

Session 3: Profiling

Anàlisi de Dades i Explotació de la Informació

Grau d'Enginyeria Informàtica.
Information System track

Prof. Josep Franquet

josep.franquet@upc.edu

Information is stored in tables

Rows of table

Represent individuals or instances, Also called sample, example, record, ... they can be repeated, at least theoretically, forming the population under study.

Thing to be classified, associated, or clustered

Characterized by a predetermined set of attributes

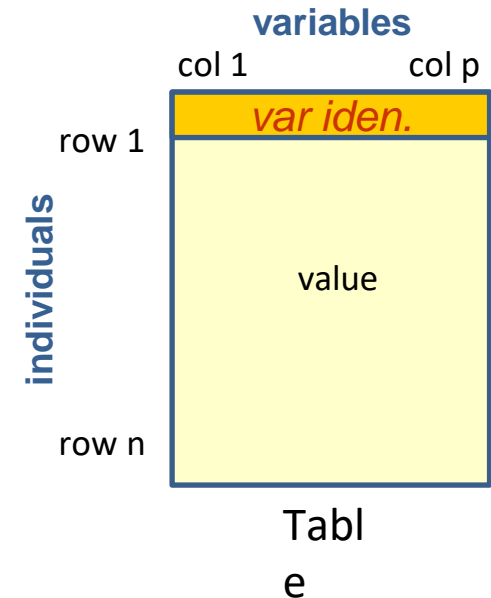
Columns of a table

Each instance is described by a predefined set of features, its variables or “attributes”. A variable is a measure of individuals which can take different values (according a probabilistic function).

Possible attribute types (“levels of measurement”): binary, nominal, ordinal, interval, ratio, textual, ...

Restriction: Same variables measured in all individuals and in the same order, but different formats are possible (fixed, csv, ...), forming a Table.

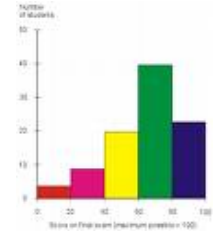
First rows usually contain the dictionary of variables (var. labels)



Coding of variables

- **Categorical**

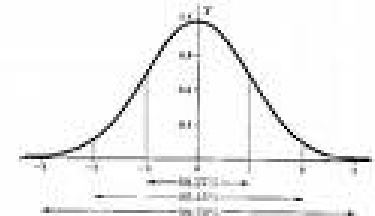
- Binary (yes/no variable, boolean attribute. Binomial distribution)
- Nominal (marital status, region, Multinomial distribution)
- Ordinal (size of clothe's size, social class,, Ordered multinomial)
- Pareto chart



Prob. distribution

- **Continuous or numeric**

- Count data (“number of words of a sentence”, “number of unemployed per country”, Poisson distribution)
- Interval (“temperature,... Laplace-Gauss distribution)
- Ratio data (age, speed, ... Laplace Gauss distribution)



But distinctions are often blurred. Ordinal data can be treated as continuous,

- **Multiresponse**

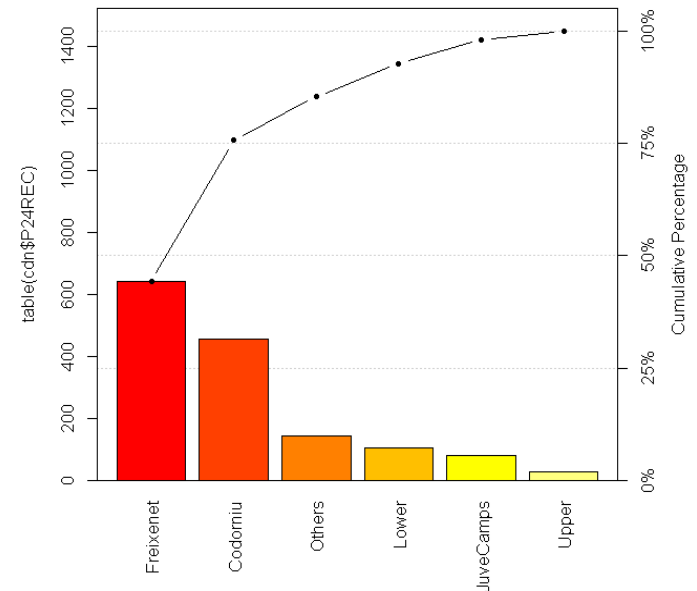
When several responses are possible.

i.e. Visited countries?. Each potential visited country becomes a binary variable

- **Textual**

Pareto chart: a classical Quality Control chart

Pareto Chart for table(cdn\$P24REC)



```
> library(qcc)
```

```
Loading required package: MASS
```

```
Package 'qcc', version 2.2
```

```
> pareto.chart(table(cdn$P24REC), ylab = "table(cdn$P24REC)")
```

Pareto chart analysis for table(cdn\$P24REC)

	Frequency	Cum.Freq.	Percentage	Cum.Percent.
Freixenet	640.000000	640.000000	44.137931	44.137931
Codorniu	457.000000	1097.000000	31.517241	75.655172
Others	142.000000	1239.000000	9.793103	85.448276
Lower	105.000000	1344.000000	7.241379	92.689655
JuveCamps	79.000000	1423.000000	5.448276	98.137931
Upper	27.000000	1450.000000	1.862069	100.000000

Role of variables

- **Response**

Variables that we want to study, by building a model, finding associations, ...
(number of products bought, passing or failing a course, income, ...)

It can be either continuous or categorical

- **Explanatory**

Variables which serve to explain the behaviour of the response variables (all the variables present in the data matrix except the response)

They can be either continuous or categorical

Types of data matrix

With or without response(s) variable

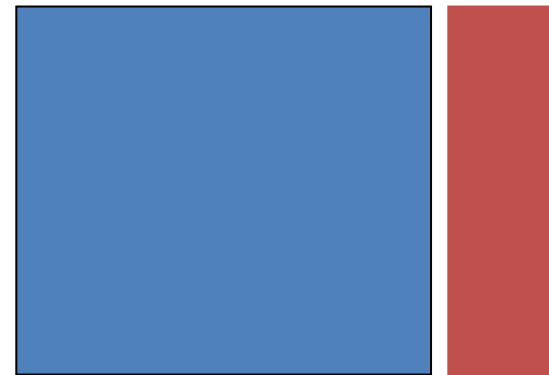
i.e. transactions data



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

Inputs

Output(s)



Idem, but we want to **find a model to predict the response**

Paradigm

Any stored data from any process always contain information about the generating phenomenon (**statistical regularity**).

Goal: **To reveal the information** (model, patterns, associations, trends, clusters, ... hidden in the data)

Data are routinely stored (and most will never be analyzed)

Data is a treasure for organizations (be aware of the data quality)

Any transactional process can be enhanced by analysis of its collected data

How? *Selecting and reporting what is interesting*

SQL queries are NOT ENOUGH. How many A products sold last month?

Profiling. What is the profile of A buyers? *Automatic detection of significant deviations*

Automatic profiling of groups of individuals

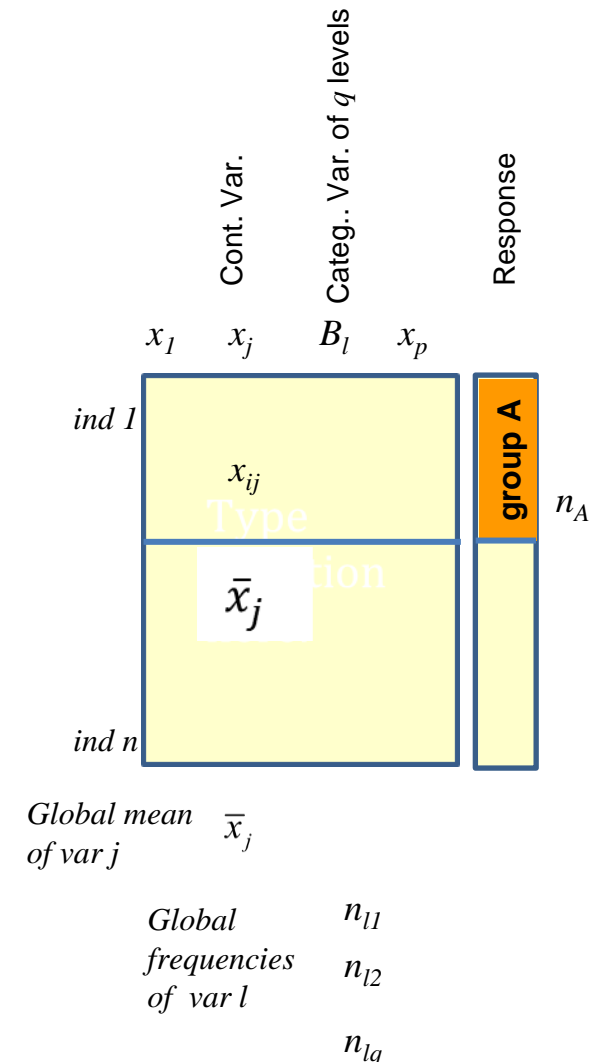
We have a group of individuals defined **by a level of a categorical variable (target).**

Problem: For every group of individuals detect which other groups of individuals (identified by the levels of the explanatory variables) or what continues variables, deviate significantly from what were expected.

- We take as response variable the variable identifying the groups that we want to find their profile.
- The explanatory variables are either categorical or continuous.

Tool: Hypothesis test

- For each group to profile, rank the modalities of the categorical explanatory variables according their p-value (ascending). Likewise, rank the continuous variables according their p-value
- Select the most significant by a threshold (0.05, 0.01, ..) defined a priori. (what matters is the ordering, actual significance depends on the number of individuals)



R function available in FactoMineR

We will use FactoMineR Package (cran R)

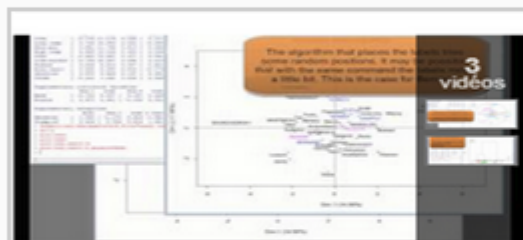
You can also consult (and download this R function from) <http://factominer.free.fr/> where a large documentation is provided, with theoretical background, examples, tutorials and so on.

The functions of this package corresponding to these sessions are:

- **Catdes:** description of the categories of a categorical variable by quantitative variables, categorical variables and categories
- **Condes:** description of a quantitative variable by quantitative and categorical variables

FACTOMINER[®]

> News bulletin



Exploratory multivariate analysis with R and FactoMineR

Videos on the use of FactoMineR (for PCA, multiple factor analysis, clustering, etc.)

The version 1.24 of FactoMineR has a new graphical module that place the labels in an "optimal" way, that allows to select some elements to draw, etc.

Four reviews on the book [Exploratory Multivariate Analysis by Example using R](#) are available in this site. To see the complete review done by Gary Evans ([for Journal of Statistical Software](#))

A new [user group](#) to ask questions on FactoMineR and on Exploratory Multivariate Data Analysis has been created. Join this group to have news about FactoMineR and to ask questions

[missMDA](#): a new package to handle missing values in PCA, MCA or MFA with FactoMineR

[English Version](#)[Version française](#)

> Top Menu

[Home](#)[Classical Methods](#)[Advanced Methods](#)[Interface](#)[Facto's best](#)[FactoMineR and Excel](#)[F.A.Q.](#)[Documents](#)[Contact](#)

> Useful Links

[Agrocampus Rennes Applied Maths Department](#)[R Project](#)[CRAN](#)

Elements of a Hypothesis test

Tool to validate a hypothesis

Hypothesis: Group k of individuals is different from the population

To test a hypothesis we need:

Null hypothesis

H_0 : denial of the hypothesis (like an evil advocate). Group k is equal to the population

Alternative hypothesis

H_1 : the hypothesis we want to validate

Test statistic:

it depends on the problem.

Reference distribution:

Distribution of the *test statistic* if the H_0 is true.

Significance threshold: Risk that we are ready to incur to reject H_0 when it is true (significance depend on the number of individuals, thus classical statistical thresholds need to be adapted).

In our case:

H_0 : Individuals of group k are taken at random

H_1 : Individuals of group k are NOT taken at random

General rationale of any test

Test

We want to see if problematic posed by the user

Null hypothesis

H_0 : conservative hypothesis.

Alternative hypothesis

H_1 : hypothesis in which the user is interested in

Data

data are observed, generally a sample

Test statistic:

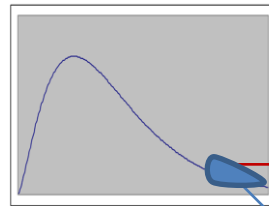
a statistic is a **function of the sample (=observed data) and thus a variate or random variable**

Reference distribution:

Distribution of the *test statistic* under H_0 (that is, if H_0 is true).

Significance threshold:

Risk of rejecting H_0 although H_0 being true



Observed value, as computed on the data

$\left[\begin{array}{ll} \leq 0.05 & H_0 \text{ is rejected} \\ > 0.05 & H_0 \text{ is not rejected} \end{array} \right.$

This area= p -value

Elements of a Hypothesis Test

The test statistic depends on the level of measurement of variables

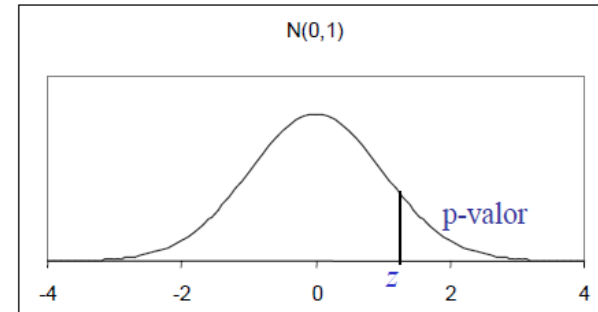
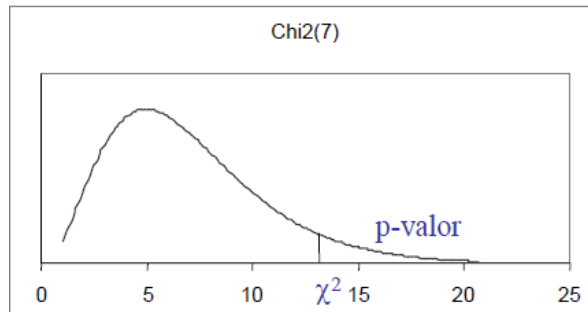
Null hypothesis

- For continuous variables, the *mean* in group should be equal to the global mean
- For categorical variables: the *proportion* in group k should be equal to the global proportion

Compute per each test, the *p-value* of the test statistic

p-value: probability to find a *test statistic* equal or farthest that the actual test statistic.

It is the probability that if the H_0 is true, of actually finding the observed test statistic



Finally rank the *p-values* in ascending order, and select the significant characteristics from those *p-values* lower than the chosen significance threshold

Response target: factor (B)
Explanatory variable: numerical (X)

$$B \sim X$$

Consider background for $X \sim B$

Profiling a categorical target from a continuous variable

$$H_0: \mu_1 = \dots = \mu_p = \mu$$

<i>groups</i>	<i>means</i>	<i>counts</i>
1	\bar{x}_1	n_1
\boxtimes	\boxtimes	\boxtimes
p	\bar{x}_p	n_p

Null Hypothesis: All group means are equal to the global mean

Global \bar{x} n

In R:

- Assuming normal distribution on X:
`oneway.test(X~B)`.
- Without normality assumption (non –parametric test): Kruskal-Wallis test
`kruskal.test(X~B)`
- Global association: Tested using a F-Fisher based-test



Ronald Fisher 1890, 1962

Profiling target categorical variables from continuous variables

$$H_0 : \mu_k = \mu \quad k = 1, \dots, p$$

Test statistic: Difference between the mean in group k and the global mean. T-Student based-test

Highlight groups with a significant different mean: level specific association tests

groups	means	counts
1	\bar{x}_1	n_1
\boxtimes	\boxtimes	\boxtimes
p	\bar{x}_p	n_p

Global \bar{x} n



William Gosset "Student",
English, 1876, 1937

$$t = \frac{\bar{x}_k - \bar{x}}{\sqrt{(1 - \frac{n_k}{n}) \frac{s^2}{n_k}}} \boxtimes t_{n-1}$$

Student's t

*Rank the continuous variables by p.value
(ascending)*

Function to compute p-values for profiling a categorical target from continuous variables – Globally and Specific Level

To Rank variables and groups according to pvalues:

```
p.xk <- function(vec,fac) {  
  nk <- as.vector(table(fac));  
  n <- sum(nk);  
  xk <- tapply(vec,fac,mean);  
  txk <- (xk-mean(vec))/(sd(vec)*sqrt((n-nk)/(n*nk)));  
  pxk <- pt(txk,n-1,lower.tail=F)}
```

Rank the continuous variables by p.value (ascending)

FactoMineR solution:

- `catdes(data.frame,num.var):` sections
 - Link between the cluster variable and the quantitative variables
 - Description of each cluster by quantitative variables

Response target: factor (B)
Explanatory variable: factor (A)

$$B \sim A$$

Profiling categories from categorical variables

– **Global Relationship between each category** of the target variable and **other categorical variables: a chi-square-test is performed**

	1... j ... q	
1	\boxtimes	
k	\boxtimes n_{kj} \boxtimes	n_k
p	\boxtimes	
	n_j	

– **Relationship between each category** of the variable target and **each category of another categorical variable: comparison of two proportions**, taking into account an hypergeometric model and normal approximations

– Descriptive tools: contingency tables (numeric) and mosaic plot (graphical)

Rank the levels of the categorical explanatory variables/ the categories by p-value (ascending)

(Global) Characterization of a target categorical variable by the other categorical variables

Test

Null hypothesis

Alternative hypothesis

H_0 : conservative hypothesis. Both variables are independent

H_1 : Both variables are not independent

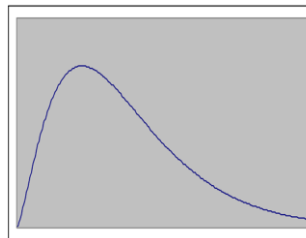
Test statistic:

$$\chi_{obs}^2 = \sum_i \sum_j \frac{\left(n_{ij} - \frac{n_{i.} n_{.j}}{n}\right)^2}{\frac{n_{i.} n_{.j}}{n}} = \sum_i \sum_j \frac{\left(n_{ij} - np_{i.} p_{.j}\right)^2}{np_{i.} p_{.j}}$$

Reference distribution:

freedom

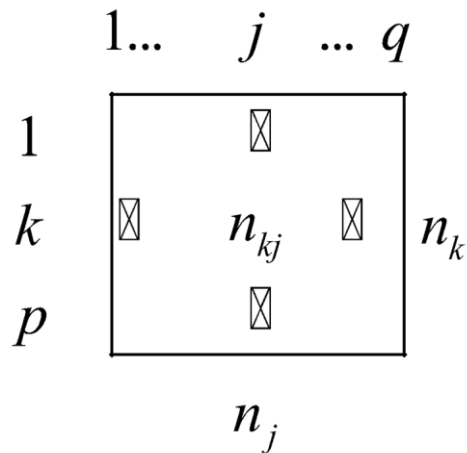
Distribution of the *test statistic* under H_0 (that is, if H_0 is true).
Chi-2 distribution, with the convenient degrees of



Significance threshold:

Risk of rejecting H_0 although H_0 being true
(significance depends on the number of individuals) **P-value**

Characterization of a target categorical variable by the levels of other categorical variables



$$H_0 : p_{j/k} = p_j \quad k = 1, \dots, p; j = 1, \dots, q$$

Assumption of normality of proportions:

$$\frac{n_{kj}}{n_k} \sim N\left(p_j = \frac{n_j}{n}, \left(1 - \frac{n_k}{n}\right) \frac{p_j(1 - p_j)}{n_k}\right)$$

Test statistic: Difference between proportion of modality j in group k and proportion of modality j in whole data

$$Z = \frac{\frac{n_{kj}}{n_k} - \frac{n_j}{n}}{\sqrt{\left(1 - \frac{n_k}{n}\right) \left(\frac{p_j(1 - p_j)}{n_k}\right)}} \sim N(0, 1)$$

Rank the levels of the categorical explanatory variables by p.value (ascending)

R function to compute the table of p-values of a categorical variable

```
p.zkj <- function(res,expl){  
  taula <- table(res,expl)  
  n <- sum(taula);  
  pk <- apply(taula,1,sum)/n;  
  pj <- apply(taula,2,sum)/n;  
  pf <- taula/(n*pk);  
  pjm <- matrix(data=pj,nrow=nrow(pf),ncol=ncol(pf), byrow=T);  
  dpf <- pf - pjm;  
  dvt <- sqrt(((1-pk)/(n*pk))%*%t(pj*(1-pj)));  
  zkj <- dpf/dvt;  
  pzkj <- pnorm(zkj,lower.tail=F);  
  list(rowpf=pf,vtest=zkj,pval=pzkj)}
```

FactoMineR solution:

- `catdes(data.frame,num.var)`
 - Link between the cluster variable and the categorical variables (chi-square test)
 - Description of each cluster by categories

Example: SwissLabor data in AER library

Usage

```
data("SwissLabor")
```

Format

A data frame containing 872 observations on 7 variables.

```
levels(SwissLabor$participation)<-  
  paste("Parti.",sep="",levels(SwissLabor$participation))  
levels(SwissLabor$foreign)<-  
  paste("Foreign.",sep="",levels(SwissLabor$foreign))
```

participation	Factor. Did the individual participate in the labor force?
income	Logarithm of nonlabor income.
age	Age in decades (years divided by 10).
education	Years of formal education.
youngkids	Number of young children (under 7 years of age).
oldkids	Number of older children (over 7 years of age).
foreign	Factor. Is the individual a foreigner (i.e., not Swiss)?

Profiling a categorical target by the categories of the other categorical variables

In SwissLabor dataset in library(AER): Participation – Yes (target) vs Foreign

H₀: The category “foreign=NO” is neither infra nor supra represented

H₁: The category “foreign=NO” is infra (versus supra) represented

```
> table(SwissLabor$foreign,SwissLabor$participation)
```

	Target-no	Target-yes
Foreign-no	402	254
Foreign-yes	69	147

```
➤ prop.table(table(SwissLabor$foreign,SwissLabor$participation),1)
```

	Target-no	Target-yes
Foreign-no	0.6128049	0.3871951
Foreign-yes	0.3194444	0.6805556

```
> prop.table(table(SwissLabor$foreign,SwissLabor$participation),2)
```

	Target-no	Target-yes
Foreign-no	0.8535032	0.6334165
Foreign-yes	0.1464968	0.3665835

```
round(prop.table
(table(SwissLabor$foreign)),dig=2)
Foreign.no Foreign.yes
0.75 0.25
```


Profiling a categorical target by other categorical variables

In SwissLabor dataset in library(AER): Participation – Yes (target) vs Foreign

H₀: Global Association is not present among Target and Explanatory Factor

H₁: Global Association is present

```
> table(df$foreign,df$participation)
```

	Target-no	Target-yes
Foreign-no	402	254
Foreign-yes	69	147

Under H0

	Target-no	Target-yes
f.For-no	354.3303	301.66972
f.For-yes	116.6697	99.33028

```
> chisq.test(table(df$foreign,df$participation))
```

Pearson's Chi-squared test with Yates' continuity correction

data: table(df\$foreign, df\$participation)

X-squared = 55.126, df = 1, p-value = 1.131e-13

```
> res.cat$test.chi2 # Global association target is factor and explanatory factors – catdes()
```

```
      p.value df
foreign 6.220116e-14 1
```

Characterization of a categorical variable by the categories of the other categorical variables

In SwissLabor dataset in library(AER): Participation – Yes (target) vs Foreign

H₀: The category “foreign=NO” is neither infra nor supra represented

H₁: The category “foreign=NO” is infra (versus supra) represented

68% of class Foreign-Yes
belongs to Category Target-Yes
 $P([B\text{-}TargetYes]/[A\text{-}Foreign\text{-}Yes])$

36.7% of Category Target-Yes belongs to class
Foreign-Yes : $P([A\text{-}Foreign\text{-}Yes] / [B\text{-}TargetYes])$

`$category$ Target=yes`` →

	Cla/Mod	Mod/Cla	Global	p.value	v.test
foreign=Foreign-yes	68.05556	36.65835	24.77064	5.591005e-14	7.517321
foreign=Foreign-no	38.71951	63.34165	75.22936	5.591005e-14	-7.517321

`prop.table(table(SwissLabor$foreign))` Foreign-no Foreign-yes
0.7522936 0.2477064

Foreign women represent 25% of the sample, but
36.7% in the target class Target-Yes

Characterization of a categorical variable by a quantitative variable

Characterization of the categorical variable “*target*” by the quantitative variables

$$Y_{ki} = \mu + \alpha_k + \varepsilon_{ki}$$

target (level k) = grand mean + effect for k level + error

H_0 (no category effect): $\alpha_1 = \dots = \alpha_k = \dots = \alpha_K = 0$

H_1 : There are at least two “target” levels k and k' such as: $\alpha_k \neq \alpha_{k'}$

```
> res.cat$quanti.var # Global association target is factor  
and explanatory variables numeric
```

	Eta2	P-value
youngkids	0.029968826	2.695567e-07
income	0.029891180	2.794460e-07
education	0.010516854	2.429641e-03
age	0.008521288	6.375401e-03
oldkids	0.006445786	1.772877e-02

Profiling categories from quantitative variables

Characterization of the categories of “*target*” by the quantitative variables

H_0 mean of the variable in the category = grand mean

H_1 : mean in the category \neq global mean

```
> res.cat$quanti # Especific association: target factor and numeric variables
$f.Par-no`
```

	v.test	Mean in category	Overall mean	sd in category	Overall sd	p.value
youngkids	5.109095	0.4097665	0.3119266	0.6770660	0.6125185	3.237063e-07
income	5.102472	10.7513327	10.6855675	0.4351131	0.4122522	3.352458e-07
education	3.026579	9.5944798	9.3073394	2.8484531	3.0345172	2.473382e-03
age	2.724342	4.0853503	3.9955275	1.1599921	1.0545623	6.442967e-03
oldkids	-2.369447	0.9023355	0.9827982	1.0622927	1.0861630	1.781471e-02

```
$`f.Par-yes`
```

	v.test	Mean in category	Overall mean	sd in category	Overall sd	p.value
oldkids	2.369447	1.0773067	0.9827982	1.1060968	1.0861630	1.781471e-02
age	-2.724342	3.8900249	3.9955275	0.9040227	1.0545623	6.442967e-03
education	-3.026579	8.9700748	9.3073394	3.2067731	3.0345172	2.473382e-03
income	-5.102472	10.6083220	10.6855675	0.3689875	0.4122522	3.352458e-07
youngkids	-5.109095	0.1970075	0.3119266	0.5029499	0.6125185	3.237063e-07

Response target: Numeric (Y)
Explanatory variable: Numeric (X)
Explanatory variable: factor (A)

$$Y \sim X$$

$$Y \sim A$$

- Target type: numeric or factor
- Type of explanatory variates
 - Global association (target, explanatory variate)
 - Specific association (target, explanatory variate)
- FactoMineR:
 - Target is a factor: `catdes()`
 - Target is numeric: `condes()`

Profiling a quantitative target from quantitative or categorical variables

– Description by quantitative variables (condes) : global association

Correlation (Pearson)-> `cor(data.frame)`

– Description by categorical variables and categories

ANOVA: test F (global association) and t-tests (level specific associations)

Response target: Numeric (Y)
Explanatory variable: Numeric (X)

$$Y \sim X$$

Relationship between a quantitative target and the other quantitative variables

H_0) no relationship (correlation is null $\rho=0$)

H_1) relationship (correlation non null $\rho \neq 0$)

Statistics:

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

Pearson linear correlation

```
> condes(SwissLabor,2) #Numeric target income
$quanti
```

	correlation	p.value
education	0.3273458	3.166132e-23
oldkids	0.1391036	3.758541e-05

Response target: Numeric (Y)
Explanatory variable: Factor (A)

$$Y \sim A$$

Global Relationship between the quantitative target “income” and the categorical variables

income (foreign group k) = mean + effect of foreign group k + error

$$Y_{ki} = \mu + \alpha_k + \varepsilon_{ki}$$

H_0 (no category effect): $\alpha_1 = \dots = \alpha_k = \dots = \alpha_K = 0$

H_1 : There are at least two “factor” levels k and k' such as: $\alpha_k \neq \alpha_{k'}$

Global association Fisher F-Based

> condes(SwissLabor,2) #Numeric target income

...

\$quali

R2 p.value

foreign 0.04389655 4.170824e-10

participation 0.02989118 2.794460e-07

Relationship between the quantitative variable “income” and the levels of categorical variables

H_0 The coefficient of category k is null $\alpha_k = 0$

H_1 : The coefficient of category k is non-null $\alpha_k \neq 0$

Level specific association t-Student based tests – Only significant levels included in the output

> condes(SwissLabor,2) #Numeric target income

\$category

Estimate p.value

Foreign.no 0.10004281 4.170824e-10

Parti.no 0.07150532 2.794460e-07

Parti.yes -0.07150532 2.794460e-07

Foreign.yes -0.10004281 4.170824e-10

Example: SwissLabor data in AER library

Usage

```
data("SwissLabor")
```

Format

A data frame containing 872 observations on 7 variables.

participation	Factor. Did the individual participate in the labor force?
income	Logarithm of nonlabor income.
age	Age in decades (years divided by 10).
education	Years of formal education.
youngkids	Number of young children (under 7 years of age).
oldkids	Number of older children (over 7 years of age).
foreign	Factor. Is the individual a foreigner (i.e., not Swiss)?

Example: SwissLabor data in AER library

```
> condes(SwissLabor,2) #Numeric target income
```

```
$quanti
```

```
correlation  p.value
```

```
education  0.3273458 3.166132e-23
```

```
oldkids    0.1391036 3.758541e-05
```

Global association

```
$quali
```

```
R2    p.value
```

```
foreign    0.04389655 4.170824e-10
```

```
participation 0.02989118 2.794460e-07
```

Global association

```
$category
```

```
Estimate  p.value
```

```
Foreign.no 0.10004281 4.170824e-10
```

```
Parti.no   0.07150532 2.794460e-07
```

```
Parti.yes  -0.07150532 2.794460e-07
```

```
Foreign.yes -0.10004281 4.170824e-10
```

Profiling on
categories