# Safe Disassociation of Set-Valued Datasets

Nancy Awad[1,2], Bechara AL Bouna[1], Jean-Francois Couchot[2], and Laurent Philippe[2]

[1] TICKET Lab., Antonine University, Hadat-Baabda, Lebanon.
nancy.awad,bechara.albouna@ua.edu.lb
[2] FEMTO-ST Institute, UMR 6174 CNRS, Université of Bourgogne Franche-Comté, France.
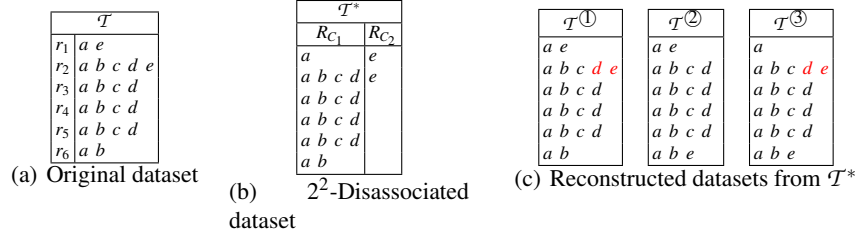jean-francois.couchot, laurent.philippe@univ-fcomte.fr

**Abstract.** Disassociation introduced by Terrovitis *et al.* is a bucketization based anonimyzation technique that divides a set-valued dataset into several clusters to hide the link between individuals and their complete set of items. It increases the utility of the anonymized dataset, but on the other side, it raises many privacy concerns, one in particular, is when the items are tightly coupled to form what is called, a cover problem. In this paper, we present safe disassociation, a technique that relies on partial-suppression, to overcome the aforementioned privacy breach encountered when disassociating set-valued datasets. Safe disassociation allows the $k^m$-anonymity privacy constraint to be extended to a bucketized dataset and copes with the cover problem. We describe our algorithm that achieves the safe disassociation and we provide a set of experiments to demonstrate its efficiency.

**Keywords:** Disassociation, cover problem, data privacy, set-valued, privacy preserving

## 1 Introduction

Privacy preservation is a key concern in data publishing where individual's personal information must remain protected under all circumstances. This sounds straightforward, but it is somehow difficult to achieve. The AOL search data leak in 2006 (BZ06) is an explicit example that shows the consequences of a unsupervised data publishing. The query logs of 650k individuals were released after omitting all explicit identifiers. They were later withdrawn due to multiple reports of attackers linking individuals to their sensitive records. Alternatively, providing a "complete" privacy over the data requires sacrifices in terms of utility, in other words, usefulness of the data (Sam01, Swe02, MGKV06, XT06, DMNS06). Hence, it is pointless to publish datasets that do not provide valuable information. A suitable trade-off between data utility and privacy must be achieved. The point is to provide not only a value anonymization technique, but instead a dataset anonymization technique, that hides/anonymizes the link between individuals and their sensitive information, and, at the same time, keeps the dataset useful for analysis. When publishing a set-valued dataset (e.g., shopping and search items) it is important to pay attention to attackers that try intentionally to link individuals to their sensitive information. These attackers may be able to single out an individual's complete record/itemset by associating data items from the dataset to their background knowledge. The example in Figure 1(a) shows a set-valued dataset $\mathcal{T}$ consisting of 6 records $r_1,\dots,r_6$, which are itemsets linked to individuals 1, ..., 6 respectively. For instance, $r_1 : \{a, e\}$ can be interpreted as individual

1 has searched for item $a$ and item $e$. If an attacker knows that individual 1 has searched for items $d$ and $e$, he/she will be able to link 1 to his record $r_2$.

**(a) Original dataset** — $\mathcal{T}$

| | |
|---|---|
| $r_1$ | a e |
| $r_2$ | a b c d e |
| $r_3$ | a b c d |
| $r_4$ | a b c d |
| $r_5$ | a b c d |
| $r_6$ | a b |

**(b) $2^2$-Disassociated dataset** — $\mathcal{T}^*$

| $R_{C_1}$ | $R_{C_2}$ |
|---|---|
| a | e |
| a b c d e | e |
| a b c d | |
| a b c d | |
| a b c d | |
| a b | |

**(c) Reconstructed datasets from $\mathcal{T}^*$**

$\mathcal{T}^{①}$

| |
|---|
| a e |
| a b c *d e* |
| a b c d |
| a b c d |
| a b c d |
| a b |

$\mathcal{T}^{②}$

| |
|---|
| a e |
| a b c d |
| a b c d |
| a b c d |
| a b c d |
| a b e |

$\mathcal{T}^{③}$

| |
|---|
| a |
| a b c *d e* |
| a b c d |
| a b c d |
| a b c d |
| a b e |

**Fig. 1.** Disassociation leading to a cover problem

Several techniques (Sam01, Swe02, XT06, LLZM12, DCdVFJ$^+$13, WWFW16, WDL18) have been defined in the literature to anonymize the dataset and cope with this particular association problem. Anonymization by disassociation (TMK08, LLGT14,LLGDT15) is a bucketisation technique (XT06,LLZM12,DCdVFJ$^+$13) that keeps the items without alternation/generalization but separates the records into clusters and chunks to hide their associations. More specifically, disassociation transforms the original data into $k^m$-anonymous clusters of chunks to ensure that an attacker who knows up to $m$ items cannot associate them with less than $k$ records. First, disassociation divides the dataset horizontally, creating clusters of records with similar frequent items. In Figure 1(a), the item $a$ has the highest number of occurrences in the records; thus, all records containing the item $a$ are added to the cluster. Next, disassociation divides the clusters vertically, creating inside each cluster $k^m$-anonymous record chunks of items, and one item chunk containing items that appear less than $k$ times. Figure 1(b) shows the result of disassociation with only one cluster, namely $\mathcal{T}^*$, containing 2 record chunks $R_{C_1}$, and $R_{C_2}$ and no item chunk. Both $R_{C_1}$ and $R_{C_2}$ are $2^2$-anonymous. That is, an attacker who knows any two items about an individual will be able link them to at least two records from the dataset.

In a previous work (BaBNG16), the disassociation technique has been evaluated and found to be vulnerable to what is called, a cover problem. This cover problem provides attackers with the ability to associate itemsets in consequent record chunks with individuals' records, and compromises the disassociated dataset. Figure 1(c) highlights all datasets $\mathcal{T}^{①}$, $\mathcal{T}^{②}$, and $\mathcal{T}^{③}$ that could be reconstructed from the disassociated dataset $\mathcal{T}^*$ by associating records from different record chunks. Suppose now the attacker knows that an individual has searched for items $d$ and $e$ , $\{d,e\}$. In such a case, he/she can remove $\mathcal{T}^{②}$ from the possible reconstructed dataset and will be able to link every record containing $d$ and $e$ with certainty to $\{a,b,c\}$ from the two remaining reconstructed datasets leading to a privacy breach.

This paper extends the previous work (BaBNG16) by providing a theoretical and practical solution to this cover problem. It introduces a privacy preserving technique to safely disassociate a set-valued dataset and address this cover problem. This method is further denoted as safe disassociation. Our contributions are summarized as follows:

- We define the privacy guarantee for safe disassociation and provide the appropriate algorithm to achieve it.
- We investigate the efficiency of safe disassociation and its impact on the utility of aggregate analysis and the discovery of association rules.

The rest of the paper is organized as follows. Section 2 presents an overview of some of the related works in set-valued dataset anonymization . In Section 3, we describe the formal model of disassociation and discuss the cover problem that makes disassociation vulnerable. We define safe disassociation in Section. 4. Experimental evaluations of this method is presented in Section 5. Finally, we conclude in Section 6 and present outlines for future work.

## 2    Related Work

Anonymization techniques can be divided into several categories, namely categorization, generalization, and bucketization, detailed hereafter.

With categorization ($JPX^+14$, $CLZ^+16$, WDL18), attributes are classified into several categories with respect to their sensitivity: identifying, sensitive or non-sensitive. On the opposite, all the data in this current work are considered to be grouped into sets whose items have the same level of sensitivity; the items combined together are sensitive and therefore only the link between the individual and their corresponding items must be protected.

Generalization techniques  (Sam01, Swe02, MGKV06) create homogeneous subsets by replacing values with ranges wide enough to create ambiguity. Reducing the extent of generalization, and thus decreasing the information loss is critical to preserving the usefulness of the data. In (HN09), only a local generalization is applied whereas in (FW10) a clustering based technique is implemented to minimize the abstraction. On the opposite, initial data in this current work are not generalized to ensure privacy or even modified.

Bucketization techniques (XT06, $CVF^+10$, BPW11) are valuable due to their ability to keep the values intact. Identifiable links between items are hidden by separating the attributes of the data. Under this category lies the disassociation technique (TMLS12, LLGT14, LLGDT15, BLLL17) that works on clustering the data and hiding identifiable links in each cluster by separating the attributes. Unfortunately, disassociation has shown to be vulnerable to a privacy breach when the items are tightly coupled. It is the ability of an attacker to link his/her partial background knowledge represented by at most $m$ items that he/she is allowed to have, according to the privacy constraint $k^m$-anonymity, with certainty to less than $k$ distinct records in a disassociated dataset.

Differential privacy  (DMNS06) is an anonymization technique that adds noise to the query results. It is based on a strong mathematical foundation that guarantees that an attacker is unable to identify sensitive data about an individual if his/her information were removed from the dataset. Unlike differential privacy, this disassociation based work does not distinguish between sensitive and non-sensitive attributes and is capable to retrieve viable and trustful information that has not

been altered nor modified. The authors in (ZHZ$^+$15) defined cocktail a framework that uses both disassociation to anonymize the data for publishing as well as differential privacy to add noise on data querying. In spite of the originality of the idea, their technique subsumes the drawbacks of both disassociation and differential privacy; as it is vulnerable to the cover problem and releases questionable information.

## 3 Background

To be self-content, this section recalls the basis of disassociation (Sec. 3.1 and Sec. 3.2) introduced in (TMLS12). Next (Sec. 3.3), it exhibits a class of privacy breach called cover problem inherent to this anonymization technique already sketched in (BaBNG16) showing its repercussion on data privacy.

### 3.1 Data model

Let $\mathcal{D} = \{x_1, ..., x_d\}$ be a set of items (*e.g.*, supermarket products, query logs, or search keywords). Any subset $I \subseteq \mathcal{D}$ is an itemset (*e.g.*, items searched together). Let $\mathcal{T} = \{r_1, ..., r_n\}$ be a dataset of records where each $r_i \subseteq \mathcal{D}$ for $1 \leq i \leq n$ is a record and $r_i$ is associated with a specific individual $i$ of a population. Let $R_{\mathcal{T}}$ be a subset of records in $\mathcal{T}$. Both $\mathcal{T}$ and $R$ have the multiset semantic, which can contain more than one instance for each of its elements.

With such notations, $s(I, \mathcal{T})$ is the number of records in $\mathcal{T}$ that contain all the elements in $I$. More formally, it is defined by the following equation

$$s(I, \mathcal{T}) = |\{r \in \mathcal{T} \mid I \subseteq r\}| \tag{1}$$

By extension $s(\mathcal{T}) = s(\emptyset, \mathcal{T})$ and $s(R_T) = s(\emptyset, R_T)$ are the number of records in $\mathcal{T}$ and $R_T$ respectively.

Table 1 recalls the basic concepts and notations used in the paper.

**Table 1.** Notations used in the paper

| | |
|---|---|
| $\mathcal{D}$ | a set of items |
| $\mathcal{T}$ | a dataset containing individuals related records |
| $\mathcal{T}^*$ | a disassociated dataset i.e. a dataset anonymized using the disassociation technique |
| $\mathcal{T}^{\odot}$ | a dataset reconstructed by cross joining the itemsets of the record chunks in a cluster from the disassociated dataset |
| $r$ | a record (of $\mathcal{T}$) which is set of items associated with a specific individual of a population |
| $I$ | an itemset included in $\mathcal{D}$ |
| $s(I, \mathcal{T})$ | support of $I$ in $\mathcal{T}$ i.e. the number of records in $\mathcal{T}$ that are superset of $I$ |
| $R$ | a cluster in a disassociated dataset, formed by the horizontal partitioning of $\mathcal{T}$ |
| $R_T$ | an item chunk in a disassociated cluster |
| $R_C$ | a record chunk in a disassociated cluster |
| $\delta$ | maximum number of records allowed in a cluster |
| $n$ | number of records in $\mathcal{T}$ |

## 3.2 Disassociation

Disassociation works under the assumption that the items should neither be altered, suppressed, nor generalized, but at the same time the resulting dataset must respect the $k^m$-anonymity privacy constraint (TMK08). Formally, $k^m$-anonymity is defined as follows:

**Definition 1** ($k^m$**-anonymity**). *Given a dataset of records $\mathcal{T}$ whose items belong to a given set of items $\mathcal{D}$. The dataset $\mathcal{T}$ is $k^m$-anonymous if $\forall I \subseteq \mathcal{D}$ such that $|I| \leq m$, the number of records in $\mathcal{T}$ that are superset of $I$ is greater than or equal to $k$, i.e., $s(I, \mathcal{T}) \geq k$.*

Given a dataset $\mathcal{T}$, applying $k^m$-disassociation[3] on $\mathcal{T}$ produces a dataset $\mathcal{T}^*$ composed of $q$ clusters, each divided into a set of record chunks and an item chunk,

$$\mathcal{T}^* = \left\{ \{R_{1_{C_1}}, \dots, R_{1_{C_t}}, R_{1_Y}\}, \dots, \{R_{q_{C_1}}, \dots, R_{q_{C_v}}, R_{q_T}\} \right\}$$

such that $\forall R_{i_{C_j}} \in \mathcal{T}^*$, $R_{i_{C_j}}$ is $k$-anonymous, where,

- $R_{i_{C_j}}$ represents the itemsets of the $i^{th}$ cluster that are contained in its $j^{th}$ record chunk.
- $R_{i_T}$ is the item chunk of the $i^{th}$ cluster containing items that occur less than $k$ times.

The example in Figure 1(b) shows that the $2^2$-disassociated dataset contains only one cluster with two $2^2$-anonymous record chunks. We thus have $\mathcal{T}^* = \{R_{C_1}, R_{C_2}\}$ with $R_{C_1} = \{\{a\}, \{a,b,c,d\}, \{a,b,c,d\}, \{a,b,c,d\}, \{a,b,c,d\}, \{a,b\}\}$, and $R_{C_2} = \{\{e\}, \{e\}\}$.

According to (TMLS12) and by construction, $k^m$-disassociation guarantees all the produced record chunks are $k^m$-anonymous. This can be better explained in this example as: any combination of two items ($m = 2$) from Figure 1(a), for example $\{a,b\}$, is found at least in two records ($k = 2$) in the record chunk, thus satisfying $k^m - anonymity$ in Figure 1(b).

However, to ensure the privacy of a disassociated dataset, $k^m$-anonymity has to be guaranteed in one of the valid reconstructed datasets of $\mathcal{T}^*$ since an attacker can produce all of them, provided he/she knows $\mathcal{T}^*$, $k$, and $m$ (TMLS12). This privacy guarantee is formally expressed as follows:

**Definition 2** (**Disassociation Guarantee**). *Let $\mathcal{G}$ be the inverse transformation of $\mathcal{T}^*$ with respect to a $k^m$ disassociation, i.e., the set of all possible datasets whose $k^m$ disassociation would yield $\mathcal{T}^*$. Disassociation guarantee is established if for any $I \subseteq \mathcal{D}$ such that $|I| \leq m$, there exists $\mathcal{T}^{\odot} \in \mathcal{G}(\mathcal{T}^*)$ with $s(I, \mathcal{T}^{\odot}) \geq k$.*

The Disassociation Guarantee ensures that for any individual with a complete record $r$, and for an attacker who knows up to $m$ items of $r$, at least one of the datasets reconstructed by the inverse transformation contains the record $r$ $k$ times or more. That is, the record $r$, as all other records, exists $k$ times in at least one of the inverse transformations.

The authors in (BaBNG16) demonstrate that this disassociation guarantee is not enough to ensure privacy. They show that whenever a disassociated dataset is subject to cover problem, a privacy breach might be encountered. In the following section, we briefly present the cover problem.

---

[3] In what follows, we use $k^m$-disassociation to denote a dataset that is disassociated and satisfies $k^m$-anonymity.

### 3.3 Cover Problem

Let us recall, from Subsection 3.2 that the disassociation technique hides itemsets that occur less than $k$ times in the original dataset for a given $m$ items, by 1) dividing them into $k^m$-anonymous sub-records in record chunks and 2) ensuring that all the records reconstructed by the inverse transformation are $k^m$-anonymous in at least one of the resulting datasets.

A cover problem is defined by the ability to associate one-to-one or one-to-many items in two distinct record chunks, from the same cluster, in the disassociated data. Without loss of generality, we focus on one cluster $R = \{R_{C_1}, \ldots, R_{C_t}, R_T\}$ of $\mathcal{T}^*$ resulting from a $k^m$-disassociation. Formally, the cover problem is defined as follows.

**Definition 3** **(Cover Problem).** *Let $I_j$ be the set of items in $R_{C_j}$, $I_j = \{x \in R_{C_j}\}$. If there exists an item $z \in I_j$ such that the support of $I_j$ is equal to the support of the singleton $\{z\}$ in $R_{C_j}$, i.e.,*

$$s(I_j, R_{C_j}) = s(\{z\}, R_{C_j}), \tag{2}$$

*the cluster $R$, and the dataset $\mathcal{T}^*$ as a consequence, are subject to a* cover prob-lem.

$\forall z \in I_j$, if $z$ satisfies equation (2), $z$ is denoted as *covered item*. The set of all the covered items is denoted as $L_j$, which is contained in $I_j$. $\forall x \in I_j$, such that $x \notin L_j$, $x$ is not a covered item, then $x$ is denoted as a *covering item*. Obviously, the set of covering items is $I_j \setminus L_j$. For instance, in Figure 1(b), $I_1 = \{a, b, c, d\}$. The support of the itemset $I_1$ in $R_{C_1}$, which is $s(I_1, R_{C_1})$, is equal to 4. In turn, it is equal to the minimum support of the items in $I_1$, which, in our example, is $s(\{c\}, R_{C_1}) = 4$. Therefore, we say that the item $c$ is covered by the items $a$ and $b$. Similarly, the item $d$ is also covered by the items $a$ and $b$.

Intuitively, a privacy breach occurs if an attacker is able to link $m$ items from his/her background knowledge, to less than $k$ records in all the datasets recon-structed by the inverse transformation. More subtle is when these records contain the same set of items in all the reconstructed datasets, thus linking more than $m$ items to the individual, or worse leading to a complete de-anonymization by linking, with certainty, the complete set of items to the individual.

We will show in the following that this privacy breach might occur whenever the dataset is subject to a cover problem. Formally speaking:

**Lemma 1.** *Let $\mathcal{T}^*$ be a $k^m$-disassociated dataset subject to a cover problem. The disassociation guarantee is thus not valid for $m \geq 2$.*

*Proof.* Let $\mathcal{T}^*$ be a $k^m$-disassociated dataset subject to a cover problem. The fol-lowing set $I_j = \{x | x \in R_{C_j}\}$ is thus not empty and there exists a covered item $z \in I_j$ such that $s(I_j, R_{C_j}) = s(\{z\}, R_{C_j})$. This means that each record $r$ of $R_{C_j}$ that con-tains $z$ includes also $I_j$. Suppose now that the attacker's background knowledge is the set $\{z, y\}$ where $y$ is an item in another record chunk $R_{C_l}$.

By contradiction, suppose that the disassociation guarantee is valid, *i.e.*, $z$ and $y$ are associated together in $k$ records in at least one of the datasets $\mathcal{T}^\odot$, recon-structed by the inverse transformation of $\mathcal{T}^*$. Since $z$ is a covered item, it appears in each record $r$ defined above. The item $y$ will also be associated $k$ times with all the items in $I_j$.

While this is correct from a privacy perspective, it cannot be considered for disassociation. Items, $y$, $z$ and any covering item $x \in I_j$ are indeed considered as $k^m$-anonymous, and, therefore, should have been allot to the same record chunk $R_{C_j}$ according to disassociation[4], whereas $z$ and $y$ are respectively items of chunks $R_{C_j}$ and $R_{C_l}$ by hypothesis.

$\square$

## 4   Safe Disassociation

In this section, we show that a safe disassociation can be achieved, to ensure that a released/published dataset is no longer subject to a cover problem. This privacy guarantee is formally defined as follows:

**Definition 4  (Safe Disassociation).** *Let $\mathcal{G}$ be the inverse transformation of $\mathcal{T}^*$ with respect to a $k^m$-disassociation, whose set of items is $\mathcal{D}$. The dataset $\mathcal{T}^*$ is safely disassociated if $\forall I \subseteq \mathcal{D}$ such that $|I| \leq m$, there exists $\mathcal{T}^{\odot} \in \mathcal{G}(\mathcal{T}^*)$ with $s(I, \mathcal{T}^{\odot}) \geq k$ and $\mathcal{T}^{\odot}$ is not subject to a cover problem.*

Safe disassociation ensures that at least a dataset $\mathcal{T}^{\odot}$ reconstructed by the inverse transformation of $\mathcal{T}^*$:
 – contains $k$ records for an itemset of size $m$ or less, abiding to the $k^m$-anonymity privacy constraint, and
 – the dataset $\mathcal{T}^{\odot}$ has no covered items.
In the following, we show how this safe disassociation can be achieved by applying a partial suppression on a disassociated dataset.

### 4.1   Achieving safe disassociation with partial suppression

In previous works dedicated to privacy preservation (JPX$^+$14), partial suppression is used to ensure that no sensitive rules can be inferred with a confidence greater than a certain threshold. Here, we assume that partial suppression can achieve disassociation safely regardless the sensitivity of the data; all items are considered with the same level of sensitivity. Moreover, using partial suppression, we minimize the information loss with respect to our privacy guarantee. Unlike global suppression that removes the items in question from all the records, partial suppression remains more efficient in terms of utility.

**Partial Suppression.** *Applying the following rules until the hypothesis is established produces a safely disassociated dataset.*
 – Hypothesis: *Let I, as defined in Definition 3, be the itemset of the record chunk $R_{C_j}$ that suffers from a cover problem.*
 – Preconditions: *Let $card = \lceil \dfrac{|I|}{2} \rceil$ be the count of records that are going to be partially suppressed from $R_{C_j}$, and let $\delta$ be the maximum number of records allowed in the cluster. Partial suppression is applicable over $R_{C_j}$ when:*

$$|R_{C_j}| \leq \delta - 2 \tag{3}$$

$$s(I, R_{C_j}) \geq k + \min(card, m) \tag{4}$$

---

[4] Vertical partitioning creates $k^m$-anonymous record chunks.

- Rules:
    1. *Create a random partition $I_1 \cup I_2 .... \cup I_{card}$ of $I$ where each set $I_j$ is composed of two items, but $I_{card}$, which is a singleton if the cardinality of $I$ is odd.*
    2. *Create two empty sets $L_1$ and $L_2$, known as the ghost records.*
    3. *For each $I_i \in I$, successively:*
        (a) *suppress $I_i$ from a record $r_i \in R_C$ such that $r_i = I$ and*
        (b) *add the items $x_1$ and $x_2$ from $I_i$ to respectively $L_1$ and $L_2$:*
        $$L_1 = L_1 \cup \{x_1\} \text{ and } L_2 = L_2 \cup \{x_2\}$$
    4. *Add the two ghost records $L_1$ and $L_2$ to $R_{C_j}$*

By definition, a *cover problem* is the ability to link one-to-one or one-to-many items in two record chunks. This arises when there exists $x \in I_j$ such that $s(I_j, R_{C_j}) = s(\{x\}, R_{C_j})$. The aim of partial suppression is to suppress items in such a way as to ensure that $s(I_j, R_{C_j})$ remains different than $s(\{x\}, R_{C_j})$, thus $s(I_j, R_{C_j}) \neq s(\{x\}, R_{C_j})$.

**4.1.1   Discussion on preconditions** Preconditions (3) and (4) play an important role in ensuring that the safe disassociation is preserved throughout the process of partial suppression.

**Precondition** (3)  is defined to keep the size of the clusters bounded by the maximum cluster size $\delta$. In fact, in the horizontal partitioning, no additional records are added to the cluster if the maximum cluster size is reached. As a consequence, the record chunks created from the vertical partitioning have a cardinality less or equal to $\delta$. Hence, to add the two ghost records $L_1$ and $L_2$ to a record chunk, its cardinality must be less than $\delta - 2$.

**Precondition** (4)  is defined to preserve the $k^m$-anonymity constraint in the record chunk. Partial suppression modifies *card* occurrences of $I$, and, therefore, leaves $s(I, R_{C_j}) - card$ unmodified. Thus, we have:

$$s(I, R'_{C_j}) = s(I, R_{C_j}) - card$$
$$\geq k + \min(m, card) - card \tag{5}$$

Let $X = \{x_1, x_2, ... x_m\}$ be a subset of $I$, $X \subset I$. There are two cases to consider, depending on whether $m$ is greater than *card* or not.

- if $card \leq m$: $s(I, R'_{C_j})$ is greater than $k$ according to inequality (5). We have $s(X, R'_{C_j}) \geq s(I, R'_{C_j}) \geq k$. This means that $k^m$-anonymity is satisfied.
- if $m < card$: $s(I, R_{C_j}) - card$ records are kept unchanged by the partial suppression. After applying step 3a in partial suppression, we notice that $s(X, R'_{C_j}) \geq card - m$. This is because the smallest value, $card - m$, is obtained when each item in $X$ belongs to a distinct pair $I_j$. In this situation, $m$ records do not associate items $x_1, x_2, ..., x_m$ together. As a result, the following applies according to hypothesis (4).

$$s(X, R'_{C_j}) \geq s(I, R_{C_j}) - card + card - m$$
$$\geq s(I, R_{C_j}) - m$$
$$\geq k + \min(card, m) - m$$
$$\geq k$$

**Lemma 2.** *Given a record chunk $R_{C_j}$, we say that partial suppression achieves safe disassociation on $R_{C_j}$ if preconditions* (3) *and* (4) *are verified.*

*Proof.* Let $R'_{C_j}$ be the result of applying the partial suppression rules on the record chunk $R_{C_j}$. The itemset $I$ is randomly partitioned in *card* (where $card = \lceil \frac{|I|}{2} \rceil$) disjoint subsets of cardinality equal to 2, and possibly one singleton if $|I|$ is odd. Those subsets are used successively to perform partial suppression. Since the preconditions are verified, it is straightforward to prove that due to the rules of partial suppression, the support of any item in $R_{C_j}$ remains the same. $\forall x \in I$, $x$ is suppressed from a record $r_i = I$ and then added randomly to one of the ghost records $L_1$ or $L_2$, therefore:

$$s(\{x\}, R'_{C_j}) = s(\{x\}, R_{C_j}). \tag{6}$$

However, what has been partially suppressed is the association between any two items $x$ and $y$ in $I$. Two case scenarios can arise:

- items $x$ and $y$ belong to the same subset $I_i$ of $I$. Thus, $x$ and $y$ are suppressed once from the same record $r_i$, and added to $L_1$ and $L_2$ respectively. This ensures that they cannot be associated again in the ghost records:

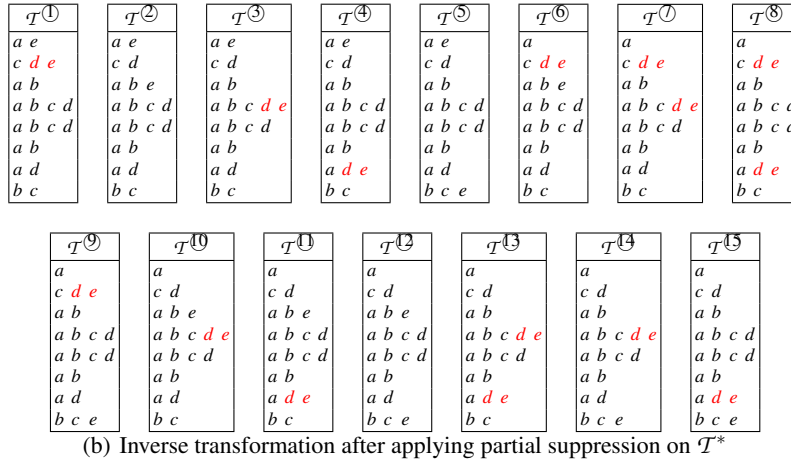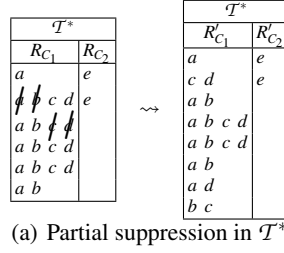$$s(\{x,y\}, R'_{C_j}) = s(\{x,y\}, R_{C_j}) - 1 = s(\{x\}, R'_{C_j}) - 1 = s(\{y\}, R'_{C_j}) - 1$$

- items $x$ and $y$ belong to two different subsets of $I$, $I_i$ and $I_j$ respectively. The association between $x$ and $y$ is lost twice since these items were suppressed from two distinct records in $R_{C_j}$, then: $s(\{x,y\}, R'_{C_j})$ is equal to $s(\{x,y\}, R_C) - 2$ (resp. to $s(\{x,y\}, R_C) - 1$) if $x$ and $y$ are added to the different ghost records $L_1$ and $L_2$ (resp. if $x$ and $y$ are added to the same ghost records).

In both cases, we have:

$$s(\{x,y\}, R'_{C_j}) < s(\{x\}, R'_{C_j})$$
$$< s(\{y\}, R'_{C_j})$$

Due to the inequality $S(I, R'_{C_j}) \leq s(\{x,y\}, R'_{C_j})$ and $\forall x, y$ in $I$ we have $s(I, R'_C) < s(\{x\}, R'_C)$, which concludes the proof.

$\square$

To illustrate how a cover problem can be eliminated through partial suppression, let us consider the example in Figure 2(a). Both items $c$ and $d$ are covered items in the disassociated dataset $\mathcal{T}^*$. After applying partial suppression, in Figure 2(a), two subsets $I_1 = \{a, b\}$ and $I_2 = \{c, d\}$ are created after randomly partitioning $I = \{a, b, c, d\}$. From $I_1$ and $I_2$, two ghost records are created $L_1 = \{a, d\}$ and $L_2 = \{b, c\}$, containing one of the items from each suppressed subsets. Next, these two ghost records are added to $R_{C_1}$. We illustrate, in Figure 2(b), all possible reconstructed datasets of the final disassociated dataset. Now, if an attacker knows that a specific individual has searched for items $\{d, e\}$ (considered as the attacker's background knowledge), he/she will be able to associate them with three possible records $\{c, d, e\}$, $\{a, d, e\}$ and $\{a, b, c, d, e\}$. While these extra associations are considered noise, for the sake of privacy, they added ambiguity to the result since the attacker cannot link $\{d, e\}$ to a particular record. In the next section, we will study and evaluate the impact of partial suppression on the utility of the dataset.

(a) Partial suppression in $\mathcal{T}^*$



(b) Inverse transformation after applying partial suppression on $\mathcal{T}^*$

**Fig. 2.** Eliminating a cover problem with partial suppression

## 5 Experimental Evaluation

In keeping with the previous work (BaBNG16), we elaborate our experiments on two datasets, the *BMS*1 and the *BMS*2, which contain click-stream E-commerce data. Table 2 shows the properties of the datasets.

The aim of the experiments can be summarized as follows:

- Evaluating the privacy breach in the disassociated dataset.
- Evaluating the utility loss w.r.t the suppression of items (or their occurrences).
- Studing the loss in associations when disassociating and safely disassociating a dataset.
- Evaluating the performance of partial suppression.

### 5.1 Privacy and utility metrics

**5.1.1 Privacy Evaluation Metric (PEM)** represents the number of vulnerable record chunks in a disassociated dataset. In fact, we consider that every

**Table 2.** Datasets properties

| Dataset | # of distinct individuals | # of distinct items | count of items' occurrences |
|---------|---------------------------|---------------------|----------------------------|
| $BMS1$ | 59602 | 497 | 149639 |
| $BMS2$ | 77512 | 3340 | 358278 |

record chunk that is subject to a cover problem is a vulnerable record chunk. We formally define our *PEM* as:

$$PEM = \frac{vRC}{RC}$$

where,

- *vRC* represents the number of vulnerable record chunks in $\mathcal{T}$, the disassociated dataset, and
- *RC* represents the total number of record chunks in $\mathcal{T}$.

**5.1.2 Relative Loss Metric (RLM)** determines the relative number of suppressed occurrences of items with partial suppression. Formally,

$$RLM = \frac{\sum_{\forall x \in \mathcal{D}}(s(\{x\}, \mathcal{T}^*) - s(\{x\}, \mathcal{T}'))}{\sum_{\forall x \in \mathcal{D}}(s(\{x\}, \mathcal{T}^*))}$$

where,

- $s(\{x\}, \mathcal{T}^*)$ represents the support of the item $x$ in the disassociated dataset $\mathcal{T}^*$, and
- $s(\{x\}, \mathcal{T}')$ represents the support of the item $x$ in the safely disassociated dataset $\mathcal{T}'$.

**5.1.3 Relative Association Error (RAE)** evaluates how likely two items remain associated together in an anonymized dataset (TMLS12). In fact, using *RAE*, we are able to evaluate the information loss due to the anonymization of the dataset (whether it is disassociated or safely disassociated). Formally, *RAE* is defined as follows:

$$RAE = \frac{s(\{x,y\}, \mathcal{T}) - s(\{x,y\}, [\mathcal{T}^*|\mathcal{T}'])}{AVG(s(\{x,y\}, \mathcal{T}), s(\{x,y\}, [\mathcal{T}^*|\mathcal{T}']))}$$

where,

- $s(\{x,y\}, \mathcal{T})$ represents the support of items $\{x,y\}$ in the original dataset $\mathcal{T}$, the disassociated dataset $\mathcal{T}^*$, or the safely disassociated dataset $\mathcal{T}'$.
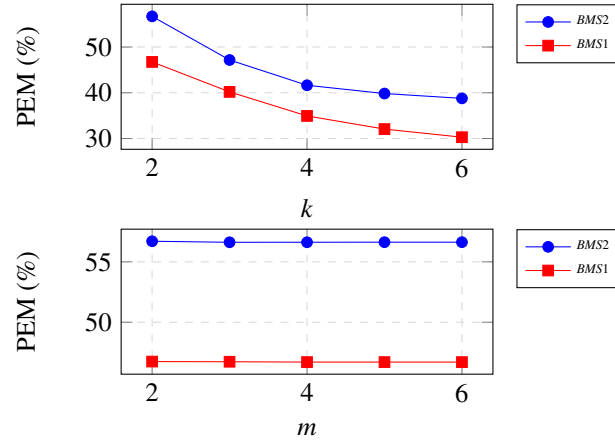
## 5.2 Experimental results

In this section, we present the results of the conducted experiments, evaluating privacy and data utility over the safely disassociated dataset.

**5.2.1   Evaluating the privacy breach**   We consider that a privacy breach occurs if an attacker is able to link $m$ items, which he/she already knows about an individual, to less than $k$ records in all the datasets reconstructed by the inverse transformation. In this test, we study the impact of the cover problem on the privacy of the disassociated dataset. In fact, we consider that a potential privacy breach exists whenever a cover problem is identified in a record chunk regardless of the background knowledge of the attacker. It is typically a strong attacker (BaBNG16) who is able to link any two items to a specific individual. Hence, we determine the following:

– the relationship between the *PEM* and both, $k$ and $m$.
– the relationship between the *PEM* and the maximum cluster size $\delta$.

**Varying $k$ and $m$:**   we vary $k$ and $m$ from 2 to 6. For each value, we compute the *PEM* to evaluate how the privacy constraint remains satisfied in the record chunks. Figure 3 shows the results of the evaluation. When $k$ increases in both datasets *BMS*1 and *BMS*2, the *PEM* decreases from 46% to 30% in *BMS*1 and from 56% to 39% in *BMS*2. While varying $k$ affects the *PEM*, varying $m$ has no noticeable impact on the *PEM* in both datasets. This is not surprising since, due to the cover problem, any two items known to the attacker can lead to a privacy breach.



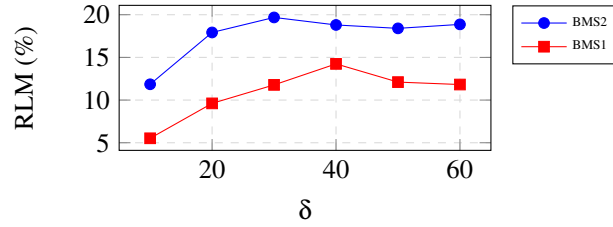**Fig. 3.** Evaluating the *PEM* while varying $k$ and $m$

**Varying $\delta$:**   we vary $\delta$ from 10 to 60. For each value, we compute the *PEM* to evaluate how the privacy constraint remains satisfied in the record chunks. The results in Figure 4 show that the *PEM* increases from 17% to 52% in *BMS*1 and from 26% to 57% in *BMS*2. This means that with more records in the cluster, the higher the chances are to compromise the dataset due to the cover problem.

**Fig. 4.** Evaluating the *PEM* while varying δ

Now, given that the privacy breach is directly related to the value of δ, we choose to vary δ and fix the values of $k$ and $m$ in the remaining tests. We use $k = 3$ to keep computational time to a minimum and $m = 2$ since only two items are sufficient to raise a privacy breach.
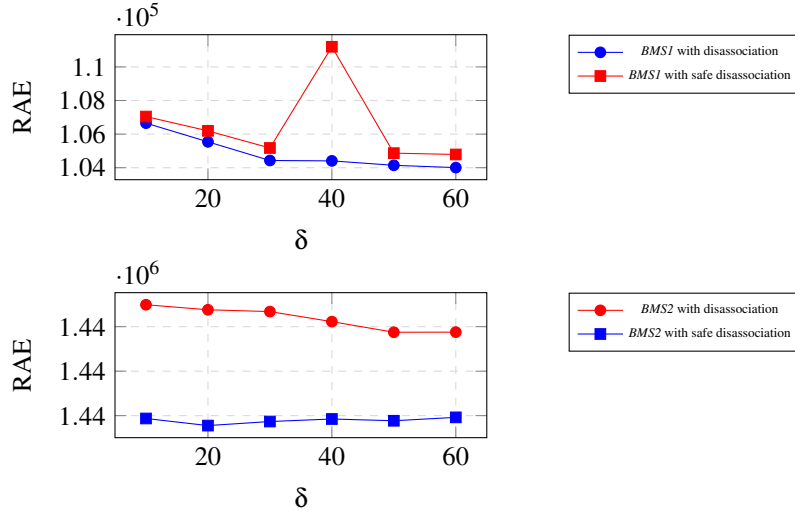
**5.2.2 Evaluating the number of suppressed items** In this test, we evaluate the number of items to be suppressed to safely disassociate a dataset. We note that the record chunks that do not comply with preconditions (3) and (4) in partial suppression, for the sake of privacy, their items are completely suppressed. We use the *RLM* to estimate this loss of items. We vary the maximum cluster size δ from 10 to 60 and compute, for each value, the *RLM*. Figure 5 shows the results of our evaluation. It is not surprising that *BMS*2 is more susceptible to partial suppression since it is more vulnerable to the cover problem due to its size, as noted in the previous test (Figure 4). In addition, in both datasets, when the maximum cluster size δ increases, we notice an increase in the *RLM*. This is actually consistent with the fact that the more vulnerable the record chunks are, the more items will be suppressed. Overall, the 20% suppressions of items represent an acceptable trade-off between privacy and utility, especially that, with safe disassociation, we release real datasets without modifying the items.



**Fig. 5.** Evaluating the *RLM* while varying δ

**5.2.3 Evaluating the association error** In this test, we evaluate the error in associations, which is the result of dividing the records into chunks after the

partitioning process. Eventually, some of the itemsets will be separated and their items will be stored into different chunks. This adds noise to the dataset since the support of associations between these separated items might be higher in the reconstructed datasets. This leads to this error in associations that can be calculated using the *RAE*, which is the relative difference between the support of the association of an itemset in the original dataset and the anonymized dataset (XT06). Again, we vary the maximum cluster size δ from 10 to 60 and compute, for each value, the *RAE* on disassociated and safely disassociated datasets. The results in Figure 6 shows that the *RAE* is higher in a safely disassociated dataset. It is not surprising because, with safe disassociation, we suppress some items to prevent the privacy breach. However, the difference between *RAE* in safe disassociation and disassociation remains acceptable; varying between 1% for *BMS*1 and 0.1% for *BMS*2.
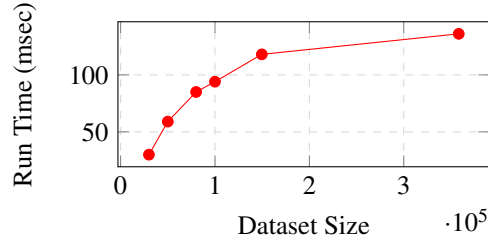


**Fig. 6.** Evaluating the *RAE* while varying δ

**5.2.4 Performance evaluation** We compare the performance of our algorithm with different dataset size. The datasets used in this experiment are formed from *BMS*1 and *BMS*2. We vary, in the $x - axis$, the number of records with {358k, 149K, 100K, 80K, 50k, 30k 10k}; and the maximum cluster size δ is fixed to 30. Figure 7 shows an increase in the run time with the increase of the dataset size.

# 6 Conclusion

Disassociation is an interesting anonymization technique that is able to hide the link between individuals and their complete set of items while keeping the items

**Fig. 7.** Performance evaluation

without generalization. In a previous work (BaBNG16), disassociation was considered vulnerable due to a cover problem. Basically, it is the ability to associate one-to-one or one-to-many items in two subsequent record chunks of the disassociated data. In this paper, we propose safe disassociation to solve this problem. We use partial suppression to achieve safe disassociation by suppressing some of the items that lead to a cover problem from subsequent record chunks. In the experiments, we show that the vulnerability of a disassociated dataset depends on the size of the cluster. We evaluate the utility of a safely disassociated dataset in terms of 1) the number of items to be suppressed to achieve safe disassociation, and 2) the additional noisy associations added due to item suppression and partitioning. The results of our evaluations showed that an acceptable trade-off between privacy and utility is met.

In future works, we aim at maximizing the utility of a safely disassociated dataset by modifying the clustering algorithm to keep user-defined itemsets associated together.

## Acknowledgments

---

[5] www.inmobiles.net

# Bibliography

Sara Barakat, Bechara al Bouna, Mohamed Nassar, and Christophe Guyeux. On the evaluation of the privacy breach in disassociated set-valued datasets. In Christian Callegari, Marten van Sinderen, Panagiotis G. Sarigiannidis, Pierangela Samarati, Enrique Cabello, Pascal Lorenz, and Mohammad S. Obaidat, editors, *Proceedings of the 13th International Joint Conference on e-Business and Telecommunications (ICETE 2016) - Volume 4: SECRYPT, Lisbon, Portugal, July 26-28, 2016.*, pages 318–326. SciTePress, 2016.

Michael Bewong, Jixue Liu, Lin Liu, and Jiuyong Li. Utility aware clustering for publishing transactional data. In Jinho Kim, Kyuseok Shim, Longbing Cao, Jae-Gil Lee, Xuemin Lin, and Yang-Sae Moon, editors, *Advances in Knowledge Discovery and Data Mining*, pages 481–494, Cham, 2017. Springer International Publishing.

Joachim Biskup, Marcel PreuB, and Lena Wiese. On the inference-proofness of database fragmentation satisfying confidentiality constraints. In *Proceedings of the 14th Information Security Conference*, Xian, China, oct 26-29 2011.

Michael Barbaro and Tom Zeller. A face is exposed for aol searcher no. 4417749, 2006.

Chen, Liuhua, Shenghai Zhong, Li-E. Wang, and Xianxian Li. A sensitivity-adaptive ρ-uncertainty model for set-valued data. In *International Conference on Financial Cryptography and Data Security, Berlin, Heidelberg, 2016.*, pages 460–473. Springer, 2016.

Valentina Ciriani, Sabrina De Capitani Di Vimercati, Sara Foresti, Sushil Jajodia, Stefano Paraboschi, and Pierangela Samarati. Combining fragmentation and encryption to protect privacy in data storage. *ACM Trans. Inf. Syst. Secur.*, 13:22:1–22:33, July 2010.

Sabrina De Capitani di Vimercati, Sara Foresti, Sushil Jajodia, Giovanni Livraga, Stefano Paraboschi, and Pierangela Samarati. Extending loose associations to multiple fragments. In *Proceedings of the 27th International Conference on Data and Applications Security and Privacy XXVII*, DBSec'13, pages 1–16, Berlin, Heidelberg, 2013. Springer-Verlag.

Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Proceedings of the Third Conference on Theory of Cryptography*, TCC'06, pages 265–284, Berlin, Heidelberg, 2006. Springer-Verlag.

Amin Milani Fard and Ke Wang. An effective clustering approach to web query log anonymization. In *Security and Cryptography (SECRYPT), Proceedings of the 2010 International Conference on*, pages 1–11. IEEE, 2010.

Yeye He and Jeffrey F. Naughton. Anonymization of set-valued data via top-down, local generalization. *Proc. VLDB Endow.*, 2(1):934–945, August 2009.

Xiao Jia, Chao Pan, Xinhui Xu, KennyQ. Zhu, and Eric Lo. ρ-uncertainty anonymization by partial suppression. In SouravS. Bhowmick, CurtisE. Dyreson, ChristianS. Jensen, MongLi Lee, Agus Muliantara, and Bernhard Thalheim, editors, *Database Systems for Advanced Applications*, volume

8422 of *Lecture Notes in Computer Science*, pages 188–202. Springer International Publishing, 2014.

Grigorios Loukides, John Liagouris, Aris Gkoulalas-Divanis, and Manolis Terrovitis. Utility-constrained electronic health record data publishing through generalization and disassociation. In Aris Gkoulalas-Divanis and Grigorios Loukides, editors, *Medical Data Privacy Handbook*, pages 149–177. Springer International Publishing, 2015.

Grigorios Loukides, John Liagouris, Aris Gkoulalas-Divanis, and Manolis Terrovitis. Disassociation for electronic health record privacy. *Journal of Biomedical Informatics*, 50:46–61, 2014.

Tiancheng Li, Ninghui Li, Jian Zhang, and Ian Molloy. Slicing: A new approach for privacy preserving data publishing. *IEEE Trans. Knowl. Data Eng.*, 24(3):561–574, 2012.

Ashwin Machanavajjhala, Johannes Gehrke, Daniel Kifer, and Muthuramakrishnan Venkitasubramaniam. *l*-diversity: Privacy beyond *k*-anonymity. In *Proceedings of the 22nd IEEE International Conference on Data Engineering (ICDE 2006)*, Atlanta Georgia, April 2006.

Pierangela Samarati. Protecting respondents' identities in microdata release. *IEEE Trans. Knowl. Data Eng.*, 13(6):1010–1027, 2001.

Latanya Sweeney. k-anonymity: a model for protecting privacy. *International Journal on Uncertainty, Fuzziness and Knowledge-based Systems*, 10(5):557–570, 2002.

Manolis Terrovitis, Nikos Mamoulis, and Panos Kalnis. Privacy-preserving anonymization of set-valued data. *PVLDB*, 1(1):115–125, 2008.

Manolis Terrovitis, Nikos Mamoulis, John Liagouris, and Spiros Skiadopoulos. Privacy preservation by disassociation. *Proc. VLDB Endow.*, 5(10):944–955, June 2012.

J. Wang, C. Deng, and X. Li. Two privacy-preserving approaches for publishing transactional data streams. *IEEE Access*, pages 1–1, 2018.

Ke Wang, Peng Wang, Ada Waichee Fu, and Raymond Chi-Wing Wong. Generalized bucketization scheme for flexible privacy settings. *Information Sciences*, 348:377 – 393, 2016.

Xiaokui Xiao and Yufei Tao. Anatomy: Simple and effective privacy preservation. In *Proceedings of 32nd International Conference on Very Large Data Bases (VLDB 2006)*, Seoul, Korea, September 12-15 2006.

Zhang, Hongli, Zhigang Zhou, Lin Ye, and D. U. Xiaojiang. Towards privacy preserving publishing of set-valued data on hybrid cloud. In *IEEE Transactions on Cloud Computing*, 2015.