



ELSEVIER

Available at
www.ComputerScienceWeb.com
POWERED BY SCIENCE @ DIRECT®

Data & Knowledge Engineering 45 (2003) 131–153

DATA &
KNOWLEDGE
ENGINEERING

www.elsevier.com/locate/datak

Bounding the cardinality of aggregate views through domain-derived constraints [☆]

Paolo Ciaccia, Matteo Golfarelli, Stefano Rizzi *

DEIS, University of Bologna, Viale Risorgimento 2, 40136 Bologna, Italy

Abstract

Accurately estimating the cardinality of aggregate views is crucial for logical and physical design of data warehouses. This paper proposes an approach based on cardinality constraints, derived a-priori from the application domain, which may bound either the cardinality of a view or the ratio between the cardinalities of two views. We face the problem by first computing satisfactory bounds for the cardinality, then by capitalizing on these bounds to determine a good probabilistic estimate for it. In particular, we propose a bounding strategy which achieves an effective trade-off between the tightness of the bounds produced and the computational complexity.

© 2002 Elsevier Science B.V. All rights reserved.

Keywords: Data warehousing; Logical design; View materialization

1. Introduction and motivation

The multidimensional model is the foundation for data representation and querying in multi-dimensional databases and data warehouses. It represents facts of interest for the decision process into *cubes* in which each cell contains numerical *measures* which quantify the fact from different points of view, while each axis represents an interesting *dimension* for analysis. For instance, within a four-dimensional cube modeling the phone calls supported by a telecommunication company, the dimensions might be the calling number, the number called, the date, and the time segment in which the call is placed; each cube cell could be associated to a measure of the total duration of the calls made from a given number to another number on a given time segment and date.

[☆]This work has been partially supported by the D2I MURST project.

* Corresponding author. Tel.: +39-051-2093542; fax: +39-051-2093540.

E-mail addresses: pciaccia@deis.unibo.it (P. Ciaccia), mgolfarelli@deis.unibo.it (M. Golfarelli), srizzi@deis.unibo.it (S. Rizzi).

The basic mechanism to extract significant information from the huge quantity of data stored in cubes is aggregation according to hierarchies of attributes rooted in dimensions [11]. In most application cases, cubes are significantly sparse (for instance, most couples of telephone numbers are never connected by a call in a given date), and so are the aggregate views.

Accurately estimating the actual cardinality of each view is crucial for logical and physical design as well as for query processing and optimization [19]. As a relevant case, consider the view materialization problem, where the aggregate views which are the most useful in answering the workload queries have to be selected for materialization (see [18] for a survey). Since the number of possible views which can be derived by aggregating a cube is exponential in the number of attributes, most approaches assume that a constraint on the total disk space occupied by materialization is posed, and attempt to find the subset of views which contemporarily satisfies this constraint and minimizes the workload cost [7,9,12]. Another case where estimation of view cardinalities is relevant is index selection [10].

If the data warehouse has already been loaded, view cardinalities can be quite accurately estimated by using statistical techniques based, say, on histograms [15] or sampling [13]. However, such techniques cannot be applied at all if the data warehouse is still under development, and the estimation of view cardinalities is needed for design purposes. To obviate this, current approaches are based on estimation models that only exploit the cardinality of the base cube and that of the single attribute domains [16,17], which however leads to significant overestimation.

In this paper we propose a novel approach to estimate the cardinality of views based on a-priori information derived from the application domain. Similarly to what is done when estimating the cardinality of projections in relational databases [6], we face the problem by first computing satisfactory bounds for the cardinality, then by capitalizing on these bounds to determine a good probabilistic estimate for it. Besides the functional dependencies (FD's) expressed by the multi-dimensional scheme, the bounds we determine also take into account additional domain-derived information expressed in the form of *cardinality constraints*, namely, bounds of the cardinality of some views and bounds (called *k-dependencies*) on the ratio between the cardinalities of two views. The main contribution of the paper is a *bounding strategy* comprising (1) a bounding function which computes effective upper bounds of cardinalities and (2) a set of formal results aimed at reducing the complexity of computation.

The paper is organized as follows. In Section 2, we provide some basic definitions, describe the basic principles of our approach and provide a motivating example. Section 3 introduces the basic properties of cardinality constraints, and Section 4 introduces the cover-based bounding strategy. Section 5 presents all major formal results on the computation of bounds. Section 6 proposes a simple probabilistic model to show how the bounds derived may be used to improve the cardinality estimates. Section 7 discusses some interesting open issues. Finally, the most complex proofs are included in the Appendix.

2. Outline of the approach

Before introducing the framework for our approach, we need to provide some basic definitions on views and on their associated lattice.

Definition 1 (*Dimensional scheme*). We call dimensional scheme \mathcal{D} a couple $(\mathcal{U}, \mathcal{F})$ where \mathcal{U} is a set of attributes and $\mathcal{F} = \{A_i \rightarrow A_j : A_i, A_j \in \mathcal{U}\}$ is a set of FD's which relate the attributes of \mathcal{U} into a set of pairwise disjoint directed trees. We call dimensions the attributes $A_k \in \mathcal{U}$ in which the trees are rooted, i.e., such that $\forall A_i \in \mathcal{U} (A_i \rightarrow A_k) \notin \mathcal{F}$; let $\dim(\mathcal{D}) \subseteq \mathcal{U}$ denote the set of dimensions of \mathcal{D} .

Definition 2 (*View*). Let $\mathcal{D} = (\mathcal{U}, \mathcal{F})$ be a dimensional scheme. We call view on \mathcal{D} any subset of attributes $V \subseteq \mathcal{U}$ such that $\forall A_i, A_j \in V (A_i \rightarrow A_j) \notin \mathcal{F}^+$ where \mathcal{F}^+ denotes the set of all FD's logically implied by \mathcal{F} .¹

Example 1. Consider an enterprise with branches in different cities. A simple dimensional scheme Transfers modeling the transfers of employees between offices might include:

$$\begin{aligned} \mathcal{U} &= \{\text{date, month, year, fromOffice, fromDept, fromCity, toOffice, toDept, toCity, employee}\} \\ \mathcal{F} &= \{\text{date} \rightarrow \text{month, month} \rightarrow \text{year, fromOffice} \rightarrow \text{fromDept,} \\ &\quad \text{fromOffice} \rightarrow \text{fromCity, toOffice} \rightarrow \text{toDept, toOffice} \rightarrow \text{toCity}\} \end{aligned}$$

thus the base cube is characterized by

$$\dim(\mathcal{D}) = \{\text{date, fromOffice, toOffice, employee}\}$$

Examples of views on the Transfers scheme are:

$$\begin{aligned} V &= \{\text{month, fromOffice, toCity, employee}\} \\ W &= \{\text{month, fromCity, fromDept, employee}\} \\ Z &= \{\text{year, fromOffice, toCity}\} \end{aligned}$$

In this work we face the problem of accurately estimating the cardinality of a view when the source data cannot be directly queried, which is the case during off-line logical design of multi-dimensional databases. We assume that a set \mathcal{S} of *cardinality constraints* is available instead, and we look for effective ways to exploit them for estimation. Without loss of generality, suppose that estimates are needed for the purpose of a view materialization algorithm. As sketched in Fig. 1, whenever the materialization algorithm requires information about a candidate view V , our approach works in two steps. First, the *bounder* uses the set \mathcal{S} of cardinality constraints supplied by the user to determine an effective upper bound for the cardinality of V ; then, the *estimator* uses this bound to derive a probabilistic estimate for the cardinality of V . Note that this two-steps approach generalizes well-known parametric models for the estimation of the cardinality of relational queries [14] and in particular those for projection size estimation [6], for which bounds are typically given as input parameters.

¹ We are using the term *view* to denote the set of grouping attributes used for aggregation, while the actual views will also include one or more measures. This slight abuse in terminology is allowable in the context of this work since the cardinality of a view only depends on its grouping attributes.

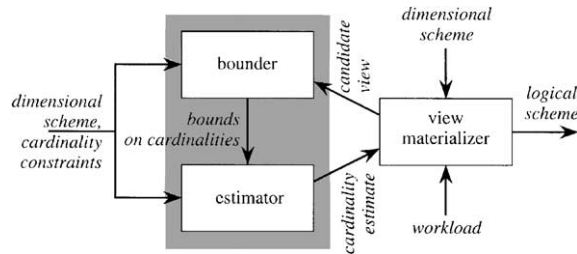


Fig. 1. Overall architecture for logical design.

We consider two different forms of cardinality constraints:

- (1) the upper bound w^+ of the cardinality w of a view W ;
- (2) k -dependencies, expressing an upper bound of the ratio between the cardinalities of two views W and Z (see Section 3.2).

We will assume that at least the upper bounds of the cardinalities of all the single attributes in the dimensional scheme are known. This assumption, which is perfectly reasonable in all application domains, is necessary in order to guarantee that at least one upper bound can be determined for each view.

The set \mathcal{S} , together with the dimensional scheme \mathcal{D} , univocally determines the *least upper bound* v^+ of the cardinality of V , meaning that:

- (1) in each instance of \mathcal{D} that does not violate any constraint in \mathcal{S} , the cardinality v of V is such that $v \leq v^+$; and
- (2) there exists an instance compatible with \mathcal{S} where $v = v^+$.

We say a constraint $c \in \mathcal{S}$ is *redundant* iff all the bounds determined by \mathcal{S} are equal to those determined by $\mathcal{S} - \{c\}$.

Definition 3 (*Sound and minimal input*). Let \mathcal{S} be a set of cardinality constraints on dimensional scheme \mathcal{D} . We say \mathcal{S} is *sound* iff there exists at least one non-empty instance of \mathcal{D} which satisfies all the constraints in \mathcal{S} . We say \mathcal{S} is *minimal* iff no constraint in \mathcal{S} is redundant.

In this paper we will assume that the input \mathcal{S} is sound and minimal. It is straightforward to derive that, in this case, all the bounds in \mathcal{S} are least upper bounds (whereas the opposite is not necessarily true).

Computing the bounds implied by \mathcal{S} turns out to be a challenging combinatorial problem, even for simple forms of cardinality constraints. For instance, it is known that the problem is NP-hard for arbitrary patterns of FD's [5]. Furthermore, the actual computational effort needed to compute these bounds might limit applicability in real-world cases. For this reason, the bounder is built around the concept of *bounding strategy*. A bounding strategy s is characterized by a bounding function that, given \mathcal{S} , \mathcal{D} , and V , computes an upper bound v_s^+ such that $v^+ \leq v_s^+$ holds.

In other terms, a bounding strategy never computes bounds which are more restrictive than the ones logically implied by the input constraints, trading-off accuracy for speed of evaluation.

Turning to the estimator, our framework supports different probabilistic models. A probabilistic model is a function that, given \mathcal{I} , \mathcal{D} , V as well as bounds computed by the bounder, provides an estimate, \bar{v} , for the cardinality of V . In general, this step can use further information from the application domain that is not suitable to derive bounds. Typically this is the case with information concerning average values (e.g., the number of transfers of each employee on each year is 1.5 on the average).

Example 2. Let 10^4 be the number of employees who have been transferred at least once, and let the enterprise consist of 10^3 offices distributed over 10 cities and belonging to one of 10 departments; let 10^3 days be the observation period. Let $V = \{\text{date, fromOffice, toOffice}\}$. Since each office is involved in transfers at most with every other office on each date, the first trivial upper bound of v is $10^3 \times 10^3 \times 10^3 = 10^9$. If a constraint states that the maximum number of transfers for an employee during one year is 2, and since we consider three years, it is derived that the cardinality of the base cube is at most six times the number of transferred employees, i.e., 6×10^4 . Thus, the upper bound of v can be improved to 6×10^4 as well (the cardinality of a view cannot exceed that of its base cube). Now, by using the probabilistic model in Section 6, the cardinality of V is estimated as $\bar{v} = 3.8 \times 10^4$.

The following compact notation is used for some examples throughout the rest of the paper. Uppercase letters from the beginning of the alphabet (A, B, \dots) denote dimensions. Attributes which are functionally determined by another attribute, i.e. attributes other than dimensions, are denoted by the corresponding primed letters (e.g., $A \rightarrow A'$, $A' \rightarrow A''$). The sets of attributes which define views are represented by omitting braces, thus writing ABC for $\{A, B, C\}$. Lowercase letters are used for the cardinalities of views and attributes (e.g., a is the cardinality of attribute A , ab is the cardinality of the view with attributes AB , and so on).

3. Basics on cardinality constraints

The possibility of exploiting cardinality constraints to bound the size of a view relies on the partial order induced on views by the FD's in the multidimensional scheme:

Definition 4 (Roll-up). Given the set $\mathcal{V}_{\mathcal{D}}$ of all possible views on \mathcal{D} , we define on $\mathcal{V}_{\mathcal{D}}$ the roll-up partial order \preceq as follows: $V \preceq W$ iff $W \rightarrow V$, i.e., iff $\forall A_i \in V \exists A_j \in W : (A_j \rightarrow A_i) \in \mathcal{F}^+$.

It is straightforward to verify that, given two views $V \in \mathcal{V}_{\mathcal{D}}$ and $W \in \mathcal{V}_{\mathcal{D}}$, both their *sup* and *inf* views always exist; we will denote them with $V \oplus W$ and $V \otimes W$, respectively.² Thus, the roll-up partial order determines a lattice, which we will call *multidimensional lattice* for \mathcal{D} , whose top and bottom elements are $\text{dim}(\mathcal{D})$ and the empty view \emptyset , respectively. The multidimensional lattice is isomorphic to the lattice of the order ideals of a partially ordered set [1]; thus, it is distributive.

² From a relational point of view, $V \oplus W$ is obtained by dropping from the natural join $V \bowtie W$ the functionally dependent attributes.

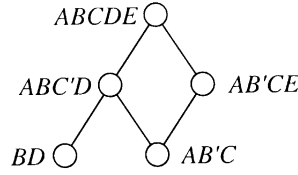


Fig. 2. Roll-up relationships of views in Example 3.

Given two views V and W , we will denote with $V \ominus W$ the least view Z such that $Z \oplus (V \otimes W) = V$. A nice property of this operator, which we will use in the proof of Lemma 7, is the following:

$$(V \ominus W) \oplus W = V \oplus W \quad (1)$$

Proof. By definition of \ominus , if we let $Z = V \ominus W$, it is $Z \oplus (V \otimes W) = V$ which implies $Z \oplus (V \otimes W) \oplus W = V \oplus W$. On the other hand, for the absorption property of distributive lattices it is $Z \oplus (V \otimes W) \oplus W = Z \oplus W$, from which $Z \oplus W = V \oplus W$. \square

Example 3. It is $ABC'D \oplus AB'CE = ABCDE$, $ABC'D \otimes AB'CE = AB'C$, $ABC'D \ominus AB'CE = BD$ (see Fig. 2).

3.1. Upper bounds

The basic observation to determine bounds for view cardinalities using bounds of the cardinalities of other views is that the multidimensional lattice induces an isomorphic structure over such cardinalities; this fact is expressed by the following two lemmas, where we use the notation $\llbracket E \rrbracket$ to denote the cardinality of a view that is the result of an expression E .

Lemma 1. If $W \preceq Z$, then $w^+ \leq z^+$.

Proof. From Definition 4 it follows that $W \preceq Z$ implies $w \leq z$ in each instance of \mathcal{D} , since $Z \rightarrow W$ holds. Now, assume that $w^+ > z^+$. Then, there is an instance of \mathcal{D} in which $z \leq z^+ < w \geq w^+$, thus $z < w$, which is a contradiction. \square

Lemma 2. Let S be a set of views and let $S' \subseteq S$ consist of all the views $W \in S$ such that for no $Z \in S$ it is $W \preceq Z$. Then, it is $\llbracket \oplus(S) \rrbracket^+ = \llbracket \oplus(S') \rrbracket^+ \leq \prod_{W \in S'} w^+$.

Proof. When $S' = S$ the result immediately follows, since the least upper bound of a set of views corresponds to their natural join, whose size can never exceed that of the Cartesian product of the views. When $S' \subset S$, the result follows from the observation that $\oplus(S) = \oplus(S')$, since each view in $S - S'$ is functionally determined by some view in S' . \square

Example 4. Let $S = \{AB, B'C, A'\}$; it is $\llbracket \oplus(S) \rrbracket^+ = \llbracket \oplus(\{AB, B'C\}) \rrbracket^+ = abc^+ \leq ab^+ \cdot b'c^+$.

3.2. The k -dependencies

A k -dependency (kD) is a relevant case of cardinality constraint which naturally generalizes a FD. In the authors' experience, kD 's are particularly useful to characterize the knowledge of the business domain held by the experts in the field. For instance, in the transfer domain, we might have some information concerning the number of destination cities for an employee, or on the number of distinct departments moved to from each department. If such information is in the form of bounds, it can be effectively used to improve the bounds of view cardinality.

Definition 5 (k -dependency). Let X and Y be two views on \mathcal{D} . We say that a kD holds between X and Y , and denote it with $X \xrightarrow{k} Y$, when k ($k \geq 1$) is an upper bound of the number of distinct tuples of Y which correspond to each distinct tuple of X within view $X \oplus Y$.

Example 5. In the Transfers scheme, assume the domain expert provides the following information: The maximum number of inter-department transfers of an employee during one year is 2. This constraint can be formalized by the following kD : $X \xrightarrow{2} Y$, where $X = \{\text{year}, \text{employee}\}$, $Y = \{\text{toDept}\}$. Intuitively, from this we can derive that the cardinality of the view $\{\text{year}, \text{employee}, \text{toDept}\}$ cannot exceed twice the cardinality of X .

The kD 's have been studied in the context of relational database theory, where they are also known as *numerical dependencies*. Grant and Minker [8] have proven that kD 's are not finitely axiomatizable, thus no fixed set of inference rules can be used to determine whether or not a given kD is logically implied by a set of kD 's. Nonetheless, a basic set of rules, which naturally extend those for FD's, was proposed in [8]. The rules we use, here generalized to work with the multi-dimensional lattice, are:

$$\begin{aligned} \text{R1: } & X \xrightarrow{k} Y \vdash X \oplus Z \xrightarrow{k} Y \oplus Z \\ \text{R2: } & X \xrightarrow{k} Y \wedge Y \xrightarrow{l} Z \vdash X \xrightarrow{k \cdot l} Y \oplus Z \\ \text{R3: } & X \xrightarrow{k} Y \oplus Z \vdash X \xrightarrow{k} Y \\ \text{R4: } & X \xrightarrow{k} Y \wedge X \xrightarrow{l} Z \vdash X \xrightarrow{k \cdot l} Y \oplus Z \\ \text{R5: } & X \oplus W \xrightarrow{k} Y \vdash X \xrightarrow{k \cdot \llbracket W \oplus Z \rrbracket^+} Y \oplus Z \end{aligned}$$

Note that the “union” rule R4 can be easily derived from R1 (“extension”), R2 (“transitivity”), and R3 (“decomposition”). As to R5, it may be proved by considering that, for each W and Z , we can write a dummy kD $X \xrightarrow{\llbracket W \oplus Z \rrbracket^+} W \oplus Z$, from which by applying R1 we obtain $X \xrightarrow{\llbracket W \oplus Z \rrbracket^+} W \oplus Z \oplus X$. On the other hand, by applying R1 to $X \oplus W \xrightarrow{k} Y$ we obtain $X \oplus W \oplus Z \xrightarrow{k} Y \oplus Z$. Applying R2 to $X \xrightarrow{\llbracket W \oplus Z \rrbracket^+} W \oplus Z \oplus X$ and $X \oplus W \oplus Z \xrightarrow{k} Y \oplus Z$ we have $X \xrightarrow{k \cdot \llbracket W \oplus Z \rrbracket^+} Y \oplus Z$, which proves R5.

The influence of kD 's on the determination of bounds is summarized by the following lemma.

Lemma 3. If $X \xrightarrow{k} Y$, then $\llbracket X \oplus Y \rrbracket^+ \leq k \cdot x^+$.

Proof. From Definition 5 it follows immediately that, if $X \xrightarrow{k} Y$, the cardinality of $X \oplus Y$ is related to the cardinality of X by inequality $\llbracket X \oplus Y \rrbracket \leq k \cdot x$. The inequalities on bounds follow immediately. \square

The following inequalities summarize Lemmas 2 and 3, thus relating upper bounds of view cardinalities to kD 's:

$$X \xrightarrow{k} Y \Rightarrow \llbracket X \oplus Y \rrbracket^+ \leq k \cdot x^+ \leq x^+ \cdot y^+ \quad (2)$$

where $k \leq y^+$ follows from the assumption of sound and minimal input.

4. The cover-based bounding strategy

In this section we introduce a bounding strategy, called *cover-based*, which relies on the concept of cover of a view to compute upper bounds.

Definition 6 (*Candidate set*). We call candidate set a couple $\mathcal{C} = (S, K)$, where S is a set of views and K is a set of kD 's. We denote with $\text{lhs}(K)$ and $\text{rhs}(K)$, respectively, the sets of views which are left- and right-hand sides for the kD 's in K ; besides, let $N_{\mathcal{C}} = S \cup \text{lhs}(K) \cup \text{rhs}(K)$.

Definition 7 (*Cover*). Let $V \in \mathcal{V}_D$ be a view on \mathcal{D} and $\mathcal{C} = (S, K)$ be a candidate set. \mathcal{C} is called a V -cover iff $V \preceq \oplus(N_{\mathcal{C}})$.

Thanks to Lemmas 1 and 2, V -covers can be used to bound from above the cardinality of V , since $v^+ \leq \llbracket \oplus(N_{\mathcal{C}}) \rrbracket^+$. More precisely, the cover-based bounding strategy cb computes v_{cb}^+ as:

$$v_{\text{cb}}^+ = \begin{cases} v^+ & \text{if } v^+ \in \mathcal{I} \\ \min\{u_{\text{cb}}(\mathcal{C}) : \mathcal{C} \text{ is a } V\text{-cover}\} & \text{if } v^+ \notin \mathcal{I} \end{cases} \quad (3)$$

where $u_{\text{cb}}(\mathcal{C})$ is the bound yielded by cover \mathcal{C} . In general, since $u_{\text{cb}}(\mathcal{C})$ depends in turn on the bounds of the views in \mathcal{C} , evaluating the cover-based bound may lead to a recursive computational flow; note however that the ‘‘case-0’’ of recursion, $v_{\text{cb}}^+ = v^+$, is correctly defined since we assumed the input \mathcal{I} to be minimal.

To start simple, assume that the V -cover \mathcal{C} does not include any kD , thus $N_{\mathcal{C}} = S$. In this case the bound provided by \mathcal{C} is directly derived from Lemma 2:

$$u_{\text{cb}}(\mathcal{C}) = u_{\text{cb}}((S, \emptyset)) = \prod_{w_j \in S'} w_{j,\text{cb}}^+ \quad (4)$$

where S' is defined as in Lemma 2.

Let us now turn to the more general (and complex) case when also kD 's are present: due to Eq. (2), it is clear that their presence can be exploited to strengthen bounds. Intuitively, if $V \preceq X \oplus Y$ and only x^+ and y^+ are known, the best we can do is to infer that $v \leq x^+ \cdot y^+$. On the other hand, if $X \xrightarrow{k} Y$ also holds, then the bound can be improved to $v \leq k \cdot x^+$, which can be much better than $x^+ \cdot y^+$. In the following we precisely characterize how sets of views and kD 's can be combined together in a common framework.

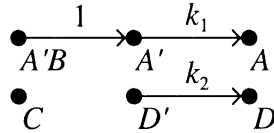


Fig. 3. The \mathcal{C} -graph associated to the candidate set in Example 6 (arrows denote kD 's).

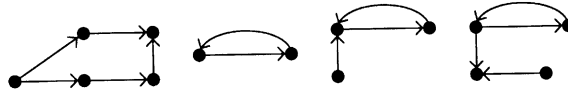


Fig. 4. Four examples of \mathcal{C} -graphs (from left to right): acyclic and reachable, cyclic and unreachable without roots, cyclic and reachable, cyclic and unreachable with root.

Definition 8 (*\mathcal{C} -graph*). The \mathcal{C} -graph $\mathcal{G}_{\mathcal{C}}$ associated with $\mathcal{C} = (S, K)$ is a labeled directed graph³ whose nodes correspond to the views in $N_{\mathcal{C}}$ and whose arcs are defined as follows:

- (1) For each $kD \ W \xrightarrow{k} Z \in K$ there is an arc, labeled k , from node W to node Z ;
- (2) For each pair of nodes W and Z such that $Z \preceq W$ there is an arc, labeled 1, from node W to node Z .

The nodes of $\mathcal{G}_{\mathcal{C}}$ with no incoming arcs are called the roots of $\mathcal{G}_{\mathcal{C}}$, and their set is denoted as $\text{root}(\mathcal{G}_{\mathcal{C}})$. The \mathcal{C} -graph $\mathcal{G}_{\mathcal{C}}$ is called *reachable* iff for each node W there exists at least one directed path from a root of $\mathcal{G}_{\mathcal{C}}$ to W .

Example 6. The reachable \mathcal{C} -graph associated to $\mathcal{C} = (\{A'B, C\}, \{A' \xrightarrow{k_1} A, D' \xrightarrow{k_2} D\})$ is depicted in Fig. 3. It is $\text{root}(\mathcal{G}_{\mathcal{C}}) = \{A'B, C, D'\}$.

Note that an acyclic \mathcal{C} -graph is always reachable, whereas the converse is not necessarily true; on the other hand, $\mathcal{G}_{\mathcal{C}}$ is reachable only if each connected component of $\mathcal{G}_{\mathcal{C}}$ has at least one root (see Fig. 4).

While each candidate set \mathcal{C} is univocally associated to one \mathcal{C} -graph, different candidate sets may be associated to the same \mathcal{C} -graph. We say that \mathcal{C} is *reduced* iff there exists no $\mathcal{C}' \subset \mathcal{C}$ such that $\mathcal{G}_{\mathcal{C}'} = \mathcal{G}_{\mathcal{C}}$ (for instance, $\mathcal{C}' = (\{C\}, \{A \xrightarrow{k} B\})$ is the reduced form of $\mathcal{C} = (\{A, B, C\}, \{A \xrightarrow{k} B\})$).

Lemma 4. A candidate set $\mathcal{C} = (S, K)$ is reduced iff $S \cap (\text{lhs}(K) \cup \text{rhs}(K)) = \emptyset$.

Proof (Only if). Trivial, since if $W \in S \cap (\text{lhs}(K) \cup \text{rhs}(K))$ then $\mathcal{C} = (S, K)$ and $\mathcal{C}' = (S - \{W\}, K)$ have the same \mathcal{C} -graph.

(If.) Assume that \mathcal{C} is not reduced. Then we can reduce either S or K and obtain the same \mathcal{C} -graph. Clearly, we cannot reduce K by dropping a kD , since the corresponding graph will have

³ Technically, $\mathcal{G}_{\mathcal{C}}$ is a *multi-graph*, since two arcs may share the same couple of nodes. This, however, does not influence the following arguments.

one arc less than $\mathcal{G}_{\mathcal{C}}$. Then, we can only remove a view W in S . This will leave unaltered the set of nodes only if W appears as the left or right side of some kD in K , as claimed. \square

The following theorem, whose proof is reported in the Appendix, precisely characterizes how the bound obtained from a candidate set is related to its \mathcal{C} -graph. For brevity, from now on we will write $\prod_E k_i$ to denote the product of all the labels of the arcs in set E .

Theorem 1. *Let V be a view, $\mathcal{C} = (S, K)$ be a V -cover, and $\mathcal{G}_{\mathcal{C}}$ be the \mathcal{C} -graph associated with \mathcal{C} . If $\mathcal{G}_{\mathcal{C}}$ is reachable it is:*

$$u_{\text{cb}}(\mathcal{C}) = u_{\text{cb}}((S, K)) = \prod_{E \in \mathcal{C}} k_i \cdot \prod_{W_j \in \text{root}(\mathcal{G}_{\mathcal{C}})} w_{j,\text{cb}}^+ \quad (5)$$

Example 7. Let $V = ABC$. Below we consider some examples of V -covers and, for each of them, show how the bound of v provided by Theorem 1 can be justified considering the lemmas proved so far and the inference rules for kD 's. In order to help the reader, Fig. 5 depicts the roll-up relationships between the views involved.

- $\mathcal{C}_1 = (\{ABCD\}, \emptyset)$ is a V -cover since $V \preceq \oplus(N_{\mathcal{C}_1}) = ABCD$. From Lemma 1 it is derived $abc \leq abcd^+$.
- $\mathcal{C}_2 = (\{AB, BC\}, \emptyset)$ is a V -cover since $V \preceq \oplus(N_{\mathcal{C}_2}) = ABC$. Since the natural join between two views is a subset of their Cartesian product, it is $abc \leq ab^+ \cdot bc^+$.
- $\mathcal{C}_3 = (\emptyset, \{AB \xrightarrow{k} C\})$. From Lemma 3 it immediately follows $abc \leq ab^+ \cdot k$.
- $\mathcal{C}_4 = (\emptyset, \{A \xrightarrow{k_1} B, B \xrightarrow{k_2} C\})$. By applying rule R2, we derive $A \xrightarrow{k_1 k_2} BC$, thus $abc \leq a^+ \cdot k_1 \cdot k_2$.
- $\mathcal{C}_5 = (\emptyset, \{A \xrightarrow{k_1} B, A \xrightarrow{k_2} C\})$. Rule R4 is now used to derive $A \xrightarrow{k_1 k_2} BC$, thus $abc \leq a^+ \cdot k_1 \cdot k_2$.
- $\mathcal{C}_6 = (\{A'B, C\}, \{A' \xrightarrow{k} A\})$. According to rule R1 it is $A'B \xrightarrow{k} AB$, and from Lemma 3 $ab^+ \leq k \cdot a'b^+$. On the other hand, $abc \leq ab^+ \cdot c^+$, thus $abc \leq k \cdot a'b^+ \cdot c^+$.

Note that Theorem 1 correctly captures also the case when no kD 's are present, thus generalizing Eq. (4); in fact, in this case $\text{root}(\mathcal{G}_{\mathcal{C}}) = S'$, where S' is as in Lemma 2. The reason why

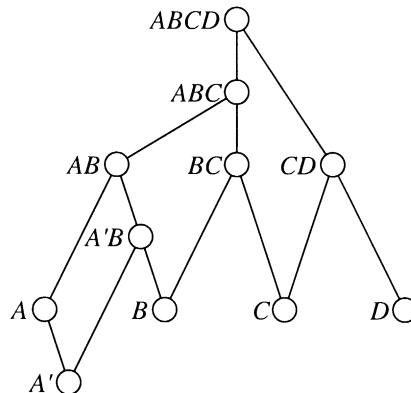


Fig. 5. Roll-up relationships of views in Example 7.

Theorem 1 applies only to reachable \mathcal{C} -graphs can be understood by means of a simple example. Let $\mathcal{C} = (\emptyset, \{A \xrightarrow{k_1} B, B \xrightarrow{k_2} A\})$, thus $\mathcal{G}_{\mathcal{C}}$ has no roots and is not reachable (see Fig. 4). If we try to apply Eq. (5) to such graph we obtain $u(\mathcal{C}) = k_1 \cdot k_2$, which is clearly wrong. On the other hand, a non-reachable \mathcal{C} -graph can be easily transformed into a reachable one by dropping one or more kD 's (for instance, $B \xrightarrow{k_2} A$ in the example above): for this reason, all the \mathcal{C} -graphs considered from now on are implicitly assumed to be reachable. Besides, since according to Theorem 1 the bound yielded by a candidate set only depends on its associated \mathcal{C} -graph, in the following we will always consider reduced candidate sets.

In order to ensure that the cover-based strategy is consistent, we need to verify that the properties expressed by Lemmas 1–3 are valid also for the cover-based upper bounds.

Lemma 5. *If $W \preceq Z$, then $w_{cb}^+ \leq z_{cb}^+$. If S is a set of views, then $\llbracket \oplus(S) \rrbracket_{cb}^+ \leq \prod_{W_i \in S'} w_{i,cb}^+$. If $X \xrightarrow{k} Y$, then $\llbracket X \oplus Y \rrbracket_{cb}^+ \leq k \cdot x_{cb}^+$.*

Proof. The first property is obvious, since the set of the W -covers includes Z and all the Z -covers. The second one derives from the fact that $\mathcal{C} = (S, \emptyset)$ is a cover for $\oplus(S) = \oplus(S')$, and $u_{cb}(\mathcal{C}) = \prod_{W_i \in S'} w_{i,cb}^+$. The third one is true since $\mathcal{C} = (\emptyset, \{X \xrightarrow{k} Y\})$ is a cover for $X \oplus Y$, and $u_{cb}(\mathcal{C}) = k \cdot x_{cb}^+$. \square

5. Domination between candidate sets

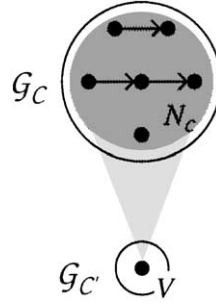
In this section we complete the presentation of the cover-based bounding strategy by presenting some formal results useful to limit the complexity of algorithms aiming to compute upper bounds of view cardinalities. The problem is inherently difficult, since the space of V -covers to be analyzed in order to determine v_{cb}^+ has exponential size. On the other hand, under some circumstances, a V -cover \mathcal{C} can be discarded from the search space without even computing its associated bound $u_{cb}(\mathcal{C})$.

Definition 9 (Domination). Let \mathcal{C} and \mathcal{C}' be two candidate sets. We say that \mathcal{C}' dominates \mathcal{C} , written $\mathcal{C}' \sqsubseteq \mathcal{C}$, iff $u_{cb}(\mathcal{C}') \leq u_{cb}(\mathcal{C})$ for every possible input \mathcal{I} .

Again, let us consider the basic case where no kD 's are present first. Note that, when $\mathcal{C}_1 = (S_1, \emptyset)$ and $\mathcal{C}_2 = (S_2, \emptyset)$, \mathcal{C}_1 can dominate \mathcal{C}_2 only if $\oplus(S_1) \preceq \oplus(S_2)$. In fact, if $\oplus(S_1) \not\preceq \oplus(S_2)$, there exists at least one attribute $A \in \oplus(S_1)$ which does not belong to $\oplus(S_2)$: intuitively, this means that the bound provided by S_2 is independent of the cardinality of A , whereas the bound yielded by S_1 is not, which is enough to show that \mathcal{C}_1 cannot dominate \mathcal{C}_2 .

The following theorem, whose proof is reported in the Appendix, provides a necessary and sufficient condition for checking dominance between two candidate sets without kD 's.

Theorem 2 (Domination between sets of views). *Let $\mathcal{C}_1 = (S_1, \emptyset)$ and $\mathcal{C}_2 = (S_2, \emptyset)$, with $S_1 = \{W_{1,1}, \dots, W_{1,i}, \dots, W_{1,m}\}$, and let S'_1 as in Lemma 2, thus $S'_1 = \{W_{1,1}, \dots, W_{1,i}, \dots, W_{1,m'}\}$ with $m' \leq m$. \mathcal{C}_1 dominates \mathcal{C}_2 iff S_2 can be partitioned into m' subsets $S_{2,1}, \dots, S_{2,m'}$ such that $W_{1,i} \preceq \oplus(S_{2,i}) \forall i = 1, \dots, m'$.*

Fig. 6. Two \mathcal{C} -graphs in Lemma 6.

Example 8. It is $\{A'B, C\} \sqsubseteq \{AB, CD, E\}$, since $A'B \preceq AB$ and $C \preceq \oplus(\{CD, E\}) = CDE$.

Domination between candidate sets including kD 's is much more complex to analyze; in the following we present a set of partial results, beginning with a simple corollary of Theorem 1.

Corollary 1. Let \mathcal{C} be a candidate set associated to a non-forest⁴ \mathcal{C} -graph $\mathcal{G}_{\mathcal{C}}$, and let \mathcal{C}' be a candidate set associated to a forest \mathcal{C} -graph $\mathcal{G}_{\mathcal{C}'}$ such that $N_{\mathcal{C}'} = N_{\mathcal{C}}$ and $E_{\mathcal{C}'} \subset E_{\mathcal{C}}$. Then $\mathcal{C}' \sqsubseteq \mathcal{C}$.

As discussed in the proof of Theorem 1, a forest \mathcal{C} -graph satisfying the conditions above exists when the reason why $\mathcal{G}_{\mathcal{C}}$ is not a forest is that two or more arcs corresponding to kD 's converge in the same node. Since Corollary 1 states that this kind of non-forest \mathcal{C} -graphs are always dominated, in the following they will not be considered. On the other hand, the other reason why $\mathcal{G}_{\mathcal{C}}$ is not a forest may be that at least two arcs labeled 1 converge in the same node: in this case, since removing such arcs would violate the very definition of \mathcal{C} -graph, Corollary 1 cannot be applied.

Next lemma expresses a sufficient condition for domination when one of the two candidate sets includes no kD 's.

Lemma 6. Let $\mathcal{C}' = (\{V\}, \emptyset)$ and $\mathcal{C} = (S, K)$ be two candidate sets. If $V \preceq \oplus(N_{\mathcal{C}})$ then $\mathcal{C}' \sqsubseteq \mathcal{C}$.

Proof. \mathcal{C} is a V -cover since $V \preceq \oplus(N_{\mathcal{C}})$, hence, from Theorem 1 it is $v_{cb}^+ \leq u_{cb}(\mathcal{C})$. Since $u_{cb}(\mathcal{C}') = v_{cb}^+$, the result follows immediately (see Fig. 6). \square

The following definition is preliminary to Lemma 7:

Definition 10 (\mathcal{C} -subgraph). Given $\mathcal{C} = (S, K)$ with associated \mathcal{C} -graph $\mathcal{G}_{\mathcal{C}} = (N_{\mathcal{C}}, E_{\mathcal{C}})$, let $N_{\mathcal{C}_1} \subset N_{\mathcal{C}}$. The \mathcal{C} -subgraph induced by $N_{\mathcal{C}_1}$ is defined as $\mathcal{G}_{\mathcal{C}_1} = (N_{\mathcal{C}_1}, E_{\mathcal{C}_1})$, where $E_{\mathcal{C}_1} = \{(X, Y) \in E_{\mathcal{C}} : X \in N_{\mathcal{C}_1} \wedge Y \in N_{\mathcal{C}_1}\}$. The \mathcal{C} -subgraph is called *proper* iff $\forall (X, Y) \in E_{\mathcal{C}}$ if $Y \in N_{\mathcal{C}_1}$ then $X \in N_{\mathcal{C}_1}$, i.e., iff there is no arc entering $\mathcal{G}_{\mathcal{C}_1}$. The (possibly empty) set

⁴ A forest is a set of disjoint directed trees.

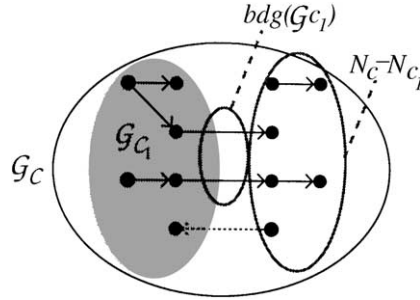


Fig. 7. A proper \mathcal{C} -subgraph.

$\text{bdg}(\mathcal{G}_{\mathcal{C}_1}) = \{(X, Y) \in E_{\mathcal{C}} : X \in N_{\mathcal{C}_1} \wedge Y \in N_{\mathcal{C}} - N_{\mathcal{C}_1}\}$, i.e., the set of the arcs exiting $\mathcal{G}_{\mathcal{C}_1}$, is called the *bridge* induced by $N_{\mathcal{C}_1}$.

Example 9. Consider the \mathcal{C} -graph shown in Fig. 7. The \mathcal{C} -subgraph induced by $N_{\mathcal{C}_1}$ is proper; if also the arc shown in dotted line were included in \mathcal{C} , $\mathcal{G}_{\mathcal{C}_1}$ would not be proper.

Lemma 7. Let $\mathcal{C}' = (\emptyset, \{V_1 \xrightarrow{k} V_2\})$ and $\mathcal{C} = (S, K)$ be two candidate sets. If $N_{\mathcal{C}}$ can be partitioned into two subsets $N_{\mathcal{C}_1}$ and $N_{\mathcal{C}_2}$ such that:

- (1) $\mathcal{G}_{\mathcal{C}_1}$ is proper; and
 - (2) $V_1 \preceq \oplus(N_{\mathcal{C}_1})$; and
 - (3) $V_2 \preceq \oplus(N_{\mathcal{C}_2})$; and
 - (4) $\oplus(\text{lhs}(\text{bdg}(\mathcal{G}_{\mathcal{C}_1}))) \ominus V_1 \preceq \oplus(N_{\mathcal{C}_u} \cup \text{rhs}(\text{bdg}(\mathcal{G}_{\mathcal{C}_1})))$, where $N_{\mathcal{C}_u} \subseteq N_{\mathcal{C}_2}$ is the subset of nodes of $\mathcal{G}_{\mathcal{C}_2}$ that are not reachable from $\text{bdg}(\mathcal{G}_{\mathcal{C}_1})$
- then $\mathcal{C}' \sqsubseteq \mathcal{C}$.

Example 10. Given $\mathcal{C}' = (\emptyset, \{AB \xrightarrow{k_1} C\})$ and $\mathcal{C} = (\{A\}, \{BE \xrightarrow{k_2} C, F \xrightarrow{k_3} D, D \xrightarrow{k_4} E\})$ it is $\mathcal{C}' \sqsubseteq \mathcal{C}$. In fact, if $N_{\mathcal{C}_1} = \{A, BE\}$ and $N_{\mathcal{C}_2} = \{C, D, E, F\}$ (see Fig. 8), it is:

- $AB = V_1 \preceq \oplus(N_{\mathcal{C}_1}) = ABE$;
- $C = V_2 \preceq \oplus(N_{\mathcal{C}_2}) = CDEF$;
- $\text{lhs}(\text{bdg}(\mathcal{G}_{\mathcal{C}_1})) = BE$, $\oplus(\text{lhs}(\text{bdg}(\mathcal{G}_{\mathcal{C}_1}))) \ominus V_1 = E$, $\oplus(N_{\mathcal{C}_u}) = DF$, $\oplus(\text{rhs}(\text{bdg}(\mathcal{G}_{\mathcal{C}_1}))) = E$, and finally $E \preceq DEF$.

Since both \mathcal{C}' and \mathcal{C} are covers for ABC , we derive that $abc^+ \leq ab^+ \cdot k_1 \leq a^+ \cdot be^+ \cdot f^+ \cdot k_2 \cdot k_3 \cdot k_4$.

From this lemma we may finally derive the following theorem, stating a sufficient condition which may be used to recursively prove domination in the general case. Both Lemma 7 and Theorem 3 are proved in the Appendix.

Theorem 3. Let $\mathcal{C}' = (S', K')$ and $\mathcal{C} = (S, K)$ be two candidate sets. Let $V \in N_{\mathcal{C}'}$ be a node of $\mathcal{G}_{\mathcal{C}'}$ with no outgoing arcs, and let $\mathcal{G}_{\mathcal{C}'_1}$ be the \mathcal{C} -subgraph induced by $N_{\mathcal{C}'} - \{V\}$. If $N_{\mathcal{C}}$ can be partitioned into two subsets $N_{\mathcal{C}_1}$ and $N_{\mathcal{C}_2}$ such that:

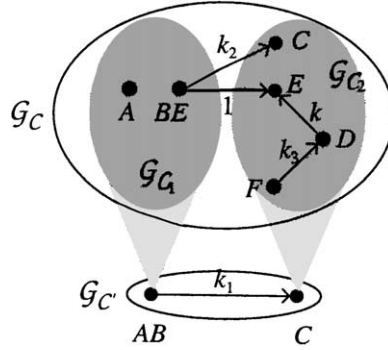


Fig. 8. The two \mathcal{C} -graphs in Example 10.

- (1) The \mathcal{C} -subgraph \mathcal{G}_{ϕ_1} induced by N_{ϕ_1} is proper; and
- (2) $\mathcal{C}'_1 \sqsubseteq \mathcal{C}_1$; and
- (3) $V \preceq \oplus(N_{\phi_2})$; and
- (4) If $\text{bdg}(\mathcal{G}_{\phi_1}) \neq \emptyset$ then $\exists(W, V) \in E_{\mathcal{C}'} : \oplus(\text{lhs}(\text{bdg}(\mathcal{G}_{\phi_1}))) \ominus W \preceq \oplus(N_{\phi_u} \cup \text{rhs}(\text{bdg}(\mathcal{G}_{\phi_u})))$, where $N_{\phi_u} \subseteq N_{\phi_2}$ is the subset of nodes of \mathcal{G}_{ϕ_2} that are not reachable from $\text{bdg}(\mathcal{G}_{\phi_1})$ then $\mathcal{C}' \sqsubseteq \mathcal{C}$.

Example 11. Given $\mathcal{C}' = (\{AB, D\}, \{A' \xrightarrow{k_1} C\})$ and $\mathcal{C} = (\{A, BE\}, \{A'B \xrightarrow{k_2} E, C \xrightarrow{k_3} BF, D' \xrightarrow{k_4} D\})$ it is $\mathcal{C}' \sqsubseteq \mathcal{C}$. In fact, as shown in Fig. 9, it is possible to apply recursively Theorem 3 to progressively smaller \mathcal{C} -graphs:

- *Step 1.* Let $V = D$ and $N_{\phi_2} = \{D, D'\}$. Conditions (1), (3) and (4) are immediately verified, we have to verify that $(\{AB\}, \{A' \xrightarrow{k_1} C\}) \sqsubseteq (\{A, BE\}, \{A'B \xrightarrow{k_2} E, C \xrightarrow{k_3} BF\})$.
- *Step 2.* Let $V = C$ and $N_{\phi_2} = \{C, E, BF\}$. Conditions (1) and (3) are immediately verified, as to condition (4) it is $W = A', \oplus(\text{lhs}(\text{bdg}(\mathcal{G}_{\phi_1}))) \ominus W = BC \preceq \oplus(N_{\phi_u} \cup \text{rhs}(\text{bdg}(\mathcal{G}_{\phi_u}))) = BCF$. It is left to verify that $(\{AB, A'\}, \emptyset) \sqsubseteq (\{A, BE, A'B\}, \emptyset)$, which is immediate due to Theorem 2.

Since both \mathcal{C}' and \mathcal{C} are covers for $ABCD$, we derive that $abcd^+ \leq ab^+ \cdot d^+ \cdot k_1 \leq a^+ \cdot be^+ \cdot d'b^+ \cdot c^+ \cdot d'^+ \cdot k_2 \cdot k_3 \cdot k_4$.

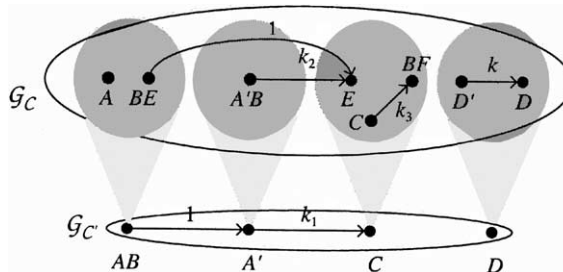


Fig. 9. The two \mathcal{C} -graphs in Example 11.

6. The estimator

Assuming that effective bounds have been derived, cardinality estimation must be based on a probabilistic model to derive an estimate, \bar{v} , of the cardinality of view V . The model we adopt here is based on the Cardenas' formula [2], which states that, when throwing N distinct objects into B buckets, the expected number of buckets in which at least one object will fall can be estimated as:

$$\Phi(B, N) = B \left(1 - \left(1 - \frac{1}{B} \right)^N \right) \leq \min\{B, N\} \quad (6)$$

Within the approach proposed in [17], (6) is used to estimate v by relying on the maximum cardinality of V , defined as the Cartesian product of the cardinalities of the attributes in V , $v_{\max} = \prod_{A_i \in V} a_i$, and on the cardinality of the base cube, $d = \llbracket \dim(\mathcal{D}) \rrbracket$, that is:

$$\bar{v}_{\text{sdnr}} = \Phi(v_{\max}, d) \leq \min\{v_{\max}, d\} \quad (7)$$

This formula turns out to significantly overestimate the cardinalities and can easily lead to violate the constraint $\bar{v}_{\text{sdnr}} \leq v^+$.

In our approach, denoted se (“safe-estimate”), the above estimate is improved in two ways: by replacing v_{\max} with the upper bound computed for v , for instance v_{cb}^+ , as a measure of the maximum cardinality of V , and by replacing the cardinality of the base cube d with an estimate, \bar{w}_{se} , of the cardinality of a view W such that $V \preceq W$. This leads to:

$$\bar{v}_{\text{se}} = \Phi(v_{\text{cb}}^+, \bar{w}_{\text{se}}) \leq \min\{v_{\text{cb}}^+, \bar{w}_{\text{se}}\} \quad (8)$$

Since both v_{cb}^+ and \bar{w}_{se} can be considerably lower than v_{\max} and d , respectively, it is usually the case that $\bar{v}_{\text{se}} \ll \bar{v}_{\text{sdnr}}$. The rationale for (8) is that we can view the problem of estimating v as the one of distributing the tuples of view W , which are estimated to be \bar{w}_{se} , over a number of v_{cb}^+ “buckets”.

Due to the need to know \bar{w}_{se} , it is obvious that our estimation process must move downward from the top of the lattice (whose cardinality d is typically known) following a path leading to V . Clearly, this represents a simplification of the correct estimation procedure, which would require to determine \bar{v} by following *all* the paths from $\dim(\mathcal{D})$ to V . On the other hand, this would lead to combinatorial explosion and necessitate of highly complex probabilistic models that are well beyond the current state-of-the-art knowledge.

From a more practical (numerical) point of view, it should be noted that moving from upper bounds to estimates leads to significant differences under specific conditions only. Two relevant cases should be considered, which arise from the limit behavior of Cardenas' formula:

- (1) When $\bar{w}_{\text{se}} \leq 0.1 v_{\text{cb}}^+$ it is $\bar{v}_{\text{se}} \approx \bar{w}_{\text{se}}$
- (2) When $\bar{w}_{\text{se}} \geq 3 v_{\text{cb}}^+$ it is $\bar{v}_{\text{se}} \approx v_{\text{cb}}^+$

The values 0.1 and 3 can thus be used to predict whether the estimator will deliver results which substantially differ from those directly obtainable from the bounder.

Example 12. In the Transfers scheme, we consider three input situations:

$$\begin{aligned} \mathcal{I}_1 = & \{ \{ \text{date} \}^+ = 10^3, \{ \text{year} \}^+ = 3, \{ \text{employee} \}^+ = 10^4, \\ & \{ \text{fromOffice} \}^+ = \{ \text{toOffice} \}^+ = 10^3, \{ \text{fromCity} \}^+ = \{ \text{toCity} \}^+ = 10, \\ & \{ \text{fromDept} \}^+ = \{ \text{toDept} \}^+ = 10 \} \end{aligned}$$

$$\mathcal{I}_2 = \mathcal{I}_1 \cup \{ \{ \text{employee}, \text{year} \} \xrightarrow{4} \{ \text{fromOffice}, \text{toOffice}, \text{date} \} \}$$

$$\begin{aligned} \mathcal{I}_3 = & \mathcal{I}_2 \cup \{ \{ \text{fromCity}, \text{fromDept} \}^+ = 40, \{ \text{toCity}, \text{toDept} \}^+ = 40, \\ & \{ \text{fromCity}, \text{fromDept} \} \xrightarrow{2} \{ \text{toCity}, \text{toDept} \}, \{ \text{fromCity}, \text{fromDept} \} \xrightarrow{30} \{ \text{fromOffice} \}, \\ & \{ \text{toCity}, \text{toDept} \} \xrightarrow{30} \{ \text{toOffice} \} \} \end{aligned}$$

Let $W = \dim(\mathcal{D}) = \{ \text{date}, \text{employee}, \text{fromOffice}, \text{toOffice} \}$ be the base cube and $V = \{ \text{fromOffice}, \text{toOffice} \}$ be the view whose cardinality is to be estimated.

- When $\mathcal{I} = \mathcal{I}_1$, the best W - and V -cover are, respectively,

$$\mathcal{C}_1 = (\{ \{ \text{date} \}, \{ \text{employee} \}, \{ \text{fromOffice} \}, \{ \text{toOffice} \} \}, \emptyset)$$

$$\mathcal{C}_2 = (\{ \{ \text{fromOffice} \}, \{ \text{toOffice} \} \}, \emptyset)$$

Clearly, $\mathcal{C}_2 \sqsubseteq \mathcal{C}_1$.

- When $\mathcal{I} = \mathcal{I}_2$, the best cover for both W and V is

$$\mathcal{C}_3 = (\emptyset, \{ \{ \text{employee}, \text{year} \} \xrightarrow{4} \{ \text{fromOffice}, \text{toOffice}, \text{date} \} \})$$

No domination relationship exists between \mathcal{C}_1 and \mathcal{C}_3 and between \mathcal{C}_2 and \mathcal{C}_3 .

- When $\mathcal{I} = \mathcal{I}_3$, the best W -cover is still \mathcal{C}_3 , while the best V -cover is

$$\begin{aligned} \mathcal{C}_4 = & (\emptyset, \{ \{ \text{fromCity}, \text{fromDept} \} \xrightarrow{2} \{ \text{toCity}, \text{toDept} \}, \{ \text{fromCity}, \text{fromDept} \} \xrightarrow{30} \{ \text{fromOffice} \}, \\ & \{ \text{toCity}, \text{toDept} \} \xrightarrow{30} \{ \text{toOffice} \} \}) \end{aligned}$$

Still, no domination relationship exists with the other candidate sets.

Table 1 shows how the upper bound w_{cb}^+ of W , the upper bound v_{cb}^+ of V and the estimate \bar{v}_{se} improve as new cardinality constraints are progressively supplied. The estimate \bar{v}_{se} is based on the estimate of w , \bar{w}_{se} , which is assumed to be equal to its upper bound w_{cb}^+ . Note how, in the absence of cardinality constraints (i.e. if only the upper bounds of the cardinalities of all the single attributes in the dimensional scheme are known) the only reasonable estimate that can be inferred is $\bar{v}_{sdnr} = 10^6$ obtained by formula (6) using $B = 10^6$ and $N = 10^{13}$. Even assuming that the cardinality of the base cube is known does not improve the estimate significantly; for example, assuming $w = 1.2 \times 10^5$, the estimate obtained using directly the Cardenas formula as in [17] is $\bar{v}_{sdnr} = 1.2 \times 10^5$, that is much worse than those obtained exploiting the constraints in \mathcal{I}_2 and \mathcal{I}_3 .

Table 1
Improving upper bounds and estimates for increasing domain-derived information

Input	w_{cb}^+	v_{cb}^+	w_{se}	\bar{v}_{se}
\mathcal{I}_1	10^{13}	10^6	10^{13}	10^6
\mathcal{I}_2	1.2×10^5	1.2×10^5	1.2×10^5	7.6×10^4
\mathcal{I}_3	1.2×10^5	7.2×10^4	1.2×10^5	5.8×10^4

7. Conclusions and open issues

In this paper we have shown how cardinality constraints derived from the application domain may be employed to determine effective bounds on the cardinality of aggregate views and how, in turn, such bounds can be used to estimate the cardinality of the views. In order to improve the approach effectiveness, some issues still need to be investigated. In the following we briefly discuss those we believe to be crucial:

- *Computation.* The utility of an approach to cardinality estimation depends on the efficiency of its computation. When no kD 's are included among the input constraints \mathcal{I} , the search for non-redundant upper bounds can be restricted to the set of *minimal* V -covers, where a minimal V -cover is one which is not dominated by any other V -cover and whose views are all constrained within \mathcal{I} [3]; in the branch-and-bound algorithm sketched in [4], a careful enumeration of minimal covers allows to significantly reduce the (otherwise exponential) search space. On the other hand, considering minimality issues when kD 's are involved is much more complex for two main reasons: firstly, a kD may be useful to compute an optimal cover even if the view on its left-hand side is not constrained in \mathcal{I} ; secondly, the utility of a kD in determining a cover can be evaluated only by considering the other kD 's included in the cover itself. Currently we are inclined to pursue a branch-and-bound approach which enumerates forest V -covers by assembling “useful” trees out of the constraints in \mathcal{I} . Once a tree covering a view $W \preceq V$ through a set of constraints $\mathcal{I}' \subset \mathcal{I}$ has been built, the problem is reduced to covering view $V \ominus W$ through the constraints in $\mathcal{I} - \mathcal{I}'$, and so on recursively.
- *Cardinality constraints.* The input knowledge may be further extended by considering other forms of cardinality constraints which are typically known to the experts of the application domain. For instance, while in this paper we have defined kD 's to express *bounds* on the ratio between the cardinalities of two views, they may also be used to denote the *average* of such ratio; while this kind of knowledge cannot be used by the bouncer, it allows the cardinality estimations to be improved. For instance, knowing that the average number of transfers for each employee on each year is 2, would allow the cardinality of the base cube to be estimated as twice the cardinality of view {employee, year}.
- *Probabilistic estimates.* Estimates based on Cardenas' formula can be improved in several ways. In particular, information on *lower bounds* of cardinalities could be considered by exploiting the results in [6], as well as information concerning the distribution of attribute values over their domains. Obviously, this requires to develop a bounding strategy for computing lower bounds; [3] presents some results in this direction.

Appendix A

Proof of Theorem 1. We will first prove that Eq. (5) holds when $\mathcal{G}_\mathcal{C}$ is a *forest*, i.e., a set of pairwise disjoint rooted trees. Then we will generalize to arbitrary \mathcal{C} -graphs.

Assume that $\mathcal{G}_\mathcal{C}$ is a forest. The set $\text{root}(\mathcal{G}_\mathcal{C})$ can be partitioned into two subsets: $\text{alone}(\mathcal{G}_\mathcal{C})$, which contains the “stand-alone” nodes, i.e. those for which the tree consists only of a root node, and $\alpha = \text{root}(\mathcal{G}_\mathcal{C}) - \text{alone}(\mathcal{G}_\mathcal{C})$, which contains the other roots (see Fig. 10). Let (N_j, E_j) be the tree rooted in $W_j \in \alpha$. Now, since $N_\mathcal{C} = \text{alone}(\mathcal{G}_\mathcal{C}) \cup (\bigcup_{j:W_j \in \alpha} N_j)$, it is

$$\llbracket \oplus (N_\mathcal{C}) \rrbracket = \llbracket \oplus (\text{alone}(\mathcal{G}_\mathcal{C})) \oplus (\oplus_{j:W_j \in \alpha} (\oplus(N_j))) \rrbracket \leq \prod_{W_j \in \text{alone}(\mathcal{G}_\mathcal{C})} w_j^+ \cdot \prod_{j:W_j \in \alpha} \llbracket \oplus (N_j) \rrbracket$$

By repeatedly applying rules R2 and R4 to (N_j, E_j) it is derived $W_j \xrightarrow{\prod_{E_j} k_i} \oplus (N_j - \{W_j\})$, hence for Lemma 3

$$\llbracket \oplus (N_j) \rrbracket \leq w_j^+ \cdot \prod_{E_j} k_i$$

which leads to

$$\llbracket \oplus (N_\mathcal{C}) \rrbracket \leq \prod_{W_j \in \text{alone}(\mathcal{G}_\mathcal{C})} w_j^+ \cdot \prod_{j:W_j \in \alpha} \left(w_j^+ \cdot \prod_{E_j} k_i \right) = \prod_{W_j \in \text{root}(\mathcal{G}_\mathcal{C})} w_j^+ \cdot \prod_{E_\mathcal{C}} k_i$$

Since by Definition 7 it is $V \preceq \oplus(N_\mathcal{C})$, we obtain $v \leq \llbracket \oplus (N_\mathcal{C}) \rrbracket$; thus, $u_{\text{cb}}(\mathcal{C})$ can be written as claimed.

When $\mathcal{G}_\mathcal{C} = (N_\mathcal{C}, E_\mathcal{C})$ is reachable but is not a forest, it is always possible to find $E'_\mathcal{C} \subset E_\mathcal{C}$ such that the graph $\mathcal{G}'_\mathcal{C} = (N_\mathcal{C}, E'_\mathcal{C})$ is a forest with the same roots as $\mathcal{G}_\mathcal{C}$. If $\mathcal{G}'_\mathcal{C}$ satisfies the definition of \mathcal{C} -graph, meaning that it is obtained by dropping no arcs labeled 1, it is obviously $v \leq u_{\text{cb}}(\mathcal{C}') \leq u_{\text{cb}}(\mathcal{C})$ (the first inequality follows from the first part of the proof, since $\mathcal{G}'_\mathcal{C}$ is a forest and $V \preceq \oplus(N_\mathcal{C})$; the second inequality follows from $\text{root}(\mathcal{G}'_\mathcal{C}) = \text{root}(\mathcal{G}_\mathcal{C})$ and from the observation that $E'_\mathcal{C} \subset E_\mathcal{C}$ implies $\prod_{E'_\mathcal{C}} k_i < \prod_{E_\mathcal{C}} k_i$).

On the other hand, if a forest $\mathcal{G}'_\mathcal{C}$ can be obtained from $\mathcal{G}_\mathcal{C}$ only by dropping one or more arcs labeled 1, Eq. (5) is still valid since arcs labeled 1 do not contribute to the bound expressed by $u_{\text{cb}}(\mathcal{C})$. \square

Proof of Theorem 2. (If.) If $W_{1,i} \preceq \oplus(S_{2,i}) \forall i = 1, \dots, m'$ the result directly follows from Lemma 2.

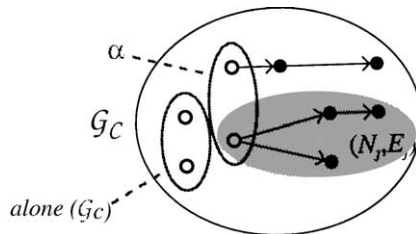


Fig. 10. A forest \mathcal{C} -graph; white circles denote views in $\text{root}(\mathcal{G}_\mathcal{C})$.

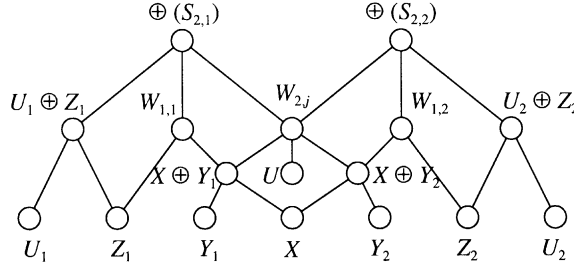


Fig. 11. Roll-up relationships of views in the Proof of Theorem 2.

(Only if.) We only provide a sketch of the complete proof. Without loss of generality assume $m = m' = 2$ and $\oplus(S_1) \preceq \oplus(S_2)$. If S_2 is not partitionable as required by the theorem there exists a view $W_{2,j}$ that is used to cover both $W_{1,1}$ and $W_{1,2}$. We can write $W_{1,1} = X \oplus Y_1 \oplus Z_1$, $W_{1,2} = X \oplus Y_2 \oplus Z_2$, $W_{2,j} = U \oplus X \oplus Y_1 \oplus Y_2$, $\oplus(S_{2,1}) = W_{2,j} \oplus (Z_1 \oplus U_1)$, and $\oplus(S_{2,2}) = W_{2,j} \oplus (Z_2 \oplus U_2)$, which guarantees that $W_{1,1} \preceq \oplus(S_{2,1})$ and $W_{1,2} \preceq \oplus(S_{2,2})$ due to the definition of \oplus (see Fig. 11). In the complete proof it is shown that there exist legal instances such that $w_{1,1}^+ \approx \llbracket X \oplus Y_1 \rrbracket^+ \cdot z_1^+$, $w_{1,2}^+ \approx \llbracket X \oplus Y_2 \rrbracket^+ \cdot z_2^+$, $w_{2,j}^+ \approx \llbracket X \oplus Y_1 \oplus Y_2 \rrbracket^+$, $\llbracket Z_1 \oplus U_1 \rrbracket^+ \approx z_1^+$, and $\llbracket Z_2 \oplus U_2 \rrbracket^+ \approx z_2^+$. Under these conditions, the inequality that should hold for domination is: $w_{1,1}^+ \cdot w_{1,2}^+ = (\llbracket X \oplus Y_1 \rrbracket^+ \cdot z_1^+) \cdot (\llbracket X \oplus Y_2 \rrbracket^+ \cdot z_2^+) \leq \llbracket X \oplus Y_1 \oplus Y_2 \rrbracket^+ \cdot z_1^+ \cdot z_2^+ = \llbracket \oplus(S_2) \rrbracket^+$, that is $\llbracket X \oplus Y_1 \rrbracket^+ \cdot \llbracket X \oplus Y_2 \rrbracket^+ \leq \llbracket X \oplus Y_1 \oplus Y_2 \rrbracket^+$, which can be easily invalidated. \square

Proof of Lemma 7. We have to prove that

$$u_{cb}(\mathcal{C}') = k' \cdot v_{1,cb}^+ \leq \prod_{E_{\mathcal{C}}} k_i \cdot \prod_{W_i \in \text{root}(\mathcal{G}_{\mathcal{C}})} w_{i,cb}^+ = u_{cb}(\mathcal{C})$$

Let $N_{\mathcal{C}_r} = N_{\mathcal{C}_2} - N_{\mathcal{C}_u}$ be the set of nodes of $\mathcal{G}_{\mathcal{C}_2}$ reachable from $\text{bdg}(\mathcal{G}_{\mathcal{C}_1})$; since $\mathcal{G}_{\mathcal{C}_1}$ is proper it is $\text{bdg}(\mathcal{G}_{\mathcal{C}_2}) = \emptyset$, thus we may write (see Fig. 12):

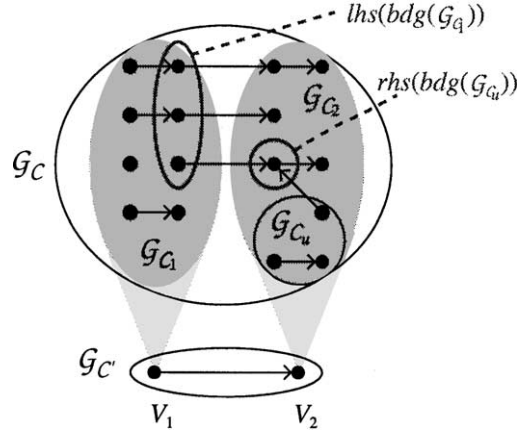
$$E_{\mathcal{C}} - E_{\mathcal{C}_u} = \text{bdg}(\mathcal{G}_{\mathcal{C}_1}) \cup E_{\mathcal{C}_2} = \text{bdg}(\mathcal{G}_{\mathcal{C}_1}) \cup E_{\mathcal{C}_r} \cup \text{bdg}(\mathcal{G}_{\mathcal{C}_u}) \cup E_{\mathcal{C}_u}$$

where $E_{\mathcal{C}_r}$ and $E_{\mathcal{C}_u}$ are, respectively, the arcs in $\mathcal{G}_{\mathcal{C}_r}$ and $\mathcal{G}_{\mathcal{C}_u}$. In fact, it is necessarily $\text{bdg}(\mathcal{G}_{\mathcal{C}_r}) = \emptyset$ since, otherwise, the nodes in $N_{\mathcal{C}_u}$ would become reachable. We can now rewrite $u_{cb}(\mathcal{C})$ as:

$$\begin{aligned} u_{cb}(\mathcal{C}) &= \prod_{E_{\mathcal{C}}} k_i \cdot \prod_{W_i \in \text{root}(\mathcal{G}_{\mathcal{C}})} w_{i,cb}^+ = \left(\prod_{E_{\mathcal{C}_1}} k_i \cdot \prod_{W_i \in \text{root}(\mathcal{G}_{\mathcal{C}_1})} w_{i,cb}^+ \right) \cdot \prod_{E_{\mathcal{C}} - E_{\mathcal{C}_u}} k_i \cdot \prod_{W_i \in \text{root}(\mathcal{G}_{\mathcal{C}_u})} w_{i,cb}^+ \\ &= u_{cb}(\mathcal{C}_1) \cdot \prod_{\text{bdg}(\mathcal{G}_{\mathcal{C}_1}) \cup E_{\mathcal{C}_r}} k_i \cdot \prod_{\text{bdg}(\mathcal{G}_{\mathcal{C}_u}) \cup E_{\mathcal{C}_u}} k_i \cdot \prod_{W_i \in \text{root}(\mathcal{G}_{\mathcal{C}_u})} w_{i,cb}^+ \end{aligned}$$

Since \mathcal{C}_1 is a V_1 -cover, it is $v_{1,cb}^+ \leq u_{cb}(\mathcal{C}_1)$. Thus, letting $h = \prod_{\text{bdg}(\mathcal{G}_{\mathcal{C}_1}) \cup E_{\mathcal{C}_r}} k_i$, it is sufficient to prove that

$$k' \leq h \cdot \prod_{\text{bdg}(\mathcal{G}_{\mathcal{C}_u}) \cup E_{\mathcal{C}_u}} k_i \cdot \prod_{W_i \in \text{root}(\mathcal{G}_{\mathcal{C}_u})} w_{i,cb}^+$$

Fig. 12. Two \mathcal{G} -graphs in Lemma 7.

From the arguments used in the proof of Theorem 1 it can be easily derived, by a repeated application of rules R2 and R4, that

$$\oplus(\text{lhs}(\text{bdg}(\mathcal{G}_{\mathcal{G}_1}))) \xrightarrow{h} \oplus(N_{\mathcal{G}_r})$$

thus, by letting $W = \oplus(\text{lhs}(\text{bdg}(\mathcal{G}_{\mathcal{G}_1})))$, $R = \oplus(N_{\mathcal{G}_r})$, and applying rules R1 and R3, we obtain:

$$V_1 \oplus W \xrightarrow{h} R$$

Due to property (1), we may write $V_1 \oplus (W \ominus V_1) \xrightarrow{h} R$. From here, rule R5 leads to

$$V_1 \xrightarrow{h \cdot \llbracket (W \ominus V_1) \oplus (V_2 \ominus R) \rrbracket^+} R \oplus (V_2 \ominus R)$$

Again for property (1), this is equivalent to

$$\begin{aligned} & V_1 \xrightarrow{h \cdot \llbracket (W \ominus V_1) \oplus (V_2 \ominus R) \rrbracket^+} R \oplus V_2 \\ \text{(R3)} \Rightarrow & V_1 \xrightarrow{h \cdot \llbracket (W \ominus V_1) \oplus (V_2 \ominus R) \rrbracket^+} V_2 \end{aligned}$$

Since $V_1 \xrightarrow{k'} V_2$ and the input is sound and minimal by hypothesis, it is necessarily

$$k' \leq h \cdot \llbracket (W \ominus V_1) \oplus (V_2 \ominus R) \rrbracket^+$$

Now, letting $U = \oplus(N_{\mathcal{G}_u} \cup \text{rhs}(\text{bdg}(\mathcal{G}_{\mathcal{G}_u})))$ and $Y = V_2 \ominus R$ (see Fig. 13), we will prove that $(W \ominus V_1) \oplus Y \preceq U$. By definition of the \ominus operator it is $Y \oplus (V_2 \otimes R) = V_2$ and for each Y' such that $Y' \oplus (V_2 \otimes R) = V_2$ it is $Y \preceq Y'$. Let $Y' = V_2 \otimes U$; using the distributive and the absorption properties, and considering that $V_2 \preceq \oplus(N_{\mathcal{G}_2}) = \oplus(N_{\mathcal{G}_r} \cup N_{\mathcal{G}_u}) = R \oplus U$ by hypothesis, we get

$$\begin{aligned} Y' \oplus (V_2 \otimes R) &= (V_2 \otimes U) \oplus (V_2 \otimes R) = ((V_2 \otimes U) \oplus V_2) \otimes ((V_2 \otimes U) \oplus R) \\ &= V_2 \otimes (R \oplus V_2) \otimes (R \oplus U) = V_2 \end{aligned}$$

Hence, $Y \preceq V_2 \otimes U \preceq U$. On the other hand, $W \ominus V_1 \preceq U$ by hypothesis. Thus:

$$k' \leq h \cdot \llbracket (W \ominus V_1) \oplus Y \rrbracket^+ \leq h \cdot \llbracket U \rrbracket^+$$

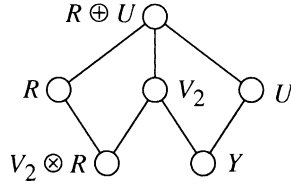


Fig. 13. Roll-up relationships of views in the Proof of Lemma 7.

Finally, because of Theorem 1 it is

$$\llbracket U \rrbracket^+ = \llbracket \oplus (N_{\mathcal{C}_u} \cup \text{rhs}(\text{bdg}(\mathcal{G}_{\mathcal{C}_u}))) \rrbracket^+ \leq \prod_{\text{bdg}(\mathcal{G}_{\mathcal{C}_u}) \cup E_{\mathcal{C}_u}} k_i \cdot \prod_{W_i \in \text{root}(\mathcal{G}_{\mathcal{C}_u})} w_{i,\text{cb}}^+$$

which proves the lemma. \square

Proof of Theorem 3. We will first examine the case $\text{bdg}(\mathcal{G}_{\mathcal{C}_1}) = \emptyset$. In this case it is:

$$\begin{aligned} u_{\text{cb}}(\mathcal{C}) &= u_{\text{cb}}(\mathcal{C}_1) \cdot u_{\text{cb}}(\mathcal{C}_2) \\ u_{\text{cb}}(\mathcal{C}'_1) &\leq u_{\text{cb}}(\mathcal{C}_1) \quad (\text{since } \mathcal{C}'_1 \sqsubseteq \mathcal{C}_1) \\ v_{\text{cb}}^+ &\leq u_{\text{cb}}(\mathcal{C}_2) \quad (\text{for Lemma 6}) \end{aligned}$$

If in $\mathcal{G}_{\mathcal{C}'}$ there is an arc, labeled k' , entering V (with $1 \leq k' \leq v^+$ since the input is sound) it is

$$u_{\text{cb}}(\mathcal{C}') = u_{\text{cb}}(\mathcal{C}'_1) \cdot k' \leq u_{\text{cb}}(\mathcal{C}_1) \cdot v_{\text{cb}}^+ \leq u_{\text{cb}}(\mathcal{C}_1) \cdot u_{\text{cb}}(\mathcal{C}_2) = u_{\text{cb}}(\mathcal{C})$$

Otherwise, it is $u_{\text{cb}}(\mathcal{C}') = u_{\text{cb}}(\mathcal{C}_1) \cdot v_{\text{cb}}^+ \leq u_{\text{cb}}(\mathcal{C})$. In both cases, we have $\mathcal{C}' \sqsubseteq \mathcal{C}$.

Now, let $\text{bdg}(\mathcal{G}_{\mathcal{C}_1}) \neq \emptyset$; we have

$$\begin{aligned} u_{\text{cb}}(\mathcal{C}') &= u_{\text{cb}}(\mathcal{C}'_1) \cdot k' \\ u_{\text{cb}}(\mathcal{C}) &= \prod_{E_{\mathcal{C}}} k_i \cdot \prod_{W_i \in \text{root}(\mathcal{G}_{\mathcal{C}})} w_{i,\text{cb}}^+ \\ &= u_{\text{cb}}(\mathcal{C}_1) \cdot \prod_{\text{bdg}(\mathcal{G}_{\mathcal{C}_1}) \cup E_{\mathcal{C}_r}} k_i \cdot \prod_{\text{bdg}(\mathcal{G}_{\mathcal{C}_u}) \cup E_{\mathcal{C}_u}} k_i \cdot \prod_{W_i \in \text{root}(\mathcal{G}_{\mathcal{C}_u})} w_{i,\text{cb}}^+ \end{aligned}$$

where the last equality follows from the proof of Lemma 7. Since $\mathcal{C}'_1 \sqsubseteq \mathcal{C}_1$ it is $u_{\text{cb}}(\mathcal{C}'_1) \leq u_{\text{cb}}(\mathcal{C}_1)$; for Lemma 7 it also is $k' \leq \prod_{\text{bdg}(\mathcal{G}_{\mathcal{C}_1}) \cup E_{\mathcal{C}_r}} k_i \cdot \prod_{\text{bdg}(\mathcal{G}_{\mathcal{C}_u}) \cup E_{\mathcal{C}_u}} k_i \cdot \prod_{W_i \in \text{root}(\mathcal{G}_{\mathcal{C}_u})} w_{i,\text{cb}}^+$. Thus, the theorem is proved. \square

References

- [1] G. Birkhoff, Lattice Theory, American Mathematical Society, 1995.
- [2] A.F. Cardenas, Analysis and performance of inverted database structures, Comm. of the ACM 18 (5) (1975) 253–263.
- [3] P. Ciaccia, M. Golfarelli, S. Rizzi, On Estimating the Cardinality of Aggregate Views, in: Proceedings of the DMDW'01, Interlaken, Switzerland, 2001, pp. 12.1–12.10.
- [4] P. Ciaccia, M. Golfarelli, S. Rizzi, Using Domain-Derived Constraints to Bound the Cardinality of Aggregate Views, in: Proceedings of the 9th SEBD, Venice, Italy, 2001, pp. 249–256.

- [5] P. Ciaccia, D. Maio, On the complexity of finding bounds for projection cardinalities in relational databases, *Information Systems* 17 (6) (1992) 511–515.
- [6] P. Ciaccia, D. Maio, Domains and active domains: What this distinction implies for the estimation of projection sizes in relational databases, *IEEE Trans. on Knowledge and Data Engineering* 7 (4) (1995) 641–655.
- [7] M. Golfarelli, S. Rizzi, View Materialization for Nested GPSJ Queries, in: *Proceedings of the DMDW'00*, Stockholm, Sweden, 2000.
- [8] J. Grant, J. Minker, Numerical Dependencies, in: H. Gallaire, J. Minker, J.-M. Nicolas (Eds.), *Advances in Database Theory*, vol. II, Plenum Publ. Co., 1983.
- [9] H. Gupta, Selection of Views to Materialize in a Data Warehouse, in: *Proceedings of the ICDDT'97*, Delphi, Greece, 1997, pp. 98–112.
- [10] H. Gupta, V. Harinarayan, A. Rajaraman, J. Ullman, Index Selection for OLAP, in: *Proceedings of the ICDE'97*, Birmingham, UK, 1997, pp. 208–219.
- [11] M. Gyssens, L.V.S. Lakshmanan, A Foundation for Multi-Dimensional Databases, in: *Proceedings of the 23rd VLDB*, Athens, Greece, 1997, pp. 106–115.
- [12] V. Harinarayan, A. Rajaraman, J. Ullman, Implementing Data Cubes Efficiently, in: *Proceedings of the ACM Sigmod Conference*, Montreal, Canada, 1996, pp. 205–216.
- [13] W. Hou, G. Özsoyoglu, Statistical estimators for aggregate relational algebra queries, *ACM Trans. on Database Systems* 16 (4) (1991) 600–654.
- [14] M.V. Mannino, P. Chu, T. Sager, Statistical profile estimation in database systems, *ACM Computing Surveys* 20 (3) (1988) 191–221.
- [15] M. Muralikrishna, D.J. DeWitt, Equi-depth Histograms for Estimating Selectivity Factors for Multi-Dimensional Queries, in: *Proceedings of the ACM Sigmod Conference*, Chicago, IL, 1988, pp. 28–36.
- [16] K. Ross, D. Srivastava, Fast Computation of Sparse Datacubes, in: *Proceedings of the 23rd VLDB*, Athens, Greece, 1997, pp. 116–125.
- [17] A. Shukla, P. Deshpande, J. Naughton, K. Ramasamy, Storage Estimation for Multidimensional Aggregates in the Presence of Hierarchies, in: *Proceedings of the 22nd VLDB*, Mumbai, India, 1996, pp. 522–531.
- [18] D. Theodoratos, M. Bouzeghoub, A General Framework for the View Selection Problem for Data Warehouse Design and Evolution, in *Proceedings of the DOLAP 2000*, Washington, DC, 2000, pp. 1–8.
- [19] P. Vassiliadis, Gulliver in the Land of Data Warehousing: Practical Experiences and Observations of a Researcher, in: *Proceedings of the DMDW'00*, Stockholm, Sweden, 2000, pp. 12/1–12/16.



Paolo Ciaccia has a “Laurea” degree in Electronic Engineering (1985) and a Ph.D. in Electronic and Computer Engineering (1992), both from the University of Bologna, Italy, where he is Full Professor of Information Systems. His research interests include similarity and preference-based query processing, data warehousing and data mining. He is one of the designers of the M-tree, an index for metric data used by many multimedia and data mining research groups in the world. He participated several international and national research projects, among which ESPRIT IV LTR HERMES on Multimedia Information Management Systems. Since 2001 he has been the responsible for DEIS within the IST/FET PANDA Network (Patterns for Next-Generation Database Systems). He has published more than 60 refereed papers in the areas of database systems, neural networks, software engineering, and autonomous systems, in major international journals (including *IEEE TKDE*, *IEEE TSE*, *ACM TODS*, *ACM TOIS*, *IS*, and *Biological Cybernetics*) and international conferences (including *VLDB*, *ACM-PODS*, *EDBT*, and *ICDE*). In 1999 he was programme co-chair of the 1st International Workshop on Similarity Search, and in 2002 he was PC chair of the SEBD Italian conference on advanced data bases. He is an Associate Editor of *IEEE TKDE*.



Matteo Golfarelli received his degree in Computer Science in 1995 and his Ph.D. for his work on autonomous agents in 1998 from the University of Bologna, Italy. In 2000 he joined the Computer Science Department as associate researcher, teaching Information Systems and Computer Architectures. He has published over 30 papers in refereed journals and international conferences in the fields of autonomous agents, pattern recognition and databases. He is currently involved in the EU PANDA thematic network concerning pattern-base management systems. His current research interests include all the aspects related to data warehouse design, in particular multidimensional modeling, view materialization and fragmentation, physical design.



Stefano Rizzi received his degree in Electronic Engineering in 1988 and his Ph.D. for his work on autonomous agents in 1996 from the University of Bologna, Italy. In 1995 he joined the Computer Science Department as a researcher. Since 1998 he is Associate Professor at the University of Bologna, where is the head of the Data Warehousing Laboratory and teaches advanced information systems and software engineering. He has published over 50 papers in refereed journals and international conferences in the fields of data warehousing, pattern recognition, mobile robotics, multi-agent systems, and visual query languages. He joined several national research projects on the above areas and is currently involved in the EU PANDA thematic network concerning pattern-base management systems. He served in the PC of several international conferences and as a reviewer in journals. His current research interests include all the aspects related to data warehouse design, in particular multidimensional modeling, view materialization and fragmentation, physical design.