# Made available by Hasselt University Library in https://documentserver.uhasselt.be

An Algebra for OLAP

Peer-reviewed author version

KUIJPERS, Bart & VAISMAN, Alejandro (2017) An Algebra for OLAP. In: Intelligent Data Analysis, 21(5), p. 1267-1300.

DOI: 10.3233/IDA-163161 Handle: http://hdl.handle.net/1942/21221

# An Algebra for OLAP

Bart Kuijpers

Hasselt University, Belgium

Alejandro Vaisman\*

Instituto Tecnológico de Buenos Aires, Argentina

# Abstract

OLAP (On Line Analytical Processing) comprises tools and algorithms that allow querying multidimensional (MD) databases. OLAP is based on the MD model, where data can be seen as a *cube*, where each cell contains one or more measures of interest, that can be aggregated along *dimensions*. Despite the extensive corpus of work in the field, a formally defined, reference language for OLAP is still needed, as there is no well-defined, accepted semantics, for many of the usual OLAP operations. In this paper, we address the problem, and present a set of operators that manipulate a data cube, clearly define their semantics, and prove that they can be composed, yielding a language powerful enough to express complex OLAP queries. We express these operations as a sequence of atomic transformations over a fixed MD matrix whose cells contain a sequence of measures. Each atomic transformation produces a new measure. When a sequence of transformations forms an OLAP operation, additionally, a flag is produced that indicates which cells must be considered as input for the next operation. In this way, an elegant algebra is defined. Our main contribution, with respect to other similar efforts in the field is that, for the first time, a formal proof to practical problems is given, and we believe the present work will serve as a basis to build more solid practical tools for data analysis.

*Keywords:* OLAP, Data Warehousing, Algebra, Data Cube, Dimension Hierarchy

### 1. Introduction

5

OLAP (On Line Analytical Processing) [1] comprises a set of tools and algorithms that allow efficiently querying multidimensional (MD) databases containing large amounts of data, usually called Data Warehouses (DW). Conceptually, in the MD model, data can be seen as a *cube*, where each cell contains one or

<sup>\*</sup>Corresponding author

*Email addresses:* bart.kuijpers@uhasselt.be (Bart Kuijpers), avaisman@itba.edu.ar (Alejandro Vaisman)

more *measures* of interest, that quantify *facts*. Measure values can be aggregated along *dimensions*, which give context to facts. At the logical level, OLAP data are typically organized as a set of *dimension and fact tables*. Typically, current database technology allows alphanumerical warehouse data to be integrated for example, with geographical or social network data, for decision making. In the

era of so-called "Big Data", the kinds of data that could be handled by data management tools, are likely to increase in the near future.

Although OLAP and Business Intelligence (BI) tools allow to manage different kinds of information, this normally requires that the user must be aware of models and query languages appropriate for each data type incorporated in the process. For example, we may have alphanumerical data coming from a local DW, spatial data (e.g., temperature) coming as rasterized images, and economical data published on the semantic web. A user who needs to integrate all of these data for analysis would need to have a knowledge not only of some

- <sup>20</sup> OLAP language (or OLAP graphic tools), but also must know how to deal with spatial data, and even with SPARQL (the standard query language for the semantic web). Ideally, this user would just like to deal with what she knows well, namely the data cube, using only the classical OLAP operators, like *Roll-up*, *Drill-down*, *Slice*, and *Dice* (among other ones), regardless the cube's underly-
- <sup>25</sup> ing data type (as we explained, spatial discrete, continuous, a graph data type, or just alphanumerical). These data types should be handled only at the logical and physical levels, not at the conceptual level. This is the idea introduced by Ciferri et al. [2], who proposed a model independent of technologies like ROLAP (for Relational OLAP), MOLAP (for Multidimensional OLAP) or HOLAP (for
- <sup>30</sup> Hybrid OLAP), and has an associated query language based exclusively on the conceptual level, thus providing high-level query operations for the user. This language, called Cube Algebra, was sketched informally in that paper. Building on this algebra, extensive examples are presented in [3], suggesting that this idea can lead to a language much more intuitive and simple than the *de facto*
- <sup>35</sup> standard MDX [4]. Nevertheless, these works do not give (since this was not their goal) any evidence of the correctness of the languages and operations proposed, other than examples with various degrees of comprehensiveness. In fact, surprisingly, and in spite of the extensive corpus of work in the field, a formally defined, reference language for OLAP is still needed [5]. There is not even a well-defined, accepted semantics, for many of the usual OLAP operations.

#### 1.1. Contributions

45

10

In this paper we address the problem introduced above. To this end, we:

- introduce a collection of operators that manipulate a data cube, and clearly define their semantics;
- prove, formally, that our operators can be composed, yielding a language powerful enough to express complex queries and cube navigation ("à la OLAP") paths.

We achieve the above representing the data cube as a fixed d-dimensional matrix, and a set of k measures, and expressing each OLAP operation as a

- <sup>50</sup> sequence of atomic transformations. Each transformation produces a new measure, and, additionally, when a sequence forms an OLAP operation, a flag that indicates which are the cells that must be considered as input for the next operation. This formalism allows us to elegantly define an algebra as a collection of operations, whose proof of correctness we provide in the paper. In this paper
- <sup>55</sup> we limit ourselves to the most usual operators, namely slice, dice, roll-up and drill-down, which constitute the core of all practical OLAP tools. This allows us to focus on our main interest, which is, to prove the feasibility of the approach. Other not-so-usual operations, and operations between two or more cubes, are left for future work.
- The main contribution of our work, with respect to other similar efforts in the field is that, for the first time, a formal proof to practical problems is given, so the present work will serve as a basis to build more solid tools for data analysis. As we show in the next section, existing work either lacks of formalism, or of applicability, and no work of any of these kinds give sound mathematical prove of its claims.

#### 1.2. Related Work

In spite that the need of an algebra for OLAP has long been acknowledged in the literature (see for example [5]), just a few works have addressed this problem so far, and in a limited way.

- The multidimensional model (MD) proposed by Gyssens and Lakshmanan [6] defines a data manipulation language that can express a so-called cube operator. The authors propose an algebra (and an equivalent calculus), which includes set operators (like selection, projection, cartesian product), operators for summarization, and re-structuring operators (fold and unfold). This model largely simplifies typical MD models for OLAP (for example, dimension hierarchies are
- rs simplifies typical MD models for OLAP (for example, dimension hierarchies are considered in a very limited way), and the operations proposed only address simple cases.

Along similar lines, Agrawal et al. [7] proposed a data model that supports multiple hierarchies along each dimension, and the possibility of performing ad-

hoc aggregates. They also define a minimal set of algebraic operators that is composed of the following operators: push and pull, destroy dimension, restriction (slice and dice), join, and associate. These operations are introduced in an informal way.

The proposals above are not appropriate for composing operations in practical cases, since they make implicit assumptions that do not apply in real world scenarios. On the contrary, we show that our approach can be applied to address typical operations composition.

Macedo and Oliveira [8] present an approach that can be considered close to ours, since they represent MD data as a matrix, with the idea of expressing OLAD are matrixed in linear shadow (LA). The many is matrixed in the idea of expressing

- <sup>90</sup> OLAP operations in linear algebra (LA). The paper's motivation is, like in our case, to fill the theoretical gap in the field. The proposal expresses some simple OLAP operations, mainly cross tabulations, as a combination of matrix multiplication, transposition, and a variant of the Kronecker product. However, this proposal is very preliminary, as the author's acknowledge, since no justification of the approach is provided. Further, as it is, the work is oriented to Excel
- <sup>95</sup> of the approach is provided. Further, as it is, the work is oriented to Excel

spreadsheets rather than to OLAP, and it has yet to be incorporated into a typical MD model for OLAP.

In a more OLAP-oriented approach, Vassiliadis [9] presented a classic MD model, which includes the concepts of dimensions, hierarchies, and cubes. The author also proposes a set of operators based on the notion of a base cube, e.g., a cube at the finest granularity level. These operations are: level climbing, packing, function application, projection, navigation, slicing, and dicing. Again, no formal language is presented. Ravat et al. [10] also propose an OLAP algebra at the conceptual level, trying to overcome their drawbacks, although no formal semantics is defined for the algebra presented by the authors.

Some works have already made use of the cube algebra proposed in [2]. Gómez el al. [11] used cube algebra to manipulate different kinds of spatial data, namely discrete data and continuous fields implemented in several different ways, like Voronoi diagrams and rasterized data. They implemented cube

- <sup>110</sup> algebra at a conceptual level, and all the needed machinery to manipulate the heterogeneous cubes at the logical and physical levels. Similar work, but with semantic web data, is presented in [12]. We envision that applications like these, are likely to grow in number and variety, for which a formalization of this algebra, like the one proposed in this paper, is clearly needed.
- <sup>115</sup> Many proposals also exist in the OLAP literature, defining several different sets of operators to handle MD data, although none of them abstract from the logical level and thus, do not provide high-level query operations for the user. For a comprehensive survey, we refer the reader to [2].

#### 1.3. Paper Organization

- <sup>120</sup> The remainder of the paper is organized as follows. In Section 2, we present our MD data model, on which we base the rest of our work. Section 3 presents the atomic transformations that we use to build the OLAP operations. In Section 4 we discuss the classical OLAP operations in terms of the transformations, show how they can be composed to address complex queries, and give proofs of
- all of our claims. We conclude in Section 5. Additional proofs are given in the appendix.

# 2. The OLAP data model

In this section, we describe the OLAP data model, which is the multidimensional data cube. Before we give the definition of the data cube, we define what we mean by its matrix. The "empty" matrix serves as a placeholder for the measures that are contained in the data cube. We also define the notions of dimension schema (with hierarchies and levels) and dimension instance (level instance, hierarchy instance and dimension graph). We end this section with a discussion of ordered domains and the representation of higher-level objects.

#### 135 2.1. Multidimensional Matrix

In this section, we give the definitions of a multidimensional matrix schema and a multidimensional matrix instance. In the following definition and throughout this paper, d is a natural number, with  $d \ge 1$ , which represents the number of dimensions of a data cube. Definition 1 (Matrix Schema). A d-dimensional matrix schema is a sequence

 $(D_1, D_2, ..., D_d)$  of d dimension names.

140

145

Dimension names can be considered to be strings. As illustrated in the following example, we use the notational convention to start dimension names with a capital letter.

**Example 1.** Our running example deals with sales information of certain products, at certain locations, at certain moments in time. For this purpose, we use the 3-dimensional matrix schema  $(D_1, D_2, D_3) = (Product, Location, Time)$ .

**Definition 2 (Matrix Instance).** A *d*-dimensional matrix instance (or a matrix, for short) over the *d*-dimensional matrix schema  $(D_1, D_2, ..., D_d)$  is a product

$$dom(D_1) \times dom(D_2) \times \cdots \times dom(D_d),$$

where, for i = 1, 2, ..., d,  $dom(D_i)$  is a non-empty, finite, ordered set, called the domain, that is associated with the dimension name  $D_i$ . For all i = 1, 2, ..., d, we denote by <, the order that we assume on the elements of  $dom(D_i)$ . For  $a_1 \in dom(D_1), a_2 \in dom(D_2), ..., a_d \in dom(D_d)$ , we call the tuple  $(a_1, a_2, ..., a_d)$ a *cell* of the matrix.

The cells of a matrix serve as placeholders for the measures that are contained in the data cube (see Definition 7). The role of the order < is discussed further in Section 2.4.

We use the notational convention to start elements of the domains  $dom(D_i)$  with a lower case letter, as it is illustrated in the following example.

**Example 2.** For the 3-dimensional matrix schema  $(D_1, D_2, D_3) = (Product, Location, Time)$  of Example 1, the non-empty sets  $dom(D_1) = \{lego, brio, apples, oranges\}, dom(D_2) = \{antwerp, brussels, paris, marseille\}, and <math>dom(D_3) = \{1/1/2014, ..., 31/1/2014\}$  give rise to the matrix instance

$$dom(D_1) \times dom(D_2) \times dom(D_3).$$

This matrix is shown in Figure 3. The cells of the matrix contain the sales figures for each combination of values in the domain. On  $dom(D_2)$ , we have, for

instance, the order antwerp < brussels < paris < marseille. On the dimension Time, we would typically have the temporal order.

2.2. Level Instance, Hierachy Instance and Dimension Graph

155

165

In this section, we define the notions of dimension schema (with hierarchies and levels) and dimension graph (or dimension instance).

Definition 3 (Dimension Schema, Hierarchy and Level). Let D be a dimension name. A dimension schema  $\sigma(D)$  for D is a lattice with a unique topnode, called All (which has only incoming edges) and a unique bottom-node, called Bottom (which has only outgoing edges), such that all maximal-length paths in the graph go from Bottom to All.

Any path from *Bottom* to *All* in a dimension schema  $\sigma(D)$  is called a *hierarchy* of  $\sigma(D)$ . Each node in a hierarchy (or in a dimension schema) is called a *level* (of  $\sigma(D)$ ).

We use the notational convention to start level names with a capital letter. We remark that the *Bottom* node is often renamed, depending on the application, as is illustrated in the following example. This example also introduces a non-graphical notation for hierarchies.

Example 3. Figure 1 gives examples of dimension schemas  $\sigma(Location)$  and  $\sigma(Time)$  for the dimensions *Location* and *Time* from Example 1.

For the dimension *Location*, we have Bottom = City and there is only one hierarchy, which we denote as

$$City \rightarrow Region \rightarrow Country \rightarrow All.$$

The node *Region* is an example of a level in this hierarchy.

For the dimension Time, we have Bottom = Day and we have two hierarchies, namely  $Day \rightarrow Month \rightarrow Semester \rightarrow Year \rightarrow All$  and  $Day \rightarrow Week \rightarrow All$ .

We remark that for the dimension *Location*, we have a linear lattice as a dimension schema. In this example, this is not the case for the dimension *Time*.



Figure 1: Dimension schemas for the dimensions Location, in (a), and Time, in (b).

Definition 4 (Level Instance, Hierachy Instance, Dimension Graph). Let D be a dimension with schema  $\sigma(D)$ , and let  $\ell$  be a level of  $\sigma(D)$ . A level instance of  $\ell$  is a non-empty, finite set  $dom(D.\ell)$ . If  $\ell = All$ , then dom(D.All) is the singleton  $\{all\}$ . If  $\ell = Bottom$ , then dom(D.Bottom) is the the domain of the dimension D, that is, dom(D) (as in Definition 2).

A dimension graph (or instance)  $I(\sigma(D))$  over the dimension schema  $\sigma(D)$ is a directed acyclic graph with node set

170

$$\bigcup_{\ell} dom(D.\ell),$$

where the union is taken over all levels in  $\sigma(D)$ . The edge set of this directed acyclic graph is defined as follows. Let  $\ell$  and  $\ell'$  be two levels of  $\sigma(D)$ , and let  $a \in dom(D.\ell)$  and  $a' \in dom(D.\ell')$ . Then, only if there is a directed edge from  $\ell$  to  $\ell'$  in  $\sigma(D)$ , there can be a directed edge in  $I(\sigma(D))$  from a to a'.

If H is a hierarchy of  $\sigma(D)$ , then the *hierarchy instance* (relative to the dimension instance  $I(\sigma(D))$ ) is the subgraph of  $I(\sigma(D))$  with nodes from  $dom(D.\ell)$ , for  $\ell$  appearing in H. This subgraph is denoted by  $I_H(\sigma(D))$ .

We use the notational convention to start the names of objects from a set  $dom(D.\ell)$  with a lower case character.

We remark that a hierarchy instance  $I_H(\sigma(D))$  is always a (directed) tree, since a hierarchy is a linear lattice. We also use the following terminology. If aand b are two nodes in a hierarchy instance  $I_H(\sigma(D))$ , such that (a, b) is in the transitive closure of the edge relation of  $I_H(\sigma(D))$ , then we say that a rolls-up to b and we denote this by  $\rho_H(a, b)$  (or  $\rho(a, b)$  if H is clear from the context).

The following example illustrates these concepts.

180

185

**Example 4.** We continue with Example 3 and focus on the dimension *Location*, whose dimension schema,  $\sigma(Location)$ , is given in Figure 1 (a). From Example 2, we have  $dom(Location) = \{antwerp, brussels, paris, marseille\}$ , which is dom(Location.Bottom), or dom(Location.City). For the levels *Region* and Country, we have  $dom(Location.Region) = \{flanders, capital, north, south\}$ , and  $dom(Location.Country) = \{belgium, france\}$ , respectively. An example of a dimension instance  $I(\sigma(Location))$  is depicted in Figure 2. This example expresses, for instance, that the city *brussels* is located in the region *capital* which is part of the country *belgium*. This means that *brussels* rolls-up to *capital* and to *belgium*, that is,  $\rho(brussels, capital)$  and  $\rho(brussels, belgium)$ . We also remark that the dimension instance of Figure 2 is indeed a tree.



Figure 2: An example of a dimension graph (or instance)  $I(\sigma(Location))$ .

In a dimension graph with multiple hierarchies, elements in some levels may be reachable from elements in the *Bottom* level, in multiple ways. However, it is important that rolling-up in different ways gives the same results. This is formalised by the concept of "sound" dimension graph.

**Definition 5 (Sound Dimension Graph).** Let  $I(\sigma(D))$  be a dimension graph (as in Definition 4). We call this dimension graph *sound*, if for any level  $\ell$  in

 $\sigma(D)$  and any two hierarchies  $H_1$  and  $H_2$  that reach  $\ell$  from the *Bottom* level and any  $a \in dom(D)$  and  $b_1, b_2 \in dom(D.\ell)$ , we have that  $\rho_{H_1}(a, b_1)$  and  $\rho_{H_2}(a, b_2)$ imply that  $b_1 = b_2$ . 

In this paper, we assume that dimension graphs are always sound (as specified in Definition 7).

#### 2.3. Multidimensional Data Cube

In this section, we give the definitions of a (multidimensional) data cube 190 schema and a data cube instance. Essentially, a data cube is a matrix in which the cells are filled with measures that are taken from some value domain  $\Gamma$ . For many applications,  $\Gamma$  will be the set of real or rational numbers. But we may also think of applications where  $\Gamma$  includes spatial regions or other geometric 195 objects, for instance.

Definition 6 (Data Cube Schema). A d-dimensional data cube schema consists of

205

- a d-dimensional matrix schema  $(D_1, D_2, ..., D_d)$ ; and
- a hierarchy schema  $\sigma(D_i)$  for each dimension  $D_i$ , with i = 1, 2, ..., d.

**Definition 7 (Data Cube Instance).** Let  $\Gamma$  be a non-empty set of "values".

- A d-dimensional, k-ary data cube instance (or data cube, for short)  $\mathcal{D}$  over the 200 d-dimensional matrix schema  $(D_1, D_2, ..., D_d)$  and hierarchy schemas  $\sigma(D_i)$  for  $D_i$ , for i = 1, 2, ..., d, with values from  $\Gamma$ , consists of
  - a d-dimensional matrix instance over the matrix schema  $(D_1, D_2, ..., D_d)$ , denoted  $M(\mathcal{D})$ ;
  - for each i = 1, 2, ..., d, a sound dimension graph  $I(\sigma(D_i))$  over  $\sigma(D_i)$ ;
    - k measures  $\mu_1, \mu_2, ..., \mu_k$ , which are functions from  $dom(D_1) \times dom(D_2) \times$  $\cdots \times dom(D_d)$  to the value domain  $\Gamma$ ; and
    - a flag  $\varphi$ , which is a function from  $dom(D_1) \times dom(D_2) \times \cdots \times dom(D_d)$ to the set  $\{0, 1\}$ .

For the remainder of this paper, we assume that  $\Gamma = \mathbf{Q}$ , the set of the rational numbers. For most applications, this suffices. Also, as a notational convention, we use calligraphic characters, like  $\mathcal{D}$ , to represent data cube instances.

The flag  $\varphi$  can be considered as a (k + 1)-st measure that is Boolean. The role of  $\varphi$  is to indicate which of the matrix cells are currently "active". The active cells have a flag value 1 and the others have a flag value 0. When we operate over a data cube, flags are used to indicate the input or output parts of the matrix of the cube. Typically, in the beginning of the operations, all cells have a flag value of 1. The role of flags will become more clear in the next sections, when we discuss OLAP transformations and operations.



Figure 3: An example of a data cube with one measure:  $\mu_1 = sales$ .

**Example 5.** We build on the previous examples. Figure 3 shows a 3-dimenional 1-ary data cube instance over the matrix schema (*Product*, *Location*, *Time*) and dimension schema  $\sigma(Product)$ ,  $\sigma(Location)$ , and  $\sigma(Time)$  (two of which were given in Example 3) with the set of the rational numbers as value domain. The matrix cells contain one measure, namely  $\mu_1 = sales$ , which expresses the sales amount per product, per location and per time instant. Initially, the flag  $\varphi$  may be, for instance, 1 for all matrix cells (not indicated in Figure 3), telling that all cells of the matrix are currently active.

#### 2.4. Ordered domains and the representation of higher-level objects

In the process of performing OLAP transformations and operations, we may need to store aggregate information about certain measures at some level above the *Bottom* level. We do nor foresee extra space for this in the data cube. We use the available cells of the original data cube to store this aggregate information,

210

215

yielding a more elegant solution, since this allows us to manipulate always the same cube schema while we perform a sequence of operations over the cube, as we will see later.

225

Recall that in Definition 2 we have assumed an order < for the domains  $dom(D_i)$ . We make use of this order for the representation of high-level objects by *Bottom*-level objects. The following definition specifies how this is achieved.

**Definition 8.** Let  $D \in \{D_1, D_2, ..., D_d\}$  be an arbitrary dimension with domain dom(D) = dom(D.Bottom). Let  $\ell$  be a level of  $\sigma(D)$ . An element  $b \in dom(D.\ell)$  is *represented* by the smallest element  $a \in dom(D)$  (according to <) for which  $\rho(a, b)$ . We denote this as rep(b) = a and say that a represents b.

We remark, that since we assume dimension graphs to be sound, this notion of representation is well defined, because the smallest element of the bottom level will always reach the same element in any level, regardless of the path traversed.

The following example illustrates the concept of representation.

**Example 6.** Continuing the previous examples, we consider the dimension Location with  $dom(Location) = \{antwerp, brussels, paris, marseille\}$ , which is dom(Location.City). On this set, we assume the order antwerp < brussels < paris < marseille. For this dimension, we have the hierarchy  $City \rightarrow Region \rightarrow Country \rightarrow All$ , and we consider the dimension instance  $I(\sigma(Location))$ , given in Figure 2.

At the Bottom = City level, cities represent themselves. On higher levels, regions and countries are represented by their "first" city in dom(Location)(according to <). This means that *flanders* and *belgium* are represented by *antwerp*, *france* is represented by *paris*, while *south* is represented by *marseille*. Let us explain further explain this. For the level *Region*, we have  $dom(Location.Region) = \{flanders, capital, north, south\}$ . At this level, *antwerp* represents *flanders* and *marseille* represents *south*, since they are the first (and, in this case only) domain elements that roll-up to these regions. So, we have rep(flanders) = antwerp. For the level *Country*, we have  $dom(Location.Country) = \{belgium, france\}$ . At this level *antwerp* represents *belgium* and *paris* represents *france*, since they are the first (but not only) domain elements that roll-up to these countries. At the level All, we have *antwerp* that represents *all*.

and a representes and

Later, we use this convention, to encode outputs of OLAP transformations and operations at different levels. That means, in our running example, if we want to represent an output at the *Country* level, we will flag *antwerp* and *paris* to represent *belgium* and *france*. The idea is to store aggregate information for higher-level objects in the cells of their *Bottom*-level representatives. In an output cube that contains this aggregate information, we have these representatives

flagged 1 and other cells flagged 0.

**Remark 1.** However, there remains a problem, as the previous example illustrates. In this example, if we have aggregate information at the level *Region*, with  $dom(Location.Region) = \{flanders, capital, north, south\}$ , then all cities of  $dom(Location) = \{antwerp, brussels, paris, marseille\}$  are flagged. At this point, it would not be clear if the cube contains information at the level *City* or at the level *Region*. If we keep a log of the OLAP operations that are performed, this log makes the level of aggregation clear.

The following property shows how the order on the *Bottom* level induces and order on higher levels. This property depends on the soundness of the dimension graph. Its proof is straightforward and we omit it.

**Property 1.** Let  $D \in \{D_1, D_2, ..., D_d\}$  be a dimension of a data cube  $\mathcal{D}$  and let  $\ell$  be a level in the dimension schema  $\sigma(D)$ . The order < on dom(D) induces an order (also denoted <) on  $dom(D.\ell)$  as follows. If  $b_1, b_2 \in dom(D.\ell)$ , then  $b_1 < b_2$  if and only if  $rep(b_1) < rep(b_2)$ .

#### 250 3. OLAP transformations and operations

A typical OLAP user manipulates a data cube by means of well-known operations. Just to mention the most popular ones, *Roll-Up* aggregates measures up to a certain level in a dimension, *Drill-Down* disaggregates measures up to a certain level in a dimension, *Slice* drops a whole dimension, and *Dice* keeps only the cells in a cube satisfying a certain Boolean condition. These operations, which we will formally define later in this paper, actually express queries over the data cube, usually submitted using some graphic tool, and translated into an underlying query language. The result is typically displayed in graphic or tabular format. An OLAP query can then be considered as a sequence of these individual operations, which receive a cube as input, and return a cube as output. For instance, using our running example, an apparently simple query

like "Total sales by region, for regions in Belgium or France", is actually expressed as a sequence of operations, whose semantics should be clearly defined, and which can be applied in different order (with the same result). For exam-

ple, we can first apply a *Roll-Up* to the *Country* level, and once at that level apply a *Dice* operation, which keeps the tuples corresponding to Belgium or France. Finally, a *Drill-Down* disaggregates the sales down to the level *Region*, returning the desired result. Note that since the sales not occurred in Belgium and France have been eliminated, this last operation must only consider the remaining members in *Country*. Thus, the *Drill-Down* operation is not a just

- an undo of the previous *Roll-Up*, as it is sometimes considered to be. Note that, in practice, this problem appears regardless of the processing type, that is, whether the operation sequence is submitted as an expression in a query language, or is processed during the user's navigation through a graphic tool.
- In what follows, we regard OLAP operations as the result of sequences of "atomic" OLAP transformations, which are measure-creating updates to a data cube. First, we give the definition of an OLAP transformation. Next, we show how these transformations can be combined into OLAP operations and how OLAP operations can be composed, along the lines explained at the beginning
  of this section. Finally, we give an overview of our arsenal of atomic OLAP transformations.

We start this section with an informal description of atomic OLAP transformations, OLAP operations and their composition.

#### 3.1. Introduction to OLAP transformations and operations

- An atomic OLAP transformation acts on a data cube instance, by adding a measure to the existing data cube measures. OLAP operations like the ones informally introduced above are defined, in our approach, as a sequence of transformations. The process of OLAP transformations starts from a given input data cube  $\mathcal{D}_{in}$ . We assume that this original data cube has k given measures  $\mu_1, \mu_2, ..., \mu_k$  (as in Definition 7). These k measures have a special status in the sense that they are "protected" and can never be altered (see Section 3.3). However, there is one exception to this protection. These original measures can be "destroyed" in some cells (see further on and Section 3.2), for instance, as the result of slice- or dice-operations, which are destructive by nature. Operations
- <sup>295</sup> of these types destroy the content of some matrix cells and remove even the protected measures in it.

Typically, the input-flag  $\varphi$  of the original data cube  $\mathcal{D}_{in}$  is set to 1 in every cell and signals that every cell of  $M(\mathcal{D}_{in})$  is part of the input cube.

On data cubes, atomic OLAP transformations can be applied. They add (or create) new measures to the sequence of existing measures by adding new measure values in each cell of the data cube's matrix. At any moment in this process, we may assume that the data cube  $\mathcal{D}$  has k + l measures  $\mu_1, \mu_2, ..., \mu_k; \tau_1, ..., \tau_l$ , where the first k are the original measures of  $\mathcal{D}_{in}$  and where the last l (with  $l \geq 0$ ) measures have been created subsequently by l OLAP transformations. A next OLAP transformation adds a new measure  $\tau_{l+1}$  to the matrix cells.

We have said that we use OLAP transformations to compute OLAP operations. In this sense, we can see that many of the measures  $\tau_1, \tau_2, ..., \tau_l$  added by

the transformations in a process, may represent the result of intermediate computations that are not really relevant to the output of an OLAP operation. We indicate that the computation of an OLAP operation O is finished by creating a 310 *m*-ary output flag  $\varphi_{O}^{(m)}$ . This output flag is a Boolean measure, that is created like other measures via atomic OLAP transformations. It indicates which of the cells of  $M(\mathcal{D})$  should be considered as belonging to the output of O. It is m-ary in the sense that it keeps the last m created measures  $\tau_{l-m+1}, \tau_{l-m+2}, ..., \tau_l$  and "trashes"  $\tau_1, \tau_2, ..., \tau_{l-m}$ . It also removes the previous flag, which it replaces. 315 The initial measures  $\mu_1, \mu_2, ..., \mu_k$  of the input data cube  $\mathcal{D}_{in}$  are never removed (unless they are "destroyed" in some cells). They are "protected" and remain in the cube throughout the process of applying one OLAP operation after another to  $\mathcal{D}_{in}$ . So, at any stage, we can use the given measures  $\mu_1, \mu_2, ..., \mu_k$  (except in destroyed cells). 320

Summarizing the above, after an OLAP operation of output arity m is completed on some cube  $\mathcal{D}$ , the measures in the cells of the output data cube  $\mathcal{D}' = O(\mathcal{D})$  are of the form

$$\mu_1, \mu_2, \dots, \mu_k; \tau_{l-m+1}, \tau_{l-m+2}, \dots, \tau_l; \varphi_O^{(m)}$$

In the previous expression, the underlining indicates the protected status of these measures. After each OLAP operation, we do a "cleaning" by renaming the unprotected measures with the symbols  $\tau_1, \tau_2, ..., \tau_m$  and the output measures become

$$\mu_1, \mu_2, ..., \mu_k; \tau_1, \tau_2, ..., \tau_m; \varphi_Q^{(m)}$$

A next OLAP operation O' can then act on  $\mathcal{D}'$  and use in its computation all the measures  $\underline{\mu_1, \mu_2, ..., \mu_k}; \tau_1, \tau_2, ..., \tau_m; \varphi_O^{(m)}$ . When O' finishes its computation after adding  $\overline{l'}$  measures  $\tau_{m+1}, \tau_{m+2}, ..., \tau_{m+l'}$  and producing a m'-ary output, the new measures in the cells will look like

$$\mu_1, \mu_2, \dots, \mu_k; \tau_{m+l'-m'+1}, \tau_{m+l'-m'+2}, \dots, \tau_{m+l'}; \varphi_{O'}^{(m')}$$

The last m' measures before the flag are renamed  $\tau_1, \tau_2, ..., \tau_{m'}$ , again. In this way, the composition of OLAP operations should be viewed.

We remark that the dimensions, the hierarchy schemas and instances of  $\mathcal{D}$  remain unaltered during the entire OLAP process.

We end this description of our view of OLAP transformations, OLAP operations and their composition, with a remark on *destructors*. Destructors are similar to flags, in the sense that they are computed by some sequence of atomic OLAP transformations and that they are Boolean. A destructor, optionally, precedes the creation of an output flag. A destructor  $\delta$  takes the value 1 for some cells of the matrix of a data cube, and 0 on other cells. When  $\delta$  is invoked (and activated by the output flag that follows it) on a data cube  $\mathcal{D}$  with measures  $\mu_1, \mu_2, ..., \mu_k; \tau_1, \tau_2, ..., \tau_m$  and flag  $\varphi_O^{(m)}$ , it empties all cells for which the value of the destructor  $\delta$  is 0 by removing all measures from them, even the protected ones, thereby effectively "destroying" these cells. This is the only case where the protected measures are altered. For example, this happens when the OLAP operation is a slice or a dice. Operations of this type destroy part of the cube and make them inaccessible for further use. In this context, the output of a destructive operation O looks like

$$\mu_1, \mu_2, ..., \mu_k; \tau_1, \tau_2, ..., \tau_l; \delta; \varphi_O^{(m)},$$

in which the destructor precedes the output flag. The effect of the presence of a destructor is the following. A cell such that  $\delta = 0$  is emptied, after which it contains no more measures and flag. For cells with  $\delta = 1$ , the sequence of measures  $\mu_1, \mu_2, ..., \mu_k; \tau_1, \tau_2, ..., \tau_l; \delta; \varphi_O^{(m)}$ ; is transformed to  $\mu_1, \mu_2, ..., \mu_k; \tau_{l-m+1}, \tau_{l-m+2}, ..., \tau_l; \varphi_O^{(m)}$ ; which is renamed as  $\mu_1, \mu_2, ..., \mu_k; \tau_1, \tau_2, ..., \tau_m; \varphi$ ; before the next transformation takes place. This transformation will act, cell per cell, on the matrix of a cube, with the understanding that it does nothing with emptied cells. That is, no new measure can ever be added to a destroyed cell.

#### 3.2. OLAP transformations

345

The following definition specifies how an OLAP transformation acts on a data cube. The atomic OLAP operations that appear in this definition are specified further on in this section.

**Definition 9 (OLAP Transformation).** Let  $\mathcal{D}$  be a *d*-dimensional, (k + l)ary data cube instance with given (or protected) measures  $\mu_1, \mu_2, ..., \mu_k$ , created <sup>350</sup> measures  $\tau_1, ..., \tau_l$  (with  $l \ge 0$ ) and flag  $\varphi$  over some value domain  $\Gamma$ . An *OLAP transformation* T, applied to  $\mathcal{D}$ , results in the creation of a new measure  $\tau_{l+1}$  in  $\mathcal{D}$ . The transformation T adds the measure  $\tau_{l+1}$  to non-empty cells of  $M(\mathcal{D})$ . The measure  $\tau_{l+1}$  is produced from

- $\mu_1, \mu_2, ..., \mu_k$  (in non-empty cells);
- $\varphi$  (in non-empty cells);
  - $\tau_1, \tau_2, ..., \tau_l$  (in non-empty cells) and
  - the hierarchy schemas and instances of  $\mathcal{D}$

and belongs to one of the following classes:

- 1. Arithmethic transformations (see Definition 11);
- <sup>360</sup> 2. Boolean transformations (see Definition 12);
  - 3. Selectors (see Definition 13);
  - 4. Counting, sum, and min-max (see Definitions 14 and 19);

5. Grouping (see Definition 18);

An OLAP transformation can also result in the creation of a measure that is an output flag  $\varphi^{(m)}$  or arity m. This should be a measure with a Boolean value and to indicate that it is a flag of arity m, we use the reserved symbol  $\varphi^{(m)}$ instead of  $\tau_{l+1}$ .

An output flag  $\varphi^{(m)}$  may (optionally) be preceded by a destructor  $\delta$  (which is created following the same rules as for other measures, but which has a special status, expressed by the reserved symbol  $\delta$ ). This should be a measure with a Boolean value (to indicate which cells are desroyed). We use the reserved symbol  $\delta$  instead of  $\tau_{l+1}$ .

The effect of output flags and destructors is discussed in Section 3.1. We remark that atomic OLAP transformations update the cells of the matrix  $M(\mathcal{D})$ <sup>370</sup> cell per cell and that empty cells of  $M(\mathcal{D})$  are unaffected by transformations.

#### 3.3. OLAP operations and their composition

Before we give the definition of an OLAP operation, we describe the *input* to the OLAP process, which may involve multiple OLAP operations. This input is a *d*-dimensional, *k*-ary data cube instance  $\mathcal{D}_{in}$ , with measures  $\mu_1, \mu_2, ..., \mu_k$  and flag  $\varphi$ . These measures are *protected* in the sense that they remain the first *k* measures throughout the entire OLAP process and are never altered or removed unless they are destroyed in some cells. The cube  $\mathcal{D}_{in}$  has also a Boolean flag  $\varphi$ , which typically is 1 in all of the cells of  $M(\mathcal{D}_{in})$ , indicating that all the matrix cells are relevant for the input. So, the measures of the input cube  $\mathcal{D}_{in}$  are denoted as follows:

#### $\underline{\mu_1, \mu_2, \dots, \mu_k}; \varphi.$

After applying some OLAP operations to  $\mathcal{D}_{in}$ , we obtain a data cube  $\mathcal{D}$ . We refer to  $\mathcal{D}$  in the following definition.

**Definition 10.** Let  $\mathcal{D}$  be a *d*-dimensional, (k+l)-ary *input* data cube instance with given measures  $\mu_1, \mu_2, ..., \mu_k$ , computed measures  $\tau_1, ..., \tau_l$  and flag  $\varphi$ . The data cube  $\mathcal{D}$  acts as the input of an OLAP operation O (of arity *m*), which consists of a sequence of *n* consecutive OLAP transformations that create the additional measures  $\tau_{l+1}, ..., \tau_{l+n}$ , followed by the creation of an *m*-ary flag  $\varphi_O^{(m)}$ (possibly preceded by a destructor  $\delta$ ). As the result of the creation of  $\varphi_O^{(m)}$ , the measures in the cells of the data cube are changed from

$$\mu_1, \mu_2, \dots, \mu_k; \tau_1, \dots, \tau_l; \varphi; \tau_{l+1}, \dots, \tau_{l+n}$$

$$\mu_1, \mu_2, \dots, \mu_k; \tau_{l+n-m+1}, \dots, \tau_{l+n}; \varphi_Q^{(m)},$$

which become

$$\mu_1, \mu_2, ..., \mu_k; \tau_1, ..., \tau_m; \varphi,$$

after renaming. The output cube  $\mathcal{D}' = O(\mathcal{D})$  has the same dimensions, hierarchy schemas and instances as  $\mathcal{D}$ , but has measures  $\mu_1, \mu_2, ..., \mu_k; \tau_1, ..., \tau_m; \varphi$ .

In the case where  $\varphi_O^{(m)}$  is preceded by a destructor  $\delta$ , the same procedure is followed, except for the cells of  $M(\mathcal{D})$  for which  $\delta$  takes the value 0. These cells of  $M(\mathcal{D})$  are emptied, contain no measures, and become inaccessible for future transformations or operations.

We remark that the output  $\mathcal{D}' = O(\mathcal{D})$  of the OLAP operation O on input  $\mathcal{D}$  can serve as input to a next OLAP operation. We illustrate the composition of two operations in Example 11 and other examples.

#### 3.4. Atomic OLAP transformations

In this section, we define our arsenal of atomic OLAP transformations, divided in five classes, as described in Definition 9.

In the remainder of this section, we use the following notational convention. For a measure  $\alpha$ , we write  $\alpha(x_1, x_2, ..., x_d)$  to indicate the value of  $\alpha$  in the cell  $(x_1, x_2, ..., x_d) \in dom(D_1) \times dom(D_2) \times \cdots \times dom(D_d)$ . We remark that

 $\alpha(x_1, x_2, ..., x_d)$  does not exist for empty cells and is therefore not considered in computations (such as sums). Also, we assume that we have protected measures  $\mu_1, \mu_2, ..., \mu_k$  and computed measures  $\tau_1, ..., \tau_l$  in the non-empty cells and that the next measure we compute is called  $\tau_{l+1}$ .

Throughout this section, we continue with the examples given in Section 2, that talk about the measure  $\mu_1 = sales$  of certain products, at certain locations, at certain moments in time, contained in a data cube  $\mathcal{D}$  over the 3-dimensional matrix schema  $(D_1, D_2, D_3) = (Product, Location, Time).$ 

3.4.1. Arithmethic transformations

Definition 11 (Arithmetic Transformations). The following creations of a

- <sup>395</sup> new measure  $\tau_{l+1}$  are arithmetic transformations:
  - 1. (Rational constant)  $\tau_{l+1} = \alpha$ , with  $\alpha \in \mathbf{Q}$ , a rational number.
  - 2. (Sum)  $\tau_{l+1} = \alpha + \beta$ , with  $\alpha, \beta \in \{\mu_1, \mu_2, ..., \mu_k, \tau_1, \tau_2, ..., \tau_l\}.$
  - 3. (**Product**)  $\tau_{l+1} = \alpha \cdot \beta$ , with  $\alpha, \beta \in \{\mu_1, \mu_2, ..., \mu_k, \tau_1, \tau_2, ..., \tau_l\}$ .

 $\operatorname{to}$ 

4. (Quotient)  $\tau_{l+1} = \alpha/\beta$ , with  $\alpha, \beta \in \{\mu_1, \mu_2, ..., \mu_k, \tau_1, \tau_2, ..., \tau_l\}$ . Here, we agree that a/0 := a for all  $a \in \mathbf{Q}$ .

**Example 7.** If  $\mu_1 = sales$  is the only measure, the following transformations 400 compute the 10% of the sales:

- $\tau_1 = 0.1$  (rational constant);
- $\tau_2 = \tau_1 \cdot \mu_1$  (product).

Next, if we want to create a Boolean measure that indicates whether a cell contains non-zero sales, we can write

405 •  $\tau_3 = \mu_1 / \mu_1$  (quotient).

The value of  $\tau_3$  is 1 if sales > 0 and 0 if sales = 0 (our definition of quotient says that 0/0 = 0).

3.4.2. Boolean transformations

**Definition 12 (Boolean Transformations).** The following creations of a new measure  $\tau_{l+1}$  are *Boolean transformations*:

1. (Equality test on measures)  $\tau_{l+1} = (\alpha = \beta)$ , with  $\alpha, \beta \in {\mu_1, \mu_2, ..., }$ 

 $\mu_k, \tau_1, \tau_2, ..., \tau_l$ . Here, the result of the comparison  $(\alpha = \beta)$  is a Boolean

- 1 or 0 (cell per cell in the non-empty cells of the matrix).
   2. (Comparison test on measures) τ<sub>l+1</sub> = (α < β), with α, β ∈ {μ<sub>1</sub>,
- $\mu_2, ..., \mu_k, \tau_1, \tau_2, ..., \tau_l$ . Here, the result of the comparison ( $\alpha < \beta$ ) is a Boolean 1 or 0 (cell per cell in the non-empty cells of the matrix).
- 3. (Equality test on levels) For  $\ell$  a level in the dimension schema  $\sigma(D_i)$  of dimension  $D_i$  and  $c \in dom(D_i.\ell)$  a constant object  $\tau_{l+1}(x_1, x_2, ..., x_d) =$  $(\ell = c)$  is an "equality" test. Here, the result of the comparison  $(\ell = c)$  is a Boolean 1 or 0 (cell per cell in the non-empty cells of the matrix) such that  $\tau_{l+1}(x_1, x_2, ..., x_d)$  is 1 if and only if  $x_i$  rolls-up to c at level  $\ell$ , that is  $\rho(x_i, c)$ .

4. (Comparison test on levels) For  $\ell$  a level in the dimension schema  $\sigma(D_i)$  of dimension  $D_i$  and  $c \in dom(D_i.\ell)$  a constant object,  $\tau_{l+1}(x_1, x_2, ..., x_d) = (\ell <_{\ell} c)$  is a "comparison" test. The result of the comparison  $(\ell <_{\ell} c)$  is a Boolean 1 or 0 (cell per cell in the non-empty cells of the matrix), such that  $\tau_{l+1}(x_1, x_2, ..., x_d)$  is 1 if and only if  $x_i$  rolls-up to an object b at level  $\ell$  for which  $b <_{\ell} c$ . The order  $<_{\ell}$  can be any order that is defined on level  $\ell$ . The transformation  $\tau_{l+1}(x_1, x_2, ..., x_d) = (c <_{\ell} \ell)$  is defined similarly.

We remark that the above equality test are superfluous since they can be expressed as a Boolean combination of comparison tests, but we include them for obvious practical reasons. Indeed, a = b is equivalent to  $\neg(a < b \lor b < a)$ . Finally, we remark that the comparison test on levels uses the order  $<_{\ell}$ , which may be the order (derived from) <, but, in practice, it will often be a lexicographical or alphabetical order, particular to the level domain.

**Example 8.** We illustrate the use of Boolean transformations by means of a sequence of transformations that implement a "dice" (see Section 4.2 for more details). The query

$$\mathsf{DICE}(\mathcal{D}, sales > 50)$$

- asks for the sales values in the matrix of  $\mathcal{D}$  which contain sales that are higher than 50. Again, we assume that  $\mu_1 = sales$  is the only available measure in the input cube. So, the measures in the cells are <u>sales</u>;  $\varphi$ . This query can be implemented by the following sequence of transformations:
  - $\tau_1 = 49.99$  (rational constant);
- $\tau_2 = (\tau_1 < sales)$  (comparison test on measures);
  - $\tau_3 = \mu_1 \cdot \tau_2$  (product);

- $\delta = \tau_2$  (destructor); and
- $\varphi^{(1)} = \tau_2$  (unary flag)

The measure  $\tau_3$  contains the *sales* values larger than or equal to 50 (and a 0 if the *sales* are lower). The flag  $\varphi^{(1)}$  selects all cells from the input as output cells and concludes the  $\mathsf{DICE}(\mathcal{D}, sales > 50)$  operation. The output of this operation is <u>sales</u>;  $\tau_3$ ;  $\varphi^{(1)}$ , which is then renamed to <u>sales</u>;  $\tau_1$ ;  $\varphi$ .

3.4.3. Selectors

- <sup>445</sup> **Definition 13 (Selector Transformations).** The following creations of a new measure  $\tau_{l+1}$  are selector transformations (or selectors) and their definition is (as always) cell per cell of  $M(\mathcal{D})$ :
  - 1. (Constant selector) For a level  $\ell$  in the dimension schema  $\sigma(D_i)$  of a dimension  $D_i$  and  $c \in dom(D_i.\ell)$ ,  $\tau_{l+1}$  can be a constant-selector for c, denoted  $\sigma_{D_i.\ell=c}$ , and it corresponds to the equality test on levels  $\tau_{l+1}(x_1, x_2, ..., x_d) = (\ell = c)$ .
  - 2. (Level selector) For a level  $\ell$  in the dimension schema  $\sigma(D_i)$  of a dimension  $D_i$ ,  $\tau_{l+1}$  can be a *level-selector for*  $\ell$ , denoted by  $\sigma_{D_i,\ell}$ , which means that we have, for all  $x_j \in dom(D_j)$  with  $j \neq i$ ,

$$\tau_{l+1}(x_1, \dots, x_{i_1}, a, x_{i+1}, \dots, x_d) = \begin{cases} 1 & \text{if } a = rep(b) \\ & \text{for some } b \in dom(D_i.\ell), \\ 0 & \text{otherwise.} \end{cases}$$

The constant selector in the Definition 13, corresponds to the equality test on levels (see 3. in Definition 12). Here, this transformation appears with a different functionality and we reserve a special notation for it. This is the reason why it was repeated. Also, we remark that the level selector selects all representatives (at the *Bottom* level) of objects at level  $\ell$  of dimension  $D_i$ .

**Example 9.** As a second example of a dice operation, we look at the query

$$\mathsf{DICE}(\mathcal{D}, Location.City = antwerp),$$

which asks for the sales in the city of *antwerp*. This operation is destructive, since it destroys all the information in cells that do not belong to *antwerp*. This query can be implemented by the following sequence of transformations:

450

- $au_1 = \sigma_{Location.City=antwerp}$  (constant selector);
  - $\tau_2 = \tau_1 \cdot \mu_1$  (product);
  - $\delta = \tau_1$  (destroys the cells outside *antwerp*);
  - $\varphi^{(1)} = \tau_1$  (unary flag creation).

The output arity of the query  $\mathsf{DICE}(\mathcal{D}, Location.City = antwerp)$  is 1. The measure  $\tau_2$  selects the sales in antwerp only. And the flag  $\varphi^{(1)}$  is a selector on the constant antwerp. The destructor  $\delta$ , that precedes the flag, empties the cells outside antwerp.

**Example 10.** As a next example, we look at the query

$$\mathsf{DICE}(\mathcal{D}, Location.City = antwerp \ OR \ Location.City = brussels),$$

which asks for the sales in the cities of *antwerp* and *brussels*. This query can be implemented by the following sequence of transformations:

- $\tau_1 = \sigma_{Location.City=antwerp}$  (constant selector);
- $\tau_2 = \sigma_{Location.City=brussels}$  (constant selector);
- $\tau_3 = \tau_1 + \tau_2 \text{ (sum)};$

470

- $\tau_4 = \tau_3 \cdot \mu_1$  (product);
- $\delta = \tau_3$  (destroys the cells outside *antwerp* and *brussels*);
  - $\varphi^{(1)} = \tau_3$  (unary flag creation).

The logical connective OR is implemented by the sum in  $\tau_3$ , which can take values 0 or 1, since the cities *antwerp* and *brussels* do not overlap. Thus, this sum implements their union. Then, measure,  $\varphi_4$  selects the sales in *antwerp* and *brussels* only. The flag  $\varphi^{(1)}$  is a selector on the constants *antwerp* and *brussels* and indicates that the cells of both these cities belong to the output. The destructor  $\delta$ , that precedes the flag, empties the cells outside *antwerp* and *brussels*. Note that in the two previous examples, the flag and the destructor do the same double work. However, this will not be the case in most situations, and, in practice, it would not have impact.

475

490

We continue with a dicing example, that also illustrates the *composition* of OLAP operations.

**Example 11.** We consider the query

 $\mathsf{DICE}(\mathcal{D}, sales > 50 \ AND \ Location.City = brussels).$ 

We can implement this by the operation  $\mathsf{DICE}(\mathcal{D}, sales > 50)$  followed by the operation  $\mathsf{DICE}(\mathcal{D}, Location.City = brussels)$ . The following implementation is a slight modification of Examples 8 and 9. Let <u>sales</u>;  $\varphi$  be the input measures.

The query  $\mathsf{DICE}(\mathcal{D}, sales > 50)$  is taken from Example 8. The output of this operation is <u>sales</u>;  $\tau_3$ ;  $\varphi^{(1)}$ , which is then renamed to <u>sales</u>;  $\tau_1$ ;  $\varphi$ . Next, the query  $\mathsf{DICE}(\mathcal{D}, Location, City = brussels)$  is implemented as

- $\tau_2 = \sigma_{Location.City=brussels}$  (constant selector);
- $\tau_3 = \tau_2 \cdot \mu_1$  (product);

•  $\delta = \tau_2$  (destroys the cells outside *brussels*);

•  $\varphi^{(1)} = \tau_2 \cdot \varphi$  (product and unary flag creation).

The output of this operation is <u>sales</u>;  $\tau_3$ ;  $\varphi^{(1)}$ , which, as above, is then renamed to <u>sales</u>;  $\tau_1$ ;  $\varphi$ .

Remark that the query can also be implemented as  $\mathsf{DICE}(\mathcal{D}, Location.City = brussels)$  followed by  $\mathsf{DICE}(\mathcal{D}, sales > 50)$ . Also in this order of operations, the appropriate cells are destroyed.

3.4.4. Counting, sum and min-max

Now, we give transformations for counting different measure values, for summing all values of a measure in a matrix, and for determining the minimum and maximum value of a measure in a matrix. Later on, in Definition 19, we give extensions of the counting and min-max transformations.

#### Definition 14 (Counting, Sum and Min-Max Transformations). The fol-

lowing creations of a new measure  $\tau_{l+1}$  are counting, sum and min-max transformations:

- 1. (Count-Distinct)  $\tau_{l+1} = \#_{\neq}(\alpha)$ , with  $\alpha \in \{\mu_1, \mu_2, ..., \mu_k, \tau_1, \tau_2, ..., \tau_l\}$  counts the number of different values of the measure  $\alpha$  in the complete matrix  $M(\mathcal{D})$  of the data cube.
- 2. (*d*-dimensional sum)

$$\tau_{l+1} = \sum_{(x_1, x_2, \dots, x_d) \in M(\mathcal{D})} \alpha(x_1, x_2, \dots, x_d),$$

with  $\alpha \in \{\mu_1, \mu_2, ..., \mu_k, \tau_1, \tau_2, ..., \tau_l\}$ , gives the sum of the measure  $\alpha$  over all non-empty matrix cells. We abbreviate this operation by writing

$$\tau_{l+1} = \mathrm{SUM}_d(\alpha)$$

and call this transformation the d-dimensional sum.

3. (min-max)  $\tau_{l+1} = \min(\alpha)$ , with  $\alpha \in \{\mu_1, \mu_2, ..., \mu_k, \tau_1, \tau_2, ..., \tau_l\}$ , gives the smallest value of the measure  $\alpha$  in non-empty cells of the matrix  $M(\mathcal{D})$ . Similarly,  $\tau_{l+1} = \max(\alpha)$ , gives the largest value of the measure  $\alpha$  in the matrix  $M(\mathcal{D})$ .

We remark that the above transformations create the same new measure value for all cells of the matrix  $M(\mathcal{D})$ .

We now give two examples of the use of *d*-dimensional sum. Examples of the use of  $\#_{\neq}(\alpha)$  are given in the following sections.

- **Example 12.** Let us consider the query "(grand total) average sales". This amount is the total sales (over all cities, products and dates), divided by the total number of cells in the matrix of the data cube. This query can be computed as follows, given  $\mu_1 = sales$ :
  - $\tau_1 = \text{SUM}_3(\mu_1)$  (this is the grand total of sales);

510

500

•  $\tau_2 = \sigma_{Location.All}$  (this puts a 1 in every cell of the matrix)

- $\tau_3 = \text{SUM}_3(\tau_2)$  (this is the grand total of cells in the matrix);
- $\tau_4 = \tau_1/\tau_3$  (this is the average);
- $\varphi^{(1)} = \sigma_{Location.Botom}$  (this flag creation selects all cells of the matrix).

The output measures are <u>sales</u>;  $\tau_4$ ;  $\varphi^{(1)}$ , which are renamed <u>sales</u>;  $\tau_1$ ;  $\varphi$ . This means that the grand total of average sales is now available in every cell of the matrix of the cube.

**Example 13.** Now, we look at the query "total sales in *antwerp*". This query <sup>515</sup> is asking for the total sales (over all products and dates) in the city of *antwerp*. The query can be computed as follows, given  $\mu_1 = sales$ :

- $\tau_1 = \sigma_{Location.City=antwerp}$  (constant selector on antwerp);
- $\tau_2 = \tau_1 \cdot \mu_1$  (product that selects the sales in *antwerp*)
- $\tau_3 = \text{SUM}_3(\tau_2)$  (this is the total sales in *antwerp* in every cell);
- $\tau_4 = \tau_3 \cdot \tau_1$  (this is the total sales in *antwerp* in the cells of *antwerp*);
  - $\varphi^{(1)} = \tau_1$  (this flag creation selects the cells of *antwerp*).

The output measures are <u>sales</u>;  $\tau_4$ ;  $\varphi^{(1)}$ , which are renamed <u>sales</u>;  $\tau_1$ ;  $\varphi$ . This means that the total of sales in *antwerp* is now available in every cell of *antwerp*. For the cells outside *antwerp* there is a 0. We remark that this example can be modified with a destructor that effectively empties cells outside *antwerp*.

**Example 14.** We look at the query "maximum sales", which should return all cells containing the maximum value in the cube. This query can be computed as follows, given  $\mu_1 = sales$ :

•  $\tau_1 = \max(\mu_1)$  (the maximum sales amount);

520

- $\tau_2 = (\tau_1 = \mu_1)$  (equality test to determine if a cell reaches the maximum);
- $\tau_3 = \tau_2 \cdot \mu_1$  (only the maximum sales remain; the others turn 0);
- $\varphi^{(1)} = \sigma_{Location.Bottom}$  (this flag creation selects all cells).

The output measures are <u>sales</u>;  $\tau_3$ ;  $\varphi^{(1)}$ , which are renamed <u>sales</u>;  $\tau_1$ ;  $\varphi$ . We remark that this example, like the previous one, can be modified with a destructor that effectively empties cells with strictly less than maximum sales.

#### 3.4.5. Grouping

The most common OLAP operations (e.g., roll-up, slice), require grouping data before aggregating them. For example, typically we will ask queries like "total sales by city", which requires grouping facts by city, and, for each group, sum all of its sales; or, we can ask "total sales by city and day", meaning that, for each city-day combination, we sum all the sales. Therefore, we need a transformation to express "grouping". We address this issue next.

To deal with grouping, we use the concept of "prime labels" for sets and products of sets. Before giving the definition of the grouping transformations, we elaborate on prime labels and product of prime labels. As we show, these prime labels work in the context of measures that take rational values (as, in practice, is often the case).

The following definition specifies our infinite supply of prime labels.

**Definition 15 (Prime Labels).** Let  $p_n$  denote the *n*-th prime number, for  $n \geq 1$ . We define the sequence of *prime labels* as follows:  $1, \sqrt{2}, \sqrt{3}, \sqrt{5}, \sqrt{7}, \sqrt{11}, \dots, \sqrt{p_n}, \dots$  We denote the set of all prime labels by  $\sqrt{\mathcal{P}}$ .

545

Now, we define a prime labeling of a finite set and of a cartesian product of finite sets.

**Definition 16 (Prime Labeling of Sets).** Let  $A, A_1, A_2, ..., A_n$  be (finite) sets. A prime labeling of the set A is an injective function  $w : A \to \sqrt{\mathcal{P}}$ . For  $a \in A$ , we call w(a) the prime label of a (for the prime labeling w).

Let I be a subset of  $\{1, 2, ..., n\}$ , which serves as an index set. A prime product I-labeling of the catesian product  $A_1 \times A_2 \times \cdots \times A_n$  consists of prime labelings  $w_i$  of the sets  $A_i$ , for  $i \in I$ , that satisfy the condition that  $w_i(A_i) \cap$  $w_j(A_j)$  is empty for  $i, j \in I$  and  $i \neq j$ . For  $(a_1, a_2, ..., a_n) \in A_1 \times A_2 \times \cdots \times A_n$ , we call  $\prod_{i \in I} w_i(a_i)$  the prime product I-label of  $(a_1, a_2, ..., a_n)$  (given the prime labelings  $w_i$ , for  $i \in I$ ). When I is a strict subset of  $\{1, 2, ..., n\}$ , we speak about a partial prime product labeling and when  $I = \{1, 2, ..., n\}$ , we speak about a full prime product labeling.

In the previous definition, whenever I is clear from the context, we can omit reference to I.

In practice, to label a set  $A_1 \times A_2 \times \cdots \times A_n$ , we use consecutive, available labels from  $\sqrt{\mathcal{P}}$  to label the sets  $A_1, A_2, \dots, A_n$ , as is illustrated by the following example. Further on, we apply this labeling to domains of dimensions (possibly at different levels).

**Example 15.** Let  $A_1 = \{a_1, a_2, a_3\}, A_2 = \{b_1, b_2\}$  and  $A_3 = \{c_1, c_2\}$ . To create a full prime product label for the elements of  $A_1 \times A_2 \times A_3$ , we can use the prime labelings  $w_1, w_2$  and  $w_3$ , defined as follows:  $w_1(a_1) = 1, w_1(a_2) = \sqrt{2}, w_1(a_3) = \sqrt{3}, w_2(b_1) = \sqrt{5}, w_2(b_2) = \sqrt{7}, w_3(c_1) = \sqrt{11}$  and  $w_3(c_3) = \sqrt{13}$ . We remark that we have used consecutive elements of  $\sqrt{\mathcal{P}}$  (with respect to the natural order or natural numbers). These labelings give the tuple  $(a_2, b_2, c_1)$  the label  $w_1(a_2) \cdot w_2(b_2) \cdot w_3(c_1) = \sqrt{2} \cdot \sqrt{7} \cdot \sqrt{11} = \sqrt{154}$ . Each cell in  $A_1 \times A_2 \times A_3$  gets a unique prime product label.

560

570

To create a partial prime product label for  $I = \{1, 2\}$ , we can use  $w_1$  and  $w_2$ , as given above. In this case, for any  $a \in A_1$  and  $b \in A_2$ , the cells  $(a, b, c_1)$  and  $(a, b, c_2)$  get the same (partial) prime product label.

If we view a cartesian product  $A_1 \times A_2 \times \cdots \times A_n$  as a finite matrix, whose cells contain rational-valued measures, we can use prime (product) labelings as follows in the aggregation process. Let us assume that the cells of  $A_1 \times A_2 \times$  $\cdots \times A_n$  contain rational values of a measure  $\mu$  and let us denote the value of this measure in the cell  $(a_1, a_2, ..., a_n)$  by  $\mu(a_1, a_2, ..., a_n)$ . If we have a full prime product labeling on  $A_1 \times A_2 \times \cdots \times A_n$ , then we can consider the sum over this cartesian product of the product of the prime product labels with the value of  $\mu$ :

$$\sum_{(a_1, a_2, \dots, a_n) \in A_1 \times A_2 \times \dots \times A_n} \mu(a_1, a_2, \dots, a_n) \cdot w_1(a_1) \cdot w_2(a_2) \cdots w_n(a_n).$$
(†1)

Since each cell of  $A_1 \times A_2 \times \cdots \times A_n$  has a unique prime product label, and since these labels are rationally independent (as we show in Property 2), this sum enables us to retrieve the values  $\mu(a_1, a_2, ..., a_n)$ .

If we have a partial prime product labeling on  $A_1 \times A_2 \times \cdots \times A_n$ , determined by an index set I, then, again, we can consider the sum over this cartesian product of the product of the partial prime product labels with the value of  $\mu$ :

$$\sum_{(a_1,a_2,\dots,a_n)\in A_1\times A_2\times\dots\times A_n}\mu(a_1,a_2,\dots,a_n)\cdot\prod_{i\in I}w_i(a_i).$$
 (†2)

Now, all cells in  $A_1 \times A_2 \times \cdots \times A_n$  above a cell in the projection of  $A_1 \times A_2 \times \cdots \times A_n$  on its components with indices in I receive the same prime label. This means that these cells are "grouped" together and the above sum allows us to retrieve the part of the sum that belongs to each group.

To make this clearer, we can write  $(\dagger_2)$  as

$$\sum_{\substack{\times_{i\in I}A_i}} \left( \sum_{\substack{\times_{i\in I^c}A_i}} \mu(a_1, a_2, ..., a_n) \right) \cdot \prod_{i\in I} w_i(a_i), \tag{\dagger}_2$$

where the outer sum ranges over the components of  $A_1 \times A_2 \times \cdots \times A_n$  whose index belongs to I and where the inner sum ranges over the components of  $A_1 \times A_2 \times \cdots \times A_n$  whose index belongs to  $I^c := \{1, 2, ..., n\} \setminus I$ . The above statement says that  $(\dagger_2)$  allows us to uniquely determine the sums

$$\sum_{A_i \in I^c A_i} \mu(a_1, a_2, \dots, a_n).$$

We remark that this last sum is the same for all cells above a cell in the projection of  $A_1 \times A_2 \times \cdots \times A_n$  on the components whose index is in I.

The following definition gives a name to the above sums.

>

**Definition 17 (Prime Sums).** We call sums of type  $(\dagger_1)$  full prime sums and sums of type  $(\dagger_2)$  partial prime sums (over I).

The following property can be derived from the well-known fact that the field extension  $\mathbf{Q}(\sqrt{2}, \sqrt{3}, ..., \sqrt{p_n}) = \{a_0 + a_1\sqrt{2} + a_2\sqrt{3} + \cdots + a_n\sqrt{p_n} \mid a_0, a_1, a_2, ..., a_n \in \mathbf{Q}\}$  has degree  $2^n$  over  $\mathbf{Q}$  and corollaries of this property (see Chapter 8 in [13]). No square root of a prime number is a rational combination of square roots of other primes.

**Property 2.** Let  $n \ge 1$  and let  $A_1 \times A_2 \times \cdots \times A_n$  be a cartesian product of finite sets. We assume that the cells  $(a_1, a_2, ..., a_n)$  of this set contain rational values  $\mu(a_1, a_2, ..., a_n)$  of a measure  $\mu$ . Let I be a subset of  $\{1, 2, ..., n\}$  and let  $w_i$ be prime labelings of the sets  $A_i$ , for  $i \in I$ , that form a prime product I-labeling (see Definition 16). Then we have that the prime sum  $(\dagger_2)$  uniquely determines the values  $\sum_{\times_{i \in I^c} A_i} \mu(a_1, a_2, ..., a_n)$  for all cells of  $A_1 \times A_2 \times \cdots \times A_n$ .

<sup>585</sup> We give the proof of this property in Appendix Appendix A.

**Remark 2.** We remark that we use these prime (product) labels in a purely *symbolic* way without actually calculating the square root values in them. The square roots are treated as symbolic entities in the computations.

Given these facts about prime (product) labels, we are ready to define atomic OLAP operations that allow us to implement grouping. In what follows, we apply these prime labels to the case where the sets  $A_i$  in  $A_1 \times A_2 \times \cdots \times A_n$  are domains of dimensions or domains of dimensions at some level.

<sup>590</sup> **Definition 18 (Grouping Transformations).** The following creations of a new measure  $\tau_{l+1}$  are grouping transformations:

- 1. (Prime labels for groups in one dimension) Let  $D_i$  be a dimension and  $\ell$  a level in the dimension schema  $\sigma(D_i)$  of a dimension  $D_i$ . Let  $dom(D_i.\ell) = \{b_1, b_2, ..., b_m\}$  with induced order  $b_1 < b_2 < \cdots < b_m$  (see Property 1). If the prime labels  $w_1, w_2, ..., w_k$  have been used by previous transformations, then for all j, with  $j \neq i$ , and all  $x_j \in dom(D_j)$ , we have  $\tau_{l+1}(x_1, ..., x_{i-1}, x_i, x_{i+1}, ..., x_d) = w_{k+l}$  if  $\rho(x_i, b_l)$ . We denote this transformation by  $\gamma_{D_i.\ell}(x_1, ..., x_{i-1}, x_i, x_{i+1}, ..., x_d)$  or  $\gamma_{D_i.\ell}$ , for short, and call the result of such a transformation a prime labeling.
- 2. (**Projection of a prime sum**) If the result of some previous transformation  $\tau_m$  is a (full or partial) prime sum  $\sum_{i=k}^{k+l} a_i \cdot w_i$  (over the complete matrix  $M(\mathcal{D})$ ) in which prime (product) labels  $w_k, w_{k+1}, ..., w_{k+l}$  (computed in a previous transformation  $\tau_n$ ) are used, then  $\tau_{l+1}$  is a new measure that "projects" on the appropriate component from the prime sum, that is,  $\tau_{l+1}(x_1, x_2..., x_d) = a_{k+l}$  if the prime (product) label  $\tau_n(x_1, x_2..., x_d) =$  $w_{k+l}$ . We denote this projection transformation by  $\tau_m \mid_{\tau_n}$ .

Now, we give some examples about (basic) counting. As always, we use *sales* information for certain products, at certain locations, at certain time moments. This is important in OLAP practice, since many times we need to compute the number of elements in a dimension level. That is, we perform an aggregation without operating on measures.

- Example 16. We look at the query "total number of cities", which asks to count the number of cities appearing at the *Bottom* level of the dimension *Location*. We can implement this query using Count-Distinct and prime labels, given  $\mu_1 = sales$ , as follows:
  - $\tau_1 = \gamma_{Location.City}$  (gives each city a different prime label);
- 610

•  $\tau_2 = \#_{\neq}(\tau_1)$  (counts the number of different prime labels and thus the number of cities);

•  $\varphi^{(1)} = \sigma_{Location.Bottom}$  (this flag creation selects all cells of the matrix).

The output measures are <u>sales</u>;  $\tau_2$ ;  $\varphi^{(1)}$ , which are renamed <u>sales</u>;  $\tau_1$ ;  $\varphi$ . This means that the total number of cities is now available in every cell of the matrix  $M(\mathcal{D})$ .

595

**Example 17.** We look at the query "total number of countries", which asks to count the number of countries appearing at the *Country* level of the dimension

615 I

Location. We can implement this query using Count-Distinct and prime labels, given  $\mu_1 = sales$ , as follows:

- $\tau_1 = \gamma_{Location.Country}$  (gives each country a different prime label);
- $\tau_2 = \#_{\neq}(\tau_1)$  (counts the number of different prime labels and thus the number of countries);
- $\varphi^{(1)} = \sigma_{Location.Bottom}$  (this flag creation selects all cells of the matrix).

The output measures are <u>sales</u>;  $\tau_2$ ;  $\varphi^{(1)}$ , which are renamed <u>sales</u>;  $\tau_1$ ;  $\varphi$ . This means that the total number of countries is now available in every cell of the matrix.

Note that, like in tghe previous example, we are operating on the dimensions, and we are not aggregating measures, which shows the generality of our approach. The next example goes further into this issue, since it uses grouping within the same dimension.

- Example 18 below, uses two prime labels on the dimension *Location*. One prime labeling is at the *Country* level and the second prime labeling is at the *City* or *Bottom* level. This construction fits in the given prime product labeling concept if we consider  $A_1 = dom(Location.Country)$  and  $A_2 = dom(Location.City)$ .
- Example 18. Consider the query "for each country, give the total number of cities". Again, we assume that  $\mu_1 = sales$  is the only available measure. This query can be implemented as follows (we explain the details below):
  - $\tau_1 = \gamma_{Location.Country}$  (this gives each country a prime label);
  - $\tau_2 = \gamma_{Location.City}$  (this gives each city a (fresh) prime label);
- $\tau_3 = \tau_1 \cdot \tau_2$  (this gives each city a product of prime labels);
  - $\tau_4 = \text{SUM}_3(\tau_3);$
  - $\tau_5 = \gamma_{Product.Bottom}$  (gives each product a different prime label);
  - $\tau_6 = \#_{\neq}(\tau_5)$  (counts the number of products—see Example 16);

- $\tau_7 = \gamma_{Time.Bottom}$  (gives each time moment a different prime label);
- $\tau_8 = \#_{\neq}(\tau_7)$  (counts the number of moments in time—see Example 16);
  - $\tau_9 = \tau_6 \cdot \tau_8$  (is the number of products times the number of time moments);
  - $\tau_{10} = \tau_4/\tau_9$  (normalisation of the sum);

640

- $\tau_{11} = \tau_{10} \mid_{\tau_2}$ ; (projection over the prime labels of city);
- $\tau_{12} = \text{SUM}_3(\tau_{11})$  (3-dimensional sum);
- $\tau_{13} = \tau_{12}/\tau_9$  (normalisation of the sum);
  - $\tau_{14} = \tau_{13} \mid_{\tau_1}$  (projection over the prime labels of country);
  - $\varphi^{(1)} = \sigma_{Location.Bottom}$  (this flag creation selects all cells of the matrix).

We now discuss this example, using the data given in Example 4. Transformation  $\tau_1$  gives each country a next available prime label. Since no labels have been used yet, *belgium* gets label 1 and *france* gets label  $\sqrt{2}$ . Transformation  $\tau_2$  gives each city a next available prime label. Since 1 and  $\sqrt{2}$  have been used, *antwerp* gets label  $\sqrt{3}$ , *brussels* gets label  $\sqrt{5}$ , *paris* gets label  $\sqrt{7}$ and *marseille* gets label  $\sqrt{11}$ .

Transformation  $\tau_3$  gives antwerp the value  $\sqrt{3}$  (i.e.,  $1.\sqrt{3}$ , brussels the value  $\sqrt{5}(1.\sqrt{5})$ , paris the value  $\sqrt{14}(\sqrt{2}.\sqrt{7})$  and marseille the value  $\sqrt{22}(\sqrt{2}.\sqrt{11})$ . If there are 10 products and 100 time moments, then  $\tau_4$  puts the value  $10 \cdot 100 \cdot (\sqrt{3} + \sqrt{5} + \sqrt{14} + \sqrt{22})$  in each cell of the matrix  $M(\mathcal{D})$ .

Transformations  $\tau_6$  and  $\tau_8$  count the number of products and the number of time moments (using fresh prime labels). In  $\tau_{10}$ ,  $\tau_3$  is divided by their product and  $\tau_{10}$  puts  $\sqrt{3} + \sqrt{5} + \sqrt{14} + \sqrt{22}$  in every cell of the matrix.

Transformation  $\tau_{11}$  is a projection on the prime labels of *City*. Since  $\sqrt{3}$ ,  $\sqrt{5}$ ,  $\sqrt{7}$  and  $\sqrt{11}$  are the prime labels for the cities, and since  $\sqrt{3} + \sqrt{5} + \sqrt{14} + \sqrt{22} = 1 \cdot \sqrt{3} + 1 \cdot \sqrt{5} + \sqrt{2} \cdot \sqrt{7} + \sqrt{2} \cdot \sqrt{11}$ , this will put 1 in the cells of *antwerp* and *brussels* and  $\sqrt{2}$  in the cells of *paris* and *marseille*.

Next,  $\tau_{12}$  puts  $10 \cdot 100 \cdot (2 \cdot 1 + 2 \cdot \sqrt{2})$  in every cell of the cube and  $\tau_{13}$  puts  $2 \cdot 1 + 2 \cdot \sqrt{2}$  in every cell of the cube. Finally,  $\tau_{14}$  projects on the prime labels of countries, which are 1 and  $\sqrt{2}$ . This puts a 2 in every cell of a Belgian city and a 2 in every cell in a French city. This is the result of the query, as the flag indicates, that is returned in every cell. Now every cell of a city in *belgium* has the count of 2 cities, as has every city in *france*.

#### 665 3.4.6. Counting and min-max revisited

Now that we know prime (product) labelings, we can give extensions of the counting and min-max transformations of Definition 14. Here, the counting, the minimum and the maximum are taken over cells which share a common prime (product) label.

- **Definition 19.** The following creations of a new measure  $\tau_{l+1}$  are generalisations of the *counting and min-max* transformations:
  - 1. (Count-Distinct) If the result of some previous transformation  $\tau_m$  is a prime (product) labeling of the cells of  $M(\mathcal{D})$ , then  $\tau_{l+1}(x_1, x_2..., x_d) =$  $\#_{\neq} \mid_{\tau_m} (\alpha)$ , with  $\alpha \in \{\mu_1, \mu_2, ..., \mu_k, \tau_1, \tau_2, ..., \tau_l\}$  counts the number of different values of the measure  $\alpha$  in cells of the matrix  $M(\mathcal{D})$  that have the same prime product label as  $\tau_m(x_1, x_2..., x_d)$ .
  - 2. (Min-Max) If the result of some previous transformation  $\tau_m$  is a prime (product) labeling of the cells of  $M(\mathcal{D})$ , then  $\tau_{l+1}(x_1, x_2..., x_d) = \min |_{\tau_m}$ ( $\alpha$ ), with  $\alpha \in \{\mu_1, \mu_2, ..., \mu_k, \tau_1, \tau_2, ..., \tau_l\}$ , gives the the smallest value of the measure  $\alpha$  in cells of the matrix  $M(\mathcal{D})$  that have the same prime product label as  $\tau_m(x_1, x_2..., x_d)$ . And  $\tau_{l+1}(x_1, x_2..., x_d) = \max |_{\tau_m} (\alpha)$  is defined similarly.

We remark that when there is only one prime label (for instance, 1) throughout the matrix  $M(\mathcal{D})$ , then the above gneralisation of the counting and min-max transformations correspond to the version of Definition 14.

#### 4. The classical OLAP operations

In this section, we show how the classical OLAP operations can be expressed using the OLAP transformations from Section 3. As we mentioned in that section, these classic operations can be combined to express complex analytical queries. The classical OLAP operations are

- Dice (see Section 4.2);
  - Slice (see Section 4.3);
  - Slice-and-Dice (see Section 4.4);
  - Roll-Up (see Section 4.5); and
  - Drill-Down (see Section 4.5).

Throughout this section, we assume that the input data cube  $\mathcal{D}_{in}$  has k given measures  $\mu_1, \mu_2, ..., \mu_k$  (as in Definition 7) and that at some point in the OLAP process this cube is transformed to a cube  $\mathcal{D}$ , having measures

$$\mu_1, \mu_2, ..., \mu_k; \tau_1, \tau_2, ..., \tau_l; \varphi,$$

where  $\tau_1, \tau_2, ..., \tau_l$ , with  $l \ge 0$ , are created measures and  $\varphi$  is an input/output flag.

4.1. Boolean cell-selection condition

In this section, we give the definition of a Boolean cell-selection condition. We also give a lemma about its expressibility that is used throughout Section 4.

**Definition 20.** Let  $M(\mathcal{D}) = dom(D_1) \times dom(D_2) \times \cdots \times dom(D_d)$  be the matrix of  $\mathcal{D}$ . A Boolean condition on the cells of  $M(\mathcal{D})$  is a function  $\phi$  from  $M(\mathcal{D})$  to  $\{0, 1\}$ . We say that the cells of  $M(\mathcal{D})$  in the set  $\phi^{-1}(\{1\})$  are selected by  $\phi$ .

We say that a Boolean condition  $\phi$  is transformation-expressible if there is a sequence of OLAP transformations  $\tau_1, \tau_2, ..., \tau_k$  such that  $\phi(x_1, x_2, ..., x_d) =$  $\tau_k(x_1, x_2, ..., x_d)$  for all  $(x_1, x_2, ..., x_d) \in M(\mathcal{D})$ .

**Lemma 1.** If  $\phi$ ,  $\phi_1$ ,  $\phi_2$  are transformation-expressible Boolean conditions on cells, then NOT  $\phi$ ,  $\phi_1$  AND  $\phi_2$  and  $\phi_1$  OR  $\phi_2$  are transformation-expressible Boolean conditions on cells.

We give the proof in Appendix Appendix B

4.2. Dice

Intuitively, the *Dice* operation selects the cells in a cube  $\mathcal{D}$  that satisfy a Boolean condition  $\phi$  on the cells. The syntax for this operation is

 $\mathsf{DICE}(\mathcal{D},\phi),$ 

where  $\phi$  is a Boolean condition over level values and measures. The resulting cube has the same dimensionality as the original cube. The dice operation is

analogous to a selection in the relational algebra. In a data cube, it selects the cells that satisfy the condition  $\phi$  by flagging them 1 in the output cube.

The *Dice* operation has been already illustrated in Examples 8, 9, 10 and 11. There, we have queries such as

 $\mathsf{DICE}(\mathcal{D}, Location.City = antwerp \ OR \ Location.City = brussels).$ 

But, we also allow equality and order constraints on objects at certain levels and in different dimensions, as illustrated by the example

 $\mathsf{DICE}(\mathcal{D}, Location. Country = belgium AND Time. Day > 15/1/2014).$ 

We also consider equality and order constraints over measures, as is illustrated by the query

$$\mathsf{DICE}(\mathcal{D}, sales > 50)$$

of Example 8. Therefore, our approach covers all typical cases in real-world OLAP [3]. We next formalize the operator's definition in terms of our transformation language.

**Definition 21 (Dice).** Given a data cube  $\mathcal{D}$ , the operation  $\mathsf{DICE}(\mathcal{D}, \phi)$ , selects all cells of the matrix  $M(\mathcal{D})$  that satisfy the Boolean condition  $\phi$  by giving them a 1 flag in the output. The Boolean condition  $\phi$  on the cells of  $M(\mathcal{D})$  is a Boolean combination of conditions of the form:

• a selector on a value b at a certain level  $\ell$  of some dimension  $D_i$ ;

710

- a comparison condition at some level  $\ell$  from a dimension schema  $\sigma(D_i)$ of a dimension  $D_i$  of the cube of the form  $\ell < c$  or  $c < \ell$ , where c is a constant (at that level  $\ell$ );
- an equality or comparison condition on some measure  $\alpha$  of the form  $\alpha = c$ ,  $\alpha < c$  or  $c < \alpha$ , where c is a (rational) constant.
- **Property 3.** Let  $\mathcal{D}$  be a data cube en let  $\phi$  be a Boolean condition on the cells of  $M(\mathcal{D})$  (as in Definition 21). The operation  $\mathsf{DICE}(\mathcal{D}, \phi)$  is expressible as an OLAP operation.

**Proof 1.** Since  $\mathsf{DICE}(\mathcal{D}, \phi)$  is a cell-selecting operation, it suffices, by Lemma 1, to show that  $\mathsf{DICE}(\mathcal{D}, \phi)$  is expressible in the OLAP algebra for an atomic Boolean cell-selection condition  $\phi$  (without logical connectives). We have to consider the three cases of Definition 21.

For the first case,  $\mathsf{DICE}(\mathcal{D}, \phi)$  is simply expressed by the selector  $\tau_{l+1} = \sigma_{D_i, \ell=b}$ , which is the output flag that indicates the appropriate cells of  $M(\mathcal{D})$ .

For the second case, if  $\ell$  is a level from a dimension schema  $\sigma(D_i)$  of a

dimension  $D_i$  and  $c \in dom(D_i.\ell)$  and  $\phi$  is of the form  $\ell < c$  or  $c < \ell$ , then the comparison test on levels  $\tau_{l+1} = (\ell <_{\ell} c)$  or  $\tau_{l+1} = (c <_{\ell} \ell)$  express  $\mathsf{DICE}(\mathcal{D}, \phi)$ . Again,  $\tau_{l+1}$  specifies the output flag.

For the third case, if  $\alpha$  is some measure, then  $\tau_{l+1} = c$  (rational constant), followed by  $\tau_{l+2} = (\alpha = c)$ ,  $\tau_{l+2} = (\alpha < c)$  or  $\tau_{l+2} = (\alpha > c)$  (equality or comparison test on a measure), respectively, express  $\mathsf{DICE}(\mathcal{D}, \phi)$ . Once again,  $\tau_{l+2}$  can serve as the output flag. This concludes the proof.

#### 4.3. Slice

Intuitively, the *Slice* operation takes as input a *d*-dimensional, *k*-ary data cube  $\mathcal{D}$  and a dimension  $D_i$  and returns as output  $\mathsf{SLICE}(\mathcal{D}, D_i)$ , which is 730 a "(d-1)-dimensional" data cube in which the original measures  $\mu_1, \ldots, \mu_k$ are replaced by their aggregation (sum) over different values of elements in  $dom(D_i)$ . In other words, dimension  $D_i$  is removed from the data cube, and, if this operation is part of a sequence of OLAP ones,  $D_i$  will not be visible in the next operations. That means, for instance, that we will not be able to dice on 735 the levels of the removed dimension. As we will see, the "removal" of dimensions are, in our approach, implemented by means of the destroyer measure  $\delta$ . We remark that the aggregation above is due to the fact that, in order to eliminate a dimension  $D_i$ , this dimension should have exactly one element [7], therefore a roll-up (which we explain later in Section 4.5) to the level All in D - i is 740 performed.

For example, if  $(D_1, D_2, D_3) = (Product, Location, Time)$ , and we consider

#### $SLICE(\mathcal{D}, Location),$

then we obtain a cube with (Product, Time)-cells which contain the sums of the given measures for certain products and times, but summed over all locations (for that product and that time).

Obviously, in our philosophy, we keep the d-dimensional data cube and store identical aggregate values for all locations, in the cells above some product-time combination. Next, we destroy all locations, except the representative for all in the Location dimension. As explained in Example 6, antwerp represents all and we only keep the cells for antwerp, where we keep the aggregate values. We formalize this next.

**Definition 22 (Slice).** Given a data cube  $\mathcal{D}$  and one of its dimensions  $D_i$ , the operation  $\mathsf{SLICE}(\mathcal{D}, D_i)$  "replaces" the measures  $\mu_1, \mu_2, ..., \mu_k$  by their aggregation (sum)  $\mu_n^{\Sigma_i}$  (for  $1 \le n \le k$ ) as follows:

$$\mu_n^{\Sigma_i}(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_d) = \sum_{x_i \in dom(D_i)} \mu_n(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_d),$$

for all  $(x_1, ..., x_{i-1}, x_i, x_{i+1}, ..., x_d) \in M(\mathcal{D})$ . The operation  $\mathsf{SLICE}(\mathcal{D}, D_i)$  destroys all cells except those of the representative of *all* for dimension  $D_i$ . We abbreviate the above 1-dimensional sum as  $\mathrm{SUM}_{D_i}(\mu_n)$ .

**Property 4.** Let  $\mathcal{D}$  be a data cube and let  $D_i$  be one of its dimensions. The operation  $\mathsf{SLICE}(\mathcal{D}, D_i)$  is expressible as an OLAP operation.

Before we give the proof of Property 4, we give a simple example that captures the idea of the proof.

**Example 19.** Consider our running example with dimensions  $(D_1, D_2, D_3) = (Product, Location, Time)$  and measure  $\mu_1 = sales$ , and consider the query

$$SLICE(\mathcal{D}, Location).$$

This query returns a cube with (product, time)-cells which contain the sums of  $\mu_1$  for each product-time combination, over all locations (for that product and that time). At the end, all cells not belonging to the representative of *all* in the dimension *Location*, that is, *antwerp*, are destroyed.

The query  $SLICE(\mathcal{D}, Location)$  is the result of the following transformations

- $\tau_{l+1} = \gamma_{Product.Bottom}$  (prime labels on products);
  - $\tau_{l+2} = \gamma_{Time.Bottom}$  (fresh prime labels on time moments);
  - $\tau_{l+3} = \tau_{l+1} \cdot \tau_{l+2}$  (product of two prime labels);
  - $\tau_{l+4} = \mu_1 \cdot \tau_{l+3}$  (product);
  - $\tau_{l+5} = \text{SUM}_3(\tau_{l+4})$  (3-dimensional sum);

•  $\tau_{l+6} = \tau_{l+5} \mid_{\tau_{l+3}}$  (projection on prime product labels);

- $\tau_{l+7} = \sigma_{Location.All}$  (selects the representative of all in the dimension Location);
- $\delta = \tau_{l+7}$  (destroys all cells apart from the representative of *all* in the dimension *Location*);
- 770
- $\varphi^{(1)} = \sigma_{Location.All}$  (this flag creation selects the relevant cells of the matrix).

The transformation  $\tau_{l+4}$  gives each (product, time)-combination a unique prime product label. This label is multiplied with the *sales* in each cell. We then make the global sum over  $M(\mathcal{D})$  in  $\tau_{l+5}$ . The transformation  $\tau_{l+6} = \tau_{l+5} |_{\tau_{l+3}}$ is the projection on the prime product labels for (product, time)-combinations. This gives each cell above some fixed (product, time)-combination the sum of the *sales* (over all locations) for that (product, time)-combination. All cells of  $M(\mathcal{D})$  that do not belong to *antwerp* (selected in  $\tau_{l+7}$ ), which represents *all*, are destroyed by  $\delta$ .

<sup>780</sup> **Proof 2 (of Property 4).** Let  $\mathcal{D}$  be a data cube and let  $D_i$  be one of its dimensions. The operation  $SLICE(\mathcal{D}, D_i)$  is expressible in the OLAP algebra by the following sequence of transformations:

- $\tau_{l+1} = \gamma_{D_1.Bottom}$  (prime labels on dimension  $D_1$ );
- ...
- $\tau_{l+1+i-2} = \gamma_{D_{i-1}.Bottom}$  (prime labels on dimension  $D_{i-1}$ );
  - $\tau_{l+1+i-1} = \gamma_{D_{i+1}.Bottom}$  (prime labels on dimension  $D_{i+1}$ );
  - ...
  - $\tau_{l+1+d-2} = \gamma_{D_d.Bottom}$  (prime labels on dimension  $D_d$ );
  - $\tau_{l+1+d-1} = \tau_{l+1} \cdot \tau_{l+1+1};$
- 790  $\tau_{l+1+d} = \tau_{l+1+d-1} \cdot \tau_{l+1+2};$

• ...

• ...

- $\tau_{l+1+2d-4} = \tau_{l+1+d-1} \cdot \tau_{l+1+d-2}$  (product of all prime labels);
- $\tau_{l+1+2d-3} = \mu_1 \cdot \tau_{l+1+2d-4}$  (product of measure with product of all prime labels):
- 795
- $\tau_{l+1+2d+k-4} = \mu_k \cdot \tau_{l+1+2d-4}$  (product of measure with product of all prime labels);
- $\tau_{l+1+2d+k-3} = \text{SUM}_d(\tau_{l+1+2d-3})$  (*d*-dimensional sum);
- ...
- $\tau_{l+1+2d+2k-4} = \text{SUM}_d(\tau_{l+1+2d+k-4})$  (*d*-dimensional sum);
  - $\tau_{l+1+2d+2k-3} = \tau_{l+1+2d+k-3} |_{\tau_{l+1+2d-4}}$  (projection on product labels);
  - ...
  - $\tau_{l+1+2d+3k-4} = \tau_{l+1+2d+2k-4} \mid_{\tau_{l+1+2d-4}}$  (projection on product labels);
  - $\tau_{l+1+2d+3k-3} = \sigma_{D_i.All}$  (selects the representative of all for dimension  $D_i$ );
- $\delta = \tau_{l+1+2d+3k-3}$  (destroyer);
  - $\varphi^{(k)} = \tau_{l+1+2d+3k-3}$  (output flag).

Transformations  $\tau_{l+1}, ..., \tau_{l+1+d-2}$  create (fresh) prime labels for each of the dimensions  $D_1, ..., D_{i-1}, D_{i+1}, ..., D_d$ . Transformation  $\tau_{l+1+2d-4}$  gives the product of all these prime labels. This means that every  $(x_1, ..., x_{i-1}, x_{i+1}, ..., x_d) \in$  $dom(D_1) \times \cdots \times dom(D_{i-1}) \times dom(D_{i+1}) \times \cdots \times dom(D_d)$  has a unique prime product label, that is shared by all cells above the projected cell  $(x_1, ..., x_{i-1}, x_{i+1}, ..., x_d)$  $\pi_{i+1}, ..., x_d)$  in the direction of the dimension  $D_i$ . Transformations  $\tau_{l+1+2d-3}, ..., \tau_{l+1+2d+k-4}$  multiply the measures  $\mu_1, \mu_2, ..., \mu_k$  with the prime product label. Transformations  $\tau_{l+1+2d+k-3}, ..., \tau_{l+1+2d+2k-4}$  make partial prime sums of the measures  $\mu_1, \mu_2, ..., \mu_k$  over the complete matrix  $M(\mathcal{D})$ . The last k transformations  $\tau_{l+1+2d+2k-3}, ..., \tau_{l+1+2d+3k-4}$  project on the prime-product-labels giving each cell above  $(x_1, ..., x_{i-1}, x_{i+1}, ..., x_d)$  the sum of the k measures above it. Finally, the destroyer  $\delta$  and the output flag  $\varphi^{(k)}$  select the representative of all for dimension  $D_i$  and make sure that the other cells of  $M(\mathcal{D})$  are destroyed.

The output, for cells that are not destroyed, is

$$\mu_1, \mu_2, \dots, \mu_k; \tau_{l+1+2d+2k-3}, \dots, \tau_{l+1+2d+3k-4}; \varphi^{(k)},$$

which is renamed to

$$\mu_1, \mu_2, ..., \mu_k; \tau_1, ..., \tau_k; \varphi.$$

For  $1 \le n \le k$ ,  $\tau_n = \mu_n^{\Sigma_i}$  is the desired aggregate value. This concludes the proof.

#### 820 4.4. Slice and Dice

825

A particular case of the *Slice* operation occurs when the dimension to be removed already contains a unique value at the bottom level. Then, we can avoid the roll-up to *All*, and define a new operation, called *Slice-and-Dice*. Although this can be seen as a *Dice* operation followed by a *Slice* one, in practice, in these situations, they are usually applied at the same time.

**Definition 23.** Given a data cube  $\mathcal{D}$ , one of its dimensions  $D_i$  and some value a in the domain  $dom(D_i)$ , the operation SLICE-DICE $(\mathcal{D}, D_i, a)$  contains all the cells in the matrix  $M(\mathcal{D})$  such that the value of the dimension  $D_i$  equals a. All other cells are destroyed.

**Property 5.** Let  $\mathcal{D}$  be a data cube,  $D_i$  on of its dimensions en let  $a \in dom(D_i)$ . The operation SLICE-DICE $(\mathcal{D}, D_i, a)$  is expressible as an OLAP operation.

**Proof 3.** Let  $\mathcal{D}$  be a data cube,  $D_i$  on of its dimensions en let  $a \in dom(D_i)$ . The selector  $\sigma_{D_i.Bottom=a}$  is the transformation that serves as destroyer and output flag and that expresses SLICE-DICE $(\mathcal{D}, D_i, a)$ . This concludes the proof.

**Example 20.** For our running example, SLICE-DICE( $\mathcal{D}$ , Location, antwerp) is implemented by the output flag  $\sigma_{Location,Citu=antwerp}$ .

#### 4.5. Roll-Up and Drill-Down

We now address two key operations in typical OLAOP practice, namely Roll - Up and Drill - Down. Intuitively, the former aggregates measure values along a dimension up to a certain level. The latter, disagregates measure values along a dimension, down to a certain level. However, as we already commented, although at first sight it may appear that Drill - Down is the inverse of Roll - Up, like stated in [7], this is not necessarily the case, particularly when we are composing several OLAP operations, and, for example, a Roll - Up is followed Roll = DLAP = DLAP

by a SLICE or a DICE. In these cases, we cannot just undo the Roll - Up, but we need to account for the cells that have been eliminated on the way.

More precisely, the *Roll-Up* operation takes as input a data cube  $\mathcal{D}$ , a dimension  $D_i$  and a subpath h of a hierarchy H over  $D_i$ , starting in a node  $\ell'$  and ending in a node  $\ell$ , and returns the aggregation of the original cube along  $D_i$ up to level  $\ell$  for some of the input measures  $\alpha_1, \alpha_2, ..., \alpha_r$ .

The roll-up operation uses one of the following classic SQL aggregation functions, applied to the indicated protected and computed measures  $\alpha_1, \alpha_2, ..., \alpha_r$ (selected from  $\mu_1, \mu_2, ..., \mu_k; \tau_1, ..., \tau_l; \varphi$ ):

#### • sum (SUM);

- average (AVG);
- minimum and maximum (MIN and MAX);
- count and count-distinct (COUNT and COUNT-DISTINCT).

We remark that, usually, measures have an associated *default* aggregation function. The typical aggregation function for the measure *sales*, for instance, is SUM.

We denote the above roll-up operation as

$$\mathsf{ROLL}\text{-}\mathsf{UP}(\mathcal{D}, D_i, H(\ell' \to \ell), \{(\alpha_i, f_i) \mid i = 1, 2, ..., r\}),\$$

where  $f_i$  is one of the above aggregation functions that is associated to  $\alpha_i$ , for i = 1, 2, ..., r. Since we are mainly interested in the expressibility of this operation as a sequence of atomic transformations, we remark that only the destination node  $\ell$  in the path h is relevant. Indeed, the result of this roll-up remains the same if the subpath h is extended to start from the *Bottom* node of dimension  $D_i$ . So, we can abbreviate the above notation to

ROLL-UP(
$$\mathcal{D}, D_i, H(\ell), \{(\alpha_i, f_i) \mid i = 1, 2, ..., r\}),\$$

and assume that the roll-up starts at the *Bottom* level.

The Drill-down operation takes as input a data cube  $\mathcal{D}$ , a dimension  $D_i$  and a subpath h of a hierarchy H over  $D_i$ , starting in a node  $\ell$  and ending in a node  $\ell'$  (at a lower level in the hierarchy), and returns the aggregation of the original cube along  $D_i$  from the bottom level up to level  $\ell'$ . The drill-down uses the same type of aggregation functions as the roll-up. Again, since we are only interested in expressibility of this operation, we remark that the drill-down operation

DRILL-DOWN
$$(\mathcal{D}, D_i, H(\ell' \leftarrow \ell), \{(\alpha_i, f_i) \mid i = 1, 2, ..., r\}),$$

has the same output as  $\mathsf{ROLL}$ -UP $(\mathcal{D}, D_i, H(\ell'), \{(\alpha_i, f_i) \mid i = 1, 2, ..., r\})$ . Therefore, for expressibility, we can limit the further discussion in this section to the roll-up.

855

We remark that, since we assume, by definition, that dimension graphs are *sound*, we can also omit reference to the hierarchy H in the above notation and simply write ROLL-UP( $\mathcal{D}, D_i, \ell, \{(\alpha_i, f_i) \mid i = 1, 2, ..., r\})$  and DRILL-DOWN( $\mathcal{D}, D_i, \ell', \{(\alpha_i, f_i) \mid i = 1, 2, ..., r\})$ , for these OLAP operations.

**Definition 24 (ROLLUP).** Given a data cube  $\mathcal{D}$ , one of its dimensions  $D_i$ , and a hierarchy H over  $D_i$ , ending in a node  $\ell$ , the operation

ROLL-UP(
$$\mathcal{D}, D_i, H(\ell), \{(\alpha_i, f_i) \mid i = 1, 2, ..., r\}$$
)

computes the aggregation of the measures  $\alpha_i$  by their aggregation functions  $f_i$ , for i = 1, 2, ..., r, as follows:

$$\begin{aligned} \alpha_i^{f_i}(x_1,...,x_{i-1},x_i,x_{i+1},...,x_d) &= \\ f_i(\{\alpha_i((x_1,...,x_{i-1},y_i,x_{i+1},...,x_d) \mid y_i \in dom(D_i) \text{ and } \rho_H(y_i,b)\}), \end{aligned}$$

for all  $(x_1, ..., x_{i-1}, x_i, x_{i+1}, ..., x_d) \in M(\mathcal{D})$ , for which  $\rho_H(y_i, b)$ , for some  $b \in dom(D_i.\ell)$ . This roll-up flags all representative *Bottom*-level objects for elements of  $dom(D_i.\ell)$  as active.

**Property 6.** Let  $\mathcal{D}$  be a data cube, let  $D_i$  be one of its dimensions, and let H be a hierarchy over  $D_i$  ending in a node  $\ell$ . Let  $\{(\alpha_i, f_i) \mid i = 1, 2, ..., r\}$  be a set of selected measures (taken from the protected measures  $\mu_1, \mu_2, ..., \mu_k$  and the computed measures  $\tau_1, ..., \tau_k$  of  $\mathcal{D}$ ), with their associated aggregation functions. The operation ROLL-UP( $\mathcal{D}, D_i, H(\ell), \{(\alpha_i, f_i) \mid i = 1, 2, ..., r\}$ ) is expressible as an OLAP operation.

**Proof 4.** Let  $\mathcal{D}$  be a data cube, let  $D_i$  be one of its dimensions, and let H be a hierarchy over  $D_i$  ending in a node  $\ell$ . Let  $\{(\alpha_i, f_i) \mid i = 1, 2, ..., r\}$  be a set of selected measures with their associated aggregation functions.

We start by remarking that the aggregations of the measures  $\alpha_i$  by the functions  $f_i$ , can be computed consecutively for i = 1, 2, ..., r. At the end their results are copied as the last r computed measures and an output flag of type  $\varphi^{(r)}$ , which is a selector  $\sigma_{D_i,\ell}$ , returns these r aggregation results as output. Now, it remains to be shown how the SQL aggregation functions SUM, AVG,
 MIN, MAX, COUNT and COUNT-DISTINCT can be implemented as sequences
 <sup>875</sup> of atomic OLAP transformations for an arbitrary measure α.

(1) SUM: We give a description of the implementation of the SUM by a sequence of atomic OLAP transformations. Since a detailed description of a similar procedure is given in the proof of Property 4, we refer to that proof for details.

- Here, we first create prime labels  $\gamma_{D_j,Bottom}$  for all  $j \neq i$  and, for dimension  $D_i$ , prime labels  $\gamma_{D_i,\ell}$  at the level  $\ell$ . Next, we create a measure that is the product of all these prime labels. This prime product label gives each cell  $(x_1, ..., x_{i-1}, y_i, x_{i+1}, ..., x_d)$  of the matrix a unique label, modulo rollingup to the same object at level  $\ell$  for the dimension  $D_i$ . This implies that  $(x_1, ..., x_{i-1}, y_i, x_{i+1}, ..., x_d)$  and  $(x_1, ..., x_{i-1}, y'_i, x_{i+1}, ..., x_d)$ , for which there is
- a  $b \in dom(D_i.\ell)$  such that  $\rho_H(y_i, b)$  and  $\rho_H(y'_i, b)$ , get the same prime product label. Then we take the *d*-dimensional sum of the product of this prime product label with  $\alpha$ . The projection on the prime product label, gives the desired result. That is, the cells  $(x_1, ..., x_{i-1}, y_i, x_{i+1}, ..., x_d)$  and  $(x_1, ..., x_{i-1}, y'_i, x_{i+1}, ..., x_d)$ , for which there is a  $b \in dom(D_i.\ell)$  such that  $\rho_H(y_i, b)$  and  $\rho_H(y'_i, b)$ , get the
- same aggregation (sum) value of  $\alpha$  over all objects that roll-up to b. For a detailed description, we refer to the proof of Property 4 and for an illustration, we refer to Example 21.

(2) COUNT: Here, we proceed in a similar way as in the case of SUM, with the modification that before taking the sum, we do not multiply the prime product labels with  $\alpha$ , but with 1.

895

If we are interested in counting the cells for which  $\alpha$  is non-zero, we can achieve this by multiplying the prime product labels by the quotient  $\alpha/\alpha$ , rather than by 1. We remark that by the definition of quotient 0/0 = 0, which implies that the cells with a zero value for  $\alpha$  are not counted.

(3) AVG: The aggregation function AVG can be implemented by the implementation of SUM, followed by the implementation of COUNT (counting all or all non-zeroes) and then computing the quotient of these two values.

(4) MIN and MAX: As in the case of SUM, we create prime product labels for all cells of  $M(\mathcal{D})$ . Let us call this prime product labels  $\tau_m$ . Then we <sup>905</sup> multiply this prime product labels by  $\alpha$ , resulting in the measure  $\tau_{m+1}$ . Next, we apply the generalised form of the maximum (or minimum) transformation max  $|\tau_m(\tau_{m+1})$  to obtain the maximum value of  $\alpha$  per prime product label. Similarly, min  $|\tau_m(\tau_{m+1})$  gives the desired minimal values.

(5) COUNT-DISTINCT: We proceed as in the case of MIN and MAX, but now we obtain the result by the transformation  $\#_{\neq} \mid_{\tau_m} (\tau_{m+1})$ , which is the generalized form of the Count-Distinct.

This concludes the proof.

Now, we illustrate the roll-up implementation, using our running example.

**Example 21.** In this example we simulate the roll-up operation, using prime (product) labels, sums and projections together with the 3-dimensional sum. We look at the query "total sales per country". We use the simplified syntax, only indicating the level to which we roll-up on the *Location* dimension (i.e., *Country*). The query

ROLL-UP(D, Location, Country, {(sales, SUM)})

is the result of the following transformations, given the measure  $\mu_1 = sales$ :

1.  $\tau_{\ell+1} = \gamma_{Product.Bottom}$  (prime labels on products);

915 2.  $\tau_{\ell+2} = \gamma_{Time.Bottom}$  (prime labels on time moments);

3.  $\tau_{\ell+3} = \gamma_{Location.Country}$  (prime labels on countries);

4.  $\tau_{\ell+4} = \tau_{\ell+1} \cdot \tau_{\ell+2} \cdot \tau_{\ell+3}$ ; (prime product label – in one step);

5.  $\tau_{\ell+5} = \mu_1 \cdot \tau_{\ell+4}$  (product of labels with *sales*);

6.  $\tau_{\ell+6} = \text{SUM}_3(\tau_{\ell+5})$  (3-dimensional sum);

920 7.  $\tau_{\ell+7} = \tau_{\ell+5} \mid_{\tau_{\ell+4}}$  (projection on prime product labels); 8.  $\varphi^{(1)} = \sigma_{Location.Country}$  (output flag on country-representatives). Transformation  $\tau_{\ell+4}$  gives every product-date-country combination a unique prime product label. Normally this product takes more steps. Above, we have abbreviated it to one transformation.

The transformation  $\tau_{\ell+7}$  gives the aggregation result and  $\varphi^{(1)}$  is the flag that says that only the cities *antwerp* and *paris*, which represent the level *Country*, are active in the output (and nothing else of the original cube).

925

930

We continue with another example of a roll-up operation. We only give high-level descriptions of its implementation as a sequence of atomic OLAP transformations.

**Example 22.** Let us consider a rather complex, although usual query in data analysis in real-world situations: "city-average sales, for cities whose average sales are above the country average".

The query can be answered using our OLAP transformations as follows:

- Compute the total sale per country (like in Example 13);
- Compute the number of sales per country (see the proof of Property 6);
- Take the quotient of these two values;
- Flag  $\sigma_{location.Country}$ ;
  - Compute the total sales over all products and all dates per city;
  - Compute the total (non-zero) sales per city;
  - Take the quotient of the two previous values;
  - Select the cities for which this quotient exceeds the "average sale per country".

940

• Use this last Boolean as an output flag.

# 4.6. The composition of classical OLAP operations

In this paper, we have proven the following theorem about the completeness of the proposed algebra. **Theorem 1.** The classical OLAP operations and their composition are expressible by OLAP operations (that is, as sequences of atomic OLAP transformations).

The proof of this theorem follows immediately from the properties in this section and the results in Section 3.

We conclude this section with an example that illustrates the power and generality of our approach, combining a sequence of OLAP operations, and expressing them as a sequence of OLAP transformations.

**Example 23.** Let us consider an OLAP user, who is analyzing sales in different countries and regions. She wants to analyze and compare sales in the north of Belgium (the Flanders region), and in the south of France (which we, generically, have denoted *south* in our running example). She starts navigating the cube (as we said, indistinctly this can be done through a query language or with a graphic tool), and first filters the cube, keeping just the cells of the two desired regions. This is done with the following expression:

 $\mathsf{DICE}(\mathcal{D}, Location.Region = flanders \ OR \ Location.Region = south).$ 

As we showed, this can be implemented as a sequence of atomic OLAP transformations. Now the user has a cube with the cells that do not have been destroyed. Next, within the same navigation process, she obtains the total sales, in France and Belgium, only considering the desired regions, by means of:

 $\mathsf{ROLL}$ -UP( $\mathcal{D}, Location, Country, \{(sales, \mathsf{SUM})\}).$ 

This will only consider the valid cells for rolling up. After this, our user only wants to keep the sales in France (since she is within the same process, she will obviously obtain the sales in the south of France). Thus, she writes (or "clicks"):

 $\mathsf{DICE}(\mathcal{D}, Location. Country = france).$ 

Finally, she wants to go back to the details, one level below in the hierarchy (that is, the sales in the south of France, the latter being the country she is at, at this stage of her navigation). For this, she does:

 $\mathsf{DRILL}$ - $\mathsf{DOWN}(\mathcal{D}, Location, Region, \{(sales, \mathsf{SUM})\})$ 

In our approach, this will be a roll-up from the bottom level to the *Region* level, but only considering the cells that have not been destroyed.  $\Box$ 

#### 5. Conclusion and discussion

We have presented a formal, mathematical approach, to solve a practical problem, which is, to provide, for the first time, a formal semantics to a collection of OLAP operations, frequently used in real-world practice. Although OLAP is a very popular field in data analytics, this is the first time a formalization like this is given. The need for this formalization is clear: in a world being flooded by data of different kinds, users must be provided with tools allowing them to have an abstract "cube view" and cube manipulation capabilities, regardless of the underlying data types. Without a solid basis and unambiguous definition of cube operations, the former could not be achieved. We claim that our work is the first one of this kind, and will serve as a basis to build more robust practical tools to address the forthcoming challenges in this field.

We have addressed the four core OLAP operations: slice, dice, roll-up, and drill-down. This does not harm the value of the work. On the contrary, this approach allows us to focus on our main interest, that is, to study the formal basis of the problem. Intuitively, although of course, we must prove this, our line of

- <sup>965</sup> work can be extended to address other kinds of OLAP queries, like queries involving more complex aggregate functions like moving averages, rankings, and the like. Further, cube combination operations, like drill-across, must be included in the picture. We believe that our contribution provides a solid basis upon which, a complete OLAP theory can be built.
- Acknowledgements: Alejandro Vaisman was supported by a travel grant from Hasselt University (Korte verblijven–inkomende mobiliteit, BOF15KV13). He was also partially supported by PICT-2014 Project 0787.

#### References

975

- R. Kimball, The Data Warehouse Toolkit: Practical Techniques for Building Dimensional Data Warehouse, Wiley, 1996.
- [2] C. Ciferri, R. Ciferri, L. Gómez, M. Schneider, A. Vaisman, E. Zimányi, Cube algebra: A generic user-centric model and query language for OLAP cubes, International Journal of Data Warehousing and Mining 9 (2) (2013) 39–65.
- [3] A. Vaisman, E. Zimányi, Data Warehouse Systems: Design and Implementation, Springer, 2014.
- [4] S. Harinath, R. Pihlgren, D.-Y. Lee, J. Sirmon, R. Bruckner, Professional Microsoft SQL Server 2012 Analysis Services with MDX and DAX, Wrox, 2012.

- [5] O. Romero, A. Abelló, On the need of a reference algebra for OLAP, in: Proceedings of the 9th International Conference on Data Warehousing and Knowledge Discovery, DaWaK'07, Regensburg, Germany, 2007, pp. 99–110.
- [6] M. Gyssens, L. Lakshmanan, A foundation for multi-dimensional databases, in: Proceedings of the 23rd International Conference on Very Large Data Bases, VLDB'97, Athens, Greece, 1997, pp. 106–115.
  - [7] R. Agrawal, A. Gupta, S. Sarawagi, Modeling multidimensional databases, in: Proceedings of the 15th International Conference on Data Engineering (ICDE'97), IEEE Computer Society, Birmingham, UK, 1997, pp. 232–243.
  - [8] H. D. Macedo, J. N. Oliveira, A linear algebra approach to OLAP, Formal Asp. Comput. 27 (2) (2015) 283–307.
  - [9] P. Vassiliadis, Modeling multidimesional databases, cubes and cube operations, in: Proceedings of the 10th International Conference on Scientific and Statistical Database Management, SSDBM'98, Capri, Italy, 1998, p. 53.
  - [10] F.Ravat, O. Teste, R. Tournier, G. Zurfluh, Algebraic and graphic languages for OLAP manipulations, IJDWM 4 (1) (2008) 17–46.
  - [11] L. Gómez, S. Gómez, A. Vaisman, A generic data model and query language for spatiotemporal OLAP cube analysis, in: Proceedings of the 15th International Conference on Extending Database Technology, EDBT 2012, Berlin, Germany, 2012, pp. 300–311.
  - [12] J. Varga, L. Etcheverry, A. Vaisman, O. Romero, T. B. Pedersen, C. Thomsen, Enabling OLAP on statistical linked open data, in: Proceedings of the 32nd International Conference of Data Engineering (ICDE 2016) (accepted, to appear).
- 1010

1005

[13] J.-P. Escofier, Galois Theory, Vol. 204 of Graduate Texts in Mathematics, Springer-Verlag, 2001.

995

## Appendix A. Proof of Property 2

1020

1025

We now give the proof of Property 2. We repeat the property here, to  $_{1015}$  facilitate reading.

**Property 2.** Let  $n \ge 1$  and let  $A_1 \times A_2 \times \cdots \times A_n$  be a cartesian product of finite sets. We assume that the cells  $(a_1, a_2, ..., a_n)$  of this set contain rational values  $\mu(a_1, a_2, ..., a_n)$  of a measure  $\mu$ . Let I be a subset of  $\{1, 2, ..., n\}$  and let  $w_i$  be prime labelings of the sets  $A_i$ , for  $i \in I$ , that form a prime product I-labelling (see Definition 16). Then we have that the prime sum  $(\dagger_2)$  uniquely

determines the values  $\sum_{\times_{i \in I^c} A_i} \mu(a_1, a_2, ..., a_n)$  for all cells of  $A_1 \times A_2 \times \cdots \times A_n$ .

**Proof 5.** First, we assume I contains one element. Without loss of generality, we may assume that  $I = \{1\}$ . Then the prime sum (over I) is

$$\sum_{(a_1,a_2,...,a_n)\in A_1\times A_2\times\cdots\times A_n} \mu(a_1,a_2,...,a_n) \cdot w_1(a_1) = \sum_{a_1\in A_1} \left( \sum_{(a_2,...,a_n)\in A_2\times\cdots\times A_n} \mu(a_1,a_2,...,a_n) \right) \cdot w_1(a_1).$$

Let us assume this sum is equal to

$$\sum_{a_1 \in A_1} \left( \sum_{(a_2, \dots, a_n) \in A_2 \times \dots \times A_n} \mu'(a_1, a_2, \dots, a_n) \right) \cdot w_1(a_1)$$

for some measure  $\mu'$  and that there exists a  $a_0 \in A_1$  such that

$$\sum_{(a_2,...,a_n)\in A_2\times\cdots\times A_n}\mu(a_0,a_2,...,a_n)\neq \sum_{(a_2,...,a_n)\in A_2\times\cdots\times A_n}\mu'(a_0,a_2,...,a_n).$$

Since all  $\mu(a_1, a_2, ..., a_n)$  and  $\mu'(a_1, a_2, ..., a_n)$  are assumed to be rational numbers, this implies that  $w_1(a_0)$  is a rational combination of the other labels  $w_1(a_1)$ , with  $a_1 \in A_1 \setminus \{a_0\}$ . Since the labels  $w_1(a_1)$ , with  $a_1 \in A_1$ , are square roots of different prime numbers, this leads to a contradiction, since the field extension  $\mathbf{Q}(\sqrt{2}, \sqrt{3}, ..., \sqrt{p_n})$ , for any n, has degree  $2^n$  over  $\mathbf{Q}$ . In other words, the square roots  $\sqrt{2}, \sqrt{3}, ..., \sqrt{p_n}$  (together with 1) are linearly independent over  $\mathbf{Q}$  (see Chapter 8 in [13]). When the cardinality of I is strictly larger than 1, we can use a similar argumentation. Then we work with prime product labels of the form  $\prod_{i \in I} w_i(a_i)$ . Because of the restrictions on these products, imposed by Definition 16 (injectivity of the labelling function per dimension and disjointness of labels between dimensions), we see that these product labels differ one from the other by at least one prime factor (under the square root). Therefore, these labels are also linearly independent over  $\mathbf{Q}$  [13]. This completes the proof.

#### Appendix B. Proof of Lemma 1

**Proof 6.** Obviously, a Boolean combination of Boolean conditions is a Boolean condition. Let us assume that  $\phi$ ,  $\phi_1$  and  $\phi_2$  are transformation-expressible by sequences of OLAP transformations that end in  $\tau_k$ ,  $\tau_{k_1}$  and  $\tau_{k_2}$ , respectively. Then  $\phi_1$  AND  $\phi_2$  can be expressed by the transformation  $\tau_m = \tau_{k_1} \cdot \tau_{k_2}$ , which is 1 on cells if and only if both  $\phi_1$  and  $\phi_2$  give 1 on those cells.

<sup>1035</sup> For the negation, we have the following sequence of additional transformations:

- $\tau_m = 1$  (rational constant);
- $\tau_{m+1} = -1$  (rational constant);
- $\tau_{m+2} = \tau_{m+1} \cdot \tau_k$  (product); and
- 1040  $\tau_{m+3} = \tau_m + \tau_{m+2}$  (sum).

Here, we simulate substraction using the sum. The transformation  $\tau_{m+3}$  equals  $\tau_m - \tau_k$  and turns  $\tau_k = 0$  into 1 and a  $\tau_k = 1$  into 0. So, the transformation  $\tau_{m+3}$  expresses NOT  $\phi$ .

Via de Morgan's law, we can express  $\phi_1$  OR  $\phi_2$  using conjunction and negation. An alternative implementation of the OR is given by  $\tau_m = \tau_{k_1} + \tau_{k_2}$  (this sum gives 0, 1 or 2); and  $\tau_{m+1} = \tau_m/\tau_m$ . This last transformation maps 1 and 2 on 1 and 0 on 0 (in Definition 11, we defined 0/0 to be 0).