



VICTORIA UNIVERSITY
MELBOURNE AUSTRALIA

Protecting Private Information for Two Classes of Aggregated Database Queries

This is the Published version of the following publication

Yang, Xuechao, Yi, Xun, Kelarev, Andrei, Rylands, Leanne, Lin, Yuqing and Ryan, Joe (2022) Protecting Private Information for Two Classes of Aggregated Database Queries. *Informatics*, 9 (3). ISSN 2227-9709

The publisher's official version can be found at
<https://www.mdpi.com/2227-9709/9/3/66>

Note that access to this version may require subscription.

Downloaded from VU Research Repository <https://vuir.vu.edu.au/47865/>



Article

Protecting Private Information for Two Classes of Aggregated Database Queries

Xuechao Yang ¹, Xun Yi ¹, Andrei Kelarev ^{1,*}, Leanne Rylands ², Yuqing Lin ³ and Joe Ryan ³¹ School of Computing Technologies, RMIT University, GPO Box 2476, Melbourne, VIC 3001, Australia² Centre for Research in Mathematics and Data Science, Western Sydney University, Locked Bay 1797, Penrith, NSW 2751, Australia³ School of Information and Physical Sciences, College of Engineering, Science and Environment, The University of Newcastle, Callaghan, NSW 2308, Australia

* Correspondence: andrei.kelarev@gmail.com

Abstract: An important direction of informatics is devoted to the protection of privacy of confidential information while providing answers to aggregated queries that can be used for analysis of data. Protecting privacy is especially important when aggregated queries are used to combine personal information stored in several databases that belong to different owners or come from different sources. Malicious attackers may be able to infer confidential information even from aggregated numerical values returned as answers to queries over large collections of data. Formal proofs of security guarantees are important, because they can be used for implementing practical systems protecting privacy and providing answers to aggregated queries. The investigation of formal conditions which guarantee protection of private information against inference attacks originates from a fundamental result obtained by Chin and Ozsoyoglu in 1982 for linear queries. The present paper solves similar problems for two new classes of aggregated nonlinear queries. We obtain complete descriptions of conditions, which guarantee the protection of privacy of confidential information against certain possible inference attacks, if a collection of queries of this type are answered. Rigorous formal security proofs are given which guarantee that the conditions obtained ensure the preservation of privacy of confidential data. In addition, we give necessary and sufficient conditions for the protection of confidential information from special inference attacks aimed at achieving a group compromise.

Keywords: privacy protection; aggregated database queries; inference attacks; nonlinear queries



Citation: Yang, X.; Yi, X.; Kelarev, A.; Rylands, L.; Lin, Y.; Ryan, J. Protecting Private Information for Two Classes of Aggregated Database Queries. *Informatics* **2022**, *9*, 66. <https://doi.org/10.3390/informatics9030066>

Academic Editor: Antony Bryant

Received: 28 July 2022

Accepted: 2 September 2022

Published: 5 September 2022

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2022 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

A large and rapidly developing area of modern informatics deals with security and privacy of data (see, for example, [1–6]). In particular, the preservation of privacy is crucial for broad adoption of digital payments [7], healthcare applications [8], location-based services [9], telemedicine [10], monitoring industrial infrastructure [11] and the Internet of things [12,13].

The investigation of formal conditions which guarantee the preservation of private information against inference attacks using aggregated database queries originates from a fundamental result obtained by Chin and Ozsoyoglu [14] in the case of linear queries and linear inference attacks. It belongs to an important research direction devoted to the protection of privacy of confidential information and provides answers to aggregated queries that can be used for analysis of data [9,15]. Protecting privacy is especially important when aggregated queries are used to combine personal information stored in several databases that belong to different owners or come from different sources [16]. Malicious attackers may be able to infer confidential information even from aggregated numerical values returned as answers to queries over large collections of data [17]. Formal proofs of security guarantees are important, because they can be used for implementing practical systems protecting privacy and providing answers to aggregated queries.

The present paper obtains novel rigorous formal conditions, which guarantee the protection of privacy of confidential information against certain possible inference attacks for two new classes of aggregated nonlinear queries motivated by the main result of [14]. Section 2 of our paper gives a review of related previous work. Section 3 contains technical details on the materials and methods used in this paper. Section 4 presents main results of our article. Section 4.1 defines MEAN and VARIANCE queries (MVQ) and introduces a new class of inference attacks, quadratic equation attacks (QEA). In order to protect confidential information from QEA attacks we design a quadratic audit system (QAS). Theorems 2 and 3 establish that QAS systems guarantee the protection of confidential data from QEA attacks. Rigorous formal security proofs are given to ensure the preservation of privacy of confidential data. Section 4.2 introduces interval inference attacks (IIA). To protect sensitive data from IIA attacks, we design an interval audit system (IAS). Theorems 4 and 5 prove that the IAS ensures protection against IIA attacks. Finally, Theorem 6 in Section 4.3 gives rigorous matrix conditions for the protection of confidential information from a group compromise. The results obtained are discussed in Section 5, where directions for future research are also proposed. A conclusion is given in Section 6.

The present paper contributes to the advancement of knowledge on the preservation of privacy of confidential information by developing formal theory, designing new formal systems for the protection against inference attacks and obtaining novel rigorous conditions that guarantee that the confidential information remains protected. In summary, a point-by-point list of the main contributions of this paper can be presented as follows:

- Formal definitions of the MVQ queries and a new class of inference attacks, the QEA attacks.
- The design of a QAS system for the protection of confidential information against the QEA attacks.
- Rigorous formal proofs of Theorems 2 and 3, which establish that QAS systems guarantee the protection of confidential data from the QEA attacks.
- Formal definition of a new class of inference attacks, the IIA attacks.
- The design of an IAS system for the protection of sensitive data from the IIA attacks.
- Rigorous formal proofs of Theorems 4 and 5, which demonstrate that IAS systems ensure protection against IIA attacks.
- Rigorous formal proof of Theorem 6, which provides stringent matrix conditions for the protection of confidential information from a group compromise.

2. Previous Work

This section is devoted to the existing literature related to the results of [14] and a brief review of other relevant research. The paper [14] investigated linear queries and designed the concept of an audit expert, which maintains a dynamic matrix for processing such queries. The paper [18] suggested using a static audit expert for arbitrary linear queries, where the query basis matrix is prepared and fixed by the system beforehand. The paper [19] proposed to apply a hybrid audit expert, which combined the advantages of the dynamic and static expert systems. The effectiveness of hybrid audit experts was further investigated in [20].

The majority of previous papers devoted to linear queries concentrated on studying the more special case of so-called SUM queries (see Section 3 for a mathematical definition). The databases where the clients are allowed to submit SUM queries, were investigated in [21–23]. The readers are referred to our survey article [24] for more details.

Wu et al. [25] used the concept of differential privacy and designed a differentially private mechanism for answering linear queries, which achieves a near-optimal data utility subject to a fixed privacy protection constraint. McKenna et al. [26] applied advanced optimisation methods to develop a mechanism for accurate answers to a user-provided set of linear queries under local differential privacy. Khalili et al. [27] proposed an incentive mechanism and a randomized response algorithm for generating differentially private answers to linear queries. Xiao et al. [28] devised a fine-grained strategy of adding Gaussian

noise to query answers in the special case of answering linear queries under differential privacy subject to per-query constraints on accuracy.

Differential privacy has also been applied for privacy protection in various more advanced scenarios recently. For example, the paper by Qu et al. [29] proposed a customizable reliable differential privacy (CRDP) model and developed a modified Laplacian mechanism that enables CRDP to simultaneously minimize background knowledge attacks and eliminate collusion attacks in cyber-physical social networks. An application of the differential privacy for the development of personalised privacy protection in cyber-physical social systems was investigated in [30].

Another important relevant direction of research deals with federated learning, which occurs when a query needs to be answered by using a large database that is a union of several separate databases that belongs to different data owners not willing to share data with others due to privacy issues. For example, Wan et al. [31] proposed to integrate differential privacy and the Wasserstein Generative Adversarial Network (WGAN) for preserving the privacy of sensitive parameters in federated learning. Cui et al. [32] introduced a blockchain-empowered decentralized and asynchronous federated learning framework and designed an improved, differentially private federated learning based on generative adversarial nets. Qu et al. [33] proposed a blockchain-enabled adaptive asynchronous federated learning paradigm (FedTwin) and designed a tailor-made consensus algorithm that uses generative adversarial network-enhanced differential privacy and an improved Markov decision process. A trade-off optimization procedure and a hybrid model were developed by Qu et al. [34] for simultaneous protection of the identity and location privacy of smart mobile devices against dynamic adversaries. Blockchain-enabled federated learning and WGAN-enabled differential privacy were applied by Wan et al. [35] in order to protect confidential model parameters in the fifth-generation broadband cellular networks and beyond fifth-generation networks.

Thus, a lot of research has been conducted that investigates related directions. However, the protection of private information for the classes of nonlinear queries examined in the present paper has never been considered in the literature before.

3. Materials and Methods

If a data repository processes aggregated numerical queries for subsets of the records and provides the outcomes of these queries without giving access to individual records, then such a repository is often called a *statistical database* (cf. [36,37]). We use standard concepts and terminology, following [36,38–42]. Our proofs also apply the main theorem of [43].

The set of all real numbers is denoted by \mathbb{R} . The cardinality of a set S is denoted by $|S|$. For positive integers $a \leq b$, the symbol $[a : b]$ stands for the set

$$[a : b] = \{a, a + 1, a + 2, \dots, b\}. \quad (1)$$

A summary of the main notation used in this paper is given in Table 1.

Let m be the number of attributes in every record of the database, and let

$$\vec{r} = (r_1, r_2, \dots, r_m) \quad (2)$$

be an arbitrary record. The attributes in the database are denoted by A_1, \dots, A_m . For $1 \leq i \leq m$, the attribute A_i is a function such that $A_i(\vec{r}) = r_i$.

Let n be the number of records stored in the database. Denote the records by $\vec{r}_1, \dots, \vec{r}_n$. We assume that the users can submit aggregated queries regarding the confidential attribute A_1 , and the attributes A_2, \dots, A_m are used to select subsets of records for these queries. Then A_1 is called a *quantitative attribute* and A_2, \dots, A_m are called *characteristic attributes* for such queries. Let x_1, x_2, \dots, x_n be the (confidential) values of the quantitative attribute A_1 in the records.

Table 1. Main terminology and notation used in the present paper.

Term	Notation
Database with confidential data	D
Number of records in D	n
All records in D	$\vec{r}_1, \vec{r}_2, \dots, \vec{r}_n$
Number of attributes in each record	m
An arbitrary record in D	$\vec{r} = (r_1, r_2, \dots, r_m)$
Quantitative attribute	A_1
Characteristic attributes	A_2, \dots, A_m
Values of attribute A_1 in $\vec{r}_1, \vec{r}_2, \dots, \vec{r}_n$	x_1, x_2, \dots, x_n
Boolean expression	$B \in \mathbb{B}$
Query	(f, B)
Query sample	$S = B(D)$
Query outcome	$f(S) = f(B(D))$

The set of records chosen for a query by specifying conditions for the characteristic attributes is called the *query sample* or *query set*. To select a sample set for a query, the users can use inequalities and Boolean expressions. Denote by \mathbb{B} the set of all Boolean expressions of inequalities involving the characteristic variables. This set can be defined inductively by the following rules:

- (B1) For any $r \in \mathbb{R}, j \in [2 : m]$, the set \mathbb{B} contains inequalities $\vec{r}_j \leq r, \vec{r}_j \geq r, \vec{r}_j < r, \vec{r}_j > r$ and equality $\vec{r}_j = r$.
- (B2) If $B_1, B_2 \in \mathbb{B}$, then $B_1 \wedge B_2 \in \mathbb{B}, B_1 \vee B_2 \in \mathbb{B}, \neg B_1 \in \mathbb{B}$, where \wedge, \vee, \neg denote the logical AND, OR and NOT operators, respectively.

Throughout, we consider a query using a Boolean expression $B \in \mathbb{B}$ to select the query sample. It specifies records \vec{r} stored in the database such that the Boolean expression holds true for these records. The query sample, i.e., the set of all records in D satisfying condition B , is denoted by $S = B(D)$.

Thorough investigation in the literature has been devoted to linear queries [14,18,19,44]. A *linear query* can be recorded as a linear combination

$$\alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_n x_n = \beta, \tag{3}$$

where β is the outcome of the query, and $\alpha_1, \dots, \alpha_n \in \mathbb{R}$. Linear queries are also called *weighted sum queries*. The *COUNT query* corresponding to the linear query (3) is defined as the number of nonzero coefficients α_i , for $i \in [1 : n]$.

A *SUM query* is defined as a linear Equation (3), where β is the outcome of the query, and where

$$\alpha_i = \begin{cases} 1 & \text{if } i\text{-th record is included in the sum,} \\ 0 & \text{otherwise.} \end{cases} \tag{4}$$

When there is a set of linear queries indexed by $j = 1, \dots, k$ with equations

$$\alpha_{j,1} x_1 + \alpha_{j,2} x_2 + \dots + \alpha_{j,n} x_n = \beta_j, \tag{5}$$

then we can collect them into the matrix $M = [\alpha_{j,i}]$ and the column vector $V = [\beta_j]$. We can represent it as the matrix equation $MX = V$. Thus, every set of SUM queries (or linear queries) can be recorded as a system of linear equations of the form

$$MX = V, \tag{6}$$

where $M = [\alpha_{j,i}]$, and where $V = [\beta_j]$ is the column vector with the values returned by the queries corresponding to the rows of the matrix M . Each query corresponds to a row of the matrix M . To derive the confidential values x_1, \dots, x_n , the user can try to solve the system of linear equations.

For linear queries, it is enough to consider one-dimensional databases, or databases with only one quantitative attribute. An arbitrary set of linear queries in a multi-dimensional database can be represented as a disjoint union of linear queries corresponding to different quantitative attributes, and each of these subsets can be viewed as a set of linear queries of the corresponding 1-dimensional database.

Every linear combination of linear queries is also a linear query. If the outcomes of several linear queries are known, then the outcomes of all their linear combinations are also known. Therefore, row and column operations can be used to simplify (6). Applying row interchange, row scaling, row addition, and column interchange, the system (6) can be reduced to a normalized basis matrix form. Therefore, without loss of generality we may assume that (6) has been simplified and is represented by a *normalized query basis matrix* $M = M_k$, where

$$M_k = (I_k \mid M'_k) \quad (7)$$

and I_k is the $(k \times k)$ identity matrix. Then the matrix M is said to be in *normalized form*. The row vectors of M_k form a basis of the space of all queries with outcomes which are known, because they can all be derived by using linear combinations of query vectors.

Inference attacks can be used to derive private information from legitimately available data. It may be possible to deduce confidential information by comparing the results of several different queries. Let x_1, x_2, \dots, x_n be the values of a protected or confidential attribute in the records. If the value x_i of a confidential attribute in one record is revealed to the user, for some $i \in [1 : n]$, then this event is called a *compromise* of the database. When it is essential to emphasize that the value in precisely one record has been revealed, then the terms *1-compromise* or *classical compromise* can also be used. *Linear inference attacks* occur when malicious adversaries try to solve the system of linear equations (6) to determine confidential values.

To provide protection against linear inference attacks, Chin and Ozsoyoglu [14] proposed a system called Audit Expert. It uses a normalized basis matrix to store all queries answered previously. When a new query is added, the Audit Expert adds it to the matrix and then reduces it to a normalized basis form again.

Theorem 1 ([14]). *A statistical database with linear queries is compromised if and only if the normalized query basis matrix M_k of the Audit Expert has a row with exactly one nonzero entry. The time complexity of the algorithm dynamically processing the query matrix of the Audit Expert and maintaining it in a normalized form for a set of k consecutive linear queries is $O(k^2)$.*

4. Results

This section presents new results obtained in this paper for the protection of confidential information against the quadratic equation attacks (Section 4.1), Interval Inference Attacks (Section 4.2), and Group Compromise (Section 4.3).

4.1. Quadratic Equation Attacks

In this subsection, we consider a new different class of nonlinear queries by using variance and mean. These notions play crucial roles in hypothesis testing, significance analysis, and other studies, see [39].

Let $S = B(D)$ be a query sample, i.e., the set of records chosen by the Boolean expression B . Denote by V the set $\{r_1 \mid (r_1, \dots, r_m) \in S\}$ of values of the confidential quantitative attribute A_1 in the records of the sample S with the corresponding probability distribution. The *mean* of the values of the quantitative attribute is also called the *expected value* of the quantitative attribute. It is denoted by $\bar{V} = E(r_1)$ and is defined by the formula:

$$\bar{V} = E(r_1) = \frac{1}{|S|} \sum_{(r_1, \dots, r_m) \in S} r_1. \quad (8)$$

The *variance* of V is the expected value $E[(r_1 - E(r_1))^2]$ of the squared differences $r_1 - E(r_1)$ of values of the quantitative attribute r_1 from the mean $E(r_1)$ (see [40]). The variance of V is denoted by σ_V and is defined by the following formula:

$$\sigma_V^2 = E[(r_1 - E(r_1))^2] = \frac{1}{|S|} \sum_{(r_1, \dots, r_m) \in S} (r_1 - \bar{V})^2, \tag{9}$$

where \bar{V} is the mean given by (8) (see [40,41]). The variance measures the variability of values of the quantitative attribute from the mean. It is explained in [40] with a complete proof (see also [41]), that formula (9) can be rewritten in the following equivalent form:

$$\sigma_V^2 = E[(r_1 - E(r_1))^2] = E(r_1^2) - (E(r_1))^2 = \frac{1}{|S|} \sum_{(r_1, \dots, r_m) \in S} r_1^2 - \left(\frac{1}{|S|} \sum_{(r_1, \dots, r_m) \in S} r_1 \right)^2. \tag{10}$$

For more explanations and worked examples, the readers are referred to [40,41].

A *MEAN and VARIANCE query*, or an *MVQ query*, can be defined as a pair (f, B) , where B is a Boolean expression and f is a function $f = (f_1, f_2)$, where f_1 is defined by (8) and f_2 is defined by (9). This means that an MVQ query submits a Boolean expression B and asks to return the values of the sample mean and variance for the sample $S = B(S)$.

Equality (10) allows us to recover the outcome of each VARIANCE query from the mean value of the squares of the values of the confidential attribute. Therefore, in order to store an MVQ query in computer memory, it is enough to keep a record of the coefficients that occur in the MEAN query, the outcome of the MEAN query, and the mean value of the squares of the values of the confidential attribute. Therefore, to store a set of MVQ queries, we can use the following pair of matrix equalities,

$$MX = V, MY = W, \tag{11}$$

where M is the matrix storing the coefficients of the MEAN queries, $X = [x_1, \dots, x_n]^T$ is the column of the confidential values x_1, \dots, x_n , $Y = [x_1^2, \dots, x_n^2]^T$ is the column of the squares of the confidential values, V is the column vector with the values returned by the MEAN queries, and where W is the column vector of the mean values of the squares of the confidential values. In concise matrix notation, the Equation (11) can be stored as the following matrix:

$$(M|V|W). \tag{12}$$

The following example illustrates our matrix notation.

Example 1. Suppose that in a dataset with two records \vec{r}_1, \vec{r}_2 the values of the confidential attribute are $x_1 = 0, x_2 = 2$. Suppose that the MVQ queries have been answered for the following three samples: $\{\vec{r}_1, \vec{r}_2\}, \{\vec{r}_1\}, \{\vec{r}_2\}$. Then we get the following matrix equalities

$$\begin{bmatrix} 1/2 & 1/2 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix}, \tag{13}$$

$$\begin{bmatrix} 1/2 & 1/2 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1^2 \\ x_2^2 \end{bmatrix} = \begin{bmatrix} 2 \\ 0 \\ 4 \end{bmatrix}. \tag{14}$$

Here (13) keeps a record of the mean values of the samples, and (14) stores the corresponding mean values of the squares of the confidential attribute. We do not have to store long records of all coefficients of the VARIANCE queries, because equality (10) makes it easy to

obtain the values of all VARIANCE queries from (14). The concise matrix notation we are going to use to keep a record of all MVQ queries is the matrix

$$\left[\begin{array}{cc|cc} 1/2 & 1/2 & 1 & 2 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 2 & 4 \end{array} \right]. \tag{15}$$

Applying the row and column operations, we can reduce M to a normalized form. Then the system (11), simplifies and reduces to the normalized form

$$M_k X = V', M_k Y = W', \tag{16}$$

where the *normalized query basis matrix* M_k has the form

$$M_k = (I_k \mid M'_k), \tag{17}$$

where I_k is the $(k \times k)$ identity matrix. In concise matrix notation equations (16) can be stored as the matrix

$$(M_k \mid V' \mid W'). \tag{18}$$

Next, we define the first type of a nonlinear inference attack, the QEA attack, which can be used by an adversary to compromise MVQ queries. Steps of the QEA attack are explained in Algorithm 1.

To protect sensitive data from QEA attacks, we design a quadratic audit system (QAS). It is described in Algorithm 2 by using the following matrix notation.

Let v be a vector with n components, and let T be a $(k \times n)$ -matrix. Denote by $|v|$ the number of nonzero components in v . For $1 \leq i \leq k$, the i -th row of T is denoted by $T(i, :)$. For $1 \leq j \leq n$, the j -th column of T is denoted by $T(:, j)$. The deletion of the j -th column from T is denoted by $T[:, j] \leftarrow []$. For $1 \leq j < \ell \leq n$, the interchanging the columns j and ℓ in T is denoted by $T(:, [j \ell]) \leftarrow T(:, [\ell j])$. The vector $[v(j), v(j + 1), \dots, v(\ell)]$ is denoted by $v(j : \ell)$. The $(k + 1 \times n)$ -matrix obtained by adding the v as the last row to T is denoted by $[T; v]$. Two vectors u and v are said to be *parallel* or *collinear* if and only if either at least one of them is a zero vector, or there exists a nonzero real number α such that $u = \alpha v$. If two vectors u, v are collinear, then we write $u \parallel v$.

A formal proof establishing that the QAS system guarantees protection of sensitive data from QEA attacks is given in Theorem 3. It relies on Theorem 2, which gives matrix conditions necessary and sufficient for QEA attack to reveal confidential data.

Theorem 2 uses the concept of c -compromise, where c is a positive integer. This concept includes as a special case the notion of a classical compromise or 1-compromise treated in Theorem 1. Namely, the disclosure of a statistic based on c or fewer records in the database is called a c -compromise. The notion of a c -compromise has already been studied in the literature (see the survey paper [24] for more references).

For any row r of the matrix M_k in (17), denote by

$$r^{(k,*)} = (r_1, \dots, r_k) \tag{19}$$

the vector of the first k components of r . Denote by

$$r^{(*,n-k)} = (r_{k+1}, \dots, r_n) \tag{20}$$

the vector of the last $n - k$ components of r . Then the row has the form

$$r = (r^{(k,*)}, r^{(*,n-k)}). \tag{21}$$

The vector $r^{(*,n-k)}$ will be called the *projection* of the row r on the matrix M'_k in (17).

Algorithm 1 Quadratic Equation Attack.**Input:** A set of MVQ queries.**Output:** A compromise of the set of queries.

- 1: First, verify whether a compromise can be achieved by using only the set of MEAN queries as in Theorem 1. If not, then proceed to the next step.
- 2: Test all combinations of $t \in [1 : n]$ and $T \subseteq [1 : n]$ to find a pair (t, T) with two properties (A1), (A2):

(A1) The set of linear equations corresponding to the MEAN queries can be used to derive equalities

$$x_i = \gamma_i x_t + \delta_i \quad (22)$$

where $\gamma_i, \delta_i \in \mathbb{R}$, for all $i \in T$.

(A2) The attackers may be able to use the outcomes of the VARIANCE queries to derive a quadratic equation of the form

$$q(x_i, i \in T) = w \quad (23)$$

depending only on $x_i, i \in T$, where $w \in \mathbb{R}$.

- 3: Substitute all expressions (22) into (23) so that it becomes a quadratic equation in one variable x_t .
- 4: Solve the resulting quadratic equation in one variable x_t to achieve a compromise.
- 5: Output t, x_t .

Algorithm 2 Quadratic Audit System.**Input:** Normalized matrix $M_k = (I_k | M'_k)$ of the answered MVQ queries and the vector v of the new MVQ query.**Output:** New normalized matrix and answer to the query, or response that the query has been rejected.

- 1: $u \leftarrow v - \sum_{i=1}^k v_i \cdot M_k(i, :); j \leftarrow k + 1$
- 2: **if** $u = 0$ **then**
- 3: Answer the query, keep the matrix unchanged.
- 4: **else if** $|u| \leq 2$ **then**
- 5: Reject the query, keep the matrix unchanged.
- 6: **else**
- 7: Let u_j be the first nonzero component of u . Set $u \leftarrow \frac{1}{u_j} u; T \leftarrow [M_k; u(k + 1 : n)]$.
- 8: **if** $j > k + 1$ **then**
- 9: $M_{k+1}(:, [(k + 1) j]) \leftarrow M_{k+1}(:, [j (k + 1)])$
- 10: **for all** $i \in [1 : k]$ **do**
- 11: $T(i, :) \leftarrow T(i, :) - T(i, k + 1)T(k + 1, :)$
- 12: **if** $T(i, :)(k + 1 : n) || T(k + 1, :)(k + 1 : n)$ **then**
- 13: Reject the query, keep the matrix unchanged.
- 14: **end if**
- 15: **end for**
- 16: **end if**
- 17: Answer the query and set $M_{k+1} = [I_{k+1}; T]$.
- 18: **end if**

Theorem 2. Let D be a database with the set of MVQ queries answered so far stored in matrix form (11) with the normalized form (16). Then the following conditions are equivalent.

- (i) The QEA attack can be used to achieve a compromise of D .
- (ii) The attackers can use the set consisting of only the MEAN queries answered so far to achieve a 2-compromise of D .

(iii) Either M_k in (16) has a row with at most two nonzero entries, or M_k has two rows with collinear projections on M'_k in (17).

Proof of Theorem 2. As in the proof of the main theorem of [14] and in other previous publications, it has been customary to assume that the attackers can gain knowledge of the COUNT query corresponding to each their query. It is important to ensure rigorous protection of privacy under this assumption, in view of the following three easy ways enabling the attackers to gain access to the outcomes of the COUNT queries.

(a) The COUNT query is a legitimate query. It can be submitted to the database and may be answered as a separate query.

(b) The COUNT query can be included as an integral part of every SUM query or linear query.

(c) It may be easy for the attackers to gain access to the values of some COUNT queries by using additional information, legal knowledge, or insider knowledge.

Theorem 1 and its proof also assume that the audit system must provide protection against database compromise even if the attackers can gain access to the COUNT queries. Without this assumption, Theorem 1 is invalid. Indeed, even if the attackers can manage to obtain an outcome of the query corresponding to the value of a confidential attribute in just one record, they will be unable to notice that they have achieved this, since without the knowledge of a COUNT query they won't know whether the outcome corresponds to just one record or many records. This is why it is a common practice to assume that the attackers can also gain access to the outcomes of the corresponding COUNT queries, and that audit system must provide protection in these circumstances.

(i)⇒(ii): Suppose that condition (i) holds, i.e., the QEA could be used to achieve a compromise of D . Let us refer to the definition of the QEA attack in Algorithm 1.

First, we consider the case where the attackers managed to achieve a compromise in Step 1 of the Quadratic Equation Attack. In this case, Step 1 results in a compromise achieved by using only the set of MEAN queries. Every classical compromise is an example of a 2-compromise required for condition (ii). Therefore in this case condition (ii) follows immediately.

Now, we assume that the attackers had to proceed to the remaining steps of the QEA. This means that they found an element $t \in [1 : n]$ and a subset $T \subseteq [1 : n]$ with properties (A1) and (A2). Let us take the equality $x_1 = \gamma_1 x_t + \delta_1$, which is the first equality of the system (22). It implies that $x_1 - \gamma_1 x_t = \delta_1$. Therefore, the attackers have managed to derive the value δ_1 of the statistic $x_1 - \gamma_1 x_t$, which depends on at most two variables. This means that the attackers have achieved a 2-compromise by using only the set of MEAN queries, and so condition (ii) holds again.

(ii)⇒(iii): Suppose that condition (ii) holds, i.e., the attackers have managed to achieve a 2-compromise of D by using only MEAN queries. This means that they derived the value η of a statistic $v_1 x_{\ell_1} + v_2 x_{\ell_2}$, for some $1 \leq \ell_1 < \ell_2 \leq n$, where $v_1^2 + v_2^2 \neq 0$. Denote the rows of the matrix M by m_1, \dots, m_k . For $i \in [1 : k]$, let us denote by λ_i the linear combination of the variables x_1, \dots, x_n corresponding to the i -th row of the matrix M . This means that

$$\lambda_i = m_i X, \tag{24}$$

where $X = [x_1, \dots, x_n]^T$. Then, as in (35) above, again it follows that there exist ξ_1, \dots, ξ_k such that

$$\eta = v_1 x_{\ell_1} + v_2 x_{\ell_2} = \xi_1 \lambda_1 + \dots + \xi_k \lambda_k. \tag{25}$$

First, we consider the case where $v_1 = 0$. Then the value $\eta = v_2 x_{\ell_2}$ provides a 1-compromise. Hence, Theorem 1 implies that the normalized basis matrix M_k of the audit system has a row with only one nonzero entry. Therefore condition (iii) is satisfied.

Second, if $v_2 = 0$, then it follows in the same way that condition (iii) holds true, as well.

Third, it remains to treat the case where $v_1, v_2 \neq 0$. Note that $M_k = [I_k \mid M'_k]$ as in (7). Let us keep in mind that because I_k is an identity matrix, it follows that every nonzero linear combination of the rows of M has at least one nonzero component in the first k columns. Applying this to the linear combination (25), we see that $\ell_1 \leq k$. Furthermore, the following two subcases are possible and we consider them separately.

Subcase 1. $\ell_2 > k$. This means that x_{ℓ_2} belongs to the columns of the matrix M'_k , which is the right block of the matrix $M_k = [I_k \mid M'_k]$ in (7). Clearly, the sum $v_1x_{\ell_1} + v_2x_{\ell_2}$ has only one nonzero component in the first k columns. More specifically, the only nonzero component of this sum in the first k columns is the ℓ_1 -th component. Because I_k is an identity matrix, it follows from (25) that $\zeta_{\ell_1} \neq 0$ and

$$\zeta_1 = \dots \zeta_{\ell_1-1} = \zeta_{\ell_1-1} = \dots = \zeta_k = 0. \tag{26}$$

Hence, $\eta = \zeta_{\ell_1} \lambda_{\ell_1}$. It follows that the ℓ_1 -th row of M_k has precisely two nonzero entries, and so condition (iii) holds.

Subcase 2. $\ell_2 \leq k$. This means that x_{ℓ_1}, x_{ℓ_2} belong to the columns of the matrix I_k in M_k . Hence, we get $\zeta_{\ell_1}, \zeta_{\ell_2} \neq 0$ and all the other coefficients x_i are equal to 0, i.e.,

$$\zeta_1 = \zeta_2 = \dots = \zeta_{\ell_1-1} = \zeta_{\ell_1+1} = \zeta_{\ell_1+2} = \dots = \tag{27}$$

$$\zeta_{\ell_2-1} = \zeta_{\ell_2+1} = \zeta_{\ell_2+2} = \dots = \zeta_k = 0. \tag{28}$$

Therefore, all entries in the last $(n - k)$ columns of η are equal to zero. Denote by p_{ℓ_1} and p_{ℓ_2} the projections of the rows m_{ℓ_1} and m_{ℓ_2} on the matrix M'_k , respectively. It follows that $\zeta_{\ell_1} p_{\ell_1} + \zeta_{\ell_2} p_{\ell_2} = 0$. This implies that the projections p_{ℓ_1} and p_{ℓ_2} are collinear, and so condition (iii) is satisfied.

(iii) \Rightarrow (i): Suppose that condition (iii) holds. The following two cases are possible.

Case 1. The matrix M_k in (16) has a row with at most two nonzero entries. Denote by ℓ the index of this row, where $1 \leq \ell \leq k$. By using the same notation m_ℓ for this row and the same linear combination λ_ℓ of the variable as in (24), we get

$$\lambda_\ell = m_\ell X. \tag{29}$$

Let ℓ_1, ℓ_2 be the indices of the two nonzero entries in m_ℓ , where $1 \leq \ell_1 < \ell_2 \leq n$. Denote these two nonzero entries of m_ℓ by v_1 and v_2 . Then it follows from (29) that

$$m_\ell X = v_1 x_{\ell_1} + v_2 x_{\ell_2}. \tag{30}$$

The ℓ -th linear equation of the system (16) shows that

$$m_\ell X = v_\ell, \tag{31}$$

where v_ℓ is the ℓ -th component of the column vector V' in (16). Therefore the value of the statistic $v_1 x_{\ell_1} + v_2 x_{\ell_2}$ is equal to v_ℓ . This establishes a 2-compromise derived by using only the set of MEAN queries. Thus, condition (ii) holds.

Case 2. The matrix M_k in (16) has two rows with collinear projections on the matrix M'_k in (17). Denote by ℓ_1, ℓ_2 the indices of these rows, where $1 \leq \ell_1 < \ell_2 \leq k$. Denote by p_{ℓ_1} and p_{ℓ_2} the projections of the rows m_{ℓ_1} and m_{ℓ_2} on the matrix M'_k , respectively. Given that p_{ℓ_1} and p_{ℓ_2} are collinear, we can multiply one of these vectors by an appropriate coefficient and obtain the second vector. Without loss of generality, we may assume that there exists a coefficient φ such that $p_{\ell_1} = \varphi p_{\ell_2}$. Because I_k is an identity matrix and the projection of the vector $m_{\ell_1} - \varphi m_{\ell_2}$ on the matrix M'_k is equal to $p_{\ell_1} - \varphi p_{\ell_2}$, it follows that

$$\lambda_{\ell_1} - \varphi \lambda_{\ell_2} = x_{\ell_1} - \varphi x_{\ell_2} = v_{\ell_1} - \varphi v_{\ell_2}. \tag{32}$$

This establishes a 2-compromise again, because equalities (32) show that the value of the statistic $x_{\ell_1} - \varphi x_{\ell_2}$ is known and is equal to the constant $v_{\ell_1} - \varphi v_{\ell_2}$. This establishes that

condition (ii) is satisfied in each of the cases, i.e., the attackers can achieve a 2-compromise by using only the set of MEAN queries.

Let us introduce notation for the set of MVQ queries answered so far. Suppose that a set of k queries consisting of the corresponding pairs of mean and variance for the set of the corresponding k samples S_1, \dots, S_k have been submitted to the audit system. Applying (8), we can record the set of MEAN queries as a system of linear equations

$$\alpha_{i1}x_1 + \alpha_{i2}x_2 + \dots + \alpha_{in}x_n = \beta_i, \tag{33}$$

where $i \in [1 : k]$, where β_i is the outcome of the MEAN query, and where

$$\alpha_{ij} = \begin{cases} 0 & \text{if } j\text{-th record is not included} \\ & \text{in } i\text{-th sample } S_i, \\ \frac{1}{|S_i|} & \text{otherwise,} \end{cases} \tag{34}$$

for $j \in [1 : n]$. Denote the left-hand-side of equality (33) by q_i .

Given that the attackers have achieved a 2-compromise by using only the queries of the system (33), they have derived the value η of a statistic $v_1x_{\ell_1} + v_2x_{\ell_2}$, for some $1 \leq \ell_1 < \ell_2 \leq n$, where $v_1^2 + v_2^2 \neq 0$. It follows that there exist coefficients ξ_1, \dots, ξ_k such that

$$\xi_1q_1 + \dots + \xi_kq_k = v_1x_{\ell_1} + v_2x_{\ell_2}, \tag{35}$$

and the value of the statistic $v_1x_{\ell_1} + v_2x_{\ell_2}$ is equal to $\eta = \xi_1\beta_1 + \dots + \xi_k\beta_k$.

For each MEAN query of the system (33), the corresponding VARIANCE query of the form (9) can be rewritten in the form (10). It follows that all VARIANCE queries can be recorded as the following system of equations expressed in terms of the quadratic variables $x_1^2, x_2^2, \dots, x_n^2$

$$\gamma_{i1}x_1^2 + \dots + \gamma_{in}x_n^2 = \delta_i, \tag{36}$$

where $i \in [1 : k]$, where $\delta_i = \sigma_i^2 + \beta_i^2$, where σ_i^2 is the outcome of the i -th VARIANCE query and β_i is the outcome from (33), and where

$$\gamma_{ij} = \begin{cases} 0 & \text{if } j\text{-th record is not included} \\ & \text{in } i\text{-th sample } S_i, \\ \frac{1}{|S_i|} & \text{otherwise.} \end{cases} \tag{37}$$

Denote the left-hand-side of equality (36) by q_i . Equalities (34) and (37) show that the coefficients $\alpha_{i1}, \dots, \alpha_{in}$ in the system (33) coincide with the corresponding coefficients $\gamma_{i1}, \dots, \gamma_{in}$ in the system (36). Therefore, it follows from (35) that

$$\xi_1q_1 + \dots + \xi_kq_k = v_1x_{\ell_1}^2 + v_2x_{\ell_2}^2 = \eta. \tag{38}$$

Because at least one of the coefficients v_1, v_2 is nonzero, without loss of generality we may assume that $v_1 \neq 0$. Hence, (35) implies that

$$x_{\ell_1} = \frac{\eta}{v_1} - \frac{v_2}{v_1}x_{\ell_2}. \tag{39}$$

Substituting (39) for x_{ℓ_1} in (38), we get

$$v_1 \left(\frac{\eta}{v_1} - \frac{v_2}{v_1}x_{\ell_2} \right)^2 + v_2x_{\ell_2}^2 = \eta. \tag{40}$$

This is a quadratic equation in one variable x_{ℓ_2} . It can be solved to determine the value of x_{ℓ_2} , which achieves a compromise of D . Thus, condition (i) is satisfied.

This completes the proof of Theorem 2. \square

Theorem 3. Let $M_k = (I_k | M'_k)$ be the normalized matrix of the answered MVQ queries, and let v be the vector of the coefficients of the mean in the next MVQ query. Then Algorithm 2 answers the next query only if it is safe to do so and the QEA attack cannot be used to disclose sensitive data. Algorithm 2 ensures that the next query is rejected if the QEA attack can reveal sensitive data after an answer to this query.

Proof. The proof establishing that QAS system guarantees protection of sensitive data from QEA attacks follows from Theorem 2. It follows immediately, because Algorithm 2 verifies condition (iii) of Theorem 2 and answers the next query only if Theorem 2 guarantees that sensitive data cannot be revealed by using the QEA attack after the query is answered. \square

4.2. Interval Inference Attacks

The class of IIA inference attacks is defined in Algorithm 3. It uses the following concepts. For a positive real number ϵ , we say that an ϵ -approximate compromise or an approximate compromise with precision ϵ has been achieved, if the attackers can determine $x \in \mathbb{R}$ such that they can deduce that the value of the confidential attribute in a record belongs to the interval $[x, x + \epsilon]$. We say that an approximate compromise occurs if there exists ϵ such that an ϵ -approximate compromise has been achieved.

To protect sensitive data from IIA attacks, we design an interval audit system (IAS). It is described in Algorithm 4.

A formal proof that the IAS system protects sensitive data from IIA attacks is presented in Theorem 4. It relies on Theorem 5, which gives necessary and sufficient conditions for an approximate compromise to occur.

Algorithm 3 Interval Inference Attack.

Input: A set of MVQ queries with query sample S_j , mean m_j , variance σ_j^2 , for $j \in [1 : \ell]$.

Output: Index s of a record and the upper and lower bounds U, L for the sensitive attribute in the record.

```

1:  $S = \cup_{j=1}^{\ell} S_j$ .
2: for all  $\vec{r} \in S$  do
3:    $L_{\vec{r}} \leftarrow -\infty; U_{\vec{r}} \leftarrow +\infty$ .
4: end for
5: for all  $j \in [1 : \ell]$  do
6:   for all  $\vec{r} \in S_j$  do
7:      $L_{\vec{r}} \leftarrow \max\{L_{\vec{r}}, m_j - \sigma_j \sqrt{|S_j| - 1}\}$ ;
8:      $U_{\vec{r}} \leftarrow \min\{U_{\vec{r}}, m_j + \sigma_j \sqrt{|S_j| - 1}\}$ .
9:   end for
10: end for
11:  $L \leftarrow -\infty; U \leftarrow +\infty; s \leftarrow -\infty$ .
12: for all  $\vec{r} \in S$  do
13:   if  $|U_{\vec{r}} - L_{\vec{r}}| < |U - L|$  then
14:      $L \leftarrow L_{\vec{r}}; U \leftarrow U_{\vec{r}}; s \leftarrow$  the index of  $\vec{r}$  in  $D$ .
15:   end if
16: end for
17: Output  $s, L, U$ .
```

Theorem 4. Let ϵ be a positive real number, let the set of already answered MVQ queries consist of ℓ queries with means m_j and variances σ_j^2 , for $j \in [1 : \ell]$. Let S be the set of all records occurring in any of these already answered queries, and let $L_{\vec{r}}, U_{\vec{r}}$ be the values defined for $\vec{r} \in S$ in Algorithm 3. Let T be the sample of records of the next submitted MVQ query. Then, Algorithm 4 answers the next query only if it is safe to do so and the IIA attack cannot result in a ϵ -approximate compromise of sensitive data. Algorithm 4 ensures that the next query is rejected if the IIA attack can result in an ϵ -approximate compromise after an answer to this query.

Proof. The proof establishing that IAS system guarantees protection of sensitive data from IIA attacks follows from Theorem 5. It follows immediately, because Algorithm 4 verifies condition (iii) of Theorem 5 and answers the next query only if Theorem 5 guarantees that ε -approximate compromise does not occur after the query is answered. \square

Algorithm 4 Interval Audit System.

Input: $\varepsilon > 0$ such that the system must protect from ε -approximate compromise. The set of ℓ already answered MVQ queries with $m_j, \sigma_j^2, j \in [1 : \ell], S$, and $L_{\vec{r}}, U_{\vec{r}}$ defined for $\vec{r} \in S$ in Algorithm 3. The new MVQ query with sample T .

Output: Reject the query if it leads to ε -compromise. Otherwise, return m and σ^2 for the new query.

- 1: Compute the mean m and variance σ^2 for T .
 - 2: **for all** $\vec{r} \in S \cap T$ **do**
 - 3: $L_{\vec{r}} \leftarrow \max\{L_{\vec{r}}, m - \sigma\sqrt{|T| - 1}\};$
 - 4: $U_{\vec{r}} \leftarrow \min\{U_{\vec{r}}, m + \sigma\sqrt{|T| - 1}\}.$
 - 5: **end for**
 - 6: **for all** $\vec{r} \in T \setminus S$ **do**
 - 7: $L_{\vec{r}} \leftarrow m - \sigma\sqrt{|T| - 1};$
 - 8: $U_{\vec{r}} \leftarrow m + \sigma\sqrt{|T| - 1}.$
 - 9: **end for**
 - 10: **if** $\min\{|U_{\vec{r}} - L_{\vec{r}}| : \vec{r} \in S \cup T\} \leq \varepsilon$ **then**
 - 11: Reject the query.
 - 12: **else**
 - 13: Output m, σ .
 - 14: **end if**
-

Theorem 5. Algorithm 3 returns the index s of a record $\vec{r} = (r_1, \dots, r_n) \in D$ and an interval $[L, U] = [L_{\vec{r}}, U_{\vec{r}}]$ such that it is guaranteed that $r_1 \in [L, U]$ and the length $|U_{\vec{r}} - L_{\vec{r}}|$ of the achieves the minimum value. There exist two databases D_L and D_U such that the record \vec{r}_L with index s_L found by Algorithm 4 in D_L has confidential attribute r_1 equal to L , and the record \vec{r}_U with index s_U in D_U has confidential attribute equal to U .

Proof. Suppose that Algorithm 3 is applied to a set of samples of MVQ queries indexed by $j \in [1 : \ell]$, with query sample S_j consisting of records $\vec{r} = (r_1, \dots, r_n) \in S_j$ such that the mean and variance of the confidential components r_1 , for $\vec{r} \in S_j$, are equal to m_j and σ_j^2 , respectively.

For each $j \in [1 : \ell]$ and each record $\vec{r} \in S = \cup_{j=1}^{\ell} S_j$, it is easily seen that lines 2 to 9 of Algorithm 3 compute the following values

$$L_{\vec{r}} = \max_{j:\vec{r} \in S_j} \left\{ m_j - \sigma_j \sqrt{|S_j| - 1} \right\}, \tag{41}$$

$$U_{\vec{r}} = \min_{j:\vec{r} \in S_j} \left\{ m_j + \sigma_j \sqrt{|S_j| - 1} \right\}. \tag{42}$$

For any sample S_j , where $j \in [1 : \ell]$, and any record $\vec{r} = (r_1, \dots, r_n) \in S_j$, the following Samuelsen’s inequalities were proven in [43]:

$$m_j - \sigma_j \sqrt{|S_j| - 1} \leq r_1 \leq m_j + \sigma_j \sqrt{|S_j| - 1}. \tag{43}$$

Combining equalities (41) and (42) with all inequalities (43) for one fixed record $\vec{r} \in S$ and all samples S_j , for $j \in [1 : \ell]$, containing $\vec{r} \in S$, we get

$$L_{\vec{r}} \leq r_1 \leq U_{\vec{r}}. \tag{44}$$

It is clear that lines 11 to 16 of Algorithm 3 find the index s of the record \vec{r} such that the length $|U_{\vec{r}} - L_{\vec{r}}|$ of the interval $[L_{\vec{r}}, U_{\vec{r}}]$ achieves the minimum value.

Let D_L be a database with n records $\vec{r}[1], \dots, \vec{r}[n]$. Suppose that there is just one sample S containing all records of D_L and that the mean μ and variance σ^2 are given and fixed. Let us define

$$\vec{r}[1]_1 = L, \tag{45}$$

$$\vec{r}[2]_1 = \dots = \vec{r}[n]_1 = \mu + (\mu - L)/(n - 1), \tag{46}$$

where L and U are defined by (41) and (42), respectively. It is routine to verify that the mean of the confidential attributes of all records in D_L is equal to μ and the variance is equal to σ^2 . Then Algorithm 4 computes

$$L_{\vec{r}[1]} = \dots = L_{\vec{r}[n]} = L, \tag{47}$$

$$U_{\vec{r}[1]} = \dots = U_{\vec{r}[n]} = U. \tag{48}$$

Therefore, Algorithm 4 returns $s_L = 1, L, U$. Because $\vec{r}[1]_1 = L$, this example shows that in full generality, the value L cannot be improved.

A dual example of database D_U with

$$\vec{r}[1]_1 = U, \tag{49}$$

$$\vec{r}[2]_1 = \dots = \vec{r}[n]_1 = \mu - (U - \mu)/(n - 1), \tag{50}$$

shows that in general the value U cannot be improved either. \square

4.3. Group Compromise

Let c, k be positive integers such that $c \leq k$, and let $M_k = (I_k \mid M'_k)$ be the normalized basis matrix of a set of linear queries as in (6) and (7). We use the following well-known definitions and facts of the matrix theory (see [38]). The rank of a matrix is equal to the dimension of the vector space spanned by the rows of the matrix. It is also equal to the maximum number of linearly independent rows of the matrix. The rank of a matrix with k rows is less than k if and only if the rows of the matrix are linearly dependent, i.e., there exists a nontrivial linear combination of the rows equal to zero. The rank of the matrix M_k is equal to k .

Theorem 6. *Let c, k be positive integers such that $c \leq k$, and let $M_k = (I_k \mid M'_k)$ be the normalized basis matrix (7) of a set of linear queries for the database D . Then the following conditions are equivalent.*

- (i) *The database D is c -compromised by the set of linear queries with the normalized basis matrix M_k .*
- (ii) *There exist c columns in M_k such that after deletion of these columns the rank of the remaining matrix becomes less than k .*
- (iii) *There exist s and t with $s + t = c$ such that it is possible to remove s columns of M'_k and in this new matrix find t rows that span a space of dimension less than t .*

Proof. Let n be the number of columns in the matrix M_k in the hypothesis of this theorem. Denote the rows of M_k by m_1, \dots, m_k , and the rows of the matrix M'_k by m'_1, \dots, m'_k . For $j \in [1 : n]$, let

$$e_j = (e_{j1}, e_{j2}, \dots, e_{jn}) \tag{51}$$

be the vector with components $e_{j\ell}$, for $\ell \in [1 : n]$, defined by

$$e_{j\ell} = \begin{cases} 1 & \text{if } j = \ell, \\ 0 & \text{otherwise.} \end{cases} \tag{52}$$

Let $X = [x_1, \dots, x_n]^T$ be the column of the confidential variables.

(i)⇒(ii) Suppose that condition (i) holds. Then there exist coefficients v_1, \dots, v_k such that the linear combination $\sum_{i=1}^k v_i m_i$ has at most c nonzero components. Therefore it can be represented in the form

$$\sum_{i=1}^k v_i m_i = \sum_{\ell=1}^c \zeta_\ell e_{i_\ell} \tag{53}$$

for some positive integers $1 \leq i_1 < \dots < i_c \leq n$ and some $\zeta_1, \dots, \zeta_c \in \mathbb{R}$. Let \bar{M}_k be the matrix obtained from the matrix M_k by deleting all columns with indices i_1, \dots, i_c . Denote by $\bar{m}_1, \dots, \bar{m}_k$ the rows obtained from the rows m_1, \dots, m_k by deleting all columns i_1, \dots, i_c . It follows from (53) that $\sum_{i=1}^k v_i \bar{m}_i = 0$. Therefore the rows of the matrix \bar{M}_k are linearly dependent. It follows that the rank of \bar{M}_k is less than k . Thus, condition (ii) is satisfied.

(ii)⇒(i) Suppose that condition (ii) holds. Then there exist c columns in the matrix M_k such that the rank of the matrix \bar{M}_k obtained by deleting these columns is less than k . Denote the indices of these columns by i_1, \dots, i_c , where $1 \leq i_1 < \dots < i_c \leq n$. Let $\bar{m}_1, \dots, \bar{m}_k$ be the rows obtained from the rows m_1, \dots, m_k by deleting all columns i_1, \dots, i_c . It follows that the rows $\bar{m}_1, \dots, \bar{m}_k$ are linearly dependent, i.e., there exist coefficients v_1, \dots, v_k such that $\sum_{i=1}^k v_i \bar{m}_i = 0$. Hence, equality (53) holds true, for some ζ_1, \dots, ζ_c . Therefore the statistic (53) produces a c -compromise of the database D . Thus, condition (i) is satisfied.

(i)⇒(iii) Suppose that there is a c -compromise. As above, then there exist coefficients v_1, \dots, v_k such that the sum $\sum_{i=1}^k v_i m_i$ can be represented in the form (53), for some $1 \leq i_1 < \dots < i_c \leq n$ and ζ_1, \dots, ζ_c . Let s be the number of the indices $1 \leq i_1 < \dots < i_c \leq n$ that are greater than k . Put $t = c - s$. Then

$$i_1 < \dots < i_t \leq k < i_{t+1} < \dots < i_c \tag{54}$$

Denote by \tilde{N} the matrix obtained from M'_k by deleting the columns with indices

$$i_{t+1} - k, i_{t+1} - k + 1, \dots, i_c - k. \tag{55}$$

Let \tilde{M} be the matrix obtained from M_k by deleting the columns with indices

$$i_{t+1}, i_{t+1} + 1, \dots, i_c. \tag{56}$$

This means that \tilde{M} is obtained from M_k by replacing M'_k with \tilde{M} . Then (7) implies that

$$\tilde{M} = [I_k \mid \tilde{N}]. \tag{57}$$

Denote the rows of the matrix \tilde{M} by $\tilde{m}_1, \dots, \tilde{m}_k$. Let $\varepsilon_1, \dots, \varepsilon_k$ be the rows of the identity matrix I_k , and let $\tilde{p}_1, \dots, \tilde{p}_k$ be the rows of the matrix \tilde{N} . Then we have

$$\tilde{m}_i = (\varepsilon_i \mid \tilde{p}_i), \tag{58}$$

for $i \in [1 : k]$. Denote by $\tilde{e}_1, \dots, \tilde{e}_n$ the vectors obtained from e_1, \dots, e_n by deleting the columns with indices (55). Clearly,

$$\tilde{e}_{i_{t+1}} = \dots = \tilde{e}_{i_c} = 0. \tag{59}$$

Therefore, equality (53) implies that

$$\sum_{i=1}^k v_i \tilde{m}_i = \sum_{\ell=1}^t \zeta_\ell \tilde{e}_{i_\ell} \tag{60}$$

It follows that the sum $\sum_{i=1}^k v_i \tilde{m}_i$ has at most t nonzero components corresponding to the t vectors \tilde{e}_{i_ℓ} in the right-hand side of (60). Therefore (58), (60) and the definition of ε_i show that

$$v_i = 0 \text{ whenever } i \notin \{i_1, \dots, i_t\}. \tag{61}$$

Hence, (58), (60) and (61) imply that

$$\sum_{\ell=1}^t v_{i_\ell} \tilde{m}_{i_\ell} = \sum_{\ell=1}^t v_{i_\ell} (\varepsilon_{i_\ell} | \tilde{p}_{i_\ell}) = \sum_{\ell=1}^t \tilde{\zeta}_\ell \tilde{e}_{i_\ell}. \tag{62}$$

It follows that $\sum_{\ell=1}^t v_{i_\ell} \tilde{p}_{i_\ell} = 0$. This means that the vectors $\tilde{p}_{i_1}, \dots, \tilde{p}_{i_t}$ are linearly dependent. Because these vectors are rows of the matrix \tilde{N} , we see that these t rows of the matrix \tilde{N} span a space of dimension less than t . Thus, condition (iii) is satisfied.

(iii) \Rightarrow (i) Suppose that condition (iii) holds. Then there exist s and t such that it is possible to remove s columns with indices $i_{t+1} - k, \dots, i_c - k$ from the matrix M'_k and in this new matrix \tilde{N} find t rows $\tilde{m}_{i_1}, \dots, \tilde{m}_{i_t}$ that span a space of dimension less than t . (For consistency, here we introduce and use the same notation as in the proof of the preceding implication, so that the numbers i_{t+1}, \dots, i_c refer to the indices of the corresponding columns in the matrix M .) Then these rows are linearly dependent, and so there exists a linear combination equal to zero,

$$\sum_{\ell=1}^t v_{i_\ell} \tilde{m}_{i_\ell} = 0 \tag{63}$$

for some v_{i_1}, \dots, v_{i_t} . Consider the following linear combination

$$\varphi = \sum_{\ell=1}^t v_{i_\ell} m_{i_\ell}. \tag{64}$$

Because I_k is the identity matrix, it follows from (58) that, if we look at the last $n - k$ components of the vector φ , then we see that all nonzero values among these components correspond to the s columns i_{t+1}, \dots, i_c of M_k of the matrix M_k corresponding to the columns of the submatrix M'_k that were deleted in the discussion above. All the other values among the last $n - k$ components of φ are equal to zero by (63). Therefore, there are at most s nonzero values among the last $n - k$ components of the vector φ .

On the other hand, because I_k is an identity matrix and φ is a sum of t rows of M_k , it follows that there are at most t nonzero coordinates among the first k components of the vector φ . In total, we see that φ has at most $s + t = c$ nonzero components. It follows that the linear combination (64) of the rows of the matrix M_k produces a c -compromise of the database D . Thus, condition (i) is satisfied. This completes the proof of Theorem 6. \square

Note that the running times of the algorithms for the detection of a c -compromise using conditions (ii) and (iii) are $O(k^2 \binom{n}{c})$ and $O(2^c c^2 \binom{n}{c})$, respectively.

5. Discussion

The results obtained in this paper advance theoretical knowledge devoted to the protection of private and confidential information and prepare a foundation for the development of future comprehensive privacy protection systems.

At the same time, the results obtained have certain limitations, which motivate future work. Next, we formulate and discuss examples of directions for future research, which are motivated by our results and will need to be addressed in separate subsequent publications.

The first limitation of our results is explained by the general approach adopted in the previous papers [14,18–24]. This approach gives only exact and correct answers to the queries submitted by the clients. However, if the system detects that a query can compromise confidential information, then it only replies that the query cannot be answered. The present paper also uses this approach.

The advantage of this approach is that in the case where it is determined that a new query submitted by a client does not lead to a disclosure of confidential information, then the client will be happy to receive an exact answer to the query. However, if it is discovered that a query leads to disclosure of confidential information, then no answer is given. Therefore, the client does not receive any helpful response in the latter case.

To tackle this issue, it may be a good idea to investigate how to supply the client with some additional information expressed, for example, in terms of evaluation of probabilities. We suggest the following direction for future research.

Direction 1. *Investigate and develop hybrid systems, which provide exact answer to a query if it does not lead to disclosure of confidential information, and which use differential privacy techniques to provide a randomised probabilistic response to a query if it leads to disclosure of confidential information.*

The second limitation of our proposed systems is their focus on the particular novel classes of attacks that have not been considered previously. However, if a system provides protection against these attacks, then it can remain vulnerable to various other types of attacks. Therefore, for practical applications it is essential to consider systems providing simultaneous protection against various types of attacks without incurring a prohibitive computational overload.

Direction 2. *Design and investigate combined comprehensive systems providing answers to aggregated queries with simultaneous protection of confidential data against various different types of attacks without incurring a prohibitive computational overhead. Consider novel approaches to the optimisation of the performance of these systems.*

The third limitation of [14,18–24] and our systems is explained by the fact that this research still remains at the theoretical stage of development, when it is paramount to develop a comprehensive theory. Clearly, useful systems can be implemented as practical software only when there is sufficient rigorous theoretical foundation and only after significant advances on Direction 2 are achieved. After that, it will become important to design software implementations and conduct experimental studies comparing their performance for various categories of practical datasets. This motivates the following direction.

Direction 3. *Design software implementations of new systems proposed during future developments of Direction 2. Conduct comprehensive experimental studies comparing their performance for various categories of practical datasets.*

The fourth limitation of our systems is in the assumption that the whole collection of data is known to the system answering queries. Therefore, the systems cannot operate in the federated learning scenario. Because federated learning is a rapidly growing area of research where aggregation techniques play significant roles (see, for example, the surveys [45,46]), we propose the following direction for future research.

Direction 4. *Develop systems for protecting the privacy of confidential information in the federated learning scenario.*

Directions 1 to 4 are recorded here in general form for arbitrary queries, even though the present article motivates the investigation of these directions with a focus on the MVQ queries as the very first option for consideration.

6. Conclusions

This paper investigated nonlinear queries, which had not been considered in the literature before. It contributed to the development of formal theory designing new systems for the protection against inference attacks and obtaining novel rigorous conditions that

guarantee that the confidential information remains protected. The paper presented the following contributions to the advancement of knowledge on the preservation of privacy of confidential information:

- Definitions of the MVQ queries (Section 4.1) and the QEA attacks (Algorithm 1).
- The design of a QAS system for the protection of confidential information against the QEA attacks (Algorithm 2).
- Theorems 2 and 3 prove that QAS systems guarantee protection against the QEA attacks.
- Definition of the IIA attacks (Algorithm 3).
- The design of an IAS system for the protection of sensitive data from the IIA attacks (Algorithm 4).
- Theorems 4 and 5 prove that IAS systems ensures protection against IIA attacks.
- Theorem 6 provides stringent matrix conditions for the protection of confidential information from a group compromise.

Four directions for future research were discussed and presented in Section 5.

Author Contributions: Conceptualization, X.Y. (Xuechao Yang), X.Y. (Xun Yi), and A.K.; methodology, X.Y. (Xun Yi) and L.R.; formal analysis, A.K., L.R., Y.L., and J.R.; investigation, X.Y. (Xuechao Yang), A.K., Y.L., and J.R.; writing—original draft preparation, A.K. and L.R.; writing—review and editing, X.Y. (Xuechao Yang), X.Y. (Xun Yi), Y.L., and J.R.; supervision, X.Y. (Xun Yi); project administration, X.Y. (Xun Yi) and L.R.; funding acquisition, X.Y. (Xun Yi) and L.R. All authors have read and agreed to the published version of the manuscript.

Funding: This research was funded by the Australian Research Council, Discovery grant DP160100913.

Institutional Review Board Statement: Not Applicable.

Informed Consent Statement: Not Applicable.

Data Availability Statement: Not Applicable.

Acknowledgments: The authors are grateful to three anonymous reviewers for thorough reports and comments that have helped to improve this paper.

Conflicts of Interest: The authors declare no conflict of interest.

Abbreviations

The following abbreviations are used in this paper and subsections where they are explained:

Abbreviation	Meaning	Subsection
IAS	Interval Audit System	Section 4.2
IIA	Interval Inference Attack	Section 4.2
MVQ	Mean and Variance Query	Section 4.1
QAS	Quadratic Audit System	Section 4.1
QEA	Quadratic Equation Attack	Section 4.1

References

1. Bartol, J.; Vehovar, V.; Petrovčič, A. Should We Be Concerned about How Information Privacy Concerns Are Measured in Online Contexts? A Systematic Review of Survey Scale Development Studies. *Informatics* **2021**, *8*, 31. [\[CrossRef\]](#)
2. Downer, K.; Bhattacharya, M. BYOD Security: A Study of Human Dimensions. *Informatics* **2022**, *9*, 16. [\[CrossRef\]](#)
3. Hirschprung, R.S.; Klein, M.; Maimon, O. Harnessing Soft Logic to Represent the Privacy Paradox. *Informatics* **2022**, *9*, 54. [\[CrossRef\]](#)
4. Antunes, M.; Oliveira, L.; Seguro, A.; Verissimo, J.; Salgado, R.; Murteira, T. Benchmarking Deep Learning Methods for Behaviour-Based Network Intrusion Detection. *Informatics* **2022**, *9*, 29. [\[CrossRef\]](#)
5. Azeez, N.A.; Odufuwa, O.E.; Misra, S.; Oluranti, J.; Damaševičius, R. Windows PE Malware Detection Using Ensemble Learning. *Informatics* **2021**, *8*, 10. [\[CrossRef\]](#)
6. Perera, S.; Jin, X.; Maurushat, A.; Opoku, D.J. Factors Affecting Reputational Damage to Organisations Due to Cyberattacks. *Informatics* **2022**, *9*, 28. [\[CrossRef\]](#)
7. Sahi, A.M.; Khalid, H.; Abbas, A.F.; Zedan, K.; Khatib, S.F.A.; Al Amosh, H. The Research Trend of Security and Privacy in Digital Payment. *Informatics* **2022**, *9*, 32. [\[CrossRef\]](#)

8. Bile Hassan, I.; Murad, M.A.A.; El-Shekeil, I.; Liu, J. Extending the UTAUT2 Model with a Privacy Calculus Model to Enhance the Adoption of a Health Information Application in Malaysia. *Informatics* **2022**, *9*, 31. [[CrossRef](#)]
9. Feng, D.; Zhou, F.; Wang, Q.; Wu, Q.; Li, B. Efficient Aggregate Queries on Location Data with Confidentiality. *Sensors* **2022**, *22*, 4908. [[CrossRef](#)]
10. Iqbal, Y.; Tahir, S.; Tahir, H.; Khan, F.; Saeed, S.; Almuhaideb, A.M.; Syed, A.M. A Novel Homomorphic Approach for Preserving Privacy of Patient Data in Telemedicine. *Sensors* **2022**, *22*, 4432. [[CrossRef](#)]
11. Sobacki, A.; Barański, S.; Szymański, J. Privacy-Preserving, Scalable Blockchain-Based Solution for Monitoring Industrial Infrastructure in the Near Real-Time. *Appl. Sci.* **2022**, *12*, 7143. [[CrossRef](#)]
12. Liu, B.; Zhang, X.; Shi, R.; Zhang, M.; Zhang, G. SEPSI: A Secure and Efficient Privacy-Preserving Set Intersection with Identity Authentication in IoT. *Mathematics* **2022**, *10*, 2120. [[CrossRef](#)]
13. Xie, Y.; Li, Y.; Ma, Y. Data Privacy Security Mechanism of Industrial Internet of Things Based on Block Chain. *Appl. Sci.* **2022**, *12*, 6859. [[CrossRef](#)]
14. Chin, F.Y.; Ozsoyoglu, G. Auditing and Inference Control in Statistical Databases. *IEEE Trans. Softw. Eng.* **1982**, *SE-8*, 574–582. [[CrossRef](#)]
15. Cellamare, M.; van Gestel, A.J.; Alradhi, H.; Martin, F.; Moncada-Torres, A. A Federated Generalized Linear Model for Privacy-Preserving Analysis. *Algorithms* **2022**, *15*, 243. [[CrossRef](#)]
16. Kelarev, A.; Yi, X.; Badsha, S.; Yang, X.; Rylands, L.; Seberry, J. A Multistage Protocol for Aggregated Queries in Distributed Cloud Databases with Privacy Protection. *Future Gener. Comput. Syst.* **2019**, *90*, 368–380. [[CrossRef](#)]
17. Ziegler, J.; Pfitzner, B.; Schulz, H.; Saalbach, A.; Arnrich, B. Defending against Reconstruction Attacks through Differentially Private Federated Learning for Classification of Heterogeneous Chest X-ray Data. *Sensors* **2022**, *22*, 5195. [[CrossRef](#)]
18. Miller, M.; Seberry, J. Audit expert and Statistical Database Security. In Proceedings of the Australian Database Research Conference, Melbourne, Australian, 6 February 1990; pp. 149–174.
19. Brankovic, L.; Miller, M.; Širán, J. Towards a Practical Auditing Method for the Prevention of Statistical Database Compromise. In Proceedings of the 7th Australasian Database Conference, Melbourne, VIC, Australia, 29–30 January 1996; pp. 177–184.
20. Brankovic, L.; Miller, M.; Širán, J. Graphs, 0-1 matrices, and usability of statistical databases. *Congr. Numer.* **1996**, *120*, 169–182.
21. Miller, M.; Roberts, I.; Simpson, J. Application of symmetric chains to an optimization problem in the security of statistical databases. *Bull. Inst. Combin. Appl.* **1991**, *2*, 47–58.
22. Brankovic, L.; Miller, M. An application of combinatorics to the security of statistical databases. *Austral. Math. Soc. Gaz.* **1995**, *22*, 173–177.
23. Griggs, J.R. Concentrating Subset Sums at k Points. *Bull. Inst. Combin. Appl.* **1997**, *20*, 65–74.
24. Kelarev, A.; Ryan, J.; Rylands, L.; Seberry, J.; Yi, X. Discrete Algorithms and Methods for Security of Statistical Databases Related to the Work of Mirka Miller. *J. Discret. Algorithms* **2018**, *52–53*, 112–121. [[CrossRef](#)]
25. Wu, G.Q.; He, Y.P.; Xia, X.Y. Near-Optimal Differentially Private Mechanism for Linear Queries. *Ruan Jian Xue Bao/J. Softw.* **2017**, *28*, 2309–2322.
26. Mckenna, R.; Maity, R.K.; Mazumdar, A.; Miklau, G. A Workloadadaptive Mechanism for Linear Queries under Local Differential Privacy. In Proceedings of the PKAW2010, Online, 31 August–4 September 2020; Volume 13, pp. 1905–1918.
27. Khalili, M.M.; Vakilinia, I. Trading Privacy through Randomized Response. In Proceedings of the IEEE Conference on Computer Communications Workshops, Vancouver, BC, Canada, 10–13 May 2021. [[CrossRef](#)]
28. Xiao, Y.; Ding, Z.; Wang, Y.; Zhang, D.; Kifer, D. Optimizing Fitness-for-Use of Differentially Private Linear Queries. In Proceedings of the 47th International Conference on Very Large Data Bases, Copenhagen, Denmark, 16–20 August 2021; Volume 14, pp. 1730–1742.
29. Qu, Y.; Yu, S.; Zhou, W.; Chen, S.; Wu, J. Customizable Reliable Privacy-Preserving Data Sharing in Cyber-Physical Social Networks. *IEEE Trans. Netw. Sci. Eng.* **2021**, *8*, 269–281. [[CrossRef](#)]
30. Qu, Y.; Gao, L.; Yu, S.; Xiang, Y. Personalized Privacy Protection of IoTs Using GAN-Enhanced Differential Privacy. In *Privacy Preservation in IoT: Machine Learning Approaches*; Springer Briefs in Computer Science; Springer: Singapore, 2022; pp. 49–76. [[CrossRef](#)]
31. Wan, Y.; Qu, Y.; Gao, L.; Xiang, Y. Differentially Privacy-Preserving Federated Learning Using Wasserstein Generative Adversarial Network. In Proceedings of the IEEE Symposium on Computers and Communications, Athens, Greece, 5–8 September 2021. [[CrossRef](#)]
32. Cui, L.; Qu, Y.; Xie, G.; Zeng, D.; Li, R.; Shen, S.; Yu, S. Security and Privacy-Enhanced Federated Learning for Anomaly Detection in IoT Infrastructures. *IEEE Trans. Ind. Inform.* **2022**, *18*, 3492–3500. [[CrossRef](#)]
33. Qu, Y.; Gao, L.; Xiang, Y.; Shen, S.; Yu, S. FedTwin: Blockchain-Enabled Adaptive Asynchronous Federated Learning for Digital Twin Networks. *IEEE Netw.* **2022**, *1–8*. [[CrossRef](#)]
34. Qu, Y.; Gao, L.; Yu, S.; Xiang, Y. Hybrid Privacy Protection of IoT Using Reinforcement Learning. In *Privacy Preservation in IoT: Machine Learning Approaches*; SpringerBriefs in Computer Science; Springer: Singapore, 2022; pp. 77–109. [[CrossRef](#)]
35. Wan, Y.; Qu, Y.; Gao, L.; Xiang, Y. Privacy-Preserving Blockchain-Enabled Federated Learning for B5G-Driven Edge Computing. *Comput. Netw.* **2022**, *204*, 108671. [[CrossRef](#)]
36. Domingo-Ferrer, J.; Muralidhar, K. *Privacy in Statistical Databases, UNESCO Chair in Data Privacy*; Springer: Cham, Switzerland, 2020.

37. Brankovic, L.; Giggins, H. Statistical Database Security. In *Security, Privacy, and Trust in Modern Data Management; Data-Centric Systems and Applications*; Springer: Berlin/Heidelberg, Germany, 2007; pp. 167–181.
38. Banerjee, S.; Roy, A. *Linear Algebra and Matrix Analysis for Statistics, Texts in Statistical Science*; Chapman and Hall/CRC: New York, NY, USA; London, UK, 2014.
39. NIST/SEMATECH. E-Handbook of Statistical Methods. 2022. Available online: <http://www.itl.nist.gov/div898/handbook/> (accessed on 15 August 2022).
40. Wikipedia. Variance. 2022. Available online: https://en.wikipedia.org/wiki/Variance#Discrete_random_variable (accessed on 22 August 2022).
41. Science Buddies. Variance and Standard Deviation. 2022. Available online: <https://www.sciencebuddies.org/science-fair-projects/science-fair/variance-and-standard-deviation> (accessed on 22 August 2022).
42. Yi, X.; Paulet, R.; Bertino, E. *Homomorphic Encryption and Applications*; Springer: New York, NY, USA, 2014.
43. Samuelson, P. How Deviant Can You Be? *J. Am. Stat. Assoc.* **1968**, *63*, 1522–1525. [[CrossRef](#)]
44. Miller, M.; Seberry, J. Relative Compromise of Statistical Databases. *Aust. Comput. J.* **1989**, *21*, 56–61.
45. Yin, X.; Zhu, Y.; Hu, J. A Comprehensive Survey of Privacy-preserving Federated Learning: A Taxonomy, Review, and Future Directions. *ACM Comput. Surv.* **2021**, *54*, 1–36. [[CrossRef](#)]
46. Liu, Z.; Guo, J.; Yang, W.; Fan, J.; Lam, K.; Zhao, J. Privacy-Preserving Aggregation in Federated Learning: A Survey. *IEEE Trans. Big Data* **2022**, 1–20. [[CrossRef](#)]