# Disclosure, Privacy Models, and Privacy Mechanisms

Vicenç Torra

January 2024

Umeå University, Sweden

# Outline

- Privacy models

  - Definition
  - Summary
  - Privacy from re-identification
  - $k$-Anonymity
  - Differential privacy
  - Homomorphic encryption
  - Secure multiparty computation
  - Result privacy
  - Summary
  - An example: Integral privacy
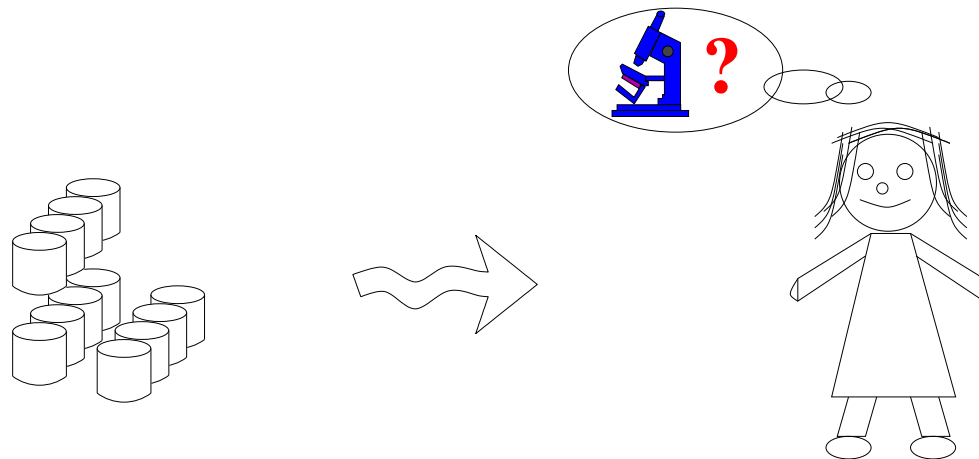
# Introduction

# Definition

# Disclosure

## Definition

- A privacy model is a computational definition of privacy.

# Summary of privacy models

# Privacy models
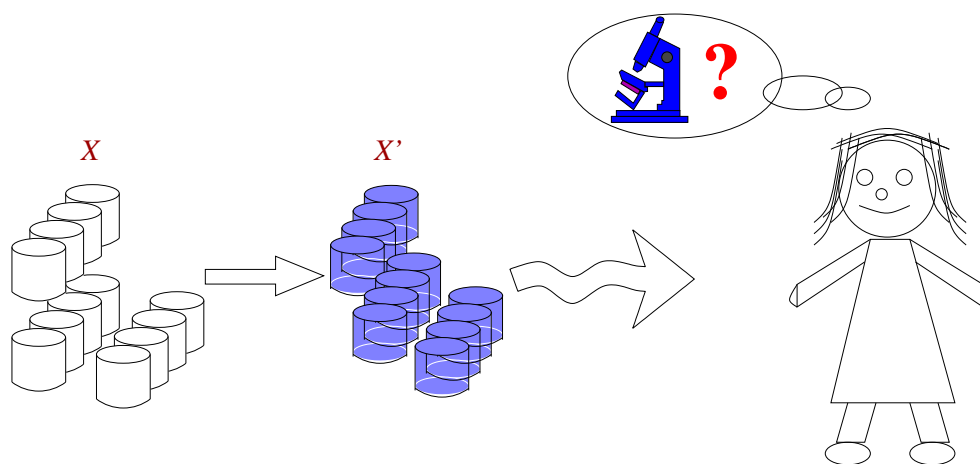
# Privacy models

**Privacy models.** A computational definition for privacy. Examples.

- **Reidentification privacy.** Avoid finding a record in a database.
- **k-Anonymity.** A record indistinguishable with $k - 1$ other records.
- **Secure multiparty computation.** Several parties want to compute a function of their databases, but only sharing the result.
- **Differential privacy.** The output of a query to a database should not depend (much) on whether a record is in the database or not.
- **Result privacy.** We want to avoid some results when an algorithm is applied to a database.
- **Integral privacy.** Inference on the databases. E.g., changes have been applied to a database.
- **Homomorphic encryption.** We want to avoid access to raw data and partial computations.

# Privacy models

**Privacy models.** A computational definition for privacy. Publish a DB

- **Reidentification privacy.** Avoid finding a record in a database.
- **k-Anonymity.** A record indistinguishable with $k - 1$ other records.
- **k-Anonymity, l-diversity.** $l$ possible categories
- **Interval disclosure.** The value for an attribute is outside an interval computed from the protected value: values different enough.
- **Result privacy.** We want to avoid some results when an algorithm is applied to a database.

# Privacy models

**Privacy models.** A computational definition for privacy. Publish a DB

- Modify DB X to obtain a DB X' compliant with the privacy model.

Original DB $X$:

| Respondent | City | Age | Illness |
|---|---|---|---|
| DRR | Barcelona | 30 | Heart attack |
| ABD | Barcelona | 32 | Cancer |
| COL | Barcelona | 33 | Cancer |
| GHE | Tarragona | 62 | AIDS |
| CIO | Tarragona | 65 | AIDS |
| HYU | Tarragona | 60 | Heart attack |

Published DB $X'$:

| | City | Age | Illness |
|---|---|---|---|
| — | Barcelona | 30 | Cancer |
| — | Barcelona | 30 | Cancer |
| — | Barcelona | 30 | Cancer |
| — | Tarragona | 60 | AIDS |
| — | Tarragona | 60 | AIDS |
| — | — | — | — |

# Privacy models

- Difficulties

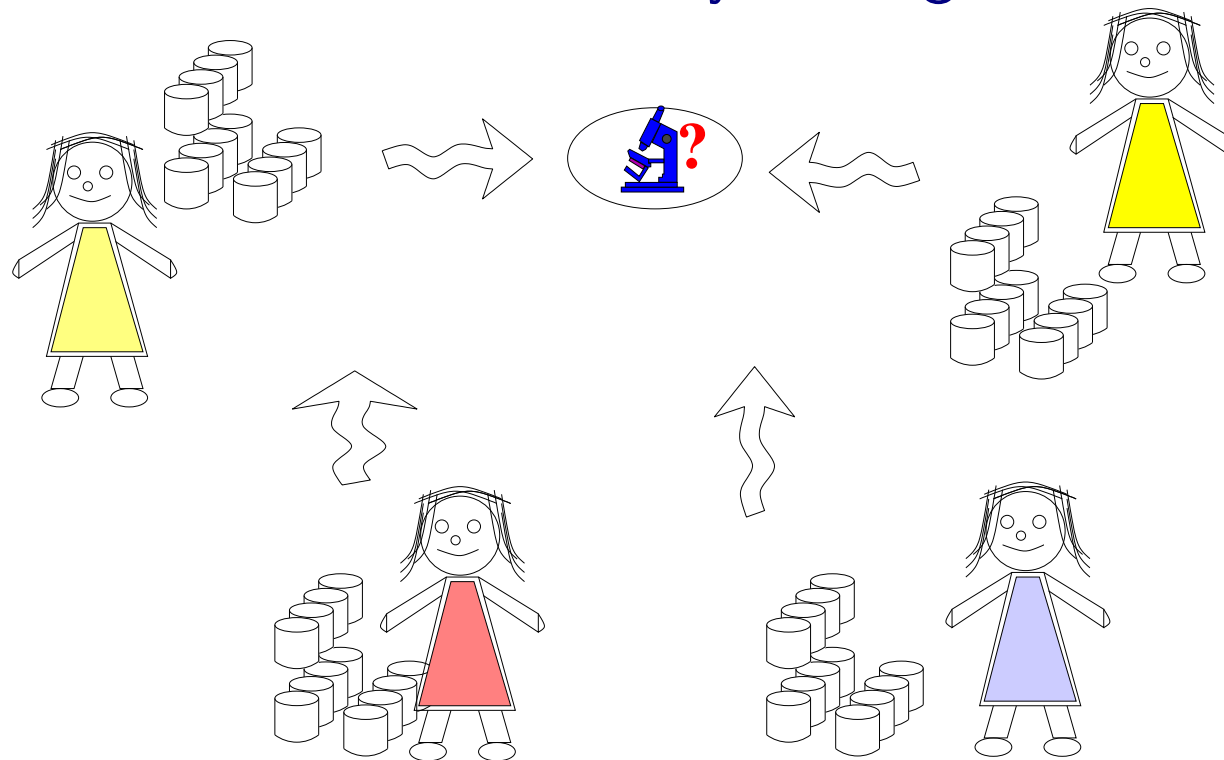  Naive anonymization does not work, highly identifiable data, high dimensional data

- Examples of successful reidentification attacks

  Sweeney analysis of USA population, data from mobile data, shopping cards, film ratings

# Privacy models

**Privacy models.** A computational definition for privacy. Share a result

- **Secure multiparty computation.** Several parties want to compute a function of their databases, but only sharing the result.

# Privacy models

**Privacy models.** A computational definition for privacy. Share a result

- Compute

$$f(DB_1, DB_2, DB_3, DB_4)$$

  without sharing $DB_1, DB_2, DB_3, DB_4$
- Example: national age mean of hospital-acquired infection patients (hospitals do not want to share the age of their infected patients!)
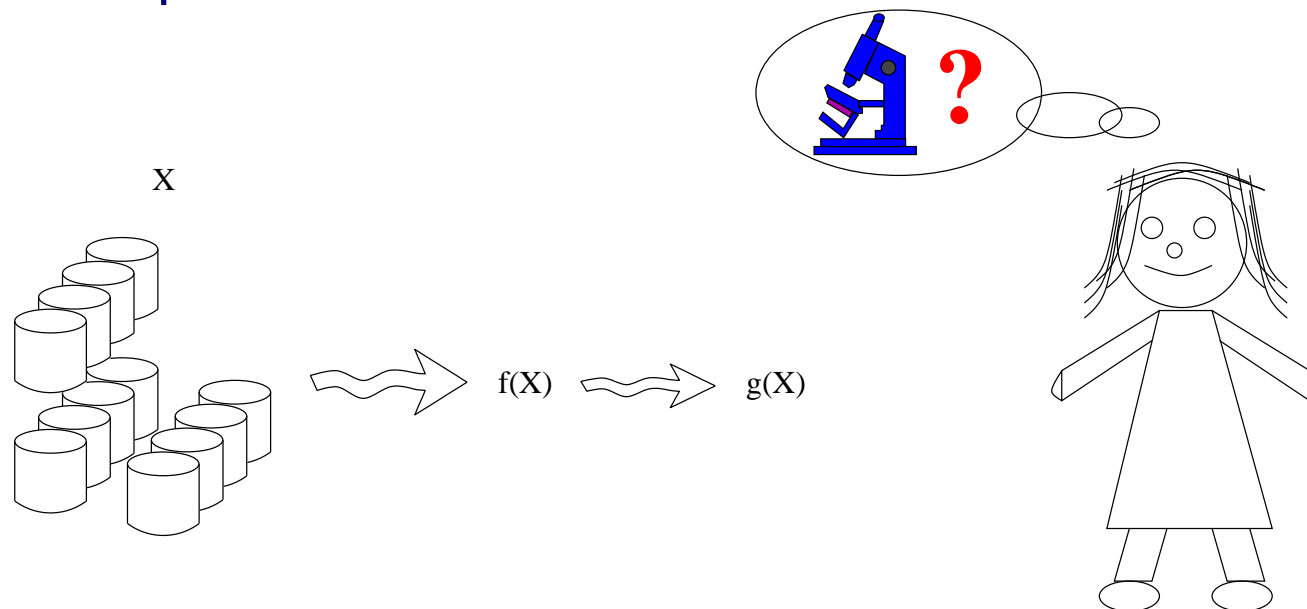
# Privacy models

- Difficulties

  Distributed approach (no trusted-third party) – computational cost of solutions

# Privacy models

**Privacy models.** A computational definition for privacy. Compute result

- **Differential privacy.** The output of a query to a database should not depend (much) on whether a record is in the database or not.
- **Integral privacy.** Inference on the databases. E.g., changes have been applied to a database.
- **Homomorphic encryption.** We want to avoid access to raw data and partial computations.

# Privacy models

- Difficulties. A simple function can give information on who is in the database

  ○ E.g., mean salary

# Privacy from re-identification

# Privacy from re-identification

- A protected database $A$ satisfies privacy from re-identification given intruder's knowledge $B$ when

$$Reid(B, A) \leq r_{R1}$$

  with a certain privacy level $r_{R1}$ (e.g., $r_{R1} = 0.25$),

- or, alternatively

$$KR.Reid(B, A) \leq (r_K, r_{R1})$$

  with certain privacy levels $r_K$ and $r_{R1}$ (e.g., $r_K = 0$ and $r_{R1} = 0.5$).

# $k$-**Anonymity**

# $k$-**Anonymity**

## Definition 3.4

- A database $A$ satisfies $k$-anonymity with respect to a set of quasi-identifiers $QI$ when the projection of $A$ in this set $QI$ results into a partition of $DB$ in sets of at least $k$ indistinguishable records.

| City | Age | Illness |
|------|-----|---------|
| Barcelona | 30 | Cancer |
| Barcelona | 30 | Cancer |
| Tarragona | 60 | AIDS |
| Tarragona | 60 | AIDS |

# $k$-**Anonymity**

- Indistinguishability w.r.t. quasi-identifiers

- $k$-Anonymity and re-identification

$$KR.Reid(B, A) \leq (0, 1/k).$$

- Plausible deniability

# $k$-**Anonymity**

- Indistinguishability w.r.t. quasi-identifiers

- $k$-Anonymity and re-identification

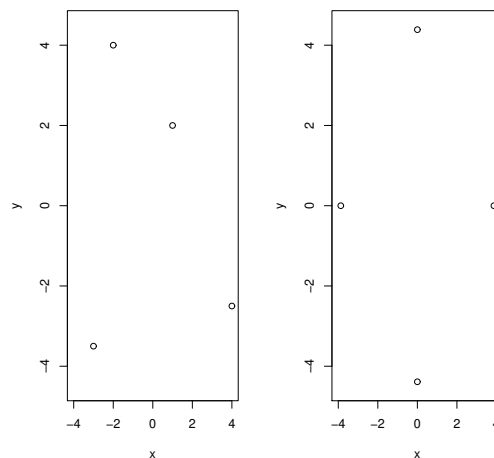$$KR.Reid(B, A) \leq (0, 1/k).$$

- Plausible deniability

  - at record level
  - but not at database level

- Records are independent

# $k$-**Anonymity**

- $k$-confusion. Drop indistinguishability

- Example

  ○ Original data: $X = \{(1,2), (-2,4), (4,-2), (-3,-3)\}$.
  ○ $k$-Anonymity: $X' = \{(0,0), (0,0), (0,0), (0,0)\}$.
  ○ $k$-Confusion: using $X'' = \{(x,0), (-x,0), (0,y), (0,-y)\}$,
         with standard deviations in $X''$ equal to the ones in $X$
    ▷ $x = \sqrt{10}/\sqrt{2/3} = 3.872983$, $y = \sqrt{12.8333}/\sqrt{2/3} = 4.387476$



- Discussion: $k$-confusion and re-identification

# $k$-**Anonymity**

- Attacks

  - Homogeneity attack (external attack)
  - External knowledge attack (internal attack)

- These are attribute disclosure attacks

  - while $k$-anonymity is for identity disclosure

- Variations of $k$-anonymity to avoid attribute disclosure

# $k$-**Anonymity**

- $p$-sensitive $k$-anonymity for $k > 1$ and $p \leq k$

  ○ if it satisfies $k$-anonymity and, for each group of records with the same combination of values for a set of quasi-identifiers, the number of distinct values for each confidential value is at least $p$ (within the same group).

- $l$-diversity

  ○ forces $l$ different categories in each set. However, in this case, categories should have to be well-represented. Different meanings have been given to what well-represented means.

- $t$-closeness.

  ○ The distribution of the attribute in any $k$-anonymous subset of the database is similar to the one of the full database. Similarity: distance

between the two distributions, distance below a given threshold $t$. The Earth Mover distance is used in the definition.

# *k*-**Anonymity**

- $k$-anonymity and computational anonymity

  - Relaxation: not-all quasi-identifiers
    "We say that unconditional anonymity is theoretical anonymity. Computational anonymity is conditioned by the assumption that the adversary has some limitation. The limitations can be (...) restricted memory or knowledge." (Stokes (2012)).
  - A data set $X$ satisfies $(k, l)$-anonymity if it is $k$-anonymous with respect to every subset of attributes of cardinality at most $l$.
    $\Rightarrow$ Intruder's knowledge limited to $l$ attributes

- Example: (2,2)-anonymity

$$D = \{(a, b, e), (a, b, f), (c, d, e), (c, d, f),$$

$$(c, b, e), (c, b, f), (a, d, e), (a, d, f)\}.$$

# Differential privacy

# Differential privacy

- Computation-driven/single database

  - Privacy model: differential privacy[1]
  - We know the function/query to apply to the database: $f$

- Example:

  compute the mean of the attribute salary of the database for all those living in Town.

---

[1]There are other models as e.g. query auditing (determining if answering a query can lead to a privacy breach), and integral privacy

# Differential privacy

- **Differential privacy** (Dwork, 2006).

  - Motivation:
    - ▷ the result of a query should not depend on the presence (or absence) of a particular individual
    - ▷ the impact of any individual in the output of the query is limited
      differential privacy ensures that the removal or addition of a single database item does not (substantially) affect the outcome of any analysis (Dwork, 2006)

# Differential privacy

- Mathematical definition of differential privacy (in terms of a probability distribution on the range of the function/query)

  ○ A function $K_q$ for a query $q$ gives $\epsilon$-differential privacy if for all data sets $D_1$ and $D_2$ differing in at most one element, and all $S \subseteq Range(K_q)$,

  $$\frac{Pr[K_q(D_1) \in S]}{Pr[K_q(D_2) \in S]} \le e^{\epsilon}.$$

  (with 0/0=1) or, equivalently,

  $$Pr[K_q(D_1) \in S] \le e^{\epsilon} Pr[K_q(D_2) \in S].$$

- $\epsilon$ is the level of privacy required (privacy budget). The smaller the $\epsilon$, the greater the privacy we have.

# Differential privacy

- Differential privacy

  - A function $K_q$ for a query $q$ gives $\epsilon$-differential privacy if ...
    - $\triangleright$ $K_q(D)$ ~~is a constant. E.g.,~~
    $$\cancel{K_q(D) = 0}$$
    - $\triangleright$ $K_q(D)$ is a randomized version of $q(D)$:
    $$K_q(D) = q(D) + \text{and some appropriate noise}$$



Kq(D)

# Differential privacy

- **Properties**

  - Plausible deniability: to an extend, in terms of $\epsilon$

# Differential privacy: Variations of differential privacy

- **Def. 3.17.** $(\epsilon, \delta)$-differential privacy (or $\delta$-approximate $\epsilon$-indistinguishability)

  ○ A function $K_q$ for a query $q$ gives $(\epsilon, \delta)$-differential privacy if for all data sets $D_1$ and $D_2$ differing in at most one element, and all $S \subseteq Range(K_q)$,

$$Pr[K_q(D_1) \in S] \leq e^\epsilon Pr[K_q(D_2) \in S] + \delta.$$
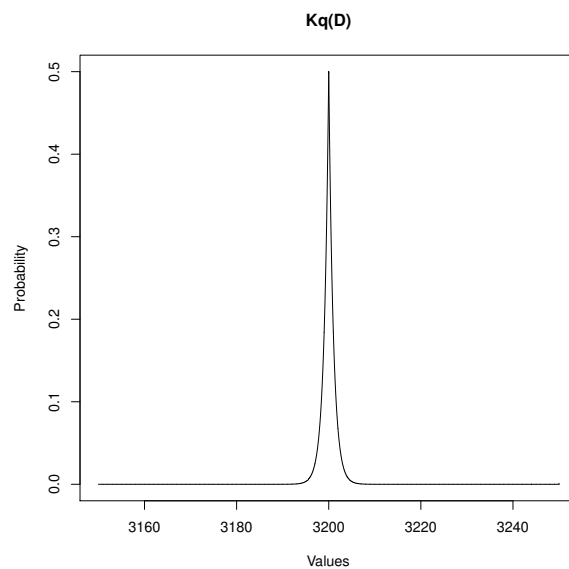
- Relaxes $\epsilon$-DP, events with a probability smaller than $\delta$ for $D_1$ are still permited even if they do not occur in $D_2$.

# Differential privacy: Variations of differential privacy

- Bounded differential privacy

  ○ The two neighboring datasets have exactly the same number of records.

# Differential privacy: Budgets

- Multiple queries and budget consumption

# Local Differential privacy

- When $q = Id$

  - That is, we want to deliver $X$, so, we provide $X' = \rho(X)$

# Local Differential privacy

- When $q = Id$

  - That is, we want to deliver $X$, so, we provide $X' = \rho(X)$

- At record level / collection level

  - Same definition applies

# Local Differential privacy

- When $q = Id$

  ○ That is, we want to deliver $X$, so, we provide $X' = \rho(X)$

- At record level / collection level

  ○ Same definition applies

- Problem

  ○ Multiple communications from the same device:
  $$x_i^1, \ldots, x_i^T$$
  provide
  $$\rho(x_i^1), \ldots, \rho(x_i^T)$$
  but, they are not independent ...

# Homomorphic encryption

# Homomorphic encryption

- Homomorphic encryption. We want to avoid access to raw data and partial computations.

  - A single database $DB$ and a function $f$. The only information learn is $f(DB)$. No other leakage is permitted. No access to the data, no acess to partial computations (i.e., similar to SMC but for a single database).
  - This allows us to store data in the cloud. No leakage during data storage, no leakage during transmissions.

# Secure Multiparty Computation

# Secure multiparty computation

- **Def. 3,18.**

  - Let $P_1, \ldots, P_n$ represent a set of parties, and let $X_1, \ldots, X_n$ be their respective databases. The parties want to compute a function $f$ of these databases (i.e., $f(X_1, \ldots, X_n)$) without revealing unnecessary information.

- After computing $f(X_1, \ldots, X_n)$ and delivering this result to each $P_i$, what $P_i$ knows is nothing more than what can be deduced from $X_i$ and the function $f$. So, the computation of $f$ has not given $P_i$ any extra knowledge.

# Secure multiparty computation

- Trivial approach: Centralized approach

  ○ Trusted third party $TTP$ that computes the analysis.
    Each $P_i$ transfers data $X_i$ using a completely secure channel (e.g., using cryptographic protocols) to the trusted third party $TTP$. Then, $TTP$ computes the result $y = f(X_1, \ldots, X_n)$, and sends $y$ to each $P_i$ in a secure way. This will satisfy the definition as each $P_i$ knows nothing more than $X_i$ and $y$.

- Secure multiparty computation provides solutions for this problem in a distributed environment (no trusted third party). Same privacy guarantees are sought.

# Secure multiparty computation

- Privacy-preserving solutions

  ○ Protocols that describe the information flow among the parties and details their computations.
  ○ Assumptions are needed on the behavior of the intruders.
    The parties themselves can be intruders trying to gain some extra knowledge from their computations. We can even consider parties that try to fool the other parties, break the protocol, and collide with others to learn relevant information from a targeted party.

# Result privacy

# Result privacy

- Result privacy. Given a database, avoid inferring knowledge $K$.

  ○ **Def.** $X$ a database, $A$ a parametric data mining algorithm. $A$ with parameter $\Theta$ is said to have ability to derive knowledge $K$ from $X$ if and only if $K$ appears either directly in the output of the algorithm or by reasoning from the output.

  $$(A, X, \Theta) \vdash K$$

  ○ $K$ is said to be derivable from $X$, if there exists any algorithm $A$ with parameter $\Theta$ such that $(A, X, \Theta) \vdash K$.

# Result privacy

- Result privacy. Given a database, avoid inferring knowledge $K$.

  - **Def.** $X$ a database, $A$ a parametric data mining algorithm. $A$ with parameter $\Theta$ is said to have ability to derive knowledge $K$ from $X$ if and only if $K$ appears either directly in the output of the algorithm or by reasoning from the output.

  $$(A, X, \Theta) \vdash K$$

  - $K$ is said to be derivable from $X$, if there exists any algorithm $A$ with parameter $\Theta$ such that $(A, X, \Theta) \vdash K$.
  - **Def.** $X$, $A$, $\Theta$ as above. $A$ with $\Theta$ is said to satisfy result privacy with respect to a set of sensitive knowledge $\mathcal{K} = \{K_1, \ldots, K_n\}$ when no $K_i$ in $\mathcal{K}$ is derivable from $X$.
    no $K_i$ is such that $(A, X, \Theta) \vdash K_i$.

# Summary

# Summary

| Privacy risk model/measure | Attribute disclosure | Identity disclosure | database release | query release | Boolean |
|---|---|---|---|---|---|
| Re-identification | | X | X | | Quantitative |
| Uniqueness | | X | X | | Quantitative |
| Result-driven | X | | X | | Boolean |
| $k$-Anonymity | | X | X | | Boolean |
| $k$-confusion | | X | X | | Boolean |
| $k$-concealment | | X | X | | Boolean |
| $p$-sensitive $k$-Anonymity | X | X | X | | Boolean |
| $k$-Anonymity, $l$-diversity | X | X | X | | Boolean |
| $k$-Anonymity, $t$-closeness | X | X | X | | Boolean |
| Interval disclosure | X | | X | | Quantitative |
| Differential privacy | X | | | X | Boolean |
| Local differential privacy | | X | X | | Boolean |
| Integral privacy | X | | | X | Boolean |
| Homomorphic encryption | X | | | X | Boolean |
| Secure multiparty computation | X | | | X | Boolean |

# An example: integral privacy

# Integral privacy

- Given a $DB$ and a function $f$, is $f(DB)$ recurrent?

  - If we consider possible databases, are we going to obtain $f(DB)$ often?
  - A $k$-anonymity flavor for $f(DB)$

# Integral privacy

Some preliminaries ...

- $P$ the population, $A$ be an algorithm that given a data set $S \subseteq P$ computes an output $A(S)$ that belongs to another domain $\mathcal{G}$.
- Given $G$ in $\mathcal{G}$, previous knowledge $S^*$ with $S^* \subset P$, the set of possible generators of $G$ is:

$$Gen(G, S^*) = \{S'|S^* \subseteq S' \subseteq P, A(S') = G\}.$$

We use $Gen^*(G, S^*) = \{S' \setminus S^*|S^* \subseteq S' \subseteq P, A(S') = G\}$
(when no information is known on $S^*$, we use $S^* = \emptyset$

# Integral privacy

Integral privacy, definition:

- $P$ data, $A : S \to \mathcal{G}$, $S^*$ background knowledge, $Gen(G, S^*)$ databases that generate $G$ and are consistent with background knowledge $S^*$.

  Then, i-integral privacy is satisfied when $Gen(G, S^*)$ is <u>large</u> and

$$\cap_{g \in Gen^*(G, S^*)} g = \emptyset.$$

Our definition of privacy has a $k$-anonymity flavor (next slides)

Requirements: why? / what?

- Empty intersection to avoid all generators sharing a record (e.g., avoiding membership attacks)
- $Gen(G, S^*)$ <u>large</u>. What is <u>large</u> ????

# Integral privacy

Integral privacy, details and the $k$-anonymity flavor

- $Gen(G, S^*)$ large . . . 1st definition
  - At least $k$ elements, $+$ empty intersection
    $= k$ different databases not sharing records

# Integral privacy

Integral privacy, details and another $k$-anonymity flavor

- $Gen(G, S^*)$ large . . . 2nd definition
  - At least $k$ different minimal sets
    Example. 10 databases:
    5 DBs only share record $r$ and 5 other DBs only share record $r'$.
    Integrally private with $k = 2$.
  - $\Rightarrow$ we formalize this notion in this paper

# Integral privacy

Integral privacy, and <span style="color:red">plausible deniability</span>

- IP satisfies plausible deniability if for any record $r$ in $P$ such that $r \notin S^*$, there is a set/database $\sigma \in Gen^*(G, S^*)$ such that $r \notin \sigma$.

Our definition satisfies plausible deniability

# Integral privacy

Integral privacy, and differential privacy

- Differential privacy, smooth function
  $A(D) \sim A(D \oplus x)$ where $D \oplus x$ means to add the record $x$ to $D$
- Integral privacy, recurrent function
  If $A^{-1}(G)$ is the set of all (real) databases that can generate the output $G$, we require $A^{-1}(G)$ to be a large and diverse set for $G$.

# Integral privacy

Integral privacy, and differential privacy

- Differential privacy, smooth function
  $A(D) \sim A(D \oplus x)$ where $D \oplus x$ means to add the record $x$ to $D$
- Integral privacy, recurrent function
  If $A^{-1}(G)$ is the set of all (real) databases that can generate the output $G$, we require $A^{-1}(G)$ to be a large and diverse set for $G$.

- Simple integrally private function:
  $A$ an algorithm that is 1 if the number of records in $D$ is even, and 0 if the number of records in $D$ is odd.
  That is, $f(D) = 1$ if and only if $|D|$ is even.

# References

# Integral privacy

- Torra, V. (2022) Guide to data privacy, Springer.
- Samarati, P. (2001) Protecting Respondents' Identities in Microdata Release, IEEE Trans. on Knowledge and Data Engineering, 13:6 1010-1027.
- Truta, T. M., Vinay, B. (2006) Privacy protection: p-sensitive k-anonymity property. Proc. 2nd Int. Workshop on Privacy Data management (PDM 2006) p. 94.
- Li, N., Li, T., Venkatasubramanian, S. (2007) T-closeness: privacy beyond k-anonymity and l-diversity, Proc. of the IEEE ICDE 2007.
- Dwork, C. (2006) Differential privacy, Proc. ICALP 2006, LNCS 4052, pp. 1-12.
- Torra, V., Navarro-Arribas, G. (2016) Integral privacy, Proc. CANS 2016.