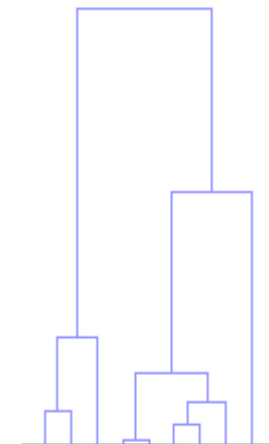


Correspondence & Hierarchical Cluster (CHIC) Analysis Software v1.0

Angelos Markos, George Menexes & Iannis Papadimitriou

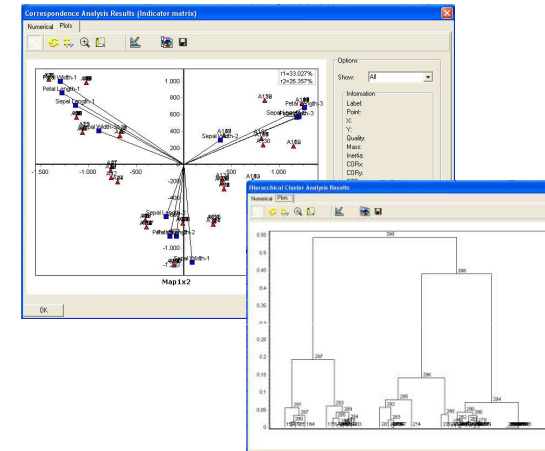
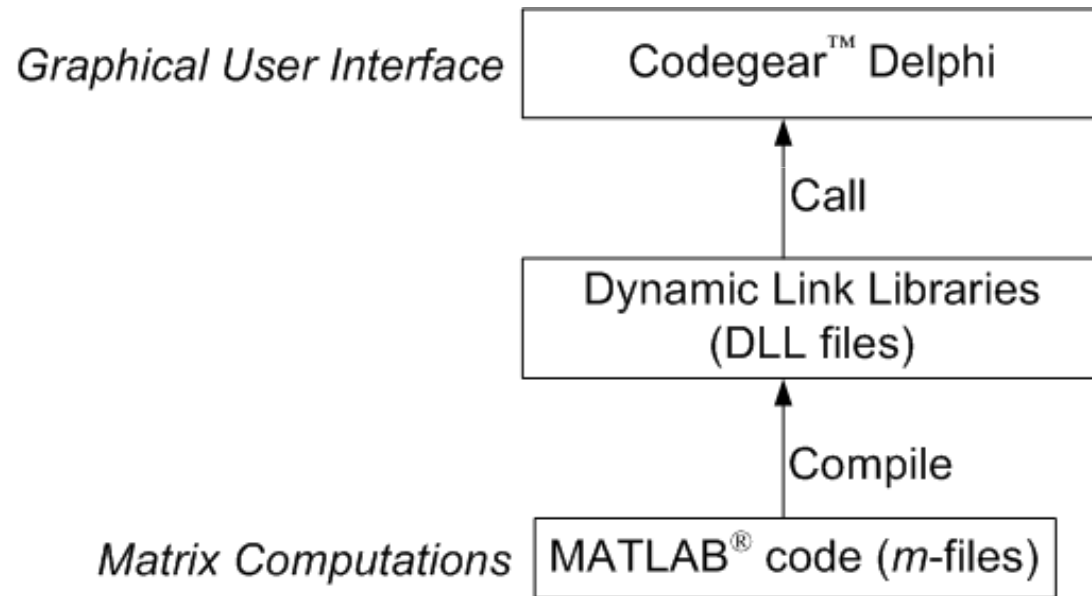
Department of Applied Informatics, University of Macedonia,
Thessaloniki, Greece



Outline

- What is CHIC Analysis?
- Motivation & Goals
- Data Entry & Data Management
- Interpretation Options
- The SVD Implementation of MCA
- Hierarchical Cluster Analysis as a complementary method
- Work in Progress

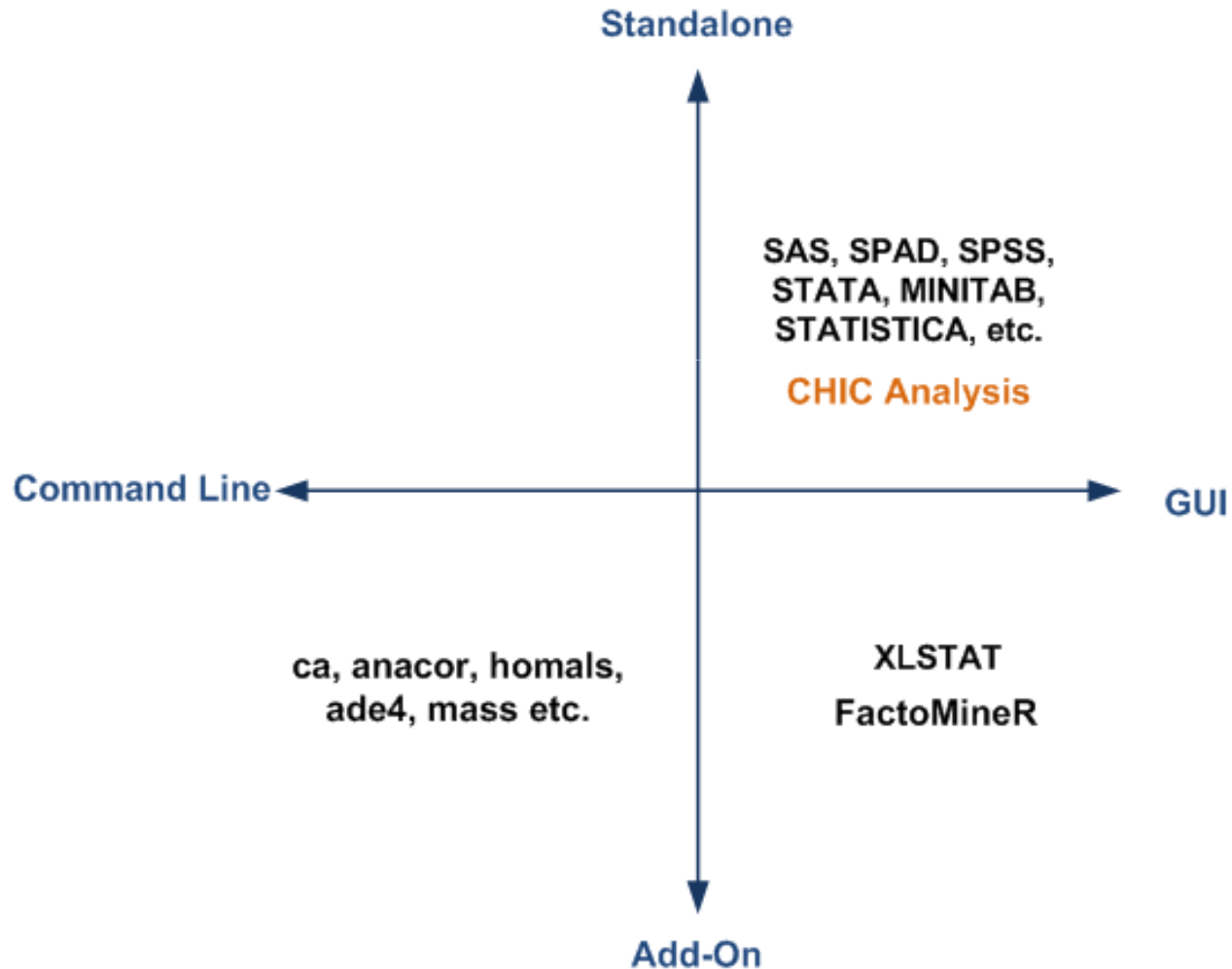
Development Tools



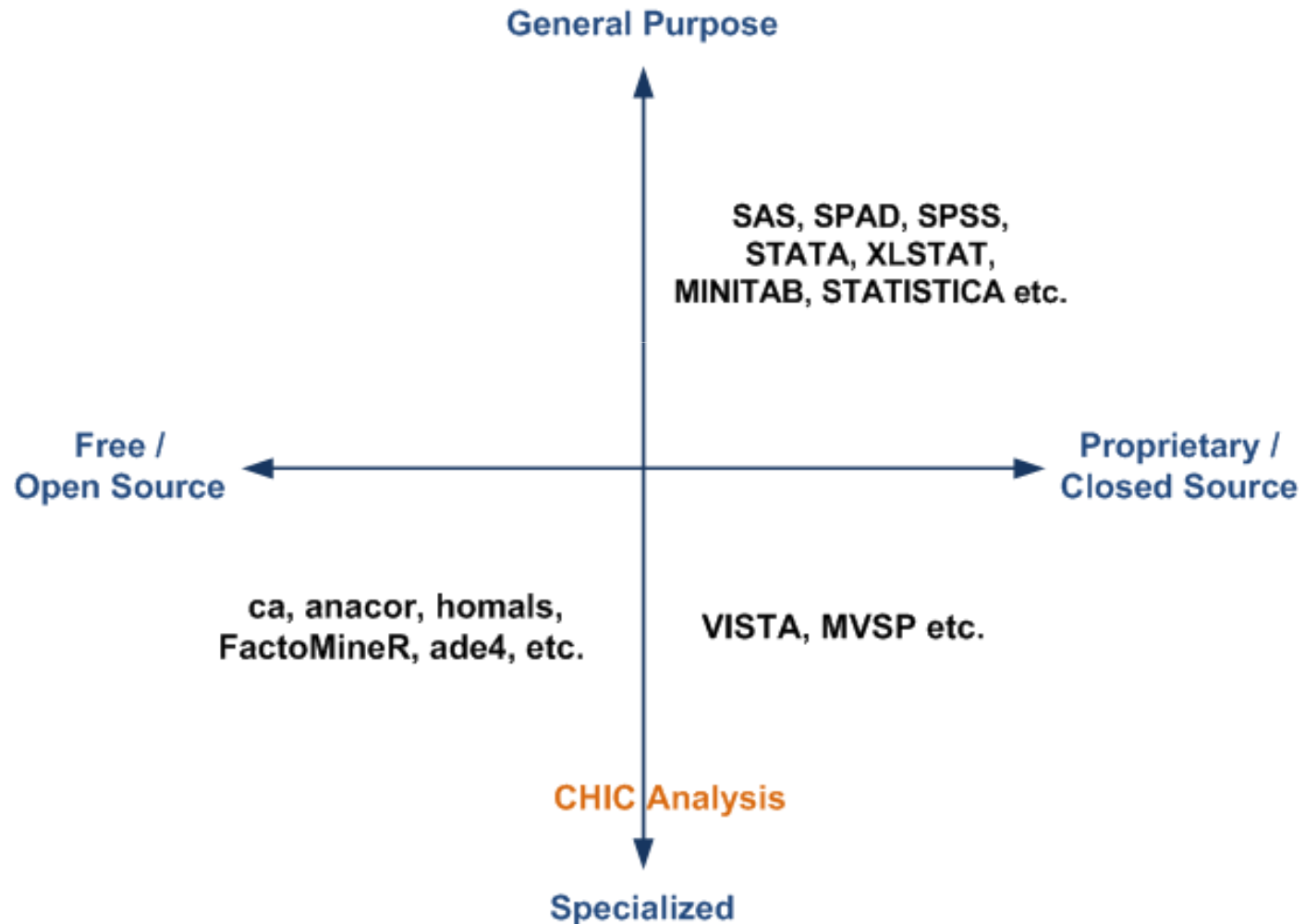
```

% Correspondence matrix
P = X./total;
% Standardized residuals matrix
%S = (P - ri'*cj)./(sqrt(ri'*cj));
S = stdresij/sqrt(total);
    
```

Corners of Statistical Software for CARME (1)



Corners of Statistical Software for CARME (2)



Motivation & Goals

- Educational purposes

- A CA software towards interpretation:
 - Options for both empirical interpretation and statistical inference
 - A synthesis of contributions (French School & the GIFI System)
 - Direct construction and analysis of Burt subtables
 - Options for Symmetric Plots and Biplots

Data Entry & Data Management (1)

■ The Data Spreadsheet

(a) Raw data table

C:\data sets\iris.ana - CHIC Analysis v1.2

C / V	Sepal Length	Sepal Width	Petal Length	Petal Width	Species
C1	5.1	3.5	1.4	0.2	1
C2	4.9	3	1.4	0.2	1
C3	4.7	3.2	1.3	0.2	1
C4	4.6	3.1	1.5	0.2	1
C5	5	3.6	1.4	0.2	1
C6	5.4	3.9	1.7	0.4	1
C7	4.6	3.4	1.4	0.3	1
C8	5	3.4	1.5	0.2	1
C9	4.4	2.9	1.4	0.2	1
C10	4.9	3.1	1.5	0.1	1
C11	5.4	3.7	1.5	0.2	1
C12	4.8	3.4	1.6	0.2	1
C13	4.8	3	1.4	0.1	1
C14	4.3	3	1.1	0.1	1
C15	5.8	4	1.2	0.2	1
C16	5.7	4.4	1.5	0.4	1
C17	5.4	3.9	1.3	0.4	1
C18	5.1	3.5	1.4	0.3	1
C19	5.7	3.8	1.7	0.3	1
C20	5.1	3.0	1.5	0.2	1

Row: 1 Column: 1 Cases:150 Variables:5 Input: Raw Data Table

(b) Contingency table

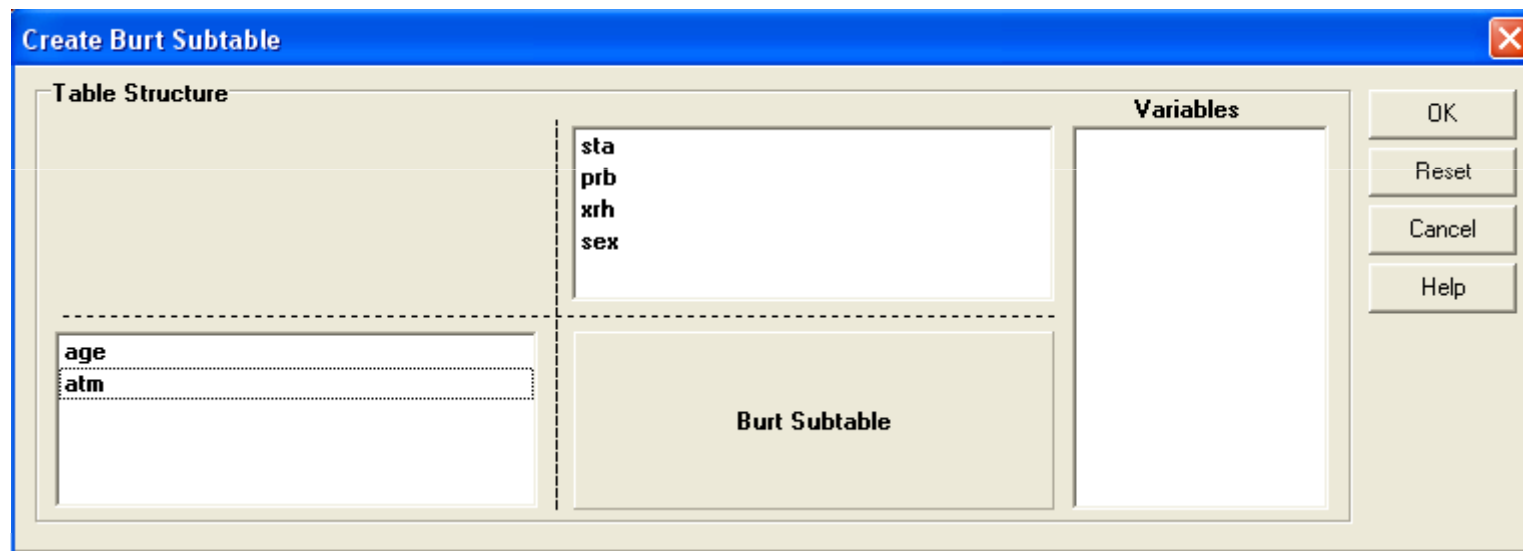
C:\data sets\science.ana - CHIC Analysis v1.2

R / C	A	B	C	D	E	Y
Geology	3	19	39	14	10	0
Biochemistry	1	2	13	1	12	1
Chemistry	6	25	49	21	29	0
Zoology	3	15	41	35	26	0
Physics	10	22	47	9	26	1
Engineering	3	11	25	15	34	1
Microbiology	1	6	14	5	11	1
Botany	0	12	34	17	23	1
Statistics	2	5	11	4	7	0
Mathematics	2	11	37	8	20	1
Museums	4	12	11	19	7	0
Math. Sciences	4	16	48	12	27	1

Row: 1 Column: 1 Rows: 12 Columns: 6 Input: Contingency Table

Data Entry & Data Management (2)

- Burt Subtable construction



Simple Correspondence Analysis Options (1)

- Significance criteria of individual points

 - Active points

 - Mass, Inertia, COR,CTR, QLT, SQCOR, Best

 - Supplementary points

 - Inertia, COR, QLT

Simple Correspondence Analysis Options (2)

- Selection of significant axes
 - Scree Plot
 - Statistical significance tests by
 - Nishisato
 - van de Geer
 - Greenacre

- Additional Options
 - Expected frequencies
 - Residuals, Standardized residuals, Adjusted standardized residuals
 - Chi-square contributions
 - Reconstruction matrix (chi-square goodness of fit)

Multiple Correspondence Analysis Options (1)

■ Significance criteria of individual points

Active

- Mass, Inertia, COR, SQCOR, CTR, SumCTR, QLT, Best, Discrimination Measures (absolute, relative)

Supplementary

- Inertia, COR, QLT, Test values

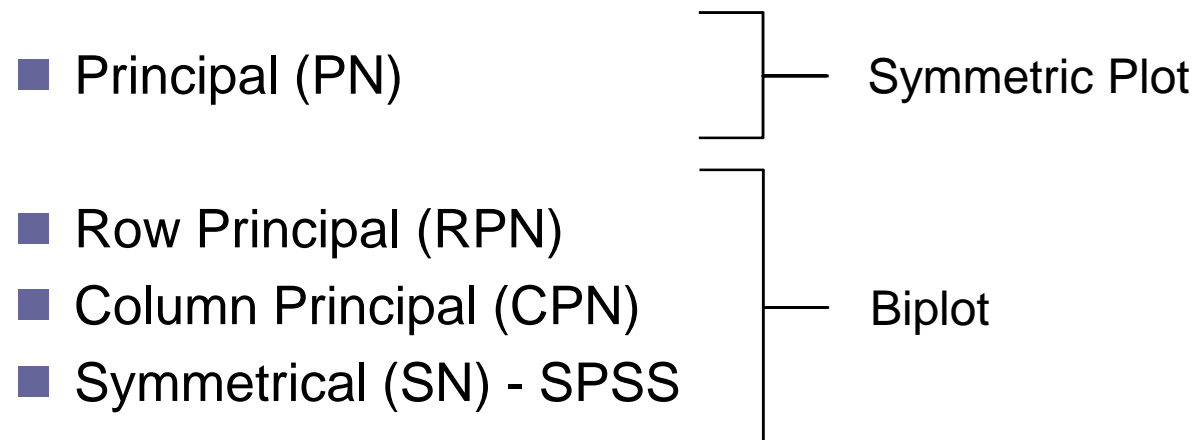
Multiple Correspondence Analysis Options (2)

- Selection of significant axes
 - Scree Plot
 - Principal inertias $> 1/p$
 - Cronbach's alpha
 - Statistical significance test by
 - Nishisato

- Inertia Adjustment
 - Greenacre
 - "Interesting" Inertia

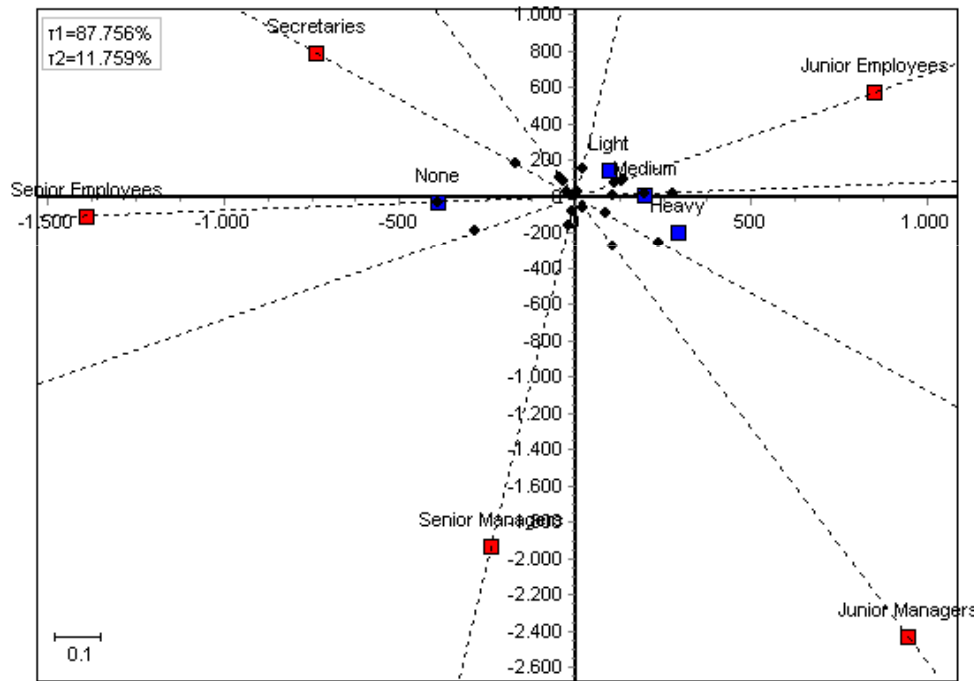
Symmetric Plots & Biplots (1)

■ Normalization Options



Symmetric Plots & Biplots (2)

■ Biplot Axes (Normalization: RPN, CPN, SN)



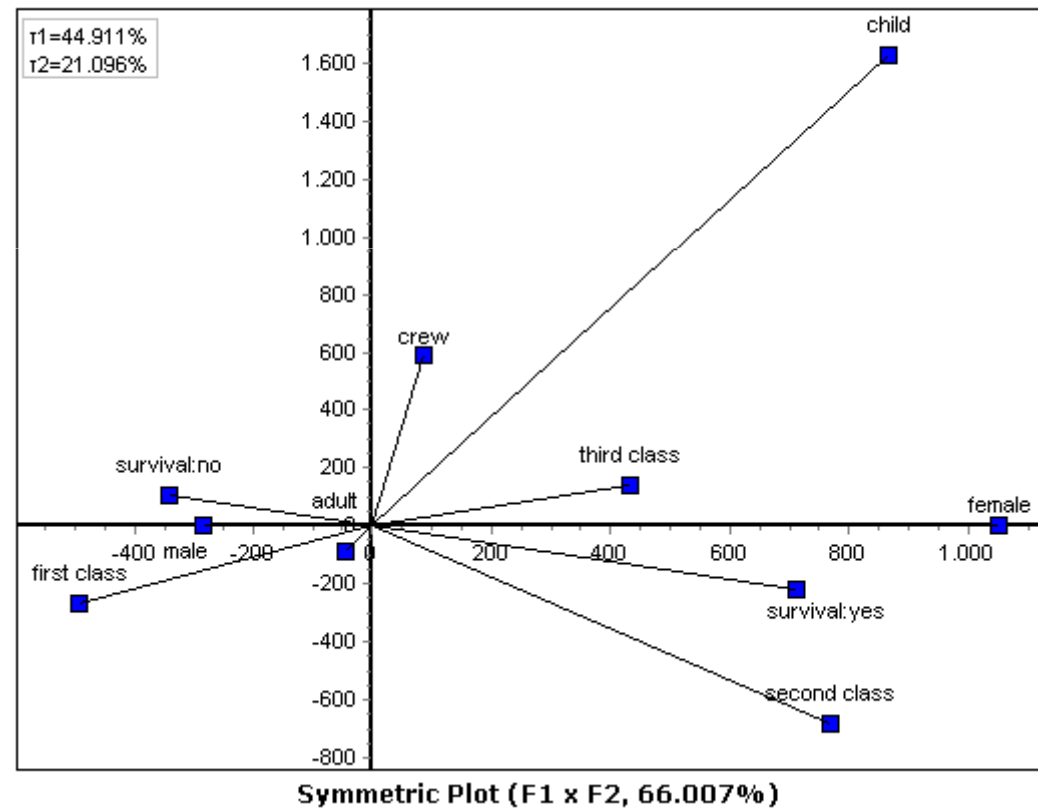
Biplot (F1 x F2, 99.515)

Table 1. Biplot Interpretation (CPN)

Senior Managers	Heavy	None	Medium	Light
Distance	1.790	1.873	1.982	2.103
Cos ²	0.140	0.039	0.998	0.307
Junior Managers	Heavy	None	Medium	Light
Distance	2.469	2.540	2.636	2.744
Cos ²	0.390	0.780	0.615	0.109
Senior Employees	Heavy	None	Medium	Light
Distance	1.373	1.383	1.403	1.434
Cos ²	0.043	0.025	0.754	0.675
Junior Employees	Light	Medium	None	Heavy
Distance	0.936	1.009	1.082	1.141
Cos ²	1.000	0.874	0.204	0.430
Secretaries	Light	Medium	None	Heavy
Distance	0.987	1.058	1.129	1.187
Cos ²	0.404	0.718	0.085	0.949

Symmetric Plots & Biplots (3)

- Position Vectors (Normalization: PN)



A note on the SVD Implementation of MCA

Step 1 - Calculate the matrix \mathbf{S}_B of standardized residuals

$$\mathbf{S}_B = \mathbf{diag}(\mathbf{r}_B)^{-1/2} (\mathbf{P}_B - \mathbf{r}\mathbf{c}^T) \mathbf{diag}(\mathbf{c}_B)^{-1/2}$$

Step 2 - Calculate the SVD of \mathbf{S}_B

$$\mathbf{S}_B = \mathbf{U}_B \mathbf{D}_B \mathbf{V}_B^T$$

Steps 3 & 4 - Standard coordinates Φ_Z of rows and Γ_Z of columns

$$\Phi_Z = \frac{1}{p} \mathbf{Z} \Gamma_Z \mathbf{D}_B^{-1/2}$$

$$\Gamma_Z = \mathbf{diag}(\mathbf{r}_Z)^{-1/2} \mathbf{U}_B = \Gamma_B$$

Steps 4 & 5 - Principal coordinates \mathbf{F}_Z of rows and \mathbf{G}_Z of columns

$$\mathbf{F}_Z = \frac{1}{p} \mathbf{Z} \Gamma_Z = \Phi_Z \mathbf{D}_B^{1/2}$$

$$\mathbf{G}_Z = \Gamma_Z \mathbf{D}_B^{1/2}$$

Hierarchical Cluster Analysis Algorithm

Input

Contingency Table, Indicator Matrix

Implementation

Ward's linkage criterion, Benzécri's chi-square distance

Output


- Dendrogram
 - Agglomeration schedule
 - Cluster Membership
 - Variable contributions to the characterization of clusters
 - Variable contributions to the division of clusters
- } VACOR

Work in Progress

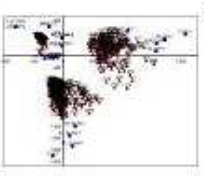
- Take advantage of the common mathematical foundation of many multivariate data analysis methods, as a basis for writing unified software
- Correspondence Analysis Variations (Canonical, Joint, Subset, Non-Symmetrical, Regularized)
- Aspects of Stability
- A forthcoming **R** version

<http://amarkos.gr/en/research/chic>

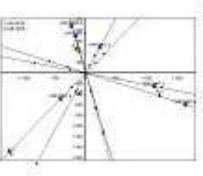
Ελληνικά ->



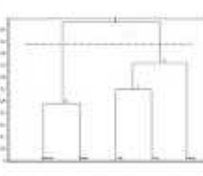
#1




#2



#3



#4



> About CHIC Analysis

CHIC (Correspondence & Hierarchical Cluster) Analysis is a standalone software designed for practitioners in the field of Correspondence Analysis and Related Methods. It features two well known exploratory data analysis methods to portray the overall structure of a data analysis session: **Correspondence Analysis-CA** (Simple and Multiple) and **Hierarchical Cluster Analysis** (Benzecri's chi-square distance, Ward's linkage criterion). The implementation of CA is in line with both the French School and the Gifi System of Data Analysis.

CHIC Analysis combines the graphical interface capabilities of Codegear Delphi with the computational power of Mathworks MatLab® (R version forthcoming). The software was implemented as an attempt to contribute to the effectiveness and reliability of Correspondence Analysis. It offers numerous aids to the results' interpretation and tools for the construction and analysis of special data tables. Special emphasis has been put on the graphical options for bi-plots, maps and dendrograms.

CHIC Analysis was designed to be reasonably comprehensive, fairly easy to use, accurate, and available for free.

> Features

> Menu

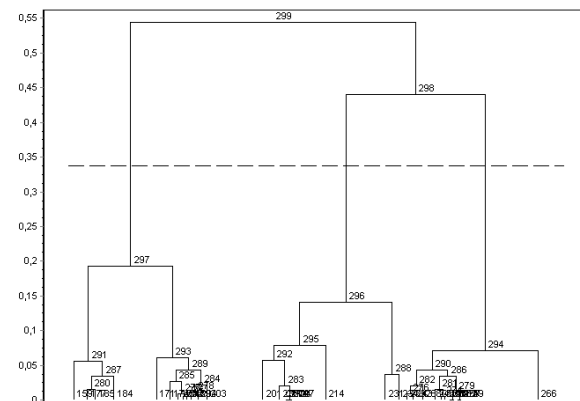
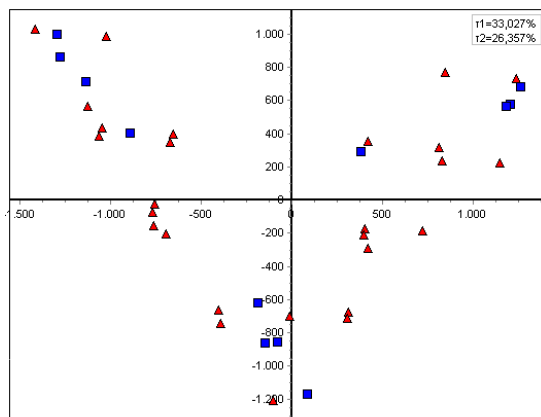
- Home
- About CHIC Analysis
- Features
- Download
- Contact

> Wise Words

"The model should follow the data, not the inverse!"

Jean-Paul Benzecri

"Science is presumably cumulative. This means that we all stand, to use Newton's beautiful phrase, "on the shoulders of giants". It also means, fortunately, that we stand on top of -



#fin