

Differentially Private Multi-Dimensional Analysis via Double-Sided Bounding and Prefix-Sum Cube

Yufei Wang[†], Xiang Cheng^{†(✉)}, Pengfei Zhang[‡], Anxing Wei[‡]

[†]State Key Laboratory of Networking and Switching Technology,
Beijing University of Posts and Telecommunications, Beijing, China,

[‡]Anhui University of Science and Technology, Huainan, China,

[†]{wyfs4321, chengxiang}@bupt.edu.cn,

[‡]{zpf, axwei}@aust.edu.cn

Abstract—Multi-dimensional analysis is a powerful tool for online analytical processing. In this paper, we focus on the process of answering multi-dimensional analytical queries over data cubes, each of which consists of a collection of cuboids, while satisfying differential privacy (DP). Existing works fail to efficiently process multi dimensional analytic queries under DP due to the high sensitivity of sum query and noise aggregation: constructing a base cuboid requires the data curator to answer a workload of linear sum queries under DP in advance, whose sensitivity will result in a large amount of DP noise, and the DP noise will finally be aggregated when constructing the remaining cuboids. To this end, we present a Differentially PRivate Multi-dimensional Analytic Approach (PRIMA). In PRIMA, we propose a Bounding-and-Debiasing Sum Query Processing Method (BAND) which first reduces the sensitivity of sum queries by bounding both the maximum and minimum contribution of each record in the data table and then debiases the data table based on a symmetric sparse vector technology. Moreover, we propose a Hypothesis Testing based Prefix Sum Computing Method (SCOPE) to compute a base prefix-sum cuboid based on hypothesis testing. By employing the base prefix-sum cuboid, any remaining cuboid can be constructed with constant pieces of DP noise aggregated. We conduct experiments on both real-world and synthetic datasets. Experimental results confirm the effectiveness of PRIMA over existing works.

I. INTRODUCTION

Multi-dimensional analysis has become an essential element of online analytical processing. By constructing a data cube that consists of a collection of cuboids in advance, the multi-dimensional analytical queries can be answered over massive amounts of data, which are stored in a data warehouse, quickly and efficiently. Specifically, the data curator first computes a base cuboid based on the data table via pre-querying. By aggregating the base cuboid in different dimensions, the data curator then constructs the remaining cuboids. Finally, any multi-dimensional analytical queries can be answered over the constructed data cube. However, since the data table usually contains private data, constructing data cubes based on such data may endanger the privacy of individuals. Therefore, the data curator must provide a rigorous privacy guarantee on how multi-dimensional analytical queries are processed.

Differential Privacy (DP) has been increasingly accepted as the state-of-the-art standard for protecting individual privacy in the central setting where a trusted data curator processes the sensitive data of all users, and has been studied extensively in

both literature [1]–[8] and industry [9]–[13]. DP guarantees that the presence or absence of any particular individual’s record has a negligible impact on the likelihood that a particular result is returned to a query. Thus, an adversary cannot make meaningful inferences about any individual’s record value, or even whether the record is present.

In this paper, we focus on the process of answering multi-dimensional analytical queries over data cubes under DP. Several existing differentially private query processing approaches can be used to give access to sensitive data with privacy guarantee. These works reduce the sensitivity of sum queries, which are essential operations in multi-dimensional analysis, by bounding users’ maximum contribution. However, all these works face a bias-variance trade-off: a loose bounding threshold cannot reduce the DP noise caused by the high sensitivity; while a tight bounding threshold will inevitably introduce systematic bias into the query results. In addition, to limit the DP noise added to the cuboids when constructing the remaining cuboids, Ding et al. [14] propose a general Noise Control Framework, which formulates an optimization problem and minimizes the utility loss. However, this framework cannot cope well with the problem of noise aggregation, i.e., the problem that the DP noise in each cell of the base cuboid will be aggregated when constructing the remaining cuboids. As the domains of attributes increase, the numbers of noisy cells aggregated in the Noise Control Framework will grow in polynomial degrees, resulting in massive DP noise aggregated in the constructed data cube.

Both the high sensitivity in answering sum queries and the noise aggregation in constructing data cubes significantly affect the accuracy of multi-dimensional analysis. To deal with these problem, we propose a Differentially PRivate Multi-dimensional Analytic Approach (PRIMA). Its main idea is that answers sum queries accurately in a double-sided bounding manner and processes multi-dimensional analysis based on the prefix-sum cube which is constructed via hypothesis testing.

In PRIMA, we first put forward a Bounding-AND-Debiasing Sum Query Processing Method (BAND) for answering sum queries under DP while alleviating the systematic bias. Different from the single-sided methods that delete or bound the records whose measures’ values are greater than a threshold, BAND carefully chooses a bounding range $[\Phi - \tau, \Phi + \tau]$

and rewrites the data table. In BAND, the data curator first bounds the local sensitivity of sum queries by rewriting the data table via a double-sided bounding mechanism, where the data curator formulates optimization problems for choosing the optimal bounding axis Φ and radio τ . After that, the data curator checks whether the lower threshold of the bounding range $\Phi - \tau$ is less than 0. If it is, the data curator further post-processing the data table again based on a symmetric sparse vector technology, which is a bottom-up mechanism that traverses the data table and outputs a new lower threshold of bounding range τ' for debiasing the rewritten data table. Specifically, we claim that although BAND is not the first method that introduces the idea of double-bounding into query processing, it is exactly the first method that take the lower threshold as an important parameter and mitigates the dilemma between systematic bias and variance by rewriting the data table with the carefully computed lower threshold.

Finally, for constructing the remaining cuboids accurately, the data curator alleviates noise aggregation by constructing the prefix-sum cubes [15]. However, computing the prefix sums require multiple accesses to any single record in the data table, which will introduce massive DP noise. Therefore, we design a HypotheSis Testing based Prefix Sum **COM**puting **M**ethod (SCOPE). Compared with computing the prefix sums directly under DP, SCOPE avoids multiple accesses to the records and computes the prefix sums based on hypothesis testing. Specifically, in SCOPE, the data curator first acquires a noisy sum vector by accessing each record once under DP. The data curator then predicts the distribution of the noisy sum vector based on hypothesis testing. Next, the data curator computes the prefix sums based on both the noisy sum vector and its predicted distribution. Finally, the data curator boosts the accuracy of the prefix sums by enforcing consistency. Based on the prefix sums, the data curator can compute a prefix sum base cuboid and constructs the remaining cuboids without noise aggregation.

The main contributions of this work are summarized as follows:

- We focus on the process of answering multi-dimensional analytical queries over data cubes under DP and present PRIMA which can bound the sensitivity of sum query while minimizing the systematic bias and alleviate noise aggregation when constructing data cubes under DP.
- We put forward a **B**ounding-**AND**-**D**ebiasing Sum Query Processing Method (BAND) which consists of a double-sided bounding mechanism and a symmetric sparse vector technology. By rewriting the table via the double-sided bounding mechanism and reducing the potential systematic bias via the symmetric sparse vector technology, BAND can answer sum queries under DP accurately.
- We design a HypotheSis Testing based Prefix Sum **COM**puting **M**ethod (SCOPE). By computing prefix sums based on hypothesis testing and boosting the accuracy of prefix sums by enforcing consistency, SCOPE can improve the utility of the noisy prefix-sum cube and alleviate the problem of noisy aggregation.

- To evaluate the performance of PRIMA, we conduct extensive experiments on both real-world and synthesized datasets. The experimental results demonstrate that PRIMA outperforms the baseline approaches.

II. RELATED WORK

Since DP [16] was introduced, many works have been proposed for preserving individuals' privacy under DP during answering queries. A straightforward way for differentially private query processing is the interactive strategy. In this strategy, when a new query comes, the data curator computes the result directly and returns the noisy version to the analyst after adding DP noise. Adopting this strategy, McSherry et al. [17] firstly propose Privacy Integrated Queries (PINQ) to process counting queries under DP. For supporting general equijoins without scaling up the noise magnitude, Proserpio et al. [18] extend PINQ by assigning a weight to each row in the data table and put forward Weighted PINQ (wPINQ), which scales down the weights of rows in a join to ensure an overall sensitivity of one. However, both PINQ and wPINQ are not compatible with standard databases since they introduce new operators that do not exist in standard SQL. Moreover, the high sensitivity of queries will result in poor utility. To process SQL queries over standard databases under DP, Zhang et al. [19] propose ϵ KTELO which is a programming framework and system that aids programmers in developing differentially private programs. To bound the sensitivity of queries, Johnson et al. [20] propose a sound approximation of local sensitivity [21] called elastic sensitivity, which is extended from the idea of smooth sensitivity [22]. Based on the elastic sensitivity, the authors implement an end-to-end DP system for SQL queries based on elastic sensitivity called FLEX. Targeting the same problem, Ge et al. [23] present APEx, which allows data analysts to pose adaptively chosen sequences of queries along with required accuracy bounds.

However, the approaches mentioned above are all designed for processing linear counting queries over a single relation with different predicates, which have a relatively small sensitivity. For answering the queries with high sensitivity, a series of works [24]–[27] bound the sensitivity by restricting the maximum contribution of each record, which rewrite the records whose attributes' values are more than a threshold τ before adding DP noise. Specifically, Kotsogiannis et al. [24] present PrivateSQL, an end-to-end differentially private relational database system that learns truncation threshold by applying the sparse vector technique [28]. Dong et al. [25] propose a truncation mechanism called R2T, which can adaptively choose a threshold that is provably close to the optimal one tuned for each particular query/dataset. Qiu et al. [26] present a mechanism for answering a set of sum queries. However, all these truncation-based works will introduce a significant systematic bias into the result. Besides, there are also a series of work [27], [29], try to bound the both maximum and minimum contribution of each user. However, these works primarily emphasize the choosing of the upper threshold while trivially setting the lower threshold as the

minimum value within the measure's domain. As a result, these works are essentially single-sided mechanisms.

Since the interactive-based approaches can only process a limited number of count queries under a fixed privacy budget. A series of approaches [14], [30]–[33] based on the non-interactive strategy are proposed. In these approaches, the data curator first constructs a differentially private data model that contains certain statistical information about the data table via a set of queries. Any subsequent queries can be then processed over the differentially private data model without spending a privacy budget. In the non-interactive strategy, noise aggregation is a major problem that leads to poor utility. To answer queries with both privacy and utility guarantees, Hay et al. [30] propose an approach that is able to deal with noise aggregation in single-dimensional range queries via hierarchical intervals. Qardaji et al. [31] optimize the hierarchical intervals by choosing a proper branching factor. Targeting the same problem, Xiao et al. [32] develop an approach called Privelet that provides accurate answers by using Haar wavelet. However, the data models in these approaches are limited to single-dimensional analysis. To tackle this problem, Ding et al. [14] present a general Noise Control Framework which carefully selects a subset of cells in a cuboid to inject Laplace noise and compute the rest cells from them. However, when the domains of attributes are large, the noise aggregation in Noise Control Framework will also cause a poor utility of query results.

III. PRELIMINARIES

A. Differential Privacy (DP)

DP requires the outputs of algorithms to be approximately the same even if any individual's record in the database is added or removed. Thus, the presence or absence of any individual's record has a statistically negligible effect on the outputs.

Formally, DP considers an instance data table T of relational schema $R(\mathbb{A}, M)$. The distance between two tables $T, T' \in D^n$ is $d(T, T') = |\{i | t_i \neq t'_i\}|$. Two tables T, T' are neighboring tables if $d(T, T') = 1$.

Definition 1 (Differential Privacy): A private algorithm \mathcal{A} preserves ϵ -differential privacy is for any pair of neighboring tables T, T' , and for all sets \mathcal{O} of possible outputs:

$$Pr[\mathcal{A}(T) \in \mathcal{O}] \leq e^\epsilon Pr[\mathcal{A}(T') \in \mathcal{O}] + \delta,$$

where the probability is taken over the randomness of \mathcal{A} .

Dwork et al. [16] propose the Laplace mechanism to achieve DP. For a function whose outputs are real, they prove that DP can be achieved by adding noise drawn randomly from Laplace distribution. The Laplace distribution with magnitude γ follows probability density function $Pr[Lap(\gamma) = x] = \frac{1}{2\gamma} e^{-|x|/\gamma}$, where $\gamma = \frac{S}{\epsilon}$ is determined by the desired privacy budget ϵ and the sensitivity S of the function. In particular, the sensitivity S is used to measure the maximum change in the outputs of a function when any individual's record in the data table is changed.

One measure of sensitivity is the global sensitivity [16], which is the maximum difference in the query's result on any two neighboring tables. However, for many functions, the global sensitivity yields high noise, which leads to a poor utility of query results. To this end, a local measure of sensitivity called local sensitivity [21] is proposed, which is the maximum difference between the query's results on a given table and any neighboring table of it.

Definition 2 (Local Sensitivity): For an algorithm \mathcal{A} and an instance data table T of relational schema $R(\mathbb{A}, M)$, the local sensitivity of \mathcal{A} at T is

$$LS_{\mathcal{A}} = \max_{y: d(T, T')=1} \|\mathcal{A}(T) - \mathcal{A}(T')\|.$$

Since the local sensitivity is a property of the given table rather than the set of all possible tables, it is much often much smaller than global sensitivity. However, the local sensitivity is not safe enough and the corresponding noise magnitude may reveals information about the database. [22]. In this paper, for utilizing local sensitivity safely, we choose to release the noisy query via Propose-Test-Release (PTR) [16].

Definition 3 ((ϵ, δ) -Propose-Test-Release): Let \mathcal{A} be the function of interest, which we wish to privately query on datasets T . A (ϵ, δ) -Propose-Test-Release consists of the following three steps:

1. Propose an upper bound β on the local sensitivity.
2. Compute the distance from T to the nearest dataset T' such that $\Delta_{LS}(T') \geq \beta$, name it γ . Distance is measured in terms of how many records must be changed to get from T to T' .
3. Compute $\hat{\gamma} = \gamma + Lap(1/\epsilon)$.
4. If $\hat{\gamma} \leq \ln 1/\delta/\epsilon$, return \perp .
5. If $\hat{\gamma} > \ln 1/\delta/\epsilon$, return $\mathcal{A}(T) + Lap(\beta/\epsilon)$.

In our work, all queries are processed via PTR to avoid the potential privacy leakage.

To support multiple differentially private queries, the sequential and parallel compositions of DP are extensively used.

Theorem 1 (Sequential Composition): Given k random algorithms $\mathcal{A}_i (1 \leq i \leq k)$ that access the same record t , each of which satisfies (ϵ_i, δ_i) -DP, then the combination of their outputs satisfies (ϵ, δ) -DP, where $\epsilon = \left(\sum_{i=1}^k \epsilon_i\right)$, $\delta = \left(\sum_{i=1}^k \delta_i\right)$.

Theorem 2 (Parallel Composition): Given k random algorithms $\mathcal{A}_i (1 \leq i \leq k)$ that access different records $t_i (1 \leq i \leq n)$, each of which satisfies (ϵ_i, δ_i) -DP, then the combination of their outputs satisfies (ϵ, δ) -DP, where $\epsilon = \max(\epsilon_1, \dots, \epsilon_k)$, $\delta = \max(\delta_1, \dots, \delta_k)$.

B. Data Cube

Data cube is a combination of cuboids. For a data table T which is an instance of relational schema $R(\mathbb{A}, M)$, a data cuboid can be viewed as the projection of T on a subset of attributes $\mathbb{B} \subseteq \mathbb{A}$, producing a multi-dimensional array. In the cuboid, each cell aggregates the measure M of records in T that match certain dimensions. Given a data cube among d

attributes in T , there are $\binom{d}{\lambda}$ λ -dimensional (' λ -D') cuboids ($0 \leq \lambda \leq d$) in a data cube, representing all possible combinations of attributes in T . In these cuboids, the only d -D cuboid is called a base cuboid. Any other cuboids can be acquired by aggregating the base cuboid. In addition, for any two cuboids C_1 and C_2 , if C_1 can be constructed from C_2 by directly aggregating one or more attributes in C_2 , then C_1 is said to be an ancestor of C_2 , and C_2 is a descendant of C_1 . We show an example of constructing a data cube based on a table T in Fig. 1.

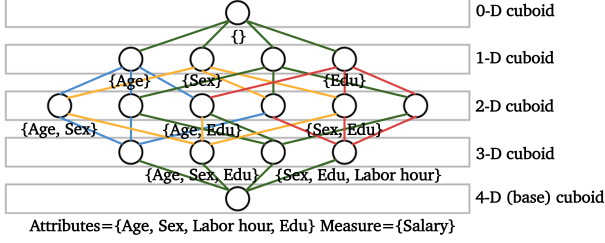


Fig. 1. Data Cube

Example 1: Suppose there is a measure *Salary* and a subset of attributes $\mathbb{A} = \{a_1, a_2, a_3, a_4\}$, where a_1, a_2, a_3 , and a_4 refer to age, gender, labor hours, and educational level, respectively. To analyze the relationship of workers' salary with their age, gender, labor hours, and educational level, the data curator can construct a data cube among these four attributes which aggregates the *Salary*. The data cube contains $\binom{4}{0} + \binom{4}{1} + \binom{4}{2} + \binom{4}{3} + \binom{4}{4} = 16$ cuboids, in which the 4-D cuboid $C_{a_1 a_2 a_3 a_4}$ is the base cuboid. In particular, for the 2-D cuboid $C_{a_1 a_2}$, the 3-D cuboid $C_{a_1 a_2 a_4}$ is the descendant of $C_{a_1 a_2}$ while $C_{a_2 a_3 a_4}$ is not, since $C_{a_1 a_2}$ cannot be constructed by aggregating $C_{a_2 a_3 a_4}$.

C. Prefix-Sum Cube

Given λ attributes, a prefix-sum cuboid [34] is a λ -D array that has the same structure as the data cuboid among these attributes. The cells in the prefix-sum cuboid contain a kind of auxiliary information called prefix sum [34], which can be computed from the distributions among λ attributes. Similarly, based on the base prefix-sum cuboid, the remaining cuboids can be constructed through the computation between two cells in the base prefix-sum cuboid. Since each cell in the prefix-sum cuboid contains a piece of DP noise, which is introduced when computing the base prefix-sum cuboid under DP, there are only constant pieces of noise being aggregated, irrespective of the domain of the aggregated attributes.

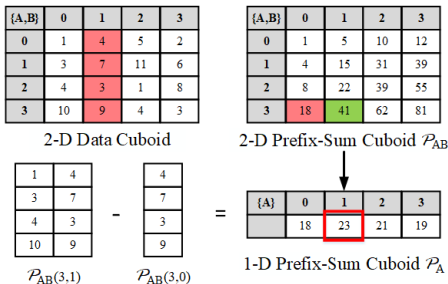


Fig. 2. Computing a cell in an ancestor cuboid based on \mathcal{P}

Suppose there is a λ -D cuboid A of size $dom(A_1) \times dom(A_2) \times \dots \times dom(A_\lambda)$. Prefix-sum cuboid \mathcal{P} is another λ -D array which is used to store various precomputed prefix sums of A . For all $0 \leq x_j < dom(A_j)$ and $j \in [0, \lambda - 1]$,

$$\begin{aligned} \mathcal{P}[x_1, x_2, \dots, x_\lambda] &= Sum(0 : x_1, 0 : x_2, \dots, 0 : x_\lambda) \\ &= \sum_{i_1=0}^{x_1} \dots \sum_{i_\lambda=0}^{x_\lambda} A[i_1, \dots, i_\lambda]. \end{aligned} \quad (1)$$

Theorem 3 provides how a cell in an ancestor cuboid can be computed from up to 2^λ appropriate cells in the prefix sum base cuboid \mathcal{P} . The left-hand side of Eq. 2 specifies a cell to be computed. The right-hand side of Eq. 2 consists of 2^λ additive terms, each of which is from a cell of \mathcal{P} with a sign '+' or '-' defined by the product of all $s(i)$. For notational convenience, let $\mathcal{P}[x_1, x_2, \dots, x_\lambda] = 0$ if $x_j = -1$ for some $j \in [0, \lambda - 1]$.

Theorem 3: For all $j \in [0, \lambda - 1]$, let

$$s(j) = \begin{cases} 1, & \text{if } x_j = h_j \\ -1, & \text{if } x_j = l_j - 1 \end{cases}.$$

Then, for all $j \in [0, \lambda - 1]$,

$$\begin{aligned} Sum(l_1 : h_1, \dots, l_d : h_d) &= \\ \sum_{\forall x_j \in \{l_j - 1, h_j\}} \left\{ \left[\prod_{i=1}^d s(i) * \mathcal{P}[x_1, x_2, \dots, x_d] \right] \right\}. \end{aligned} \quad (2)$$

We show an example of how to compute a cell in an ancestor cuboid based on a prefix-sum cuboid in Example 2.

Example 2: Fig. 2 shows an example of a prefix-sum cuboid \mathcal{P} with $\lambda = 2$, $dom(M_A) = dom(M_B) = 4$, and the geometrical explanation for the process of computing the value of a certain cell in the 1-D prefix-sum cuboid $\mathcal{P}_A(1) = Sum(0 : 1, 0 : 3) - Sum(0 : 0, 0 : 3) = \mathcal{P}_{AB}(3, 1) - \mathcal{P}_{AB}(3, 0)$.

D. Multi-Dimensional Analysis

Given a table T which is an instance of relational schema $R(\mathbb{A}, M)$, where \mathbb{A} is a set of d attributes and M is a measure, we focus on the following class of multi-dimensional analytical queries:

SELECT SUM(M) FROM T WHERE

$A_1 \in [l_1, r_1]$ **AND** ... **AND** $A_d \in [l_d, r_d]$.

We denote the answer of a multi-dimensional analytical query q over data table T as \mathcal{A}_q^T .

E. Linear Sum Queries and Query Matrix

The base cuboid can be constructed by answering a set of linear sum queries. For a data table T which is an instance of relational schema $R(\mathbb{A}, M)$, we denote by $dom(\mathbb{A})$ the cross-product of the domains of attributes in \mathbb{A} . The data curator chooses a set of attributes $\mathbb{B} \subseteq \mathbb{A}$ relevant to the base cuboid.

For example, if the data analyst is interested in a subset of two dimensional multi-dimensional analytic queries over attributes A_1 and A_2 , they would set $\mathbb{B} = \{A_1, A_2\}$. The data curator then forms a sum vector \mathbf{x} with one entry for each element of $\text{dom}(\mathbb{B})$. For simplicity, we assume $\text{dom}(\mathbb{B}) = \{1, 2, \dots, n\}$ and for each $i \in \text{dom}(\mathbb{B})$, x_i is the sum of records' measures that satisfy the constraints that $|\mathbb{B}| = i$. We represent \mathbf{x} as a column vector of sums: $\mathbf{x} = [x_1, x_2, \dots, x_n]^t$. A linear query computes a linear combination of the counts in \mathbf{x} .

Definition 4 (Linear Sum Query): A linear query is a length- n row vector $\mathbf{q} = [q_1, \dots, q_n]$ with each $q_i \in \mathbb{R}$. The answer to a linear sum query \mathbf{q} on \mathbf{x} is the vector product $\mathbf{q}\mathbf{x} = q_1x_1 + \dots + q_nx_n$.

A series of linear sum queries can be organized into the rows of a pre-query matrix [4], [5] whose definition is given as follows:

Definition 5 (Pre-Query Matrix): A pre-query matrix is a collection of m linear sum queries, arranged by rows to form an $m \times n$ matrix.

If \mathbf{Q} is an $m \times n$ pre-query matrix, the query answer for \mathbf{Q} is a length m column vector of query results, which can be computed as the matrix product $\mathbf{Q}\mathbf{x}$. The sensitivity of \mathbf{Q} [4], [5] can be computed as follows.

Proposition 1 (Pre-Query Matrix Sensitivity): The sensitivity of matrix \mathbf{Q} , denoted as $\Delta_{\mathbf{Q}}$, is

$$\Delta_{\mathbf{Q}} = \max_{\|\mathbf{x}-\mathbf{x}'\|=1} \|\mathbf{Q}\mathbf{x} - \mathbf{Q}\mathbf{x}'\|_1 = \max_j \sum_{i=1}^n |q_{ij}|.$$

Thus the sensitivity of a pre-query matrix is the maximum L_1 norm of a column.

We show the pre-query matrices for computing a data cuboid and a prefix-sum cuboid in Example 3.

Example 3: Fig. 3 shows two query matrices. Given a subset of attributes $\mathbb{B} = A, B$, where $\text{dom}(M_A) = \text{dom}(M_B) = 4$, $\text{dom}(\mathbb{B}) = \text{dom}(M_A) \times \text{dom}(M_B) = 16$. \mathbf{I}_{16} is the pre-query matrix for computing a 2-D data cuboid, and \mathbf{P}_{16} is the pre-query matrix for computing a 2-D prefix-sum cuboid. Specifically, for a data table where the sensitivity of the sum query is LS_f , \mathbf{I}_{16} is an identity matrix of size 16×16 . This matrix consists of 16 queries and each row of \mathbf{I}_{16} represents a linear sum query that asks for the value of a cell in the data cuboid. According to Proposition 1, the sensitivity of each linear sum query in \mathbf{I}_{16} is LS_f . \mathbf{P}_{16} is a lower triangular matrix of size 16×16 , which consists of 16 queries, and each row of \mathbf{P}_{16} represents a linear sum query that asks for the value of a cell in the prefix-sum cuboid. According to Proposition 1, the sensitivity of each linear sum query in \mathbf{P}_{16} is $16 \cdot LS_f$.

$$\begin{aligned} \text{(a) } \mathbf{I}_{16} &= \begin{bmatrix} LS_f & 0 & \dots & 0 \\ 0 & LS_f & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & LS_f \end{bmatrix} & \text{(b) } \mathbf{P}_{16} &= \begin{bmatrix} LS_f & 0 & \dots & 0 \\ LS_f & LS_f & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ LS_f & LS_f & \dots & LS_f \end{bmatrix} \end{aligned}$$

Fig. 3. Two Query Matrices

IV. BASELINE APPROACHES

Although there is no work that can be directly used to answer multi-dimensional analytical queries over data cubes under DP, in this section, we provide three baseline approaches, two of which are combinations of related works.

A. BASE

For answering multi-dimensional analytical queries over data cubes under DP, one straightforward approach is to compute a base cuboid from the data table under DP and construct the other cuboids from the noisy base cuboid. We denote this approach as BASE.

Limitations. Obviously, BASE will introduce excessive DP noise when answering the linear sum queries and such noise will be aggregated when constructing the remaining cuboids, leading to poor utility of results.

B. SVTC

Sparse Vector Technology based Data Cube Constructing (SVTC) is a combination of PrivateSQL [24] and the Noise Control Framework [14]. In SVTC, for a data table $T = \{t_i | 1 \leq i \leq n\}$, the data curator first truncates the data table T via SVT [28]. For choosing a proper truncation threshold τ , the data curator first assigns a parameter τ_s and a candidate set of truncation threshold $\mathcal{B} = \{b_i = i | 1 \leq i \leq \text{dom}(M)\}$ in advance. Then, the data curator traverses \mathcal{B} in a bottom-up manner. For each b_i , the data curator truncates the data table and outputs a truncated data table T^{b_i} where any record t_i that satisfies $|M^{t_i}| > b_i$ is deleted. Based on T^{b_i} , the data curator answers a stream of queries $Q = \{q_i | 1 \leq i \leq \text{dom}(M)\}$ under DP, where

$$q_i : \text{SELECT SUM}(M) \text{ FROM } T^{b_i} \text{ WHERE } M = i,$$

and the noisy result is $\mathcal{A}_{q_i}^{T^{b_i}} + \text{Lap}(\frac{b_i}{\epsilon})$. During the process of traversing \mathcal{B} in a bottom-up manner, for each b_i , the data curator can acquire a query result $\mathcal{A}_{q_i}^{T^{b_i}} = \sum_{q_i \in Q} (\mathcal{A}_{q_i}^{T^{b_i}} + \text{Lap}(\frac{b_i}{\epsilon}))$. Since a larger b_i indicates less records in T are truncated, there is an increment in $\mathcal{A}_{q_i}^{T^{b_i}}$ compared with $\mathcal{A}_{q_i}^{T^{b_i-1}}$. For the first $\mathcal{A}_{q_i}^{T^{b_i}}$ whose increment is less than the threshold τ_s , the data curator identifies the corresponding b_i as the optimal truncation threshold τ . After getting τ , the data curator traverses every record in the data table and deletes the records whose measures' values exceed τ . After truncating the data table via SVT, the data curator chooses an initial set \mathcal{L}_{pre} of cuboids according to the Noise Control Framework and computes these cuboids by answering the linear sum queries in the pre-query matrix under DP. Based on \mathcal{L}_{pre} , the data curator constructs the remaining cuboids.

Limitations. SVTC attempts to improve the utility of the sum query results by truncating the data table. However, the truncation operator will introduce a systematic error into the result. Moreover, since SVT requires a manually assigned parameter as the metric for choosing the optimal truncation threshold τ , the effectiveness of SVTC is unstable. Furthermore, although the Noise Control Framework can bound the

maximal noise over all cuboids, when the domain of attributes $\text{dom}(A)$ is large (e.g., $\text{dom}(A) > 10$), the utility loss caused by noise aggregation is still enormous.

C. R2TC

Race-to-the-Top based Data Cube Constructing (R2TC) is a combination of R2T [25] and the Noise Control Framework [14]. In R2TC, the data curator first chooses an initial set \mathcal{L}_{pre} of cuboids according to the Noise Control Framework. The data curator then initials a candidate set $\mathcal{B} = \{b_i = 2^i | 1 \leq i \leq \log_2 \text{dom}(M)\}$ and a series of truncated data table $\{T^{b_i} | 1 \leq i \leq \log_2 \text{dom}(M)\}$. For a certain linear sum query q in the pre-query matrix, the data curator computes a list of noisy results $\mathcal{A}_q^{T^{b_i}} = \mathcal{A}_q^{T^{b_i}} + \text{Lap}(B_i, \epsilon') - \ln(\frac{\log(\text{dom}(M))}{\beta}) \cdot \frac{B_i}{\epsilon'}$, where $\epsilon' = \frac{\epsilon}{\log \text{dom}(M)}$ is the privacy budget allocated for q and β is an assigned probability parameter. The data curator chooses $\max(\mathcal{A}_q^{T^{b_i}})$ as the final result of q . After answering all linear sum queries in the pre-query matrix under DP, the data curator computes the initial set \mathcal{L}_{pre} of cuboids and constructs remaining cuboids based on \mathcal{L}_{pre} .

Limitations. Although R2TC does not need to manually assign τ_s like SVTC, it still requires a probability parameter β given in advance, which has a direct impact on the utility. Moreover, for choosing the optimal threshold, R2TC must shift the result down by an amount that equals to the scale of the noise. As $\text{dom}(M)$ increases, it will make the effectiveness of R2TC drop rapidly.

V. SOLUTION

In this section, we present PRIMA for answering multi-dimensional analytic queries over data cubes under DP.

A. Overview

As shown in Section II, the existing approaches fail to reduce the systematic bias while bounding the high sensitivity of sum queries and alleviate the problem of noise aggregation. To this end, we present a Differentially **PR**ivate **Multi-dimensional Analytic Approach** (PRIMA). PRIMA contains two key building blocks: **B**ounding-**AND**-**D**ebiasing Sum Query Processing Method (BAND) and **H**ypothe**S**is Testing based Prefix Sum **Co**mputing **M**ethod (SCOPE). We show the overview of PRIMA in Fig. 4 and introduce the main phases of PRIMA as follows:

1. Data Rewriting Phase. Given a class of multi-dimensional analytical queries that associate with d attributes and a measure M . The data curator first determines the optimal bounding axis Φ and bounding radius τ via a double-bounding mechanism, which forms the bounding range $[\Phi - \tau, \Phi + \tau]$. Then, the data curator traverses the data table and clips the measure M of each record whose value is beyond the bounding range $[\Phi - \tau, \Phi + \tau]$. The details of the double-bounding mechanism are shown in Section V-B.

2. Data Debiasing Phase. After rewriting the data table, if the lower threshold of bounding range $\Phi - \tau < 0$, the data curator determines a new lower threshold of bounding range τ' for the bounded values of M via a symmetric sparse vector technology and clips the measure M of each record whose

value is beyond the new bounding range $[\tau', \Phi + \tau]$. The details of the symmetric sparse vector technology are described in Section V-C.

3. Pre-querying Phase. The data curator then organizes a workload of linear sum queries into a pre-query matrix to acquire the noisy sum vectors for constructing the base prefix-sum cuboid. Since the sum queries in the pre-query matrix access different records in the table, according to the parallel composition, a piece of privacy budget ϵ can be shared when answering these queries. To prevent the privacy leakage, all queries are processed via PTR.

4. Cuboid Constructing Phase. Based on the noisy sum vectors, the data curator computes the base prefix-sum cuboid via SCOPE and constructs the remaining prefix-sum cuboids by aggregating the base prefix-sum cuboid in different dimensions. The details of SCOPE are given in Section V-D.

5. Query Processing Phase. Finally, based on the constructed data cube, the multi-dimensional analytical queries can be processed privately and efficiently.

In the following, we first introduce the double-sided bounding mechanism and the symmetric sparse vector technology in Section V-B and Section V-C, respectively. We then give the details of SCOPE in Section V-D. Finally, we show the privacy guarantee and utility analysis of PRIMA in Section V-E.

B. Double-sided Bounding Mechanism

The high sensitivity of the sum query is the underlying cause of the poor result utility. A series of works [24], [25] bound the sensitivity of sum queries by bounding users' maximum contribution. Given records $\{t_i | 1 \leq i \leq n\}$ in a data table T , if each record t_i has a measure whose measure's value is $|M^{t_i}|$, the data curator deletes the record t_i that $|M^{t_i}| > \tau$ [24], [25] or rewrites its measure $|M^{t_i}| = \tau$, where τ is a carefully chosen threshold. However, only bounding users' maximum contribution introduces systematic bias to the data table, leading to the poor utility of the query results.

In BAND, for answering sum queries under DP while alleviating the systematic bias, we put forward a double-sided bounding mechanism. Rather than deleting or rewriting the records whose measures' values are larger than a specified threshold, the data curator carefully chooses a bounding range $[\Phi - \tau, \Phi + \tau]$ and rewrites $|M^{t_i}|$ such that any $|M^{t_i}|$ beyond the bounding range only contributes $\Phi - \tau$ or $\Phi + \tau$ to the query result. In particular, by bounding the records whose measures' values larger than $\Phi + \tau$, the sensitivity of sum queries is bounded and less DP noise is added; by bounding the records whose measures' values smaller than $\Phi - \tau$, the systematic bias in query results is alleviated.

The double-sided bounding mechanism takes a data table T which is an instance of relational schema $R(\mathbb{A}, M)$ as input and outputs a rewritten data table T' . In the double-sided bounding mechanism, the data curator first computes the optimal bounding axis Φ by traversing the records in the data table T , then determines the optimal bounding radius τ by formulating an optimization problem. Φ and τ form a bounding range $[\Phi - \tau, \Phi + \tau]$.

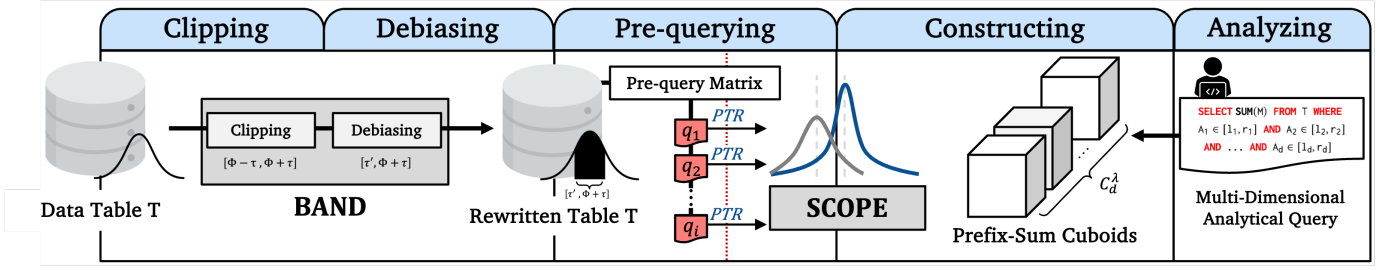


Fig. 4. The overview of PRIMA

After determining the bounding range, the data curator clips the records in data table T such that the contribution of record t_i is $\max \min |M^{t_i}|, \Phi - \tau, \Phi + \tau$. We show the details of how to determine the optimal Φ and τ as follows.

1) *Bounding Axis Φ* : Given the original records $\{t_i | 1 \leq i \leq n\}$ and the bounded records $\{t'_i | 1 \leq i \leq n\}$, in order to reduce the systematic bias in the rewritten data table, we aim to find the optimal Φ that makes the mathematic expectation of the bounded records' measures converge to the mathematic expectation of the original records' measures. We show the optimal Φ in Eq. 3.

Theorem 4: The optimal bounding axis is

$$\Phi = \frac{\sum_{i=1}^n |M^{t_i}|}{n}. \quad (3)$$

Due to space limitation, we put the proof of Theorem 4 in Appendix A.

2) *Bounding Radius τ* : To determine an optimal bounding radius τ , we consider two types of errors. Specifically, we denote the error caused by adding Laplace noise to the linear sum query results as the DP noise and denote the error caused by the deviation between the original values of the records and the bounded values of the records as the data distortion. A small τ indicates there are more values of records bounded, resulting in a large data distortion. While a large τ indicates the sum query has a large sensitivity, leading to a large DP noise. Therefore, we formalize an optimization problem. Before giving the details of the optimization problem, we first quantify the DP noise and the data distortion.

The magnitude of DP noise can be directly represented by the variance of the Laplace distribution:

$$2\left(\frac{\Delta f}{\epsilon}\right)^2 = 2[\min(\Phi + \tau, \text{dom}(M))/\epsilon]^2, \quad (4)$$

where $\text{dom}(M)$ is the domain of the measure M . As for the data distortion, we denote the data distortion as

$$\sum_{i=1}^n \frac{(|M^{t_i}| - |M^{t'_i}|)^2}{n}. \quad (5)$$

where t_i and t'_i refer to the original and bounded records, respectively.

In Eq. 5, the value of $|M^{t'_i}|$ depends on the relationship between the domain of measure $\text{dom}(M)$ and the bounding

range $[\Phi - \tau, \Phi + \tau]$. We first consider the case where the bounding range is included in the domain of measure: $0 < \Phi - \tau < \Phi + \tau < \text{dom}(M)$. For any original record t_i , the value of its measure $|M^{t_i}|$ belongs to either $[0, \Phi - \tau]$, $[\Phi - \tau, \Phi + \tau]$, or $(\Phi + \tau, \text{dom}(M)]$. After rewriting the data table, the value of the bounded record' measure $|M^{t'_i}|$ is $\Phi - \tau$, $|M^{t_i}|$, or $\Phi + \tau$. Therefore, we can convert Eq. 5 to

$$\sum_{i=1}^{\alpha} \frac{(\Phi - \tau - |M^{t_i}|)^2}{n} + \sum_{i=1}^{\beta} \frac{(\Phi + \tau - |M^{t_i}|)^2}{n}. \quad (6)$$

where α is the number of records that satisfy $|M^{t_i}| < \Phi - \tau$, and β is the number of records that satisfy $|M^{t_i}| > \Phi + \tau$. For the case where the bounding range and the domain of measure only partially overlap, we discuss in Section V-C.

Based on Eq. 4 and Eq. 6, the objective function is shown as follows:

$$\begin{aligned} \arg \min & \sum_{i=1}^{\alpha} \frac{(\Phi - \tau - |M^{t_i}|)^2}{n} + \sum_{i=1}^{\beta} \frac{(\Phi + \tau - |M^{t_i}|)^2}{n} \\ & + 2\left[\frac{\min(\Phi + \tau, \text{dom}(M))}{\epsilon}\right]^2. \end{aligned} \quad (7)$$

s.t. $\tau > 0$

The problem is a convex optimization problem. By taking the partial derivative of Eq. 7 with respect to τ , the optimal τ is shown in Theorem 5.

Theorem 5: Suppose $B = \sum_{j=1}^{\beta} |M^{t_j}| - \sum_{i=1}^{\alpha} |M^{t_i}|$. The optimal bounding radius is

$$\tau = \frac{\epsilon^2 B - [\epsilon^2(\beta - \alpha) + 2n]\Phi}{\epsilon^2(\alpha + \beta) + 2n}. \quad (8)$$

Due to space limitation, we put the proof of Theorem 5 in Appendix B.

C. Symmetric Sparse Vector Technology

Given a data table T where each record $t_i (1 \leq i \leq n)$ has a measure whose value is $|M^{t_i}|$, when the distribution of $|M^{t_i}|$ is dispersed and the mathematic expectation of $|M^{t_i}|$ is small, a large bounding radius τ is required to avoid significant data distortion. In this case, the optimal bounding radius τ determined by Theorem 5 may be larger than the bounding axis Φ , making the bounding range and the domain of measure only partially overlap. As a result, the double-sided bounding mechanism will degrade to single-sided bounding mechanism and lead to a significant systematic bias.

Algorithm 1: Symmetric Sparse Vector Technique

Input: The domain of the measure $dom(M)$, the data table $T = \{t_i | 1 \leq i \leq n\}$, a stream of difference queries $Q = \{q_i | 1 \leq i \leq n\}$.

Output: Bounding threshold B .

```
1 Initialize a candidate set  $\mathcal{T} = \{1, 2, \dots, dom(M)\}$ ;  
2 for  $\tau \in \mathcal{T}$  do  
3   The rewritten data table  
    $T' = \{t'_i, 1 \leq i \leq n | t'_i = \min(\max(t_i, \tau), \Phi + \tau)\}$ ;  
4   if  $\sum_{i=1}^n \mathcal{A}_{q_i}^T - \sum_{i=1}^n \mathcal{A}_{q_i}^{T'} \leq 0$  then  
5      $\text{Return } \tau$ .
```

To alleviate such systematic bias, we propose a Symmetric Sparse Vector Technique (SSVT). In SSVT, when $\Phi - \tau < 0$, the data curator will choose a new lower threshold τ' for the bounding range and bound the records that are beyond $[\tau', \Phi + \tau]$ once more.

Inspired by the Sparse Vector Technique (SVT) [28], SSVT is a bottom-up mechanism that outputs the first candidate value that satisfies the constraint as the new lower threshold τ' for the bounding range. The details of SSVT are shown in Algorithm 1. For a data table T which is an instance of relational schema $R(\mathbb{A}, M)$. To compute τ' , the data curator first initializes a candidate set $\mathcal{T} = \{1, 2, \dots, dom(M)\}$, which contains all possible bounding lower bounds. The data curator traverses \mathcal{T} in ascending order. For each candidate $\tau_i \in \mathcal{T}$, the data curator first bounds the records in T with the bounding range $[\tau_i, \Phi + \tau]$ and gets a temporary rewritten table T' . Based on T and T' , the data curator answers a differential query stream $Q = \{q_i | 1 \leq i \leq n\}$ and compares the sum of results $\sum_{i=1}^n \mathcal{A}_{q_i}^T - \sum_{i=1}^n \mathcal{A}_{q_i}^{T'}$. Finally the data curator treats the first τ_i that makes $\sum_{i=1}^n \mathcal{A}_{q_i}^T - \sum_{i=1}^n \mathcal{A}_{q_i}^{T'} < 0$ as the new lower threshold of bounding range.

There are three differences between SVT and SSVT. Firstly, SSVT is customized to alleviate potential systematic bias caused by the double-sided bounding mechanism and cannot be replaced by SVT. Secondly, SSVT can maintain a stable performance since there is no manually assigned threshold required. Thirdly, since the output τ' of SSVT is obtained by tuning it on rewritten datasets without accessing the original data table, there is no privacy leakage.

D. Hypothesis Testing based Prefix Sum Computing Method

After clipping the data table via BAND, the data curator can acquire the noisy sum vector by answering the linear sum queries in the pre-query matrix via PTR. Based on the noisy sum vector, the data curator first computes the base cuboid and then constructs the remaining cuboids by aggregating the base cuboid in different dimensions. However, constructing the remaining cuboids based on the base cuboid directly will lead to noise aggregation, which refers to the problem that the DP noise in each cell of the base cuboid will be aggregated when constructing the remaining cuboids.

As a baseline approach, the Noise Control Framework [14] alleviates noise aggregation by computing a subset of cuboids from the data table in a manner that reduces the overall noise and aggregates the remaining cuboids from the descendant cuboids in the subset. However, noise aggregation is still inevitable in this framework. Since the amount of aggregated DP noise is polynomial in the cardinality of $dom(M)$, a large $dom(M)$ will lead to a massive amount of DP noise aggregated in the constructed cuboids.

We consider alleviating noise aggregation by using the prefix-sum cube, which is a kind of structure where any ancestor cuboids can be constructed by the computation between pieces of prefix sums. However, the prefix sum cannot be applied to construct the base cuboid directly. As shown in Fig. 3(b), computing prefix sums by answering a workload of linear sum queries under DP will result in multiple DP noise. Therefore, to take the advantage of the prefix sum without introducing multiple DP noise, we propose a Hypothesis Testing based Prefix Sum Computing Method (SCOPE). Rather than directly computing the prefix sums under DP, SCOPE first computes the prefix sums based on the distribution of the noisy sum vector, which is predicted via hypothesis testing, and then boosts the accuracy of prefix sums by enforcing consistency.

We first describe how to compute prefix sums based on the noisy sum vector. The idea for computing prefix sums is based on an observation: for an attribute A in data table T , if its marginal follows a distribution \mathcal{X} with probability density function f , then the corresponding prefix sums follows a distribution that is the integral of \mathcal{X} with a probability distribution function F . Inspired by this observation, we can first predict the distribution \mathcal{X} that the marginal of involved attributes follows by using hypothesis testing, and then compute the prefix sums based on the query results. Specifically, the hypothesis testing is a statistical method used to make inferences and draw conclusions about the distribution of a population based on a sample set drawn from them. For ease of illustration, we take a 1-D prefix-sum cuboid \mathcal{P} of an attribute A as an example. Suppose the data curator acquires a noisy sum vector $\mathbf{x} = [x_1, x_2, \dots, x_{dom(A)}]$, where $x_j = \sum_{|A^{t_i}|=j, 1 \leq i \leq n} |M^{t_i}| (1 \leq j \leq dom(A))$, by answering a workload of linear sum queries in the pre-query matrix. Since the input of the hypothesis testing is a sample set, the data curator first converts the noisy sum vector \mathbf{x} to a synthetic sample set \mathbf{x}' . The details of the conversion procedure are shown in Algorithm 2. In Algorithm 2, the data curator takes the noisy sum vector \mathbf{x} as input and initials a synthetic sample set \mathbf{x}' . For any element x_i in \mathbf{x} , the data curator appends i into \mathbf{x}' for x_i times. For \mathcal{X} which represents the distribution of A , \mathbf{x}' can be considered as a sample set drawn from the sample space $\Omega = \{1, 2, \dots, dom(A)\}$ following a probability distribution P .

To compute the prefix-sum cuboid \mathcal{P} , the data curator predicts the probability distribution P of the random variable \mathcal{X} through hypothesis testing. There are various hypothesis testing methods that can be used, such as Komogorov-Smirnov test [35], Anderson-Darling test [36], and Pearson's chi-

Algorithm 2: Converting

Input: The noisy sum vector \mathbf{x} .**Output:** The synthetic sample set \mathbf{x}' .

```

1 Initialize a synthetic sample set  $\mathbf{x}'$ ;
2 for  $i \in \{1, 2, \dots, \text{dom}(A)\}$  do
3   for  $j \in \text{range}(1, x_i)$  do
4     Append  $i$  into  $\mathbf{x}'$ 
5 Return  $\mathbf{x}'$ .
```

squared test [37]. Without loss of generality, here we adopt the Komogorov-Smirnov test to predict the distribution P , which is one of the most useful and general hypothesis testing methods for quantifying the distance between the empirical distribution function of the sample and the cumulative distribution function of the reference distribution. Suppose there exist a series of probability distributions in a hypothesis set $\mathbb{P} = P_1, P_2, \dots, P_i, \dots$ (e.g., normal distribution, Weibull distribution, exponential distribution, etc.). For any distribution P_i in \mathbb{P} , the data curator first computes the parameter of P_i based on \mathbf{x}' and gets the corresponding cumulative distribution function $F_{exp}^i(x)$ of P_i . After getting $F_{exp}^i(x)$, the data curator then computes the empirical distribution function $F_{obv}(x)$:

$$F_{obv}^i(x) = \frac{1}{n} \sum_{j=1}^n \mathbb{1}_{(-\infty, x]}(X_j), \quad (9)$$

where X_j is the element in \mathbf{x}' and $\mathbb{1}_{(-\infty, x]}(X_j)$ is the indicator function:

$$\mathbb{1}_{(-\infty, x]}(X_i) = \begin{cases} 1 & \text{if } X_j \leq x \\ 0 & \text{if } X_j > x \end{cases}. \quad (10)$$

Given $F_{exp}^i(x)$ and $F_{obv}^i(x)$, the Kolmogorov-Smirnov statistic is $D_i = \max |F_{exp}^i(x) - F_{obv}^i(x)|$. and the p -value is $p_i = 2e^{-2(D_i)^2 \frac{n}{2}}$. For all p_i ($1 \leq i \leq |\mathbb{P}|$), the data curator chooses the largest p -value $p_{max} = \max(p_1, p_2, \dots)$ and identifies the corresponding probability distribution of random variable \mathcal{X} as P_{max} .

We find that according to the definition, prefix-sum cuboid \mathcal{P} can be regarded as the integral of the sum vector \mathbf{x} . If we denote the discretized probability density function and the discretized cumulative distribution function of P_{max} as $f(i)$ and $F(i)$ respectively, we have $\mathcal{P}(i) = \frac{\mathbf{x}_i}{f(i)} F(i)$. Based on the noisy sum vector \mathbf{x} and the predicted distribution P_{max} , the prefix-sum cuboid \mathcal{P} can be computed.

However, the process of computing the prefix sums will inevitably introduce errors: 1) the predicted probability distribution P_{max} may deviate from the actual distribution; 2) the sum vector \mathbf{x} contains DP noise. These errors may make the noisy prefix sum $\mathcal{P}(i_2)$ smaller than $\mathcal{P}(i_1)$ ($0 \leq i_1 < i_2 \leq \text{dom}[M]$), which leads to inconsistency among cells and violates the prior knowledge that the values in the prefix-sum cube are in ascending order. To improve the utility, the data curator post-processes the noisy prefix-sum cube to remove the inconsistency. Specifically, for any cell $\tilde{P}(i_i)$ in

the noisy prefix-sum cube \tilde{P} , the data curator makes sure $0 \leq \mathcal{P}(i-1) \leq \mathcal{P}(i) \leq \mathcal{P}(i+1) \leq n$.

E. Privacy and Utility Analysis

In this section, we first show the privacy guarantee of PRIMA, then analyze the effectiveness of PRIMA and compare it with baseline approaches from two aspects: pre-query result error, noise aggregation. In addition, we analyze the prediction error in PRIMA, which is caused when predicting the distribution via hypothesis testing.

1) *Privacy Guarantee:* In PRIMA, the data table T will be accessed three times. The data curator first traverses the data table twice by using BAND to determine the bounding threshold. Then the data curator answers the linear sum queries via PTR to acquire the noisy sum vectors. In the first two times, all access occurs on the data curator's side, and no sensitive information will be released. When answering the linear sum queries in the pre-query matrix, according to the parallel composition, a piece of privacy budget ϵ is shared by all queries the pre-query matrix. Overall, according to the sequential composition, PRIMA satisfies (ϵ, δ) -DP.

2) *Pre-query Result Error:* We first focus on the pre-query result error while answering the linear sum queries under DP. For the truncation-based approaches, i.e., SVTC and R2TC, the pre-query result error is affected by both the DP noise and the systematic bias. Both SVTC and R2TC face the dilemma that a large truncation threshold cannot alleviate the DP noise, while a small truncation threshold will result in massive systematic bias. In particular, the systematic bias of SVTC is hard to quantify since it is highly impacted by the manually assigned parameter for SVT. While the systematic bias of R2TC is proportional to the logarithm of $\text{dom}(M)$. Compared with these baseline approaches, PRIMA resolves the dilemma by bounding the values of records from both sides and minimizes the systematic bias.

3) *Noise Aggregation:* Noise aggregation refers to that the DP noise added into the pre-query results will aggregate into the cells when constructing the ancestor cuboids from the descendant cuboids. Both the number of attributes aggregated and the domains of the aggregated attributes will affect the magnitude of noise aggregation. For an attribute A with domain $\text{dom}(A)$, the DP noise in the cells of ancestor cuboids will be amplified by $\text{dom}(A)$ times. In addition, if there is more than one attribute aggregated, noise aggregation will occur more than once. In BASE, all ancestor cuboids are constructed from the base cuboid directly. The DP noise will be amplified by $\text{dom}(A)^{d-\lambda}$ times when constructing λ -D cuboids, which leads to a massive noise aggregation. SVTC and R2TC alleviate noise aggregation by using the Noise Control Framework. By carefully choosing an initial set \mathcal{L}_{pre} , the ancestor cuboids can be constructed from arbitrary cuboids in \mathcal{L}_{pre} with a minimum number of attributes aggregated. As a result, in SVTC and R2TC, the DP noise will be amplified by $\text{dom}(A)^{\lambda'} \cdot 2|\mathcal{L}_{pre}|^2/\epsilon^2$ times, where λ' is typically set to 1 or 2. However, when $\text{dom}(A)$ is large, there is still a large amount of DP noise aggregated. By contrast, PRIMA

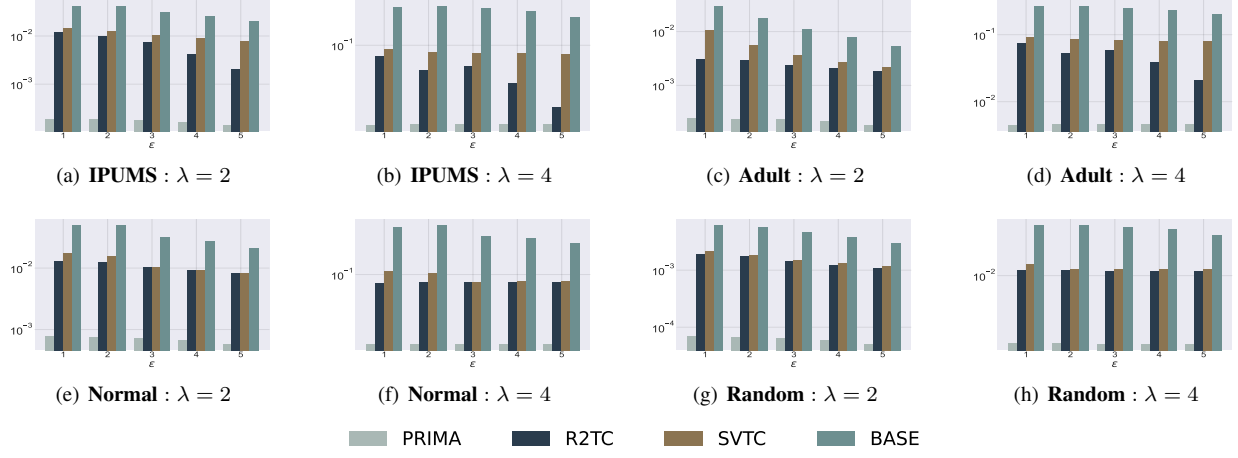


Fig. 5. Varying ϵ on all datasets under setting of $\text{dom}(A) = 30, \text{dom}(M) = 30, n = 10^5, d = 4$. MAEs are shown in log scale.

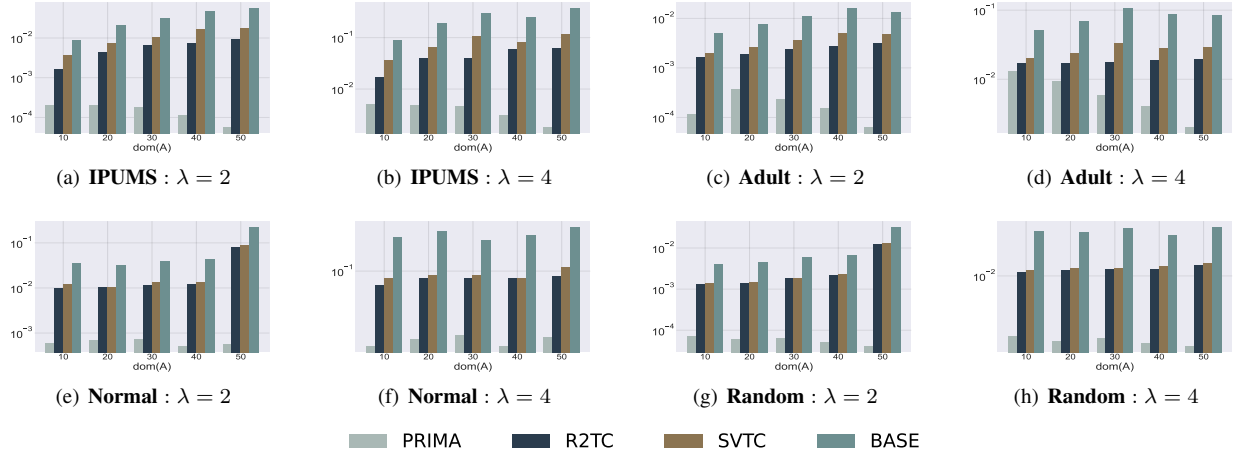


Fig. 6. Varying $\text{dom}(A)$ on all datasets under setting of $\epsilon = 2, \text{dom}(M) = 30, n = 10^5, d = 4$. MAEs are shown in log scale.

constructs the cuboids from a base prefix-sum cuboid with an amplification factor $2^{d-\lambda}$, which is a constant regardless of $\text{dom}(A)$ and smaller than any baseline approaches and leads to a better utility.

4) *Prediction Error*: The prediction error in PRIMA refers to the distance between the true distribution of the noisy sum vector and the distribution that is predicted by hypothesis testing. A major factor that impacts the accuracy of the predicted distribution is the sample size. Larger samples contain more characteristics about the distribution, which can improve the accuracy of the predicted distribution. Since PRIMA is designed for online analytical processing in a data warehouse, which stores massive amounts of data, it is able to take large enough samples to reduce the prediction error. In addition, the prior knowledge possessed by the data curator, e.g., people’s ages typically follow a normal distribution, can also boost the accuracy of the predicted distribution.

VI. EXPERIMENTS

In this section, we first evaluate the performance of PRIMA and its building blocks by varying different parameters.

A. Experimental Settings

Datasets. We make use of two real datasets and two synthetic datasets in our experiments:

- *IPUMS* [38]: A US census dataset from the IPUMS repository, which contains around 3 million records.
- *Adult* [39]: A dataset from the UCI machine learning repository. After removing missing values, the dataset contains around 50 thousands records.
- *Normal*: A dataset which is synthesized from the multivariate normal distribution with mean 0, standard deviation 1, contains 50 ordinal attributes. The covariance between every two attributes is 0.5.
- *Random*: A dataset which is synthesized from the uniform distribution, contains 50 ordinal attributes.

To experiment with different sizes of data table, we generate multiple test datasets from the four datasets with the number of records ranging from 10k to 1M. In addition, for evaluating different numbers of attributes and domain sizes, we generate multiple versions of these four datasets with the domain sizes of attributes ranging from 10 to 50.

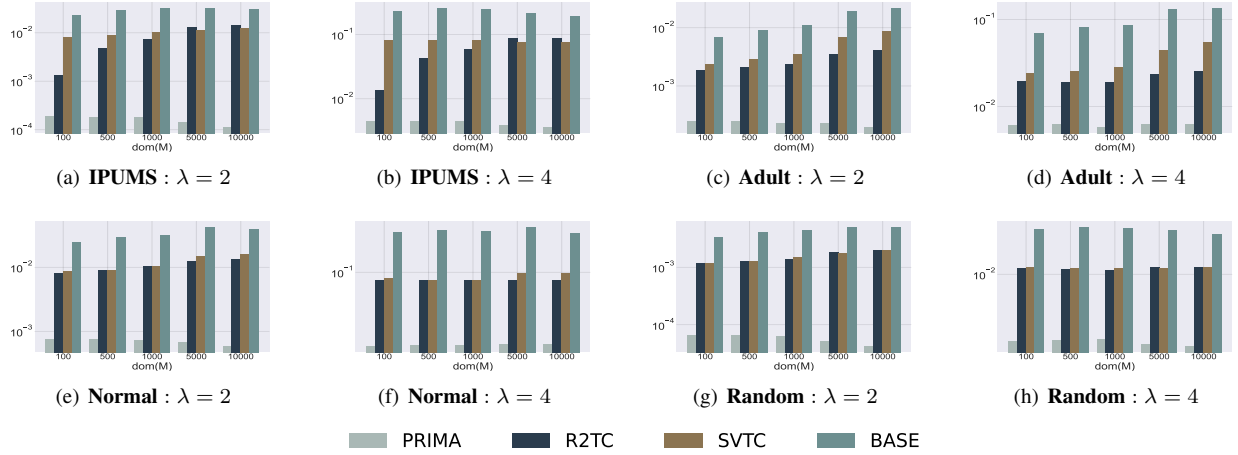


Fig. 7. Varying $dom(M)$ on all datasets under setting of $\epsilon = 2, dom(A) = 30, n = 10^5, d = 4$. MAEs are shown in log scale.

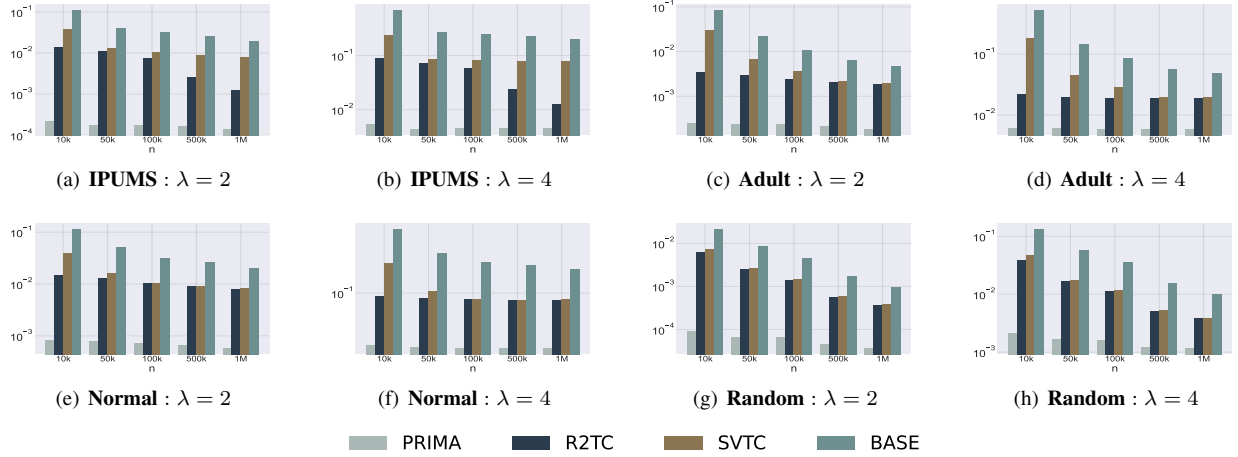


Fig. 8. Varying n on all datasets under setting of $\epsilon = 2, dom(A) = 30, dom(M) = 30, d = 4$. MAEs are shown in log scale.

Competitors. To evaluate the overall effectiveness of PRIMA, we compare PRIMA against the baseline approaches including Base, SVTC and R2TC. In particular, in SVTC and R2TC, we compute an optimal initial set \mathcal{L}_{pre} which contains a 4-D cuboid, two 3-D cuboids, two 2-D cuboids and a 1-D cuboids according to the guideline in [14].

Utility Metric. We use the Mean Absolute Error to measure the accuracy of constructed cuboids: $MAE = AVG(|\frac{\tilde{c}-c}{n}|)$, where n is the number of cells in a certain cuboids, \tilde{c} and c are the noisy and true values stored in the cells, respectively.

Methodology. For each dataset and each approach, we construct data cube between any λ attributes and measure the MAE of each cell in the m -D cuboids. We repeated all experiments 10 times to obtain stable results. In all experiments, unless explicitly stated, we use the following default values for other relevant parameters: $\epsilon = 2.0, dom(A) = 30, dom(M) = 30, n = 10^5, d = 4, \lambda = \{2, 4\}$.

Environment. We implement all the approaches in Python. All experiments are conducted on an Intel Core i5 2.50GHz PC with 8GB RAM.

B. Performance of PRIMA

In this part, we evaluate the performance of different methods under various privacy budget ϵ , the domains of attributes $dom(A)$, the domain of measure $dom(M)$ and the total number of records n .

Varying privacy budget ϵ . Fig. 5 shows the results varying ϵ from 1.0 to 5.0. It is obvious that PRIMA has a clear advantage over other approaches. We also find that, compared with other approaches, the MAEs of PRIMA are robust to the growth of ϵ . The reason can be explained as follows. Since PRIMA reduces the DP noise and the systematic bias via BAND, these two types of errors have only a minor impact on PRIMA. As a result, the effectiveness of PRIMA is not sensitive to the variations in the privacy budget ϵ .

Varying the domains of attributes $dom(A)$. Fig. 6 presents the results varying $dom(A)$ from 10 to 50 on all datasets. We observe that PRIMA consistently outperforms all other approaches. We can find that as $dom(A)$ increases, the MAEs of PRIMA are stable while the MAEs of the baseline approaches increase. The reason can be explained as follows. A larger

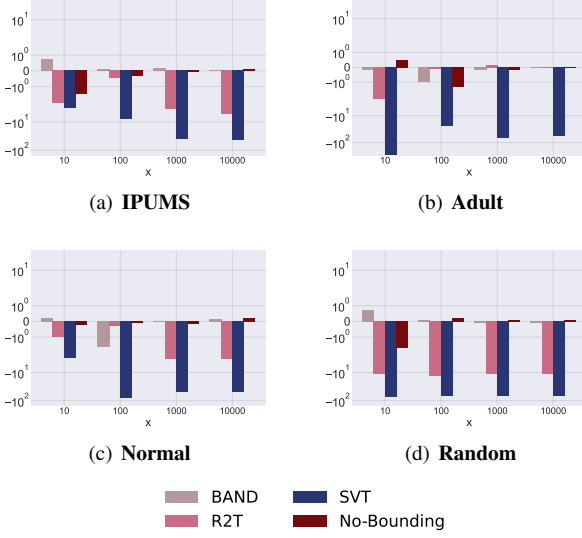


Fig. 9. MBs of x queries on all datasets.

$dom(A)$ indicates there are more cells in each cuboid. For the baseline approaches, as the number of cells in each cuboid increases, more noise aggregation will take place, resulting in higher MAEs. While PRIMA alleviates noise aggregation by constructing prefix-sum cubes. As a result, the effectiveness of PRIMA is not sensitive to the increase of $dom(A)$.

Varying the domain of attributes $dom(M)$. Fig. 7 studies the impact of $dom(M)$. It is obvious that PRIMA performs the best among all approaches. Moreover, we observe that as $dom(M)$ increases, the MAEs of PRIMA are stable while the MAEs of the baseline approaches increase. The reason can be explained as follows. A large $dom(M)$ leads to a higher sensitivity of the sum query. In BASE, a high sensitivity makes the magnitude of DP noise grow. In SVTC and R2TC, massive systematic bias is introduced when truncating the records. While PRIMA can guarantee a good query result utility by bounding the records from double sides.

Varying the number of records n . Fig. 8 shows the results varying n from 10K to 1M on all datasets. We can observe that the MAEs of all approaches decrease with the increase in the number of records, and PRIMA always performs the best.

C. Performance of Key Building Blocks

We first evaluate the advantage of BAND by comparing BAND with SVT, R2T and No-Bounding, which is a straightforward method where the data curator answers sum queries under DP without rewriting the data table. For each dataset and each method, we randomly choose x sum queries and measure the Mean Bias $MB = \frac{1}{x} \sum_{i=1}^x \tilde{c} - c$. In all experiments, we use the following default values: $\epsilon = 2.0, dom(A) = 30, dom(M) = 30, n = 10^5, d = 4, \lambda = 4$.

The MBs of four methods are shown in Fig. 9. For the four methods, the MBs of BAND and No-Bounding are close to zero as the number of queries x increases. While the MBs of SVT and R2T are negative in all cases. It indicates that SVTC and R2TC will result in a greater systematic bias.

We evaluate the advantages of SCOPE by comparing SCOPE with NCF, which refers to the Noise Control Framework, and CDB, which refers to the method for constructing a data cube in BASE. In this experiment, to construct the m -D cuboids via each method, we first obtain the required noisy sum vectors by answering linear sum queries in the pre-query matrix based on the original data table. We evaluate the effectiveness of these methods by measuring the MAE of m -D cuboids. In all experiments, we use the following default values: $\epsilon = 2.0, dom(A) = 30, dom(M) = 30, n = 10^5, d = 4$.

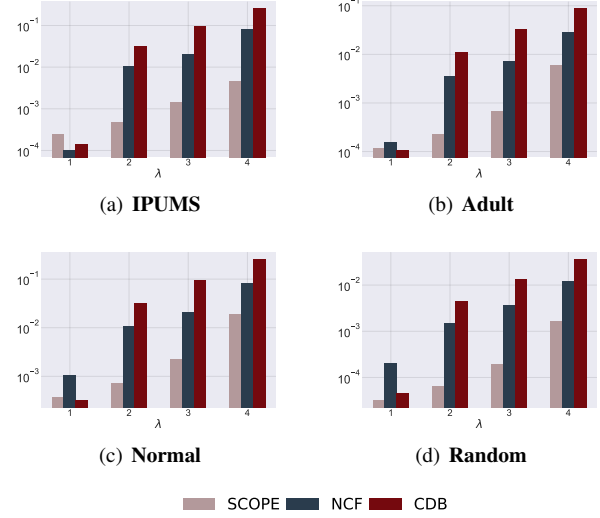


Fig. 10. MAEs of λ -D cuboids on all datasets.

According to the results shown in Fig. 10, we find SCOPE outperforms the other methods in all cases except for $m = 1$. The reason is that: since SCOPE alleviates the problem of noise aggregation by constructing prefix-sum cubes, the ancestor prefix-sum cuboids can be constructed without excessive DP noise aggregation. As m increases, the DP noise in the constructed base cuboid generally aggregates in the ancestor cuboids that are constructed via NCF and CDB, while SCOPE can maintain lower MAEs. We also observe that when $m = 1$, the MAEs of CDB are smaller than those of SCOPE and NCF. This is because, in SCOPE, the base prefix-sum cuboid contains prediction error; in NCF, there is more than one cuboid in the initial set \mathcal{L}_{pre} required to be computed, leading to a higher sensitivity when answering linear sum queries.

VII. CONCLUSION

In this paper, we propose PRIMA for processing multi-dimensional analytical queries over data cubes under DP. To answer sum queries accurately while alleviating systematic bias, we put forward BAND, which contains a double-sided bounding mechanism for rewriting the data table and a symmetric sparse vector technology for debiasing the rewritten data table. In addition, to improve the utility of noisy prefix-sum cube and alleviate noisy aggregation, we design SCOPE, where the distribution of noisy sum vector is predicted via hypothesis testing and the prefix-sum cube is constructed based on the predicted distribution. Sufficient experiments are conducted to prove the effectiveness of PRIMA.

REFERENCES

- [1] A. Edmonds, A. Nikolov, and J. R. Ullman, “The power of factorization mechanisms in local and central differential privacy,” in *Proceedings of Theory of Computing*, K. Makarychev, Y. Makarychev, M. Tulsiani, G. Kamath, and J. Chuzhoy, Eds. ACM, 2020, pp. 425–438.
- [2] M. Hardt, K. Ligett, and F. McSherry, “A simple and practical algorithm for differentially private data release,” in *Proceedings of NIPS*, P. L. Bartlett, F. C. N. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, Eds., 2012, pp. 2348–2356.
- [3] M. Hay, V. Rastogi, G. Miklau, and D. Suciu, “Boosting the accuracy of differentially private histograms through consistency,” *Proc. VLDB Endow.*, vol. 3, no. 1, pp. 1021–1032, 2010.
- [4] C. Li, M. Hay, V. Rastogi, G. Miklau, and A. McGregor, “Optimizing linear counting queries under differential privacy,” in *Proceedings of PODS*, J. Paredaens and D. V. Gucht, Eds. ACM, 2010, pp. 123–134.
- [5] C. Li, G. Miklau, M. Hay, A. McGregor, and V. Rastogi, “The matrix mechanism: optimizing linear counting queries under differential privacy,” *VLDB J.*, vol. 24, no. 6, pp. 757–781, 2015.
- [6] V. Rastogi and S. Nath, “Differentially private aggregation of distributed time-series with transformation and encryption,” in *Proceedings of SIGMOD*, A. K. Elmagarmid and D. Agrawal, Eds. ACM, 2010, pp. 735–746.
- [7] G. Yuan, Y. Yang, Z. Zhang, and Z. Hao, “Convex optimization for linear query processing under approximate differential privacy,” in *Proceedings of KDD*, B. Krishnapuram, M. Shah, A. J. Smola, C. C. Aggarwal, D. Shen, and R. Rastogi, Eds. ACM, 2016, pp. 2005–2014.
- [8] G. Yuan, Z. Zhang, M. Winslett, X. Xiao, Y. Yang, and Z. Hao, “Optimizing batch linear queries under exact and approximate differential privacy,” *ACM Trans. Database Syst.*, vol. 40, no. 2, pp. 11:1–11:47, 2015.
- [9] A. Bittau, Ú. Erlingsson, P. Maniatis, I. Mironov, A. Raghunathan, D. Lie, M. Rudominer, U. Kode, J. Tinnés, and B. Seefeld, “Prochlo: Strong privacy for analytics in the crowd,” in *Proceedings of Operating Systems Principles*. ACM, 2017, pp. 441–459.
- [10] Ú. Erlingsson, V. Pihur, and A. Korolova, “RAPPOR: randomized aggregatable privacy-preserving ordinal response,” in *Computer and Communications Security*. ACM, 2014, pp. 1054–1067.
- [11] A. D. P. Team., “Local privacy and statistical minimax rates,” 2017.
- [12] A. Machanavajjhala, D. Kifer, J. M. Abowd, J. Gehrke, and L. Vilhuber, “Privacy: Theory meets practice on the map,” in *Proceedings of International Conference on Data Engineering*. IEEE Computer Society, 2008, pp. 277–286.
- [13] S. Haney, A. Machanavajjhala, J. M. Abowd, M. Graham, M. Kutzbach, and L. Vilhuber, “Utility cost of formal privacy for releasing national employer-employee statistics,” in *Proceedings of on Management of Data*. ACM, 2017, pp. 1339–1354.
- [14] B. Ding, M. Winslett, J. Han, and Z. Li, “Differentially private data cubes: optimizing noise sources and consistency,” in *Proceedings of SIGMOD*. ACM, 2011, pp. 217–228.
- [15] Y. Wang and X. Cheng, “PRISM: prefix-sum based range queries processing method under local differential privacy,” in *Proceedings of ICDE*. IEEE, 2022, pp. 433–445.
- [16] C. Dwork, F. McSherry, K. Nissim, and A. Smith, “Calibrating noise to sensitivity in private data analysis,” in *Theory of Cryptography Conference*, 2006.
- [17] F. McSherry, “Privacy integrated queries: an extensible platform for privacy-preserving data analysis,” in *Proceedings of SIGMOD*. ACM, 2009, pp. 19–30.
- [18] D. Proserpio, S. Goldberg, and F. McSherry, “Calibrating data to sensitivity in private data analysis,” *Proc. VLDB Endow.*, vol. 7, no. 8, pp. 637–648, 2014.
- [19] D. Zhang, R. McKenna, I. Kotsogiannis, G. Bissias, M. Hay, A. Machanavajjhala, and G. Miklau, “ektelo: A framework for defining differentially private computations,” *ACM Trans. Database Syst.*, vol. 45, no. 1, pp. 2:1–2:44, 2020.
- [20] N. M. Johnson, J. P. Near, and D. Song, “Towards practical differential privacy for SQL queries,” *Proc. VLDB Endow.*, vol. 11, no. 5, pp. 526–539, 2018.
- [21] C. Dwork and J. Lei, “Differential privacy and robust statistics,” in *Proceedings of Theory of Computing*. ACM, 2009, pp. 371–380.
- [22] K. Nissim, S. Raskhodnikova, and A. D. Smith, “Smooth sensitivity and sampling in private data analysis,” in *Proceedings of Theory of Computing*. ACM, 2007, pp. 75–84.
- [23] C. Ge, X. He, I. F. Ilyas, and A. Machanavajjhala, “Apex: Accuracy-aware differentially private data exploration,” in *Proceedings of SIGMOD*. ACM, 2019, pp. 177–194.
- [24] I. Kotsogiannis, Y. Tao, X. He, M. Fanaeepour, A. Machanavajjhala, M. Hay, and G. Miklau, “Privatesql: A differentially private SQL query engine,” *Proc. VLDB Endow.*, vol. 12, no. 11, pp. 1371–1384, 2019.
- [25] W. Dong, J. Fang, K. Yi, Y. Tao, and A. Machanavajjhala, “R2T: instance-optimal truncation for differentially private query evaluation with foreign keys,” in *Proceedings of SIGMOD*. ACM, 2022, pp. 759–772.
- [26] Y. Qiu, W. Dong, K. Yi, B. Wu, and F. Li, “Releasing private data for numerical queries,” in *Proceedings of KDD*, A. Zhang and H. Rangwala, Eds. ACM, 2022, pp. 1410–1419.
- [27] R. J. Wilson, Z. C. Yuxin, L. William, D. Damien, S. Daniel, and G. Bryant, “Differentially private SQL with bounded user contribution,” *Proc. Priv. Enhancing Technol.*, vol. 2020, no. 2, pp. 230–250, 2020.
- [28] J. P. Near and C. Abuah, *Programming Differential Privacy*, 2021, vol. 1. [Online]. Available: <https://uvm-plaid.github.io/programming-dp/>
- [29] Y. Liu, A. T. Suresh, W. Zhu, P. Kairouz, and M. Gruteser, “Algorithms for bounding contribution for histogram estimation under user-level privacy,” 2023.
- [30] M. Hay, V. Rastogi, G. Miklau, and D. Suciu, “Boosting the accuracy of differentially private histograms through consistency,” *Proceedings of the VLDB Endowment*, vol. 3, no. 1, pp. 1021–1032, 2010.
- [31] W. H. Qardaji, W. Yang, and N. Li, “Understanding hierarchical methods for differentially private histograms,” *Proc. VLDB Endow.*, vol. 6, no. 14, pp. 1954–1965, 2013.
- [32] X. Xiao, G. Wang, and J. Gehrke, “Differential privacy via wavelet transforms,” in *Proceedings of ICDE*. IEEE, 2010, pp. 225–236.
- [33] R. McKenna, G. Miklau, M. Hay, and A. Machanavajjhala, “Optimizing error of high-dimensional statistical queries under differential privacy,” *Proceedings of the VLDB Endowment*, vol. 11, no. 10, pp. 1206–1219, 2018.
- [34] C. Ho, R. Agrawal, N. Megiddo, and R. Srikant, “Range queries in OLAP data cubes,” in *Proceedings of SIGMOD*. ACM, 1997, pp. 73–88.
- [35] W. Daniel, *Applied Nonparametric Statistics*. PWS-KENT Pub., 1990.
- [36] T. W. Anderson and D. A. Darling, “Asymptotic Theory of Certain “Goodness of Fit” Criteria Based on Stochastic Processes,” *The Annals of Mathematical Statistics*, vol. 23, no. 2, pp. 193–212, 1952.
- [37] K. Pearson, “X. On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling,” Jul. 1900.
- [38] Sobek, Matthew, Ruggles, and Steven, “The ipums project,” *Historical Methods*, vol. 32, no. 3, pp. 102–102, 1999.
- [39] D. Dua and C. Graff, “UCI machine learning repository,” 2017. [Online]. Available: <http://archive.ics.uci.edu/ml>

APPENDIX

A. Proof of Theorem 4

Proof. The systematic bias E_b between the original records $\{t_i | 1 \leq i \leq n\}$ and the clipped records $\{t'_i | 1 \leq i \leq n\}$ is

$$E_b = \frac{\sum_{i=1}^n |M^{t_i}| - \sum_{i=1}^n |M^{t'_i}|}{n}. \quad (11)$$

Suppose there are α records that satisfy $|M^{t_i}| < \Phi - \tau$ ($1 \leq i \leq \alpha$), β records that satisfy $|M^{t_j}| > \Phi + \tau$ ($1 \leq j \leq \beta$), γ records that satisfy $\Phi - \tau \leq |M^{t_k}| \leq \Phi + \tau$ ($1 \leq k \leq \gamma$), Equation (11) can be converted to

$$\begin{aligned} E_b &= \frac{\sum_{i=1}^{\alpha} |M^{t_i}| - (\Phi - \tau)}{n} + \frac{\sum_{k=1}^{\gamma} |M^{t_k}| - |M^{t_k}|}{n} \\ &\quad + \frac{\sum_{j=1}^{\beta} |M^{t_j}| - (\Phi + \tau)}{n} \\ &= \frac{\sum_{i=1}^n |M^{t_i}|}{n} - \frac{n - \gamma}{n} \Phi + \frac{\alpha - \beta}{n} \tau - \frac{\sum_{i=1}^{\gamma} |M^{t_i}|}{n}. \end{aligned} \quad (12)$$

As shown in Equation (12), the systematic bias is affected by both Φ and τ . For determining the optimal Φ , we ignore the effect of τ by considering the situation where τ approaches 0. Then, we have

$$\lim_{\tau \rightarrow +0} E_b = \frac{\sum_{i=1}^n |M^{t_i}|}{n} - \Phi. \quad (13)$$

To minimize the systematic bias, we can claim the optimal clipping axis is

$$\Phi = \frac{\sum_{i=1}^n |M^{t_i}|}{n}. \quad (14)$$

B. Proof of Theorem 5

Proof. We use $Err(\tau)$ to represent the objective function in Equation (7). We can convert $Err(\tau)$ to

$$\begin{aligned} Err(\tau) &= \sum_{i=1}^{\alpha} \frac{[(\Phi - \tau) - |M^{t_i}|]^2}{n} + \sum_{j=1}^{\beta} \frac{[(\Phi + \tau) - |M^{t_j}|]^2}{n} \\ &\quad + 2\left[\frac{\Phi + \tau}{\epsilon}\right]^2 \\ &= \frac{\alpha + \beta}{n} \tau^2 + \frac{2(\beta - \alpha)\Phi + 2(\sum_{i=1}^{\alpha} |M^{t_i}| + \sum_{j=1}^{\beta} |M^{t_j}|)}{n} \tau \\ &\quad + \frac{\alpha\Phi^2 + \beta\Phi^2 - 2\Phi(\sum_{i=1}^{\alpha} |M^{t_i}| + \sum_{j=1}^{\beta} |M^{t_j}|)}{n} \\ &\quad + 2\left[\frac{\Phi + \tau}{\epsilon}\right]^2. \end{aligned} \quad (15)$$

The problem is a convex optimization problem. By taking the partial derivative of Equation (15) with respect to τ , the derivative of $Err(\tau)$ is

$$\begin{aligned} Err(\tau)' &= 2\frac{\alpha + \beta}{n} \tau + \frac{2(\beta - \alpha)\Phi + 2(\sum_{i=1}^{\alpha} |M^{t_i}| + \sum_{j=1}^{\beta} |M^{t_j}|)}{n} \\ &\quad + 4\left[\frac{\Phi + \tau}{\epsilon}\right]. \end{aligned} \quad (16)$$

Let $Err(\tau)' = 0$, we have

$$\begin{aligned} Err(\tau)' &= 0 \Rightarrow \\ &2\frac{\alpha + \beta}{n} \tau + \frac{2(\beta - \alpha)\Phi + 2(\sum_{i=1}^{\alpha} |M^{t_i}| + \sum_{j=1}^{\beta} |M^{t_j}|)}{n} \\ &\quad + 4\left[\frac{\Phi + \tau}{\epsilon}\right] = 0 \Rightarrow \\ \tau &= \frac{\epsilon^2 B - [\epsilon^2(\beta - \alpha) + 2n]\Phi}{\epsilon^2(\alpha + \beta) + 2n} \end{aligned} \quad (17)$$