

# pCube: Update-Efficient Online Aggregation with Progressive Feedback and Error Bounds

Mirek Riedewald, Divyakant Agrawal, and Amr El Abbadi  
Department of Computer Science \*  
University of California, Santa Barbara  
Santa Barbara, CA 93106  
{mirek, agrawal, amr}@cs.ucsb.edu

## Abstract

Multidimensional data cubes are used in large data warehouses as a tool for online aggregation of information. As the number of dimensions increases, supporting efficient queries as well as updates to the data cube becomes difficult. Another problem that arises with increased dimensionality is the sparseness of the data space. In this paper we develop a new data structure referred to as the pCube (data cube for progressive querying), to support efficient querying and updating of multidimensional data cubes in large data warehouses. While the pCube concept is very general and can be applied to any type of query, we mainly focus on range queries that summarize the contents of regions of the data cube. pCube provides intermediate results with absolute error bounds (to allow trading accuracy for fast response time), efficient updates, scalability with increasing dimensionality, and pre-aggregation to support summarization of large ranges. We present both a general solution and an implementation of pCube and report the results of experimental evaluations.

## 1. Introduction

Data warehouses typically accumulate and store operational data from a collection of OLTP (On-Line Transaction Processing) databases. Statistical and scientific databases (e.g., geographical information systems, census databases) contain huge amounts of “raw” information. Human analysts are interested in extracting that part of the data that is of interest for a certain task, e.g., discovering population development trends or analyzing the success of a campaign of TV commercials. That data is described by a set of attributes. One of them is chosen as the *measure* attribute (the

attribute of interest, e.g., number of people). The others are the functional attributes, or simply *dimensions* (e.g., age, income). The data can conceptually be viewed as a hyper-rectangle, where each single *cell* is described by a unique combination of dimension values and contains the corresponding value of the measure attribute. We will henceforth refer to this hyper-rectangle as the *data cube*. We should point out that in this paper the data cube model only represents a way of *viewing* the data. It does not imply or exclude a certain physical data organization, like in relational tables (ROLAP) or in multidimensional databases (MOLAP).

The more complex dependencies an analyst wants to discover, the more functional attributes have to be included into the analysis. This obviously leads to a high dimensionality of the data cube. It is, however, important to note that even though the data itself is massive, the high dimensionality causes the data cube to be very *sparse*. For example, a U.S. census data cube with the dimensions shown in Table 1 contains more than  $2 * 10^{14}$  cells. Assuming that the U.S. census bureau stores data about 200 million Americans, at most 0.0001% of the cells will be non-empty. Assuming further, that each value of the measure attribute (in the example the number of people) can be encoded with a single byte, storing all cells will require more than 180 Terabytes.

Attribute	Domain	Attribute	Domain
Age	0 - 150	Income	0 - 99999
Weight	1 - 500	Race	1 - 5
Family type	1 - 5	Marital status	1 - 7
Class of worker	0 - 8	Education	1 - 17

**Table 1. Attributes and their domains for a census data set**

An important tool for analyzing data is a range query. Range queries apply aggregate operators (like the SQL op-

\*This research was supported by NSF grants EIA-9818320, IIS-98-17432, and IIS-99-70700.

erators SUM, COUNT, MAX) to data cube cells that are selected by the given ranges on the dimensions. A range query could ask to “Return the number of U.S. citizens that are younger than 30 years and earn more than \$50,000 per year”. In the data cube model a range query selects a hyper-rectangle of cells. We will use the term *query-cube* to refer to this hyper-rectangle. To answer a range query exactly, the values stored in the cells inside the query-cube must be known. For large ranges this can lead to a situation where answering the range query requires a considerable number of accesses to slow secondary storage. Online Analytical Processing (OLAP) queries often are more complex than “simple” range queries and can contain subqueries. Nevertheless OLAP is an interactive process where analysts expect fast responses in the order of seconds at most. When *exact* results are not required, e.g., for exploratory data analysis when it suffices to get an idea of the order of magnitude of some phenomenon in the database, queries can be sped up by returning an *approximate* answer based on summary information about the contents of the database.

OLAP applications typically have to deal with updates, e.g., data warehouses collect constantly changing data from a company's databases. Often it is regarded to be sufficient if those updates can be efficiently processed in batches. There are, however, applications that require *fast updates* to single cells. For instance, business leaders will want to construct interactive *what-if* scenarios using their data cubes, in much the same way they construct *what-if* scenarios using spreadsheets now. These applications require real-time (or even hypothetical) data to be integrated with historic data for the purpose of instantaneous analysis and subsequent action. In a system where updates are only propagated in batches (possibly during the night), such *what-if* analysis will not be possible. Also, while some applications do not suffer in the presence of stale data, in many emerging applications like decision support and stock trading the instant availability of the latest information plays a crucial role. Finally, the greater the number of updates, the longer batch updates will make the data in the data cube inaccessible to an analyst. Finding suitable time slots for batch update windows becomes increasingly harder when businesses demand flexible work hours and 24 hour availability of their important data. Especially data collections that are accessible from all over the world (e.g., for multinational companies) do not follow a simple “high load at daytime, no accesses at nighttime” pattern. Consequently even systems that apply updates in large batches benefit from reduced update costs which shrink the size of the update window and enable the system to process more updates in the same time.

In this paper we develop the concept of a *progressive data cube*, abbreviated as *pCube*. The main idea is to build a data structure that is able to provide the user with fast approximate answers for aggregation queries. The answer

can then be progressively refined. Since an approximate answer without quality guarantees is of no use, pCube provides *absolute* error bounds for the query result. Compared to probabilistic error bounds (confidence intervals) absolute bounds offer the following advantages. First, probabilistic bounds do not guarantee that the final result will be inside the confidence interval; especially for small answer sets with high variance and for queries that look for extremal values. Second, while non-expert users can easily understand the concept of absolute error bounds, this is not necessarily the case with confidence intervals. This argument is very important for data exploration tools that are designed for a diverse group of users.

To be more specific, pCube has the following properties.

- pCube can deal with high-dimensional sparse data cubes as well as with dense data cubes. Its access, update and resource costs depend on the number of non-empty cells, not on the size of the complete data cube (the product of the domain sizes of all dimensions).
- pCube's feedback mechanism provides the user with quick responses to queries. A response consists of an approximate answer with absolute error bounds which are progressively refined.
- The feedback mechanism enables the user
  - to interrupt the query execution as soon as the result is good enough, or
  - to set a time limit for the query execution and get the best possible result according to this limit (with the possibility to further refine the result by allowing additional execution time).
- pCube combines pre-aggregation and hierarchical data organization to speed up query execution and enable fast updates. It does not depend on batch updates. On the other hand, pCube benefits from applying updates in batches. The batch size therefore can be used as a tuning parameter.

For popular aggregate functions like SUM (sum over all selected non-empty cells), COUNT (number of data items that satisfy the selection condition) and others, we state explicitly which values have to be stored to enable an update-efficient approximation.

The rest of the paper is organized as follows. Section 2 gives an overview of relevant related work. In Section 3 we develop the general pCube concept. Experimental results for a pCube implementation are reported in Section 4. Section 5 concludes this article.

## 2. Previous Work

To support the process of extracting and analyzing data from a data collection, [13] proposes a new SQL operator — “CUBE” or “data cube”. It is a  $d$ -dimensional generalization of the GROUP BY operator and computes all possible group-bys for  $d$  attributes (given in the SELECT clause). We will refer to the result of the CUBE operator as the *extended data cube*. For a selection of  $d$  attributes, there exist  $2^d$  possible group-bys. Due to this exponential number of results, computing the complete extended data cube efficiently is one of the major challenges (see [5] for an overview and references). Especially for high-dimensional and sparse data sets computing and storing the complete extended data cube is infeasible ([23] describes the problem of database explosion). Updates to a single item in one of the source data sets would cause an aggregate in each of the  $2^d$  views to be updated as well. To overcome the above problems, efficient algorithms were proposed to choose those views (and indexes) to be stored that optimize a certain cost/benefit metric (e.g., [14, 15]). [19] and [20] describe index structures especially designed for the extended data cube. The above techniques do not provide fast responses for all aggregation queries.

In [18], Ho et al. develop techniques to guarantee an access time for range SUM queries that is independent of the size of the query-cube (constant cost) by pre-aggregating information in the data cube. Their prefix-sum approach suffers from high update costs (single cell update cost is exponential in the domain sizes). [6], [9] and [10] address this issue by balancing the cost of queries and updates. None of the above approaches can, however, efficiently deal with high-dimensional sparse data sets. In fact, applying the proposed algorithms to sparse data sets leads to dense sets of the size of the complete data cube, including empty cells. Suggested techniques for dealing with sparse data sets are rather sketchy and not effective for arbitrary data distributions.

The main idea behind most approaches for handling sparseness is to find dense chunks or clusters in the data set. Then dense and sparse regions of the data set are each handled by a specialized technique. Usually these approaches try to combine the advantages of special “dense” and “sparse” storage and access techniques, e.g., [7, 12]. Techniques for sparse data sets mainly focus on space efficiency. Aspects like early feedback for queries and update-efficiency are not examined.

An established technique to support fast accesses to databases is to create index structures on them. [8] and [22] provide an excellent overview of the field. Index structures that only index non-empty cells (i.e., the data items that are in the database) can deal with sparseness. On the other hand just providing fast access to all selected data items does not

suffice. Retrieving and aggregating each selected item on-the-fly is still too slow for large data collections. In [2] it is suggested to augment indexes by aggregate data to exploit summary information for a data set. This idea up to now was used to support query cost estimation and approximate answers to queries [29], but to the best of our knowledge never for *progressive feedback* with *absolute error bounds* on data cubes. Also, in contrast to our approach, the solution presented in [29] does not consider efficient updates. Another feature that distinguishes pCube from index structures is that it is primarily a way of *organizing and storing* data for aggregation and summarization. The indexing comes for “free” by exploiting the data organization.

In an early paper Rowe [25] presents an approach to approximate query results on databases which is called “antismpling”. The main idea is to maintain statistics about the contents of a database (a database abstract), e.g., mean, median, number of data items, and other distribution information. Instead of returning the (expensive) exact result to a query, an approximate answer is returned based on the statistical information. Antismpling techniques provide *absolute* error bounds. For instance, the number of non-empty cells in the data cube is an upper bound for the number of non-empty cells in any region of the data cube. Combining different statistics, each possibly providing loose bounds only, can lead to much tighter bounds. It requires expert knowledge to determine how much and which information to store and how to combine it. Guaranteeing tight bounds for “simple” queries that count the number of non-empty cells in an arbitrary range of the data cube is a hard task. How detailed should the information about the data distribution be and how should it be organized for fast accesses and updates? We will show how pCube answers these questions.

The New Jersey Data Reduction Report [2] provides an overview of various techniques that can be useful for online aggregation. Barbará and Wu [4] describe a technique to use loglinear models for compressing the data cube and obtaining approximate answers. For the technique to be efficient, dense clusters in the data cube have to be identified and are then approximated by loglinear models (other approximation techniques could be used as well, e.g., regression models [3]). The approach is mainly applicable to dense low-dimensional data cubes. It is possible to return absolute error bounds based on the approximation model and a guaranteed bound for the approximation error per cell. However, for queries with large query-cubes this approach leads to slow responses. Maintaining the approximation models in the presence of updates requires batch updates that access all cells in a chunk. [26] presents a compression technique that achieves very high compression ratios (1000 and more) by using density estimations for dense clusters. This approach, however, is not efficient in the presence of high-

dimensional sparse data sets and frequent updates. With increasing sparseness, the estimation will become much less accurate. Higher dimensionality reduces the accuracy of the query result which is obtained by multiple integrations over the density functions. Also, the technique can not provide absolute error bounds for its approximate answers and a query has to access the data sources if the user requests an exact result. Compressing sparse data cubes using wavelets is described in [28]. The more wavelet coefficients are accessed, the more accurate the result. Theoretically one could select the optimal time-accuracy tradeoff for a query. But since no error bounds are returned to the user, he/she can not make use of this possibility. To get an exact answer, the data sources have to be accessed, because the wavelet decomposition algorithm drops “unimportant” coefficients. Even though wavelets provide good approximations for most practical applications, they do not fulfill our requirements of error bounds for approximations and of efficient updates. Another direction of research aims at using histograms for fast approximate answers [24]. Histograms are an efficient technique to summarize data. Systems like AQUA [1] use small pre-computed statistics, called synopses [11] (which usually are histograms and samples), to provide fast approximate answers. Synopses can be used to answer queries. The AQUA system, for instance, returns fast answers with probabilistic error bounds [1] (confidence intervals, no absolute bounds). Up to a certain degree the user can theoretically trade off accuracy for faster responses by allowing smaller samples to be used. The accuracy of the query response is limited by the information stored in the histograms and samples. For further refinements (up to the exact answer) the query has to access the data in the sources, which is very expensive. Techniques to efficiently maintain synopses typically require batch updates, which makes them inaccessible during the update. pCube on the other hand can easily deal with concurrent single cell updates and queries by using locks on the tree nodes.

Recent developments show that even though the issue of fast approximate answering has gained increasing attention, there are still not many solutions that provide progressive refinement of the approximate result with error bounds. An exception is the CONTROL project [16]. Users are provided with *probabilistic* error bounds (confidence intervals) that are progressively refined. CONTROL is based on algorithms especially designed for online aggregation (e.g., [17]). While confidence intervals give an idea about the expected result for a query, they do not guarantee that the final result will be inside the interval. This holds especially for skewed data and for queries that search for extremum statistics (outliers). In contrast to the CONTROL approach, pCube returns *absolute* error bounds.

### 3. The pCube Concept

For large data collections accessing and aggregating the data that is relevant for a query can consume time in the order of minutes. Take for example a range query that asks for the number of Americans that are younger than 30 years. Obtaining the answer from the “raw” data requires accessing at least all those data items (i.e., data cube cells) that satisfy the selection criterion (which could be some millions of items), not counting additional index accesses. In general, for an *exact* answer the value of each single cell in the query-cube must be found out and included in the aggregation on-the-fly. Index structures are of little help when a query-cube contains a large number of non-empty cells. The alternative way to speed up the query process is by pre-computing the aggregation values for groups of cells and to use those results when answering a query. For instance, if the number of Americans between 0 and 14 and the number of Americans between 15 and 29 years of age is known, i.e., was pre-computed and stored, then two accesses to these pre-computed values provide the answer, thus saving a large number of disk accesses and CPU time. Pre-computation (pre-aggregation) leads to the following trade-offs. The smaller the targeted response time for a query and the more queries are to be sped up, the more storage is necessary for pre-computed values. On the other hand, the more pre-computed values depend on a certain cell, the more expensive it is to update this cell's value (since the pre-computed values have to be updated as well). High dimensionality and sparseness of the data cube affect the benefit-cost ratio of these tradeoffs negatively. More dimensions dramatically increase the number of possible range queries and storing pre-computed values for sparse data cubes leads to the problem of database explosion [23]. Approaches which are effective for low-dimensional and dense data cubes therefore are not feasible for high-dimensional sparse data cubes.

#### 3.1. Requirements

The only way to guarantee fast responses for any query on high-dimensional sparse and massive data cubes is to weaken the requirement to return exact answers. Instead, pCube returns *approximate* answers with absolute error bounds. To get a quick response, the content of the query-cube has to be approximated with a small number of parameters. The difficult task is to choose which parameters to store in order to get tight error bounds (see for instance [25] for an overview). Different queries and applications make different demands on the approximation technique.

(1) When the query-cube is large and contains a large number of non-empty cells, it can only be approximated quickly if the approximation model's granularity is coarse.

Queries with small query-cubes on the other hand are best approximated by choosing a finer approximation model. Popular approximation techniques like (multidimensional) histograms only offer a single granularity to all queries, i.e., favor either large or small query-cubes.

(2) Applications require different degrees of accuracy. The degree of accuracy should be chosen by the user or application, not by the system (designer). In some cases the required accuracy is not determined *a priori* (e.g., when it suffices to find out if the result is a positive number, a result range  $100 \leq x \leq 10000$  is acceptable, but the range  $-1 \leq x \leq 10$  is not).

(3) The requirement of efficient updates rules out techniques where an update to a single cell can cause expensive recomputations (e.g., with costs in the order of the data cube size or the number of non-empty cells in the data cube).

### 3.2. The pCube Structure

In the following a pCube structure that fulfills the above requirements is described. pCube is organized as a tree. Each node covers a region of the data cube. We will use the term *node-region* to denote this region (the node-region could be a hyper-rectangle, hyper-sphere, etc.). The root node covers the complete data cube (or any region that contains all non-empty cells in the data cube). Leaf nodes cover single *non-empty* cells. Any other node covers a part of the node-region of its parent, such that all children together completely cover their parent's node-region (or at least all non-empty cells covered by the parent).<sup>1</sup> Each leaf node stores the value of the cell it covers. Internal nodes store information about the cells in their node-region, which will henceforth be referred to as *node-information*. For each aggregate function supported by pCube, the node-information contains the pre-computed aggregate value for the cells in the node-region or values that enable the computation of the aggregate value in constant time. Additional parameters are stored that allow for the fast computation of absolute error bounds for each of the supported aggregate functions. To guarantee space and update efficiency, the size of the node-information is bounded by a (*a priori* known) constant.

Figure 1 shows an instance of pCube that supports the aggregate function COUNT. For a combination of dimension values the measure attribute stores the number of data items with this value combination. COUNT range queries return the number of data items that fall into a certain range. The node-information simply consists of the sum of the values of the non-empty cells in the node-region. Since there are overlapping node-regions of sibling nodes, it has to be ensured that no data item is counted more than once. Non-

<sup>1</sup>Note that there might be non-empty intersections, e.g., overlapping hyper-rectangles for structures similar to the R-tree.

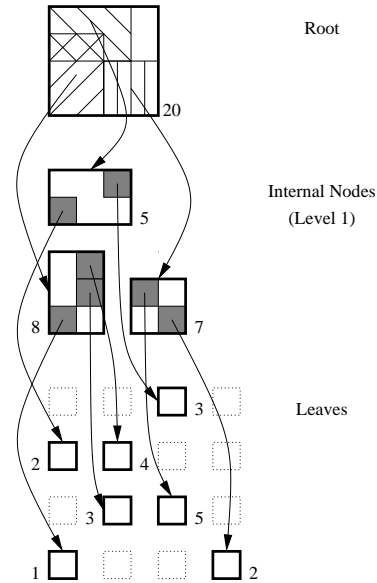


Figure 1. pCube with R-tree-like structure for the aggregate function COUNT

empty cells that fall into more than one node-region of sibling nodes have to be assigned to only one of the nodes.

### 3.3. Querying pCube

A query is answered by descending the tree. At each node the query result is approximated based on the node-information (e.g., as described in [25]). The further the query descends, the smaller the node-regions, i.e., the better the approximation based on the node-information. Each time the query descends to a new node, a finer intermediate result (with better error bounds) can be output. Thus the user receives continuous feedback of improving quality. Both small and large query-cubes can be supported efficiently. For large query-cubes the coarse granularity in the beginning allows for fast results which quickly improve. Small query-cubes typically intersect a smaller number of tree nodes, and therefore the query quickly reaches lower level nodes that provide more accurate answers. The pre-computed aggregate value at a tree node can dramatically reduce the response time. If a query-cube completely contains a node-region, this value represents the *exact* aggregate value for all query-cube cells in the node-region. Consequently, no descendents of the node need to be accessed and the aggregate function does not need to be computed on-the-fly. In our census example for the aggregate function COUNT, a tree node whose node-region contains all cells that satisfy  $10 \leq \text{age} < 30$  and  $0 \leq \text{income} \leq 30000$  (all other attributes range over their complete domain) stores

the number of people that fall into this range in its node-information. If, for instance, a query asks for the number of people that are younger than 30 and earn less than \$40,000, the pre-computed value in the node contributes the exact answer for all people between 10 and 30 that earn less than \$30,000. Thus, a single access to the node provides an exact answer for a large number of cells, without requiring to descend further in that subtree.

The cost of a query (measured as the number of disk/memory accesses, or response time) in the worst case depends linearly on the number of node-regions the query-cube intersects. This cost decreases if node-regions are completely included in the query-cube and if the user is satisfied with approximate results.

### 3.4. Obtaining Absolute Error Bounds

The computation of the error bounds is based on the following observation. When a query-cube and a node-region have no cells in common, the answer is empty (which is an exact answer). If the query-cube contains the complete node-region, the aggregate value in the node already provides the exact answer for those cells (no approximation error at all). In the non-trivial case when the query-cube contains only some of the node-region's cells, error bounds are obtained in constant time by only using the information in the nodes visited so far and properties of the aggregate function. A node can either store bounds on the extreme values or store enough information to allow their fast computation. For the COUNT function a node only needs to store the *sum* of the values of the non-empty cells in its node-region. Storing the number of non-empty cells would not lead to a correct result, because there might be more than one data item in the database with a certain combination of dimension values (i.e., a cell can store a value greater than 1). The stored sum provides the aggregate over all cells in the node-region and the upper bound. The value 0 is the trivial lower bound. Table 2 contains selected aggregate functions (including the SQL standard aggregate functions) together with possible auxiliary values (explained in Table 3) for the computation of absolute error bounds. The last column contains the formula to compute the aggregate value for the non-empty cells in the node region.

Note that AVG, Variance and Standard deviation are not approximated by combining the aggregate values of different nodes directly. We illustrate this for the aggregate function AVG. If a query contains two pCube nodes such that for one node an average value  $v_1$  and for the other a value  $v_2$  is computed, the overall average value is still not determined. It can be any value between  $v_1$  and  $v_2$ , depending on the number of base values that contributed to each of the two partial results. A simple solution to this problem is to compute the auxiliary values for the nodes, i.e., the

sum and the count value that produce the corresponding average, and to combine those in order to obtain the query result. Thus, instead of combining the aggregate values directly, the corresponding auxiliary values are computed for the pCube nodes and then combined to get the value for the aggregate function.

Figure 2 illustrates the technique for the aggregate function COUNT for a non-overlapping tree structure. The shaded region indicates the intersection between node-region and query-cube for all nodes that have to be accessed to get the exact result. Note that the node-region of one of the internal nodes is completely included in the query-cube, therefore the corresponding leaf nodes do not need to be accessed, even though their node-regions are covered by the query-cube. Each node stores the sum of the values of the non-empty cells it covers. In case the query-cube completely contains the node-region, upper and lower bounds for the result are equal to the count value stored in the node. Otherwise, when query-cube and node-region intersect, at least none and at most all non-empty cells of the node-region can fall into the intersection, providing the error bounds. Even though the approximation is very coarse at the upper nodes near the root, it is refined quickly when the query descends the tree. In the example, the approximate value is obtained based on the assumption that the non-empty cells are uniformly distributed over the node-region. Thus the query result is approximated as 9.4 at the root (the query contains 6 out of 16 cells, the aggregate value for the 16 cells is 25) and as 10.5 at the next level of nodes. Finally, after accessing the values in the leaves the exact result of 12 is returned.

### 3.5. Updating pCube

We define an *update* to be any operation that changes the value of a data cube cell (e.g., an insert is an update that changes a cell's value from "empty" to a certain value; a delete is defined similarly). An update to a cell affects the node-information of all nodes that cover this cell. It takes advantage of the tree structure by first descending the tree until the corresponding leaf node is reached. The leaf's value is changed. Then, knowing the former and the new value of the updated cell, all affected nodes on the path from the root to the leaf are updated (bottom-up). The overall update cost depends not only on the number of nodes that are affected, but also on the cost of updating a node's information about the cells in its node-region. Ideally the node-information can be updated in constant time without additional disk accesses. Then, for trees with non-overlapping node-regions of sibling nodes, the update cost is linear in the height of the tree. Otherwise, in the worst case all internal tree nodes might have to be accessed. Also, updates can lead to node/page overflow or underflow. How these cases

Function	Auxiliary values	Upper	Lower	Aggregate for node-region
COUNT	sum	sum	0	sum
SUM	sum <sub>+</sub> , sum <sub>-</sub>	sum <sub>+</sub>	sum <sub>-</sub>	sum <sub>+</sub> + sum <sub>-</sub>
AVG (average)	cnt, sum <sub>+</sub> , sum <sub>-</sub>	sum <sub>+</sub>	sum <sub>-</sub>	(sum <sub>+</sub> + sum <sub>-</sub> )/cnt
AVG, alternative implementation	cnt, sum, max, min	max	min	sum/cnt
MAX	max, min	max	min	max
MIN	max, min	max	min	min
Variance	cnt, sum, sum <sub>2</sub>	sum <sub>2</sub>	0	sum <sub>2</sub> /cnt - (sum/cnt) <sup>2</sup>
Standard deviation	cnt, sum, sum <sub>2</sub>	$\sqrt{\text{sum}_2}$	0	$\sqrt{\text{sum}_2/\text{cnt} - (\text{sum}/\text{cnt})^2}$

**Table 2. Values that provide absolute error bounds for aggregate functions**

Symbols	Meaning
sum	Sum of the values of all non-empty cells in the node-region
sum <sub>2</sub>	Sum of the squares of the values of all non-empty cells in the node-region
sum <sub>+</sub>	Sum of the positive values in the node-region
sum <sub>-</sub>	Sum of the negative values in the node-region
cnt	Number of non-empty cells in the node-region
max	Maximum value in the node-region
min	Minimum value in the node-region

**Table 3. Node-information values and their meaning**

are handled depends on the specific tree structure used to implement the pCube. Basically techniques similar to those known from index structures can be applied. Figure 3 illustrates the insertion of a new value into a pCube; deleting the cell would reverse the changes.

### 3.6. Applicability of pCube

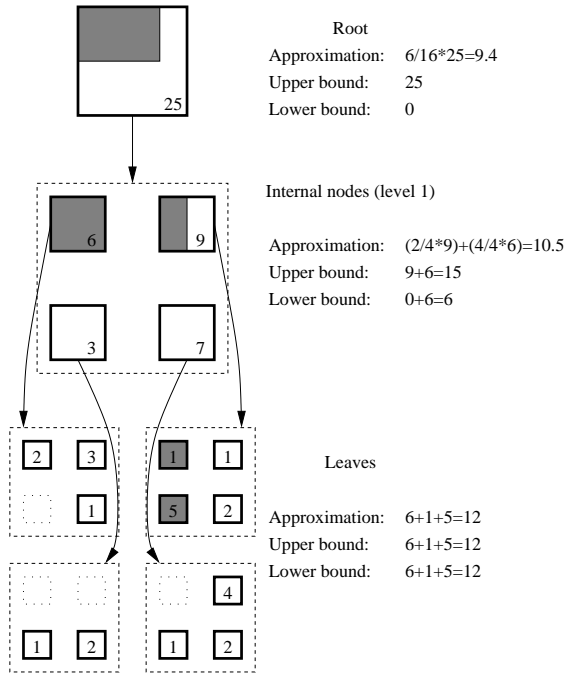
The way pCube's nodes divide the space determines pCube's structure and therefore heavily influences the query and update costs. In general, techniques have to be chosen that guarantee that both costs depend on the number of *non-empty* cells, rather than on the dimensionality or the size of the data cube (which includes empty cells). Possible choices include structures that are similar to the R-tree, X-tree, kdB-tree, etc. (see [8] for references). In addition to the structure, it must be determined which kind of information to store at the nodes for supporting efficient aggregation and good approximation.

Generally a pCube can be constructed for any aggregate function. There are, however, functions that can not be supported efficiently. The main problem is caused by nodes in the upper tree levels (close to the root) that aggregate a large number of values. We require the cost of recomputing the aggregate to be independent of the number of aggregated values, i.e., to be bounded by a constant. It should be possible to compute the new aggregate value only from the knowledge about the update operation (previous contents of the updated data cube cell, new contents of that cell) and a small amount of auxiliary information about the values. Note that any set of aggregate functions that is *self-maintainable* with respect to insertions and updates can be

efficiently implemented by pCube. Such a set of aggregate functions is self-maintainable if the new values of the functions can be computed solely from the old values and from the changes caused by the update operation [21].

The SQL aggregate functions COUNT and SUM are self-maintainable. COUNT can be supported efficiently by storing the number of data items that are in the node-region; SUM is handled similarly. For AVG this direct approach fails. Knowing the former average for the cells in the node-region and the update on a single cell does not provide enough information for computing the new average. However, if instead the sum and the number of values (i.e., the COUNT) are known, the average can easily be computed and incrementally updated. Consequently, even though AVG is not self-maintainable, it can be made self-maintainable by adding a small amount of additional information.

There are other aggregate functions like MAX, MIN and Median, that can not be made self-maintainable. It is not possible to efficiently maintain them based only on the information stored at single nodes, e.g., for MAX, when the current maximum value is deleted. pCube's hierarchical structure provides a powerful tool that solves this problem for some aggregate functions. We describe a technique for the function MAX. Since updates have to descend the tree, one can take advantage of the combined information of the visited nodes. Let each node store the maximum of the values in its node-region and assume that an update deletes the current maximum of a node. Then the new maximum can be computed by a function that returns the maximum over the maximum values of all child nodes. Since the update



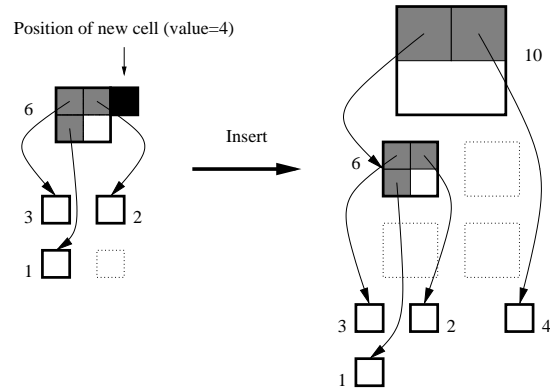
**Figure 2. Obtaining error bounds for a pCube with quadtree-like structure**

will affect the maximum value of at least one child node, this function proceeds recursively. A similar technique can be applied to MIN. This way both functions are maintainable at a cost which is at most  $b$  times higher than maintaining COUNT,  $b$  being the largest number of children of a node. Similar “workarounds” as for MAX and MIN are possible for other aggregate functions which are not self-maintainable. Details are not provided here and are part of our future work.

Table 2 contains some popular aggregate functions and formulas for obtaining absolute error bounds and aggregate values from auxiliary values that can be efficiently incrementally maintained (see Table 3). Note that the auxiliary values are chosen with update-efficiency being the priority. By storing more information, the error bounds can be tightened. Essentially the choice of which kind of information to store at a tree node depends on the application (which types of queries are supported, how sensitive the query result is to approximation errors).

## 4. Experimental Results

To evaluate the pCube concept we performed extensive experimental evaluations. We start with a short description of a pCube implementation for aggregation over hyper-rectangular ranges of the data cube. Then experimental re-



**Figure 3. Effects of an update on a pCube with quadtree-like structure**

sults for this implementation are discussed.

### 4.1. Implementation Issues

The concept of pCubes is very general, and leaves much freedom for an implementation to be optimized for a particular application. The only assumption we make for this implementation is that the values of the dimensions and the measure attribute can be expressed with data types of a small bounded size, e.g., integer, long, float, double (32 or 64 bits, respectively). Most attributes of real-world data collections have this property. For attributes with large data types, like strings, functions can often be found that compress the attribute's range.

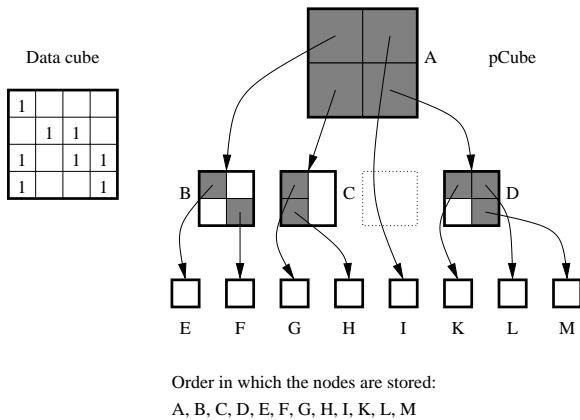
For tree structures where the regions of sibling nodes are allowed to overlap, accessing a single cell of the data cube possibly results in multiple search paths from the root down to the leaves. Also, counting data items in the overlap twice has to be avoided. We therefore decided to implement pCube such that the node-regions of sibling nodes do not overlap. A simple way to ensure this property is to use spatial decomposition. This can be done based on a certain division factor (e.g., partition the space along dimension  $X$  in halves) or based on attribute hierarchies (e.g., partition year into 12 months). To take the data distribution into account, the spatial partitioning is “optimized” to avoid storing empty nodes. We performed the experiments on a pCube with quadtree-like structure, which is described in the next section.

### 4.2. The pCube Implementation

The structure used for this pCube implementation is based on spatial decomposition. The children of a node partition the node-region of their parent into non-overlapping



regions of equal size by dividing *each* dimension in half. Among all possible tree nodes, we materialize in the tree only those nodes that contain non-empty cells and have two or more children containing non-empty cells. Figure 4 shows an example. Thus it is guaranteed that the number of internal tree nodes is less than the number of leaves. Since the number of leaves equals the number of non-empty cells in the data cube, our pCube implementation's *space requirement* depends linearly on the number on non-empty cells.



**Figure 4. Example for the pCube structure**

Depending on which aggregate functions pCube supports, the nodes store values that summarize the information in their node-region (see Tables 2 and 3). For instance, to support the aggregate function SUM, each internal node stores the sum of the positive and the sum of the negative values of its node-region.

Queries are performed as described for pCube in general. Their cost in the worst case depends linearly on the number of node-regions that intersect the query-cube. Consequently the *cost of a query* is in the worst case linear in the number of non-empty cells. On the average, however, better behavior can be expected because queries can take advantage of pre-computed values.

To update a cell, first the tree is descended to find the region where the cell is located. If the update inserts a new cell or deletes a cell, it could be necessary to insert a new node or shrink/delete a node, respectively. These operations, however, due to the spatial decomposition approach, do not propagate for more than one level of the tree. After changing the cell's value, the node-information is updated for all nodes on the path from the cell up to the root. For efficiently maintainable functions (see Section 3.6) this results in an overall worst case update cost that is linear in the height of the tree. Since we assumed that the data types of the attributes are of bounded size, the spatial decomposition guarantees that the height of the tree is also bounded by a constant. Consider an example where no data type uses more than 32 bits. At each level of the spatial decomposi-

tion each dimension is divided in half. Hence the height of the tree is bounded by  $32^2$ .<sup>2</sup> Thus the overall *cost of an update* is bounded by a constant.

For large data sets the complete pCube structure will not fit into a computer's main memory and has to be stored on disk (secondary memory). We store pCube in a packed format onto disk pages where a page can contain more than one node. Approximate results with error bounds can be output after accessing all relevant nodes in a page, while the next page is retrieved from disk.

The nodes on the disk pages are ordered by their level (i.e., how often the space was divided), nodes on the same level are stored in Z-order (space-filling curve, see for instance [8]). See Figure 4 for an example. Thus for queries that need to access a high percentage of the internal nodes, the corresponding pages can efficiently be accessed with sequential I/O. By keeping the first pages that contain the high-level (close to the root) nodes in fast storage (cache/memory) pCube quickly returns results without accessing the disk. Note, that maintaining the node order on the disk pages results in extra costs for insert and delete operations. Whenever an insert creates a new node or a delete causes a node to be shrunken, the correct page for inserting the node must be located. By indexing the page numbers with a B-tree and using the "greatest" node (according to the ordering of the nodes) in the page as the index attribute, the page can be found at an extra cost which is logarithmic in the number of non-empty cells.

### 4.3. Results

We ran various experiments for real-world and artificial data. Artificial data sets included dense and sparse data cubes. We examined the effect of skew (Zipf distributed values [30]), clusters and dimensions with different domain sizes.

The results presented here were obtained by executing 30,000 random range queries of varying selectivity on different data cubes (Table 4). In each dimension, a range was independently and uniformly selected based on selection predicates. Let  $[a \dots b]$  denote the set  $\{a; a + 1; \dots; b\}$  and let dimension  $D$ 's domain be  $[\min_D \dots \max_D]$ . Then the following ranges were selected for dimension  $D$  (where  $x$  and  $y$  are the randomly determined bounds):

1.  $[\min_D \dots x]$  with 30% probability
2.  $[x \dots y]$  with 10% probability
3.  $[x \dots x]$  (single value) with 10% probability

<sup>2</sup>For floating point data types that encode a large range with a comparably small number of bits, the floating point numbers are mapped to corresponding bit strings (order-preserving). Then the spatial decomposition and the accesses are performed based on the "encoded" numbers.

4.  $[\min_D \dots \max_D]$  (complete domain) with 50% probability

Note, that this selection favors queries that are disadvantageous for our pCube implementation. It is very likely that range queries are generated that cut through many dimensions, i.e., that intersect a large number of tree nodes, without completely covering their node-regions. Furthermore, in many high-dimensional applications, it is realistic to assume users to specify a range over a limited number of dimensions and leave other dimension attributes unspecified.

In each experiment we measured the number of page accesses for obtaining results of different approximation qualities for queries with a selectivity between 0% and 100%. Those numbers are compared to two benchmark algorithms. For our experiments we set pCube's parameters such that the cells are stored in Z-order. We used the same ordering to pack the data onto disk pages for the benchmark algorithms.

The first comparison benchmark algorithm ("Ideal Index") behaves like an optimal index structure where an index on each dimension exists. It is motivated by real index structures in database systems. If an index exists on an attribute  $A$ , then pointers to all data items whose  $A$ -value satisfies the selection condition of the query can be obtained from the index. Instead of scanning the whole database, only the selected data items are retrieved. By creating an index on each dimension attribute, this technique generalizes to selections where conditions on more than one attribute are given in the query. First, for each dimension the set of data items that satisfy the selection condition for the dimension are obtained from the corresponding (one-dimensional) index. Then, by intersecting the result sets, the set of all data items that satisfy the whole selection condition is computed. Finally only the selected data items are retrieved from the database. Our Ideal Index benchmark simulates this behavior. In contrast to the real index, we assume that it gets the pointers to the selected data items (i.e., cells) "for free". Thus the Ideal Index only accesses those disk pages that contain the selected data items. Note that this gives the Ideal Index benchmark a great advantage over any real technique. It not only saves the accesses to the index pages, but also the cost of combining the result sets. For high-dimensional data sets those costs can be very high. For instance, a range query with a selectivity of 80% in each dimension only selects about  $(80\%)^{15} = 3.5\%$  of the cells of a 15-dimensional data cube. Thus the costs of obtaining the large intermediate result sets for each dimension are high, while the final result set will be rather small.

When using index structures like B-trees, the index pages are accessed with random I/O (each access incurs seek and rotational latency). Random I/O are much slower than sequential I/O where disk pages are read sequentially. We therefore implemented a second benchmark algorithm

("Ideal Scan") that is based on sequential I/O. We assume that this algorithm "magically" knows about the first and the last page on disk that contain data items (cells) that are selected by the query. The access cost is then measured as the number of sequential page accesses to pages between and including the given pages.

Note that both benchmark algorithms are highly idealized. They constitute a lower bound for the corresponding classes of techniques that do not use pre-aggregation and do not provide any early feedback with error bounds. Also, it should be kept in mind that all page accesses of the Ideal Scan are sequential I/O. For pCube and the Ideal Index some accesses will be sequential, some random. pCube can take advantage of the order in which the nodes are stored on disk. For queries that are expected to intersect a high percentage of tree nodes, it can be decided to perform a sequential scan on the pages that contain the internal nodes.

Figure 5 contain the experimental results for a real-world data set (CENSUS, see Table 4). The data was obtained from the U.S. Census Bureau database [27]. We extracted the 1995 Current Population Survey's person data file. The data cube has 10 dimensions (age, class of worker, family type, age recode, educational attainment, marital status, race, Hispanic origin, sex, household summary). The data cube is extremely sparse (less than 0.001% non-empty cells). Figures 5(a) and 5(b) show that queries with large query-cubes benefit from storing aggregate information at internal nodes, irrespective of whether the query runs till completion (exact result) or is aborted earlier (when the approximation is good enough). Recall that in the case of large query-cubes when the query-cube is expected to intersect a high percentage of tree nodes, the disk accesses are performed sequentially. If a user is satisfied with a relative approximation error of 20%, the savings in access costs are substantial. In the case of the Ideal Index (Figure 5(c)), the larger the query-cube, the more non-empty cells are included and the more pages have to be accessed. Note, that real-world algorithms require additional accesses to index pages to find the relevant data pages. Index accesses are slow random I/O. In contrast, when simply scanning the data sequentially, one can take advantage of faster sequential I/O. The Ideal Scan benchmark (Figure 5(d)) scans the smallest possible set of data pages, such that the relevant data is retrieved. It has to access almost the entire set of data pages, even for small ranges. Thus it is not efficient for queries that select a small number of data items only. The results generally show that the user can be provided with fast approximate answers at almost no extra costs. When the application does not require exact results, savings in query costs can be expected.

For a dense data set pCube performs even better (Figure 6). The results shown here were obtained for a four-dimensional data cube where each dimension has 16 pos-

Data set	Dimensions	Aggregate function	Non-empty cells	Internal nodes (pCube only)
CENSUS	10	COUNT	45988 (352 pages)	11598 (161 pages)
DENSE	4	SUM	39319 (151 pages)	4369 (80 pages)
ZIPF	11	SUM	49996 (404 pages)	10698 (170 pages)

**Table 4. Data sets used in the experiments**

sible values (see Table 4 for more details). Especially for large query-cubes pCube clearly outperforms the benchmark algorithms. The main reason is that for a dense data cube the spatial decomposition produces a balanced tree where almost all possible internal nodes are created. Thus the number of leaf nodes accessed is dramatically reduced.

Since our pCube implementation is based on a regular spatial decomposition technique, we expected a performance degradation with increasing skew in the data distribution. The worst results were obtained for an 11-dimensional data cube (dimension domains of size 2, 4, 8, 16, 32, 64 and the others 128) where the values for eight dimensions are Zipf distributed (skew parameter between 1.8 and 1.2), and the values for the other three are uniformly distributed. Due to the skew the tree is poorly balanced. Queries that access nodes in the dense region of the data cube intersect many tree nodes, resulting in a large number of disk accesses. However, even for this “worst case scenario” the pCube implementation provides fast answers with error bounds at moderate extra costs (Figure 7).

## 5. Conclusion

The experiments show that pCube is a powerful concept for supporting efficient online aggregation. In contrast to earlier approaches, pCube provides early feedback with absolute error bounds for queries on data cubes. Since the worst case update and query costs depend on the number of *non-empty* cells, rather than on the dimensionality or the size of the data cube (including empty cells), our technique can efficiently handle high-dimensional and sparse data cubes. The combination of a hierarchical data structure and simple approximation models at the tree nodes enables fast updates and progressively improving error bounds. Also, by taking advantage of the hierarchical structure, some aggregate functions that are not self-maintainable and can not be made self-maintainable, e.g., MAX and MIN, can be handled efficiently by pCube.

An important goal of statistical analysis and decision support queries lies in finding unusual information in the data set, for instance extremum statistics. Techniques that primarily capture “typical” or “average” data items (e.g., most types of histograms, sampling) and only provide probabilistic error bounds, are not well-suited for this type of analysis. In contrast to that, pCube’s absolute error bounds

guarantee that the unusual information is contained in the approximate answer.

In real-world applications attributes often form hierarchies. The selection of ranges then typically corresponds to those hierarchies (e.g., one can expect more queries that summarize over the range 01/01/1999-12/31/1999, than for the range 01/14/1999-01/13/2000). Incorporating this knowledge into the pCube implementation would lead to a better “match” of query boundaries and node-regions. For pCubes that are based on spatial decomposition a fixed hierarchy can be chosen for an attribute. Then the partitioning is performed based on the hierarchy (e.g., instead of partitioning a week into two “half-weeks” it will be split into seven days). We expect such a partitioning to result in performance improvements.

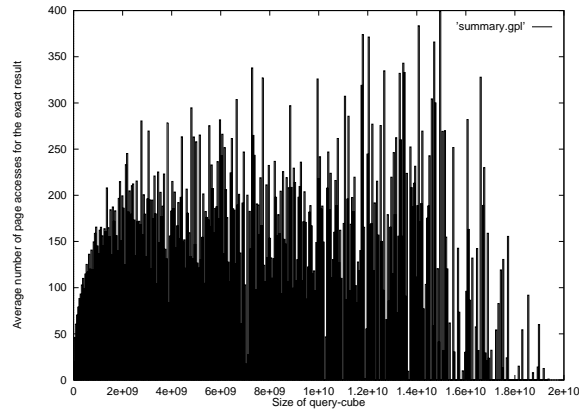
The concept of pCube is very general. Our experiments were performed for just one of its many possible implementations. Still the results are very encouraging and show, that early feedback with absolute error bounds can be added to data analysis tools at a low extra cost. There are, however, aggregate functions like MostFrequentValue, for which pCube in its current version is not applicable. Our future work includes the development of concepts that provide pCube’s properties for diverse functions and hierarchy based partitioning.

## 6. Acknowledgements

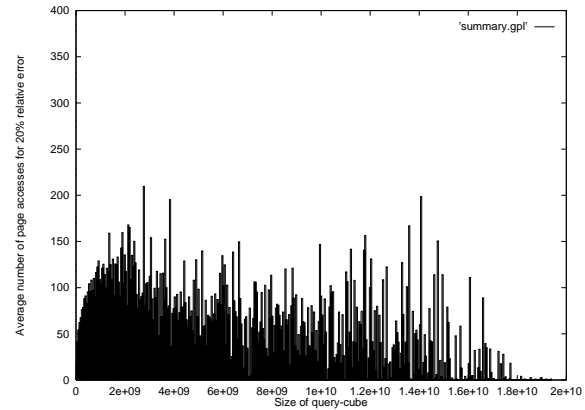
The authors would like to thank Christian Lang and Kristian Kvilekval for productive discussions. We are also grateful for the useful comments made by the referees that helped in improving this work.

## References

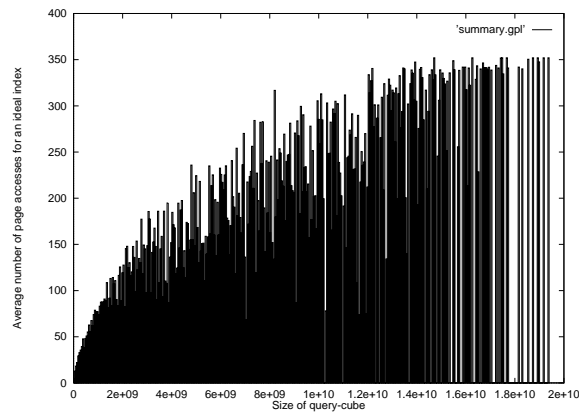
- [1] S. Acharya, P. B. Gibbons, V. Poosala, and S. Ramaswamy. Join synopsis for approximate query answering. In *Proc. ACM SIGMOD*, 1999.
- [2] D. Barbará et al. The new jersey data reduction report. *Data Engineering Bulletin*, 20(4), 1997.
- [3] D. Barbará and M. Sullivan. Quasi-cubes: Exploiting approximations in multidimensional databases. *SIGMOD Record*, 26(3), 1997.
- [4] D. Barbará and X. Wu. Using loglinear models to compress datacubes. Technical report, George Mason University, 1999.



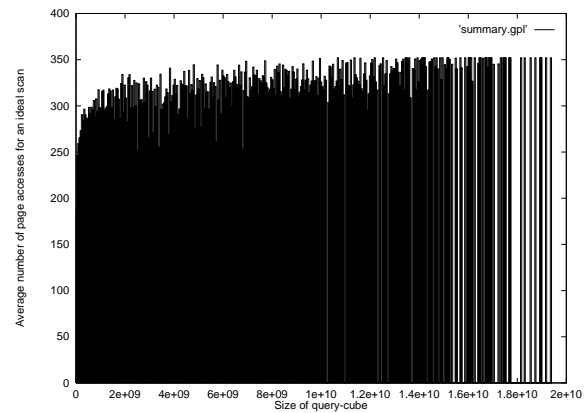
(a) pCube, page accesses for the exact result



(b) pCube, page accesses for a relative error of 20%



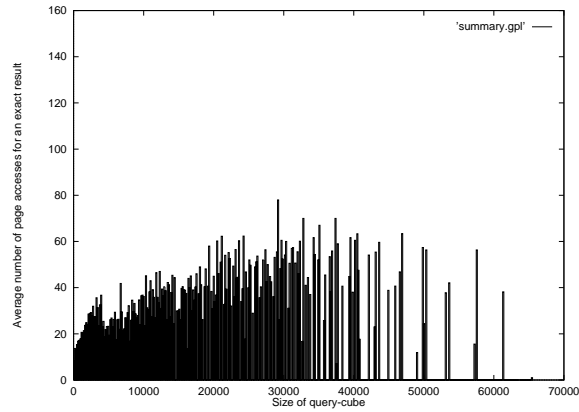
(c) Ideal Index, page accesses for the exact result



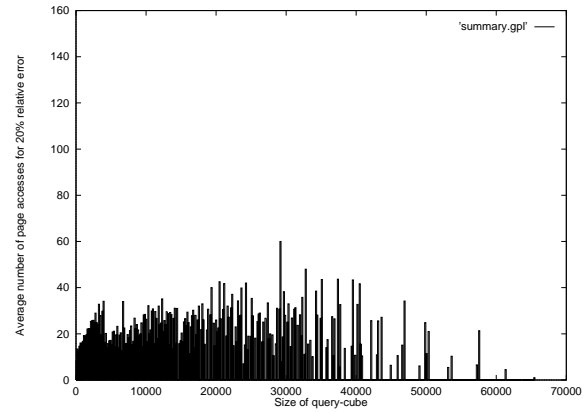
(d) Ideal Scan, page accesses for the exact result

**Figure 5. CENSUS dataset**

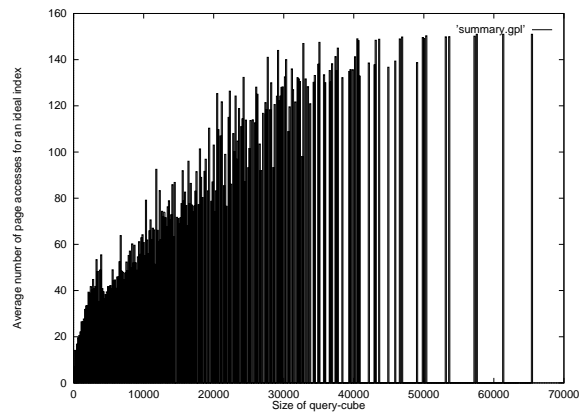
- [5] K. Beyer and R. Ramakrishnan. Bottom-up computation of sparse and iceberg CUBEs. In *Proc. ACM SIGMOD*, 1999.
- [6] C.-Y. Chan and Y. E. Ioannidis. Hierarchical cubes for range-sum queries. In *Proc. 25th VLDB*, 1999.
- [7] D. W. Cheung, B. Zhou, B. Kao, K. Hu, and S. D. Lee. DRO-LAP — a dense-region based approach to on-line analytical processing. Technical report, Univ. of Hong Kong, 1999.
- [8] V. Gaede and O. Günther. Multidimensional access methods. *ACM Computing Surveys*, 30(2), 1998.
- [9] S. Geffner, D. Agrawal, and A. E. Abbadi. The dynamic data cube. In *Proc. EDBT*, 2000.
- [10] S. Geffner, M. Riedewald, D. Agrawal, and A. E. Abbadi. Data cubes in dynamic environments. *Data Engineering Bulletin*, 22(4), 1999.
- [11] P. B. Gibbons and Y. Matias. Synopsis data structures for massive data sets. Technical report, Bell Labs, 1998.
- [12] S. Goil and A. Choudhary. BESS: Sparse data storage of multi-dimensional data for OLAP and data mining. Technical report, Northwestern University, 1997.
- [13] J. Gray, S. Chaudhuri, A. Bosworth, A. Layman, D. Reichart, M. Venkatrao, F. Pellow, and H. Pirahesh. Data cube: A relational aggregation operator generalizing group-by, cross-tab, and sub-totals. *Data Mining and Knowledge Discovery*, pages 29–53, 1997.
- [14] H. Gupta, V. Harinarayan, A. Rajaraman, and J. D. Ullman. Index selection for OLAP. In *Proc. ICDE*, 1997.
- [15] V. Harinarayan, A. Rajaraman, and J. D. Ullman. Implementing data cubes efficiently. In *Proc. ACM SIGMOD*, 1996.
- [16] J. M. Hellerstein et al. Interactive data analysis: The CONTROL project. *IEEE Computer*, August 1999.
- [17] J. M. Hellerstein, P. J. Haas, and H. J. Wang. Online aggregation. In *Proc. ACM SIGMOD*, 1997.
- [18] C. Ho, R. Agrawal, N. Megiddo, and R. Srikant. Range queries in OLAP data cubes. In *Proc. ACM SIGMOD*, 1997.
- [19] T. Johnson and D. Shasha. Some approaches to index design for cube forests. *IEEE Data Engineering Bulletin*, 20(1), March 1997.



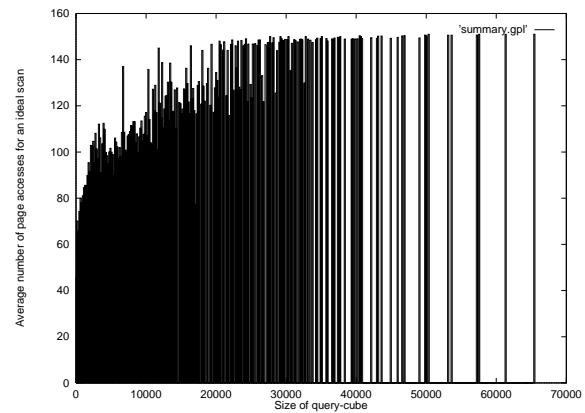
(a) pCube, page accesses for the exact result



(b) pCube, page accesses for a relative error of 20%



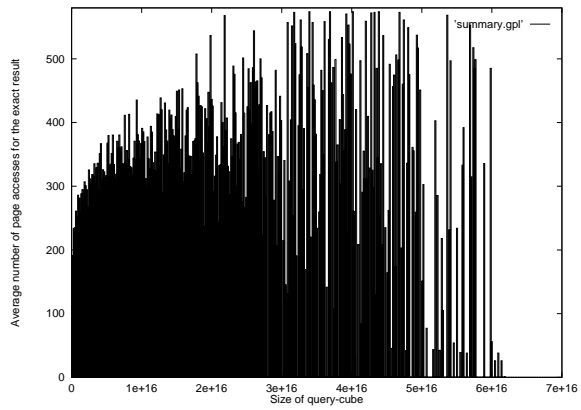
(c) Ideal Index, page accesses for the exact result



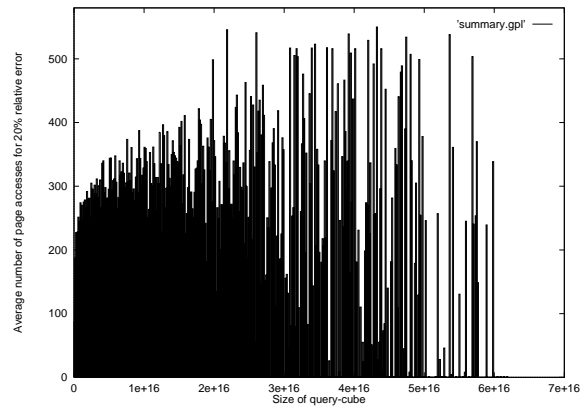
(d) Ideal Scan, page accesses for the exact result

**Figure 6. DENSE dataset**

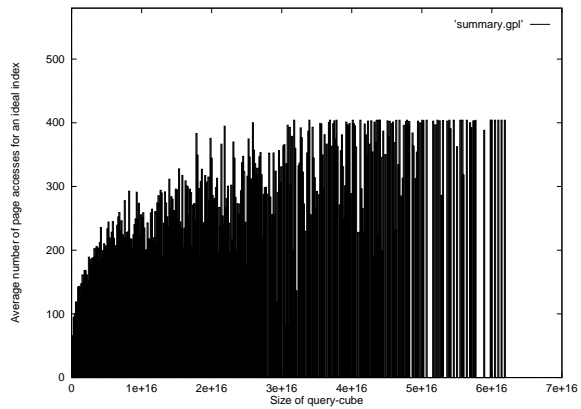
- [20] Y. Kotidis and N. Roussopoulos. An alternative storage organization for ROLAP aggregate views based on cubetrees. In *Proc. ACM SIGMOD*, 1998.
- [21] I. S. Mumick, D. Quass, and B. S. Mumick. Maintenance of data cubes and summary tables in a warehouse. In *Proc. ACM SIGMOD*, 1997.
- [22] P. E. O'Neil and D. Quass. Improved query performance with variant indexes. In *Proc. ACM SIGMOD*, 1997.
- [23] N. Pendse and R. Creeth. The OLAP report. <http://www.olapreport.com/Analyses.htm>, 1999. Parts available online in the current edition.
- [24] V. Poosala and V. Ganti. Fast approximate answers to aggregate queries on a data cube. In *Proc. SSDBM*, 1999.
- [25] N. C. Rowe. Antisampling for estimation: An overview. *IEEE Transactions on Software Engineering*, 11(10), 1985.
- [26] J. Shanmugasundaram, U. Fayyad, and P. S. Bradley. Compressed data cubes for OLAP aggregate query approximation on continuous dimensions. In *Proc. ACM SIGKDD*, 1999.
- [27] U.S. Census Bureau. <http://www.census.gov/main/www/access.html>.
- [28] J. S. Vitter and M. Wang. Approximate computation of multidimensional aggregates of sparse data using wavelets. In *Proc. ACM SIGMOD*, 1999.
- [29] W. Wang, J. Yang, and R. Muntz. STING: A statistical information grid approach to spatial data mining. In *Proc. 23rd VLDB*, 1997.
- [30] G. K. Zipf. *Human Behaviour and the Principle of Least Effort*. Addison-Wesley, Reading, MA, 1949.



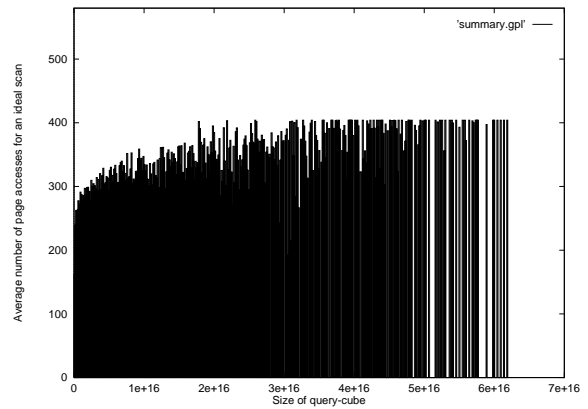
(a) pCube, page accesses for the exact result



(b) pCube, page accesses for a relative error of 20%



(c) Ideal Index, page accesses for the exact result



(d) Ideal Scan, page accesses for the exact result

**Figure 7. ZIPF dataset**