



present a technique that tackles this problem. Similar issues arise in transaction recording systems [13] as well as in statistical databases [1, 16].

Even if the original base data exist, the ability to reconstruct the original data from the summaries is of great value. Computing summaries is essentially a data reduction process in which great information loss takes place. In order to reconstruct the data, various assumptions have to be made about the statistical properties of the reduced data. Given the reconstructed and the original data at hand, we can test how strong our assumptions about the original data were, just by comparing the two. This is useful in reasoning about the properties of the underlying data set and could be of great value in data mining. It can help detect correlations in the data, and identify deviations, that is, values that do not conform to the underlying model. Such results are of great interest to the analyst, because they indicate local or global abnormalities.

In this paper, we propose the use of an information theoretic principle for the reconstruction of the original data from the summarized forms. Our reconstruction technique is based on the well recognized and widely applicable information theoretic principle of *maximum entropy* [9]. We present algorithms for the efficient reconstruction of data from the aggregates. Moreover, using an information theoretic formalism we identify and describe in detail an alternate benefit of the proposed reconstruction techniques, namely the ability to ‘rank’ each reconstructed value by its potential ‘interest’ to the user, as a means of aiding data analysis. The notion of interest in data mining is a difficult problem and several approaches have been proposed for its formal definition [22]. Although all the measures of interest are subjective and heavily dependent on the application, we argue that an information theoretic approach to this problem, besides being mathematically rigorous, appears conceptually appealing as well.

Our contributions can be summarized as follows. We propose a method for reconstructing multidimensional values from the already computed aggregated data, which does not require any additional special data structures. We describe an extension to the above method, in which we are able to provide quality (error) guarantees for the reconstruction. We present a method to identify and rank deviations in multidimensional datasets. This method does not depend on any a priori or domain knowledge for the problem at hand, and it also does not require any parameter settings or calibrations. Finally, we perform an extensive experimental evaluation, using both synthetic and real datasets.

The rest of the paper is organized as follows. Section 2 reviews related work. In Section 3 we present some necessary background material and Section 4 presents the reconstruction algorithm. Section 5 discusses the two particular application scenarios of the algorithm. In Section 6 we present

experimental results evaluating the performance and the utility of the proposed algorithms, and we present our conclusions in Section 7.

## 2 Related Work

The principle of maximum entropy [9] has been successfully applied in different domains, including linguistics [6] and databases [11]. Faloutsos et al. apply maximum entropy in addition to other techniques, for one dimensional data reconstruction [11]. In a sense our work, generalizes the work of Faloutsos et al. in multiple dimensions.

There exists a sizeable bibliography in histogramming techniques and approximate query answering [19, 12, 23, 2, 21, 4, 5]. Our approach is fundamentally different. Previous work focused on the problem of data reconstruction by constructing specialized summarized representations of the data. We argue, that since aggregated data are already stored in the warehouse, it is imperative to examine the quality of reconstruction one can attain from the aggregates. This approach could potentially have a great deal of benefit, if the quality of the reconstruction is good.

The problem of identifying interesting values in a dataset is related to deviation detection. Arning et al. [3] try to identify the subset of a database that is most dissimilar to the rest of the data. Other approaches discuss algorithms specialized in metric spaces that scale to large datasets [14]. The drawback of the above approaches is that the user is required to come up with the right selection of functions and parameters which requires a great deal of effort. Our algorithm does not need such input, making the whole procedure less cumbersome and more robust. In that sense, our work is closer to the framework proposed by Sarawagi et al. [20]. They describe an algorithm that mines the data in a data cube for exceptions. However, this method is computationally expensive, depends on the computation of the entire data cube, and cannot accommodate updates.

## 3 Background

Consider a schema  $R = (A_1, A_2, \dots, A_n, Y)$  and  $r$  an instance of  $R$ .  $A_1, \dots, A_n$  are attributes having some specified type and  $Y$  is a measure attribute. Attribute  $Y$  could represent volume of sales, dollar amount, number of calls, etc. Attributes  $A_i, 1 \leq i \leq n$  include dimension attributes and hierarchies possibly defined on them. For example in Figure 1(a) the schema is  $R(\text{STATE}, \text{CITY}, \text{BRAND}, \text{GENDER}, \text{SALES})$ . Attributes  $\text{STATE}$  and  $\text{BRAND}$  define hierarchies on  $\text{CITY}$  and  $\text{GENDER}$  respectively.

With a schema of  $n$  attributes, assume that  $h$  attributes define hierarchies on the remaining  $d = n - h$  attributes. We can enumerate the nodes of hierarchy  $i, 1 \leq i \leq d$  starting bottom up, such that  $h_i = j$  denotes the  $j$ th node

of hierarchy  $i$ . Viewing the table of Figure 1(a) with its dimensions and hierarchies as a multidimensional grid, we refer to the  $d$ -dimensional vector  $(h_1, \dots, h_d)$  as a *grid query*. For example, the grid query  $(2, 1)$  is a point query asking for the sales of women's Levi's in Queens, while the grid query  $(7, 5)$  is a range query referring to the entire upper half of the dataset. In our experience, grid queries are the most common in warehouse environments. The reason we are interested in this special class of queries is that it significantly reduces the space of possible queries. For the rest of this paper we will refer to grid queries simply as queries, and to the subset of the entire dataset they involve as dataset. Note though, that by answering a grid query using our techniques, we can also answer any other query asking for a subset of the answer we computed.

A basic observation about any instance  $r$  of  $R$  is that it can be viewed as a discrete  $n$  dimensional probability distribution,  $P_r(A_1, \dots, A_n)$ . This can be accomplished by normalizing the value  $Y$  on each row of  $r$  by the sum of all  $Y$  values. The analogy between  $r$  and  $P_r$  can be extended further. Consider for example the following query:

```
SELECT  A1, A2, . . . , An-1, sum(Y)
FROM    r
GROUPBY A1, A2, . . . , An-1.
```

The outcome of this query is one of the  $n$  marginal distributions of  $P_r$  of order (number of variables)  $n - 1$ . All we need to do, is normalize  $sum(Y)$  by the sum of all  $Y$  values. (This normalization is needed for mathematical correctness, but is not required by the techniques discussed in this paper.) Reasoning similarly, one can draw the analogy between all the group by's on  $r$  and all the marginal distributions of  $P_r$ . Thus, we can view the problem of reconstructing  $r$  from its aggregates as analogous to the problem of reconstructing an  $n$  dimensional probability distribution from a number of its marginal distributions. Based on this analogy, in the rest of this paper, we use the terms group by on instance  $r$  and marginal distribution of  $P_r$  interchangeably.

### 3.1 Maximum Entropy Distributions

Let  $P(A_1, \dots, A_n)$  be an  $n$  dimensional discrete probability distribution, to be estimated from a number of its marginals. With  $n$  variables, there are  $2^n - 2$  marginals (excluding the grand total and the base data) of  $P$  in total. Moreover, there are  $\binom{n}{k}$  marginals with  $k$  variables (equivalently of order  $k$ ). Let  $S$  be an arbitrary subset of the powerset of  $X = \{A_1, \dots, A_n\}$ . The problem of *maximum entropy estimation* of  $P$  is defined as follows:

**Problem 1** Maximize the entropy  $H(P)$  of  $P$  over all probability distributions satisfying:

$$P(X) \geq 0, \text{ with equality outside the support of } P, \\ \sum P(X) = 1, \text{ and}$$

$$\forall i \in S, \sum_{j \in S-i} P(j) = P(i), \\ \text{where } P(j) \text{ is a marginal distribution of } P, \text{ with } j \in S.$$

The maximum entropy estimation of  $P$  is a model fitting technique. It finds the model with the 'least' information or fewest assumptions given the specified constraints. The method appears ideal for this estimation problem as it is designed for lack of data rather than an excess thereof. The overriding principle in maximum entropy is that when nothing is known the distribution should be as uniform as possible. The constraints (the marginal distributions in our case) specify the regions where the estimated distribution should be minimally non-uniform in order to satisfy these constraints.

### 3.2 Properties of Maximum Entropy

Let  $P(X)$  be an  $n$ -dimensional probability distribution. We denote by  $P_i, 1 \leq i \leq n - 1$  the maximum entropy approximation to  $P$  using only marginals of order  $i$  as constraints. Then the following theorems hold.

**Theorem 1** [15] Let  $p$  be a member of the class of discrete probability distributions that can be constructed using a specified set of lower order marginal distributions of order  $i$ . Then, among all distributions  $p, P_i$  minimizes:

$$D(p, P) = \sum_x p(X) \log \frac{p(X)}{P(X)} \quad (1)$$

The measure  $D$  is known in the literature as the relative entropy [9] and measures the similarity of two probability distributions. More precisely  $D$  is a measure of the inefficiency of assuming that the distribution is  $p$  when the true distribution is  $P$ . In statistics  $D$  arises as an expected logarithm of the likelihood ratio, and one can show that by minimizing  $D(p, P)$  we also minimize the  $\chi^2$  test between  $p$  and  $P$  [15].

**Theorem 2** [15] Let  $P_i, 1 \leq i \leq n - 1$  be the maximum entropy approximation of  $P$  using only marginals of order  $i$ . Then the following inequality holds:

$$D(P_1, P) \geq D(P_2, P) \dots \geq D(P_{n-1}, P) \quad (2)$$

Theorem 2 states that a better estimation of  $P$  can be performed by using marginals of order  $i + 1$  than marginals of order  $i$ , for  $1 \leq i \leq n - 1$ .

## 4 Algorithmic Solution

Problem 1 is a constraint optimization problem that is not amenable to a general closed form solution. The standard technique for solving this maximization problem is the method of *Lagrange Multipliers* [7]. This method calls for the solution of a rather complex system of equations with

high computational cost. Moreover, a different system of equations has to be derived each time, depending on the specified constraints. Ideally, we are after an algorithmic approach that can operate directly on the specified constraints and derive the estimation with minimal human intervention.

In this section, we propose the use of an algorithmic approach. The technique is called *Iterative Proportional Fitting* (IPF), and was introduced in the statistics literature [10]. It is an iterative algorithm that converges to the maximum entropy solution. IPF has the following properties [8]: it always converges monotonically to the required unique maximum entropy estimation, given a number of marginals; a stopping rule may be used that ensures accuracy to any desired degree in the solution; the estimates depend only on a particular small set of marginals; convergence, as well as the speed of convergence, are not directly affected by the starting values; finally, in some cases, convergence to the exact estimates is achieved after only a single iterative cycle.

#### 4.1 Description of IPF

Let  $r(A_1, A_2, \dots, A_n, Y)$  be a relational instance. We specify a multidimensional range of interest  $[d_{1_s}, d_{1_e}] \times \dots \times [d_{d_s}, d_{d_e}]$ , defined by a  $d$ -dimensional grid query  $Q = (h_1, \dots, h_d)$ . We assume that all the marginals of order  $k$ ,  $1 \leq k \leq d - 1$  needed for the reconstruction of  $Q$  have been located<sup>1</sup> and are given as constraints on Problem 1. Without loss of generality, assume that the set  $\{A_1, \dots, A_d\}$  contains all the attributes with hierarchies defined on them. Let  $S_k$  be the set containing all subsets of the powerset of  $\{A_1, \dots, A_d\}$  with  $k$  elements. Clearly  $|S_k| = \binom{d}{k}$ . Let  $P_i(j)$ ,  $1 \leq i \leq \binom{d}{k}$ ,  $j \in S_k$ , be the marginals of  $P$  of order  $k$  corresponding to  $Q$ . On each  $P_i(j)$ , the aggregation has been computed on all attributes not present in  $j$ .

We denote as  $Y_{A_1 A_2 \dots A_n}$ , the value of attribute  $Y$  for the specific combination of the  $A_i$ , attribute values ( $d$  of those are specified by the grid query and the remaining  $n - d$  from the hierarchy). We denote as  $\hat{Y}^{(t)}_{A_1 A_2 \dots A_n}$  the estimate of the value of  $Y_{A_1 A_2 \dots A_n}$  during the  $t$ -th iteration of the algorithm. IPF starts the reconstruction by initializing a  $d$ -dimensional grid,  $G$ , of size  $d_{i_e} - d_{i_s} + 1$  per dimension  $i$ , to one. We refer to each element of the  $d$ -dimensional grid as a *cell*. In addition, given the initialized  $G$ , it computes its marginals of order  $k$ . Let  $\hat{P}_i^{(t)}(j)$ , denote the marginals computed from  $G$  in the  $t$ -th iteration of the algorithm. Denote  $\hat{P}_i^0(j)$ , the marginals after the initialization.

At each iteration, IPF loops over all marginals  $i$ ,  $1 \leq i \leq \binom{d}{k}$ , and all grid cells,  $l_1 \in [d_{1_s}, d_{1_e}], \dots, l_n \in [d_{d_s}, d_{d_e}]$ , and adjusts the values of the grid cells according to the formula:  $\hat{Y}_{l_1 l_2 \dots l_d}^{(t+1)} = \hat{Y}_{l_1 l_2 \dots l_d}^{(t)} \frac{P_i(j)}{\hat{P}_i^{(t)}(j)}$ .

<sup>1</sup> We defer our discussion of the issues related to the identification of the suitable marginals to Section 5.

This procedure is guaranteed to converge to the maximum entropy estimation of  $P$  given the specified collection of marginals. The estimates converge in a monotonically decreasing fashion, and we commonly choose to terminate the iterations when the change in each individual estimate becomes smaller than some user specified  $\delta$  value. For all the experiments presented herein we set  $\delta$  to 10% of the median of the values designated by  $Q$ . We chose the median because it is not affected by any extreme values, and it can be efficiently computed [17]. In addition, we have shown that the algorithm is not very sensitive to the value of  $\delta$  [18], making the specific choice less important. A skeleton of the IPF algorithm is given in Figure 2. Due to lack of space, an example of applying the IPF algorithm is presented in the full version of the paper [18].

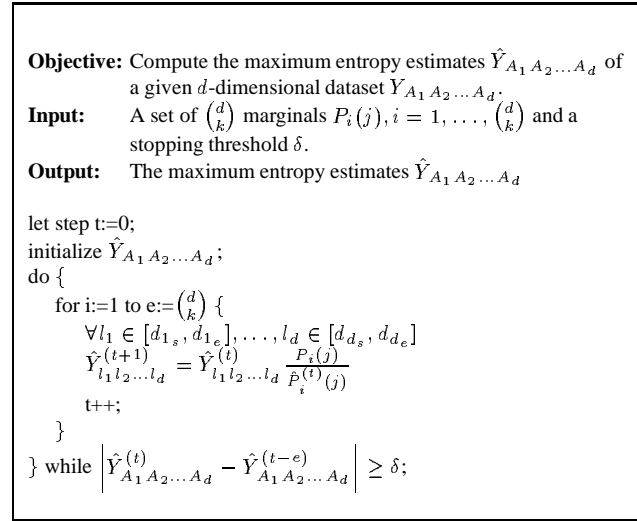


Figure 2. The IPF algorithm.

#### 4.2 Algorithmic Complexity

The IPF algorithm requires as input the marginals corresponding to a query  $Q$ . It then iterates over a grid  $G$  which is the same size as the result set of  $Q$ . We can safely assume that the marginals fit in main memory, but this may not be true for  $G$ .

If the available memory is large enough to fit the grid  $G$  then the algorithm needs to read from disk only the marginals. All the subsequent operations take place in main memory. Considering the large sizes of memory that are commonplace nowadays, we expect that the algorithm will be able to provide fast answers to a significant number of queries by operating on memory-resident data.

In the case where  $G$  does not fit in memory, the algorithm during each iteration makes  $\binom{n}{k}$  passes over  $G$ , where  $n$  is the dimensionality of  $Q$ , and  $k$  is the order of the marginals

used for the estimation. After having computed the estimates based on a single marginal for all the values in  $G$ , the algorithm has to make one more disk pass over  $G$  in order to calculate the new estimated marginals. In our implementation, we cut the number of passes in half by incorporating the update of the estimated marginals with the computation of the base values, reducing the cost of the algorithm to  $\binom{n}{k}t$  passes over  $G$ , where  $t$  is the number of iterations. The results that we report in the experimental section illustrate the **worst case** scenario, where the dataset does not fit in main memory.

### 4.3 Error Guarantees for the Approximation

For various applications, being able to provide error guarantees for the reconstruction of individual values is imperative. We can safely assume that at the time of the computation of the aggregates the original data are still available. Let  $W$  be the set of grid queries of interest. One approach to provide error guarantees for each query in  $W$  would be the following. Estimate the values of each query in  $W$ , while the original data are still available, compute the largest difference (between the actual and the estimated value) for each query in  $W$ , and store them separately. This would incur a storage overhead of  $O(|W|)$ . More formally, let us denote by  $Y_i$  the original value of some cell  $i$ , and with  $\hat{Y}_i$  its estimated value. We can use the absolute difference  $d_i = |\hat{Y}_i - Y_i|$  in order to provide an upper bound for the error. Assuming that a grid query  $Q$  encompasses  $N$  cells from the base data, the upper bound for the total error for  $Q$  can be calculated using the formula<sup>2</sup>  $N \cdot \max\{d_i\}$ ,  $1 \leq i \leq N$  (albeit this is an over-estimation of the real error).

We may extend the above approach in order to provide tighter error guarantees as follows. Store a number  $k$  (user defined) of the largest estimation errors for each query in  $W$ . Given a query in  $W$  that involves a number of the cells whose correct values have been explicitly stored, the reconstruction algorithm uses these correct values, and thus induces no error for the specific cells. Consequently, the overall error for the query is dramatically reduced. If the error guarantee per query should be bounded by a user specified value, we can choose  $k$  (the number of values to store) such that the overall error of reconstruction satisfies the error bound.

In many cases the number of cells for which the approximation is really poor will be relatively small. In the experimental section we evaluate this argument and we present graphs depicting the distribution of errors for the real datasets we used. The graphs show that only a minor percentage of the cells exhibit high errors, and thus, can be ef-

<sup>2</sup> Without loss of generality, we assume that the error metric is the Root Mean Square Error. Similar arguments hold if we choose other error metrics as well.

ficiently stored, inducing a small storage overhead.

## 5 Using IPF

In the following sections, we discuss issues related to the applications of IPF, both for query answering and knowledge discovery.

### 5.1 Computing the Marginals

Given a  $d$ -dimensional grid query  $Q$  the algorithm has to determine the order of the marginals to use for reconstruction as well as the ranges of these marginals pertinent to  $Q$ . In light of Theorem 2, marginals of order  $k$  will produce more accurate estimates than marginals of order  $k - 1$ . In the experimental section we present graphs exploring the time and accuracy tradeoffs related to this choice.

Assuming one decides to use the marginals of order  $k$ , the exact marginals relevant to  $Q$  have to be determined. A simple rewriting of  $Q$  can produce the  $\binom{d}{k}$  marginals of order  $k$  that should be queried in order to retrieve the values for IPF to operate on. Any grid query on relational instance  $r(A_1, \dots, A_n)$ , where  $A_1, \dots, A_d$  are dimension attributes and  $A_{d+1} \dots A_n$  define hierarchies on them, can be expressed in SQL as:

```
SELECT A1, A2, ... Ad
FROM r
WHERE Ad+1 = a1 and ... An = an-d,
```

where the values  $a_1, \dots, a_{n-d}$  designate the multidimensional range of  $Q$ . Let  $S = \{A_1, \dots, A_d\}$ . Then, the marginals of order  $d - 1$  relevant to the reconstruction of the grid query can be retrieved by issuing the following SQL query:

```
SELECT S - Ai
FROM r
WHERE Ad+1 = a1 and ... An = an-d,
```

for  $1 \leq i \leq d$ . This will give us all the  $\binom{d}{d-1} = d$  marginals. Similarly, the marginals of order  $d - 2$  relevant to the reconstruction of the grid query can be retrieved by:

```
SELECT S - {Ai, Aj}
FROM r
WHERE Ad+1 = a1 ... An = an-d,
```

for  $1 \leq i \leq d, i \leq j \leq d$ . Reasoning similarly we can construct the expressions for the choice of marginals of any order.

### 5.2 Mining Interesting Patterns

In the case that the base data are available, the proposed reconstruction technique has a different utility. Maximum entropy reconstruction from a number of marginals is performed based on the assumption that the marginals of interest are pairwise independent. By reconstructing the data

and comparing them with the base data, the validity of the pairwise independence assumption can be tested. Any data value that violates the pairwise independence assumption, will induce a larger reconstruction error than one that does not. This provides an automatic way to reason about the underlying correlations among attributes. For example if the volume of sales of men’s Levi’s jeans in Queens incurs the largest estimation error than all sales of Levi’s in the cities of NY state, we know that the volume of sales at Queens represents the strongest violation of the independence assumption among sales in NY state. We term such values *deviations*, because they deviate from the estimation model.

The basis for terming a specific value as a deviant can be a measure of the distance between the actual and the estimated value, i.e., the estimation error, such as absolute difference. However, the aforementioned metric may not always produce qualitative results. In order to remedy this situation, we can use a formula which normalizes the estimation error of a value with respect to the standard deviation  $\sigma$  of all the estimation errors returned by the algorithm  $s = \frac{|\hat{Y}_i - Y_i|}{\sigma}$ , where with  $Y_i$  we denote the original value of cell  $i$ , and with  $\hat{Y}_i$  its estimation. Then, we choose a cutting threshold for  $s$ , that can effectively prune all the normal perturbations in the dataset, leaving us with only the large deviations. The above technique splits the sorted set of deviations into two regions: it assigns the statistically large deviations to the first region, and the rest to the second one. We refer to the boundary point between those two regions as the *cutoff point*, and we demonstrate its role in the experimental evaluation. A commonly used threshold is for  $s = 2$ , which will prune 95% of the approximation errors as trivial, leaving only the largest 5% for consideration (the values follow from the properties of Normal distributions). The system can subsequently sort those deviating values, and pick the top- $k$  among them. In the experimental section, we present graphs that visualize the ability of the algorithm to spot deviations, using both synthetic and real datasets.

### 5.3 Optimization Problems

So far we have assumed that we have enough space to materialize all the marginals needed to reconstruct each query in our workload. In the case where only a subset of the marginals can be materialized, one would be interested in obtaining the most benefit in reconstructing queries while satisfying the imposed space constraints. The benefit in query reconstruction is defined as the number of “important” grid queries that can be reconstructed with the greatest accuracy, where importance may be user-defined, or determined by the query workload. The above situation gives rise to an optimization problem that one can show is NP-Hard, but is amenable to an efficient solution using heuristics. Moreover, assuming that we wish to utilize the reconstruction for iden-

tifying values of interest, ideally one would be interested in materializing the most informative marginals, that is the ones that are most likely to contain the values with the largest deviations. This is a similar optimization problem, and can be solved using the same basic algorithms.

The aforementioned problems constitute ongoing work.

## 6 Experimental Evaluation

In order to test the IPF algorithm we used a mixture of synthetic and real datasets. The synthetic datasets produced are derived by sampling uniform and Gaussian data distributions.

**Uniform:** We produced datasets of dimensionality 2, 3, and 4, with different variance values. The size of the datasets varied from 1,000 to 20,000 tuples.

**Gaussian:** We produced datasets of dimensionality 2, 3, and 4. The values were sampled from Gaussian distributions with 3 different sigma  $\sigma$  values. In the experiments we refer to these values for  $\sigma$  as *small*, *medium*, and *large*. The size of the datasets varied from 1,000 to 20,000 tuples.

We also experimented with three Gaussian datasets which had additional random noise uniformly distributed. The third dataset was derived from a mixture of two multidimensional Gaussian distributions. We refer to these datasets as *g\_small\_s*, *g\_large\_s*, and *g\_mix* respectively. The statistical properties of those datasets are reported in Table 1.

**Real:** Two of the real datasets, *calls* and *calls3*, are derived from AT&T proprietary data. They represent aggregated telephone calls in certain regions and over time. They are 2- and 3-dimensional, and their size is 10,000 tuples.

We also used *census*, a dataset from the U.S. Census Bureau, containing information about the age, education, and income of individuals. It is a 4-dimensional dataset from which we extracted instances of 10,000-50,000 tuples.

Table 1 summarizes the statistical properties of *calls* and *census\_10K*. The rest of the real datasets exhibit similar characteristics.

<i>dataset</i>	<i>min</i>	<i>max</i>	<i>mean</i>	<i>median</i>	<i>std.dev.</i>	<i>skew</i>
<i>g_large_s</i>	9.0	76.72	28.15	28.00	5.80	0.40
<i>g_small_s</i>	0.0	72.14	20.72	20.00	12.45	0.36
<i>g_mix</i>	0.0	106.68	24.66	22.00	13.57	0.53
<i>calls</i>	1.0	729.00	18.01	4.00	37.79	5.37
<i>census_10K</i>	1000	196623	24736	20000	23450	3.67

**Table 1. Statistical properties for some of the datasets used in the experiments.**

The error metric that we report in the experiments is the *Root Mean Square Error (RMSE)*, defined as  $RMSE = \sqrt{\frac{\sum_{i=1}^N (Y_i - \hat{Y}_i)^2}{N}}$ , where  $Y_i$  represents the original values in

the dataset,  $\hat{Y}_i$  the corresponding estimated values, and  $N$  is the total number of values in the dataset (remember that this refers to the dataset determined by some grid query).

Note that all the above datasets also have an additional measure attribute. Therefore, the total number of attributes for the datasets we used ranges from 3 to 5.

### 6.1 Exploring the Properties of the Algorithm

In this section we present experiments concerning the run-time of the algorithm, and the property stated in Theorem 2.

Figure 3(a) shows how the run-time of the algorithm changes when the dataset size increases. The graph demon-

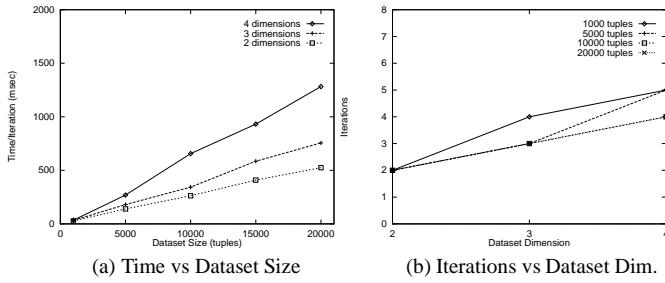


Figure 3. Scalability of the algorithm.

strates the outcome of experiments for datasets of different dimensionalities. As expected, there exists a linear relationship between the size of the dataset and the run-time of the algorithm. Moreover, when dimensionality increases, the number of marginals that the algorithm uses increases as well. This explains the steeper curves of the graph for the higher dimensions.

Other experiments that we performed indicate that there is no correlation between the dataset size and the number of iterations that the algorithm needs to perform in order to converge. Nevertheless, as Figure 3(b) depicts, the number of iterations increases with the dimensionality of the dataset.

The above results demonstrate that this approach can be effectively used by the analyst in real time, and in an interactive fashion, for middle-sized queries even when they do not fit in main memory.

The following graph is the experimental verification of Theorem 2. We used 4-dimensional datasets, and measured the error of the estimation when the algorithm operates with marginals of orders 1 to 3. Figure 4(a) displays the fact that when we use higher order marginals for the reconstruction of the same dataset, the error is diminishing. It is interesting to note here, that the relative benefit of employing marginals of higher order is increasing. Thus, when we use marginals of order 3 the reduction of the error is more acute. The greater accuracy that we are gaining by using marginals of high or-

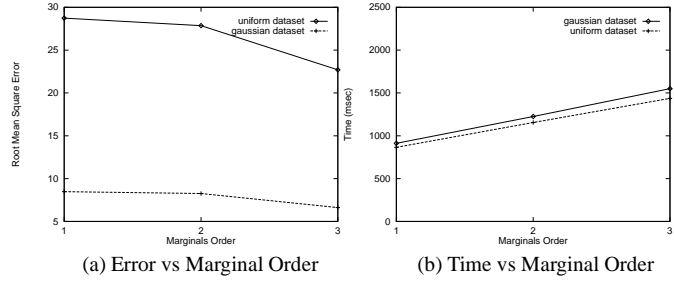


Figure 4. The effect of the order of the marginals used for the reconstruction.

der comes at the expense of time. The run-time of the algorithm increases with the order of the marginals (Figure 4(b)).

### 6.2 Evaluating the Accuracy of Reconstruction

A number of experiments try to investigate the behavior of the algorithm when reconstructing an unknown dataset.

The graph in Figure 5(a) shows how the error changes when the dataset size increases. We used three different 3-dimensional datasets drawn from Gaussian distributions with the same mean  $\mu$ , but different sigma  $\sigma$ , as described in the previous section. As is evident from figure 5(a), the

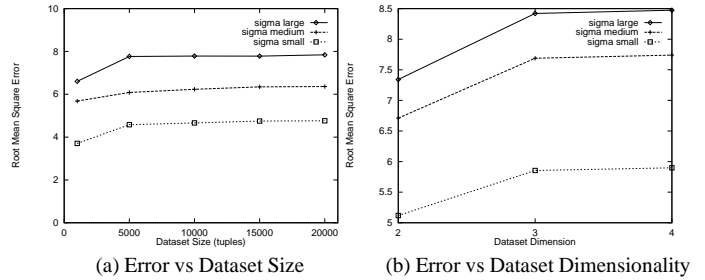


Figure 5. The effect of dataset size and dimensionality on the reconstruction error.

error of reconstruction is related to the variance of the underlying dataset and it increases as the variance increases. By increasing the size of the dataset, the error increases (but not dramatically) as more values are included in the computation of the error. The next graph, Figure 5(b), depicts how the dataset dimensionality affects the accuracy of the estimations. During this experiment we instructed the algorithm to use the marginals of the highest order common to all the datasets. It is evident that the reconstruction error increases with dimensionality; however, the increase seems correlated to the variance of the underlying dataset, since the increase of the error as the dimensionality increases is only nominal for datasets with similar variance.

The experiments with the uniform datasets gave results similar to the above, and we do not present them here.

In the following experiment we used the real datasets to verify the trends that the reconstruction error follows as indicated by the previous experiments. Figure 6(a) depicts a graph of the error for increasing dimensionality, while Figure 6(b) demonstrates the behavior of the error as the dataset size grows. These graphs show a deterioration in the ac-

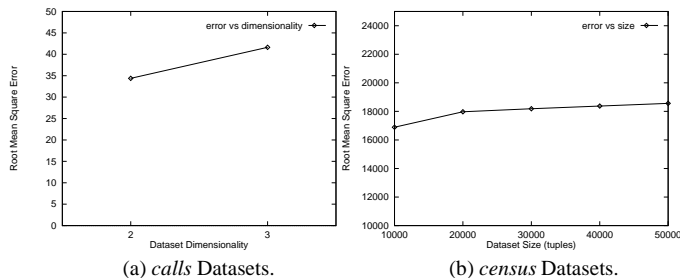


Figure 6. Error vs dimensionality and error vs dataset size for the real datasets.

curacy of reconstruction both when dimensionality goes up and when the size of the dataset increases, which is in accordance with the results of the experiments using the synthetic datasets.

### 6.3 Reconstruction with Error Guarantees

In the following experiments, we try to assess the benefit of providing error guarantees. The method we propose in order to achieve this goal (as discussed in Section 4.3) explicitly stores a small number of deviating values, which are subsequently used during the reconstruction phase to diminish the error.

In the first set of experiments we explore the distribution of the size of the estimation errors (i.e., the absolute error between the real and the estimated value for an individual cell). The graph in Figure 7(a) depicts the distributions for

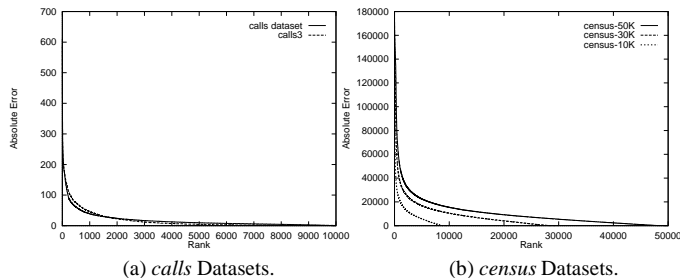


Figure 7. Distribution of absolute error for real datasets.

the datasets *calls* and *calls3*. Both curves indicate that the error sizes follow a Zipf-ian like distribution, with a very small

fraction of the instances having large values. This fact indicates that the choice to store the largest estimation errors as extra information is very likely to pay off during reconstruction. Figure 7(b) shows the same graph for different sizes of the *census* dataset. The results show that for all the sizes the error follows a highly skewed distribution.

The next experiments evaluate the relative benefit of storing a number of deviating values per query in order to guarantee a specified reconstruction error level. Figure 8 shows the error of the reconstruction when the number of deviations that the algorithm is using increases from 0 to the size of the dataset. The user may choose between those two ends according to the application requirements for error guarantees, and the space restrictions on the number of values that can be stored.

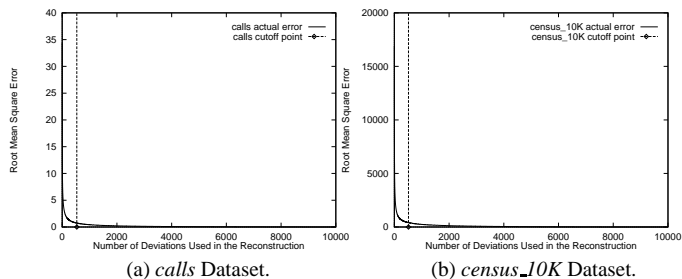


Figure 8. The reconstruction error when a varying number of deviations is used by the algorithm.

In every case, storing only a very small number of deviations is enough to dramatically decrease the error. In both cases of the real datasets we used, storing only the few largest deviations decreased the error by two orders of magnitude. As expected, at the other end of the spectrum, when the algorithm has knowledge of all the errors, it uses this information to achieve a perfect reconstruction.

It is interesting to note here the role that the *cutoff point* can play in this situation (see Section 5.2). In the graphs, the cutoff point, is marked with a vertical line, and it can be used to determine the point (and subsequently the number of values to materialize) at which the relative benefit of storing additional deviating values becomes negligible.

In Figure 9, we use a logarithmic scale for the y-axis to present the same graphs as before. In addition to the actual reconstruction error, we plot a theoretical upper-bound for the error, following the discussion in Section 4.3. Even though this bound is not very tight, it may still be useful for certain kind of applications. These graphs also enforce our argument about the *cutoff point*. It is clear that the *cutoff point* separates the initial region of dramatic decrease of the error from the plateau that follows. The added benefit of storing extra values from the latter region is quite small.



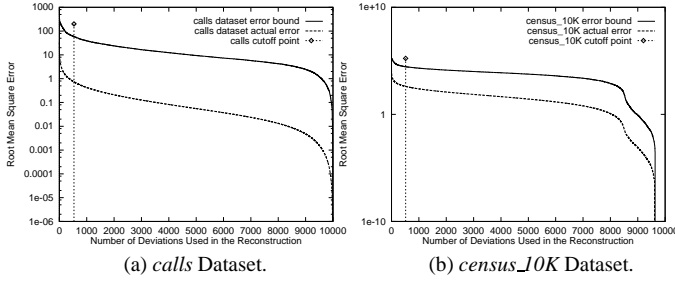


Figure 9. The actual reconstruction error and an upper bound when a varying number of deviations is used by the algorithm. The y-axis is plotted in logarithmic scale.

## 6.4 Mining Interesting Patterns

We evaluate the ability of the IPF algorithm to mine the underlying general structure of the data and report any deviations with the following experiments with synthetic and real datasets. Note, that the graphs we present involve datasets in two dimensions only, for illustration purposes.

We produced two datasets (namely  $g\_large\_s$  and  $g\_small\_s$ ) drawn from Gaussian distributions, one with large (Figure 10(a)), and one with small sigma value (Figure 10(b)). We then added some uniform noise on top of the

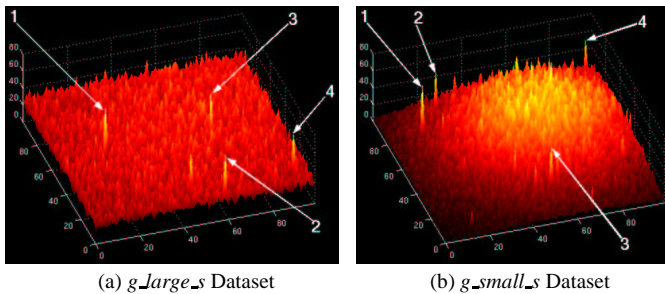


Figure 10. Identifying deviants using synthetic datasets.

Gaussian distributions. The algorithm captured the general trends of the data, and was able to report the values that deviate the most from the norm. The top-4 of these values are presented in Table 2. Manual inspection of the results reveals that these are indeed the predominant deviations in the datasets.

The third synthetic dataset ( $g\_mix$ ) we tested is a combination of two multidimensional Gaussian distributions with different mean and sigma values, and some noise on top (Figure 11(a)). Once more, the algorithm correctly identified the base distributions, and singled out the most significant deviating values (reported in Table 2). Note that the algorithm does not merely identify global phenomena, e.g., reporting the maximum value along a dimension. Instead, it

$g\_large\_s$		$g\_small\_s$		$g\_mix$	
cell	diff.	cell	diff.	cell	diff.
(50,21)	47.34	(74,14)	56.80	(60,67)	62.78
(2,59)	47.27	(89,25)	44.01	(48,6)	58.33
(48,68)	42.15	(19,54)	42.74	(28,2)	52.26
(11,93)	41.24	(95,95)	36.98	(23,60)	51.54

Table 2. The top-4 deviations reported for each of the three synthetic datasets. The metric used is the *difference* of the real value from the estimated.

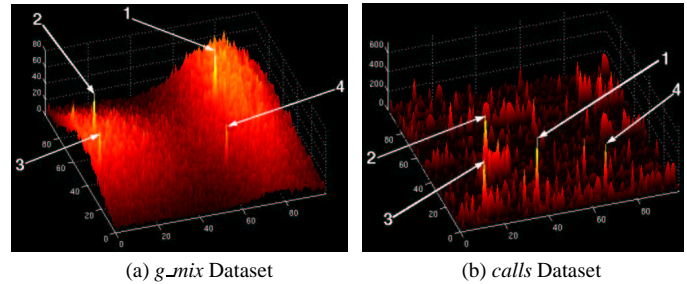


Figure 11. Identifying deviants using a synthetic and a real dataset.

takes into account the local neighborhood in which a particular value appears, and reports any incongruities therein.

In the following experiments, we instructed the algorithm to find the most deviating values in two of the real datasets, the *calls*, and the *census\_50K* dataset.

Figure 11(b) depicts the *calls* dataset along with the top-4 deviating values, which are also listed in Table 3. All the

<i>calls</i>	
cell	diff.
(7,37)	611.45
(38,24)	572.39
(7,13)	506.32
(7,68)	434.07

Table 3. The top-4 deviations reported for the *calls* datasets. The metric used is the *difference* of the real value from the estimated.

marked values are instances of unusually high volume of calls. This information is important to the analyst since it indicates exceptional behavior which can either be fraudulent, or mark special cases in the dataset.

The outcome of the second experiment, with *census\_50K*, cannot be graphically depicted, because the dataset is 5-dimensional. The attributes of the dataset are age, command of English, number of children, level of education, and in-

come. As expected, the above attributes are not independent. For example, the income tends to get larger with age, and when the level of education is higher. Nevertheless, there exist values that do not follow these patterns. Among the top deviations are a middle-aged person with high level of education who earns less than 20K, a person with a PhD degree who earns merely 3K, and a 24-year old who earns 200K. These are certainly results that deviate from the norm, and therefore are interesting.

Note that the algorithm is able to identify all the above results as interesting even though it has no domain knowledge, and it gets no user input.

## 7 Conclusion

In this paper we considered the problems of using just the aggregate information in order to provide approximate answers to queries and identify interesting values in multidimensional datasets. Each problem is of particular interest in the field of data analysis and approximate query answering respectively, especially since the volume of data stored in warehouses is huge. The techniques we discussed are based on the theoretic principle of *maximum entropy*. We also proposed an extended framework that allows the user to choose specific error guarantees for the reconstruction process. Finally, we presented a detailed performance study using both real and synthetic data, highlighting the applicability and benefits of the approach, as well as the efficiency of the proposed algorithms.

## Acknowledgements

We would like to thank Alberto Mendelzon and the anonymous referees for their comments that helped improve the content of this paper.

## References

- [1] S. Abad-Mota. Approximate Query Processing with Summary Tables in Statistical Databases. In *EDBT*, pages 499–515, Vienna, Austria, Mar. 1992.
- [2] S. Acharya, P. B. Gibbons, V. Poosala, and S. Ramaswamy. Join Synopses for Approximate Query Answering. In *ACM SIGMOD*, pages 275–286, Philadelphia, PA, USA, June 1999.
- [3] A. Arning, R. Agrawal, and P. Raghavan. A Linear Method for Deviation Detection in Large Databases. In *International Conference on Knowledge Discovery and Data Mining*, pages 164–169, Portland, OR, USA, Aug. 1996.
- [4] D. Barbará and M. Sullivan. Quasi-Cubes: Exploiting Approximations in Multidimensional Databases. *26(3):12–17*, 1997.
- [5] D. Barbará and X. Wu. Using Loglinear Models to Compress Datacubes. In *Web-Age Information Management*, pages 311–322, Shanghai, China, June 2000.
- [6] A. Berger, S. Pietra, and V. Pietra. A Maximum Entropy Approach to Natural Language Modelling. *Computational Linguistics*, *22(1)*, May 1996.
- [7] D. Bertsekas. *Constrained Optimization and Lagrange Multiplier Methods*. Academic Press, 1982.
- [8] Y. M. M. Bishop, S. E. Fienberg, and P. W. Holland. *Discrete Multivariate Analysis: Theory and Practice*. The MIT Press, 1975.
- [9] T. Cover and J. Thomas. *Elements of Information Theory*. Wiley, 1991.
- [10] W. E. Deming and F. F. Stephan. On a Least Square Adjustment of a Sampled Frequency Table When the Expected Marginal Totals Are Known. *Annals of Mathematical Statistics*, *11:427–444*, 1940.
- [11] C. Faloutsos, H. V. Jagadish, and N. Sidiropoulos. Recovering Information from Summary Data. *VLDB, Athens, Greece*, pages 36–45, Aug. 1997.
- [12] H. V. Jagadish, N. Koudas, S. Muthukrishnan, V. Poosala, K. C. Sevcik, and T. Suel. Optimal Histograms with Quality Guarantees. *VLDB*, pages 275–286, Aug. 1998.
- [13] H. V. Jagadish, I. S. Mumick, and A. Silberschatz. View Maintenance Issues in the Chronicle Data Model. *ACM PODS*, pages 113–124, June 1995.
- [14] E. M. Knorr and R. T. Ng. Algorithms for Mining Distance-Based Outliers in Large Datasets. In *VLDB*, pages 392–403, New York, NY, USA, Aug. 1998.
- [15] S. Kullback. *Information Theory and Statistics*. John Wiley and Sons, 1968.
- [16] F. Malvestuto. A Universal Scheme Approach to Statistical Databases Containing Homogeneous Summary Tables. *ACM TODS*, *18(4)*, pages 678–708, Dec. 1993.
- [17] G. S. Manku, S. Rajagopalan, and B. G. Lindsay. Random Sampling Techniques for Space Efficient Online Computation of Order Statistics of Large Datasets. In *ACM SIGMOD*, pages 251–262, Philadelphia, PA, USA, June 1999.
- [18] T. Palpanas and N. Koudas. Entropy Based Approximate Querying and Exploration of Datacubes. Technical Report CSRG-409, Dept. of Computer Science, University of Toronto, May 2000.
- [19] V. Poosala, Y. Ioannidis, P. Haas, and E. Shekita. Improved Histograms for Selectivity Estimation of Range Predicates. *ACM SIGMOD, Montreal Canada*, pages 294–305, June 1996.
- [20] S. Sarawagi, R. Agrawal, and N. Megiddo. Discovery-driven Exploration of OLAP Data Cubes. In *EDBT*, pages 168–182, Valencia, Spain, Mar. 1998.
- [21] J. Shanmugasundaram, U. M. Fayyad, and P. S. Bradley. Compressed Data Cubes for OLAP Aggregate Query Approximation on Continuous Dimensions. In *International Conference on Knowledge Discovery and Data Mining*, pages 223–232, San Diego, CA, USA, Aug. 1999.
- [22] A. Silberschatz and A. Tuzhilin. What Makes Patterns Interesting in Knowledge Discovery Systems. *IEEE Transactions on Knowledge and Data Engineering*, *Vol. 8, No 6.*, pages 970–974, Dec. 1996.
- [23] J. S. Vitter, M. Wang, and B. Iyer. Data Cube Approximation and Histograms via Wavelets. In *ACM CIKM*, pages 96–104, Washington, DC, USA, 1998.