

# Cardinality-Based Inference Control in Sum-Only Data Cubes<sup>\*</sup>

Lingyu Wang, Duminda Wijesekera, and Sushil Jajodia

Center for Secure Information Systems, George Mason University  
Fairfax, VA 22030-4444, USA  
{lwang3,dwijesek,jajodia}@gmu.edu

**Abstract.** This paper addresses the inference problems in data warehouses and decision support systems such as on-line analytical processing (OLAP) systems. Even though OLAP systems restrict user accesses to predefined aggregations, inappropriate disclosure of sensitive attribute values may still occur. Based on a definition of non-compromiseability to mean that any member of a set of variables satisfying a given set of their aggregations can have more than one value, we derive sufficient conditions for non-compromiseability in sum-only data cubes. Under this definition, (1) the non-compromiseability of multi-dimensional aggregations can be reduced to that of one dimensional aggregations, (2) full or dense core cuboids are non-compromiseable, and (3) there is a tight lower bound for the cardinality of a core cuboid to remain non-compromiseable. Based on these results, taken together with a three-tier model for controlling inferences, we provide a divide-and-conquer algorithm that uniformly divides data sets into chunks and builds a data cube on each such chunk. The union of these data cubes are then used to provide users with inference-free OLAP queries.

## 1 Introduction

Decision support systems such as On-line Analytical Processing (OLAP) are becoming increasingly important in industry. These systems are designed to answer queries involving large amounts of data and their statistical averages in near real time. It is well known that access control alone is insufficient in eliminating all forms of disclosures, as information not released directly may be inferred indirectly from answers to legitimate queries. This is known as *inference problem*. An OLAP query typically consists of multiple aggregations, and hence vulnerable to unwanted inferences. Providing inference-free answers to sum-only data cube style OLAP queries while not adversely impacting performance or restricting availability in an OLAP system is the subject matter of this paper.

The inference problem has been investigated since 70's and many inference control methods have been proposed for statistical databases. However, most of

---

<sup>\*</sup> This work was partially supported by the National Science Foundation under grant CCR-0113515.

these methods become computationally infeasible if directly applied to OLAP systems. OLAP applications demand short response times although queries usually aggregate large amounts of data [20]. Because most existing inference control algorithms have run times proportional to the size of queries or aggregated data sets, their impact upon performance renders them impractical for OLAP systems.

While arbitrary queries are common in statistical databases, OLAP queries usually comprise of well-structured operations such as group-by, cross-tab and sub-totals. These operations can conveniently be integrated with various data cube operations, such as slicing-dicing, roll up and drill down [19]. We will show that the limited formats and predictable structures of OLAP queries as well as the multi-dimensional hierarchical data model of OLAP datasets can be exploited to simplify inference control.

Table 1 shows a small two-dimensional data set consisting of monthly employee salaries. Individual salaries should be hidden from users, and hence have been replaced with the symbol ?. The symbol N/a denotes null value for inapplicable combinations of months and employees that are known to users.<sup>1</sup> Assume subtotals are allowed to be released to users through OLAP queries. Inference problem occurs if any of the values represented by symbol ? can be derived from the subtotals. No value in the first two quarters can be inferred, because infinitely many different values may satisfy each ? symbol with the subtotals satisfied. For the third quarter, Mary’s salary in September can be inferred from the subtotals of September salaries because she is the only employee with a valid salary for September. For the fourth quarter, by subtracting Bob’s and Jim’s fourth quarter salaries (\$4300 and \$3000 respectively) from the subtotals in October and November (\$7100 and \$4100 respectively) Alice’s salary for October can be computed to be \$3900.

Based on a definition of non-compromiseability to mean that any member of a set of variables satisfying a given set of their aggregations can have more than one value<sup>2</sup>, we derive sufficient conditions for non-compromiseability in sum-only data cubes: (1) the non-compromiseability of multi-dimensional aggregations can be reduced to that of one dimensional aggregations, (2) full or dense core cuboids are non-compromiseable, and (3) there is a tight lower bound on the cardinality of a core cuboid for it to remain non-compromiseable. Based on these results, and a three-tier model for controlling inferences, we provide a divide-and-conquer algorithm that uniformly divides data sets into chunks and builds a data cube on each such chunk. The union of these data cubes are then used to provide users with inference-free OLAP queries.

The rest of the paper is organized as follows. Section 2 reviews existing inference control methods proposed in statistical databases and OLAP systems. Section 3 formalizes sum-only data cube and proves sufficient conditions for its

---

<sup>1</sup> In general, data values are known through various kinds of *external knowledge* (knowledge obtained through channels other than queries.)

<sup>2</sup> In the settings of this paper, each variable can have either one value or infinitely many values.

non-compromiseability. On the basis of a three-tier model these conditions are integrated into an inference control algorithm in Section 4. Section 5 concludes the paper.

## 2 Related Work

Inference control has been extensively studied in statistical databases as surveyed in [14, 1, 15]. Inference control methods proposed in statistical databases are usually classified into two main categories; *restriction based* techniques and *perturbation based* techniques. Restriction based techniques [18] include restricting the size of a *query set* (i.e., the entities that satisfy a single query), overlap control [16] in query sets, auditing all queries in order to determine when inferences are possible [11, 8, 22, 24], suppressing sensitive data in released tables containing statistical data [12], partitioning data into mutually exclusive sets [9, 10], and restricting each query to range over at most one partition. Perturbation based technique includes adding noise to source data [29], outputs [5, 25], database structure [27], or size of query sets (by sampling data to answer queries) [13]. Some variations of the inference problem have been studied lately, such as the inferences caused by arithmetic constraints [7, 6], inferring approximate values instead of exact values [24] and inferring intervals enclosing exact values [23].

The inference control methods proposed for statistical databases generally sacrifice efficiency for the control of inferences caused by arbitrary queries, which

**Table 1.** An example data cube

Quarter	Month	Alice	Bob	Jim	Mary	Sub Total
1	January	?	?	?	?	5500
	February	?	?	?	?	5500
	March	?	?	?	?	5500
	Sub Total	3000	3000	4500	6000	
2	April	?	?	?	?	6100
	May	?	N/a	?	?	6100
	June	?	?	?	?	4100
	Sub Total	4500	3300	4500	4000	
3	July	?	?	?	?	6100
	August	?	?	?	?	6100
	September	N/a	N/a	N/a	?	2000
	Sub Total	3500	2200	2500	6000	
4	October	?	?	?	N/a	7100
	November	N/a	?	?	N/a	4100
	December	?	N/a	N/a	?	4100
*	Bonus	?	N/a	N/a	?	6000
	Sub Total	7000	4300	3000	7000	

is essential for general purpose databases. However, this sacrifice is not desirable for OLAP systems, because in OLAP systems near real time response takes priority over the generality of answerable queries. Hence most of these methods are computationally infeasible in OLAP systems. As an example, Audit Expert [11] models inference problem with a linear system  $Ax = b$  and detects the occurrence of inference by transforming the  $m$  by  $n$  matrix  $A$  (corresponding to  $m$  queries on  $n$  attribute values) to its reduced row echelon form. The transformation has a well-known complexity of  $O(m^2n)$ , which is prohibitive in the context of data warehouses and OLAP systems since  $m$  and  $n$  can be as large as a million.

Our work shares similar motivation with that of [16], i.e., to efficiently control inference with the cardinality of data and queries, which can be easily obtained, stored and maintained. Dobkin et al. gave sufficient conditions for the non-compromiseability of arbitrary sum only queries [16]. The conditions are based on the smallest number of queries that suffices to compromise the individual data. Our work addresses multi-dimensional data cube queries. The fact that data cube queries are a special case of arbitrary queries implies better results.

To the best of our knowledge, inference control for OLAP systems and data warehouses are limited to [3, 2, 17, 26]. They all attempt to perturb sensitive values while preserving the data distribution model, such that classification (or association rules) results obtained before and after the perturbation do not change. These approaches are application-specific, that is, the sole purpose of data analysis is limited to classification (or association rule mining). We do not have this restriction. Moreover, we do not use perturbation in this paper.

### 3 Cardinality-Based Non-compromiseability Criteria for Data Cubes

This section defines our model for sum-only data cubes and formalizes compromiseability. We then derive cardinality-based sufficient conditions for non-compromiseability based on the model and definitions.

#### 3.1 Problem Formulation

In our model, a  $k$ -dimensional *data cube* consists of one *core cuboid* and several *aggregation cuboids*. In addition, we use an *aggregation matrix* to abstract the relationship between them. Each *dimension* is modeled as a closed integer interval. The Cartesian product of the  $k$  dimensions is referred to as *full core cuboid*. Each integer vector in the full core cuboid is a *tuple*. A *core cuboid* is a subset of the full core cuboid, which contains at least one tuple for every value chosen from every dimension. This condition allows us to uniquely identify the size of each dimension for a given core cuboid. Definition 1 formalizes these concepts.

#### Definition 1 (Core Cuboids and Slices).

Given a set of  $k$  integers  $D_i$  satisfying  $D_i > 1$  for all  $1 \leq i \leq k$ . A  $k$ -dimensional core cuboid is any subset  $S$  of  $\Pi_{i=1}^k [1, D_i]$  satisfying the property that for any  $x_i \in [1, D_i]$  there exist  $(k - 1)$  integers  $x_j \in [1, D_j]$  for all

$1 \leq j \leq k$  and  $j \neq i$ , such that  $(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_k) \in S$ .  $C_c$  denotes a core cuboid. Each vector  $t \in C_c$  is referred to as a tuple. Further, the  $i^{\text{th}}$  element of vector  $t \in C_c$ , denoted by  $t[i]$ , is referred to as the  $i^{\text{th}}$  dimension of  $t$ . We say that  $\Pi_{i=1}^k [1, D_i]$  is the full core cuboid denoted by  $C_f$ . We say a tuple  $t$  is missing from the core cuboid  $C_c$  if  $t \in C_f \setminus C_c$ . The subset of  $C_c$  defined by  $\{t \mid t \in C_c, t[i] = j\}$  for each  $j \in [1, D_i]$  is said to be the  $j^{\text{th}}$  slice of  $C_c$  on the  $i^{\text{th}}$  dimension, denoted by  $P_i(C_c, j)$ . If  $P_i(C_c, j) = \{t \mid t \in C_f, t[i] = j, j \in [1, D_i]\}$ , we say that  $P_i(C_c, j)$  is a full slice.

For example, the fourth quarter salaries in the sample data cube of Table 1 is depicted in Table 2. It has two dimensions: month (dimension 1) and employee name (dimension 2). Both dimensions have four different values that are mapped to the integer interval  $[1, 4]$ . Therefore, the full core cuboid  $C_f$  is  $[1, 4] \times [1, 4]$ . The core cuboid  $C_c$  contains nine tuples and seven missing tuples are shown as N/a.

To define aggregations of a data cube, we follow [19] to augment each dimension with a special value *ALL*, for which we use the symbol  $*$ . Each *aggregation vector* is similar to a tuple except that it is formed with the augmented dimensions. An aggregation vector selects a set of tuples in core cuboids with its  $*$  values, which form its *aggregation set*. All aggregation vectors having  $*$  value in the same dimensions form an *aggregation cuboid*. The concepts of aggregation vector, aggregation cuboid and aggregation set are formalized in Definition 2.

**Definition 2 (j-\* Aggregation Vectors, Cuboids and Data Cubes).**

A  $j$ -\* aggregation vector  $t$  is a  $k$  dimensional vector satisfying  $t \in \Pi_{i=1}^k ([1, D_i] \cup \{*\})$  and  $|\{i : t[i] = * \text{ for } 1 \leq i \leq k\}| = j$ . If  $t[i] = *$ , then we say that the  $i^{\text{th}}$  element is a  $*$ -elements, and others are called non  $*$ -elements. A  $j$ -\* aggregation cuboid is a set of aggregation vectors  $C$  such that for any  $t, t' \in C$ ,  $\{i : t[i] = *\} = \{i : t'[i] = *\}$  and  $|\{i : t[i] = *\}| = j$ . The aggregation set of an  $j$ -\* aggregation vector  $t$  is defined as  $\{t' : t' \in C_c \text{ such that } t'[i] = t[i], \forall i \text{ } t[i] \neq *\}$ . We use the notation  $Qset(t)$  for the aggregation set of  $t$ . The aggregation set of a set of aggregation vectors  $S$  is defined as the union of  $Qset(t)$  for all  $t \in S$ . We use notation  $Qset(S)$  for the aggregation set of  $S$ . A data cube is defined as a pair  $\langle C_c, S_{all} \rangle$ , where  $C_c$  is a core cuboid, and  $S_{all}$  is the set of all  $j$ -\* aggregation cuboids, for all  $1 \leq j \leq k$ .

For example, subtotals of the fourth quarter data in the data cube of Table 1 is depicted in Table 2. Each subtotal is represented as an aggregation vector with  $*$  value. For example,  $(1, *)$  represents the subtotal in October. The aggregation set of  $(1, *)$  is  $\{(1,1), (1,2), (1,3)\}$ . The set of four aggregation vectors  $\{(1, *), (2, *), (3, *), (4, *)\}$  form an aggregation cuboid because all of them have  $*$  value in the second dimension.

To abstract the relationship between a core cuboid and its aggregation cuboids in a given data cube, we define a binary matrix referred to as *aggregation matrix*. Each element of an aggregation matrix is associated with a tuple and an aggregation vector. An element 1 means the tuple is contained in the aggregation set of the aggregation vector, 0 otherwise. We assign the tuples in  $C_f$

**Table 2.** Illustration of data cube

	1 (Al)	2 (Bob)	3 (Jim)	4 (Ma)	5 (SubT)
1 (Oct)	(1,1)	(1,2)	(1,3)	N/a	(1,*)
2 (Nov)	N/a	(2,2)	(2,3)	N/a	(2,*)
3 (Dec)	(3,1)	N/a	N/a	(3,4)	(3,*)
4 (Bonus)	(4,1)	N/a	N/a	(4,4)	(4,*)
5 (SubT)	(*,1)	(*,2)	(*,3)	(*,4)	(*,*)

and any  $C$  in dictionary order, the aggregation cuboids in  $S_{all}$  in ascending order on the number of \*-elements and descending order on the index of the \*-element. This assignment enables us to refer to the  $i^{th}$  tuple in  $C_f$  as  $C_f[i]$  (similarly for  $C_c$ ,  $S_{all}$  or their subsets). We use  $M[i, j]$  for the  $(i, j)^{th}$  element of matrix  $M$ . Aggregation matrix is formalized in Definition 3.

**Definition 3 (Aggregation Matrix).**

The aggregation matrix of the aggregation cuboid  $C$  on the core cuboid  $C_c$  is defined as the following  $(m \times n)$  matrix  $M_{C_c, C}$  ( or simply  $M$  when  $C_c$  and  $C$  are clear from context).

$$M_{C_c, C}[i, j] = \begin{cases} 1, & \text{if } C_f[j] \in Qset(C[i]); \\ 0, & \text{otherwise.} \end{cases}$$

We define the aggregation matrix of  $S$  on  $C_c$  as the row block matrix with the  $i^{th}$  row block as the aggregation matrix of the  $i^{th}$  aggregation cuboid in  $S$ .

We use  $S_1$  to represent the set of all 1-\* aggregation cuboids for a given  $C_c$ , and  $M_1$  the aggregation matrix of  $S_1$  on  $C_c$  (that is  $M_{C_c, S_1}$ ), referred to as the 1-\* aggregation matrix.

Aggregation matrix and compromiseability are illustrated in Table 3. By representing individual salaries with variables  $x_i$ , we get a linear system  $M_{C_c, S_1} \cdot \vec{X} = \vec{B}$ . It has at least one solution, since  $\vec{B}$  is calculated from the “real” salaries, which must satisfy the stated linear system. From linear algebra [21], each  $x_i$  can have either a unique value or infinitely many different values among all the solutions to  $M_{C_c, S_1} \cdot \vec{X} = \vec{B}$ . This depends on  $M_{C_c, S_1}$  but not on  $\vec{B}$  (this is not valid if additional knowledge about  $\vec{X}$  is known to users, for example, salaries are non-negative [23, 24, 22]). If an  $x_i$  has a unique value among all the solutions, then clearly the sensitive value represented by  $x_i$  was compromised. In this example,  $x_1$  has the value of 3900 in any solution, so Alice’s salary for October is compromised. In Definition 4, we formalize the definition of compromiseability. We distinguish between two cases, that is, the trivial case illustrated by the third quarter data of Table 1, and the complementary case illustrated by the fourth quarter data.



**Theorem 1.** 1. A full core cuboid  $C_f$  cannot be trivially compromised by any set of aggregation cuboids  $S$ .

2.  $C_c$  is trivially compromised by  $S_1$  if  $|C_c| < 2^{k-1} \cdot \max(D_1, D_2, \dots, D_k)$  for  $k \geq 2$

**Proof:** Given in [30] (omitted here due to space limitation).

Theorem 1 provides cardinality-based criteria for the two extreme cases, i.e., the core cuboid is either full or sparse. However, cardinality-based criteria are ineffective for cases in between. As an example, consider the third quarter data in Table 1, which is trivially compromised. Without changing the cardinality, evenly distributing the three “N/a” in three months makes the core cuboid free of trivial compromise. This invalidates any cardinality based criteria because trivial compromiseability varies for core cuboids with exactly the same cardinality.

### 3.3 Non-trivial Compromiseability

In this section, we derive cardinality-based criteria for determining compromiseability in the non-trivial case. We have two results. The first is that full core cuboids cannot be non-trivially compromised. The second is a lower bound on the cardinality of the core cuboid such that it remains safe from non-trivial compromise. First we have Lemma 1.

**Lemma 1.** 1.  $C_c$  can not be non-trivially compromised by any single cuboid.

2. If  $C_c$  cannot be compromised by  $S_1$ , then it cannot be compromised by  $S_{all}$ .

3. For any integers  $k$  and  $D_1, D_2, \dots, D_k$  that satisfy  $D_i \geq 4$  for  $1 \leq i \leq k$ , there is a  $k$ -dimensional data cube  $\langle C_c, S_{all} \rangle$ , with integer boundaries  $D_i$ , such that  $C_c$  is non-trivially compromised by  $S_1$ .

**Proof:** Given in [30] (omitted here due to space limitation).

Because of the second claim of Lemma 1, it is sufficient to safeguard the core cuboid from the collection of 1-\* aggregation cuboids. The last condition in Lemma 1 shows that it is impossible to obtain a criteria for preventing non-trivial compromiseability by only looking at the dimension cardinalities.

**Theorem 2 (Non-trivial Compromiseability).**

1.  $C_f$  cannot be non-trivially compromised by  $S_1$ .

2. For any integers  $k$  and  $D_i \geq 4$ , there exists a  $k$ -dimensional data cube  $\langle C_c, S_{all} \rangle$  satisfying  $|C_f - C_c| = 2D_l + 2D_m - 9$  such that  $C_c$  is non-trivially compromised by  $S_1$ , where  $D_l$  and  $D_m$  are the least two among  $D_i$ .

3. If  $|C_f - C_c| < 2D_l + 2D_m - 9$ , then  $C_c$  cannot be non-trivially compromised.

**Proof:** See the Appendix.

The first claim in Theorem 2 guarantees the non-trivial compromiseability of full core cuboid. Second and third claims give a tight lower bound on the cardinality of a core cuboid for it to remain free of non-trivial compromises. The second claim also implies that no cardinality based criteria can be derived for core cuboids with a cardinality below the computed lower bound.



**Corollary 1 (Non-trivial Compromiseability).**

If for any  $i \in [1, k]$ , there exists  $j \in [1, D_i]$  such that  $|P_i(C_f, j) - P_i(C_c, j)| = 0$ ,  $C_c$  cannot be non-trivially compromised.

**Proof:** Follows from the proof of Theorem 2. □

By Corollary 1, full slice on every dimension guarantees non-compromiseability in the non-trivial case.

## 4 A Cardinality-Based Inference Control Algorithm for Data Cubes

This section describes an algorithm to control inferences in data cube style OLAP queries using the results on compromiseability developed in Section 3. Our algorithm is based on a three-tier model consisting of core data, pre-computed aggregations and answerable queries.

### 4.1 Three-Tier Inference Control Model

Our three-tier model consists of three basic components and two abstract relations in between as given below and illustrated in Figure 1. In addition we enforce three properties on the model.

**1. Three Tiers:**

- (a) A set of data items  $D$ .
- (b) A set of aggregations  $A$ .
- (c) A set of queries  $Q$ .

**2. Relations Between Tiers:**

- (a)  $R_{AD} \subseteq A \times D$ .
- (b)  $R_{QA} \subseteq Q \times A$ .

**3. Properties:**

- (a)  $|A| \ll |Q|$ .
- (b) There are partitions  $P_D$  on  $D$  and  $P_A$  on  $A$ , such that for any  $(a, d) \in R_{AD}$  and  $(a', d') \in R_{AD}$ ,  $d$  and  $d'$  are in the same chunk of  $P_D$  if and only if  $a$  and  $a'$  are in the same chunk of  $P_A$ .
- (c)  $D$  is not compromised by  $A$ .

Three-tier inference control model simplifies inference control problem in several ways. Firstly, because all queries in  $Q$  are derived from aggregations in  $A$ , it suffices to ensure the non-compromiseability  $A$  instead of  $Q$ . This reduces the complexity of inference control due to the first characteristic of  $A$ . Secondly, the second characteristic of  $A$  allows us to adopt a divide-and-conquer approach to further reduce the complexity of inference control. Thirdly, inference control is embedded in the off-line design of  $A$  and  $R_{AD}$ , so the overhead of on-line inference control is eliminated or reduced. Although the restriction of  $Q$  to be derived from  $A$  reduces the total number of answerable queries,  $A$  can be designed in such a way that it contains most semantics required by the application. Hence the restricted queries are mostly arbitrary and meaningless with respect to application requirements.

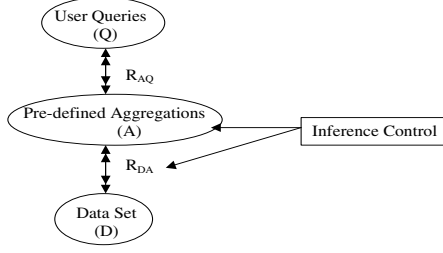


Fig. 1. Three-tier model for controlling inferences

**Algorithm** *Ctrl\_Inf\_Cube*

**Input:** Data Cube  $\langle C_c, S_{all} \rangle$ , integers  $D_i^j \in [1, D_i]$  for  $1 \leq i \leq k$  and  $0 \leq j \leq m_i$ , where  $m_i$  is fixed for each  $i$  and  $D_i^j$  satisfy  $D_i^0 = 1$ ,  $D_i^{m_i} = D_i$ , and  $D_i^{j-1} < D_i^j$  for all  $1 \leq j \leq m_i$ .

**Output:** a set of aggregations  $S_A$  that do not compromise  $C_c$

**Method:**

1. Let  $S_A = \phi$ ;
2. **For** each  $k$  dimensional vector  $v \in \prod_{i=1}^k [1, m_i]$   
     Let  $C_{tmp} = \{t : t \in C_c, \forall i \in [1, k] \ t[i] \in [D_i^{v[i]-1}, D_i^{v[i]}]\}$ ;  
     Let  $C_p = \{t : \exists t' \in C_{tmp}, \forall i \in [1, k] \ t[i] = t'[i] - D_i^{v[i]-1} + 1\}$ ;  
     Let  $S_A = S_A \cup Ctrl\_Inf\_Chunk(C_p, [1, D_i^{v[i]} - D_i^{v[i]-1} + 1])$ ;
3. **Return**  $S_A$ .

**Subroutine** *Ctrl\_Inf\_Chunk*

**Input:**  $k$  dimensional core cuboid  $C'_c$  and the  $k$  dimension domains  $[1, D'_i], i = 1, 2, \dots, k$

**Output:**  $S'_{all}$  if it does not compromise  $C'_c$  according to the cardinality-based criteria,  $\phi$  otherwise, where  $S'_{all}$  is the set of all aggregation cuboids defined on  $C'_c$

**Method:**

1. **If**  $|C'_c| = 0$  or  $|C'_c| = |C'_f|$ , where  $C'_f$  is the full core cuboid of  $C'_c$   
     **Return**  $S'_{all}$ ;
2. **If**  $C'_c$  is trivially compromised by  $S'_{all}$   
     **Return**  $\phi$ ;
3. Let  $D_l, D_m$  be the two smallest among  $D'_i$ ;
4. **If**  $|C'_f - C'_c| < 2D_l + 2D_m - 9$   
     **Return**  $S'_{all}$ ;
5. **If** for all  $i \in [1, k]$  there exists  $j \in [1, D'_i]$  such that  $|P_i(C'_f, j) - P_i(C'_c, j)| = 0$   
     **Return**  $S'_{all}$ ;
6. **Return**  $\phi$ .

Fig. 2. Inference control algorithm for data cubes

## 4.2 Inference Control Algorithm

The inference control algorithm shown in Figure 2 applies the results given in Section 3 using our three-tier model. The algorithm first partitions the core cuboid into disjoint chunks, each of which is then passed to the subroutine *Ctrl\_Inf\_Chunk*. The subroutine checks the non-compromisability of the *sub-*

*data cube* defined on this chunk of data, using the cardinality based criteria. If it is compromised the subroutine returns an empty set, indicating no aggregation is allowed on the data. Otherwise, the subroutine returns the set of all aggregation cuboids of the sub-data cube. The final outcome is the union of all partial results returned by the subroutine (this set of aggregations can then be used to answer data cube style OLAP queries without inference problem).

**Correctness** The correctness of the algorithm, that is, the non-compromiseability of the final result is straight-forward. The subroutine *Ctrl\_Inf\_Chunk* guarantees the non-compromiseability of each sub-data cube respectively. In addition, the sub-data cubes are disjoint, making the non-compromiseability of each of them independent of others.

**Runtime Analysis:** The main routine of the algorithm partitions  $C_c$  by evaluating the  $k$  dimensions of each tuple. Let  $n = |C_c|$ , so the runtime of the main routine is  $O(nk)=O(n)$  (suppose  $k$  is fixed with respect to  $n$ ). The subroutine *Ctrl\_Inf\_Chunk* is called for each of the  $N = \prod_{i=1}^k m_i$  chunks ( $m_i$  are defined in the algorithm). It evaluates the cardinality of each input chunk  $C'_c$ , which has the same complexity as establishing its 1-\* aggregation matrix  $M'_1$ .

Let  $n' = \prod_{i=1}^k D'_i$  be the number of columns in  $M'_1$  ( $D'_i$  are defined in the algorithm), then  $m' = n' \sum_{i=1}^k \frac{1}{D'_i}$  is the number of rows. Let  $D_i^{max}$  be the maximum value among  $D'_i$ . Out of the  $(m'n')$  elements,  $O(m' \cdot D_i^{max})$  elements must be considered to compute  $M'_1$ . Suppose  $(\sum_{i=1}^k \frac{1}{D'_i})D_i^{max} = O(k)$ . Then the runtime of the subroutine is  $O(k \cdot \prod_{i=1}^k D'_i)$ . It is called  $N$  times so the total runtime is  $O(k \cdot \prod_{i=1}^k m_i \cdot \prod_{i=1}^k D'_i) = O(k \cdot \prod_{i=1}^k m_i \cdot \prod_{i=1}^k \frac{D_i}{m_i})$ , which is  $O(k \cdot \prod_{i=1}^k D_i) = O(n)$ . We note that by definition, determining non-compromiseability has a complexity of  $O(n^3)$  and the maximum non-compromiseable subset of aggregations cannot be found in polynomial time [11].

**Enhancing the Algorithm:** The algorithm demonstrates a simple application of the cardinality based criteria in Section 3, which can be improved in many aspects. The dimension hierarchies inherent to most OLAP datasets can be exploited to increase the semantics included in the final output of the algorithm. For example, assume time dimension has a hierarchy comprised of day, week, month and year. Instead of partitioning the dataset arbitrarily, each week can be used for a chunk. Queries about weeks, months and years can then be answered with only the aggregations in algorithm output.

Notice that the key to cardinality-based non-compromiseability is that each chunk in the partition of a core cube must be either empty or dense (full). The row shuffling [4] technique proposed by Barbara et al. increases the subspace density of data sets by shuffling tuples along those categorical, unordered dimensions. Row shuffling can be integrated into the inference control algorithm as a pre-processing step prior to partitioning.

**Data Cube Operations:** We briefly discuss how our algorithm may address common data cube operations such as slicing, dicing, rolling up, drilling down and range query. Slicing, dicing and range query require aggregations to be defined on a subspace formed by intervals in dimension domains. Our algorithm also partitions the data set into small chunks. Therefore, in order to enhance our algorithm to address these operations, the subspace required by these data cube operations should be formed as the union of multiple chunks. Rolling up and drilling down require aggregations to be defined at different granularities than those in the original data cube. Rolling up does not directly create an inference threat because with coarser granulated queries include less information about individual data. Our ongoing work is addressing these details.

Although update operations are uncommon in decision support systems, data stored in data warehouses need to be updated over time. Our algorithm is suitable for update operations in two aspects. Firstly, changing values has no effect on cardinality, which determines the non-compromiseability in our algorithm. Secondly, because we have *localized* protection by partitioning data set into small disjoint chunks, the effect of an insertion or deletion is restricted to only the chunks containing updated tuples.

## 5 Conclusions

Based on a definition of non-compromiseability to mean that each unknown sensitive variable has more than one choices of value to fit a given set of their aggregations, we have derived sufficient conditions for non-compromiseability in sum-only data cubes. We have proved that full core cuboids can not be compromised, and that there is a tight lower bound on the cardinality of a non-compromiseable core cuboid. To apply our results to the inference control of data cube style OLAP queries, we have shown a *divide and conquer* algorithm based on a three-tier model. Future work includes enhancing our results and algorithm to include data cube operations and consider other variations of OLAP queries.

## References

- [1] N.R. Adam and J.C. Wortmann. Security-control methods for statistical databases: a comparative study. *ACM Computing Surveys*, 21(4):515–556, 1989. [57](#)
- [2] D. Agrawal and C. C. Aggarwal. On the design and quantification of privacy preserving data mining algorithms. In *Proceedings of the Twentieth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems*, pages 247–255, 2001. [58](#)
- [3] R. Agrawal and R. Srikant. Privacy-preserving data mining. In *Proceedings of the 2000 IEEE Symposium on Security and Privacy*, pages 439–450, 2000. [58](#)
- [4] D. Barbará and X. Wu. Using approximations to scale exploratory data analysis in datacubes. In *Proceedings of the Fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 382–386, 1999. [65](#)

- [5] L. L. Beck. A security mechanism for statistical databases. *ACM Trans. on Database Systems*, 5(3):316–338, 1980. 57
- [6] A. Brodsky, C. Farkas, and S. Jajodia. Secure databases: Constraints, inference channels, and monitoring disclosures. *IEEE Trans. Knowledge and Data Engineering*, 12(6):900–919, 2000. 57
- [7] A. Brodsky, C. Farkas, D. Wijesekera, and X.S. Wang. Constraints, inference channels and secure databases. In *the 6th International Conference on Principles and Practice of Constraint Programming*, pages 98–113, 2000. 57
- [8] F. Y. Chin, P. Kossowski, and S. C. Loh. Efficient inference control for range sum queries. *Theoretical Computer Science*, 32:77–86, 1984. 57
- [9] F. Y. Chin and G. Özsoyoglu. Security in partitioned dynamic statistical databases. In *Proc. of IEEE COMPSAC*, pages 594–601, 1979. 57
- [10] F. Y. Chin and G. Özsoyoglu. Statistical database design. *ACM Trans. on Database Systems*, 6(1):113–139, 1981. 57
- [11] F. Y. Chin and G. Özsoyoglu. Auditing and inference control in statistical databases. *IEEE Trans. on Software Engineering*, 8(6):574–582, 1982. 57, 58, 65
- [12] L. H. Cox. Suppression methodology and statistical disclosure control. *Journal of American Statistic Association*, 75(370):377–385, 1980. 57
- [13] D. E. Denning. Secure statistical databases with random sample queries. *ACM Trans. on Database Systems*, 5(3):291–315, 1980. 57
- [14] D. E. Denning and P. J. Denning. Data security. *ACM computing surveys*, 11(3):227–249, 1979. 57
- [15] D. E. Denning and J. Schlörer. Inference controls for statistical databases. *IEEE Computer*, 16(7):69–82, 1983. 57
- [16] D. Dobkin, A. K. Jones, and R. J. Lipton. Secure databases: protection against user influence. *ACM Trans. on Database Systems*, 4(1):97–106, 1979. 57, 58
- [17] A. Evfimievski, R. Srikant, , R. Agrawal, and J. Gehrke. Privacy preserving mining of association rules. In *Proceedings of the 8th Conference on Knowledge Discovery and Data Mining (KDD'02)*, 2002. 58
- [18] L. P. Fellegi. On the gestion of statistical confidentiality. *Journal of American Statistic Association*, 67(337):7–18, 1972. 57
- [19] J. Gray, A. Bosworth, A. Layman, and H. Pirahesh. Data cube: A relational operator generalizing group-by, crosstab and sub-totals. In *Proceedings of the 12th International Conference on Data Engineering*, pages 152–159, 1996. 56, 59
- [20] V. Harinarayan, A. Rajaraman, and J. D. Ullman. Implementing data cubes efficiently. In *Proceedings of the 1996 ACM SIGMOD international conference on Management of data*, pages 205–227, 1996. 56
- [21] K. Hoffman. *Linear Algebra*. Prentice-Hall, 1961. 60, 61
- [22] J. Kleinberg, C. Papadimitriou, and P. Raghavan. Auditing boolean attributes. In *Proc. of the 9th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*, pages 86–91, 2000. 57, 60
- [23] Y. Li, L. Wang, X. S. Wang, and S. Jajodia. Auditing interval-based inference. In *Proceedings of the 14th Conference on Advanced Information Systems Engineering (CAiSE'02)*, 2001. 57, 60
- [24] F. M. Malvestuto and M. Moscarini. Computational issues connected with the protection of sensitive statistics by auditing sum-queries. In *Proc. of IEEE Scientific and Statistical Database Management*, pages 134–144, 1998. 57, 60

- [25] J.M. Mateo-Sanz and J. Domingo-Ferrer. A method for data-oriented multivariate microaggregation. In *Proceedings of the Conference on Statistical Data Protection'98*, pages 89–99, 1998. 57
- [26] S. Rizvi and J.R. Haritsa. Maintaining data privacy in association rule mining. In *Proceedings of the 28th Conference on Very Large Data Base (VLDB'02)*, 2002. 58
- [27] J. Schlörer. Security of statistical databases: multidimensional transformation. *ACM Trans. on Database Systems*, 6(1):95–112, 1981. 57
- [28] R.P. Tewarson. *Sparse Matrices*. Academic Press, 1973. 69
- [29] J.F. Traub, Y. Yemini, and H. Woźniakowski. The statistical security of a statistical database. *ACM Trans. on Database Systems*, 9(4):672–679, 1984. 57
- [30] L. Wang, D. Wijesekera, and J. Sushil. Cardinality-based inference control in sum-only data cubes (extended version). Technical Report, 2002. Available at <http://ise.gmu.edu/techrep/>. 62

## Appendix

### Proof(Theorem 2):

- Without loss of generality, we show that  $t_0 = (1, 1, \dots, 1)$  cannot be nontrivially compromised by  $S_1$ . Let  $C'_f = \{t : \forall i \in [1, k], t[i] = 1 \vee t[i] = 2\}$ . Since  $(D_1, D_2, \dots, D_k) > 1$ , we have that  $C'_f \subseteq C_f$  and  $|C'_f| = 2^k$ . It suffices to prove that  $t_0$  cannot be compromised by  $S_1$  in the data cube  $\langle C'_f, S_{all} \rangle$ . Let  $M'_1$  be the 1-\* aggregation matrix of  $S_1$  on  $C'_f$ . According to Definition 3, there are  $2^k$  non-zero column vectors in  $M'_1$ , corresponding to the  $2^k$  tuples in  $C'_f$ . In the rest of the proof we formally show that each of the  $2^k$  non-zero column vectors can be represented by the linear combination of the left  $2^k - 1$  column vectors. Then, it follows from linear algebra that  $t_0$  cannot be compromised by  $S_1$  in data cube  $\langle C'_f, S_{all} \rangle$  (and hence in  $\langle C_f, S_{all} \rangle$ ). In order to prove our informally stated claim, we define the *sign assignment vector* as an  $n$  dimensional column vector  $t_{sign}$  where  $n$  is  $|C_f|$ , as follows:

- $t_{sign}[1] = 1$
- $t_{sign}[2^i + j] = -t_{sign}[j]$  for all  $0 \leq i \leq k - 1$  and  $1 \leq j \leq 2^i$
- $t_{sign}[j] = 0$  for all  $j > 2^k$

**Claim:**  $M'_1 \cdot t_{sign} = 0$ , where 0 is the  $n$  dimensional zero column vector.

**Justification:**

Let  $t = S_1[i], t[l] = *$  for  $l \in [1, k]$ .

Let  $v$  be  $M'_1[i, -]$ .

Suppose  $t[j] = 1$  or  $t[j] = 2$  for all  $j \neq l$ .

Then  $|Qset(t)| = 2$ , and as a consequence we get  $Qset(t) = \{t_1, t_2\}$

where  $t_1, t_2 \in C_f$ ,  $t_1[l] = 1, t_1[l] = 2$

and  $t_1[j] = t_2[j] = t[j]$  for all  $j \neq l$

Hence, there are two integers  $j_1, j_2 \in [1, n]$  satisfying

$v[j_1] = v[j_2] = 1$  and  $v[j] = 0$  for any  $j \neq j_1, j_2$ .

By Definition 3,  $M'_1[-, j_1]$  (the  $j_1^{th}$  column of  $M'_1$ ) and  $M'_1[-, j_2]$

correspond to  $t_1$  and  $t_2$  respectively.

Because  $C'_f$  is formed in dictionary order, we get  $j_2 = j_1 + 2^{l-1}$ .

Hence, we have  $v \cdot t_{sign} = 0$ .

Otherwise,  $|Qset(t)| = 0$ ; and hence  $Qset(t) = \phi$ .

Hence,  $v = 0$ , and hence,  $0 \cdot t_{sign} = 0$ .

This justifies our claim. Hence, as stated earlier, the justification concludes the main proof.

2. Without loss of generality we assume  $D_1, D_2$  are the least two among  $D_i$ 's. For an arbitrary but fixed value of  $D_2$ , we show by induction on  $D_1$  that  $C_c$  as constructed in the proof of Lemma 1 satisfies  $|C_f - C_c| = 2D_1 + 2D_2 - 9$ . **Inductive Hypothesis:**  $C_c$  as constructed in the proof of Lemma 1 satisfies:
  - $|C_f - C_c| = 2j + 2D_2 - 9$  for any  $j \geq 4$ .
  - $|P_1(C_f, j) - P_1(C_c, j)| = 2$  for any  $j \in [1, D_1]$ .

**Base Case:** In the base case of the proof of Lemma 1, the core cuboid  $C_c$  satisfies  $|C_f - C_c| = 2D_1 + 2D_2 - 9$ . Notice that the core cuboid,  $D_1 = 4$ , and  $|P_1(C_f, j) - P_1(C_c, j)| = 2$ . This validates the base case of our inductive hypothesis.

**Inductive Case:** Suppose for  $D_1 = j$  we have  $|C_f - C_c| = 2j + 2D_2 - 9$  and  $|P_1(C_f, j) - P_1(C_c, j)| = 2$ . Let  $C'_f$  be the full core cuboid corresponding to  $C'_c$  for  $D_1 = j + 1$ . By the definition of  $C$  in the proof of Lemma 1, we have  $|C| = |P_1(C_c, j)|$  and as a consequence  $|C'_f - C'_c| = |C_f - C_c| + 2 = 2(j + 1) + 2D_2 - 9$ . Since  $P_1(C'_c, j + 1) = C$ . Hence,  $|P_1(C'_f, j) - P_1(C'_c, j)| = 2$ . This validates the inductive case of our inductive argument and consequently concludes our proof of the tightness of the cardinality lower bound for avoiding nontrivial compromiseability.

3. **Lower Bound:** We show that if  $C_c$  is nontrivially compromised then we have  $|C_f - C_c| \geq 2D_1 + 2D_2 - 9$ . First we make following assumptions.
  - (a) The tuple  $t = (1, 1, \dots, 1) \in C_c$  is nontrivially compromised by  $S_1$
  - (b) No tuple in  $C_c$  is trivially compromised
  - (c) There exists  $S \subseteq S_1$  such that  $S$  non-trivially compromises  $t$ , but for any  $C \in S$ ,  $S \setminus C$  does not non-trivially compromise  $t$
  - (d) For any  $t' \in C_f \setminus C_c$ ,  $t$  cannot be nontrivially compromised by  $S_1$  in data cube  $\langle C_c \cup \{t'\}, S_{all} \rangle$ .

Assumption (b) holds by Definition 4. Assumption (c) is reasonable considering the case  $S$  contains only a single aggregation cuboid. Assumption (d) is reasonable considering the case  $|C_f \setminus C_c| = 1$ .

**Claim:** Suppose Assumption (a),(b),(c), and (d) hold. Furthermore assume that there is a  $C \in S$  where  $t \in C$  satisfies  $t[i] = *$ . Then  $|P_i(C_f, 1) - P_i(C_c, 1)| \geq 1$ , and  $|P_i(C_f, j) - P_i(C_c, j)| \geq 2$  holds for any  $j \in [2, D_i]$ .

**Justification:** The proof is by contradiction. Without loss of generality, we only justify the claim for  $i = k$  and  $j = 2$ . That is, given a  $C \in S$  satisfying  $t[k] = *$  for any  $t \in C$ , we prove that  $|P_k(C_f, 2) - P_k(C_c, 2)| \geq 2$ .

First we transform the aggregation matrix of  $S$  on  $C_c$  by row permutation into a singly bordered block diagonal form (SBBDF) [28], denoted by  $M_{m \times n}$ . The  $i^{th}$  diagonal block of  $M$  corresponds to  $P_k(C_c, i)$  and  $\{t : t \in S \setminus C, t[k] =$

$i\}$ , and the border of  $M$  denotes the aggregation cuboid  $C$ . We call the columns of  $M$  corresponding to the  $i^{\text{th}}$  diagonal block as the  $i^{\text{th}}$  slice of  $M$ . Due to Assumption (a), there exists a row vector  $a$  satisfying  $a \cdot M = e_1$ . Let  $r_i$  be  $M[i, -]$  then we get  $e_1 = \sum_{i=1}^m a[i] \cdot r_i$ . Suppose each diagonal block of  $M$  has size  $m' \times n'$ . Use  $r_i^j$ , for  $1 \leq j \leq D_k$  to represent the row vector composed of the elements of  $r_i$  that falls into the  $j^{\text{th}}$  slice of  $M$ . Notice that there are  $n'$  elements in  $r_i^j$ . We also use  $e'_1$  and  $0'$  to represent the  $n'$  dimensional unit row vector and  $n'$  dimensional zero row vector, respectively.

Then the following are true:

**i.**  $e'_1 = \sum_{i=1}^{m'} a[i]r_i^1 + \sum_{i=m-m'+1}^m a[i]r_i^1$

**ii.**  $0' = \sum_{i=m'+1}^{2m'} a[i]r_i^2 + \sum_{i=m-m'+1}^m a[i]r_i^2$

First we suppose  $|P_k(C_f, 2) - P_k(C_c, 2)| = 0$ , that is, the second slice of  $M$  contains no zero column. We then derive contradictions to our assumptions. Since  $|P_k(C_f, 2) - P_k(C_c, 2)| = 0$  the first slice of  $M$  contains no more non-zero columns than the second slice of  $M$  does. Intuitively if the latter is transformed into a zero vector then applying the same transformation on the former leads to a zero vector, too. This is formally represented as:

**iii.**  $0' = \sum_{i=1}^{m'} a[m' + i]r_i^1 + \sum_{i=m-m'+1}^m a[i]r_i^1$ .

Subtracting (iii) from (i) gives  $e'_1 = \sum_{i=1}^{m'} (a[i] - a[m' + i])r_i^1$ . That implies  $C_c$  is nontrivially compromised by  $S \setminus \{C_k\}$ , contradicting Assumption (c). Thus  $|P_k(C_f, 2) - P_k(C_c, 2)| \neq 0$ .

Next we assume  $|P_k(C_f, 2) - P_k(C_c, 2)| = 1$  and derive a contradiction to our assumptions.

First the row vector  $r_i^3$  satisfies the following condition:

**iv.**  $0' = \sum_{i=2m'+1}^{3m'} a[i]r_i^3 + \sum_{i=m-m'+1}^m a[i]r_i^3$ .

Let  $t' \in P_k(C_f, 2) \setminus P_k(C_c, 2)$ . Notice that (i), (ii) still hold. Suppose  $t'$  corresponds to  $M[-, y] = 0$ . Now assume we add  $t'$  to  $P_k(C_c, 2)$ , consequently we have  $M[-, y] \neq 0$ . Due to Assumption (d), we have that the left side of (ii) becomes  $e'_1$ , that is,  $a \cdot M[-, y] = 1$ . There is also an extra 1-element  $M[x, y]$  in the border of  $M$ .

Now let  $t''$  be the tuple corresponding to  $M[-, y + n']$  in the third slice of  $M$ . Suppose  $t'' \in P_k(C_c, 3)$  and consequently  $M[-, y + n'] \neq 0$ . We have that  $M[-, y + n'] = M[-, y]$  and consequently  $a \cdot M[-, y + n'] = 1$ .

By removing  $t'$  from  $P_k(C_c, 2)$  we return to the original state that all our assumption hold. Now we show by contradiction that  $t'' \in P_k(C_c, 3)$  cannot hold any longer. Intuitively, since  $t'$  is the only missing tuple in the second slice of  $M$ , the third slice of  $M$  contains no more non-zero vectors than the second slice of  $M$  does, except  $t''$ . Because  $a \cdot M[-, y + n'] = 1$ , elements of  $a$  transform the second slice of  $M$  to a zero vector, as shown by (ii), also transform the third slice of  $M$  to a unit vector. This is formally represented in (v):

**v.**  $e'' = \sum_{i=2m'+1}^{3m'} a[i - m']r_i^3 + \sum_{i=m-m'+1}^m a[i]r_i^3$

Subtracting (iv) from (v) we get that  $e'' = \sum_{i=2m'+1}^{3m'} (a[i - m'] - a[i])r_i^3$ ; implying  $C_c$  is compromised by  $S \setminus \{C_i\}$ . Hence, Assumption (c) is false. Consequently,  $t'' \notin C_c$ .



Similar proof exists for the  $i^{th}$  slice of  $C_c$ , where  $i = 4, 5, \dots, D_k$ . However,  $M[x, -] \neq 0$  because if so, we can let  $a_x$  be zero and then decrease the number of missing tuples in  $C_c$ , contradicting Assumption (d). Hence  $M[x, -]$  is a unit vector with the 1-element in the first slice of  $M$ . However, this further contradicts Assumption (b), that no trivial compromise is possible. Hence we have that  $|P_k(C_f, 2) - P_k(C_c, 2)| = 1$  is false.

Now consider  $|P_k(C_f, 1) - P_k(C_c, 1)|$ . Suppose all the assumptions hold, and  $|P_k(C_f, 1) - P_k(C_c, 1)| = 0$ . Let  $t_1, t_2 \in P_k(C_f, 2) \setminus P_k(C_c, 2)$ . Now define  $C'_c = C_c \setminus \{t\} \cup \{t_1\}$  and  $M'$  be the aggregation matrix of  $S$  on  $C'_c$ . From  $a \cdot M = e_1$ , and Assumption (d) we get  $a \cdot M' = e_i$ , where  $M[-, i]$  corresponds to  $t_1$ . This implies the nontrivially compromise of  $t_1$  in  $\langle C'_c, S_{all} \rangle$ , with  $|P_k(C_f, 1) - P_k(C'_c, 1)| = 1$ , which contradicts what we have already proved. Hence, we get  $|P_k(C_f, 1) - P_k(C_c, 1)| \geq 1$ . This concludes the justification of our claim.

The claim implies that the number of missing tuples in  $C_c$  increases monotonically with the following:

- The number of aggregation cuboids in  $S$ .
- $D_i$ , provided there is  $C \in S$  satisfying  $t[i] = *$  for any  $t \in C$ .

Hence  $|C_f - C_c|$  reaches its lower bound when  $S = \{C_1, C_2\}$ , which is equal to  $2D_1 + 2D_2 - 9$ , as shown in the first part of the current proof - concluding the proof of Theorem 2.

□