

Classification of procedures

Vicenç Torra

October, 2013

Institut d'Investigació en Intel·ligència Artificial (IIIA-CSIC), Bellaterra, Catalonia

Outline

1. Dimensions
2. 1st dimension
3. 2nd dimension
4. 3rd dimension

Dimensions

Data Privacy: Dimensions

Data Privacy

- Dimensions for classification
 - Classification 1:
 - ★ On whose privacy is being sought
 - Classification 2:
 - ★ On the computations to be done
 - Classification 3:
 - ★ On the number of data sources

Data Privacy

Classification 1: On whose privacy is being sought

Data Privacy

Classification 1: On whose privacy is being sought

Subjects involved: Respondent, owner and user

Data Privacy

Classification 1: On whose privacy is being sought

Subjects involved: Respondent, owner and user

- Respondent privacy
- Owner privacy
- User privacy

Data Privacy

Classification 1: On whose privacy is being sought

- Respondent privacy
 - “Preventing re-identification of the respondents to which the records of a database correspond”
 - An issue when data has to be made available to third parties
- Owner privacy
 - Prevent the disclosure of the database.
- User privacy
 - Avoid the disclosure of information related to the user

Data Privacy

Classification 2: On the computations to be done

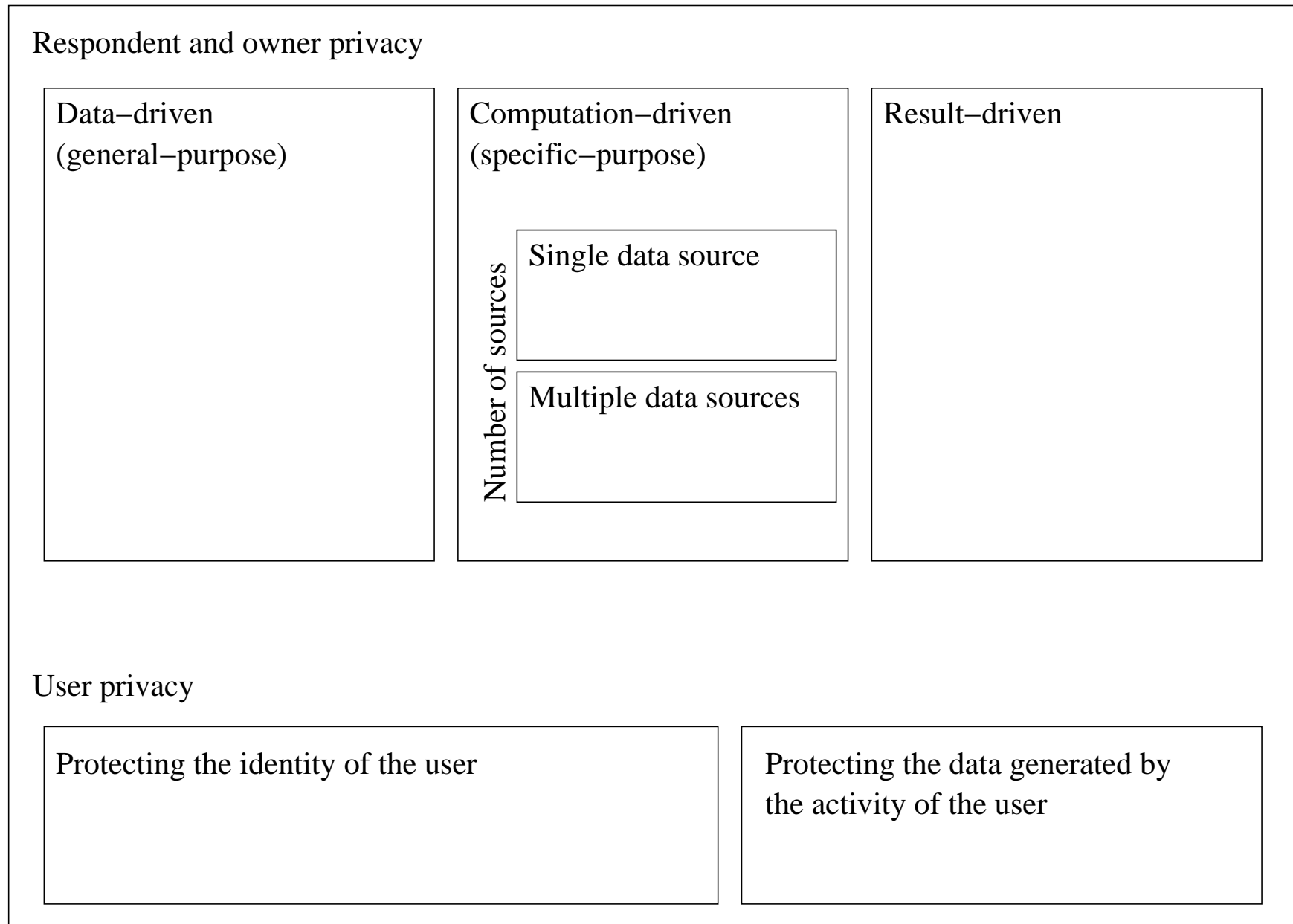
- Data-driven or general purpose protection procedures
- Computation-driven or specific purpose protection procedures
- Result-driven

Data Privacy

Classification 3: On the number of data sources

- Single data source
- Multiple data sources

Data Privacy



Dimensions. 1st classification
On whose privacy is being sought

Data Privacy

Classification 1: On whose privacy is being sought

- Respondent privacy
- Owner privacy
- User privacy

Data Privacy

User privacy

- Protecting the identity of the user
- Protecting the data generated by the activity of the user

Data Privacy

User privacy

- Protecting the identity of the user
- Protecting the data generated by the activity of the user

Tools for anonymous communications belong to user privacy

Data Privacy

User privacy

- Protecting the identity of the user
- Protecting the data generated by the activity of the user

Tools for anonymous communications belong to user privacy

Other examples with users querying databases

Data Privacy

User privacy in database search

- Protecting the identity of the user
 - Protect who is making a query

Data Privacy

User privacy in database search

- Protecting the identity of the user
 - Protect who is making a query
 - Anonymous database search
- Protecting the data generated by the user

Data Privacy

User privacy in database search

- Protecting the identity of the user
 - Protect who is making a query
 - Anonymous database search
- Protecting the data generated by the user
 - Protect the query of the user

Data Privacy

User privacy in database search

- Protecting the identity of the user
 - Protect who is making a query
 - Anonymous database search
- Protecting the data generated by the user
 - Protect the query of the user
 - Private Information Retrieval (PIR)

Data Privacy

User privacy

- Private Information Retrieval (PIR)
- Anonymous database search

Data Privacy

User privacy

- Private Information Retrieval (PIR)
 - How a user should retrieve an element from a DB or a search engine, without the system or the server being able to deduce **which element** is the object of the user's interest.

Data Privacy

User privacy

- **Private Information Retrieval (PIR)**
 - (Information Theoretic) Private Information Retrieval (PIR)
 - Computational PIR (cPIR)
 - Trusted-hardware PIR
 - Other approaches
 - ★ Goopir
 - ★ TrackMeNot

Data Privacy

User privacy

- (Information Theoretic) Private Information Retrieval (PIR)
 - Information theoretic: cannot be broken with **unlimited computing power**

Data Privacy

User privacy

- (Information Theoretic) Private Information Retrieval (PIR)
 - Information theoretic: cannot be broken with **unlimited computing power**
 - Every (information theoretic) PIR scheme with a single-database (with n bits) **requires $\Omega(n)$ bits of communication.**

Data Privacy

User privacy

- (Information Theoretic) Private Information Retrieval (PIR)
 - Information theoretic: cannot be broken with **unlimited computing power**
 - Every (information theoretic) PIR scheme with a single-database (with n bits) **requires $\Omega(n)$ bits of communication.**
 - **It can be proven (Chor et al. 1998)** that if a user wants to keep its privacy (in the information theoretic sense), then essentially the only thing he can do is to ask for a copy of the whole database.

Data Privacy

User privacy

- (Information Theoretic) PIR:
 - Communication complexity is reduced: sublinear in n by assuming that the **data is replicated**.

Data Privacy

User privacy

- (Information Theoretic) PIR:
 - Communication complexity is reduced: sublinear in n by assuming that the **data is replicated**.
 - ★ k **copies** of the database are considered
 - ★ DB copies **do not** collaborate

Data Privacy

User privacy

- (Information Theoretic) PIR:
 - Communication complexity is reduced: sublinear in n by assuming that the **data is replicated**.
 - ★ k **copies** of the database are considered
 - ★ DB copies **do not** collaborate
 - **Example.** Scheme in (Chor et al., 1999) with **communication complexity** $O(n^{1/3})$ for $k = 2$

Data Privacy

User privacy

- (Information Theoretic) PIR: k copies of the database (not being intercommunicated)
 - **Problem.**
 - ★ Database. A binary string $x = x_1 \cdots x_n$ of length n
(Identical copies of this string are stored in $k \geq 2$ servers)
 - ★ User. Given index i , is interested in obtaining the value of bit x_i
 - ★ *Solution:* The user queries each of the servers and gets replies from which the desired bit x_i can be computed.
The server does not gain any information about i from the query.

Data Privacy

Definition of the problem. (Information Theoretic) PIR (I)

- Input
 - $i \in [n]$ where $[n] = \{1, \dots, n\}$
 - r random input of length ℓ_{rnd}
- Overview of the process
 - k queries $Q_1(i, r), \dots, Q_k(i, r)$ of length ℓ_q each
 - Servers respond according to strategies A_1, \dots, A_k with replies of length ℓ_a according to the content of the DB x
 - The user reconstructs the desired bit x_i from the k replies, together with i and r

Data Privacy

Definition of the problem. (Information Theoretic) PIR (I)

- Formalization
 - A k -server PIR scheme for database length n consists of
 - ★ k query functions $Q_1, \dots, Q_k : [n] \times \{0, 1\}^{\ell_{rnd}} \rightarrow \{0, 1\}^{\ell_q}$
 - ★ k answer functions, $A_1, \dots, A_k : \{0, 1\}^n \times \{0, 1\}^{\ell_q} \rightarrow \{0, 1\}^{\ell_a}$
 - ★ a reconstruction function $R : [n] \times \{0, 1\}^{\ell_{rnd}} \times (\{0, 1\}^{\ell_a})^k \rightarrow \{0, 1\}$
 - These functions should satisfy

Data Privacy

Definition of the problem. (Information Theoretic) PIR (I)

- Formalization
 - A k -server PIR scheme for database length n consists of
 - ★ k query functions $Q_1, \dots, Q_k : [n] \times \{0, 1\}^{\ell_{rnd}} \rightarrow \{0, 1\}^{\ell_q}$
 - ★ k answer functions, $A_1, \dots, A_k : \{0, 1\}^n \times \{0, 1\}^{\ell_q} \rightarrow \{0, 1\}^{\ell_a}$
 - ★ a reconstruction function $R : [n] \times \{0, 1\}^{\ell_{rnd}} \times (\{0, 1\}^{\ell_a})^k \rightarrow \{0, 1\}$
 - These functions should satisfy
 - ★ **Correctness.** For every $x \in \{0, 1\}^n$, $i \in [n]$, and $r \in \{0, 1\}^{\ell_{rnd}}$

$$R(i, r, A_1(x, Q_1(i, r)), \dots, A_k(x, Q_k(i, r))) = x_i$$

Data Privacy

Definition of the problem. (Information Theoretic) PIR (I)

- Formalization

- A k -server PIR scheme for database length n consists of

- ★ k query functions $Q_1, \dots, Q_k : [n] \times \{0, 1\}^{\ell_{rnd}} \rightarrow \{0, 1\}^{\ell_q}$

- ★ k answer functions, $A_1, \dots, A_k : \{0, 1\}^n \times \{0, 1\}^{\ell_q} \rightarrow \{0, 1\}^{\ell_a}$

- ★ a reconstruction function $R : [n] \times \{0, 1\}^{\ell_{rnd}} \times (\{0, 1\}^{\ell_a})^k \rightarrow \{0, 1\}$

- These functions should satisfy

- ★ **Correctness.** For every $x \in \{0, 1\}^n$, $i \in [n]$, and $r \in \{0, 1\}^{\ell_{rnd}}$

$$R(i, r, A_1(x, Q_1(i, r)), \dots, A_k(x, Q_k(i, r))) = x_i$$

- ★ **Privacy.** For every $i, j \in [n]$, $s \in [k]$, and $q \in \{0, 1\}^{\ell_q}$

$$\Pr(Q_s(i, r) = q) = \Pr(Q_s(j, r) = q)$$

where the probabilities are taken over uniformly chosen $r \in \{0, 1\}^{\ell_{rnd}}$

Data Privacy

User privacy

- (Information Theoretic) PIR: k copies of the database (not being intercommunicated)
 - Variations.
 - ★ Protocols can be defined to coalitions of up to $t < k$ servers

Data Privacy

User privacy

- Computational PIR (cPIR): privacy against one single database
 - The server has limited computational capacity
 - ★ **The computations** the server has to perform in order to gather enough information on the searches of a user to vulnerate her privacy, **exceeds the capacity** of the server.

Data Privacy

User privacy

- Computational PIR (cPIR): privacy against one single database
 - First approaches:
 - (Chor, Gilboa, 1997) For every $0 < c < 1$ there is a cPIR scheme for $k = 2$ DB with **communication complexity** $O(n^c)$.
 - (Kushilevitz, Ostrovsky, 1997) For every $c > 0$ there exists a single-database cPIR scheme with **communication complexity** $O(n^c)$, assuming the hardness of deciding quadratic residuosity¹. Linear time for the DB with respect to the number of rows.
 - They present a basic scheme and a recursive scheme

¹Given (x, N) where N is a composite number, it is difficult to determine whether x is a quadratic residue modulo N (i.e., $x = y^2 \pmod N$ for a certain y).

Data Privacy

User privacy

- Trusted-hardware Private Information Retrieval
(hardware-based Private Information Retrieval)
 - PIR protocols based on the assumption of a trusted hardware

Data Privacy

User privacy

- Other systems
 - Goopir: A user masks the query with $k - 1$ fake queries (example: change w_1 by $w_1 \text{ or } w_2 \text{ or } \dots \text{ or } w_k$) and submit the query to the search engine
 - ★ It assumes that **frequencies** of keywords and phrases that appear in a query **are known in advance**.
 - the frequencies of the target and the fake queries should be similar
 - so that the uncertainty of the search engine about the real target query is maximum
 - maximum privacy

Data Privacy

User privacy

- Other systems
 - TrackMeNot: A plugin for Firefox that periodically issues search queries
 - it hides the users actual search trails in a cloud of *ghost* queries.
 - ★ Generalization of its use: overhead of ghost queries
 - ★ Automatic ghost queries might be distinguishable and provide clues

Data Privacy

User privacy

- Private Information Retrieval (PIR)
- **Anonymous database search**

Data Privacy

User privacy

- Anonymous database search
 - How a user should retrieve an element from a database or a search engine without the system or the server being able to deduce who the retrieving user is.
 - It does not hide the content of the query, but obstructs the possibilities for the database of profiling users.

Data Privacy

User privacy

- P2P UPIR: Peer-to-peer User-Private Information Retrieval
 - Users submit queries on behalf of other users
 - The way in which users share communication spaces (memory sectors and cryptographic keys) is defined using combinatorial configurations
 - P2P UPIR offers privacy versus peer users

Data Privacy

P2P UPIR: Peer-to-peer User-Private Information Retrieval

- Communities of users and communication space: case 1
 - one memory sector and one cryptographic key
 - ★ all write and read
 - ★ the DB cannot know who is asking what: no profiling (except for the group)
 - but, no privacy between users
- The user does not know who made the query, but all queries are known

Data Privacy

P2P UPIR: Peer-to-peer User-Private Information Retrieval

- Communities of users and communication space: case 2
 - each user shares a different communication space with every other user
 - ★ every user only reads requests from “neighbours”
 - The user knows who requested a query, and its content
 - Not all the queries are known

Data Privacy

P2P UPIR: Peer-to-peer User-Private Information Retrieval

- Communities of users and communication space: case 3
 - different communication spaces for different users
 - ★ n_c communication spaces
 - with a memory sector and a cryptographic key
 - ★ n_u a set of users
 - all of them having access to a subset of d_u communication spaces
 - so that every communication space is shared by d_c users
 - and every pair of users share at most one communication space

case1

$$n_c = 1 \text{ (one space)}$$

$$d_u = 1 \text{ (one space per user)}$$

$$d_c = n_u \text{ (the only space is shared by all users)}$$

case2

$$n_c = \frac{n_u(n_u-1)}{2} \text{ (one space for each pair)}$$

$$d_u = n_u - 1 \text{ (for each user, one space for each other user)}$$

$$d_c = 2 \text{ (each space: only two users)}$$

Data Privacy

Respondent and owner privacy

Data-driven
(general-purpose)

Computation-driven
(specific-purpose)

Result-driven

Number of sources

Single data source

Multiple data sources

User privacy

Protecting the identity of the user

Anonymous database search

P2PUPIR

Protecting the data generated by
the activity of the user

PIR, ThPIR, cPIR

GooPIR
TrackMeNot

Data Privacy

Classification 1: On whose privacy is being sought

- Respondent privacy
- Owner privacy
- User privacy
 - Private Information Retrieval (PIR)
 - Anonymous database search

Dimensions. 2nd classification

On the computations to be done

Data Privacy

Respondent and owner privacy

- Data-driven or general-purpose
- Computation-driven or specific-purpose
- Result-driven

Data Privacy

Result-driven

- **Prevent** data mining procedures **infer some knowledge** that is valuable for the database owner
- Other uses: avoid discriminatory knowledge inferred from databases

Data Privacy

Result-driven

- **Formalization.** Database \mathcal{D} , A data mining algorithm, with parameters Θ is said to have ability to derive knowledge K from \mathcal{D} if and only if K is obtained from the output of the algorithm. Notation: $(A, \mathcal{D}, \Theta) \vdash K$.
- Any knowledge K such that $(A, \mathcal{D}, \Theta) \vdash K$ is in $KSet_{\mathcal{D}}$.

Data Privacy

Result-driven

- **Formalization.** Database \mathcal{D} , A data mining algorithm, with parameters Θ is said to have ability to derive knowledge K from \mathcal{D} if and only if K is obtained from the output of the algorithm. Notation: $(A, \mathcal{D}, \Theta) \vdash K$.
- Any knowledge K such that $(A, \mathcal{D}, \Theta) \vdash K$ is in $KSet_{\mathcal{D}}$.

Definition. \mathcal{D} a database, $\mathcal{K} = \{K_1, \dots, K_n\}$ sensitive knowledge to be hidden. The problem of hiding knowledge \mathcal{K} from \mathcal{D} consists on transforming \mathcal{D} into a database \mathcal{D}' such that

1. $\mathcal{K} \cap KSet_{\mathcal{D}} = \emptyset$
2. the information loss from \mathcal{D} to \mathcal{D}' is minimal

Data Privacy

Result-driven for **association rules mining**: Association rule hiding

- Recall that rules are mined when

$$\textit{Support}(R) \geq \textit{thr} - s$$

and

$$\textit{Confidence}(R) \geq \textit{thr} - c$$

for certain thresholds $\textit{thr} - s$ and $\textit{thr} - c$.

Data Privacy

Result-driven for **association rules mining**: Association rule hiding

- Recall that rules are mined when

$$\textit{Support}(R) \geq \textit{thr} - s$$

and

$$\textit{Confidence}(R) \geq \textit{thr} - c$$

for certain thresholds $\textit{thr} - s$ and $\textit{thr} - c$.

Two approaches:

- To reduce the support of the rule.
- To reduce the confidence of the rule.

Data Privacy

Result-driven for association rules mining: example

- **A formalization.** \mathcal{D} a database; $thr - s$ threshold. Let $\mathcal{K} = \{K_1, \dots, K_n\}$ sensitive itemsets, \mathcal{A} non-sensitive itemsets.

Data Privacy

Result-driven for association rules mining: example

- **A formalization.** \mathcal{D} a database; $thr - s$ threshold. Let $\mathcal{K} = \{K_1, \dots, K_n\}$ sensitive itemsets, \mathcal{A} non-sensitive itemsets.
- Transform $\mathcal{D} \rightarrow \mathcal{D}'$ such that
 1. $Support_{\mathcal{D}}(K) < thr - s$ for all $K_i \in \mathcal{K}$
 2. The number of itemsets K in \mathcal{A} such that $Support_{\mathcal{D}}(K) < thr - s$ is minimized.

This problem is NP-hard (Atallah et al., 1999)

Because of this: heuristic approaches

Data Privacy

Result-driven for association rules mining: heuristic algorithm

- **Algorithm.**

While HI is not hidden **do**

$HI' = HI$;

While $|HI'| > 2$ **do**

$P =$ subsets of HI with cardinality $|HI'| - 1$;

$HI' = \arg \max_{hi \in P} Support(hi)$;

$T_s =$ transaction in T supporting HI that affects the minimum number of itemsets of cardinality 2;

Set $HI' = 0$ in T_s ; Propagate results forward;

Data Privacy

Result-driven for association rules mining: heuristic algorithm

- **Algorithm.**

While HI is not hidden **do**

$HI' = HI$;

While $|HI'| > 2$ **do**

$P =$ subsets of HI with cardinality $|HI'| - 1$;

$HI' = \arg \max_{hi \in P} Support(hi)$;

$T_s =$ transaction in T supporting HI that affects the minimum number of itemsets of cardinality 2;

Set $HI' = 0$ in T_s ; Propagate results forward;

The algorithm does not cause false positives,
only false negatives (rules no longer inferred)

Data Privacy

Result-driven for association rules mining: heuristic algorithm

- **Example.** Computation of the algorithm to hide $HI = \{a, b, c\}$.

Data Privacy

Result-driven for association rules mining: heuristic algorithm

- **Example.** Computation of the algorithm to hide $HI = \{a, b, c\}$.

Transaction number	Items
T1	a, b, c, d
T2	a, b, c
T3	a, c, d

Data Privacy

Result-driven for association rules mining: heuristic algorithm

- **Example.** Computation of the algorithm to hide $HI = \{a, b, c\}$.

Transaction number	Items
T1	a, b, c, d
T2	a, b, c
T3	a, c, d

- Subsets of HI with cardinality $|HI| - 1$: $\{a, b\}$, $\{b, c\}$, $\{a, c\}$.

Data Privacy

Result-driven for association rules mining: heuristic algorithm

- **Example.** Computation of the algorithm to hide $HI = \{a, b, c\}$.

Transaction number	Items
T1	a, b, c, d
T2	a, b, c
T3	a, c, d

- Subsets of HI with cardinality $|HI| - 1$: $\{a, b\}$, $\{b, c\}$, $\{a, c\}$.
- $Support(\{a, b\}) = Support(\{b, c\}) = 2$, and $Support(\{a, c\}) = 3$

Data Privacy

Result-driven for association rules mining: heuristic algorithm

- **Example.** Computation of the algorithm to hide $HI = \{a, b, c\}$.

Transaction number	Items
T1	a, b, c, d
T2	a, b, c
T3	a, c, d

- Subsets of HI with cardinality $|HI| - 1$: $\{a, b\}$, $\{b, c\}$, $\{a, c\}$.
- $Support(\{a, b\}) = Support(\{b, c\}) = 2$, and $Support(\{a, c\}) = 3$
→ We select $HI' = \{a, c\}$.

Data Privacy

Result-driven for association rules mining: heuristic algorithm

- **Example.** Computation of the algorithm to hide $HI = \{a, b, c\}$.

Transaction number	Items
T1	a, b, c, d
T2	a, b, c
T3	a, c, d

- Subsets of HI with cardinality $|HI| - 1$: $\{a, b\}$, $\{b, c\}$, $\{a, c\}$.
- $Support(\{a, b\}) = Support(\{b, c\}) = 2$, and $Support(\{a, c\}) = 3$
 → We select $HI' = \{a, c\}$.
- Set of transactions in T that support HI (and HI'): $\{T1, T2\}$.

Data Privacy

Result-driven for association rules mining: heuristic algorithm

- **Example.** Computation of the algorithm to hide $HI = \{a, b, c\}$.

Transaction number	Items
T1	a, b, c, d
T2	a, b, c
T3	a, c, d

- Subsets of HI with cardinality $|HI| - 1$: $\{a, b\}$, $\{b, c\}$, $\{a, c\}$.
- $Support(\{a, b\}) = Support(\{b, c\}) = 2$, and $Support(\{a, c\}) = 3$
 → We select $HI' = \{a, c\}$.
- Set of transactions in T that support HI (and HI'): $\{T1, T2\}$.
- T 's transaction in $\{T1, T2\}$ that affects the minimum number of itemsets of cardinality 2: $T2$ affects less itemsets than $T1$.

Data Privacy

Result-driven for association rules mining: heuristic algorithm

- **Example.** Computation of the algorithm to hide $HI = \{a, b, c\}$.
- Remove one of the items in $HI' = \{a, c\}$ that are in $T2$:

Data Privacy

Result-driven for association rules mining: heuristic algorithm

- **Example.** Computation of the algorithm to hide $HI = \{a, b, c\}$.
- Remove one of the items in $HI' = \{a, c\}$ that are in $T2$:
Both have the same support, we select one of them at random.
- Propagate the results forward: recompute supports

Data Privacy

Respondent and owner privacy

- Data-driven or general-purpose
- Computation-driven or specific-purpose
- Result-driven

Data Privacy

Computation-driven approaches: centralized

- **Example.** Parties P_1, \dots, P_n own databases DB_1, \dots, DB_n . The parties want to compute a function, say f , of these databases (i.e., $f(DB_1, \dots, DB_n)$) without revealing unnecessary information. In other words, after computing $f(DB_1, \dots, DB_n)$ and delivering this result to all P_i , what P_i knows is nothing more than what can be deduced from his DB_i and the function f .
- So, the computation of f has not given P_i any extra knowledge.

Data Privacy

Computation-driven approaches: distributed

- Vertically partitioned data
- Horizontally partitioned data

Data Privacy

Computation-driven approaches

- The centralized approach as a reference

Data Privacy

Computation-driven approaches.

Privacy leakage for the distributed approach is usually analyzed considering two types of **adversaries**.

Data Privacy

Computation-driven approaches.

Privacy leakage for the distributed approach is usually analyzed considering two types of **adversaries**.

- **Semi-honest adversaries.** Data owners follow the cryptographic protocol but they analyse all the information they get during its execution to discover as much information as they can.
- **Malicious adversaries.** Data owners try to fool the protocol (e.g. aborting it or sending incorrect messages on purpose) so that they can infer confidential information.

Data Privacy

Respondent and owner privacy

- Data-driven or general-purpose
- Computation-driven or specific-purpose
- Result-driven

Data Privacy

Data-driven approaches.

Data Privacy

Data-driven approaches.

- Reduction of quality
 - **Information loss** measures or utility measures

Data Privacy

Data-driven approaches.

- Reduction of quality
 - **Information loss** measures or utility measures
- Not ensuring complete privacy
 - **Disclosure risk** measures

Data Privacy

Data-driven approaches.

- Reduction of quality
 - **Information loss** measures or utility measures
- Not ensuring complete privacy
 - **Disclosure risk** measures
- Trade-off and visualization
 - Tools for visualization and finding the trade-off

Data Privacy

Data-driven approaches. Topics

- (a) which methods are available ?
 - **Methods:** Perturbative methods, non-perturbative methods, synthetic data generators
- (b) do methods keep the most relevant information in the data?
 - **Information Loss (IL):** in what extent data is still useful?
- (c) is the resulting data still acceptable w.r.t. privacy?
 - **Disclosure Risk (DR):** in what extent data does not lead to disclosure?
- (d) How to visualize the contradiction between DR and IL?
 - **Visualization tools:** R-U maps, scores

Dimensions. 3rd classification

On the number of data sources

Data Privacy

Number of data sources

- Single data source
- Multiple data source
 - Usually computation-driven
 - Centralized approach vs. distributed approach
 - ★ Trusted Third Party (TTP) vs. PPDM

Dimensions

Data Privacy

Respondent and owner privacy

Data-driven
(general-purpose)

Computation-driven
(specific-purpose)

Result-driven

Number of sources

Single data source

Multiple data sources

User privacy

Protecting the identity of the user

Protecting the data generated by
the activity of the user

Data Privacy

Yet other classifications

- Perturbative vs. Cryptographic approaches
 - Mainly data-driven vs. Computation-driven
- Type of data
 - Files and databases
 - Tabular data (aggregates)
 - Logs
 - ...