

# View Discovery in OLAP Databases through Statistical Combinatorial Optimization

Cliff Joslyn<sup>1</sup>, John Burke<sup>1</sup>, Terence Critchlow<sup>1</sup>,  
Nick Hengartner<sup>2</sup>, and Emilie Hogan<sup>1,3</sup>

<sup>1</sup> Pacific Northwest National Laboratory

<sup>2</sup> Los Alamos National Laboratory

<sup>3</sup> Mathematics Department, Rutgers University

**Abstract.** The capability of OLAP database software systems to handle data complexity comes at a high price for analysts, presenting them a combinatorially vast space of views of a relational database. We respond to the need to deploy technologies sufficient to allow users to guide themselves to areas of local structure by casting the space of “views” of an OLAP database as a combinatorial object of all projections and subsets, and “view discovery” as a search process over that lattice. We equip the view lattice with statistical information theoretical measures sufficient to support a combinatorial optimization process. We outline “hop-chaining” as a particular view discovery algorithm over this object, wherein users are guided across a permutation of the dimensions by searching for successive two-dimensional views, pushing seen dimensions into an increasingly large background filter in a “spiraling” search process. We illustrate this work in the context of data cubes recording summary statistics for radiation portal monitors at US ports.

## 1 Introduction and Related Work

OnLine Analytical Processing (OLAP) [6,7] is a relational database technology providing users with rapid access to summary, aggregated views of a single large database, and is widely recognized for knowledge representation and discovery in high-dimensional relational databases. OLAP technologies provide intuitive and graphical access to the massively complex set of possible summary views available in large relational (SQL) structured data repositories [21]. But the ability of OLAP database software systems, such as the industry-leading Hyperion<sup>1</sup> and ProClarity<sup>2</sup> platforms, to handle data complexity comes at a high price for analysts. The available portions and projections of the overall data space present a bewilderingly wide-ranging, combinatorially vast, space of options. There is an urgent need for knowledge discovery techniques that guide users’ knowledge

---

<sup>1</sup> <http://www.oracle.com/technology/products/bi/essbase/visual-explorer.html>

<sup>2</sup> <http://www.microsoft.com/bi/products/ProClarity/proclarity-overview.aspx>

discovery tasks; to find relevant patterns, trends, and anomalies; and to do so within the intuitive interfaces provided by “business intelligence” OLAP tools.

For example, consider a decision-maker responsible for analyzing a large relational database of records of events of personal vehicles, cargo vehicles, and others passing through radiation portal monitors (RPMs) at US ports of entry. In our data cubes include dimensions for multiple time representations, spatial hierarchies of collections of RPMs at different locations, and RPM attributes such as vendor. In this context, a vast collection of different views, focusing on different combinations of dimensions, and different subsets of records, are available to the user. How can the user assess the relative significance of different views? Given, for example, an initial view focusing on RPM type and date, is it more significant to focus on passenger or cargo RPMs, or a particular month in the year? And given such a selection, is it more significant to next consider a spatial dimension, or any of more than a dozen independent dimensions available?

Through the Generalized Data-Driven Analysis and Integration (GDDAI) Project [11], our team has been developing both pure and hybrid OLAP data analysis capabilities for a range of homeland security applications. The overall GDDAI goal is to provide a seamless integration of analysis capabilities, allowing analysts to focus on understanding the *data* instead of the *tools*. We describe GDDAI’s approach to knowledge discovery in OLAP data cubes using information-theoretical combinatorial optimization, and as applied in the ProClarity platform on databases of surveillance data from radiation monitors at US ports of entry. We aim at a formalism for user-assisted knowledge discovery in OLAP databases around the fundamental concept of **view chaining**. Users are provided with analytical feedback to guide themselves to areas of high local structure within view space: that is, to significant collections of dimensions and data items (columns and rows, respectively), in an OLAP-structured database.

OLAP is fundamentally concerned with a collection of  $N$  variables  $X^i$  and a multi-dimensional data relation over their Cartesian product  $\mathbf{X} := \times_{i=1}^N X^i$ . Thus formalisms for OLAP data analysis are naturally rooted in relational database theory, and OLAP formalisms [1,10,28] extend relational calculi and algebras, for example extending the SQL language to its multi-dimensional analog MDX<sup>3</sup>. But OLAP shares mathematical connections with a range of multivariate analytical approaches operable on the space  $\mathbf{X}$ , for example statistical databases [26]; the analysis of contingency tables [3]; hierarchical log-linear modeling [18]; grand tour methods in multi-variate data visualization [2]; projection pursuit [19]; data tensor analysis [14]; and reconstructibility analysis [13,20].

Our ultimate goal is to place OLAP knowledge discovery methods within a mathematical context of combinatorial optimization in such a manner as to be realizable within existing industry-standard database platforms. Specifically:

---

<sup>3</sup> <http://msdn.microsoft.com/en-us/library/ms145506.aspx>

- Given a foundational OLAP database engine platform (e.g. EssBase<sup>4</sup>, SSAS<sup>5</sup>);
- And an OLAP client with technology for graphical display and an intuitive interface (e.g. ProClarity, Hyperion);
- Cast the space of **views** (sub-cubes) of an OLAP database as a combinatorial, lattice-theoretical [9] structure;
- Equipped with statistical measures reflecting the structural relations among views (their dimensional scope, depth, etc.) in the context of the data observed within them;
- To support both automated search to areas of high local structure;
- And user-guided exploration of views in the context of these measures.

While we believe that our emphasis on a combinatorial approach is distinct, our work resonates with that of a number of others. Our “view chaining” (moving from one projected subset of a data cube to another intersecting in dimensionality) is similar to the navigational processes described by others [23,24,25], and anticipated in some of our prior work [12]. But approaches which seek out “drill-down paths” [5] only “descend” the view lattice along one “axis” of views with increasing dimensional extension, sequentially adding variables to the view at each step. In contrast, our “hop-chaining” technique chains through a sequence of two-dimensional views, affecting a permutation of the variables  $X^i$ .

Our overall approach is consistent with an increasingly large body of similar work drawing on information theoretical statistical measures in data cubes to provide quantities for making navigational choices [20,22]. However, some other researchers have used different statistical approaches, for example variance estimation [25] or skewness measures [16]. Our primary departure from traditional OLAP analysis is the extension to conditioning and conditional probability measures over views. This not only provides the basis for optimization and navigation, it also creates a strong connection to graphical or structural statistical models [4,27], graphoid logics [17,27], as well as systems-theory based structural model induction methodologies [13,15].

We begin by establishing concepts and notation for (non-hierarchical) OLAP databases over data tensors, and then define the **view lattice** of projected subsets over such structures. This brings us to a point where we can explicate the fundamental (again non-hierarchical) OLAP operations of projection, extension, filtering, and “flushing” (decreasing a filter). We introduce **conditional views** and the complex combinatorial object which is the **conditional view space**. This prepares us to introduce “hop-chaining” as a particular view discovery algorithm over this combinatorial object, wherein users are guided by conditional information measures across a permutation of the dimensions by searching for successive two-dimensional views, pushing seen dimensions in a “spiraling” search process into an increasingly large background filter. We then consider how to move to the fully hierarchical case, before illustrating hop-chaining on databases of surveillance data from radiation monitors at US ports of entry.

---

<sup>4</sup> <http://www.oracle.com/appserver/business-intelligence/essbase.html>

<sup>5</sup> [http://msdn.microsoft.com/en-us/library/ms175609\(SQL.90\).aspx](http://msdn.microsoft.com/en-us/library/ms175609(SQL.90).aspx)

## 2 OLAP Formalism

Although the mathematical tools required to analyze OLAP databases are relatively simple, their notational formalisms are inevitably, and regrettably, not [1,10,28]. While our formalism is similar, it differs in a number of ways as well:

- We combine projections  $I$  on dimensions and restrictions  $J$  on records into a lattice-theoretical object called a **view**  $\mathcal{D}_{I,J}$ .
- OLAP concerns databases organized around collections of variables which can be distinguished as: dimensions, which have a hierarchical structure, and whose Cartesian product forms the data cube’s schema; and measures, which can be numerically aggregated within different slices of that schema. For this work we consider cubes with a single integral measure, which in our application is the count of a number of records in the underlying database. While in principle any numerical measure could yield, through appropriate normalization, frequency distributions for use in our view discovery technique, these count measures do so directly and naturally. In future work we will consider the generalization to arbitrary numerical measures.
- Our view discovery method is currently only available on flat dimensions which are not hierarchically-structured to support roll-up aggregation and drill-down disaggregation operations. Future directions to extend to fully hierarchical OLAP data cubes will be indicated in Sec. 5.3.

### 2.1 Chaining Operations in the View Lattice of Data Tensor Cubes

Let  $\mathbb{N} = \{1, 2, \dots\}$ ,  $\mathbb{N}_N := \{1, 2, \dots, N\}$ . For some  $N \in \mathbb{N}$ , define a **data cube** as an  $N$ -dimensional tensor  $\mathcal{D} := \langle \mathbf{X}, \mathcal{X}, c \rangle$  where:

- $\mathcal{X} := \{X^i\}_{i=1}^N$  is a collection of  $N$  **variables** or **columns** with  $X^i := \{x_{k^i}\}_{k^i=1}^{L^i} \in \mathcal{X}$ ;
- $\mathbf{X} := \times_{X^i \in \mathcal{X}} X^i$  is a **data space** or **data schema** whose members are  $N$ -dimensional vectors  $\mathbf{x} = \langle x_{k^1}, x_{k^2}, \dots, x_{k^N} \rangle = \langle x_{k^i} \rangle_{i=1}^N \in \mathbf{X}$  called **slots**;
- $c : \mathbf{X} \rightarrow \{0, 1, \dots\}$  is a **count** function.

Let  $M := \sum_{\mathbf{x} \in \mathbf{X}} c(\mathbf{x})$  be the total number of records in the database. Then  $\mathcal{D}$  also has relative frequencies  $f$  on the cells, so that  $f : \mathbf{X} \rightarrow [0, 1]$ , where  $f(\mathbf{x}) = \frac{c(\mathbf{x})}{M}$ , and thus  $\sum_{\mathbf{x} \in \mathbf{X}} f(\mathbf{x}) = 1$ . An example of a data tensor with simulated data for our RPM cube is shown in Table 1, for  $\mathcal{X} = \{X^1, X^2, X^3\} = \{\text{RPM Manufacturer, Location, Month}\}$ , with RPM Mfr = { Ludlum, SAIC }, Location = { New York, Seattle, Miami }, and Month = { Jan, Feb, Mar, Apr }, so that  $N = 3$ . The table shows the counts  $c(\mathbf{x})$ , so that  $M = 74$ , and the frequencies  $f(\mathbf{x})$ .

At any time, we may look at a projection of  $\mathcal{D}$  along a sub-cross-product involving only certain dimensions with indices  $I \subset \mathbb{N}_N$ . Call  $I$  a **projector**, and denote  $\mathbf{x} \downarrow I = \langle x_{k^i} \rangle_{i \in I} \in \mathbf{X} \downarrow I$  where  $\mathbf{X} \downarrow I := \times_{i \in I} X^i$ , as a projected vector

**Table 1.** An example data tensor. Blank entries repeat the elements above, and rows with zero counts are suppressed.

RPM Mfr	Location	Month	$c(\mathbf{x})$	$f(\mathbf{x})$	RPM Mfr	Location	Month	$c(\mathbf{x})$	$f(\mathbf{x})$
Ludlum	New York	Jan	1	0.014	SAIC	New York	Jan	1	0.014
		Mar	3	0.041			Feb	4	0.054
		Apr	7	0.095		Seattle	Mar	3	0.041
	Seattle	Jan	9	0.122			Apr	3	0.041
		Apr	15	0.203			Miami	Jan	6
	Miami	Jan	2	0.027		Feb		2	0.027
		Feb	8	0.108		Mar		4	0.054
		Mar	4	0.054		Apr		1	0.014
Apr		1	0.014						

and data schema. We write  $\mathbf{x} \downarrow i$  for  $\mathbf{x} \downarrow \{i\}$ , and for projectors  $I \subseteq I'$  and vectors  $\mathbf{x}, \mathbf{y} \in \mathbf{X}$ , we use  $\mathbf{x} \downarrow I \subseteq \mathbf{y} \downarrow I'$  to mean  $\forall i \in I, \mathbf{x} \downarrow i = \mathbf{y} \downarrow i$ .

Count and frequency functions convey to the projected count and frequency functions denoted  $c[I] : \mathbf{X} \downarrow I \rightarrow \mathbb{N}$  and  $f[I] : \mathbf{X} \downarrow I \rightarrow [0, 1]$ , so that

$$c[I](\mathbf{x} \downarrow I) = \sum_{\mathbf{x}' \downarrow \mathbb{N}_N \supseteq \mathbf{x} \downarrow I} c(\mathbf{x}') \quad (1)$$

$$f[I](\mathbf{x} \downarrow I) = \sum_{\mathbf{x}' \downarrow \mathbb{N}_N \supseteq \mathbf{x} \downarrow I} f(\mathbf{x}'), \quad (2)$$

and  $\sum_{\mathbf{x} \downarrow I \in \mathbf{X} \downarrow I} f[I](\mathbf{x} \downarrow I) = 1$ . In words, we add the counts (resp. frequencies) over all vectors in  $\mathbf{y} \in \mathbf{X}$  such that  $\mathbf{y} \downarrow I = \mathbf{x} \downarrow I$ . This is just the process of building the  $I$ -marginal over  $f$ , seen as a joint distribution over the  $X^i$  for  $i \in I$ .

Any set of record indices  $J \subseteq \mathbb{N}_M$  is called a **filter**. Then we can consider the filtered count function  $c^J : \mathbf{X} \rightarrow \{0, 1, \dots\}$  and frequency function  $f^J : \mathbf{X} \rightarrow [0, 1]$  whose values are reduced by the restriction in  $J \subseteq \mathbb{N}_M$ , now determining

$$M' := \sum_{\mathbf{x} \in \mathbf{X}} c^J(\mathbf{x}) = |J| \leq M. \quad (3)$$

We renormalize the frequencies  $f^J$  over the resulting  $M'$  to derive

$$f^J(\mathbf{x}) = \frac{c^J(\mathbf{x})}{M'}, \quad (4)$$

so that still  $\sum_{\mathbf{x} \in \mathbf{X}} f^J(\mathbf{x}) = 1$ .

Finally, when both a selector  $I$  and filter  $J$  are available, then we have  $c^J[I] : \mathbf{X} \downarrow I \rightarrow \{0, 1, \dots\}, f^J[I] : \mathbf{x} \downarrow I \rightarrow [0, 1]$  defined analogously, where now  $\sum_{\mathbf{x} \downarrow I \in \mathbf{X} \downarrow I} f^J[I](\mathbf{x} \downarrow I) = 1$ . So given a data cube  $\mathcal{D}$ , denote  $\mathcal{D}_{I,J}$  as a **view** of  $\mathcal{D}$ , restricting our attention to just the  $J$  records projected onto just the  $I$  dimensions  $\mathbf{X} \downarrow I$ , and determining counts  $c^J[I]$  and frequencies  $f^J[I]$ .

In a lattice theoretical context [9], each projector  $I \subseteq \mathbb{N}_N$  can be cast as a point in the Boolean lattice  $\mathcal{B}^N$  of dimension  $N$  called a **projector lattice**. Similarly, each filter  $J \subseteq \mathbb{N}_M$  is a point in a Boolean lattice  $\mathcal{B}^M$  called a **filter**

**lattice.** Thus each view  $\mathcal{D}_{I,J}$  maps to a unique node in the **view lattice**  $\mathcal{B} := \mathcal{B}^N \times \mathcal{B}^M = 2^N \times 2^M$ , the Cartesian product of the projector and filter lattices.

We then define **chaining** operations as transitions from an initial view  $\mathcal{D}_{I,J}$  to another  $\mathcal{D}_{I',J}$  or  $\mathcal{D}_{I,J'}$ , corresponding to a move in the view lattice  $\mathcal{B}$ :

**Projection:** Removal of a dimension so that  $I' = I \setminus \{i\}$  for some  $i \in I$ . This corresponds to moving a single step down in  $\mathcal{B}^N$ , and to marginalization in statistical analyses. We have  $\forall \mathbf{x}' \downarrow I' \in \mathbf{X} \downarrow I'$ ,

$$c^J [I'](\mathbf{x}' \downarrow I') = \sum_{\mathbf{x} \downarrow I \supseteq \mathbf{x}' \downarrow I'} c^J [I](\mathbf{x}). \tag{5}$$

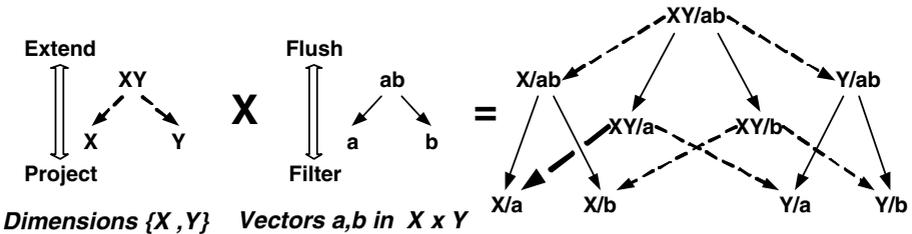
**Extension:** Addition of a dimension so that  $I' = I \cup \{i\}$  for some  $i \notin I$ . This corresponds to moving a single step up in  $\mathcal{B}^N$ . We're now *disaggregating* or *distributing* information about the  $I$  dimensions over the  $I' \setminus I$  dimensions. Notationally, we have the converse of (5), so that  $\forall \mathbf{x} \downarrow I \in \mathbf{X} \downarrow I$ ,

$$\sum_{\mathbf{x}' \downarrow I' \supseteq \mathbf{x} \downarrow I} c^J [I'](\mathbf{x}') = c^J [I](\mathbf{x} \downarrow I).$$

**Filtering:** Removal of records by strengthening the filter, so that  $J' \subseteq J$ . This corresponds to moving potentially multiple steps down in  $\mathcal{B}^M$ .

**Flushing:** Addition of records by weakening (reversing, flushing) the filter, so that  $J' \supseteq J$ . Corresponds to moving potentially multiple steps up in  $\mathcal{B}^M$ .

Repeated chaining operations thus map to trajectories in  $\mathcal{B}$ . Consider the very small example shown in Fig. 1 for  $N = M = 2$  with dimensions  $\mathcal{X} = \{X, Y\}$  and two  $N$ -dimensional data vectors  $\mathbf{a}, \mathbf{b} \in X \times Y$ , and denote e.g.  $X/\mathbf{ab} = \{\mathbf{a} \downarrow \{X\}, \mathbf{b} \downarrow \{X\}\}$ . The left side of Fig. 1 shows the separate projector and selector lattices (bottom nodes  $\emptyset$  not shown), with extension as a transition to a higher rank in the lattice and projection as a downward transition. Similarly, filtering and flushing are the corresponding operations in the filter lattice. The view lattice is shown on the right, along with a particular chain operation  $\mathcal{D}_{\{X,Y\},\{\mathbf{a}\}} \mapsto \mathcal{D}_{\{X\},\{\mathbf{a}\}}$ , which projects the subset of records  $\{\mathbf{a}\}$  from the two-dimensional view  $\{X, Y\} = \mathcal{X}$  to the one-dimensional view  $\{X\} \subseteq \mathcal{X}$ .



**Fig. 1.** The lattice theoretical view of data views. (Left) The projector and filter lattices  $\mathcal{B}^N, \mathcal{B}^M$  (global lower bounds  $\emptyset$  not shown). (Right) The view lattice  $\mathcal{B}$  as their product. The projection chain operation  $\mathcal{D}_{\{X,Y\},\{\mathbf{a}\}} \mapsto \mathcal{D}_{\{X\},\{\mathbf{a}\}}$  is shown as a bold link.

## 2.2 Relational Expressions and Background Filtering

Note that usually  $M \gg N$ , so that there are far more records than dimensions (in our example,  $M = 74 > 3 = N$ ). In principle, filters  $J$  defining which records to include in a view can be specified arbitrarily, for example through any SQL or MDX **where** clause, or through OLAP operations like **top  $n$** , including the  $n$  records with the highest value of some feature. In practice, filters are specified as relational expressions in terms of the dimensional values, as expressed in MDX **where** clauses. In our example, we might say **where RPM Mfr = "Ludlum" and ( Month <= "Feb" and Month >= "Jan")**, using chronological order on the Month variable to determine a filter  $J$  specifying just those 20 out of the total possible 74 records. For notational purposes, we will therefore sometimes use these relational expressions to indicate the corresponding filters.

Note that each relational filter expression references a certain set of variables, in this case RPM Mfr and Month, denoted as  $R \subseteq \mathbb{N}_N$ . Compared to our projector  $I$ ,  $R$  naturally divides into two groups of variables:

**Foreground:** Those variables in  $R^f := R \cap I$  which appear in both the filter expression and are included in the current projection.

**Background:** Those variables in  $R^b := R \setminus I$  which appear only in the filter expression, but are not part of the current projection.

The portions of filter expressions involving foreground variables restrict the rows and columns displayed in the OLAP tool. Filtering expressions can have many sources, such as **Show Only** or **Hide**. It is common in full (hierarchical) OLAP to select a collection of siblings within a particular sub-branch of a hierarchical dimension. For example for a spatial dimension, the user within the ProClarity tool might select **All -> USA -> California**, or its children **California -> Cities**, all siblings. But those portions of filter expressions involving background variables do not change which rows or columns are displayed, but only serve to reduce the values shown in cells. In ProClarity, these are shown in the Background pane.

## 2.3 Example

Table 2 shows the results of four chaining operations from our original example in Table 1, including a projection  $I = \{1, 2, 3\} \mapsto I' = \{1, 2\}$ , a filter using relational expressions, and a filter using a non-relational expression. The bottom right shows a hybrid result of applying both the projector  $I' = \{1, 2\}$  and the relational filter expression **where RPM Mfr = "Ludlum" and ( Month <= "Feb" and Month >= "Jan")**. Compare this to the top left, where there is only a quantitative restriction for the same dimensionality because of the use of a background filter. Here  $I = \{ \text{RPM Mfr, Location} \}$ ,  $R = \{ \text{RPM Mfr, Month} \}$ ,  $R^f = \{ \text{RPM Mfr} \}$ ,  $R^b = \{ \text{Month} \}$ ,  $M' = 20$ .

**Table 2.** Results from chaining operations  $\mathcal{D}_{\mathbb{N}_N, \mathbb{N}_M} \mapsto \mathcal{D}_{I', J'}$  from the data cube in Table 1. (Top Left) Projection:  $I' = \{1, 2\}, M' = M = 74$ . (Top Right) Filter:  $J' =$  where RPM Mfr = "Ludlum" and ( Month <= "Feb" and Month >= "Jan"),  $M' = 20$ . (Bottom Left) Filter:  $J'$  determined from top 5 most frequent entries,  $M' = 45$ . (Bottom Right)  $I' = \{1, 2\}$  and  $J'$  determined by the relational expression where RPM Mfr = "Ludlum" and ( Month <= "Feb" and Month >= "Jan"),  $M' = 20$ .

RPM Mfr	Location	$c[I'](x)$	$f[I'](x)$	RPM Mfr	Location	Month	$c^{J'}(x)$	$f^{J'}(x)$
Ludlum	New York	11	0.150	Ludlum	New York	Jan	1	0.050
	Seattle	24	0.325		Seattle	Jan	9	0.450
	Miami	15	0.203		Miami	Jan	2	0.100
SAIC	New York	1	0.014			Feb	8	0.400
	Seattle	10	0.136					
	Miami	13	0.176					
RPM Mfr	Location	Month	$c^{J'}(x)$	$f^{J'}(x)$	RPM Mfr	Location	$c^{J'}[I'](x)$	$f^{J'}[I'](x)$
Ludlum	Seattle	Apr	15	0.333	Ludlum	New York	1	0.050
		Jan	9	0.200		Seattle	9	0.450
	Miami	Feb	8	0.178		Miami	10	0.500
SAIC	New York	Apr	7	0.156				
	New York	Jan	6	0.133				

### 3 Conditional Views

In this section consider the filter  $J$  to be fixed, and suppress the superscript on  $f$ . We have seen that the frequencies  $f : \mathbf{X} \rightarrow [0, 1]$  represent joint probabilities  $f(\mathbf{x}) = f(x_{k^1}, x_{k^2}, \dots, x_{k^N})$ , so that from (2) and (5),  $f[I](\mathbf{x} \downarrow I)$  expresses the  $I$ -way marginal over a joint probability distribution  $f$ . Now consider two projectors  $I_1, I_2 \subseteq \mathbb{N}_N$ , so that we can define a conditional frequency  $f[I_1|I_2] : \mathbf{X} \downarrow I_1 \cup I_2 \rightarrow [0, 1]$  where  $f[I_1|I_2] := \frac{f[I_1 \cup I_2]}{f[I_2]}$ . For individual vectors, we have

$$f[I_1|I_2](\mathbf{x}) = f[I_1|I_2](\mathbf{x} \downarrow I_1 \cup I_2) := \frac{f[I_1 \cup I_2](\mathbf{x} \downarrow I_1 \cup I_2)}{f[I_2](\mathbf{x} \downarrow I_2)}.$$

$f[I_1|I_2](\mathbf{x})$  is the probability of the vector  $\mathbf{x} \downarrow I_1 \cup I_2$  restricted to the  $I_1 \cup I_2$  dimensions given that we know we can only choose vectors whose restriction to  $I_2$  is  $\mathbf{x} \downarrow I_2$ . We note that  $f[I_1|\emptyset](\mathbf{x}) = f[I_1](\mathbf{x}), f[\emptyset|I_2] \equiv 1$ , and since  $f[I_1|I_2] = f[I_1 \setminus I_2|I_2]$ , in general we can assume that  $I_1$  and  $I_2$  are disjoint.

We can now extend our concept of a view to a **conditional view**  $\mathcal{D}_{I_1|I_2, J}$  as a view on  $\mathcal{D}_{I_1 \cup I_2, J}$ , which is further equipped with the conditional frequency  $f^J[I_1|I_2]$ . Conditional views  $\mathcal{D}_{I_1|I_2, J}$  live in a different combinatorial structure than the view lattice  $\mathcal{B}$ . To describe  $I_1|I_2$  and  $J$  in a conditional view, we need three sets  $I_1, I_2 \in \mathbb{N}_N$  and  $J \in \mathbb{N}_M$  with  $I_1$  and  $I_2$  disjoint. So define  $\mathcal{A} := 3^{[N]} \times 2^M$  where  $3^{[N]}$  is a graded poset [9] with the following structure:

- $N + 1$  levels numbered from the bottom  $0, 1, \dots, N$ .
- The  $i^{\text{th}}$  level contains all partitions of each of the sets in  $\binom{[N]}{i}$ , that is the  $i$ -element subsets of  $\mathbb{N}_N$ , into two parts where

1. The order of the parts is significant, so that  $[\{1, 3\}, \{4\}]$  and  $[\{4\}, \{1, 3\}]$  of  $\{1, 3, 4\}$  are not equivalent.
  2. The empty set is an allowed member of a partition, so  $[\{1, 3, 4\}, \emptyset]$  is in the third level of  $3^{[N]}$  for  $N \geq 4$ .
- We write the two sets without set brackets and with a  $|$  separating them.
  - The partial order is given by an extended subset relation: if  $I_1 \subseteq I'_1$  and  $I_2 \subseteq I'_2$ , then  $I_1|I_2 \prec I'_1|I'_2$ , e.g.  $1\ 2|3 \prec 1\ 2\ 4|3$ .

An element in the poset  $3^{[N]}$  corresponds to an  $I_1|I_2$  by letting  $I_1$  (resp.  $I_2$ ) be the elements to the left (resp. right) of the  $|$ . We call this poset  $3^{[N]}$  because it's size is  $3^N$  and it really corresponds to partitioning  $\mathbb{N}_N$  into three disjoint sets, the first being  $I_1$ , the second being  $I_2$  and the third being  $\mathbb{N}_N \setminus (I_1 \cup I_2)$ . The structure  $3^{[2]}$  is shown in Fig. 2.

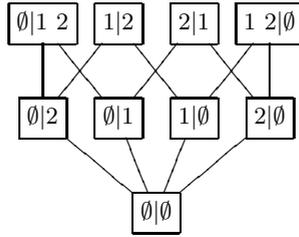


Fig. 2. The structure  $3^{[2]}$

## 4 Information Measures on Conditional Views

For a view  $\mathcal{D}_{I,J} \in \mathcal{B}$  which we identify with its frequency  $f^J[I]$ , or a conditional view  $\mathcal{D}_{I_1|I_2;J} \in \mathcal{A}$  which we identify with its conditional frequency  $f^J[I_1|I_2]$ , we are interested in measuring how “interesting” or “unusual” it is, as measured by departures from a null model. Such measures can be used for combinatorial search over the view structures  $\mathcal{B}, \mathcal{A}$  to identify noteworthy features in the data. The entropy of an unconditional view  $\mathcal{D}_{I,J}$

$$H(f^J[I]) := - \sum_{\mathbf{x} \in \mathbf{X}_{\downarrow I}} f^J[I](\mathbf{x}) \log(f^J[I](\mathbf{x})).$$

is a well-established measure of the information content of that view. A view has maximal entropy when every slot has the same expected count. Given a conditional view  $\mathcal{D}_{I_1|I_2;J}$ , we define the **conditional entropy**,  $H(f^J[I_1|I_2])$  to be the expected entropy of the conditional distribution  $f^J[I_1|I_2]$ , which operationally is related to the unconditional entropy as

$$H(f^J[I_1|I_2]) := H(f^J[I_1 \cup I_2]) - H(f^J[I_2]).$$

Given two views  $\mathcal{D}_{I,J}, \mathcal{D}_{I,J'}$  of the same dimensionality  $I$ , but with different filters  $J$  and  $J'$ , the **relative entropy** (Kullback-Leibler divergence)

$$D(f^J[I] \| f^{J'}[I]) := \sum_{\mathbf{x} \in \mathbf{X}_{\downarrow I}} f^J[I](\mathbf{x}) \log \left( \frac{f^J[I](\mathbf{x})}{f^{J'}[I](\mathbf{x})} \right)$$

is a well-known measure of the similarity of  $f^J[I]$  to  $f^{J'}[I]$ .  $D$  is zero if and only if  $f^J[I] = f^{J'}[I]$ , but it is not a metric because it is not symmetric, i.e.,  $D(f^J[I] \| f^{J'}[I]) \neq D(f^{J'}[I] \| f^J[I])$ .

$D$  is a special case of a larger class of  $\alpha$ -divergence measures between distribution introduced by Csiszár [8]. Given two probability distributions  $P$  and  $Q$ , write the density with respect to the dominating measure  $\mu = P + Q$  as  $p = dP/d(P + Q)$  and  $q = dQ/d(P + Q)$ . For any  $\alpha \in \mathbb{R}$ , the  $\alpha$ -divergence is

$$D_\alpha(P \| Q) = \int \frac{\alpha p(x) + (1 - \alpha)q(x) - p(x)^\alpha q(x)^{1-\alpha}}{\alpha(1 - \alpha)} \mu(dx).$$

$\alpha$ -divergence is convex with respect to both  $p$  and  $q$ , is non-negative, and is zero if and only if  $p = q$   $\mu$ -almost everywhere. For  $\alpha \neq 0, 1$ , the  $\alpha$ -divergence is bounded. The limit when  $\alpha \rightarrow 1$  returns the relative entropy between  $P$  and  $Q$ . There are other special cases that are of interest to us:

$$\begin{aligned} D_2(P \| Q) &= \frac{1}{2} \int \frac{(p(x) - q(x))^2}{q(x)} \mu(dx) \\ D_{-1}(P \| Q) &= \frac{1}{2} \int \frac{(q(x) - p(x))^2}{p(x)} \mu(dx) \\ D_{1/2}(P \| Q) &= 2 \int \left( \sqrt{p(x)} - \sqrt{q(x)} \right)^2 \mu(dx). \end{aligned}$$

In particular the **Hellinger metric**  $\sqrt{D_{1/2}}$  is symmetric in both  $p$  and  $q$ , and satisfies the triangle inequality. We prefer the Hellinger distance over the relative entropy because it is a bonified metric and remains bounded. In our case and notation, we have the **Hellinger distance** as

$$G(f^J[I], f^{J'}[I]) := \sqrt{\sum_{\mathbf{x} \in \mathbf{X}_{\downarrow I}} \left( \sqrt{f^J[I](\mathbf{x})} - \sqrt{f^{J'}[I](\mathbf{x})} \right)^2}.$$

## 5 Hop-Chaining View Discovery

Given our basic formalism on data views, conditional views, and information measures on them, a variety of possible user-guided navigational tasks become possible. For example, above we discussed Cariou *et al.* [5], who develop methods for discovering “drill-down paths” in data cubes. We can describe this as creating a series of views with projectors  $I_1 \supseteq I_2 \supseteq I_3$  of increasingly specified dimensional structure.

Our approach is motivated by the idea that most users will be challenged by complex views of high dimensionality, while still needing to explore many

possible data interactions. We are thus interested in restricting our users to two-dimensional views only, producing a sequence of projectors  $I_1, I_2, I_3$  where  $|I_k| = 2$  and  $|I_k \cap I_{k+1}| = 1$ , thus affecting a permutation of the variables  $X^i$ .

### 5.1 Preliminaries

We assume an arbitrary permutation of the  $i \in \mathbb{N}_N$  so that we can refer to the dimensions  $X^1, X^2, \dots, X^N$  in order. The choice of the initial variables  $X^1, X^2$  is a free parameter to the method, acting as a kind of “seed”.

One thing that is critical to note is the following. Consider a view  $\mathcal{D}_{I,J}$  which is then filtered to include only records for a particular member  $x_0^{i_0} \in X^{i_0}$  of a particular dimension  $X^{i_0} \in \mathcal{X}$ ; in other words, let  $J'$  be determined by the relational expression **where**  $X^{i_0} = x_0^{i_0}$ . Then in the new view  $\mathcal{D}'_{I,J'}, f^{J'}[I]$  is positive only on the fibers of the tensor  $\mathbf{X}$  where  $X^{i_0} = x_0^{i_0}$ , and zero elsewhere. Thus the variable  $X^{i_0}$  is effectively removed from the dimensionality of  $\mathcal{D}'$ , or rather, it is removed from the *support* of  $\mathcal{D}'$ .

Notationally, we can say  $\mathcal{D}_{I, X^{i_0}=x_0^{i_0}} = \mathcal{D}_{I \setminus \{i_0\}, X^{i_0}=x_0^{i_0}}$ . Under the normal convention that  $0 \cdot \log(0) = 0$ , our information measures  $H$  and  $G$  above are insensitive to the addition of zeros in the distribution. This allows us to compare the view  $\mathcal{D}_{I, X^{i_0}=x_0^{i_0}}$  to any other view of dimensionality  $I \setminus \{i_0\}$ .

This is illustrated in Table 3 through our continuing example, now with the filter **where** `Location="Seattle"`. Although formally still an RPM Mfr  $\times$  Location  $\times$  Month cube, in fact this view lives in the RPM Mfr  $\times$  Month plane, and so can be compared to the RPM Mfr  $\times$  Month marginal.

**Table 3.** Our example data tensor from Table 1 under the filter **where** `Location="Seattle"`;  $M' = 34$

RPM Mfr	Location	Month	$c(\mathbf{x})$	$f(\mathbf{x})$
Ludlum	Seattle	Jan	9	0.265
		Apr	15	0.441
SAIC		Feb	4	0.118
		Mar	3	0.088
		Apr	3	0.088

Finally, some caution is necessary when the relative entropy  $D(f^J[I] \| f^{J'}[I])$  or Hellinger distance  $G(f^J[I], f^{J'}[I])$  is calculated from data, as their magnitudes between empirical distributions is strongly influenced by small sample sizes. To counter spurious effects, we supplement each calculated entropy with the probability that under the null distribution that the row has the same distribution as the marginal, of observing an empirical entropy larger or equal to actual value. When that probability is large, say greater than 5%, then we consider its value spurious and set it to zero before proceeding with the algorithm.

## 5.2 Method

We can now state the hop-chaining methodology.

1. Set the initial filter to  $J = \mathbb{N}_M$ . Set the initial projector  $I = \{1, 2\}$ , determining the initial view  $f^J[I]$  as just the initial  $X^1 \times X^2$  grid.
2. For each row  $x_{k^1} \in X^1$ , we have the marginal distribution  $f^{X^1=x_{k^1}}[I]$  of that individual row, using the superscript to indicate the relational expression filter. We also have the marginal  $f^J[I \setminus \{X^1\}]$  over all the rows for the current filter  $J$ . In light of the discussion just above, we can calculate all the Hellinger distances between each of the rows and this row marginal as

$$G(f^{X^1=x_{k^1}}[I], f^J[I \setminus \{X^1\}]) = G(f^{X^1=x_{k^1}}[I \setminus \{X^1\}], f^J[I \setminus \{X^1\}]),$$

and retain the maximum row value  $G^1 := \max_{x_{k^1} \in X^1} G(f^{X^1=x_{k^1}}[I], f^J[I \setminus \{X^1\}])$ . We can dually do so for columns against the column marginal:

$$G(f^{X^2=x_{k^2}}[I], f^J[I \setminus \{X^2\}]) = G(f^{X^2=x_{k^2}}[I \setminus \{X^2\}], f^J[I \setminus \{X^2\}]),$$

retaining the maximum value  $G^2 := \max_{x_{k^2} \in X^2} G(f^{X^2=x_{k^2}}[I], f^J[I \setminus \{X^2\}])$ .

3. The user is prompted to select either a row  $x_0^1 \in X^1$  or a column  $x_0^2 \in X^2$ . Since  $G^1$  (resp.  $G^2$ ) represents the row (column) with the largest distance from its marginal, perhaps selecting the global maximum  $\max(G^1, G^2)$  is most appropriate; or this can be selected automatically. Letting  $x'_0$  be the selected value from the selected variable (row or column)  $i' \in I$ , then  $J'$  is set to **where**  $X^{i'} = x'_0$ , and this is placed in the *background* filter.
4. Let  $i'' \in I$  be the variable *not* selected by the user, so that  $I = \{i', i''\}$ .
5. For each dimension  $i''' \in \mathbb{N}_N \setminus I$ , that is, for each dimension which is neither in the background filter  $R^b = \{i'\}$  nor retained in the view through the projector  $\{i''\}$ , calculate the conditional entropy of the retained view  $f^{J'}[\{i''\}]$  against that variable:  $H(f^{J'}[\{i''\}|\{i'''\}])$ .
6. The user is prompted to select a new variable  $i''' \in \mathbb{N}_N \setminus I$  to add to the projector  $\{i''\}$ . Since  $\underset{i''' \in \mathbb{N}_N \setminus I}{\operatorname{argmin}} H(f^{J'}[\{i''\}|\{i'''\}])$  represents the variable with the most constraint against  $i''$ , that may be the most appropriate selection, or it can be selected automatically.
7. Let  $I' = \{i'', i'''\}$ . Note that  $I'$  is a sibling to  $I$  in  $\mathcal{B}^N$ , thus the name ‘‘hop-chaining’’.
8. Let  $I', J'$  be the new  $I, J$  and go to step 2.

Keeping in mind the arbitrary permutation of the  $X^i$ , then the repeated result of applying this method is a sequence of hop-chaining steps in the view lattice, building up an increasing background filter:

1.  $I = \{1, 2\}, J = \mathbb{N}_M$
2.  $I' = \{2, 3\}, J' = \mathbf{where} X^1 = x_0^1$
3.  $I'' = \{3, 4\}, J'' = \mathbf{where} X^1 = x_0^1, X^2 = x_0^2$
4.  $I''' = \{4, 5\}, J''' = \mathbf{where} X^1 = x_0^1, X^2 = x_0^2, X^3 = x_0^3$

### 5.3 Extension to Hierarchical Data Cubes

In Sec. 2.1 we introduced data cubes as flat data tensors, while in general in OLAP the dimensions are hierarchically structured. While a full development of a hierarchical approach to hop-chaining awaits future work, we can indicate those directions here. First, we extend our definition of a data cube to be  $\mathcal{D} := \langle \mathbf{X}, \mathcal{X}, \mathcal{Q}, c \rangle$ , where now  $\mathcal{Q} = \{\mathcal{P}^i\}_{i=1}^N$  is a collection of  $N$  partially-ordered hierarchical [9] **dimensions**  $\mathcal{P}^i = \langle P^i, \leq^i \rangle$  with **members**  $p^i \in P^i$ . Each partially-ordered set (poset)  $\mathcal{P}^i$  is isomorphic to a sub-poset of the Boolean lattice  $\mathcal{B}^i = \langle 2^{X^i}, \subseteq \rangle$  which is the power set of the values of the variable  $X^i$  ordered by set inclusion. While in principle each poset  $\mathcal{P}^i$  could be as large as  $\mathcal{B}^i$ , in practice, they are trees, with  $X^i \in P^i$  and  $\forall x^i \in X^i, \{x^i\} \in P^i$ .

We can identify the **cube schema** as  $\mathbf{P} := \times_{i=1}^N \mathcal{P}^i$ , so that each  $\mathbf{p} \in \mathbf{P}$  is a **cell**. We then have the hierarchical count function  $\hat{c} : \mathbf{P} \rightarrow \mathbb{N}$ , where  $\hat{c}(\mathbf{p}) := \sum_{\mathbf{x} \leq \mathbf{p}} c(\mathbf{x})$ , and  $\leq := \times_{i=1}^N \leq^i$ , the product order of the hierarchies. There is also the hierarchical frequency function  $\hat{f} : \mathbf{P} \rightarrow [0, 1]$ , with  $\hat{f}(\mathbf{p}) := \frac{\hat{c}(\mathbf{p})}{M}$ . In a real OLAP view, the entries actually correspond not to slots  $\mathbf{x} \in \mathbf{X}$ , but to cells  $\mathbf{p} \in \mathbf{P}$ ; and the rows and columns not to collections of data items  $Y^i \subseteq X^i$ , but of members  $Q^i \subseteq P^i$ . If  $X^1 = \text{“Location”}$ , and  $p_0^1 = \text{“California”} \in P^1$ , then classical drilldown might take a row like **California** from a view, restrict  $J$  with the relational expression **where Location = California**, and then replace  $Q^1$  with all the children of  $p_0^1$ , so that  $Q^{1'} = \{p^1 \leq^1 p_0^1\}$ .

We are experimenting with view discovery and hop-chaining formalisms which operate over these member collections  $Q^i$ , and in general over their Cartesian products  $\times_{i \in I} Q^i \subseteq \mathbf{P} \downarrow I$ . But in the current formulation, it is sufficient to consider each dimension  $X^i$  involved in a foreground view to be drilled-down to the immediate children of the top of  $\mathcal{P}^i$ , that is, the children of **All**.

## 6 Implementation

We have implemented the hop-chaining methodology in a prototype form for experimentation and proof-of-principle. ProClarity 6.2 is used in conjunction with SQL Server Analysis Services (SSAS) 2005 and the R statistical platform v. 2.7<sup>6</sup>. ProClarity provides a flexible and friendly GUI environment with extensive API support which we use to gather current display contents and query context for row, column and background filter selections. R is currently used in either batch or interactive mode for statistical analysis and development. Microsoft Visual Studio .Net 2005 is used to develop plug-ins to ProClarity to pass ProClarity views to R for hop-chain calculations.

A first view of the data set used for demonstrating this method is shown in Fig. 3, a screenshot from the ProClarity tool. The database is a collection of 1.9M records of events of RPM events. The 15 available dimensions are shown on the left of the screen (e.g. “day of the month”, “RPM hierarchy”), tracking such

<sup>6</sup> <http://www.r-project.org>

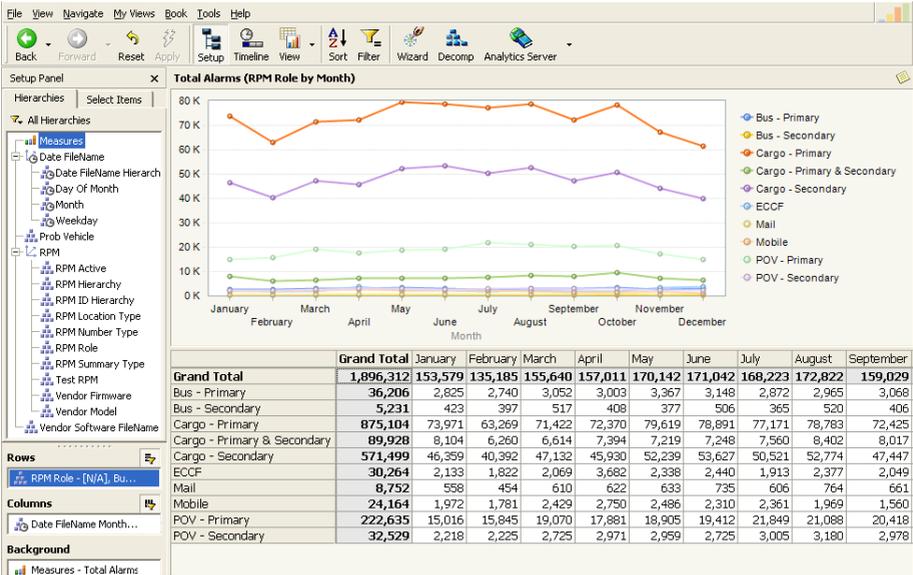


Fig. 3. Initial 2-D view of the alarm summary data cube, showing count distribution of RPM Role by months

things as the identities and characteristics of particular RPMs, time information about events, and information about the hardware, firmware, and software used at different RPMs.

## 7 Examples

Space limitations will allow showing only a single step for the hop-chaining procedure against the alarm summary data cube.

Fig. 3 shows the two-dimensional projection of the  $X^1 = \text{“RPM Role”} \times X^2 = \text{“Month”}$  dimensions within the 15-dimensional overall cube, drilled down to the first level of the hierarchies (see Sec. 5.3). Its plot shows the distributions of

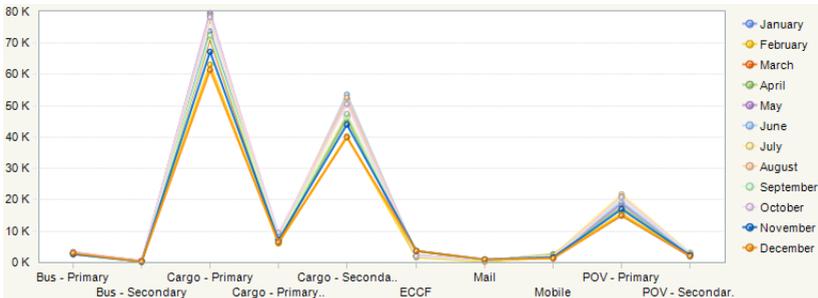


Fig. 4. Count distribution of months

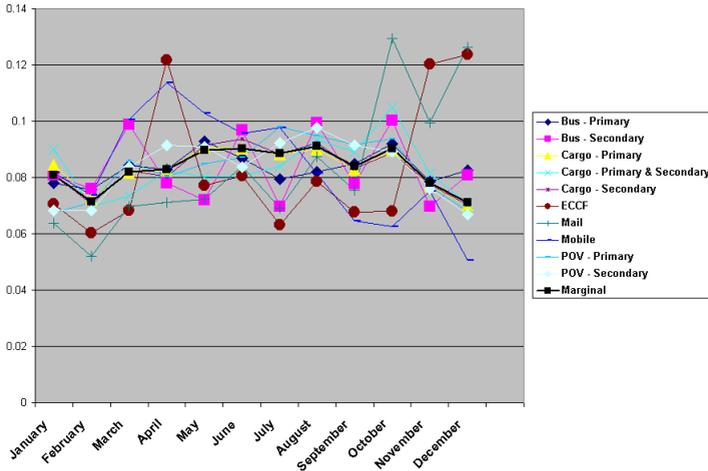


Fig. 5. Frequency distributions of RPM roles

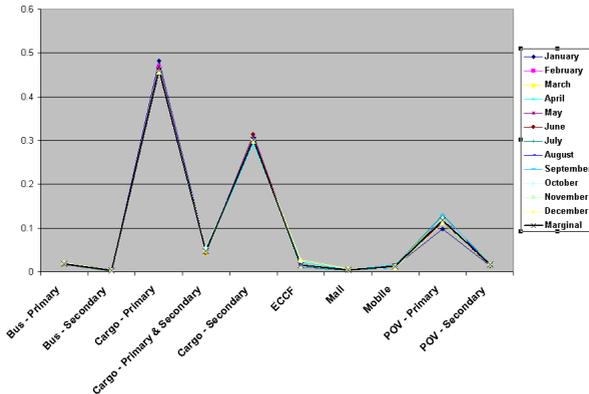
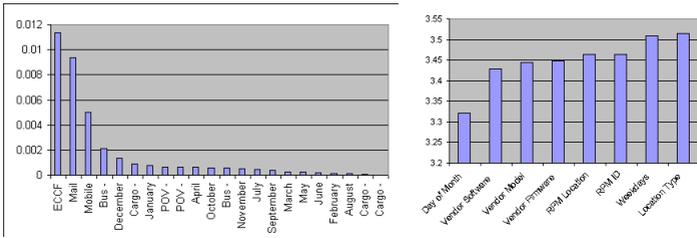


Fig. 6. Frequency distributions of months

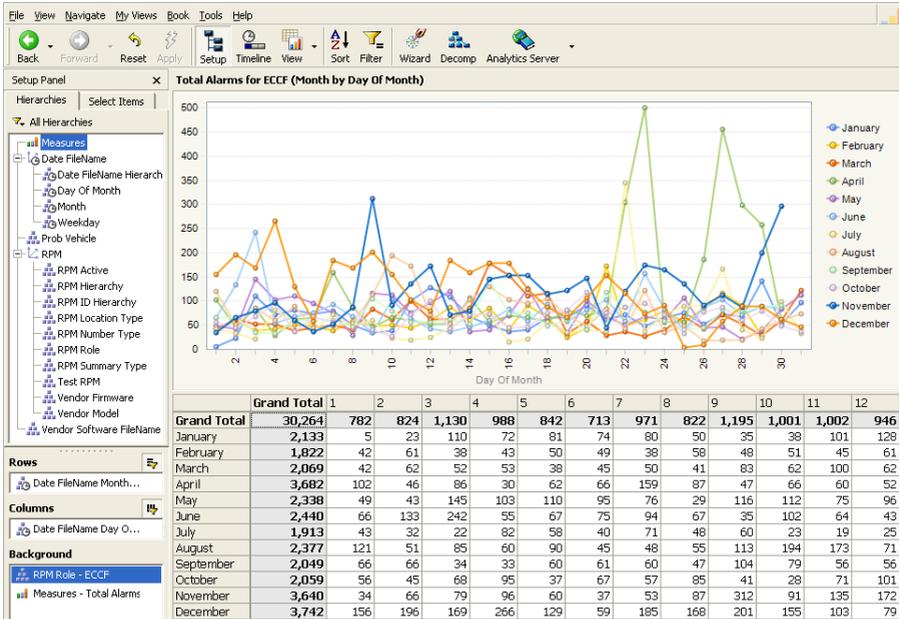
count  $c$  of alarms by RPM role (Busses Primary, Cargo Secondary, etc.)  $X^1$ , while Fig. 4 shows the distribution by Month  $X^2$ .

The distributions for roles seem to vary at most by overall magnitude, rather than shape, while the distributions for months appear almost identical. However, Fig. 5 and Fig. 6 show the same distributions, but now in terms of their frequencies  $f$  relative to their corresponding marginals, allowing us to compare the shapes of the distributions normalized by their absolute sizes. While the months still seem identical, the RPM roles are clearly different, although it is difficult to see which is most unusual with respect to the marginal (bold line).

The left side of Fig. 7 shows the Hellinger distances  $G(f^{X^i=x_{ki}}[I], f^J[I \setminus \{X^i\}])$  for  $i \in \{1, 2\}$  for each row or column against its marginal. The RPM roles “ECCF” and “Mail” are clearly the most significant, which can be verified by examining the anomalously shaped plots in Fig. 5. The most significant



**Fig. 7.** (Left) Hellinger distances of rows and columns against their marginals. (Right) Relative entropy of months against each other significant dimension, given that RPM Role = ECCF.



**Fig. 8.** Subsequent view on the  $X^2 = \text{Months} \times X^3 = \text{Day of Month}$  projector. Note the new background filter where RPM Role = ECCF.

month is December, although this is hardly evident in Fig. 6. We select the maximal row-wise Hellinger value  $G^1 = .011$  for ECCF, so that  $i' = 1, x_0^1 = \text{ECCF}$ .  $X^{i'} = X^1 = \text{“RPM Role”}$  is added to the background filter,  $X^{i''} = X^2 = \text{Months}$  is retained in the view, and we calculate  $H(f^{j'}[\{2\}|\{i'''\}])$  for all  $i''' \in \{3, 4, \dots, 15\}$ , which are shown on the right of Fig. 7 for all significant dimensions. On that basis  $X^3$  is selected as Day of Month with minimal  $H = 3.22$ .

The subsequent view for  $X^2 = \text{Months} \times X^3 = \text{Day of Month}$  is then shown in Fig. 8. Note the strikingly divergent plot for April: it in fact does have the

highest Hellinger distance at .07, an aspect which is completely invisible from the overall initial view, e.g. in Fig. 5.

## 8 Discussion, Analysis, and Future Work

In this paper, we have provided the fundamentals necessary to express view discovery in OLAP databases as a combinatorial search and optimization operation in general, aside from the specific hop-chaining method. What remains to be addressed is a precise formal expression of this optimization problem. This is dependent on the mathematical properties of our information measures  $H$ ,  $D$ , and  $G$  over the lattices  $\mathcal{B}$ ,  $\mathcal{A}$ . It is well known, for example, that  $H$  is a monotonic function in  $\mathcal{B}$  [9], in that  $\forall I_1 \subseteq I_2, H(f^J[I_1]) \geq H(f^J[I_2])$ . There should be ample literature (e.g. [27]) to fully explicate the behavior of these functions on these structures, and guide development of future search algorithms.

Also as mentioned above, we are restricting our attention to OLAP cubes with a single “count” measure. Frequency distributions are available from other quantitative measures, and exploring the behavior of these algorithms in those contexts is of interest.

Information measures are used because of their mathematical properties in identifying unusual distributions. It remains to be demonstrated that these measures have high utility for real analysts of these databases, and which mix of statistical measures, whether including our precise hop-chaining algorithm or not, is ideal for their needs.

The value of our implementation within commercial front-end database tools provides the opportunity for this validation. Generally, software implementations provide a tremendous value in performing this research, not only for this practical validation by sponsors and users, but also for assisting with the methodological development itself. As our software platform matures, we look forward to incorporating other algorithms for view discovery [5,16,20,23,24,25], for purposes of comparison and completeness.

## Acknowledgements

This work is supported in part by the U.S. Department of Homeland Security under Grant Award 2007-ST-104-000006. The views and conclusions contained in this document are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of the U.S. Department of Homeland Security. We are grateful for the support of Judith Cohn at the Los Alamos National Laboratory; David Gillen, Eric Stephan, Bryan Olsen, and Elena Peterson at the Pacific Northwest National Laboratory; Matt Damante and Robert Fernandez at the Lawrence Livermore National Laboratory; and our sponsors and colleagues at DHS/S&T, DHS/ICE, and DHS/CBP.

## References

1. Agrawal, R., Gupta, A., Sarawagi, S.: Modeling Multidimensional Databases. In: Proc. 13th Int. Conf. on Data Engineering (1997)
2. Asimov, D.: The Grand Tour: A Tool for Viewing Multidimensional Data. *SIAM J. Statistical Computing* 6, 1 (1985)
3. Barak, B., Chaudhuri, K., Dwork, C., Kale, S., McSherry, F., Talwar, K.: Privacy, Accuracy, and Consistency Too: A Holistic Solution to Contingency Table Release. In: Proc. 2007 Conf. Principles of Database Systems (PODS 2007) (2007)
4. Borgelt, C., Kruse, R.: Graphical Models. Wiley, New York (2002)
5. Cariou, V., Cubillé, J., Derquenne, C., Goutier, S., Guisnel, F., Klajnmic, H.: Built-In Indicators to Discover Interesting Drill Paths in a Cube. In: Song, I.-Y., Eder, J., Nguyen, T.M. (eds.) *DaWaK 2008. LNCS*, vol. 5182, pp. 33–44. Springer, Heidelberg (2008)
6. Chaudhuri, S., Umeshwar, D.: An Overview of Data Warehousing and OLAP Technology. *ACM SIGMOD Record* 26(1), 65–74 (1997)
7. Codd, E.F., Codd, S.B., Salley, C.T.: Providing OLAP (On-Line Analytical Processing) to User-Analysts: An IT Mandate” (1993), [www.cs.bgu.ac.il/~dbm031/dw042/Papers/olap\\_to\\_useranalysts\\_wp.pdf](http://www.cs.bgu.ac.il/~dbm031/dw042/Papers/olap_to_useranalysts_wp.pdf)
8. Csiszár, I.: Information-type measures of divergence of probability distributions and indirect observations. *Studia Sci. Math. Hung* 2, 299–318 (1967)
9. Davey, B.A., Priestly, H.A.: Introduction to Lattices and Order, 2nd edn. Cambridge University Press, Cambridge (1990)
10. Gyssens, M., Lakshmanan, L.V.S.: A Foundation for Multi-Dimensional Databases. In: Proc. 23rd VLDB Conf., pp. 106–115 (1997)
11. Joslyn, C.A., Gillen, D., Fernandes, R., Damante, M., Burke, J., Critchlow, T.: Hybrid Relational and Link Analytical Knowledge Discovery for Law Enforcement. In: 2008 IEEE Int. Conf. on Technologies for Homeland Security (HST 2008), pp. 161–166. IEEE, Piscataway (2008)
12. Joslyn, C., Mniszejski, S.: DEEP: Data Exploration through Extension and Projection. Los Alamos Technical Report LAUR 02-1330 (2002)
13. Klir, G., Elias, D.: Architecture of Systems Problem Solving, 2nd edn. Plenum, New York (2003)
14. Kolda, T.G., Bader, B.W.: Tensor Decompositions and Applications. *SIAM Review* (in press) (2008)
15. Krippendorff, K.: Information Theory: Structural Models for Qualitative Data. Sage Publications, Newbury Park (1986)
16. Kumar, N., Gangopadhyay, A., Bapna, S., Karabatis, G., Chen, Z.: Measuring Interestingness of Discovered Skewed Patterns in Data Cubes. *Decision Support Systems* 46, 429–439 (2008)
17. Lauritzen, S.L.: Graphical Models, Oxford UP (1996)
18. Malvestuto, F.M.: Testing Implication of Hierarchical Log-Linear Models for Probability Distributions. *Statistics and Computing* 6, 169–176 (1996)
19. Montanari, A., Lizzani, L.: A Projection Pursuit Approach to Variable Selection. *Computational Statistics and Data Analysis* 35, 463–473 (2001)
20. Palpanas, T., Koudas, N.: Using Database Aggregates for Approximating Querying and Deviation Detection. *IEEE Trans. Knowledge and Data Engineering* 17(11), 1–11 (2005)
21. Pedersen, T.B., Jensen, C.S.: Multidimensional Database Technology. *IEEE Computer* 34(12), 40–46 (2001)

22. Sarawagi, S.: Explaining Differences in Multidimensional Aggregates. In: Proc. 25th Int. Conf. Very Large Databases, VLDB 1999 (1999)
23. Sarawagi, S.: User-Adaptive Exploration of Multidimensional Data. In: Proc. 26th Very Large Database Conf., pp. 307–316 (2000)
24. Sarawagi, S.: iDiff: Informative Summarization of Differences in Multidimensional Aggregates. *Data Mining and Knowledge Discovery* 5, 255–276 (2001)
25. Sarawagi, S., Agrawal, R., Megiddo, N.: Discovery-driven Exploration of OLAP Data Cubes. In: Schek, H.-J., Saltor, F., Ramos, I., Alonso, G. (eds.) EDBT 1998. LNCS, vol. 1377, pp. 168–182. Springer, Heidelberg (1998)
26. Soshani, A.: OLAP and Statistical Databases: Similarities and Differences. In: Proc. PODS 1997, pp. 185–196 (1997)
27. Studeny, M.: Probabilistic Conditional Independence Structures. Springer, London (2005)
28. Thomas, H., Datta, A.: A Conceptual Model and Algebra for On-Line Analytical Processing in Decision Support Databases. *Information Systems Research* 12(1), 83–102 (2001)