



HAL
open science

Missing Data in Hierarchical Classification of Variables, a simulation study

Ana Lorga da Silva, Helena Bacelar-Nicolau, Gilbert Saporta

► **To cite this version:**

Ana Lorga da Silva, Helena Bacelar-Nicolau, Gilbert Saporta. Missing Data in Hierarchical Classification of Variables, a simulation study. IFCS 2002, International Federation of Classification Societies, Jul 2002, Krakow, Poland. pp.121-128. hal-01124802

HAL Id: hal-01124802

<https://hal.science/hal-01124802>

Submitted on 26 Mar 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Missing Data in Hierarchical Classification of variables - a simulation study^{*}

Ana Lorga da Silva¹, Helena Bacelar-Nicolau², and Gilbert Saporta³

¹ Department of Mathematics, ISEG ,Tecnico University,
Lisbon,Portugal

² FPCE LEAD - Lisbon University,
Lisbon,Portugal

³ Statistics Department,CNAM,
Paris, FRANCE

Abstract. Here we develop from a first work the effect of missing data in hierarchical classification of variables according to the following factors: amount of missing data, imputation techniques, similarity coefficient, and aggregation criterion. We have used two methods of imputation, a regression method using an OLS method and an EM algorithm. For the similarity matrices we have used the basic affinity coefficient and the Pearson's correlation coefficient. As aggregation criteria we apply average linkage, single linkage and complete linkage methods. To compare the structure of the hierarchical classifications the Spearman's coefficient between the associated ultrametrics has been used. We present here simulation experiments in five multivariate normal cases.

1 Introduction

The missing data problem has been dealt in a large number of papers and books where several methods to minimise missing data effect have been developed (Rubin(1974), Little and Rubin(1987), Dempster, Laird and Rubin(1977), Orchard and Woodbury(1972), Beale and Little(1975) among others).

When one wants to classify variables, for instance in marketing analysis and social sciences, one frequently finds missing data. We are interested in analysing the effect of missing data in some particular (originally complete) hierarchical classification structures of variables, as well the results of imputation methods in those cases. In the present work we consider hierarchical clustering models based on two similarity coefficients - basic affinity, (Bacelar-Nicolau(1981,1988,2000), Matusita(1955), Nicolau(1998) among others) and Pearson's correlation - and three classical aggregation criteria. We use two types of imputation methods in simulation studies with different percentage

^{*} This work has been partially supported by the Franco-Portuguese Scientific Programme MSPLDM-542-B2 (Embassy of France and Portuguese Ministry of Science and Technology - ICCTI) and the Multivariate Data Analysis research team of CEAUL/FCUL.

of missing data at random. The data are issued from multinormal populations (Saporta(1990)).

2 Hierarchical cluster analysis

In this work we are interested in the classification of variables. We use the following hierarchical aggregation criteria:

Average linkage (AL): $C(A, B) = \frac{1}{(\#A) \times (\#B)} \sum c(X_j, X_{j'})$, $X_j \in A, X_{j'} \in B$

Single linkage (SL): $C(A, B) = \max \{ c(X_j, X_{j'}) \}$, $X_j \in A, X_{j'} \in B$ }

Complete linkage (CL): $C(A, B) = \min \{ c(X_j, X_{j'}) \}$, $X_j \in A, X_{j'} \in B$ }

where A and B represent two clusters and c is a similarity coefficient between two variables ($X_j, X_{j'}$ are $(n \times 1)$ variables) which can be one of the two following:

The (unweighted) basic affinity coefficient $\sum_{i=1}^n \sqrt{\frac{x_{ij} x_{ij'}}{x_{.j} x_{.j'}}$, where $x_{.j} = \sum_{i=1}^n x_{ij}$ and $x_{.j'} = \sum_{i=1}^n x_{ij'}$, as defined in Bacelar-Nicolau(2000).

The Bravais-Pearson correlation coefficient $c_p = \frac{\sum_{i=1}^n (x_{ij} - \bar{x}^j)(x_{ij'} - \bar{x}^{j'})}{s_{x_j} s_{x_{j'}}$

and C is the respective extension to the clusters. In order to compare hierarchical classification models, we will use the Spearman's coefficient- c_s -between the ultrametric matrices, based on pairs of observations with the usual correction for ties.

3 The missing data - MAR

The data are said that missing at random if its missingness does not depend of the values assumed on the variables having missing values, but depends on the values observed in other completely observed variables. The expression of the general notion of MAR can be then written as: $Prob(R|X_{obs}, X_{miss}) = Prob(R|X_{obs})$, were X_{obs} represents the observed values of $X_{n \times p}$, X_{miss} the missing values of $X_{n \times p}$ and $R = [R_{ij}]$ is a missing data indicator,

$$R_{ij} = \begin{cases} 1, & \text{if } x_{ij} \text{ observed} \\ 0, & \text{if } x_{ij} \text{ missing} \end{cases}$$

4 The imputation methods

An ordinary least square regression method (OLS) is used: is defined as usually, β_1, β_2 are estimated over the observed values of the dependent variable, $X_{obs} = \beta_1 + \beta_2 X'_{obs}$ (X'_{obs} is a "sample" of X corresponding to the observed values of $X - X_{obs}$) and then the missing values of X (X_{miss}) are imputed by the regression on X'_{miss} (those are observed values corresponding to the missing dependent values of under the estimated model $X_{miss} = \beta_1 + \beta_2 X'_{miss}$).

An EM algorithm has been used as follows:

At the E step of the algorithm (at the t th iteration),

$$x_{ij}^t = \begin{cases} x_{ij}, & \text{if } x_{ij} \text{ is observed} \\ \hat{x}_{ij}, & \text{if } x_{ij} \text{ is missing} \end{cases}$$

"The E step imputes the best linear predictors of the missing values, using current estimates of the parameters available so that a suitable choice can be made. It also calculates the adjustments c_{jk} to the estimated covariance matrix needed to allow for imputation of missing values" Little and Rubin(1987) At the M step

$$\mu^{(t+1)} = [n^{-1} \sum_{i=1}^n x_{ij}],$$

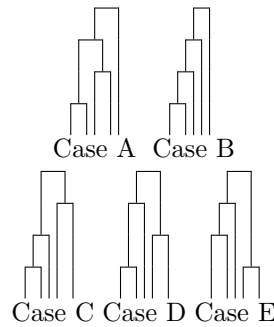
$\sigma_{jk}^{(t+1)} = (n-1)^{-1} E(\sum_{i=1}^n x_{ij}x_{ik} | X_{obs}) - \mu_j^{t+1} \mu_k^{t+1}, j, k = 1, \dots, p$. We consider missing values over a dependent variable.

5 Numerical Experiments

In order to study the performance of the affinity and the Pearson's correlation coefficients as measures of similarity between variables, in hierarchical classification and in presence of missing data, we use here the three classical hierarchical clustering methods AL, SL and CL: in the cases of complete data; in MAR case - 10%, 15% and 20% (over the total of the data - each 1000×5 matrices); and when the missing data are filled-in using the two imputation methods as mentioned in 4..

One hundred samples have been generated of each type of simulated data set, from five normal multivariate populations: Cases A, B, C, D and E, such as: $X_i \sim N(\mu_i, \Sigma_i), i = 1, \dots, 5$, and are 1000×5 matrices ($\mu_1 = \mu_2 = \mu_3 = \mu_4 = \mu_5$).

The values of the variance-covariance matrices have been chosen with the aim of obtaining specific hierarchical structures:



(order of variables: X_1, X_2, X_3, X_4, X_5)

Note that cases C, D and E, have the same topology but they have different aggregation levels. In order to have missing data at random MAR we

have deleted (10%, 15% and 20%) values at random from variables X_1 and X_2 .

In the following we present the results of the simulations, respectively in all the cases, by increasing order of the percentages of missing data (MD), according to the similarity coefficients, the agglomerative methods and the imputation methods. In each case we compare the ultrametrics associated to the originally complete data with the ultrametric matrices associated to the incomplete and the reconstructed data respectively(imputed data - ID).

The comparison between ultrametrics is obtained using a 5% Spearman's bilateral test (the critical value is $c'_s = 0,684$, see for instance Saporta (1990)). In presence of MD, the classification is obtained by a listwise method i.e. we have only considered for the analysis the complete rows (by eliminating the rows with MD). In analysis of cases A,B,C,D and E, $c_s = 1, |c_s| > c'_s, |c_s| < c'_s$ mean that:

$c_s = 1$ the general "structure" of the two hierarchical classifications being compared is the same, that is the two associated ultrametrics are "ordinal equivalent" (each pair of ranked trees give the same "ordinal" structure).

$|c_s| > c'_s$ the two hierarchical classification structures are not the same, but the two ultrametrics are "significantly correlated" (at 5%).

$|c_s| < c'_s$ the two hierarchical classification structures are "significantly different"

The percentages of cases and $c_s = 1, |c_s| > c'_s$ and $|c_s| < c'_s$ are also indicated in each cell of the tables.

All the simulated complete data reproduced the same general hierarchical structure using both coefficients - Affinity and Pearson's correlation - and the three hierarchical methods, AL, SL and CL.

MD	Affinity coefficient			Pearson's coefficient		
	AL	SL	CL	AL	SL	CL
10%	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	89% $c_s = 1$ 11% $ c_s > c'_s$	100% $c_s = 1$
15%	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	50% $c_s = 1$ 50% $ c_s > c'_s$	100% $c_s = 1$
20%	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	69% $c_s = 1$ 31% $ c_s > c'_s$	100% $c_s = 1$

Table 1. describes the results in presence of MD - Case A

6 Conclusions

In all the studied cases the affinity coefficient performs better than Pearson's correlation coefficient in presence of data missing at random.

MD	Affinity coefficient			Pearson's coefficient		
	AL	SL	CL	AL	SL	CL
10%	100% $c_s = 1$	69% $c_s = 1$ 31% $ c_s > c'_s$	100% $c_s = 1$	99% $c_s = 1$ 1% $ c_s > c'_s$	92% $c_s = 1$ 8% $ c_s > c'_s$	99% $c_s = 1$ 1% $ c_s > c'_s$
15%	100% $c_s = 1$	17% $c_s = 1$ 83% $ c_s > c'_s$	100% $c_s = 1$	99% $c_s = 1$ 1% $ c_s > c'_s$	41% $c_s = 1$ 39% $ c_s > c'_s$	99% $c_s = 1$ 1% $ c_s > c'_s$
20%	79% $c_s = 1$ 21% $ c_s > c'_s$	3% $c_s = 1$ 97% $ c_s > c'_s$	97% $c_s = 1$ 17% $ c_s > c'_s$	12% $c_s = 1$ 88% $ c_s > c'_s$	3% $c_s = 1$ 97% $ c_s > c'_s$	12% $c_s = 1$ 88% $ c_s > c'_s$

Table 2. results after using OLS method - Case A

MD	Affinity coefficient			Pearson's coefficient		
	AL	SL	CL	AL	SL	CL
10%	98% $c_s = 1$ 2% $ c_s > c'_s$	81% $c_s = 1$ 19% $ c_s > c'_s$	98% $c_s = 1$ 2% $ c_s > c'_s$	93% $c_s = 1$ 7% $ c_s > c'_s$	70% $c_s = 1$ 30% $ c_s > c'_s$	93% $c_s = 1$ 7% $ c_s > c'_s$
15%	58% $c_s = 1$ 42% $ c_s > c'_s$	14% $c_s = 1$ 86% $ c_s > c'_s$	95% $c_s = 1$ 5% $ c_s > c'_s$	91% $c_s = 1$ 9% $ c_s > c'_s$	49% $c_s = 1$ 51% $ c_s > c'_s$	94% $c_s = 1$ 6% $ c_s > c'_s$
20%	75% $c_s = 1$ 25% $ c_s > c'_s$	8% $c_s = 1$ 92% $ c_s > c'_s$	80% $c_s = 1$ 20% $ c_s > c'_s$	31% $c_s = 1$ 69% $ c_s > c'_s$	1% $c_s = 1$ 99% $ c_s > c'_s$	31% $c_s = 1$ 69% $ c_s > c'_s$

Table 3. results after using EM method - Case A

MD	Affinity coefficient			Pearson's coefficient		
	AL	SL	CL	AL	SL	CL
10%	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$
15%	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	98% $c_s = 1$ 2% $ c_s > c'_s$	100% $c_s = 1$	100% $c_s = 1$
20%	100% $c_s = 1$	99% $c_s = 1$ 1% $ c_s > c'_s$	100% $c_s = 1$	94% $c_s = 1$ 6% $ c_s > c'_s$	97% $c_s = 1$ 3% $ c_s > c'_s$	81% $c_s = 1$ 18% $ c_s > c'_s$ 1% $ c_s < c'_s$

Table 4. results in presence of MD - Case B

MD	Affinity coefficient			Pearson's coefficient		
	AL	SL	CL	AL	SL	CL
10%	99% $c_s = 1$ 1% $ c_s > c'_s$	99% $c_s = 1$ 1% $ c_s > c'_s$	99% $c_s = 1$ 1% $ c_s > c'_s$	99% $c_s = 1$ 1% $ c_s > c'_s$	99% $c_s = 1$ 1% $ c_s > c'_s$	99% $c_s = 1$ 1% $ c_s > c'_s$
15%	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$
20%	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$

Table 5. results after using OLS method -Case B

Better results are obtained in presence of MD, than after the application of both imputation methods.

When using the imputation methods in case C both imputation methods gave the same results, and also that the affinity coefficient performs better than the Pearson's coefficient. In cases A, D and E the results are different when using the two imputation methods, some times the least squares method performs better, others, is with the EM algorithm that we obtain better performance. We have obtained similar results in a study with Personality development data, in presence of seven variables(Silva et al.(2001))

MD	Affinity coefficient			Pearson's coefficient		
	AL	SL	CL	AL	SL	CL
10%	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$
15%	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$
20%	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$

Table 6. results after using EM method - Case B

MD	Affinity coefficient			Pearson's coefficient		
	AL	SL	CL	AL	SL	CL
10%	100% $c_s = 1$	99% $c_s = 1$ 1% $ c_s > c'_s$	100% $c_s = 1$	5% $c_s = 1$ 95% $ c_s < c'_s$	49% $c_s = 1$ 51% $ c_s < c'_s$	100% $ c_s < c'_s$
15%	100% $c_s = 1$	96% $c_s = 1$ 4% $ c_s > c'_s$	99% $c_s = 1$ 1% $ c_s > c'_s$	5% $c_s = 1$ 95% $ c_s < c'_s$	41% $c_s = 1$ 59% $ c_s < c'_s$	100% $c_s < 1$
20%	100% $c_s = 1$	87% $c_s = 1$ 13% $ c_s > c'_s$	100% $c_s = 1$	3% $c_s = 1$ 97% $ c_s < c'_s$	21% $c_s = 1$ 69% $ c_s < c'_s$	100% $ c_s < c'_s$

Table 7. results in presence of MD - Case C

MD	Affinity coefficient			Pearson's coefficient		
	AL	SL	CL	AL	SL	CL
10%	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	98% $c_s = 1$ 2% $ c_s < c'_s$	100% $c_s = 1$	98% $c_s = 1$ 2% $ c_s < c'_s$
15%	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	94% $c_s = 1$ 5% $ c_s > c'_s$ 1% $ c_s < c'_s$	96% $c_s = 1$ 4% $ c_s > c'_s$	95% $c_s = 1$ 4% $ c_s > c'_s$ 1% $ c_s < c'_s$
20%	94% $c_s = 1$ 2% $ c_s > c'_s$ 4% $ c_s < c'_s$	94% $c_s = 1$ 2% $ c_s > c'_s$ 4% $ c_s < c'_s$	94% $c_s = 1$ 2% $ c_s > c'_s$ 4% $ c_s < c'_s$	21% $c_s = 1$ 58% $ c_s > c'_s$ 21% $ c_s < c'_s$	21% $c_s = 1$ 58% $ c_s > c'_s$ 21% $ c_s < c'_s$	21% $c_s = 1$ 58% $ c_s > c'_s$ 21% $ c_s < c'_s$

Table 8. results after using both imputation methods - Case C

MD	Affinity coefficient			Pearson's coefficient		
	AL	SL	CL	AL	SL	CL
10%	99% $c_s = 1$ 1% $ c_s > c'_s$	100% $c_s = 1$	100% $c_s = 1$	97% $c_s = 1$ 3% $ c_s > c'_s$	96% $c_s = 1$ 4% $ c_s > c'_s$	99% $c_s = 1$ 1% $ c_s > c'_s$
15%	98% $c_s = 1$ 1% $ c_s > c'_s$ 1% $ c_s < c'_s$	99% $c_s = 1$ 1% $ c_s > c'_s$	99% $c_s = 1$ 1% $ c_s > c'_s$	95% $c_s = 1$ 5% $ c_s < c'_s$	90% $c_s = 1$ 10% $ c_s > c'_s$	95% $c_s = 1$ 5% $ c_s > c'_s$
20%	99% $c_s = 1$ 1% $ c_s < c'_s$	100% $c_s = 1$	100% $c_s = 1$	92% $c_s = 1$ 8% $ c_s < c'_s$	85% $c_s = 1$ 15% $ c_s > c'_s$	93% $c_s = 1$ 7% $ c_s > c'_s$

Table 9. results in presence of MD - Case D

The following developments of this work are related to other similarity coefficients and hierarchical structures, namely concerning a probabilistic classification approach, the use of the probabilistic weighted coefficient (without ignore objects or impute the missing data) and different types of missing data and imputation methods.

MD	Affinity coefficient			Pearson's coefficient		
	AL	SL	CL	AL	SL	CL
10%	100% $ c_s > c'_s$	99% $c_s = 1$ 1% $ c_s > c'_s$	96% $c_s = 1$ 4% $ c_s > c'_s$	59% $c_s = 1$ 40% $ c_s > c'_s$ 1% $ c_s < c'_s$	47% $c_s = 1$ 63% $ c_s > c'_s$	60% $c_s = 1$ 40% $ c_s > c'_s$
15%	9% $c_s = 1$ 90% $ c_s > c'_s$ 1% $ c_s < c'_s$	6% $c_s = 1$ 94% $ c_s > c'_s$	18% $c_s = 1$ 82% $ c_s > c'_s$	70% $c_s = 1$ 29% $ c_s > c'_s$ 1% $ c_s < c'_s$	52% $c_s = 1$ 47% $ c_s > c'_s$	76% $c_s = 1$ 23% $ c_s > c'_s$
20%	99% $ c_s > c'_s$ 1% $ c_s < c'_s$	100% $ c_s > c'_s$	100% $ c_s > c'_s$	31% $c_s = 1$ 68% $ c_s > c'_s$ 1% $ c_s < c'_s$	18% $c_s = 1$ 82% $ c_s > c'_s$	68% $c_s = 1$ 32% $ c_s > c'_s$

Table 10. results after using OLS method - Case D

MD	Affinity coefficient			Pearson's coefficient		
	AL	SL	CL	AL	SL	CL
10%	98% $c_s = 1$ 2% $ c_s > c'_s$	83% $c_s = 1$ 12% $ c_s > c'_s$ 1% $ c_s < c'_s$	92% $c_s = 1$ 7% $ c_s > c'_s$ 1% $ c_s < c'_s$	93% $c_s = 1$ 7% $ c_s > c'_s$	99% $c_s = 1$ 1% $ c_s > c'_s$	93% $c_s = 1$ 7% $ c_s > c'_s$
15%	58% $c_s = 1$ 42% $ c_s > c'_s$	21% $c_s = 1$ 79% $ c_s > c'_s$	87% $c_s = 1$ 1% $ c_s > c'_s$ 12% $ c_s < c'_s$	91% $c_s = 1$ 9% $ c_s > c'_s$	80% $c_s = 1$ 19% $ c_s > c'_s$ 1% $ c_s < c'_s$	93% $c_s = 1$ 6% $ c_s > c'_s$ 1% $ c_s < c'_s$
20%	75% $c_s = 1$ 25% $ c_s > c'_s$	1% $c_s = 1$ 99% $ c_s > c'_s$	97% $ c_s > c'_s$ 3% $ c_s < c'_s$	31% $c_s = 1$ 69% $ c_s > c'_s$	19% $c_s = 1$ 81% $ c_s > c'_s$	98% $ c_s > c'_s$ 2% $ c_s < c'_s$

Table 11. results after using EM method - Case D

MD	Affinity coefficient			Pearson's coefficient		
	AL	SL	CL	AL	SL	CL
10%	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$
15%	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	97% $c_s = 1$ 3% $ c_s < c'_s$	98% $c_s = 1$ 2% $ c_s < c'_s$
20%	100% $c_s = 1$	98% $c_s = 1$ 1% $ c_s > c'_s$ 1% $ c_s < c'_s$	99% $c_s = 1$ 1% $ c_s > c'_s$	100% $c_s = 1$	87% $c_s = 1$ 13% $ c_s < c'_s$	92% $c_s = 1$ 8% $ c_s < c'_s$

Table 12. results in presence of MD - Case E

MD	Affinity coefficient			Pearson's coefficient		
	AL	SL	CL	AL	SL	CL
10%	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$
15%	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$	100% $c_s = 1$
20%	99% $c_s = 1$ 1% $ c_s > c'_s$	99% $c_s = 1$ 1% $ c_s > c'_s$	99% $c_s = 1$ 1% $ c_s > c'_s$	73% $c_s = 1$ 25% $ c_s > c'_s$ 2% $ c_s < c'_s$	73% $c_s = 1$ 25% $ c_s > c'_s$ 2% $ c_s < c'_s$	73% $c_s = 1$ 25% $ c_s > c'_s$ 2% $ c_s < c'_s$

Table 13. results after using OLS method - Case E

References

- BACELAR-NICOLAU, H.(1981): Contributions to the Study of Comparison Coefficients in Cluster Analysis, Univ. Lisbon.
- BACELAR-NICOLAU, H.(1988), Two probabilistic models for classification of variables in frequency tables, *Classif. and Relat. Meth. of Data Analysis*, H. .H.

MD	Affinity coefficient			Pearson's coefficient		
	AL	SL	CL	AL	SL	CL
10%	$100\%c_s = 1$	$100\%c_s = 1$	$100\%c_s = 1$	$100\%c_s = 1$	$100\%c_s = 1$	$100\%c_s = 1$
15%	$100\%c_s = 1$	$100\%c_s = 1$	$100\%c_s = 1$	$100\%c_s = 1$	$100\%c_s = 1$	$100\%c_s = 1$
20%	$99\%c_s = 1$	$99\%c_s = 1$	$99\%c_s = 1$	$78\%c_s = 1$	$76\%c_s = 1$	$76\%c_s = 1$
	$1\% c_s > c'_s$	$1\% c_s > c'_s$	$1\% c_s > c'_s$	$22\% c_s > c'_s$	$24\% c_s > c'_s$	$24\% c_s > c'_s$

Table 14. results after using EM method - Case E

- Bock (ed.), North Holland, pp. 181-186
- BACELAR-NICOLAU(2000) The Affinity Coefficient in Analysis of Symbolic Data Exploratory Methods for Extracting Statistical Information from Complex Data. H.H. Bock and E.Diday (Eds.), Springer,160-165.
- BEALE, E. M. L. and LITTLE, R. J. A.(1975) Missing values in multivariate data analysis. *J. R. Statist. Soc. B*, **37**, 129-145.
- DEMPSTER, A. P., LAIRD, N. M. and RUBIN, D. B.(1977) Maximum likelihood estimation from incomplete data via the EM algorithm (with discussion). *J. R. Statist. Soc. B*, **39**, 1-38.
- LITTLE, R. J. A. and RUBIN, D. B.(1987) Statistical Analysis With Missing Data, John Wiley & Sons, New York.
- MATUSITA,K.(1955) Decision rules, based on distance for problems of fit, two samples and estimation. *Ann. Math. Stat.*, vol26, n4,631-640.
- NICOLAU F.C., BACELAR-NICOLAU, H. (1998), Some Trends in the Classification of variables, Data Science, Classification, and Related Methods, C. Hayashi, N. Ohsumi, K. Yajima, Y. Tanaka, H. H. Bock, Y. Baba (eds.), Springer, pp. 89-98
- ORCHARD, T. and WOODBURY, M. A.(1972) A missing information principle: theory and applications. Proceedings 6th Berkley Symposium on Mathematical Statistic and Probability, **1**, 697-715.
- RUBIN, D. B.(1974) Characterising the estimation of parameters in the estimation of parameters in incomplete-data problems. *JASA*, **69**,467-474.
- SAPORTA, G.(1990) Probabilités, Analyse des Données et Statistique, Editions Technip, Paris.
- SILVA,A.L, BACELAR-NICOLAU, SAPORTA, G. and GEADA, M.(2001) Missing Data in Hierarchical Classification – a study with Personality development data, - 32nd European Mathematical Psychology /EMPG 2001, 109-110.