

A Multiple Correspondence Analysis to Organize Data Cubes

Riadh BEN MESSAOUD, Omar BOUSSAID and Sabine LOUDCHER RABASÉDA

*Laboratory ERIC – University of Lyon 2
5 avenue Pierre Mendès-France, 69676, Bron Cedex – France*

Abstract. On Line Analytical Processing (OLAP) is a technology basically created to provide users with tools in order to explore and navigate into data cubes. Unfortunately, in huge and sparse data, exploration becomes a tedious task and the simple user's intuition or experience does not lead to efficient results. In this paper, we propose to exploit the results of the Multiple Correspondence Analysis (MCA) in order to enhance data cube representations and make them more suitable for visualization and thus, easier to analyze. Our approach addresses the issues of organizing data in an interesting way and detects relevant facts. Our purpose is to help the interpretation of multidimensional data by efficient and simple visual effects. To validate our approach, we compute its efficiency by measuring the quality of resulting multidimensional data representations. In order to do so, we propose an homogeneity criterion to measure the visual relevance of data representations. This criterion is based on the concept of geometric neighborhood and similarity between cells. Experimental results on real data have shown the interest of using our approach on sparse data cubes.

Keywords. OLAP, Data cubes, data representation, MCA, test-values, arrangement of attributes, characteristic attributes, homogeneity criterion

Introduction

On-Line Analytical Processing (OLAP) is a technology supported by most data warehousing systems [8,11]. It provides a platform for analyzing data according to multiple dimensions and multiple hierarchical levels. Data are presented in multidimensional views, commonly called data cubes [3]. A data cube can be considered as a space representation composed by a set of cells. A cell is associated with one or more measures and identified by coordinates represented by one attribute from each dimension. Each cell in a cube represents a precise fact. For example, if dimensions are *products*, *stores* and *months*, the measure of a particular cell can be the *sales* of one *product* in a particular *store* on a given *month*. OLAP provides the user with visual based tools to summarize, explore and navigate into data cubes in order to detect interesting and relevant information. However, exploring a data cube is not always an easy task to perform. Obviously, in large cubes containing sparse data, the whole analysis process becomes tedious and complex. In such a case, an intuitive exploration based on the user's experience does not quickly lead to efficient results. More generally, in the case of a data cube with more than three dimensions, a user is naturally faced to a hard task of navigation and explo-

ration in order to detect relevant information. Current OLAP provides query-driven and visual tools to browse data cubes, but does not deeply assist the user and help him/her to investigate interesting patterns.

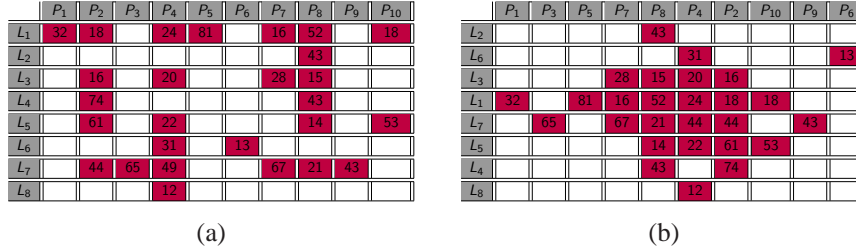


Figure 1. Example of two representations of a 2-dimensional data cube.

For example, consider the cube of Figure 1. On the one hand, representation 1(a) displays sales of products (P_1, \dots, P_{10}) crossed by geographic locations of stores (L_1, \dots, L_8). In this representation, full cells (gray cells) are displayed randomly according to the lexical ordering of the *attributes* – also called *members* – in each dimension. The way the cube is displayed does not provide an attractive representation that visually helps a user to easily interpret data.

On the other hand, Figure 1(b) contains the same information as Figure 1(a). However, it displays a data representation which is visually easier to analyze. In fact, the cube of Figure 1(b) expresses important relationships by providing a visual representation that gathers full cells together and separates them from empty ones. In a natural way, such a representation is more comfortable to the user and allows to drive easy and efficient analysis.

Nevertheless, note that the representation (b) of Figure 1 can be interactively constructed by the user from representation (a) via some classic OLAP operators. This suppose that the user intuitively knows how to arrange the attributes of each dimension. Hence, we propose to provide the user with an automatic assistance to identify interesting facts and arrange them in a suitable visual representation [19,18]. As shown in Figure 1, we propose an approach that allows the user to get relevant facts expressing relationships and displays them in an appropriate way that enhances the exploration process independently of the cube's size. Thus, we suggest to carry out a Multiple Correspondence Analysis [9] (MCA) on a data cube as a preprocessing step. Basically, MCA is a powerful describing method even for huge volumes of data. It factors categorical variables and displays data in a factorial space constructed by orthogonal system of axes that provides relevant views of data. These elements motivate us to exploit the results of the MCA in order to better explore large data cubes by identifying and arranging its interesting facts. The first constructed factorial axis summarizes the maximum of information contained in the cube. We focus on relevant OLAP facts associated with characteristic attributes (variables) given by the factorial axes. These facts are interesting since they reflect relationships and concentrate a significant information. For a better visualization of these facts, we highlight them and arrange their attributes in the data space representation by using the *test-values* [14].

In order to evaluate the visual relevance of multidimensional data representations, we also propose in this paper a novel criterion to measure the homogeneity of cells distribution in the space representation of a data cube [17]. This criterion is based on geometric neighborhood of data cube cells, and also takes into account the similarity of cells' measures and provides a scalar quantification for the homogeneity of a given data cube representation. It also allows to evaluate the performance of our approach by comparing the quality of the initial data representation and the arranged one.

This paper is organized as follows. In section 1, we present related work to our approach. We provide in section 2 the problem formalization and present the general context of this work. The section 4 introduces the *test-values* and details the steps of our approach. We define in the next section our quality representation criterion. The section 6 presents a real world case study on a huge and sparse data cube. We propose experimental results in the section 7. Finally, we conclude and propose some future researches directions.

1. Related Work

Several works have already treated the issue of enhancing the space representation of data cubes. These works were undertaken following different motivations and adopted different ways to address the problem. We note that while many efforts are interested to computational aspects of data cubes (optimization of storage space, compression strategies, queries response time, etc.), a small number of studies have focused on OLAP aspects. Our present work fits into the second category. In our work, we focus on assisting OLAP users in order to improve and help the analysis process on large and sparse data cubes. We use a factorial approach to highlight relevant facts and provide an interesting visual data representations. Nevertheless, we dress an overview of main studies as well in the first as in the second category.

Some studies approximate computation of compressed data cube. In [23], Vitter *et al.* proposed to build compact data cubes by using approximation through wavelets. **Quasi-Cube** [1] compresses data representation by materializing only sufficient parts of a data cube. In [21] approximation is performed by estimating the density function of data. Other efforts address the issue of computing data cubes with index structure. For instance, **Dwarf** [22] uses indexes to reduce the storage space of a cube by identifying and factoring redundant tuples. Wang *et al.* propose to factorize data redundancies with **BST** [24] (*Base Single Tuple*). In [6], Feng *et al.* introduce **PrefixCube**, a data structure based on only one **BST**. The **Quotient Cube** [12] summarizes the semantic contents of a data cube and partitions it into similar cells. In [13], **QC-Tree** is directly constructed from the base table in order to maintain it under updates. Some other studies optimize storage spaces by partitioning the initial cube. **Range CUBE** [7] identifies correlations between attributes and compresses the data cube. **Partitioned-Cube** [20] partitions large relations into fragments. Operations on the cube are, therefore, performed in memory-sized fragments independently. In [15], high dimensional data are transformed into small local cubes and used to for online queries.

Finally, our approach shares already the same motivation of Choong *et al* [4]. The authors address the problem of high dimensionality of data cubes. They try to enhance analysis processes by preparing the dataset into appropriate representation. Thus, the

user can explore it in a more effective manner. The authors use an approach that combines association rules algorithm and a fuzzy subsets method. Their approach consists in identifying blocks of similar measures in the data cube. However, this approach does not take into account the problem of data sparsity. Furthermore, it does not provide a quality evaluation of the resulting multidimensional representations.

We emphasize that our approach does not deal with the issues of data cube compression, reduction of dimensionality or optimization of storage space. Through this study, we try to act on sparsity in huge multidimensional representations. Not to reduce it, but to reduce its negative effects on the interpretations and OLAP analysis of data [19,18]. Thus, we use the MCA to arrange differently the facts and highlight their relevant relationships in a data cube within a visual effect that gathers them as well as possible in the space representation.

2. Problem Formalization

Let \mathcal{C} denote a data cube. Note that, our approach can be applied directly on \mathcal{C} or on a data view (a sub-cube) extracted from \mathcal{C} . It is up to the user to select dimensions, fix one hierarchical level per dimension and select measures in order to create a particular data view (s)he wishes to visualize. Thus, to enhance the data representation of the constructed view, the user can apply on it our proposed approach. In order to lighten the formalization, in the followings of the paper, we assume that a user has selected a data cube \mathcal{C} , with d dimensions $(D_1, \dots, D_t, \dots, D_d)$, m measures $(M_1, \dots, M_q, \dots, M_m)$ and n facts. We also assume that the user has fixed one hierarchical level with p_t categorical attributes per dimension. Let a_j^t the j^{th} attribute of the dimension D_t and $p = \sum_{t=1}^d p_t$ the total number of attributes in \mathcal{C} . For each dimension D_t , we note $\{a_1^t, \dots, a_j^t, \dots, a_{p_t}^t\}$ the set of its attributes.

In a first step, the aim of our approach is to organize the space representation of a given data cube \mathcal{C} by arranging the attributes of its dimensions. For each dimension D_t , our approach establishes a new arrangement of its attributes a_j^t in the data space (see subsection 4.2). This arrangement provides a data representation visually easier to interpret and displays multidimensional information in a more suitable way for analysis. In a second step, our approach detects from the resulted representation relevant facts expressing interesting relationships. To do that, we select from each dimension D_t a subset Φ_t of significant attributes, also called characteristic attributes (see subsection 4.3). The crossing of these particular attributes allows to identify relevant cells in the cube.

Our approach is based on the MCA [9,14]. The MCA is a factorial method that displays categorical variables in a property space which maps their associations in two or more dimensions. From a table of n observations and p categorical variables, describing a p -dimensional cloud of individuals ($p < n$), the MCA provides orthogonal axes to describe the most variance of the whole data cloud. The fundamental idea is to reduce the dimensionality of the original data thanks to a reduced number of variables (factors) which are a combination of the original ones. The MCA is generally used as an exploratory approach to unearth empirical regularities of a dataset.

In our case, we assume the cube's facts as the individuals of the MCA, the cube's dimensions as its variables, and the attributes of a dimension as values of their corresponding variables. We apply the MCA on the n facts of the cube \mathcal{C} and use its results

Id	D ₁	D ₂	D ₃	M ₁
1	L1	T2	P1	9
2	L2	T2	P3	5
3	L2	T1	P2	6
4	L1	T1	P3	7

(a)

Id	Z						
	Z ₁		Z ₂		Z ₃		
	L1	L2	T1	T2	P1	P2	P3
1	1	0	0	1	1	0	0
2	0	1	0	1	0	0	1
3	0	1	1	0	0	1	0
4	1	0	1	0	0	0	1

(b)

Figure 2. Example of a conversion of a data cube to a complete disjunctive table.

to build *test-values* (see subsection 4.1) for the attributes a_j^t of the dimensions D_t . We exploit these *test-values* to arrange attributes and detect characteristic ones in their corresponding dimensions.

3. Applying the MCA on a Data Cube

Like all statistical methods, MCA needs a tabular representation of data as input. Therefore, we can not apply it directly on multidimensional representations like data cubes. Therefore, we need to convert \mathcal{C} to a *complete disjunctive table*. For each dimension D_t , we generate a binary matrix Z_t with n rows and p_t columns. Rows represent facts, and columns represent dimension's attributes. The i^{th} row of Z_t contains $(p_t - 1)$ times the value 0 and one time the value 1 in the column that fits with the attribute taken by the fact i . The general term of Z_t is:

$$z_{ij}^t = \begin{cases} 1 & \text{if the fact } i \text{ takes the attribute } a_j^t \\ 0 & \text{else} \end{cases} \quad (1)$$

By merging the d matrices Z_t , we obtain a complete disjunctive table $Z = [Z_1, Z_2, \dots, Z_t, \dots, Z_d]$ with n rows and p columns. It describes the d positions of the n facts of \mathcal{C} through a binary coding. For instance, Figure 2 shows an simple example of a data cube (a), with 3 dimensions $D_1 : \{L_1, L_2\}$, $D_2 : \{T_1, T_2\}$, and $D_3 : \{P_1, P_2, P_3\}$. This cube is converted to a complete disjunctive table Z in Figure 2(b). In the case of a large data cube, we naturally obtain a very huge matrix Z . Recall that MCA is a factorial method perfectly suited to huge input dataset with high numbers of rows and columns.

Once the complete disjunctive table Z is built, MCA starts by constructing a matrix $B = Z'Z$ – called *Burt table* –, where Z' is the transposed matrix of Z . *Burt table* B is a (p, p) symmetric matrix which contains all the category marginal on the main diagonal and all possible cross-tables of the d dimensions of \mathcal{C} in the off-diagonal. Let X be a (p, p) diagonal matrix which has the same diagonal elements of B and zeros otherwise. We construct from Z and X a new matrix S according to the formula:

$$S = \frac{1}{d} Z'ZX^{-1} = \frac{1}{d} BX^{-1} \quad (2)$$

By diagonalizing S , we obtain $(p - d)$ diagonal elements, called *eigenvalues* and denoted λ_α . Each eigenvalue λ_α is associated to a directory vector u_α and corresponds to a factorial axis F_α , where $Su_\alpha = \lambda_\alpha u_\alpha$.

An eigenvalue represents the amount of inertia (variance) that reflects the relative importance of its axis. The first axis always explains the most inertia and has the largest eigenvalue. Usually, in a factorial analysis process, researchers keep only the first, two or three axes of inertia. Other researchers give complex mathematical criterion [2,10,16,5] to determine the number of axes to keep. In [9], Benzecri suggests that this limit should be fixed by user's capacity to give a meaningful interpretation to the axes he keeps. It is not because an axis has a relatively small eigenvalue that we should discard it. It can often help to make a fine point about the data. It is up to the user to choose the number k of axis to keep by checking eigenvalues and the general meaning of axes.

4. Organizing Data Cubes and Detecting Relevant Facts

Usually in a factorial analysis, relative contributions of variables are used to give sense to the axes. A relative contribution shows the percentage of inertia of a particular axis which is explained by an attribute. The largest relative contribution of a variable to an axis is, the more it gives sense to this axis. In our approach, we interpret a factorial axis by characteristic attributes detected through the use of the *test-values* proposed by Lebart *et al.* in [14]. In the followings, we present the theoretical principle of test-values applied to the context of our approach.

4.1. Test-Values

Let $I(a_j^t)$ denotes the set of facts having a_j^t as attribute in the dimension D_t . We also note $n_j^t = \text{Card}(I(a_j^t)) = \sum_{i=1}^n z_{ij}^t$ the number of elements in $I(a_j^t)$. It corresponds to the number of facts in \mathcal{C} having a_j^t as attribute (weight of a_j^t in the cube). $\varphi_{\alpha j}^t = \frac{1}{n_j^t \sqrt{\lambda_\alpha}} \sum_{i \in I(a_j^t)} \psi_{\alpha i}$ is the coordinate of a_j^t on the factorial axis F_α , where $\psi_{\alpha i}$ is the coordinate of the facts i on F_α . Suppose that, under a null hypothesis H_0 , the n_j^t facts are selected randomly in the set of the n facts, the mean of their coordinates in F_α can be represented by a random variable $Y_{\alpha j}^t = \frac{1}{n_j^t} \sum_{i \in I(a_j^t)} \psi_{\alpha i}$, where $E(Y_{\alpha j}^t) = 0$ and $\text{VAR}_{H_0}(Y_{\alpha j}^t) = \frac{n-n_j^t}{n-1} \frac{\lambda_\alpha}{n_j^t}$.

Remark that $\varphi_{\alpha j}^t = \frac{1}{\sqrt{\lambda_\alpha}} Y_{\alpha j}^t$. Thus, $E(\varphi_{\alpha j}^t) = 0$, and $\text{VAR}_{H_0}(\varphi_{\alpha j}^t) = \frac{n-n_j^t}{n-1} \frac{1}{n_j^t}$. Therefore, the test-value of the attribute a_j^t is:

$$V_{\alpha j}^t = \sqrt{n_j^t \frac{n-1}{n-n_j^t}} \varphi_{\alpha j}^t \quad (3)$$

$V_{\alpha j}^t$ measures the number of standard deviations between the attribute a_j^t (the gravity center of the n_j^t facts) and the center of the factorial axis F_α . The position of an attribute is interesting for a given axis F_α if its cloud of facts is located in a narrow zone in the direction α . This zone should also be as far as possible from the center of the axis. The test-value is a criterion that quickly provides an appreciation if an attribute has a *significant* position on a given factorial axis or not.

4.2. Arrangement of Attributes

In a classic OLAP representation of data cubes, attributes are usually organized according to a lexical order such as alphabetic order for *geographic* dimensions or chronological order for *times* dimensions. In our approach, we propose to exploit the test-values of attributes in order to organize differently the data cube's facts. The new organization will display a relevant data representation easier to analyze and to interpret, especially in the case of large and sparse cubes. For each dimension, we sort its attributes according to the increasing order of their test-values. Actually, a test-value indicates the position of an attribute on a given axis. The relative geometric position of an attribute is more significant to factorial axes when these axes are important (have the greatest eigenvalues). For this, we propose to sort attributes according to the k first axes selected by the user. We sort the p_t test-values $V_{\alpha_j}^t$ of the attributes a_j^t on the axis F_α . This will provide a new order of indices j . According to this order, we arrange attributes a_j^t in the dimension D_t .

In general, we assume that all attributes of a dimension D_t are geometrically ordered in the data cube space representation according to the order of indices j_t , i.e, the attribute $a_{j_t-1}^t$ precedes $a_{j_t}^t$ and $a_{j_t}^t$ precedes $a_{j_t+1}^t$ (see the example of Figure 3). Indices j_t are ordered according to the arrangement of the attributes in the space representation of the dimension D_t .

4.3. Characteristic Attributes

In general, an attribute is considered significant for an axis if the absolute value of its test-value is higher than $\tau = 2$. This roughly corresponds to an error threshold of 5%. We note that, the lower error threshold is, the greater τ is. In our case, for one attribute, the test of the hypothesis H_0 can induce a possible error. This error will inevitably be increased when we perform the test p times for all the attributes of the cube. To minimize this accumulation of errors, we propose to fix for each test an error threshold of 1% which correspond to $\tau = 3$. We also note that, when a given axis can be characterized by too much attributes according to their test-values, instead of taking them all, we can restrict the selection by considering only a percentage of the most characteristic ones. i.e, those having the highest absolute test-values. Finally to detect interesting facts in a data cube, for each dimension D_t , we select the following set of characteristic attributes.

$$\Phi_t = \left\{ \begin{array}{l} a_j^t, \text{ where } \forall j \in \{1, \dots, p_t\}, \\ \exists \alpha \in \{1, \dots, k\} \text{ such as } |V_{\alpha_j}^t| \geq 3 \end{array} \right\} \quad (4)$$

5. Quality of a Data Representation

We provide a quality criterion of data cube representations [17]. It measures the homogeneity of the geometric distribution of cells in a data cube. One cell contains one or more measures of an OLAP fact. The attributes of a cell are coordinates of a fact according to dimensions in the data space representation. Let $A = (a_{j_1}^1, \dots, a_{j_t}^t, \dots, a_{j_d}^d)$ be a cell in \mathcal{C} , where $t \in \{1, \dots, d\}$ and $j_t \in \{1, \dots, p_t\}$. j_t is the index of the attribute that takes the cell A according to D_t . We denote $|A|$ the value of a measure contained in A which is equal to NULL if A is empty. For example, in Figure 3, $|A| = 5.7$ whereas

$|Y| = \text{NULL}$. We say that a cell $B = (b_{j_1}^1, \dots, b_{j_t}^t, \dots, b_{j_d}^d)$ is neighbor of A , denoted $B \dashv A$, if $\forall t \in \{1, \dots, d\}$, the coordinates of B satisfy:

$$b_{j_t}^t = a_{j_t-1}^t \text{ or } b_{j_t}^t = a_{j_t}^t \text{ or } b_{j_t}^t = a_{j_t+1}^t$$

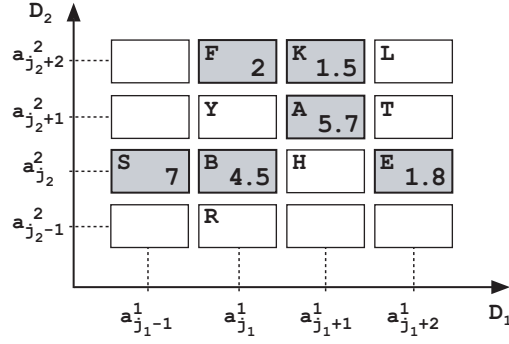


Figure 3. A 2-dimensional example of a data cube.

In Figure 3, cell B is neighbor of A ($B \dashv A$). Y is also neighbor of A ($Y \dashv A$). Whereas cells S and R are not neighbors of A . For a cell A of a cube \mathcal{C} , we define the neighborhood of A , denoted $\mathcal{N}(A)$, by the set of all cells B of \mathcal{C} neighbors of A .

$$\mathcal{N}(A) = \{B \in \mathcal{C} \text{ where } B \dashv A\}$$

For example, in Figure 3, the neighborhood of A corresponds to the set $\mathcal{N}(A) = \{F, K, L, T, E, H, B, Y\}$. To evaluate similarities between neighbor cells, we define a similarity function δ .

Definition The similarity δ of two cells A and B from a cube \mathcal{C} is defined as follows:

$$\delta : \mathcal{C} \times \mathcal{C} \longrightarrow \mathbb{R}$$

$$\delta(A, B) \longmapsto \begin{cases} 1 - \left(\frac{||A| - |B||}{\max(\mathcal{C}) - \min(\mathcal{C})} \right) & \text{if } A \text{ and } B \text{ are full} \\ 0 & \text{else} \end{cases}$$

Where $||A| - |B||$ is the absolute difference of measures contained in cells A and B , and $\max(\mathcal{C})$ (respectively, $\min(\mathcal{C})$) is the maximum (respectively, the minimum) measure value in cube \mathcal{C} .

In the cube of Figure 3, where grayed cells are full and white ones are empty, $\max(\mathcal{C}) = 7$, which matches with the cell S and $\min(\mathcal{C}) = 1.5$, which matches with the cell K . For instance, $\delta(A, B) = 1 - \left(\frac{|5.7 - 4.5|}{7 - 1.5} \right) \simeq 0.78$ and $\delta(A, Y) = 0$.

Now, let consider a function Δ from \mathcal{C} to \mathbb{R} such as $\forall A \in \mathcal{C}$, $\Delta(A) = \sum_{B \in \mathcal{N}(A)} \delta(A, B)$. It corresponds to the sum of the similarities of A with all its full neighbor cells. For instance, according to Figure 3, $\Delta(A)$ is computed as follows: $\Delta(A) = \delta(A, F) + \delta(A, K) + \delta(A, L) + \delta(A, T) + \delta(A, E) + \delta(A, H) + \delta(A, B) + \delta(A, Y) \simeq 1.64$.

We introduce the crude homogeneity criterion of a data cube \mathcal{C} according to:

$$chc(\mathcal{C}) = \sum_{\substack{A \in \mathcal{C} \\ |A| \neq \text{NULL}}} \sum_{B \in \mathcal{N}(A)} \delta(A, B) = \sum_{\substack{A \in \mathcal{C} \\ |A| \neq \text{NULL}}} \Delta(A)$$

The crude homogeneity criterion computes the sum of similarities of every couple of full and neighbor cells in a data cube. For instance, in Figure 3, the crude homogeneity criterion is computed as $chc(\mathcal{C}) = \Delta(F) + \Delta(K) + \Delta(A) + \Delta(S) + \Delta(B) + \Delta(E) \simeq 6.67$. Note that, the crude homogeneity criterion of a data cube touches its maximum when all the cells of the cube are full and have equal measures. We denote $chc_{max}(\mathcal{C}) = \sum_{A \in \mathcal{C}} \sum_{B \in \mathcal{N}(A)} 1$.

Definition The homogeneity criterion of a data cube is defined as:

$$hc(\mathcal{C}) = \frac{chc(\mathcal{C})}{chc_{max}(\mathcal{C})} = \frac{\sum_{\substack{A \in \mathcal{C} \\ |A| \neq \text{NULL}}} \Delta(A)}{\sum_{A \in \mathcal{C}} \sum_{B \in \mathcal{N}(A)} 1}$$

The homogeneity criterion evaluates the quality of a multidimensional data representation. This quality is rather better when full and similar cells are neighbors. Indeed, when similar cells are gathered in specific regions of the space representation of a data cube, this cube is easier to visualize and so, a user can directly focus his/her data interpretation on these regions.

For example, in Figure 3, $chc_{max}(\mathcal{C}) = 84$. So, the homogeneity criterion of this representation is: $hc(\mathcal{C}) = \frac{6.67}{84} \simeq 0.08$. Nevertheless, such a criterion can not make real sense for a single situation of a data representation. In all cases, we should rather compare it to other data representations of the same cube. In fact, recall that the aim of our method is to organize the facts of an initial data cube representation by arranging attributes in each dimensions according to the order of test-values. Let us denote the initial cube \mathcal{C}_{ini} and the organized one \mathcal{C}_{org} . To measure the relevance of the organization provided by our method, we compute the gain $g = \frac{hc(\mathcal{C}_{org}) - hc(\mathcal{C}_{ini})}{hc(\mathcal{C}_{ini})}$ realized by the homogeneity criterion.

We also note that, for the same cube, its organized representation does not depend on the initial representation because the results of the MCA are insensitive to the order of input variables.

6. A Case Study

To test and validate our approach, we apply it on a 5-dimensional cube ($d = 5$) that we have constructed from the *Census-Income Database*¹ of the *UCI Knowledge Discovery in Databases Archive*². This data set contains weighted census data extracted from the 1994 and 1995 current population surveys conducted by the *U.S. Census Bureau*. The data contains demographic and employment related variables. The constructed cube contains 199 523 facts and one fact represents a particular profile of a sub population measured by the *Wage per hour*. The dimensions of the cube are : *Education level* (D_1 ,

¹<http://kdd.ics.uci.edu/databases/census-income/census-income.html>

²<http://kdd.ics.uci.edu/>

$p_1 = 17$), *Professional category* (D_2 , $p_2 = 22$), *State of residence* (D_3 , $p_3 = 51$), *Household situation* (D_4 , $p_4 = 38$), and *Country of birth* (D_5 , $p_5 = 42$).

We generate a complete disjunctive table $Z = [Z_1, Z_2, Z_3, Z_4, Z_5]$ according to a binary coding of the cube dimensions. Z contains 199523 rows and $p = \sum_{t=1}^5 p_t = 170$ columns. By applying the MCA on Z we obtain $p - d = 165$ factorial axes F_α . Each axis is associated to an eigenvalue λ_α . Suppose that, according to the histogram of eigenvalues, a user chooses the three first axes ($k = 3$). These axes explain 15.35% of the total inertia of the facts cloud. This contribution does not seem very important at a first sight. But we should note that in a case of a uniform distribution of eigenvalues, we get normally a contribution of $\frac{1}{p-d} = 0.6\%$ per axis, i.e. the three first axes represent an inertia already 25 times more important than a uniform distribution.

j	Attributes	Test-values		
		V_{1j}^1	V_{2j}^1	V_{3j}^1
9	<i>Hospital services</i>	-99.90	-99.90	-99.90
14	<i>Other professional services</i>	-99.90	-99.90	99.90
17	<i>Public administration</i>	-99.90	-99.90	99.90
12	<i>Medical except hospital</i>	-99.90	99.90	-99.90
5	<i>Education</i>	-99.90	99.90	99.90
7	<i>Finance insurance</i>	-99.90	99.90	99.90
19	<i>Social services</i>	-99.90	99.90	99.90
8	<i>Forestry and fisheries</i>	-35.43	-8.11	83.57
3	<i>Communications</i>	-34.05	-99.90	99.90
15	<i>Personal services except private</i>	-21.92	-5.50	10.28
13	<i>Mining</i>	-6.59	-99.64	-5.25
16	<i>Private household services</i>	7.77	51.45	11.68
6	<i>Entertainment</i>	40.04	99.90	96.23
1	<i>Agriculture</i>	68.66	3.39	-27.38
4	<i>Construction</i>	99.90	-99.90	-99.90
10	<i>Manufact. durable goods</i>	99.90	-99.90	-99.90
11	<i>Manufact. nondurable goods</i>	99.90	-99.90	-99.90
21	<i>Utilities and sanitary services</i>	99.90	-99.90	-99.90
22	<i>Wholesale trade</i>	99.90	-99.90	-24.37
20	<i>Transportation</i>	99.90	-99.90	99.90
18	<i>Retail trade</i>	99.90	99.90	-99.90
2	<i>Business and repair</i>	99.90	99.90	99.90

Table 1. Attribute's test-values of *Professional category* dimension.

The organized *Census-Income* data cube is obtained by sorting the attributes of its dimensions. For each dimension D_t its attributes are sorted by the increasing values of V_{1j}^t , then by V_{2j}^t , and then by V_{3j}^t . Table 1 shows the new attributes' order of the *Professional category* dimension (D_2). Note that j is the index of the original alphabetic order of the attributes. This order is replaced by a new one according to the sort of test-values. In the Figures 4(a) and 4(b), we can clearly see the visual effect of this arrangement of attributes. These figures display views of data by crossing the *Professional category* dimension on columns (D_2) and the *Country of birth* dimension on rows (D_5). The representation 4(a) displays the initial view according to the alphabetic order of attributes, whereas representation 4(b) displays the same view where attributes are rather sorted according to their test-values.

Remember that the aim of our current approach is not to compress or reduce the dimensions of a data cube. We do not also reduce sparsity of a data representation. Nevertheless, we act on this sparsity and reduce its negative effect on OLAP interpretation. Thus, we arrange differently original facts within a visual effect that gathers them as

well as possible in the space representation of the data cube. At a first sight, the visual representation 4(b) is more suitable to interpretation than 4(a). We clearly distinguish in Figure 4(b) four dense regions of full cells. In this regions, the homogeneity is higher than the rest of the space representation of the data cube.

This is confirmed by the measure of homogeneity criterion. Indeed, for a sparsity ratio of 63.42%, the homogeneity criterion for the organized cube of representation 4(b) is $hc(C_{org}) = 0.17$; whereas it measures $hc(C_{ini}) = 0.14$ for the initial cube of representation 4(a), i.e, we release a gain $g = 17.19\%$ of homogeneity when arranging the attributes of the cube according to test-values.

According to the test of the Equation (4), for each $t \in \{1, \dots, 5\}$, we select from D_t the set of characteristic attributes for the three selected factorial axes. These characteristic attributes give the best semantic interpretation of factorial axes and express strong relationships for their corresponding facts. To avoid great number of possible characteristic attributes per axis, we can consider, for each axis, only the first 50% of attributes having the highest absolute test-values. For instance, in the *Professional category* dimension D_2 , the set Φ_2 of characteristic attributes correspond to grayed rows in table 1.

In the same way, we apply the test of the Equation (4) on the other dimensions of the cube. In the representation of Figure 4(b), we clearly see that the zones of facts corresponding to characteristic attributes of the dimensions D_2 and D_5 seem to be more interesting and denser than other regions of the data space representation. These zones contains relevant information and reflect interesting association between facts. For instance, we can easily note that industrial and physical jobs, like construction, agriculture and manufacturing are highly performed by *Native Latin Americans* from Ecuador, Peru, Nicaragua and Mexico for example. At the opposite, *Asians* people from India, Iran, Japan and China are rather concentrated in commerce and trade.

7. Experimental Results

We have realized some experiments on the *Census-Income* data cube presented in section 6. The aim of these experiments is to appreciate the efficiency of our approach by measuring the homogeneity gain realized by our MCA-based organization on data representations with different sparsity ratios. To vary sparsity we proceeded by a random sampling on the initial dataset of the 199 523 facts from the considered cube.

According to Figure 5, the homogeneity gain has an increasing general trend. Nevertheless, we should note that for low sparsity ratios, the curve is rather oscillating around the null value of the homogeneity gain. In fact, when sparsity is less then 60%, the gain does not have a constant variation. It sometimes drops to negative values. This means that our method does not bring a value added to the quality of data representation. For dense data cubes, the employment of our method is not always significant. This is naturally due to the construction of the homogeneity criterion which closely depends on the number of empty and full cells. It can also be due to the structure of the random data samples that can generate data representations already having good qualities and high homogeneity values.

Our MCA-based organization method is rather interesting for data representations with high sparsity. In Figure 5, we clearly see that curve is rapidly increasing to high positive values of gain when sparsity is greater than 60%. Actually, with high relative

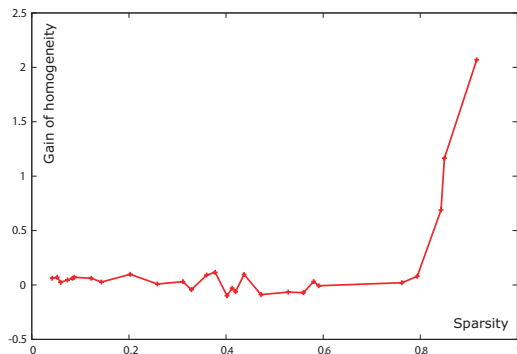


Figure 5. Evolution of the homogeneity gain according to sparsity.

number of empty cells in a data cube, we have a large manoeuvre margin for concentrating similar full cells and gathering them in the space representation. This shows the vocation of using our approach in order to enhance the visual quality representation, and thus the analysis of huge and sparse data cubes.

8. Conclusion and Future Work

In this paper, we introduced a MCA-based approach to enhance the space representation of large and sparse data cubes. This approach aims to provide an assistance to the OLAP user and helps him/her to easily explore huge volumes of data. For a given data cube, we compute the test-values of its attributes. According to these test-values, we arrange attributes of each dimension and so display in an appropriate way the space representation of facts. This representation provides better property for data visualization since it gather full cells expression interesting relationships of data. We also identify relevant regions of facts in this data representation by detecting characteristic attributes of factorial axes. This solve the problem of high dimensionality and sparsity of data and allows the user to directly focus his exploration and data interpretation on these regions. We have also proposed an homogeneity criterion to measure the quality of data representations. This criterion is based on the notion of geometric neighborhood of cells and their measures' similarities. Through experiments we led on real world data, our criterion proved the efficiency of our approach for huge and sparse data cubes.

Currently, we are studying some possible extensions for this work. We consider the problem of optimizing complexity of our approach. We also try to involve our approach in order to take into account the issue of data updates. Finally, we project to implement this approach under a Web environment that offers an interesting on-line aspect and a good user interaction context.

References

- [1] D. Barbará and M. Sullivan. Quasi-Cubes: Exploiting Approximations in Multidimensional Databases. *SIGMOD Record*, 26(3):12–17, 1997.
- [2] R. Cattell. The Scree Test for the Number of Factors. *Multivariate Behavioral Research*, 1:245–276, 1966.

- [3] S. Chaudhuri and U. Dayal. An Overview of Data Warehousing and OLAP Technology. *SIGMOD Record*, 26(1):65–74, 1997.
- [4] Y. W. Choong, D. Laurent, and P. Marcel. Computing Appropriate Representations for Multidimensional Data. *Data & knowledge Engineering Journal*, 45(2):181–203, 2003.
- [5] B. Escofier and B. Leroux. Etude de trois problèmes de stabilité en analyse factorielle. *Publication de l'Institut Statistique de l'Université de Paris*, 11:1–48, 1972.
- [6] J. Feng, Q. Fang, and H. Ding. PrefixCube: Prefix-sharing Condensed Data Cube. In *7th ACM International Workshop on Data warehousing and OLAP (DOLAP 2004)*, pages 38–47, Washington D.C., U.S.A., November 2004.
- [7] Y. Feng, D. Agrawal, A. E. Abbadi, and A. Metwally. Range CUBE: Efficient Cube Computation by Exploiting Data Correlation. In *20th International Conference on Data Engineering (ICDE 2004)*, pages 658–670, Boston, Massachusetts, U.S.A., March–April 2004.
- [8] W. H. Inmon. *Building the Data Warehouse*. John Wiley & Sons, 1996.
- [9] J.P. Benzecri. *Correspondence Analysis Handbook*. Marcel Dekker, hardcover edition, January 1992.
- [10] H. Kaiser. A note on Guttman's lower bound for the number of common factors. *Brit. J. Statist. Psychol.*, 14:1–2, 1961.
- [11] R. Kimball. *The Data Warehouse toolkit*. John Wiley & Sons, 1996.
- [12] L. V. Lakshmanan, J. Pei, and J. Han. Quotient Cube: How to Summarize the Semantics of a Data Cube. In *28th International Conference of Very Large Data Bases (VLDB 2002)*, Hong Kong, China, August 2002.
- [13] L. V. Lakshmanan, J. Pei, and Y. Zhao. QC-Trees: An Efficient Summary Structure for Semantic OLAP. In A. Press, editor, *ACM SIGMOD International Conference on Management of Data (SIGMOD 2003)*, pages 64–75, San Diego, California, U.S.A., 2003.
- [14] L. Lebart, A. Morineau, and M. Piron. *Statistique exploratoire multidimensionnelle*. Dunold, Paris, 3^e édition edition, 2000.
- [15] X. Li, J. Han, and H. Gonzalez. High-Dimensional OLAP: A Minimal Cubing Approach. In *30th International Conference on Very Large Data Bases (VLDB 2004)*, pages 528–539, Toronto, Canada, August 2004.
- [16] E. Malinvaud. Data Analysis in Applied Socio-Economic Statistics with Special Consideration of Correspondence Analysis. In *Marketing Science Conference*, Jouy en Josas, France, 1987.
- [17] R. B. Messaoud, O. Boussaid, and S. L. Rabaséda. Evaluation of a MCA-Based Approach to Organize Data Cubes. In *Proceedings of the 14th ACM International Conference on Information and Knowledge Management (CIKM'2005)*, pages 341–342, Bremen, Germany, October – November 2005. ACM Press.
- [18] R. B. Messaoud, O. Boussaid, and S. L. Rabaséda. Efficient Multidimensional Data Representation Based on Multiple Correspondence Analysis. In *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'2006)*, pages 662–667, Philadelphia, PA, USA, August 2006. ACM Press.
- [19] R. B. Messaoud, O. Boussaid, and S. L. Rabaséda. Using a Factorial Approach for Efficient Representation of Relevant OLAP Facts. In *Proceedings of the 7th International Baltic Conference on Databases and Information Systems (DB&IS'2006)*, pages 98–105, Vilnius, Lithuania, July 2006. IEEE Communications Society.
- [20] K. A. Ross and D. Srivastava. Fast Computation of Sparse Datacubes. In *23rd International Conference on Very Large Data Bases (VLDB 1997)*, pages 116–125, Athens, Greece, August 1997.
- [21] J. Shanmugasundaram, U. M. Fayyad, and P. S. Bradley. Compressed Data Cubes for OLAP Aggregate Query Approximation on Continuous Dimensions. In *5th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (SIGKDD 1999)*, pages 223–232, San Diego, California, U.S.A., August 1999.
- [22] Y. Sismanis, A. Deligiannakis, N. Roussopoulos, and Y. Kotidis. Dwarf: Shrinking the PetaCube. In *ACM SIGMOD International Conference on Management of Data (SIGMOD 2002)*, pages 464–475, Madison, Wisconsin, U.S.A., 2002.
- [23] J. S. Vitter and M. Wang. Approximate Computation of Multidimensional Aggregates of Sparse Data Using Wavelets. In *ACM SIGMOD International Conference on Management of Data (SIGMOD 1999)*, pages 193–204, Philadelphia, Pennsylvania, U.S.A., June 1999. ACM Press.
- [24] W. Wang, H. Lu, J. Feng, and J. X. Yu. Condensed Cube: An Effective Approach to Reducing Data Cube Size. In *18th IEEE International Conference on Data Engineering (ICDE 2002)*, pages 155–165, San Jose, California, U.S.A., February–March 2002.