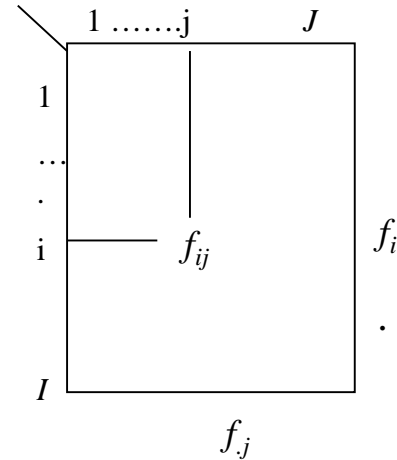
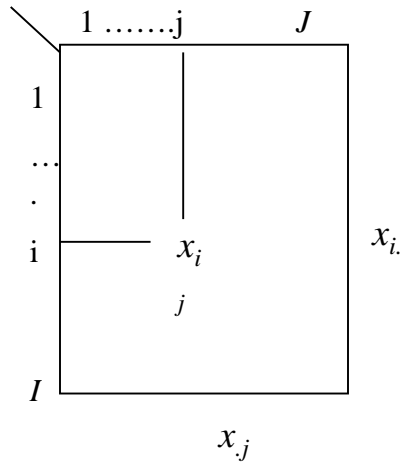


Column-profiles table

Tabla  $\mathbf{FD_J^{-1}}$

$$\frac{f_{ij}}{f_{.j}} = f_{i.}$$

## Observed data



## Estimation of the independence model

$$\hat{f}_{ij} = f_{i.} \cdot f_{.j}$$

Expected counted under independence hypothesis

$$\hat{x}_{ij} = n \cdot f_{i.} \cdot f_{.j}$$

Significance of the relationship between the two variables:  $\chi^2$  test

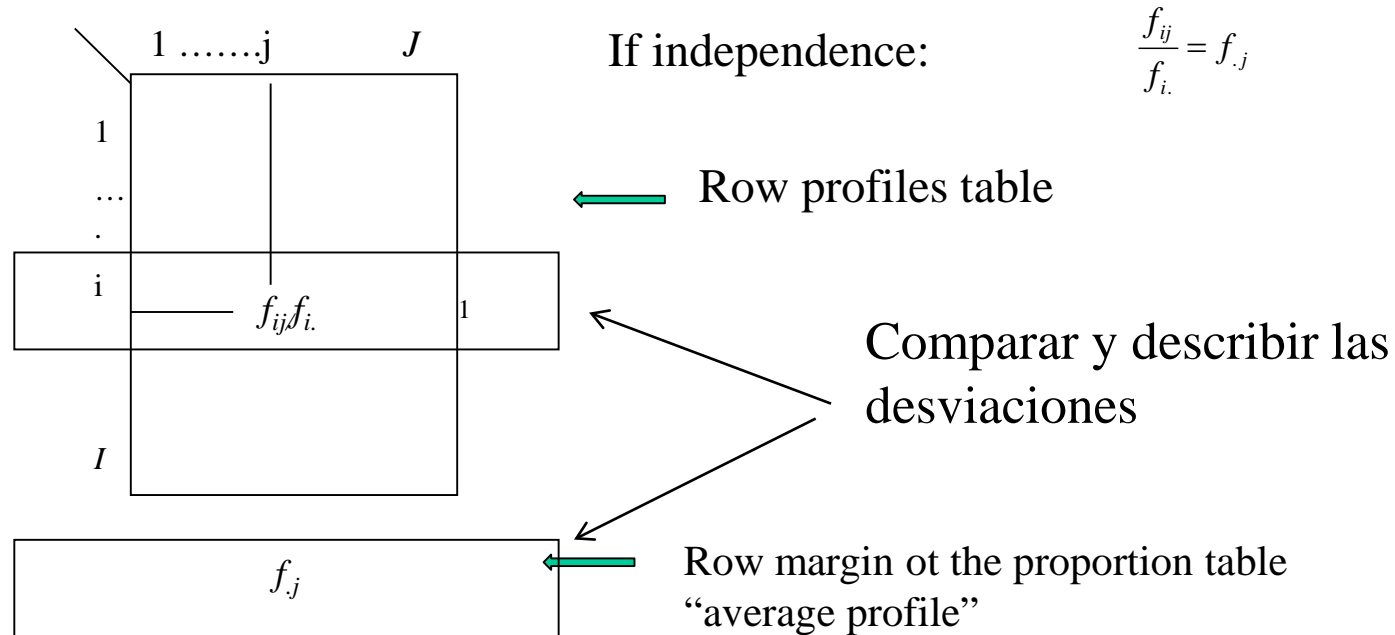
$$\chi^2 = \sum_{i,j} \frac{(x_{ij} - \hat{x}_{ij})^2}{\hat{x}_{ij}} \quad \dots \text{distribution...}p\text{-value}$$

Intensity of the relationship between the two variables

$$\Phi^2 = \sum_{i,j} \frac{(f_{ij} - \hat{f}_{ij})^2}{\hat{f}_{ij}} = \frac{\chi^2}{n}$$

CA simultaneously analyzes the row and column profiles tables. It does not inform about the significance of the relationship but about the intensity of the relationship and visualizes the structure of the relationship.

# Comparing profiles

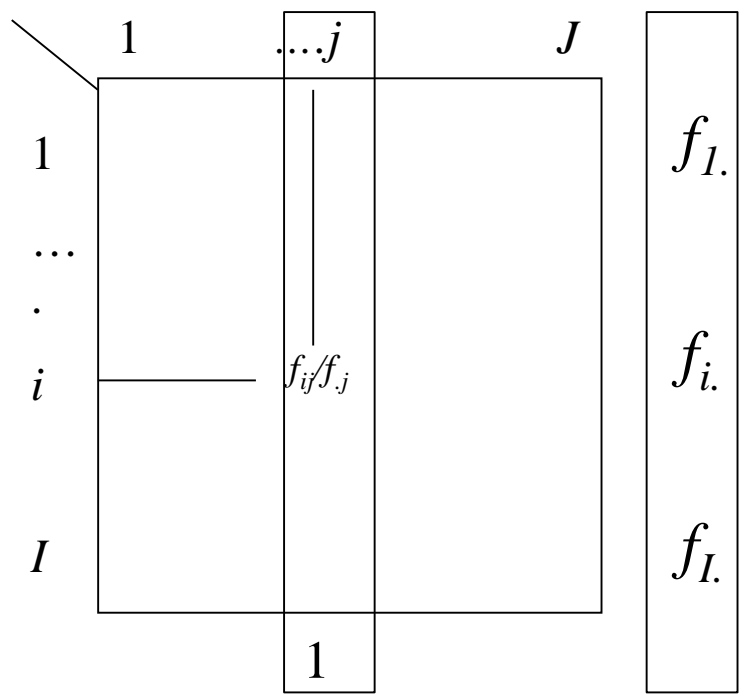


|                     | health-excellent | health-very good | health-good  | health-fair  | health-poor  |
|---------------------|------------------|------------------|--------------|--------------|--------------|
| 18-25 años          | 0.283            | 0.338            | 0.252        | 0.108        | 0.019        |
| 26-35 años          | 0.173            | 0.316            | 0.311        | 0.155        | 0.046        |
| 36-45 años          | 0.081            | 0.196            | 0.347        | 0.262        | 0.114        |
| 46-55 años          | 0.044            | 0.132            | 0.327        | 0.301        | 0.195        |
| 56-65 años          | 0.033            | 0.054            | 0.238        | 0.352        | 0.323        |
| 66-75 años          | 0.021            | 0.046            | 0.203        | 0.346        | 0.384        |
| 76 y más años       | 0.018            | 0.046            | 0.167        | 0.308        | 0.460        |
| <b>Perfil-medio</b> | <b>0.093</b>     | <b>0.167</b>     | <b>0.272</b> | <b>0.263</b> | <b>0.207</b> |

Column profiles table

If there were independence  $\frac{f_{ij}}{f_{i.}} = f_{.j}$

Column-margin of the proportion table  
 "Average column profile"



To compare and describe the deviations



> profil.col

|               | health-excell | health-very good | health-good | health-fair | health-poor |              |
|---------------|---------------|------------------|-------------|-------------|-------------|--------------|
| 18-25 años    | 0.383         | 0.259            | 0.118       | 0.052       | 0.012       | <b>0.127</b> |
| 26-35 años    | 0.305         | 0.316            | 0.189       | 0.098       | 0.036       | <b>0.166</b> |
| 36-45 años    | 0.131         | 0.180            | 0.195       | 0.152       | 0.083       | <b>0.152</b> |
| 46-55 años    | 0.074         | 0.126            | 0.190       | 0.181       | 0.149       | <b>0.158</b> |
| 56-65 años    | 0.055         | 0.052            | 0.139       | 0.213       | 0.247       | <b>0.159</b> |
| 66-75 años    | 0.036         | 0.046            | 0.121       | 0.214       | 0.301       | <b>0.162</b> |
| 76 y más años | 0.015         | 0.022            | 0.048       | 0.091       | 0.172       | <b>0.078</b> |

# 3. CA

CA= First, analysis of the cloud of rows

Cloud of rows described by their profile  $\frac{f_{ij}}{f_{i.}}$  Matrix  $\mathbf{D_I^{-1}F}$

Weights of rows

$f_{i.}$  stored into diagonal matrix

$\mathbf{D_I}$

chi.2 metric

with generic term  $\frac{1}{f_{.j}}$

$\mathbf{D_J^{-1}}$

$$d^2(i,l) = \sum_{j=1}^J \frac{1}{f_{.j}} \left( \frac{f_{ij}}{f_{i.}} - \frac{f_{lj}}{f_{l.}} \right)^2$$

→ distributional equivalence

CA= Then, analysis of the cloud of columns

Cloud of columns described by their profile  $\frac{f_{ij}}{f_{.j}}$  Matrix  $\mathbf{D}_J^{-1}\mathbf{F}'$

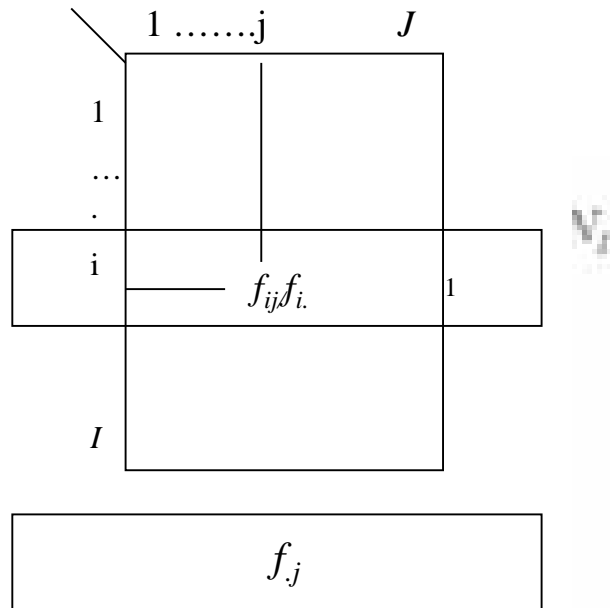
Weighted of the columns  $f_{.j}$  stored into diagonal matrix  $\mathbf{D}_J$

Métrica del chi.2 with generic term  $\frac{1}{f_{.i}}$   $\mathbf{D}_I^{-1}$

$$d^2(j, h) = \sum_{i=1}^I \frac{1}{f_{.i}} \left( \frac{f_{ij}}{f_{.j}} - \frac{f_{ih}}{f_{.h}} \right)^2$$

→ distributional equivalence

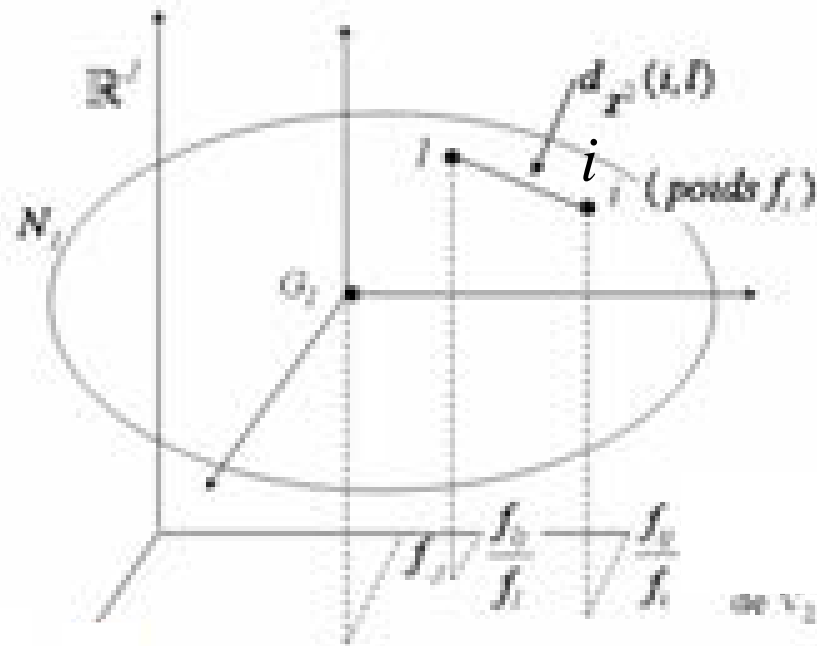
## 4. Geometrical approach: cloud of rows and cloud of columns



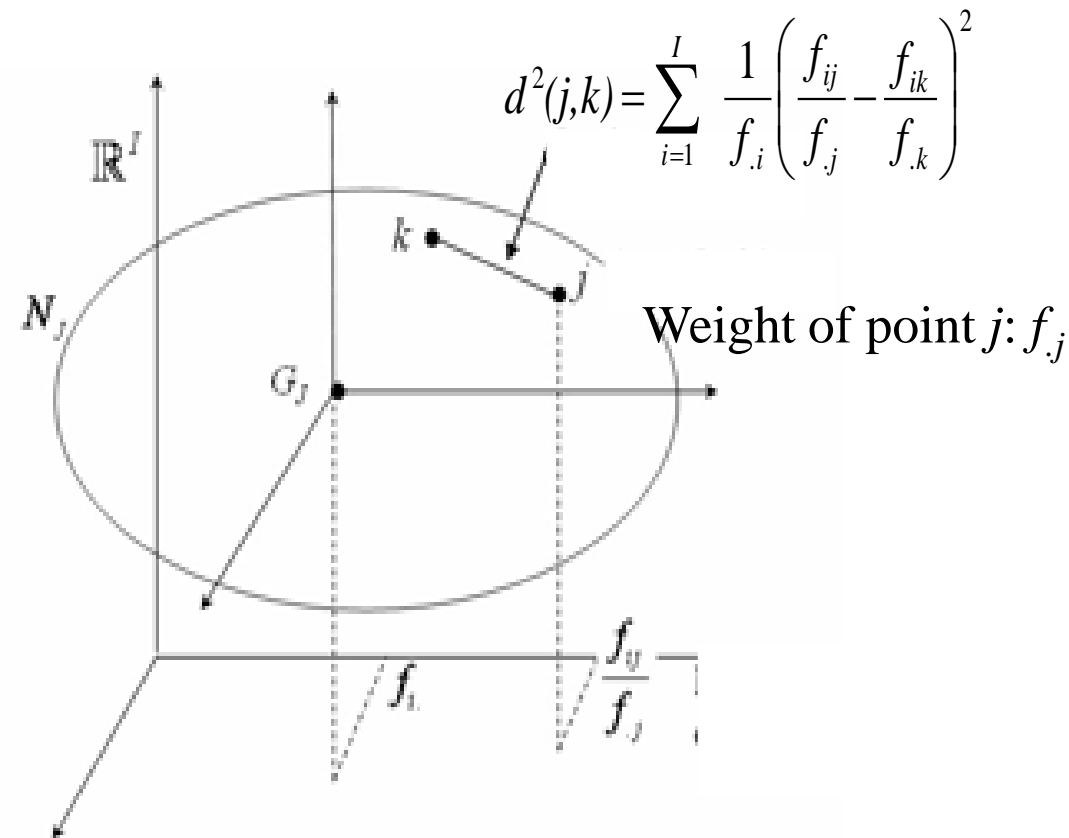
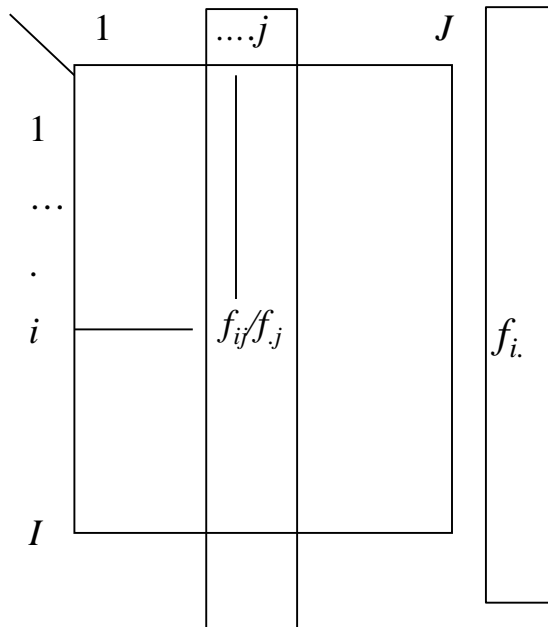
Weight of point  $i$ :  $f_i$ .

Cloud of rows

$$d^2(i, l) = \sum_{j=1}^J \frac{1}{f_j} \left( \frac{f_{ij}}{f_i} - \frac{f_{lj}}{f_l} \right)^2$$



## Cloud of column profiles



If there were independent

$$\frac{f_{ij}}{f_{i.}} = f_{.j}$$

$$\frac{f_{ij}}{f_{i.}} = f_{.j}$$

Both clouds have a null inertia

$$Inercia(N_I|G_I) = Inercia(N_J|G_J)$$

The intensity of the relationship is higher insofar as the inertia is higher

$$\begin{aligned} Inercia(N_I|G_I) &= \sum_i Inercia(i|G_I) = \sum_i f_{i.} d^2(i, G_I) = \sum_j f_{.j} d^2(j, G_J) = \\ &= \sum_i \sum_j \frac{1}{f_{i.} f_{.j}} (f_{ij} - f_{i.} \cdot f_{.j})^2 = \\ &= \Phi^2 = \frac{\chi^2}{n} = Inercia(N_J|G_J) \end{aligned}$$

Representation is a low dimension space

Find the subspace which better sums up the data

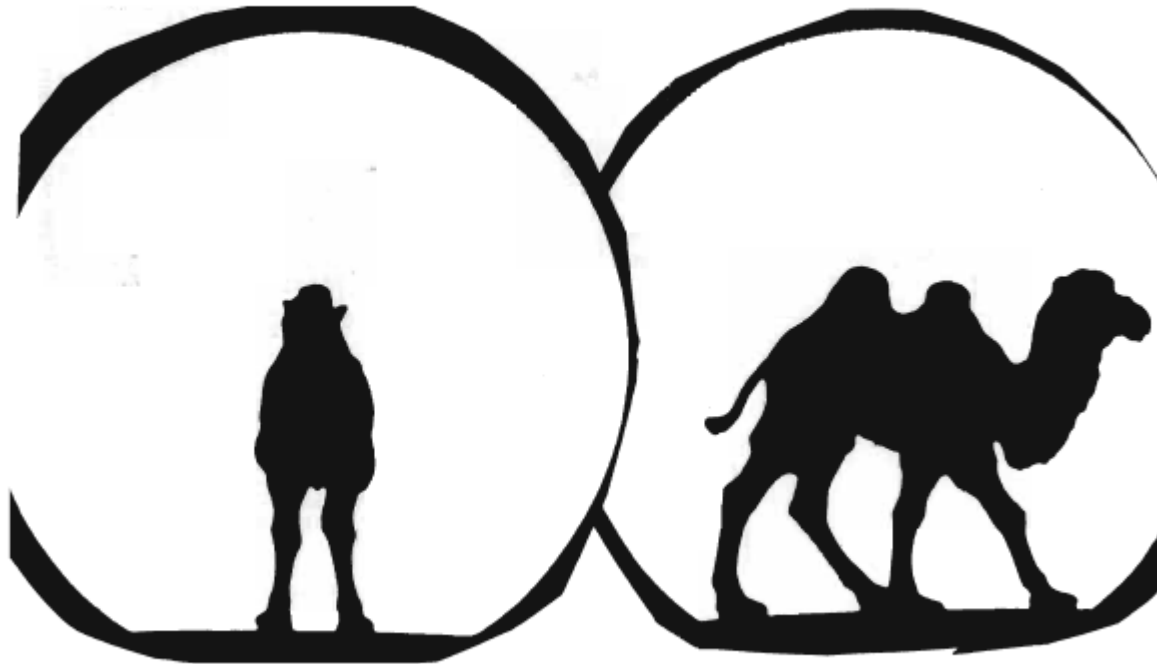
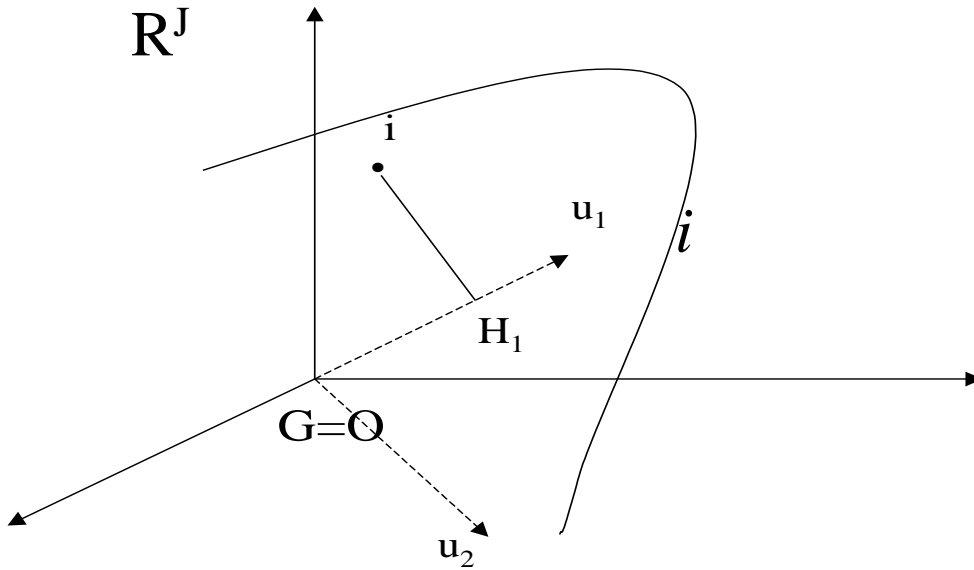


Figure: Camel vs dromedary?



## Row cloud



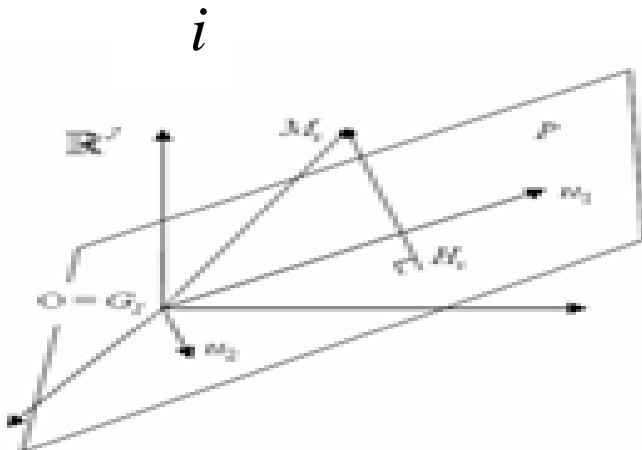
$$\text{Max} \sum_i f_i \cdot O H_i^2$$

$$\begin{array}{ll} u_1 & \lambda_1 \\ u_2 & \lambda_2 \\ u_3 & \lambda_3 \end{array}$$

... ..

$$u_{\min(I-1, J-1)} \lambda_{\min(I-1, J-1)}$$

-



In the column cloud.....

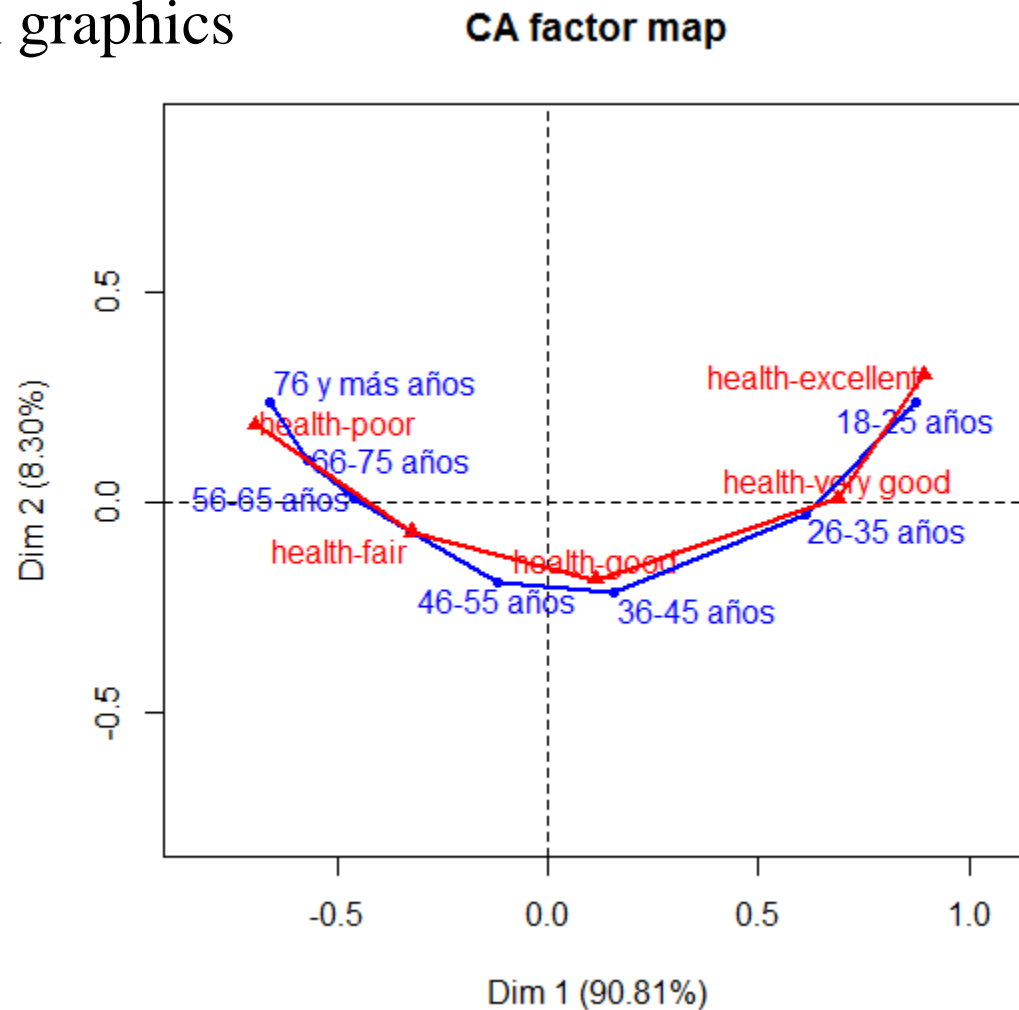
$$\begin{array}{cc} v_1 & \lambda_1 \\ v_2 & \lambda_2 \\ v_3 & \lambda_3 \end{array}$$

$$v_{\min(I-1,J-1)} \lambda_{\min(I-1,J-1)}$$

-

$$\Phi^2 = \sum_i \lambda_i = \sum_i f_{i.} d^2(i, G_I) = \sum_j f_{.j} d^2(j, G_J)$$

Graphical results: in this case, it is legitimate to superpose the row and column graphics



Guttman effect

## 5. Helps to interpretation

### Global quality of the representation

$$\frac{\text{Conserved inertia}}{\text{Total inertia}}$$

On the first plane:

$$\frac{\lambda_1 + \lambda_2}{\sum_{s=1}^S \lambda_s}$$

```
> round(res.ca$eig,2)
```

|       | eigenvalue | percentage of variance | cumulative percentage of variance |
|-------|------------|------------------------|-----------------------------------|
| dim 1 | 0.29       | 90.81                  | 90.81                             |
| dim 2 | 0.03       | 8.30                   | 99.11                             |
| dim 3 | 0.00       | 0.83                   | 99.94                             |
| dim 4 | 0.00       | 0.06                   | 100.00                            |
| dim 5 | 0.00       | 0.00                   | 100.00                            |

```
> FI2
[1] 0.3142015=chi2/n
```

V de Cramer

```
> sqrt(sum(res.ca$eig[,1])/4)
[1] 0.2802684
```

Particularities of the eigenvalues in CA  $0 \leq \lambda_s \leq 1$

What does it mean to observe an eigenvalue equal to 1?

Maximum number of axes?

How many axes we have to interpret?

## Contributions and quality of representation

= what we have seen in PCA

But do not forget that the elements (rows and columns) are endowed with weights

## 6. Transition relationships also called barycentric relationships

$$F_s(i) = \frac{1}{\sqrt{\lambda_s}} \sum_j \frac{f_{ij}}{f_{i.}} G_s(j)$$

$$G_s(j) = \frac{1}{\sqrt{\lambda_s}} \sum_i \frac{f_{ij}}{f_{.j}} \cdot F_s(i)$$



## 6. Complements

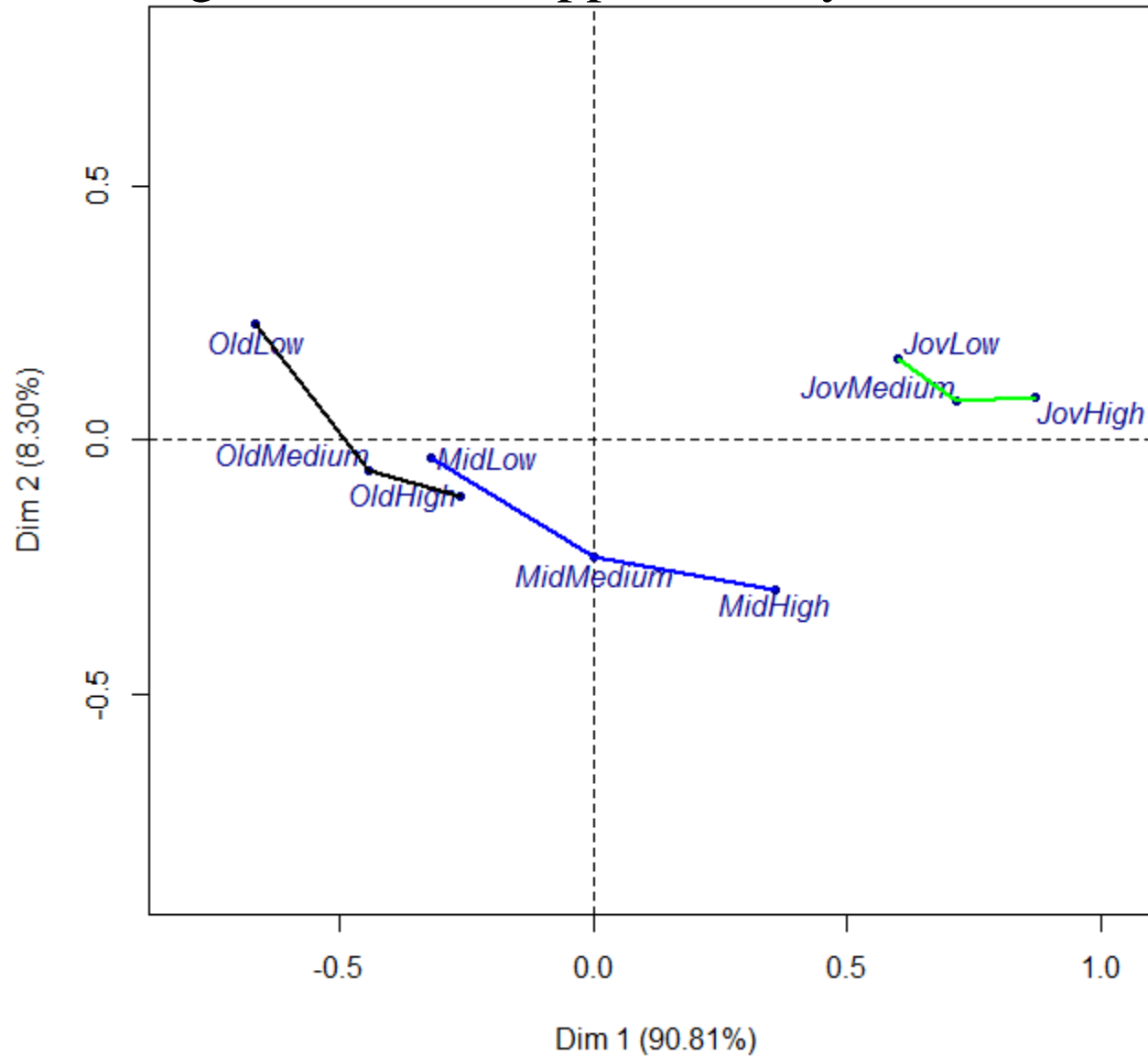
### 6.1 Supplementary columns and/or rows

$$F_s(i^+) = \frac{1}{\sqrt{\lambda_s}} \sum_j \frac{f_{i^+j}}{f_{i^+ \cdot}} G_s(j)$$

$$G_s(j^+) = \frac{1}{\sqrt{\lambda_s}} \sum_i \frac{f_{ij^+}}{f_{\cdot j^+}} \cdot F_s(i)$$

CA factor map

Categories of Age  $\times$  Educ as supplementary rows



## 6.2 Nature of the relation. Intensity of the relation. Cramer V

The graph informs about the nature of the relationship between the variables through the visualisation of the associations between the categories of one variable and these of the other

The eigenvalues– and their sum– inform about the intensity of the relationship.

The Cramer V allows for comparing the intensity of the relationship with its maximum (and therefore, between cross-tables with different dimensions)

$$V = \sqrt{\frac{\phi^2}{\text{Max}(\phi^2)}} = \sqrt{\frac{\phi^2}{\text{Min}(I-1, J-1)}}$$