

Query Auditing for Protecting Sensitive Attributes in Statistical Databases

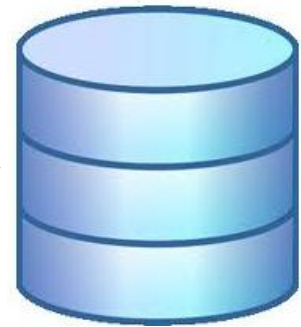
Vinh-Thong Ta

INRIA Lyon, PRIVATICS Team

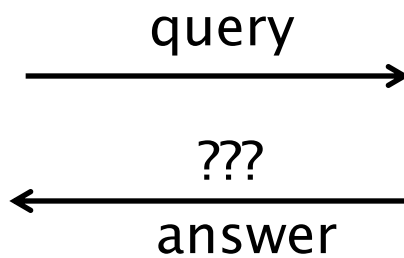
vinh-thong.ta@inria.fr

Query Auditing

Statistical
database



auditor



querier

(statistician, doctor,...)

Prevent or detect unintentional disclosure of sensitive information after a series of queries and corresponding answers

Typical Setting

- n denotes the total number of records in the DB
- $X = \{x_1, x_2, \dots, x_n\}$ are the sensitive attribute values in the records.
- $q = (Xq, f)$ is an aggregate query, where
 - Xq specifies a subset of records, called the query set
 - f is aggregation function such as MAX, MIN, SUM, AVG, MEDIAN.
- $a = f(Xq)$ is the result of applying f to Xq .

*Salary
(sensitive info)*

<i>employer_1</i>	x_1
<i>employer_2</i>	x_2
	⋮	⋮	⋮
<i>employer_{n-1}</i>	x_{n-1}
<i>employer_n</i>	x_n

E.g., $q = (\{x_1, x_2\}, \text{SUM})$

Offline vs. Online Auditing

- Offline auditing: Detection
 - Given t queries q_1, \dots, q_t and their answers a_1, \dots, a_t over X
 - Goal: determine if sensitive information has been disclosed.

$$\begin{array}{l} q_1: \quad x_i + \dots + x_n = a_1 \\ q_2: \quad x_j + \dots + x_m = a_2 \\ \dots \\ q_{t-1}: \quad x_k + \dots + x_p = a_{t-1} \\ q_t: \quad x_1 + x_2 = a_t \end{array} \left. \vphantom{\begin{array}{l} q_1 \\ q_2 \\ \dots \\ q_{t-1} \\ q_t \end{array}} \right\} \text{can the value of } x_i \text{ be determined?}$$

Offline vs. Online Auditing

- Online auditing: Prevention (and Detection)
 - Given $(t-1)$ queries $q1, \dots, qt-1$ and their $(t-1)$ answers $a1, \dots, at-1$ over X
 - **Goal:** determine if answering the new query qt would cause disclosure of sensitive information, and in case yes it denies to answer.

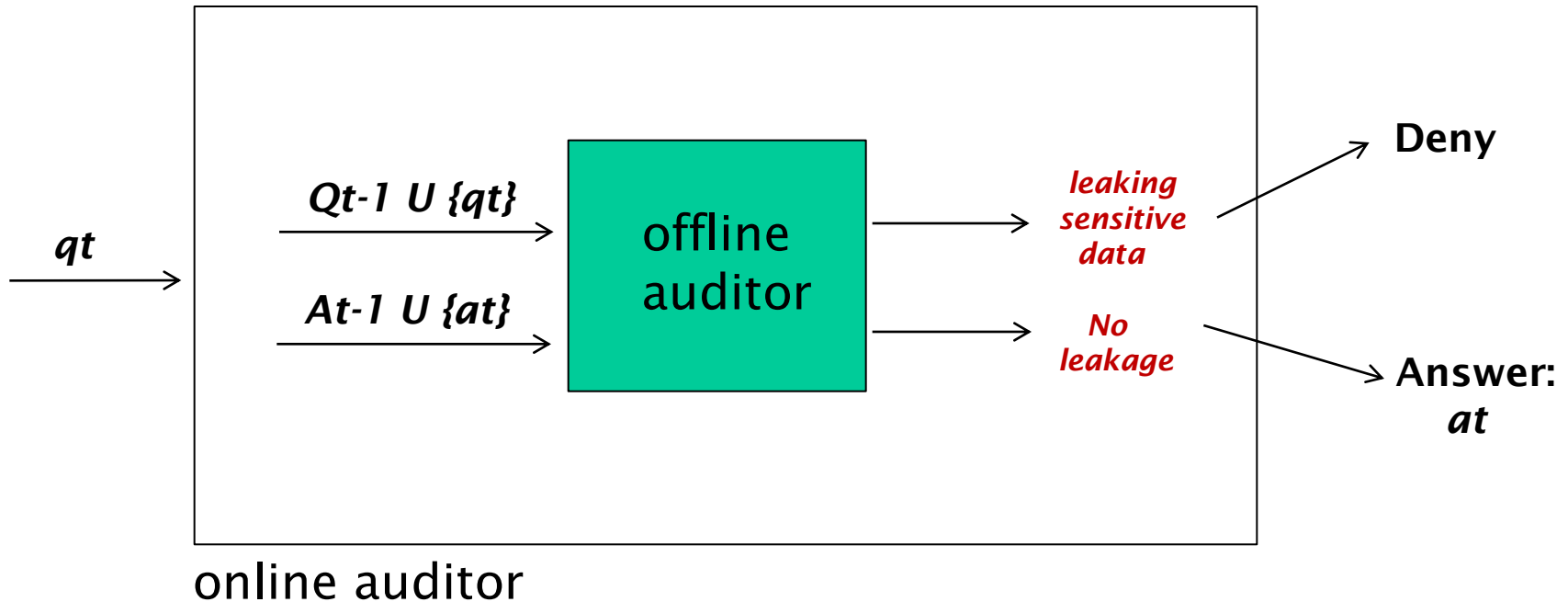
$$\left. \begin{array}{l} q1: \quad x_i + \dots + x_n = a1 \\ q2: \quad x_j + \dots + x_m = a2 \\ \dots \\ qt-1: \quad x_k + \dots + x_p = at-1 \end{array} \right\} \begin{array}{l} \text{no value of } x_i \\ \text{can be determined yet.} \end{array}$$

$$qt: \quad x_1 + x_2$$

Should we provide the answer at ?

Can we apply an offline auditor directly to solve the on-line auditing problem?

- Past queries $Q_{t-1} = \{q_1, \dots, q_{t-1}\}$ and the corresponding $A_{t-1} = \{a_1, \dots, a_{t-1}\}$ over set X .
- Determine whether to answer the *new* query q_t .



Can we apply an offline auditor directly to solve the on-line auditing problem?

This approach does not work in general, because denials also leak information!

Example:

- let $n = 3$ and $X = \{5, 5, 5\}$, SUM and MAX queries are allowed.
- let $q1 = \text{SUM}(x1, x2, x3)$ then $A = \{15\}$
- let $q2 = \text{MAX}(x1, x2, x3) \rightarrow$ this is *denied*,
 - because answering with 5 would lead to leakage: $x1=x2=x3= 5$.
- however, based on the deny it can be logically deduced that $x1=x2=x3= 5$
 - $\text{MAX}(x1, x2, x3)$ cannot be smaller than 5, otherwise the SUM cannot be 15
 - if $\text{MAX}(x1, x2, x3) > 5$ then the query would have been answered

Hence: $\text{MAX}(x1, x2, x3)$ must be 5

Another bad approach that applies offline auditor

Deny whenever the offline auditor does, and in addition, randomly deny some queries that would normally be answered by offline auditor.

- Now denials leak less information, but leakage is not generally prevented
- The auditing algorithm needs to remember which queries were randomly denied, since otherwise an attacker can repeatedly pose the same query until it is answered
- A difficulty is then to define whether two queries are equivalent

Better Approach: Simulatable Online Auditor

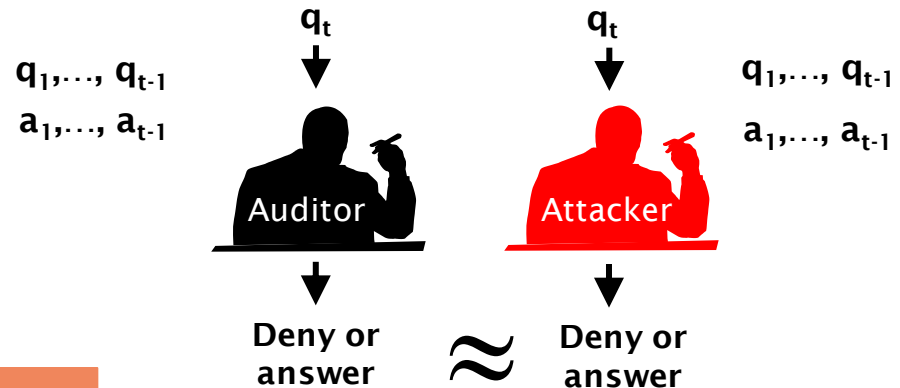
Crucial observation:

- query denials have the potential to leak information if when deciding to deny, the auditor uses information that is unavailable to the attacker (i.e., the answer ***at*** to the current query ***qt***)

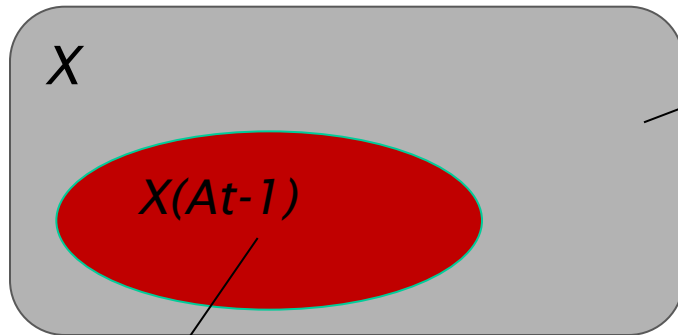
The idea behind simulatable auditors:

- the attacker should be able to ***simulate or mimic*** the auditors decisions to answer or deny a given query.
- as the attacker can equivalently determine for himself when his queries will be denied, denials leak no more information than what the attacker already knows.

→ ***provable privacy***



A sufficient condition for simulatability



Set of all possible data sets
i.e. all possible values of (x_1, \dots, x_n)
(from the attacker's point of view)

Set of all data sets consistent with past answers, $A_{t-1} = \{a_1, \dots, a_{t-1}\}$
(from the attacker's point of view)

If there is one X_i in $X(A_{t-1})$ that leads to data leakage then the auditor denies q_t , otherwise it answers with a_t .

Simulatable Auditing Example

Example revisited:

- let $n = 3$ and $X = \{x1, x2, x3\}$
- $q1 = (\{x1, x2, x3\}, SUM)$, which can be responded, $a1 = 15$.
- then $q2 = (\{x1, x2, x3\}, MAX)$ should always be denied

(even if for the values (2, 5, 8) of (x1, x2, x3), it would be safe to respond $q2$)

Because the data set (5, 5, 5), consistent with $a1=15$, would lead to data leakage. This is also known by attacker.

x1	2
x2	5
x3	8

database

- ☹ Simulatability has bad usability
 - deny too much

Full disclosure vs. Partial disclosure

Full disclosure model

- A value of x_i is fully disclosed by (Q, A) if it can be **uniquely** determined (*i.e.*, x_i is the same in all possible data sets consistent with (Q, A))

Partial disclosure model

- x_i cannot be uniquely determined, but it may be deduced to lie in a tiny interval, or in a large interval with a heavily skewed distribution.

The auditor is randomized

- it's decision to answer or deny needs not be deterministic.
- x_i is drawn from some distribution D on $(-\infty, \infty)^n$ known to both the attacker and the auditor

State-of-the-art

Auditing/Disclosure	Online Auditing	Offline Auditing	Handling Update
Probabilistic Disc.	Sim. MAX, MIN, MAX & MIN - real, unbound values Sim. SUM - real, unbound values		
Full Disc.	Sim. MAX, MIN, MAX & MIN - real, unbound values Sim. SUM - real, unbound values	MAX & SUM (NP-hard) - real, unbound values SUM - real, unbound values - boolean values MAX, MIN, MAX & MIN - real, unbound values	MAX, MIN, MAX & MIN - delete, modify, insert SUM - delete, modify, insert
Interval Disc.		SUM (real, unbound values)	

TABLE I

SUMMARY OF QUERY AUDITING PROBLEMS AND RELATED WORKS. THE ABBREVIATION SIM. MEANS SIMULATABLE.

Challenges and Open questions

- While there has been some investigation into auditing SUM, MAX, MIN, MEDIAN queries, intermingling these queries has proven to be a greater challenge.
- Auditing SQL style, Select-Project-Join queries
- Existing works are concerning with protecting individual values (e.g., salary of one person), it would be interesting to protect aggregated values (e.g., MAX, MIN).
- Collusion is a largely unaddressed issue in most interactive data sharing mechanisms today. Users can cooperate and share info.
- Utility, usability measurement
 - e.g; how to define and measure?

Our work

- While there has been some investigation into auditing SUM, MAX, MIN, MEDIAN queries, intermingling these queries has proven to be a greater challenge.
- Auditing SQL style, Select-Project-Join queries
- Existing works are concerning with protecting individual values (e.g., salary of one person), it would be interesting to protect aggregated values (e.g., MAX, MIN).
- Collusion, information sharing
 - Vinh-Thong Ta and Levente Buttyán. *Query Auditing for Protecting Max/Min Values of Sensitive Attributes in Statistical Databases*. In 9th International Conference on Trust, Privacy, Security in Digital Business (Trustbus 2012), pp. 171-186, Springer LNCS, July 2012.
- Utility, *information loss*
 - e.g; how to define and measure?

Related Works vs. Our Work

Related works:

- Detect or prevent the disclosure of the sensitive fields of *individual records* in the database
 - e.g., the salary of a given employee (x_i)
- $x_i \in (-\infty, \infty)$, real number

	<i>salary</i>	
<i>employer_1</i> x_1
<i>employer_2</i> x_2
	⋮	⋮
<i>employer_{n-1}</i> x_{n-1}
<i>employer_n</i> x_n

Our work:

- No solution proposed previously
- Detect or prevent the disclosure of *aggregate values* in the database
 - e.g., the maximum salary, $\text{MAX}(x_1, \dots, x_n)$

• $x_i \in [\alpha, \beta]$, real number

→ yields new problems which cannot be solved with existing methods!!!

	<i>salary</i>	
<i>empl_1</i> x_1
<i>empl_2</i> x_2
	⋮	⋮
<i>empl_{n-1}</i> x_{n-1}
<i>empl_n</i> x_n

MAX

- Use body mounted WSNs to collect medical data from a patient
 - e.g., ECG signals, blood pressure measurements, temperature samples
- Use a personal device (e.g., a smart phone) to collect data
- Provide controlled access to the data for external parties
 - e.g., hospital personnel, personal coach services, and health insurance companies



- The records in database all belong to the same patient.
- Individual values (i.e., sensor readings) may not be sensitive,
- **Aggregates computed over those values** can disclose the health status of the patient
 - e.g., the maximum of the blood pressure in a given time interval
- Some of the accessing parties (e.g., health insurance companies) should be prevented to learn that information.

Contributions

- Query: $AVG(Q)$, $Q \subseteq \{x_1, \dots, x_n\}$, $\forall x_i \in [\alpha, \beta]$
- Goal: detecting/preventing disclosure of $MAX\{x_1, \dots, x_n\}$ (MIN)

- We proposed three query auditors for three different settings
 - An offline query auditor for the full disclosure model
 - Based on Linear optimization problem
 - An online query auditor for the full disclosure model
 - Based on Linear optimization problem
 - An online simulatable auditor for the partial disclosure model
 - Application of the random sampling method proposed by other researchers (Lovász et. al.).

- In each case, we proved that our proposed query auditor is secure, and detects/prevents the MAX (MIN) value from disclosure.
- We also showed that they are polynomial-time algorithms.

Conclusions and future directions

- Query auditing is a broadly investigated problem, and relevant in database data protection.
- Investigating the case of combined queries (e.g., MAX & MIN & SUM).
- Investigating the case of SQL style queries.
- Addressing the protection of other types of aggregated values than MAX and MIN.
- Examining the impact of collusion attackers.
- Defining and measuring the precise degree of utility, usability.
- Proposing something with better usability than simulatable auditing.

Offline Auditor^{max}_{avg} in the Full Disclosure Model

t queries
 $b_{i,j} \in \{0,1\}$

$$\bar{A} = \begin{pmatrix} b_{1,1} & b_{1,2} & \dots & b_{1,n} \\ b_{2,1} & b_{2,2} & \dots & b_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ b_{t,1} & b_{t,2} & \dots & b_{t,n} \end{pmatrix}, \bar{a} = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_t \end{pmatrix}$$

corr. answers

$$L_j = \left\{ \begin{array}{l} \alpha \leq x_i \leq \beta, \forall x_i : x_i \in \{x_1, \dots, x_n\} \\ \bar{A}x = \bar{a}, \text{ where } \bar{x} = (x_1, \dots, x_n)^T \\ \text{ob}_j = \text{maximize}(x_j) \end{array} \right\} F$$

Lin. Eq. System
 (Feasible set)

We have n linear programming problems $P_j, j \in \{1, \dots, n\}$

Maximum of $x_1, \dots, x_n : \max\{\text{ob}_1, \dots, \text{ob}_n\}$

problem P

Simulatable Online $Auditor_{avg}^{\max}$ in the Partial Disclosure Model

λ -safe :

A sequence of queries and answers, q_1, q_2, \dots, q_t and a_1, a_2, \dots, a_t is said to be λ -safe with respect to MAX and an interval $J \subseteq [\alpha, \beta]$ if

$$Safe_{\lambda, J}(q_1, \dots, q_t, a_1, \dots, a_t) = \begin{cases} 1, & \frac{1}{1+\lambda} \leq \frac{\Pr_D(MAX \in J \mid \bigwedge_{j=1}^t (avg(Q_j) = a_j))}{\Pr_D(MAX \in J)} \leq 1+\lambda \\ 0, & \text{otherwise} \end{cases}$$

i.e., the attacker's confidence that $MAX \in I$ does not change significantly upon seeing the queries and answers.

AllSafe : $AllSafe_{\lambda, \omega}(q_1, \dots, q_t, a_1, \dots, a_t) = \begin{cases} 1, & \text{if } Safe_{\lambda, J}(q_1, \dots, q_t, a_1, \dots, a_t), \forall J \\ 0, & \text{otherwise} \end{cases}$

where interval J is ω -significant, namely, $P(MAX \in J) \geq \frac{1}{\omega}$

Simulatable Online $Auditor_{avg}^{\max}$ in the Partial Disclosure Model

(λ, ω, T) – privacy game : there are up to T rounds

in each round t :

- the attacker (adaptively) poses a query $q_t = (Q_t, AVG)$
- the auditor determines whether q_t should be answered; the auditor responds with $a_t = AVG_X(Q_t)$ if q_t is allowed, and “denies” otherwise
- the attacker wins if $AllSafe_{\lambda, \omega}(q_1, \dots, q_t, a_1, \dots, a_t) = 0$

$(\lambda, \omega, T, \delta)$ – private auditor :

for any attacker A

$$\Pr\{A \text{ wins the } (\lambda, \omega, T)\text{-privacy game}\} \leq \delta$$

where the probability is taken over the distribution D that the data comes from and the coin tosses of the randomized auditor and the attacker