

Minimality Attack in Privacy Preserving Data Publishing

Technical Report TR 2006-28

Raymond Chi-Wing Wong¹ Ada Wai-Chee Fu¹ Ke Wang² Jian Pei²

¹ The Chinese University of Hong Kong, {cwwong, adafu}@cse.cuhk.edu.hk

² Simon Fraser University, Canada, {jpei, wangk}@cs.sfu.ca

Abstract

Data publishing generates much concern over the protection of individual privacy. In the well-known k -anonymity model and the related models such as l -diversity and (α, k) -anonymity, the adversary is assumed to possess knowledge about an external table with information of the quasi-identifiers of individuals. In this paper, we show that knowledge of the mechanism or algorithm of anonymization for data publication can also lead to extra information that assists the adversary and jeopardizes individual privacy. In particular, all known mechanisms try to minimize information loss and such an attempt provides for a loophole for attacks. We call such an attack a minimality attack. In this paper, we propose a model called m -confidentiality which deals with the individual privacy issue with the consideration of minimality attacks. Though the problem of optimal m -confidentiality anonymization is NP-hard, we propose an algorithm which generates m -confidential data sets efficiently. We also conducted experiments to show how such an attack can succeed on a real dataset and that our algorithm suffers almost no penalty in information loss.

1 Introduction

Although data mining is potentially useful, many data holders are reluctant to provide their data for data mining for fear of violating individual privacy. In recent years, study has been made to ensure that the sensitive information of individuals cannot be identified easily [19, 20, 12, 16, 11]. One well-studied approach is the k -anonymity model which in turn led to other models such as confidence bounding [22], l -diversity [14], (α, k) -anonymity [24], t -closeness [13], (k, e) -anonymity [27] and (c, k) -safety [15]. These models assume that the data or table T contains (1) a *quasi-identifier (QID)*, which is a set of attributes (e.g. a QID may be {Date of birth, Zipcode, Sex}) in T which can be used to identify an individual, and (2) *sensitive attributes*, attributes in T which may contain some sensitive values (e.g. HIV of attribute Disease) of individuals. Often, it is also assumed that each tuple in T corresponds to an individual and no two tuples refer to the same individual. All tuples with the same QID value form an *equivalence class*, which we call QID-EC. The table T is said to satisfy k -anonymity if the size of every equivalence class is greater than or equal to k .

In this paper, we study the case where the adversary has some additional knowledge about the mechanism involved in the anonymization and launches an attack based on this knowledge. We focus on the protection of the relationship between the quasi-identifier and a single sensitive attribute. In a simplified setting of l -diversity model [14], a QID-EC is said to be *l -diverse* or satisfy *l -diversity* if the proportion of each sensitive value is at most $1/l$. A table satisfies l -diversity (or it is l -diverse) if all QID-EC's in it are l -diverse. In the following discussion, when we refer to l -diversity, we refer to this simplified setting.¹ The complex l -diversity model is discussed in Section 5, in which we show that our results can be extended to other anonymization models.

¹This simplified model is a special case of the confidence bounding in [22] and is the same as (α, k) -anonymity [24] when $k = 1$ and $\alpha = 1/l$.

QID	Disease	QID	Disease	QID	QID
$q1$	HIV	$q1$	HIV	Q	Q
$q1$	non-sensitive	$q1$	HIV	Q	Q
$q2$	HIV	$q2$	non-sensitive	Q	Q
$q2$	non-sensitive	$q2$	non-sensitive	Q	Q
$q2$	non-sensitive	$q2$	non-sensitive	Q	$q2$
$q2$	non-sensitive	$q2$	non-sensitive	Q	$q2$
$q2$	non-sensitive	$q2$	non-sensitive	Q	$q2$

(a) good table (b) bad table (c) global (d) local

Table 1: 2-diversity: global and local recoding

1.1 Minimality Attack

In Table 1(a), assume that the QID values of $q1$ and $q2$ can be generalized to Q and assume only one sensitive attribute “disease”, in which HIV is a sensitive value. For example, $q1$ may be $\{Nov\ 1930, Z3972, M\}$, $q2$ may be $\{Dec\ 1930, Z3972, M\}$ and Q is $\{Nov/Dec\ 1930, Z3972, M\}$. (Note that $q1$ and $q2$ may also be generalized values.) A tuple associated with HIV is said to be a *sensitive* tuple. For each equivalence class, at most half of the tuples are sensitive. Hence, the table satisfies 2-diversity.

As observed in [12], existing approaches of anonymization for data publishing have an implicit principle: “For any anonymization mechanism, it is desirable to define some notion of minimality. Intuitively, a k -anonymization should not generalize, suppress, or distort the data more than it is necessary to achieve k -anonymity.” Based on this minimality principle, Table 1(a) will not be generalized.² In fact the above notion of minimality is too strong since almost all known anonymization problems for data publishing are NP-hard, many existing algorithms are heuristical and only attain local minima. We shall later give a more relaxed notion of the minimality principle in order to cover both the optimal as well as the heuristical algorithms. For now, we assume that minimality principle means that a QID-EC will not be generalized unnecessarily.

Next, consider a slightly different table, Table 1(b). Here, the set of tuples for $q1$ violates 2-diversity because the proportion of the sensitive tuples is greater than 1/2. Thus, this table will be anonymized to a *generalized* table by generalizing the QID values as shown in Table 1(c) by *global recoding* [26, 21]. The tuples in this table contains the generalized values of the QID arranged in the same tuple ordering as the corresponding tuples in Table 1(b). This is a convention we shall use for all examples in this paper showing the anonymization of a table. In global recoding, all occurrences of an attribute value are recoded to the same value. If *local recoding* [19, 1] is adopted, occurrences of the same value of an attribute may be recoded to different values. Such an anonymization is shown in Table 1(d). These anonymized tables satisfy 2-diversity. However, do these tables really protect individual privacy?

In most previous works [20, 11, 12, 26], the knowledge of the adversary involves an external table T^e such as a voter registration list that maps QIDs to individuals.³ As in many previous works, we assume that each tuple in T^e maps to one individual and no two tuples map to the same individual. The same is also assumed in the table T to be published. Let us first consider the case when T and T^e are mapped to the same set of individuals. Table 2(a) is an example of T^e .

Assume further that the adversary knows the goal of 2-diversity, s/he also knows whether it is a global or local recoding, and Table 2(a) is available as the external table T^e . With the notion of minimality in anonymization, the adversary reasons as follows: From the published Table 1(c), there are 2 sensitive tuples in total. From T^e , there are 2 tuples with $QID=q1$ and 5 tuples with $QID=q2$. Hence, the equivalence

²This is the case for each of the anonymization algorithms in [14, 22, 24].

³There are many sources of such an external table T^e . *Most municipalities sell population registers that include the identifiers of individuals along with basic demographics; examples include local census data, voter lists, city directories, and information from motor vehicle agencies, tax assessors, and real estate agencies* [17]. In the voter list, 97% of the voters were identifiable with just the full postal code and birth date [18]. From [20], it is reported that a city’s voter list in two diskettes was purchased for twenty dollars, and was used to re-identify medical records.

Name	QID	QID	Name	QID	QID
Andre	$q1$	$q1$	Andre	$q1$	$q1$
Kim	$q1$	$q1$	Kim	$q1$	$q1$
Jeremy	$q2$	$q2$	Jeremy	$q2$	$q2$
Victoria	$q2$	$q2$	Victoria	$q2$	$q2$
Ellen	$q2$	$q2$	Ellen	$q2$	$q2$
Sally	$q2$	$q2$	Sally	$q2$	$q2$
Ben	$q2$	$q2$	Ben	$q2$	$q2$
			Tim	$q4$	$q4$
			Joseph	$q4$	$q4$

(a) individual QID (b) multiset (c) individual QID (d) multiset

Table 2: T^e : external table available to the adversary

QID	Disease	QID	Disease	QID	QID
$q1$	HIV	$q1$	HIV	Q	Q
$q1$	Lung Cancer	$q1$	HIV	Q	Q
$q2$	Gallstones	$q2$	Gallstones	Q	Q
$q2$	HIV	$q2$	Lung Cancer	Q	Q
$q2$	Ulcer	$q2$	Ulcer	Q	$q2$
$q2$	Alzheimer	$q2$	Alzheimer	Q	$q2$
$q2$	Diabetes	$q2$	Diabetes	Q	$q2$
$q4$	Ulcer	$q4$	Ulcer	$q4$	$q4$
$q4$	Alzheimer	$q4$	Alzheimer	$q4$	$q4$

(a) good table (b) bad table (c) global (d) local

Table 3: 2-diversity (where all values in Disease are sensitive): global and local recoding

class for $q2$ in the original table *must* already satisfy 2-diversity, because even if both sensitive tuples have QID= $q2$, the proportion of sensitive values in the class for $q2$ is only $2/5$. Since *generalization* has taken place, at least one equivalence class in the original table T must have violated 2-diversity, because otherwise no generalization will take place according to minimality. The adversary concludes that $q1$ has violated 2-diversity, and that is possible only if both tuples with QID= $q1$ have a disease value of “HIV”. The adversary therefore discovers that Andre and Kim are linked to “HIV”.

In some previous works, it is assumed that the set of individuals in the external table T^e can be a superset of that for the published table. Table 2(c) shows such a case, where there is no tuple for Tim and Joseph in Table 1(a) and Table 1(b). If it is known that $q4$ cannot be generalized to Q (e.g. $q4 = \{Nov\ 1930, Z3972, F\}$ and $Q = \{Jan/Feb\ 1990, Z3972, M\}$), then the adversary can be certain that the tuples with QID= $q4$ are not in the original table. Thus, the extra $q4$ tuples in T^e do not have any effect on the above reasoning of the adversary and, therefore, the same conclusion can be drawn. We call such an attack based on the minimality principle a *minimality attack*.

Observation 1 *If a table T is anonymized to T^* which satisfies l -diversity, it can suffer from a minimality attack. This is true for both global and local recoding and for the cases when the set of individuals related to T^e is a superset of that related to T .*

In the above example, some values in the sensitive attribute Disease are not sensitive. Would it help if all values in the sensitive attributes are sensitive? In the tables in Table 3, we assume that all values for Disease are sensitive. Table 3(a) satisfies 2-diversity but Table 3(b) does not. Suppose anonymization of Table 3(b) results in Table 3(c) by global recoding and Table 3(d) by local recoding. The adversary is armed with the external table Table 2(c) and the knowledge of the goal of 2-diversity, s/he can launch

an attack by reasoning as follows: with 5 tuples for $QID=q2$ and each sensitive value appearing at most twice, there cannot be any violation of 2-diversity for the tuples with $QID=q2$. There must have been a violation for $QID=q1$. For a violation to take place, both tuples with $QID=q1$ must be linked to the same disease. Since HIV is the only disease that appears twice, Andre and Kim must have contracted HIV.

Observation 2 *Minimality attack is possible whether the sensitive attribute contains non-sensitive values or not.*

The intended *objective* of 2-diversity is to make sure that an adversary cannot deduce with a probability above $1/2$ that an individual is linked to any sensitive value. Thus, the published tables violate this objective.

Observation 3 *The above attacks to Andre would also be successful if the knowledge of the external table Table 2(a) is replaced by that of a multiset of the QID values as shown in Table 2(b) plus the QID value of Andre; or if Table 2(c) is replaced by the multiset in Table 2(d) plus the QID value of Andre.*

Note that the multisets in Tables 2(b) and (d) are inherently available in the published data if the bucketization technique as in [25, 27, 15] is used.

1.2 Contributions

In this paper, we introduce the problem of minimality attacks in privacy preservation for data publishing. Our contributions include the following.

1. To the best of our knowledge, this is the first work to study the attack by minimality in privacy preserving data publishing. We propose an m -confidentiality model to capture the privacy preserving requirement under the additional adversary knowledge of the minimality of the anonymization mechanisms.
2. Since almost all known anonymization methods for data publishing attempt to minimize information loss, we show in Section 5 how minimality attack can be successful in a variety of known anonymization models.
3. We propose a solution to generate a published data set which satisfies m -confidentiality. Our method makes use of the existing mechanisms for k -anonymity with additional precaution steps. Interestingly, although it has been discovered by recent research works that k -anonymity is incapable of handling sensitive values in some cases, it is precisely this feature that makes it a useful component in our method to counter attacks by minimality for protecting sensitive data. Since k -anonymization does not consider the sensitive values, its result is not related to whether some tuples need to be anonymized due to the sensitive values. Without this relationship, an attack by minimality becomes infeasible.
4. We have conducted a comprehensive empirical study on both the problem and our method. We show how such a minimality attack can succeed on a real data set in our experiment. Compared to the most competent existing algorithms for k -anonymity, our method introduces very minor computation overhead. The information loss generated by our method is also comparable to those resulting in known algorithms for k -anonymity.

2 Related Work

Since the introduction of k -anonymity, there have been a number of enhanced models such as confidence bounding [22], l -diversity [14], (α, k) -anonymity [24], t -closeness [13], (k, e) -anonymity [27] and personalized privacy [26], which additionally consider the privacy issue of disclosure of the *relationship* between the quasi-identifier and the sensitive attributes. Confidence bounding is to bound the confidence by which a QID can be associated with a sensitive value. T is said to satisfy (α, k) -anonymity if T is k -anonymous and the proportion of each sensitive value in every equivalence class is at most α , where $\alpha \in [0, 1]$ is a user parameter. If we set $\alpha = \frac{1}{k}$ and $k = 1$, then the (α, k) -anonymity model becomes the simplified model of l -diversity.

An adversary may also have some additional knowledge about the individuals in the dataset or some knowledge about the data involved [14, 9, 15]. [14] considers the possibility that the adversary can exclude some sensitive values. For example, Japanese have an extremely low incidence of heart disease. Thus, the adversary can exclude heart disease in a QID-EC for a Japanese individual. [9] considers that additional information may be available in terms of some statistics on some of the attributes, such as age statistics and zip code statistics. More recently, [15] tries to protect sensitive data against background knowledge in the form of implications, e.g., if an individual A has HIV then another individual B also has HIV, and proposes a model called (c, k) -*safety* to protect against such attacks. However, none of the above work considers the knowledge of the anonymization mechanism discussed in this paper. In Section 5 we shall show that the above previous works are vulnerable to minimality attacks. Other than generalization, more general distortion can be applied to data before publishing. The use of distortion has been proposed in earlier works such as [2, 5].

The idea of attack by minimality has been known for some time in cryptographic attack where the adversary makes use of the knowledge of the underlying cryptographic algorithm. In particular, a timing attack [10] in a public-key encryption system, such as RSA, DSS and SSL, is a practical and powerful attack that exploits the timing factor of the implemented algorithm, with the assumption that the algorithm will not take more time than necessary. Measuring response time for a specific query might give away relatively large amounts of information. To defend timing attack, the same algorithm can be implemented in such a way that every execution returns in exactly x seconds, where x is the maximum time it ever takes to execute the routine. In this extreme case, timing does not give an attacker any helpful information. In 2003, Boneh and Brumley [4] demonstrated a practical network-based timing attack on SSL-enabled web servers which recovered a server private key in a matter of hours. This led to the widespread deployment and use of blinding techniques in SSL implementations.

3 Problem Definition

Let T be a table. A record of T is a tuple. An attribute defines all the possible values in a column. A *quasi-identifier* (QID) is a set of attributes of T that may serve as identifications for some individuals. We assume that one of the attributes is a sensitive attribute where some values in this attribute should not be linkable to any individual.

Assumption 1 *Each tuple in the table T is related to one individual and no two tuples are related to the same individual.*

We assume that each attribute has a corresponding conceptual *taxonomy* \mathcal{T} . A lower level domain in the taxonomy \mathcal{T} provides more details than a higher level domain. For example, Figure 1 shows a generalization taxonomy of “Education” in the “Adult” dataset [3]. Values “undergrad” and “postgrad” can be generalized to “university”.⁴ Generalization replaces lower level domain values in the taxonomy with higher level domain values.

Some previous works consider taxonomies only for QID attributes while some consider also taxonomies for the sensitive attributes. In earlier works on anonymization, the taxonomy for an attribute in the QID or the sensitive attribute is a tree. However, in general, a taxonomy may be a directed acyclic graph (DAG). For example, “lung cancer” may have two parents “cancer” and “respiratory disease”. Also, “day” can be generalized to “week”, or via “month” to “year”, or via “season” to “year”. Therefore, we extend the meaning of a taxonomy to any *partially order set* with a partial order \prec where $v' \prec v$ if v is an ancestor of v' . An attribute may have more than one taxonomy, where a certain value can belong to two or more taxonomies.⁵

Let \mathcal{T} be a taxonomy for an attribute in QID. We call the leaf nodes of the taxonomy \mathcal{T} the *ground values*.

In Figure 1, values “1st-4th”, “undergrad” and “vocational” are some ground values in \mathcal{T} . As “university” is an ancestor of “undergrad”, we obtain “undergrad” \prec “university”.

⁴Such hierarchies can also be created for numerical attributes by generalizing values to value range and to wider value ranges. The ranges can be determined by users or a machine learning algorithm [6].

⁵Note that a taxonomy may not be a lattice. For example, consider attribute disease. “Nasal cancer” and “lung cancer” may both be under two parents of “cancer” and “respiratory disease”.

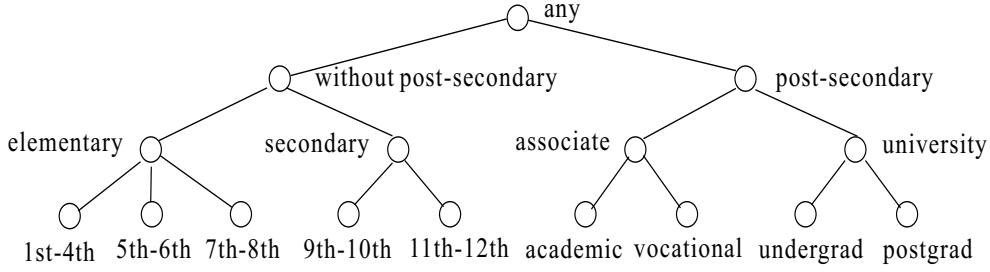


Figure 1: Generalization taxonomy of “Education” in the “Adult” dataset

When a record contains the sensitive value of “lung cancer”, it can be generalized to either “respiratory disease” or “cancer”. It can be reasoned that, while “cancer” and “lung cancer” are possibly considered sensitive, “respiratory disease” may not be. Therefore, we can assume the following property.

Assumption 2 (Taxonomy property): *In a taxonomy for a sensitive attribute, the ancestor nodes of a non-sensitive node are also non-sensitive. The ancestor of a sensitive node may be either sensitive or non-sensitive.*

In a faithful anonymization, a value can be generalized to any ancestor. For example, “lung cancer” may be generalized to “cancer” or “respiratory disease”. With the above assumption, if a node is sensitive, all ground values in its descendants are sensitive.

With a taxonomy for the sensitive attribute, such as the one in Figure 1, in general, the protection is not targeting on a single ground value. In Figure 1, all the values under “elementary” may be sensitive in the sense that there should not be linkage between an individual and the set of values {1st-4th, 5th-6th, 7th-8th}. That is, the adversary must not be able to deduce with confidence that an individual has education between 1st to 8th grade. In general, a group of sensitive values may not be under one subtree. For example, for diseases, it is possible that cancer and HIV are both considered sensitive. So, a user should not be linked to the set {HIV, cancer} with a high probability. However, HIV and Cancer are not under the same category in the taxonomy. For this more general case, we introduce the *sensitive value set*, which is a set of ground values in the taxonomy for the sensitive attribute. In such a taxonomy, there can be multiple sensitive value sets.

A major technique used in previous works is to recode the QID values in such a way that a set of individuals will be matched to the same generalized QID value and, in the set, the occurrence of values in any sensitive value set is not frequent. Hence, the records with the same QID value (which could be a generalized value) is of interest. In a table T , the equality of the QID values determines an equivalence relation on the set of tuples in T . A QID equivalence class, or simply QID-EC, is a set of tuples in T with identical QID value. For simplicity, we also refer to a QID-EC by the identical QID value.

Definition 1 (Anonymization) *Anonymization is a mapping from a table T to an anonymized table T^* , such that a one-to-one mapping function f maps each tuple t in T to a tuple t^* in T^* , where $f(t) = t^*$. We define $t^*.A$ (or $f(t).A$) to be the value of attribute A of tuple t^* (or $f(t)$). Given a set of taxonomies $\tau = \{\mathcal{T}_1, \dots, \mathcal{T}_u\}$, we say that an anonymization defined by f conforms to τ iff either $t.A \prec f(t).A$ or $t.A = f(t).A$ holds for any t and A .*

For instance, Table 1(b) is anonymized to Table 1(c). The mapping function f is mapping the tuples with $q1$ and $q2$ to Q .

Let K_{ad} be the knowledge of the adversary. In most previous work [20, 11, 12, 26], in addition to the published data set T^* , K_{ad} involves an external table T^e such as Voter registration list that maps QIDs to individuals. In the literature, two possible cases of T^e have been considered: (1) *Worst Case*: the set of individuals in the external table T^e is equal to the set of individuals in the original table T ; (2) *Superset Case*: the set of individuals in the external table T^e is a proper superset of the set of individuals in the original table T . Assuming the worst case scenario is the safest stance and it has been the assumption in most previous works. We have shown in our first two examples that, in either of the above 2 cases, minimality attack is possible.

The objective of privacy preservation is to limit the probability of the linkage from any individual to any sensitive value set s in the sensitive attribute. We define this probability or *credibility* as follows.

Definition 2 (Credibility) *Let T^* be a published table which is generated from T . Consider an individual $o \in O$ and a sensitive value set s in the sensitive attribute. $Credibility(o, s, K_{ad})$ is the probability that an adversary can infer from T^* and background knowledge K_{ad} that o is associated with s .*

The particular background knowledge we are concern here is about the minimality principle described in Assumption 3, and is formally defined as follows.

Definition 3 (Minimality Principle) *Suppose \mathcal{A} is an algorithm for anonymization for a privacy requirement \mathcal{R} which follows the minimality principle. Let table T^* be a table generated by \mathcal{A} and satisfies \mathcal{R} . Then, for any QID-EC X in T^* , there is no specialization (reverse of generalization) of the QID's in X which results in another table T' which also satisfies l -diversity.*

Note that this minimality principle holds for both global recoding and local recoding. If \mathcal{A} is for global recoding (local recoding), both T^* and T' are global recoding (local recoding). So far we focus on the privacy requirement of l -diversity. However, in Section 5, we shall consider cases where \mathcal{R} is equal to other requirements.

Assumption 3 (Adversary knowledge K_{ad}^{min}) *In the above definition of $Credibility(o, s, K_{ad})$, we consider the cases where K_{ad} includes T^* , the multiset T^q containing all QID occurrences in the table T , the QID values of a target individual in T , a set of taxonomies τ and whether the anonymization \mathcal{A} conforms to the taxonomies τ , the target privacy requirement \mathcal{R} , and whether \mathcal{A} follows the minimality principle. We refer to this knowledge as K_{ad}^{min} .*

If Table 1(a) is the result generated from an anonymization mechanism (e.g. the adapted Incognito algorithm in [14]) for l -diversity that follows the minimality principle, suppose the multiset in Table 2(b) is known and the QID value of individual o is known to be $q1$, then $Credibility(o, \{HIV\}, K_{ad}^{min}) = 1/2$. When the same K_{ad}^{min} is applied to Table 1(c), $Credibility(o, \{HIV\}, K_{ad}^{min}) = 1$. Section 4 describes how to compute the credibility.

The above minimality principle is very general and does not demand that \mathcal{A} minimizes the overall information loss, nor does it depend on how the information loss is defined. Almost all known anonymization algorithms (including Incognito based methods [12, 14, 15, 13] and top-down approaches [7, 26, 24, 21]) try to reduce information loss of one form or another, and they all follow the above principle.

In the examples in Section 1, the value of l (for l -diversity) is used by the adversary. However, l is not included in K_{ad}^{min} . This is because, in many cases, it can be deduced from the published table T^* . For example, for the anonymization in Table 1(d), the adversary can deduce that l must be 2. For Table 3(c), l must be 2 since $q4$ is linked to 2 different values only.

Definition 4 (m -confidentiality) *A table T is said to satisfy m -confidentiality (or T is m -confidential) if, for any individual o and any sensitive value set s , $Credibility(o, s, K_{ad})$ does not exceed $1/m$.*

For example, Tables 1(a) and 3(a) satisfy 2-confidentiality.

When a table T is anonymized to a more generalized table T^* , it is of interest to measure the information loss that is incurred. There are different ways to define information loss. Since we will measure the effectiveness of our method based on the method in [26], we also adopt a similar measure of information loss.

Definition 5 (Coverage and Base) *Let \mathcal{T} be the taxonomy for an attribute in QID. The coverage of a generalized QID value v^* , denoted by $coverage[v^*]$, is given by the number of ground values v' in \mathcal{T} such that $v' < v^*$. The base of the taxonomy \mathcal{T} , denoted by $base(\mathcal{T})$, is the number of ground values in the taxonomy.*

For example, in Figure 1, $coverage[\text{“university”}] = 2$ since “undergrad” and “postgrad” can be generalized to “university”, $base(\mathcal{T}) = 9$.

A weighting can be assigned for each attribute A , denoted by $weight(A)$, to reflect the users' opinion on the significance of information loss in different attributes. Let $t.A$ denote the value of A in tuple t .

QID	Disease
$q1$	HIV
$q1$	HIV
$q2$	HIV
$q2$	non-sensitive
$q3$	HIV
$q3$	HIV
$q3$	non-sensitive
$q3$	non-sensitive
$q3$	non-sensitive
...	...
$q3$	non-sensitive

Table 4: A table which violates 2-diversity

QID	Disease
Q	HIV
Q	HIV
Q	HIV
Q	non-sensitive
Q	HIV
Q	HIV
Q	non-sensitive
Q	non-sensitive
Q	non-sensitive
...	...
Q	non-sensitive

Table 5: A 2-diverse table by global recoding of Table 4

Definition 6 (Information Loss) Let T^* be an anonymization of table T by means of a mapping function f . Let \mathcal{T}_A be the taxonomy for attribute A which is used in the mapping and v^* be the nearest common ancestor of $t.A$ and $f(t).A$ in \mathcal{T}_A . The information loss of a tuple t^* in T^* introduced by f is given by

$$\mathcal{IL}(t^*) = \sum_{A \in QID} \left\{ \frac{\text{coverage}[v^*] - 1}{\text{base}(\mathcal{T}_A) - 1} \times \text{weight}(A) \right\}$$

The information loss is given by $\text{Dist}(T, T^*) = \frac{\sum_{t^* \in T^*} \mathcal{IL}(t^*)}{|T^*|}$

If $f(t).A = t.A$, then $f(t).A$ is a ground value, the nearest common ancestor $v^* = t.A$, and $\text{coverage}[v^*] = 1$. If this is true for all A 's in QID , then $\mathcal{IL}(t^*)$ is equal to 0, which means there is no information loss. If $t.A$ is generalized to the root of taxonomy \mathcal{T}_A , then the nearest common ancestor $v^* = \text{the root of } \mathcal{T}_A$. Thus, $\text{coverage}[v^*] = \text{base}(\mathcal{T}_A)$ and, if this is the case for all A 's in QID , then $\mathcal{IL}(t^*) = 1$. Note that we have modified the definition in [26] in order to achieve the range of $[0,1]$ for $\mathcal{IL}(t^*) = 1$ and also for $\text{Dist}(T, T^*)$.

Although minimizing information loss poses a loophole for attack by minimality, one cannot completely ignore information loss since, without such a notion, we allow for complete distortion of the data which will also render the published data useless.

Definition 7 (PROBLEM) *Optimal m -confidentiality: Given a table T , generate an anonymized table T^* from T which satisfies m -confidentiality where the information loss $\text{Dist}(T, T^*)$ is minimized.*

4 Credibility: Source of Attack

In this section, we characterize the nature of minimality attack. Minimality attack is successful if the adversary can compute the credibility values and find a violation of m -confidentiality when the privacy requirement is l -diversity. This computation depends on a combinatorial analysis on the possibilities given the knowledge of K_{ad}^{min} . In particular, the adversary attacks by excluding some possible scenarios, tilting the probabilistic balance towards privacy disclosure.

4.1 Global Recoding

The derivation of credibility is better illustrated with the example as shown in Table 5 which is a global recoding of Table 4 to achieve 2-diversity. In Table 4, {HIV} is the only sensitive value set and the

	Number of sensitive tuples			Total number of cases
	$q1$	$q2$	$q3$	
(a)	2	0	3	120
(b)	2	1	2	90
(c)	2	2	1	10
(d)	1	2	2	90
(e)	0	2	3	120

Table 6: Possible combinations of number of sensitive tuples

goal is 2-diversity. Assume that T and T^e have *matching cardinality* on Q . From T^e , the adversary can determine that there are two tuples in $q1$, two tuples in $q2$ and 10 tuples in $q3$. Since there are 10 tuples with a QID value of $q3$, and there are in total 5 sensitive tuples, $q3$ trivially satisfies 2-diversity. As T^* (Table 5) is generalized, the adversary decides that at least one of the QID-EC's $q1$ and $q2$ contains two sensitive tuples. With this in mind, the adversary lists all the possible combinations of the number of sensitive tuples among the three classes $q1$, $q2$ and $q3$ in which either $q1$ or $q2$ or both contain 2 sensitive tuples as shown in Table 6. There are only five possible combinations as shown. We call this table as the *sensitive tuple distribution table*.

In scenario (a), there are $C_2^2 \times C_0^2 \times C_3^{10} = 120$ different possible ways to assign the sensitive values to the tuples. In scenario (b), there are $C_2^2 \times C_1^2 \times C_2^{10} = 90$ different *assignments* or *cases*⁶. Similarly, there are 10 cases, 90 cases and 120 cases in scenarios (c), (d) and (e), respectively. The total number of cases is equal to $120 + 90 + 10 + 90 + 120 = 430$. Without any additional knowledge about the assignments, one must assume that each of these cases occurs with the same probability $1/430$. Consider the credibility that an individual o with value $q1$ is linked to HIV given K_{ad}^{min} . There are two possible cases.

- *Case 1:* There are two sensitive tuples in $q1$.

The total number of cases where there are two sensitive tuples in $q1$ is equal to $120 + 90 + 10 = 220$. The probability that Case 1 occurs given K_{ad}^{min} is equal to $220/430 = 0.5116$.

- *Case 2:* There is one sensitive tuple in $q1$.

The total number of cases where there is one sensitive tuple in $q1$ is equal to 90. The probability that Case 2 occurs given K_{ad}^{min} is equal to $90/430 = 0.2093$.

In the following, we use $Prob(E)$ to stand for the probability that event E occurs.

Thus, the credibility that an individual o with QID value $q1$ is linked to HIV given K_{ad}^{min} is equal to

$$\begin{aligned}
& Prob(\text{Case 1}) \times Prob(q1 \text{ is linked to HIV in Case 1}) \\
& + Prob(\text{Case 2}) \times Prob(q1 \text{ is linked to HIV in Case 2})
\end{aligned}$$

$Prob(q1 \text{ is linked to HIV in Case 1})$ is equal to $2/2=1$.

$Prob(q1 \text{ is linked to HIV in Case 2})$ is equal to $1/2=0.5$.

$$\begin{aligned}
Credibility(o, \{HIV\}, K_{ad}^{min}) &= 0.5116 \times 1 + 0.2093 \times 0.5 \\
&= 0.616,
\end{aligned}$$

which is greater than 0.5. This result shows that the published table violates 2-confidentiality.

General Formula

The general formula of the computation of the credibility is based on the idea illustrated above. We have a probability space $(\Omega, \mathcal{F}, \mathcal{P})$, where Ω is the set of all possible assignments of the sensitive values to the

⁶In the following, we refer assignments as cases.

tuples. \mathcal{F} is the power set of Ω . P is a probability mass function from \mathcal{F} to the real numbers in $[0,1]$ which gives the probability for each element in \mathcal{F} . Given K_{ad}^{min} , there will be a set of assignments \mathcal{G} in Ω which are impossible or $P(\mathcal{G}) = 0$ and if $x \in \mathcal{G}$ then $P(\{x\}) = 0$. Without any other additional knowledge, we assume that the probability of the remaining assignments are equal. That is, $\mathcal{G}' = \Omega - \mathcal{G}$, $P(\mathcal{G}') = 1$ and for $x \in \mathcal{G}'$, $P(\{x\}) = 1/|\mathcal{G}'|$.

Definition 8 Let Q be a QID-EC in T^* . Tables T^* and T^e have matching cardinality on Q if the number of tuples in T^e with QID that can be generalized to Q is the same as the that in T^* .

Let \mathcal{X} be a maximal set of QID-EC's in T which are generalized to the same QID-EC Q in the published table T^* . Suppose T^* and T^e have matching cardinality on Q . Let C_1, C_2, \dots, C_u be the QID-EC's in \mathcal{X} sorted in ascending order of the size of the QID-EC's. Let n_i be the number of tuples in class C_i . Hence, $n_1 \leq n_2 \leq \dots \leq n_u$. Let n_s be the total number of tuples with values in sensitive value set s in the data set.

In Table 4, there are three classes, namely $q1, q2$ and $q3$. Thus, $u = 3$. C_1 corresponds to $q1$, C_2 corresponds to $q2$ and C_3 corresponds to $q3$. Also, $n_1 = 2, n_2 = 2$ and $n_3 = 10$.

Suppose the published table is generalized in order to satisfy the l -diversity requirement.

If $n_s \leq \lfloor \frac{n_i}{l} \rfloor$, then C_i in the original data set must satisfy the l -diversity requirement without any generalization. Class C_i may violate the l -diversity requirement only if $n_s > \lfloor \frac{n_i}{l} \rfloor$. Let \mathcal{C} be the set of all classes C_i where $n_s > \lfloor \frac{n_i}{l} \rfloor$. Let \mathcal{C}' be the set of the remaining classes. Let p be the total number of classes in \mathcal{C} . Since the classes are sorted, $\mathcal{C} = \{C_1, C_2, \dots, C_p\}$ and $\mathcal{C}' = \{C_{p+1}, C_{p+1}, \dots, C_u\}$.

Lemma 1 If the classes $\mathcal{X} = \{C_1, \dots, C_u\}$ have been generalized to their parent class in \mathcal{T} , the adversary can deduce that at least one class (in the original table) violates l -diversity among \mathcal{C} and all classes in \mathcal{C}' (in the original table) do not violate l -diversity.

Obviously, the credibility of individuals in a class in \mathcal{C}' is smaller than or equal to $\frac{1}{l}$. However, the credibility of individuals in a class in \mathcal{C} may be greater than $\frac{1}{l}$. Thus, the adversary tries to compute $Credibility(o, s, K_{ad}^{min})$, where $o \in C_i$, for $i = 1, 2, \dots, p$. Suppose there are j tuples with the sensitive value set s in C_i . Let $|C_i(s)|$ denote the number of occurrences of the tuples with s in C_i . The probability that o is linked to a sensitive value set is $\frac{j}{n_i}$, where n_i is the class size of C_i . Let $Prob(|C_i(s)| = j | K_{ad}^{min})$ be the probability that there are exactly j occurrences of tuples with s in C_i given K_{ad}^{min} . By considering all possible number j of occurrences of tuples with s from 1 to n_i in C_i , the general formula for credibility is given by:

$$\begin{aligned} & Credibility(o, s, K_{ad}^{min}), \text{ where } o \in C_i, 1 \leq i \leq p \\ &= Prob(o \text{ is linked to } s \text{ in } C_i | K_{ad}^{min}) \\ &= \sum_{j=1}^{n_i} Prob(|C_i(s)| = j | K_{ad}^{min}) \times \frac{j}{n_i} \end{aligned}$$

In the above formula, $Prob(|C_i(s)| = j | K_{ad}^{min})$ can be calculated by considering all possible cases. Conceptually, a table such as Table 6 will be constructed, in which some possible combinations will be excluded due to the minimality notion in K_{ad}^{min} .

More specifically, in the above formula, the calculation of the term $Prob(|C_i(s)| = j | K_{ad}^{min})$ can be done by a dynamic programming approach. Before describing how to make use of a dynamic programming approach, we define the following events.

F_i	$0 \leq n_{i,s} \leq \lfloor \frac{n_i}{l} \rfloor$
G_i	$\lfloor \frac{n_i}{l} \rfloor + 1 \leq n_{i,s} \leq n_i$
H_i	$0 \leq n_{i,s} \leq n_i$

We illustrate the events in Figure 2. We can see that $F_i \cup G_i = H_i$.

Now, we like to evaluate $Prob(|C_i(s)| = j | K_{ad}^{min})$. Let E_1 be the event that at least one C_j among C_1, C_2, \dots, C_p violates l -diversity. Let E_2 be the event that there are j sensitive values in C_i (i.e. $|C_i(s)| =$

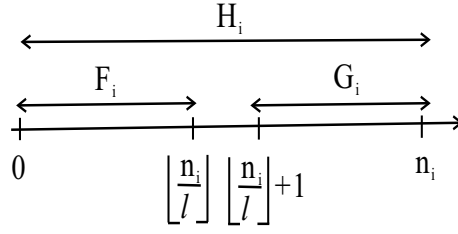


Figure 2: Illustration of Some Events

j). We illustrate with Figure 3. As the published table is generalized, the adversary should know that event E_1 should occur. That means the adversary does not need to guess the sensitive value of an individual from all possible cases. Instead, he/she can deduce the sensitive value of an individual from the total number of cases that E_1 occurs. Now, we consider $\text{Prob}(|C_i(s)| = j \mid K_{ad}^{min})$. Thus, among the total number of cases that E_1 occurs, we want to know the total number of cases that E_2 occurs. Thus, we obtain the following formula.

$$\begin{aligned}
& \text{Prob}(|C_i(s)| = j \mid K_{ad}^{min}) \\
= & \text{Prob}(E_2|E_1) \\
= & \frac{\text{total no. of cases that both } E_1 \text{ and } E_2 \text{ occur}}{\text{total no. of cases that } E_1 \text{ occurs}}
\end{aligned}$$

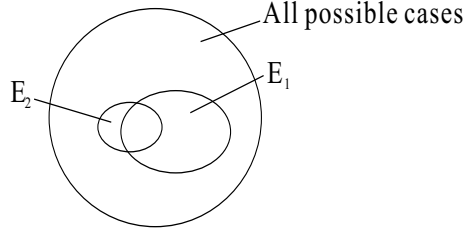


Figure 3: Illustration for Probability with K_{ad}^{min}

Let E_3 be the event at least one G_j occurs among G_1, G_2, \dots, G_p . As $E_1 = E_3$, we have the following formula.

$$\begin{aligned}
& \text{Prob}(|C_i(s)| = j \mid K_{ad}^{min}) \\
= & \frac{\text{total no. of cases that both } E_3 \text{ and } E_2 \text{ occur}}{\text{total no. of cases that } E_3 \text{ occurs}}
\end{aligned}$$

E_1	at least one C_j among C_1, C_2, \dots, C_p violates l -diversity
E_2	there are j sensitive values in C_i (i.e. $ C_i(s) = j$)
E_3	at least one G_j occurs among G_1, G_2, \dots, G_p

Let A be total no. of cases that both E_3 and E_2 occur. Let B be total no. of cases that E_3 occurs.

A	total no. of cases that both E_3 and E_2 occur
B	total no. of cases that E_3 occurs

In the following, we consider the total number of cases in the records in C_1, C_2, \dots, C_p only. We do not consider the total number of cases in the records in $C_{p+1}, C_{p+2}, \dots, C_u$ in the following for the sake

of illustration. The reason is described as follows. Assume there are x sensitive values in C_1, C_2, \dots, C_p . Suppose we obtain the total number of cases in the records in C_1, C_2, \dots, C_p equal to Q by our dynamic programming which will be described as follows. We can easily obtain the total number of cases in the records in all classes (i.e. $C_1, C_2, \dots, C_p, C_{p+1}, \dots, C_u$) by multiplying Q by $C_{n_s-x}^N$, where N is the total number of records in $C_{p+1}, C_{p+2}, \dots, C_u$ and n_s is the total number of sensitive values in the total data set.

In our dynamic programming, we make use of three variables, namely $n([a, b], x)$, $m([a, b], x)$ and $u([a, b], x)$, for the computation of A and B used in $\text{Prob}(|C_i(s)| = j \mid K_{ad}^{min})$.

1. $n([a, b], x)$ is the total no. of cases that at least one G_j occurs among G_a, G_{a+1}, \dots, G_b when there are x sensitive values in C_a, C_{a+1}, \dots, C_b , for $a, b = 1, 2, \dots, p$ and $x = 1, 2, \dots, n_s$.
2. $m([a, b], x)$ is the total no. of cases that H_a and H_{a+1} and ... and H_b occur when there are x sensitive values in C_a, C_{a+1}, \dots, C_b , for $a, b = 1, 2, \dots, p$ and $x = 1, 2, \dots, n_s$.
3. $u([a, b], x)$ is the total no. of cases that F_a and F_{a+1} and ... and F_b occur when there are x sensitive values in C_a, C_{a+1}, \dots, C_b , for $a, b = 1, 2, \dots, p$ and $x = 1, 2, \dots, n_s$.

$n([a, b], x)$	total no. of cases that at least one G_j occurs among G_a, G_{a+1}, \dots, G_b when there are x sensitive values in C_a, C_{a+1}, \dots, C_b
$m([a, b], x)$	total no. of cases that H_a and H_{a+1} and ... and H_b occur when there are x sensitive values in C_a, C_{a+1}, \dots, C_b
$u([a, b], x)$	total no. of cases that F_a and F_{a+1} and ... and F_b occur when there are x sensitive values in C_a, C_{a+1}, \dots, C_b

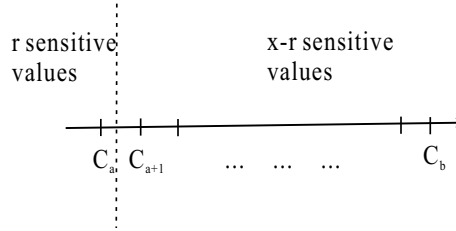


Figure 4: Illustration of $n([a, b], x)$, $m([a, b], x)$ and $u([a, b], x)$

Consider $m([a, b], x)$. Among C_a, C_{a+1}, \dots, C_b , we divide the classes into two parts. The first part contains the first class C_a and the second part contains the rest of the classes, C_{a+1}, \dots, C_b . See Figure 4. Suppose we allocate r sensitive values in G_a and $x - r$ sensitive values in C_{a+1}, \dots, C_b , among x sensitive values. The number of cases that there are r sensitive values in class C_a of size n_a is equal to $C_r^{n_a}$. The total number of cases that G_{a+1}, \dots, G_b occur when the number of sensitive values allocated to them is equal to $x - r$ is equal to $m([a + 1, b], x - r)$. Thus, for a given r , the total number of cases is equal to $C_r^{n_a} \times m([a + 1, b], x - r)$.

$$m([a, b], x) = \sum_{r=0}^{n_a} C_r^{n_a} \times m([a + 1, b], x - r)$$

We define the base cases of $m([a, b], x)$ as follows. The base case happens when $a = b$. It is impossible that the number of sensitive values allocated to C_a is greater than the class size of C_a or smaller than 0. Thus the term should be set to 0. If the number of sensitive values allocated to C_a ranges from 0 to n_a , the term should be equal to the number of possible combinations that there are x sensitive values in class C_a of size n_a (i.e. $C_x^{n_a}$).

$$m([a, a], x) = \begin{cases} 0 & \text{if } x > n_a \\ 0 & \text{if } x < 0 \\ C_x^{n_a} & \text{if } 0 \leq x \leq n_a \end{cases}$$

Consider $u([a, b], x)$. As $u([a, b], x)$ is same as $m([a, b], x)$ except the upper boundary of term F_i is equal to $\lfloor \frac{n_a}{l} \rfloor$, instead of n_a . Similarly, we obtain the following formula.

$$u([a, b], x) = \sum_{r=0}^{\lfloor \frac{n_a}{l} \rfloor} C_r^{n_a} \times u([a+1, b], x-r)$$

$$u([a, a], x) = \begin{cases} 0 & \text{if } x \geq \lfloor \frac{n_a}{l} \rfloor + 1 \\ 0 & \text{if } x < 0 \\ C_x^{n_a} & \text{if } 0 \leq x \leq \lfloor \frac{n_a}{l} \rfloor \end{cases}$$

Consider $n([a, b], x)$. Among C_a, C_{a+1}, \dots, C_b , we select the first class C_a to consider whether G_a occurs or not. Suppose we allocate r sensitive values in G_a and $x-r$ sensitive values in C_{a+1}, \dots, C_b , among x sensitive values. See Figure 4. There are two cases.

Case (1): $\lfloor \frac{n_a}{l} \rfloor + 1 \leq r \leq n_a$. Thus, G_a occurs. In this case, as event G_a occurs, the number of sensitive values in each C_i of C_{a+1}, \dots, C_b can range from 0 to n_i , where n_i is the class size of C_i . The number of cases that the number of sensitive values in each C_i of C_{a+1}, \dots, C_b can range from 0 to n_i is equal to $m([a+1, b], x-r)$. As the number of cases that there are r sensitive values in C_a of size n_a is equal to $C_r^{n_a}$, the total number of cases is equal to $C_r^{n_a} \times m([a+1, b], x-r)$.

Case (2): $0 \leq r \leq \lfloor \frac{n_a}{l} \rfloor$. Thus, G_a does not occur. In this case, as event G_a does not occur, at least one G_j occurs among G_{a+1}, \dots, G_b . The number of cases that at least one G_j occurs among G_{a+1}, \dots, G_b is equal to $n([a+1, b], x-r)$. Similarly, as the number of cases that there are r sensitive values in C_a of size n_a is equal to $C_r^{n_a}$, the total number of cases is equal to $C_r^{n_a} \times n([a+1, b], x-r)$.

Combining the above two cases, we obtain the following formula.

$$n([a, b], x) = \sum_{r=\lfloor \frac{n_a}{l} \rfloor + 1}^{n_a} C_r^{n_a} \times m([a+1, b], x-r) + \sum_{r=0}^{\lfloor \frac{n_a}{l} \rfloor} C_r^{n_a} \times n([a+1, b], x-r)$$

We define the base cases of $n([a, b], x)$ as follows. The base case happens when $a = b$. It is impossible that the number of sensitive values allocated to C_a is greater than the class size of C_a . The term should be set to 0. It is impossible that the number of sensitive values allocated to C_a is smaller than or equal to $\lfloor \frac{n_a}{l} \rfloor$, because G_a does not occur in this case. Thus, the term should also be set to 0. If the number of sensitive values allocated to C_a ranges from $\lfloor \frac{n_a}{l} \rfloor + 1$ to n_a , the term should be equal to the number of combinations that there are x sensitive values in class C_a of size n_a (i.e. $C_x^{n_a}$).

$$n([a, a], x) = \begin{cases} 0 & \text{if } x > n_a \\ 0 & \text{if } 0 \leq x \leq \lfloor \frac{n_a}{l} \rfloor \\ C_x^{n_a} & \text{if } \lfloor \frac{n_a}{l} \rfloor + 1 \leq x \leq n_a \end{cases}$$

Now, consider A . Recall that A is the total number of cases that (1) at least one G_j occurs among G_1, G_2, \dots, G_p (which corresponds to E_3) and (2) there are j sensitive values in C_i (which corresponds to E_2).

Let $\mathcal{A}(x)$ be the total number of above cases provided that there are x sensitive values in C_1, C_2, \dots, C_p .

Here, we consider the number of cases involving all classes (i.e. $C_1, \dots, C_p, C_{p+1}, \dots, C_u$). Suppose we allocate x sensitive values in C_1, C_2, \dots, C_p and $n_s - x$ sensitive values in $C_{p+1}, C_{p+2}, \dots, C_u$. Among C_1, C_2, \dots, C_p , we divide into three parts: (1) C_1, C_2, \dots, C_{i-1} , (2) C_i and (3) $C_{i+1}, C_{i+2}, \dots, C_p$. Recall that C_i contains j sensitive values. Within C_1, C_2, \dots, C_p , we further allocate (1) r sensitive values to C_1, C_2, \dots, C_{i-1} , (2) j sensitive values to C_i , and (3) $x - r - j$ sensitive values to $C_{i+1}, C_{i+2}, \dots, C_p$. See Figure 5.

There are two cases.

Case 1: $\lfloor \frac{n_i}{l} \rfloor + 1 \leq j \leq n_i$. That is, G_i occurs. This means that the number of sensitive values in a class C_y of C_1, C_2, \dots, C_{i-1} ranges from 0 to n_y . Also, the number of sensitive values in a class C_y of $C_{i+1}, C_{i+2}, \dots, C_p$ ranges from 0 to n_y .

The number of cases that the number of sensitive values in a class C_y of C_1, C_2, \dots, C_{i-1} ranges from 0 to n_y is equal to $m([1, i-1], r)$. Similarly, the number of cases that the number of sensitive values in a class C_y of $C_{i+1}, C_{i+2}, \dots, C_p$ ranges from 0 to n_y is equal to $m([i+1, p], x-r-j)$. Thus,

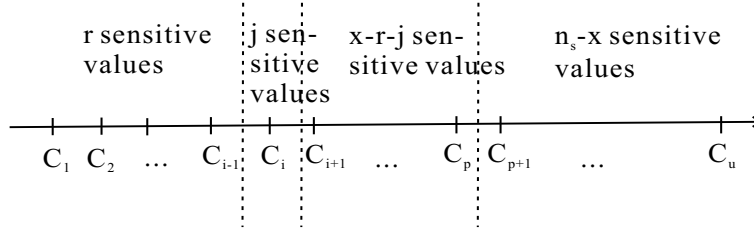


Figure 5: Illustration of $\mathcal{A}(x)$

if we consider all possible values of r from 0 to $x - j$, the total number of these cases is equal to $\sum_{r=0}^{x-j} m([1, i-1], r) \times m([i+1, p], x-r-j)$.

Note that the number of cases that there are j sensitive values in C_i of size n_i is equal to $C_j^{m_i}$. Also, the number of cases that there are $n_s - x$ sensitive values in N tuples in classes $C_{p+1}, C_{p+2}, \dots, C_u$ is equal to $C_{n_s-x}^N$. Thus, $\mathcal{A}(x) = C_{n_s-x}^N \times C_j^{m_i} \times \sum_{r=0}^{x-j} m([1, i-1], r) \times m([i+1, p], x-r-j)$ in this case.

Case 2: $0 \leq j \leq \lfloor \frac{n_i}{t} \rfloor$. That is, G_i does not occur. There are the following sub-cases.

Case 2(a): At least one G_z occurs among G_1, G_2, \dots, G_{i-1} . In this case, the number of sensitive values in a class C_y of $C_{i+1}, C_{i+2}, \dots, C_p$ ranges from 0 to n_y .

The number of cases that at least one G_z occurs among G_1, G_2, \dots, G_{i-1} is equal to $n([1, i-1], r)$. The number of cases that the number of sensitive values in a class C_y of $C_{i+1}, C_{i+2}, \dots, C_p$ ranges from 0 to n_y is equal to $m([i+1, p], x-r-j)$. Thus, the total number of these cases is equal to $n([1, i-1], r) \times m([i+1, p], x-r-j)$.

Case 2(b): All G_z among G_1, G_2, \dots, G_{i-1} does not occur. In other words, all F_1, F_2, \dots, F_{i-1} occur. Besides, we should also know that there is at least one G_z occurs among $G_{i+1}, G_{i+2}, \dots, G_p$.

The number of cases that all F_1, F_2, \dots, F_{i-1} occur is equal to $u([1, i-1], r)$. The number of cases that at least one G_z occurs among $G_{i+1}, G_{i+2}, \dots, G_p$ is equal to $n([i+1, p], x-r-j)$. Thus, the total number of these cases is equal to $u([1, i-1], r) \times n([i+1, p], x-r-j)$.

By combining Case 2(a) and Case 2(b) and considering all possible values r from 0 to $x-j$, we obtain the total number of cases equal to $\sum_{r=0}^{x-j} [n([1, i-1], r) \times m([i+1, p], x-r-j) + u([1, i-1], r) \times n([i+1, p], x-r-j)]$.

Similarly, the number of cases that there are j sensitive values in C_i of size n_i is equal to $C_j^{m_i}$. Also, the number of cases that there are $n_s - x$ sensitive values in N tuples in classes $C_{p+1}, C_{p+2}, \dots, C_u$ is equal to $C_{n_s-x}^N$. Thus, the total number of cases in Case (2) is equal to $C_{n_s-x}^N \times C_j^{m_i} \times \sum_{r=0}^{x-j} [n([1, i-1], r) \times m([i+1, p], x-r-j) + u([1, i-1], r) \times n([i+1, p], x-r-j)]$.

We obtain $\mathcal{A}(x)$ as follows.

$\mathcal{A}(x)$

$$= \begin{cases} C_{n_s-x}^N \times C_j^{m_i} \times \sum_{r=0}^{x-j} m([1, i-1], r) \times m([i+1, p], x-r-j) & \text{if } \lfloor \frac{n_i}{t} \rfloor + 1 \leq j \leq n_i \\ C_{n_s-x}^N \times C_j^{m_i} \times \sum_{r=0}^{x-j} [n([1, i-1], r) \times m([i+1, p], x-r-j) + u([1, i-1], r) \times n([i+1, p], x-r-j)] & \text{if } 0 \leq j \leq \lfloor \frac{n_i}{t} \rfloor \end{cases}$$

By considering all possible values x from $\lfloor \frac{n_1}{t} \rfloor + 1$ to n_s , A is equal to the following. Note that it is impossible that $x < \lfloor \frac{n_1}{t} \rfloor + 1$. This is because, if $x < \lfloor \frac{n_1}{t} \rfloor + 1$, E_3 does not occur.

$$A = \sum_{x=\lfloor \frac{n_1}{t} \rfloor + 1}^{n_s} \mathcal{A}(x)$$

Consider B where B is the total number of cases that at least one G_j occurs among G_1, G_2, \dots, G_p . By considering all possible values x from $\lfloor \frac{n_1}{t} \rfloor + 1$ to n_s , we obtain the following formula.

$$B = \sum_{x=\lfloor \frac{n_1}{t} \rfloor}^{n_s} n([1, p], x) \times C_{n_s-x}^N$$

4.1.1 Generalization Hierarchy of Height > 1

In the last section, we derived a general formula of credibility with the generalization hierarchy of height equal to 1. In the above discussion, the generalization hierarchy we discussed so far has height of 1.

Let us summarize the derivation of the general formula as follows. Suppose we have u classes in the original data set T , namely C_1, C_2, \dots, C_u . The number of tuples in class C_i is equal to n_i . After the traditional anonymization, we obtain a generalized data set T^* . In T^* , there are u' equivalence classes, namely $E_1, E_2, \dots, E_{u'}$. If the data set is generalized and we have the generalization hierarchy of height = 1 as shown in Figure 6, we have $u' = 1$. That is, there is only one equivalence class E_1 in T^* . We call this P . Let n_P be the number of tuples in P . Let $n_{P,s}$ be the number of sensitive tuples in P . Let $v(P)$ be the value of class P . In this case, n_P is equal to the number of tuples with value equal to $v(P)$ in T^* (i.e. the number of tuples in T^* because there is only one equivalence class). $n_{P,s}$ is equal to the number of sensitive tuples with value equal to $v(P)$ in T^* .

Note that we obtain the following information.

1. n_P
2. $n_{P,s}$
3. n_1, n_2, \dots, n_u

n_P	the number of tuples in P
$n_{P,s}$	the number of sensitive tuples in P
$v(P)$	the value of class P
n_i	the number of tuples in class C_i (which is a child of P)

After we obtain the above information, we compute the credibility for each class C_i , denoted by p_i , according to the formula we derived in the last section. Let us illustrate with a black-box as shown in Figure 7. We call the step of the computation of the credibility in this simple scenario as function V .

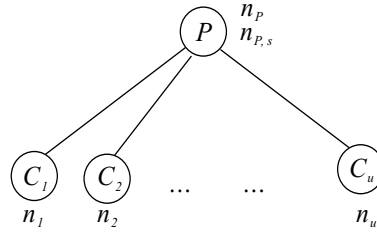


Figure 6: Generalization hierarchy of classes C_1, C_2, \dots, C_u

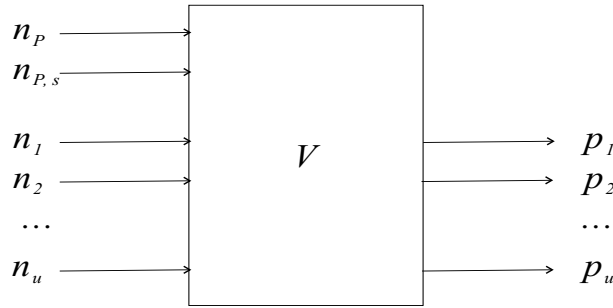


Figure 7: Function V

However, in general, the generalization hierarchy has height greater than one. The generalization hierarchy of our running example can be modified as in Figure 1. Suppose the original table T contains different education levels such as “undergrad” and “postgrad”. The generalized table T^* contains multiple

equivalence classes, instead of only one equivalence class. Let the equivalence classes be $E_1, E_2, \dots, E_{u'}$, where u' is the number of equivalence classes in T^* .

Let us consider the simple case that T^* contains the values which are generalized one level upwards. Then, we discuss a complicated case that T^* contains the values which are generalized multiple levels upwards next.

1. Value Generalized One Level Upwards

Let us consider the case that the generalized value in T^* is just one level upwards compared with the original value in T according to the generalization hierarchy. Suppose we have the original table T as follows.

ID	T	$Sensitive$
1	1st-4th	HIV
2	1st-4th	none
3	5th-6th	HIV
4	5th-6th	HIV
5	7th-8th	none
6	7th-8th	none
7	9th-10th	HIV
8	9th-10th	none

We may generalize the table as follows. The following generalized table contains the value generalized one level upwards. This is the case we are studying in this section.

ID	T	$Sensitive$
1	elementary	HIV
2	elementary	none
3	elementary	HIV
4	elementary	HIV
5	elementary	none
6	elementary	none
7	secondary	HIV
8	secondary	none

We can also make use of the function V to calculate the credibility. Suppose we calculate the credibility for (o, s, K_{ad}^{min}) , where $o \in C_i$. As the table is generalized, there exists an equivalence class E_j such that E_j is the generalized class of C_i . For example, if we want to know the credibility of an individual in “1st-4th”, $C_i = 1st - 4th$. Then, the equivalence class E_j in T^* is “elementary”. Let us map this situation to the function V . P is mapped with E_j . n_P is the number of tuples in P (i.e. E_j). $n_{P,s}$ is the number of sensitive tuples in P (i.e. E_j).

Note that, in this case, n_P is equal to the number of tuples with value equal to $v(P)$ in T^* , and $n_{P,s}$ is equal to the number of sensitive tuples with value equal to $v(P)$ in T^* . For example, n_P is the number of tuples with “elementary” in T^* (which is equal to 6) and $n_{P,s}$ is the number of sensitive tuples with “elementary” in T^* (which is equal to 3).

The adversary wants to find the number of tuples in each child (or specialized value) of P (or E_j). For example, if $P = elementary$, there are three children, 1st-4th, 5th-6th and 7th-8th. Let u'' be the number of children of E_j . In order to obtain the number of tuples of each child c_y of E_j , we perform the following operations. As the generalization is global, all tuples with the same original value in T should belong to only one of $E_1, E_2, \dots, E_{u'}$. For example, 1st-4th only belongs to elementary but not secondary.

Note that the adversary knows two tables T^e and T^* with matching cardinality.

Let D_X^e be the set of tuples in T^e which can be generalized to $v(X)$, where X is a (generalized/ungeneralized) class. For example, if $X = \text{"1st - 4th"}$, $|D_X^e| = 2$. If $X = \text{"elementary"}$, $|D_X^e| = 6$. If $X = \text{"without post-secondary"}$, $|D_X^e| = 8$. The number of tuples in child c_y is equal to $|D_{c_y}^e|$, where $y = 1, 2, \dots, u''$.

T^e	the external table
T^*	the published table
D_X^e	the set of tuples in T^e which can be generalized to $v(X)$, where X is a (generalized/ungeneralized) class

Then, we can call function V to calculate the credibility by mapping $n_y = |D_{c_y}^e|$ for $y = 1, 2, \dots, u''$. That is, we have the following mappings.

Parameter	Mapping Value
n_P	the number of tuples in E_j
$n_{P,s}$	the number of sensitive tuples in E_j
$n_1, n_2, \dots, n_{u''}$	$ D_{c_1}^e , D_{c_2}^e , \dots, D_{c_{u''}}^e $

In the above example, we have the following mappings.

Parameter	Mapping Value
n_P	the number of tuples in "elementary"
$n_{P,s}$	the number of sensitive tuples in "elementary"
$n_1, n_2, \dots, n_{u''}$	the number of tuples in 1st-4th (which corresponds to C_1) the number of tuples in 5th-6th (which corresponds to C_2) the number of tuples in 7th-8th (which corresponds to C_3)

Suppose we want to know the credibility for 1st-4th (which corresponds to C_1). Then, we can obtain the value by p_1 from function V .

2. Value Generalized Multiple Levels Upwards

Suppose we calculate the credibility for (o, s, K_{ad}^{min}) , where $o \in C_x$. Now, we consider the general case. For example, the original value in T may be generalized more than one level upwards. Suppose C_x is generalized to E_i in T^* . Let \mathcal{P} be the path from E_i to C_x in the generalization hierarchy. Let $\mathcal{P} = \{N_1, N_2, \dots, N_h\}$, where h is the length of the path.

For example, we have the original table as follows.

ID	T	$Sensitive$
1	1st-4th	HIV
2	1st-4th	none
3	5th-6th	HIV
4	5th-6th	HIV
5	7th-8th	none
6	7th-8th	none
7	9th-10th	HIV
8	9th-10th	HIV

After the generalization, we may obtain the following published table. We are considering this case in this section.

<i>ID</i>	<i>T</i>	<i>Sensitive</i>
1	without post-secondary	HIV
2	without post-secondary	none
3	without post-secondary	HIV
4	without post-secondary	HIV
5	without post-secondary	none
6	without post-secondary	none
7	without post-secondary	HIV
8	without post-secondary	HIV

If $C_x = \text{"1st-4th"}$, $E_i = \text{"without post-secondary"}$. Then, path $\mathcal{P} = \{\text{without post-secondary, elementary, 1st-4th}\}$.

We can perform a top-down approach for calculating the credibility of C_x . We first consider the parent-child relationship between N_1 and N_2 . Then, we calculate the credibility of node N_2 with function V . After that, we consider the parent-child relationship between N_2 and N_3 . Then, we calculate the credibility for node N_3 with function V . We continue calculating the credibility of node N_4 with function V until we finish calculating all nodes in \mathcal{P} .

For example, we first consider the parent-child relationship between “without post-secondary” and “elementary”. We calculate the credibility of an individual in “elementary” (which is the child of this relationship). After we obtain this credibility, we further consider the parent-child relationship between “elementary” and “1st-4th”. Then, we calculate the credibility of an individual in “1st-4th” (which is the child of this relationship).

Now, we map to function V . “without post-secondary” has two children, “elementary” and “secondary”. From T^* , we find that there are 8 tuples (containing 5 sensitive values) for “without post-secondary”. Thus, we get n_P and $n_{P,s}$.

For each child c_j , we obtain the number of tuples with value equal to $v(c_j)$ in the following way. Recall that T^e is the external table, and D_X^e is the set of tuples in T^e which can be generalized to $v(X)$. The set of tuples with value equal to $v(c_j)$ is equal to $D_{c_j}^e$. For example, if $c_j = \text{elementary}$, we find that the tuples with ID=1, 2, 3, 4, 5 and 6 can be generalized to “elementary”. Thus, $|D_{c_j}^e| = 6$.

After collecting the values of all parameters of V , we call function V . Then, we find the correspondence child to be processed next (i.e. $\mathcal{P}[i+1]$). Let c_j be this child. In our example, the next processing child is “elementary”. Then, the number of tuples for this child c_j is n_j and the credibility for this child c_j is p_j .

For the next step, we change c_j to the parent P in the next parent-child relationship. Also, we need to map to the function V . n_P is mapped to n_j and $n_{P,s}$ is mapped to $\lceil n_j \times p_j \rceil$. The reason of the assignment of $n_{P,s}$ by $\lceil n_j \times p_j \rceil$ is described as follows. The credibility that an individual in new parent P is linked to a sensitive value is equal to p_j . As there are n_j tuples in new parent P , the expected number of sensitive tuples in P is equal to $n_j \times p_j$. We choose the ceiling of $n_j \times p_j$ (i.e. $\lceil n_j \times p_j \rceil$), instead of a floor function, because this gives a higher privacy of an individual.

After performing all nodes in \mathcal{P} , we obtain the final credibility in p_j in the final loop, where p_j is the credibility for the child c_j in the last parent-child relationship.

4.1.2 QID of size > 1

We have described QID of size 1. Now, we describe how the formula can be extended to the QID of size greater than one.

Suppose there is an equivalence class E_i . In the previous section, we consider the children of E_i for an attribute.

Suppose QID contains q attributes. We can calculate the credibility of C_i as follows. We consider the table containing the tuples which are C_i or the generalized values of C_i . We perform a top down approach by considering the first attribute. After considering the first attribute, we continue the computation of the top down approach for the remaining attributes until all attributes in QID are processed.

QID	Disease
$q1$	HIV
$q1$	HIV
$q1$	non-sensitive
$q1$	non-sensitive
$q1$	HIV
$q2$	non-sensitive
$q2$	non-sensitive
...	...
$q2$	non-sensitive
$q2$	HIV

Table 7: Another table which violates 2-diversity

QID	Disease
$q1$	HIV
$q1$	HIV
$q1$	non-sensitive
$q1$	non-sensitive
Q	HIV
Q	non-sensitive
$q2$	non-sensitive
...	...
$q2$	non-sensitive
$q2$	HIV

Table 8: A 2-diverse table of Table 7 by local recoding

For example, we have the following table.

ID	Education	Postcode	Sensitive
1	1st-4th	4361	HIV
2	1st-4th	4362	none
3	5th-6th	4351	HIV
4	5th-6th	4352	none

The generalized table T^* is shown as follows.

ID	Education	Postcode	Sensitive
1	elementary	436*	HIV
2	elementary	436*	none
3	elementary	435*	HIV
4	elementary	435*	none

Suppose we want to calculate the credibility that an individual with (1st-4th, 4361) is linked to s . Then, we just focus the tuples with ID=1 and 2 in T^* because these two tuples are the generalized tuples with (1st-4th, 4361).

4.2 Local Recoding

An example is shown in Table 7 to illustrate the derivation of the credibility with local recoding for l -diversity. For the QID, assume that only $q1$ and $q2$ can be generalized to Q . Assume that Table 7 and the corresponding T^e have matching cardinality on Q . The proportion of the sensitive tuples in the set of tuples with $q1$ is equal to $3/5 > 1/2$. Thus, the set of tuples with $q1$ does not satisfy 2-diversity. Table 7 is generalized to Table 8, which satisfies 2-diversity, while the distortion is minimized.

Assume the adversary has knowledge of K_{ad}^{min} . From the external table T^e , there are 5 tuples with $q1$ and 8 tuples with $q2$. These are the only tuples with QID that can be generalized to Q . The adversary reasons in this way. There are four sensitive tuples in T^* . Suppose they all appear in the tuples containing $q2$, $q2$ still satisfies 2-diversity. The generalization in T^* must be caused by the set of tuples in $q1$. In T^* , the QID-EC for Q contains one sensitive tuple and one non-sensitive tuple. The sensitive tuple should

come from $q1$ because if this sensitive tuple does not come from $q1$, there will have been no need for the generalization.

Consider the credibility that an individual o with QID $q1$ is linked to HIV given K_{ad}^{min} . There are two cases.

- *Case 1:* the tuple of o appears in the QID-EC of $q1$ in T^* .

There are four tuples with value $q1$ in T^* . From T^e , there are five tuples with $q1$. The probability that Case 1 occurs is $4/5$.

- *Case 2:* the tuple of o appears in the QID-EC of Q in T^* .

There are totally five tuples with $q1$ and there are four tuples with value $q1$ in T^* . Hence, one such tuple must have been generalized and is now in the QID-EC of Q in T^* . The probability of Case 2 is $1/5$.

$Credibility(o, \{HIV\}, K_{ad}^{min})$ is equal to

$$= Prob(\text{Case 1}) \times Prob(o \text{ is linked to HIV in Case 1} \mid K_{ad}^{min}) \\ + Prob(\text{Case 2}) \times Prob(o \text{ is linked to HIV in Case 2} \mid K_{ad}^{min})$$

Since 2 out of 4 tuples in the QID-EC of $q1$ in T^* contain HIV, and the HIV tuple in the QID-EC of Q in T^* must be from $q1$, Thus,

$$Prob(o \text{ is linked to HIV in Case 1} \mid K_{ad}^{min}) = \frac{2}{4} = \frac{1}{2}. \\ Prob(o \text{ is linked to HIV in Case 2} \mid K_{ad}^{min}) = 1. \\ Credibility(o, \{HIV\}, K_{ad}^{min}) = \frac{4}{5} \times \frac{1}{2} + \frac{1}{5} \times 1 = \frac{3}{5},$$

which is greater than 0.5. Thus, the anonymized table violates 2-confidentiality.

General Formula

The above example shows the basic idea of the derivation of the general formula.

Suppose there are u QID-EC's in the original data set, namely C_1, C_2, \dots, C_u , which can be generalized to the same value \mathcal{C}_G . After the generalization, some tuples in some C_i are generalized to \mathcal{C}_G while some are not. We define the following symbols which will be used in the derivation of the credibility.

n_i	number of tuples with class C_i in T^e
$n_{i,g}$	number of generalized tuples in T^* whose original QID is C_i
$n_{i,u}$	number of ungeneralized tuples in T^* with QID = C_i
$n_{i,u(s)}$	number of sensitive ungeneralized tuples in T^* with QID = C_i

The value of $n_{i,u}$ can be easily obtained by scanning the tuples in T^* . $n_{i,g}$ can be obtained by subtracting $n_{i,u}$ from n_i . Similarly, it is easy to find $n_{i,u(s)}$. For example, in Table 8, C_i corresponds to $q1$ and \mathcal{C}_G corresponds to Q . Thus, $n_{i,u} = 4$, $n_i = 5$, $n_{i,g} = 1$ and $n_{i,u(s)} = 2$.

In order to calculate $Credibility(o, s, K_{ad}^{min})$, where o has QID of C_i , the adversary needs to consider two cases. The first case is that the tuple of o is generalized to \mathcal{C}_G . The second case is that the tuple of o is not generalized in T^* . Let $t^*(o)$ be the tuple of individual o in T^* . By considering these two cases,

$$Credibility(o, s, K_{ad}^{min}), \text{ where } o \in C_i \\ = Prob(o \text{ is linked to } s \text{ in } T^* \mid K_{ad}^{min}) \\ = Prob(t^*(o) \in \mathcal{C}_G \text{ in } T^*) \\ \times Prob(o \text{ is linked to } s \text{ in } \mathcal{C}_G \text{ in } T^* \mid K_{ad}^{min})$$

$$\begin{aligned}
& + Prob(t^*(o) \in C_i \text{ in } T^*) \\
& \times Prob(o \text{ is linked to } s \text{ in } C_i \text{ in } T^* | K_{ad}^{min}) \\
= & \frac{n_{i,g}}{n_i} \times Prob(o \text{ is linked to } s \text{ in } \mathcal{C}_G \text{ in } T^* | K_{ad}^{min}) \\
& + \frac{n_{i,u}}{n_i} \times \frac{n_{i,u}(s)}{n_{i,u}}
\end{aligned}$$

The term $Prob(o \text{ is linked to } s \text{ in } \mathcal{C}_G \text{ in } T^* | K_{ad}^{min})$ can be computed by using the formula in global-recoding, which takes into account of the minimality of the anonymization.

For the case when a set of QID-EC's are generalized to more than one values, the above analysis is extended to include more possible combinations of outcomes. However, the basic ideas remain similar.

More specifically, the question is how to compute $Prob(o \text{ is linked to } s \text{ in } \mathcal{C}_G \text{ in } T^* | K_{ad})$. We can also make use of function V . We can regard this generalized data set as follows. We just consider \mathcal{C}_G but do not consider C_i in T^* . For example, suppose we compute $Prob(o \text{ is linked to } s \text{ in elementary in } T^* | K_{ad})$. We just consider tuples with ID=3, 4, 5 and 6.

As \mathcal{C}_G has the same generalized value in the published data set for different values of i , n_P is mapped to the total number of generalized tuples in T^* (i.e. $\sum_{i=1}^u n_{i,1}$). $n_{P,s}$ is mapped to the total number of sensitive generalized tuples in T^* (i.e. $\sum_{i=1}^u n_{i,1,s}$). Then, n_1 in V is mapped to $n_{1,1}$, n_2 in V is mapped to $n_{2,1}$ and so on. Thus, we have the following mapping.

Parameter	Mapping Value
n_P	$\sum_{i=1}^u n_{i,1}$
$n_{P,s}$	$\sum_{i=1}^u n_{i,1,s}$
n_1, n_2, \dots, n_u	$n_{1,1}, n_{2,1}, \dots, n_{u,1}$

Then, we can perform the function V to obtain the credibility.

4.2.1 Generalization Hierarchy of Height > 1

Now, we consider the general case with the generalization hierarchy of height greater than one.

Let T be the original data set. Let T^* be the generalized data set. Let T^e be the external data set. Suppose we consider an individual $o \in C_i$. Note that C_i may be generalized to different values in T^* . Let us illustrate with the following example.

<i>ID</i>	<i>T</i>	<i>Sensitive</i>
1	1st-4th	HIV
2	1st-4th	HIV
3	5th-6th	none
4	5th-6th	none
5	5th-6th	none
6	5th-6th	none
7	7th-8th	HIV
8	7th-8th	none
9	9th-10th	HIV
10	9th-10th	HIV

We may generalize the table as follows.

<i>ID</i>	<i>T*</i>	<i>Sensitive</i>
1	elementary	HIV
2	elementary	HIV
3	elementary	none
4	elementary	none
5	without post-secondary	none
6	without post-secondary	none
7	7th-8th	HIV
8	7th-8th	none
9	without post-secondary	HIV
10	without post-secondary	HIV

Suppose the adversary wants to calculate the credibility that an individual in 1st-4th is linked to a sensitive value s in the published data set T^* . He/she may guess “elementary” or “without post-secondary” in the published data set. In other words, the tuples with “elementary” and “without post-secondary” in T^* are suspected to be the tuples with “1st-4th”.

In general, let C_i be the class to be considered. In T^* , some tuples t with C_i may be generalized while some tuples T^e may not be generalized. Besides, the generalized tuples t may be generalized to different values. For example, some tuples in “5th-6th” are generalized to “elementary” and some tuples are generalized “without post-secondary”. Let \mathcal{W}_{C_i} be the set of values in T^* to which C_i can be generalized. For example, if $C_i = \text{“5th-6th”}$, $\mathcal{W}_{C_i} = \{\text{elementary, without post-secondary}\}$.

We have the following derivation.

$$\begin{aligned}
& \text{Credibility for } (o, s, K_{ad}), \text{ where } o \in C_i \\
&= \text{Prob}(o \text{ is linked to } s \text{ in } T^* | K_{ad}) \\
&= \text{Prob}(o \in \mathcal{W}_{C_i} \text{ in } T^*) \times \text{Prob}(o \text{ is linked to } s \text{ in } \mathcal{W}_{C_i} \text{ in } T^* | K_{ad}) \\
&\quad + \text{Prob}(o \in C_i \text{ in } T^*) \times \text{Prob}(o \text{ is linked to } s \text{ in } C_i \text{ in } T^* | K_{ad}) \\
&= \text{Prob}(o \in \mathcal{W}_{C_i} \text{ in } T^*) \times \\
&\quad \sum_{W \in \mathcal{W}_{C_i}} [\text{Prob}(o \in W \text{ in } \mathcal{W}_{C_i} \text{ in } T^*) \times \text{Prob}(o \text{ is linked to } s \text{ in } W \text{ in } T^* | K_{ad})] \\
&\quad + \text{Prob}(o \in C_i \text{ in } T^*) \times \text{Prob}(o \text{ is linked to } s \text{ in } C_i \text{ in } T^* | K_{ad})
\end{aligned}$$

In the example, suppose $C_i = 5th - 6th$. As there is no 5th-6th in T^* , $\text{Prob}(o \in C_i \text{ in } T^*) = 0$ and $\text{Prob}(o \in \mathcal{W}_{C_i} \text{ in } T^*) = 1$. We know that $\mathcal{W}_{C_i} = \{\text{elementary, without post-secondary}\}$. If $W = \text{elementary}$, then $\text{Prob}(o \in W \text{ in } \mathcal{W}_{C_i} \text{ in } T^*) = 4/8=1/2$. If $W = \text{“without post-secondary”}$, then $\text{Prob}(o \in W \text{ in } \mathcal{W}_{C_i} \text{ in } T^*) = 4/8=1/2$. However, we need to compute $\text{Prob}(o \text{ is linked to } s \text{ in } W \text{ in } T^* | K_{ad})$. We can make use of the function V . The step to find the credibility is also similar. But, it involves some complicated issues.

Suppose $W \in \mathcal{W}_{C_i}$. In function V , we map P to W . Let \mathcal{P} be the path from P to C_i .

There is a parent node P and a number of children of node P in the generalization hierarchy. In order to use the function V , we need to know the number of tuples of the parent node P and the number of tuples of the children of P .

There are the following cases. Suppose $W \in \mathcal{W}_{C_i}$. If \mathcal{W}_{C_i} contains more than one value, some values in \mathcal{W}_{C_i} is generalized than others. For example, if $C_i = 5th-6th$, “without post-secondary” is more generalized than “elementary” in \mathcal{W}_{C_i} . We have the following cases. After considering the following cases, we can map P and the correspondence children in function V .

1. W is the most generalized value in \mathcal{W}_{C_i}
2. W is not the most generalized value in \mathcal{W}_{C_i}
1. **Case 1:** W is the most generalized value in \mathcal{W}_{C_i} .

Suppose P (or W) equals “without post-secondary” and one of the children of P equals “elementary”.

- (a) We consider how the adversary obtains n_P and $n_{P,s}$. It is easy for the adversary to obtain this information from T^* . For example, if P = “without post-secondary”, we can obtain $n_P = 4$ and $n_{P,s} = 2$ by scanning the table T^* .
- (b) Now, we consider how the adversary calculates the number of tuples of each child of P (or W). Let us illustrate with an example. The adversary wants to know how many tuples with “without post-secondary” in T^* come from “elementary” from his knowledge. His knowledge contains the statistics of the education levels and the published table T^* . Let us denote the statistics of the education levels by T^e for simplicity.

It is easy to know that we can derive the number of tuples with “without post-secondary” in T^* which come from “elementary” is 2. In T^e , tuples with ID=1, 2, 3, 4, 5, 6, 7 and 8 belong to elementary. In T^* , tuples with ID=1, 2, 3, 4, 7 and 8 can be generalized to “elementary”. Let D^e be the set of tuples which can be generalized to “elementary” in T^e . Let D^* be the set of tuples which are equal to or can be generalized to “elementary” in T^* . $D^e = \{1, 2, 3, 4, 5, 6, 7, 8\}$ and $D^* = \{1, 2, 3, 4, 7, 8\}$. Then, we know that the number of tuples with “without post-secondary” in T^* which come from “elementary” is 2 by subtracting $|D^*|$ from $|D^e|$.

In general, we define the following notations. Let X be a generalized/ungeneralized value in D^e or T^* . Let $D_X^* = \{t | t \text{ is equal to or can be generalized to } X \text{ in } T^*\}$. Let $D_X^e = \{t | t \text{ is equal to or can be generalized to } X \text{ in } T^e\}$.

D_X^*	$\{t t \text{ is equal to or can be generalized to } X \text{ in } T^*\}$
D_X^e	$\{t t \text{ is equal to or can be generalized to } X \text{ in } T^e\}$

If we set $X = \text{elementary}$, we can obtain the number of tuples with “without post-secondary” in T^* which come from X is 2 by subtracting $|D_X^*|$ from $|D_X^e|$.

In summary, for each child c_j of P (or W), we can also obtain the number of tuples by subtracting $|D_{c_j}^*|$ from $|D_{c_j}^e|$. After obtaining this information, we can also make use of function V to derive the credibility of an individual in C_i .

2. **Case 2:** W is not the most generalized value in \mathcal{W}_{C_i} .

Suppose P equals “elementary” in T^* and one of the children of P equals “1st-4th”.

- (a) Now, we first consider how the adversary obtains n_P and $n_{P,s}$ in function V . Similarly, we can scan T^* to obtain n_P and $n_{P,s}$. For example, if P = “elementary”, $n_P = 4$ and $n_{P,s} = 2$. Note that we are considering the tuples with “elementary” in T^* . Thus, we are considering the tuples with ID=1, 2, 3 and 4 because we are calculating the probability that the tuples with “elementary” in T^* are linked to s . Thus, we do not need to consider the tuples with ID=5 and 6 where “without post-secondary” is a generalized value of “elementary”.
- (b) We consider how the adversary calculates the number of tuples in each child of P (or W). Also, we illustrate with an example. The adversary wants to know how many tuples with “elementary” in T^* come from “1st-4th” from his knowledge. His knowledge contains the statistics of the education levels and the published table T^* .

It is not trivial for the adversary to deduce the number of tuples with “elementary” in T^* come from “1st-4th”. The tuples with “elementary” in T^* have ID=1, 2, 3 and 4. However, 1st-4th may appear in “without post-secondary” and “elementary”.

We know that there are four tuples with “elementary” in T^* (i.e. ID=1, 2, 3 and 4). Thus, the adversary may think in one of these ways.

- i. “1st-4th” does not appear in tuples 1, 2, 3 and 4
- ii. “1st-4th” appears in one of these four tuples
- iii. “1st-4th” appears in two of these four tuples

For each of the above cases, we obtain the number of tuples with 1st-4th. Similarly, we obtain the number of tuples with 5th-6th and 7th-8th such that the total number of tuples for “elementary” education levels, 1st-4th, 5th-6th and 7th-8th, is equal to 4. In other words, the adversary considers all possible combinations of these four tuples into these three education levels.

For each combination, the adversary makes use of the function V to calculate the credibility. Finally, we obtain the average credibility among all combinations.

4.2.2 QID of size > 1

The idea is similar to Section 4.1.2.

4.3 Attack Conditions

We have seen in the above that a minimality attack is always accompanied by some exclusion of some possibilities by the adversary because of the minimality notion. We can characterize this attack criterion in the following.

Theorem 1 *An attack by minimality is possible only if the adversary can exclude some possible combinations of the number of sensitive tuples among the QID-EC's in the sensitive tuple distribution table based on the knowledge of K_{ad}^{min} .*

Proof Sketch: The above theorem is true because, if there is no exclusion from the table, then the credibility as computed by the formulae for credibility is exactly the ratio of the sensitive tuples to the total number of tuples in the generalized QID-EC. \square

Lemma 2 *An attack by minimality is not always successful even when there are some excluded combination(s) in the sensitive tuple distribution table based on K_{ad}^{min} .*

To prove this lemma, we give an example where 2 QID's $q1$ and $q2$ are generalized to Q . There are 4 tuples of $q1$ and 2 tuples of $q2$. In total, there are 3 occurrences of the sensitive value set s in the 6 tuples. If 2-diversity is the goal, then we can exclude the case of 2 sensitive $q1$ tuple and 1 sensitive $q2$ tuple. After the exclusion, the credibility of any linkage between any individual to s still does not exceed 0.5.

5 General Model

In this section, we show that minimality attack can be successful on a variety of anonymization models. In Tables 9 to 13, we show *good tables* that satisfy the corresponding privacy requirements in different models, *bad tables* that do not, and *global* and *local* recodings of the bad tables which follow the minimality principle and unfortunately suffer from minimality attacks.

5.1 Recursive (c, l) -diversity

With recursive (c, l) -diversity [14], in each QID-EC, let v be the most frequent sensitive value, if we remove the next $l - 2$ most frequent sensitive values, the frequency of v must be less than c times the total count of the remaining values. Table 9(c) is a global recoding for Table 9(b). With the knowledge of minimality in the anonymization, the adversary deduces that the QID-EC for $q2$ must satisfy (3,3)-diversity and that the QID-EC for $q1$ must contain two HIV values. Thus, the intended obligation that an individual should be linked to at least 3 different sensitive values is breached. Similar arguments can be applied to Table (d).

5.2 t -closeness

Recently, t -closeness [13] was proposed. If table T satisfies t -closeness, the distribution \mathbb{P} of each equivalence class in T is roughly equal to the distribution \mathbb{Q} of the whole table T with respect to the sensitive attribute. More specifically, the difference between the distribution of each equivalence class in T and the distribution of the whole table T , denoted by $D[\mathbb{P}, \mathbb{Q}]$, is at most t . Let us use the definition in [13]:

$$D[\mathbb{P}, \mathbb{Q}] = 1/2 \sum_{i=1}^m |p_i - q_i|$$

Consider Table 10(c). For each possible sensitive value distribution \mathbb{P} for QID-EC $q2$, the adversary computes $D[\mathbb{P}, \mathbb{Q}]$. S/he finds that $D[\mathbb{P}, \mathbb{Q}]$ is always smaller than 0.2. Hence the anonymization is due to

QID	Disease
$q1$	Diabetics
$q1$	HIV
$q1$	Lung Cancer
$q2$	HIV
$q2$	Ulcer
$q2$	Alzhema
$q2$	Gallstones

(a) good table

QID	Disease
$q1$	Diabetics
$q1$	HIV
$q1$	HIV
$q2$	Lung Cancer
$q2$	Ulcer
$q2$	Alzhema
$q2$	Gallstones

(b) bad table

QID
Q
Q
Q
Q
Q
Q
Q

(c) global

QID
Q
Q
Q
Q
$q2$
$q2$
$q2$

(d) local

Table 9: Anonymization for (3,3)-diversity

QID	Disease
$q1$	HIV
$q1$	non-sensitive
$q2$	non-sensitive
$q2$	non-sensitive
$q2$	HIV
$q2$	HIV

(a) good table

QID	Disease
$q1$	HIV
$q1$	HIV
$q2$	non-sensitive
$q2$	non-sensitive
$q2$	non-sensitive
$q2$	HIV

(b) bad table

QID
Q
Q
Q
Q
Q
Q

(c) global

QID
Q
Q
Q
$q2$
$q2$
$q2$

(d) local

Table 10: 0.2-closeness anonymization

$q1$. S/he concludes that both tuples with QID= $q1$ are sensitive. Similar arguments can also be made to Table (d).

5.3 (k, e) -anonymity

The model of (k, e) -anonymity [27] considers the anonymization of tables with numeric sensitive attributes. It generates a table where each equivalence class is of size at least k and has a range of the sensitive values at least e . In the tables in Table 11 and Table 12, we show the bucketization in terms of QID values, the individuals with the same QID value are in the same bucket. Consider the tables in Table 11 (where Income is a sensitive numeric attribute). From Table (c), the adversary deduces that the tuples with QID= $q1$ must violate (k, e) -anonymity and must be linked with two $30k$ incomes. We obtain a similar conclusion from Table (d) for local recoding.

QID	Income
$q1$	30k
$q1$	20k
$q2$	30k
$q2$	20k
$q2$	40k

(a) good table

QID	Income
$q1$	30k
$q1$	30k
$q2$	20k
$q2$	10k
$q2$	40k

(b) bad table

QID
Q
Q
Q
Q
Q

(c) global

QID
Q
Q
Q
$q2$
$q2$

(d) local

Table 11: (k, e) -anonymity for $k = 2$ and $e = 5k$

QID	Disease	QID	Disease	QID	QID
$q1$	HIV	$q1$	HIV	Q	Q
$q1$	non-sensitive	$q1$	HIV	Q	Q
$q1$	non-sensitive	$q1$	non-sensitive	Q	Q
$q2$	HIV	$q2$	non-sensitive	Q	Q
$q2$	non-sensitive	$q2$	non-sensitive	Q	Q
$q2$	non-sensitive	$q2$	non-sensitive	Q	Q
$q2$	non-sensitive	$q2$	non-sensitive	Q	Q
$q2$	non-sensitive	$q2$	non-sensitive	Q	Q
$q2$	non-sensitive	$q2$	non-sensitive	Q	Q
$q2$	non-sensitive	$q2$	non-sensitive	Q	$q2$
$q2$	non-sensitive	$q2$	non-sensitive	Q	$q2$
$q2$	non-sensitive	$q2$	non-sensitive	Q	$q2$

(a) good table (b) bad table (c) global (d) local

Table 12: Anonymization for (0.6, 2)-safety

QID	Education	Guarding Node	QID	QID
$q1$	1st-4th	elementary	Q	Q
$q2$	undergrad	none	Q	Q
$q2$	undergrad	none	Q	$q2$

(a) bad table (b) global (c) local

Table 13: Anonymization for personalized anonymity

5.4 (c, k) -safety

(c, k) -safety [15] considers the worst-case background knowledge. The background knowledge defined in [15] is a set of k implications whose simplest form is: if individual p_1 is linked to sensitive value v_1 , then individual p_2 is linked to sensitive value v_2 , where v_1 may equal v_2 . This model requires that the probability that any individual is linked to a sensitive value given a background knowledge containing any k implications is at most c . Consider the tables in Table 12. Table (c) is a global recoding of Table (b). From Table (c), the adversary can eliminate the cases where there is 0 or 1 HIV value in the bucket $q1$, since in such cases Table (b) would have been (0.6,2)-safe. Since Table (b) is anonymized to Table (c), there must be 2 HIV occurrences in the bucket for $q1$ in Table (b). Hence an individual with $q1$ is linked to HIV with a probability of $2/3$, higher than the intended threshold of 0.6. Similar attack can be launched against the local recoding in Table (d).

5.5 Personalized Privacy

[26] proposed a personalized privacy model where each individual can provide his/her preference on the protection of his/her sensitive value, denoted by a *guarding node*. For example, an individual o with a value “1st-4th” may specify “elementary” as a guarding node in order that any QID-EC that may contain o should contain at most $1/l$ tuples with “elementary” values. For $l = 2$, in the tables in Table 13, Tables (b) and (c) are global and local recodings for Table (a), respectively. Suppose the adversary knows that everyone with undergraduate degree do not mind to disclose his/her education. Based on the minimality principle, if the “1st-4th” belongs to a $q2$ tuple, then Table (a) will not be anonymized, so the tuple with QID= $q1$ must be linked to “1st-4th”. Similar attack will be successful on Table (c).

5.6 General Attack by Minimality

In the proposed anonymization mechanism for each of the above cases in the respective references, the Minimality Principle in Definition 3 holds if we set \mathcal{R} to the objective at hand, such as recursive (c, l) -diversity, t -closeness and (k, e) -anonymity. By including the knowledge related to minimality attack to the background knowledge, the adversary can derive the probabilistic formulae for computing the corresponding credibility in each case, where the idea of eliminating impossible cases as shown in Section 4 is a key to the attack.

6 Algorithm

The problem of optimal m -confidentiality is a difficult problem. In most data anonymization methods, if a generalization step does not reach the privacy goal, further generalization can help. However, further generalizations will not solve the problem of m -confidentiality. If we further generalize Q to $*$ in Table 1(c) or further generalize $q2$ to Q in Table 1(d), it does not deter the minimality attack. The result still reveals the linkage of $q1$ to HIV as before. We show below optimal m -confidentiality is NP-hard for global recoding.

Optimal global m -confidentiality: Given a table T and a non-negative cost e , can we generate a table T^* from T by global recoding which satisfies m -confidentiality and where the information loss of $Dist(T, T^*)$ is less than or equal to e ?

Theorem 2 *Optimal m -confidentiality under global recoding is NP-hard.*

Proof: We shall transform the problem of EXACT COVER BY 3-SETS (X3C) [8] to the m -confidentiality anonymization problem. X3C is defined by: Given a set X with $|X| = 3q$ and a collection C of 3-element subsets of X . Does C contain an exact cover for X , i.e. a subcollection $C' \subseteq C$ such that every element of X occurs in exactly one member of C' ?

Given an instance of X3C, we transform it to an instance of Optimal m -confidentiality under global recoding as follows. Create a table T with two attributes Q and S , where Q is a QID attribute and S is a sensitive attribute that may contain sensitive values. For S , there is only one sensitive value s_v and one non-sensitive value s_n . We set $weight(Q) = 1$. For each element x in X , create a tuple with $Q = x$ and $S = s_v$. Hence, each value of x appears in exactly one tuple. Let the elements in C be c_1, \dots, c_m . For each element $c_i = (x, y, z)$ in C , create a taxonomy \mathcal{T}_i . \mathcal{T}_i contains ground elements of x, y, z, n_{i1}, n_{i2} and n_{i3} , which are children of a root node r_i . Create 3 tuples with $Q = n_{ij}$ and $S = s_n$, for $j = 1, 2, 3$.

Recall that with global recoding all occurrences of an attribute value are recoded to the same value. Note that the adversary cannot launch a minimality attack since each QID value appears only in one tuple in the set of tuples. The adversary cannot exclude any possible combination of the table of sensitive tuple distribution. From Theorem 1, minimality attack is not possible. Hence C contains an exact cover for X if and only if there is a solution T^* for the 2-confidentiality problem with $Dist(T, T^*) = 6q$. \square

However, the following lemma gives us a good pointer to a possible solution.

Lemma 3 *Given the adversary background knowledge K_{ad}^{min} , if the anonymization mechanism that generates T^* does not follow the minimality principle, then T^* satisfies m -confidentiality iff it satisfies l -diversity, where $l = m$.*

The above lemma holds because in this case K_{ad}^{min} does not give any extra information for the adversary compared to the problem setting for l -diversity.

From the above lemma the minimality principle and l -diversity are the crucial keys. As the adversary relies on the minimality assumption, we can tackle the problem at its source by removing the minimality notion from the anonymization. The main idea is that, even if some QID-EC's in a given table T originally do not violate l -diversity, we can still generalize the QID. Since the anonymization does not play according to the minimality rule, the adversary cannot launch the minimality attack directly. However, a question is: how much shall we generalize or anonymize? It is not desirable to lose on data utility.

A naive method to generalize everything in an excessive manner would not work well, since the information loss will also be excessively large. From the formula for information loss, if every QID

attribute value must go at least one level up the taxonomies, then for typical taxonomies, the information loss will be a sizeable fraction.

Here we propose a feasible solution for the m -confidentiality problem. Although some problems are uncovered that questions the utility of k -anonymity in protecting sensitive values, k -anonymity has been successful in some practical applications. This indicates that when a data set is k -anonymized for a given k , the chance of a large proportion of a sensitive value set s in any QID-EC is very likely reduced to a safe level. Since k -anonymity does not try to anonymize based on the sensitive value set, it will anonymize a QID-EC even if it satisfies l -diversity. This is the blinding effect we are targeting for. However, there is no guarantee of m -confidentiality by k -anonymity alone, where $m = l$.

Hence, our solution is based on k -anonymity, with additional precaution steps taken to ensure m -confidentiality. Let us call our solution Algorithm MASK (Minimality Attack Safe K-anonymity), which involves four steps.

Algorithm 1 – MASK

- 1: From the given table T , generate a k -anonymous table T^k where k is a user parameter.
 - 2: From T^k , determine the set \mathcal{V} containing all QID-EC's which violate l -diversity in T^k , and a set \mathcal{L} containing QID-EC's which satisfy l -diversity in T^k . How to select \mathcal{L} will be described below.
 - 3: For each QID-EC Q_i in \mathcal{L} , find the proportion p_i of tuples containing values in the sensitive value set s . The distribution \mathcal{D} of the p_i values is determined.
 - 4: For each QID-EC $E \in \mathcal{V}$, randomly pick a value of p_E from the distribution \mathcal{D} . The sensitive values in E are distorted in such a way that the resulting proportion of the sensitive value set s in E is equal to p_E .
-

After Step 1, some QID-EC's may not satisfy l -diversity. Steps 2 to 4 above will ensure that all QID-EC's in the result are l -diverse. In Step 2 above, we select a QID-EC set \mathcal{L} from T^k . The purpose is to disguise the distortion so that the adversary cannot tell the difference between a distorted QID-EC and many un-distorted QID-EC's. We set the size of \mathcal{L} , denoted by u , to $(l - 1) \times |\mathcal{V}|$. Among all the QID-EC's in T^k that satisfies l -diversity, we pick u QID-EC's with the highest proportions of the sensitive value set s .

Theorem 3 *Algorithm MASK generates m -confidential data sets.*

The above follows from Lemma 3 because MASK does not follow the minimality principle. It is easy to find an l -diverse table T^* generated by MASK with a QID-EC X in T^* so that a specialization of the QID's in X results in another table T' which also satisfies l -diversity.

The use of \mathcal{L} for the distortion of \mathcal{V} is to make the distribution of s proportions in \mathcal{V} look indistinguishable from that of a large QID-EC set (\mathcal{L}). This is an extra safeguard for the algorithm in case the adversary knows the mechanism of anonymization. If the QID-EC's in \mathcal{V} simply copy the s proportion from an l -diverse QID-EC in T_k with the greatest s proportion, the repeated pattern may become a source of attack. In our setting the probability that some QID-EC in \mathcal{V} has the same s proportion as a QID-EC in \mathcal{L} is $1/l$. Therefore, for l repeated occurrences of an s proportion, the probability that any one belongs to a QID-EC in \mathcal{V} is only $1/l (= 1/m)$.

Generation of Two Tables - Bucketization

Conventional anonymization methods produce a single generalized table T as shown in Table 5. Recently [25] proposed to generate two separate tables from T with the introduction of an attribute called GID that is shared by the two tables. The first table T_{QID} contains the attributes of QID and GID, and the second table T_{sen} contains GID and the sensitive attribute(s). The two tables are created from T^* by assigning each QID-EC in T^* a unique GID. The advantage is that we can keep the original values in T of the QID in T_{QID} and hence reduce information loss. However, the single table T has the advantage of

	Attribute	Distinct Values	Generalizations	Height
1	Age	74	5-, 10-, 20-year ranges	4
2	Work Class	7	Taxonomy Tree	3
3	Marital Status	7	Taxonomy Tree	3
4	Occupation	14	Taxonomy Tree	2
5	Race	5	Taxonomy Tree	2
6	Sex	2	Suppression	1
7	Native Country	41	Taxonomy Tree	3
8	Salary Class	2	Suppression	1
9	Education	16	Taxonomy Tree	4

Table 14: Description of Adult Data Set

clarity and requiring no extra interpretation on the data receiver’s part. In our experiments, we will try both the approach of generating a single table T and the approach of generating two tables (also known as bucketization) as in [25, 27, 15].

7 Empirical Study

A Pentium IV 2.2GHz PC with 1GM RAM was used to conduct our experiment. The algorithm was implemented in C/C++. In our experiment, we adopted the publicly available data set, Adult Database from the UCIrvine Machine Learning Repository [3]. This data set (5.5MB) was also adopted by [12, 14, 23, 7]. We used a configuration similar to [12, 14]. The records with unknown values were first eliminated resulting in a data set with 45,222 tuples (5.4MB). Nine attributes were chosen in our experiment, as shown in Table 14. By default, we chose the first eight attributes and the last attribute in Table 14 as the quasi-identifier and the sensitive attribute, respectively. As discussed in the previous sections, attribute “Education” contains a sensitive value set containing all values representing the education levels before “secondary” (or “9th-10th”) such as “1st-4th”, “5th-6th” and “7th-8th”.

7.1 Analysis of the minimality attack

We are interested to know how successful the minimality attack can be in a real data set with existing minimality-based anonymization algorithms. We adopted the Adult data set and the selected algorithm was the (α, k) -anonymity algorithm [24]. We set $\alpha = 1/l$ and $k = 1$, so that it corresponds to the simplified l -diversity. We have implemented an algorithm based on the general formulae in Section 4 to compute the credibility values. We found that minimality attack successfully uncovered QID-EC’s which violates m -confidentiality, where $m = l$. We use m and l exchangeably in the following. Let us call the tuples in such QID-EC’s the *problematic tuples*. Figure 8(a) shows the proportion of problematic tuples among all sensitive tuples under the variation of m , where the total number of sensitive tuples is 1,566. The general trend is that the proportion increases when m increases. When m increases, there is higher chance that problematic tuples are generalized with more generalized tuples. Also, it is more likely that those generalized tuples are easily uncovered for the minimality attack.

In Figure 8(b), when m increases, it is obvious that the average credibility of problematic tuples decreases. When m increases, $1/m$ decreases. Thus, each QID-EC contains at most $1/m$ occurrences of the sensitive value set. Thus, this lowers the credibility of the tuples in QID-ECs.

Figure 8(c) shows that the proportion of problematic tuples increases with QID size. This is because, when QID size is larger, the size of each QID-EC is smaller. It is more likely that a QID-EC violates the privacy requirement. Thus, more tuples are vulnerable for the minimality attack. Figure 8(d) shows that the average credibility of problematic tuples remain nearly unchanged when the QID size increases. This is because the credibility is based on m . It is noted that the average credibility in Figure 8(d) is about 0.9, which is greater than 0.5 ($=1/2$).

We also examined some cases obtained in the experiment. Suppose we adopt the QID attributes as (age, workclass, martial status) with sensitive attribute Education. The original table contains one tuple with QID=(80, self-emp-not-inc, married-spouse-absent) and two tuples with QID=(80, private, married-spouse-absent).

Age	Workclass	Martial Status	Education
80	self-emp-not-inc	married-spouse-absent	7th-8th
80	private	married-spouse-absent	HS-grad
80	private	married-spouse-absent	HS-grad

Suppose $m = 2$. Recall that 7th-8th is in the sensitive value set. Since the first tuple violates 2-diversity, the Workclass of tuple 1 and tuple 2 are generalized to “with-pay” as follows. Thus, the published table contains the following tuples.

Age	Workclass	Martial Status	Education
80	with-pay	married-spouse-absent	7th-8th
80	with-pay	married-spouse-absent	HS-grad
80	private	married-spouse-absent	HS-grad

In this case, it is easy to check that the credibility for an individual with QID= (80, self-emp-not-inc, married-spouse-absent) is equal to 1.

Another uncovered case involves more tuples. The original table contains one tuple with QID=(33, self-emp-not-inc, married-spouse-absent) and 17 tuples with QID=(33, private, married-spouse-absent).

Similarly, when $m = 2$, the first tuple violates 2-diversity. Thus, Workclass of tuple 1 and tuple 2 are generalized to “with-pay” in the published table. Similarly, the adversary can deduce that the individual with QID=(33, self-emp-not-inc, married-spouse-absent) is linked with a low education (i.e. Education=“1st-4th”) since this credibility is equal to 1.

Consider the default QID size = 8. When $m = 2$, the execution time of the computation of the credibility of each QID-ECs in the original table is about 173s. When $m = 10$, the execution time is 239s. It is not costly for an adversary to launch a minimality attack.

7.2 Analysis of the proposed algorithm

We compared our proposed algorithm with a local recoding algorithm for (α, k) -anonymity [24] ((α, k) -A). Let us refer to our proposed algorithm MASK described in Section 6 by m -conf. (α, k) -A does not guarantee m -confidentiality, but it is suitable for comparison since it considers both k -anonymity and l -diversity, where $l = m$. We are therefore interested to know the overhead required in our approach in order to achieve m -confidentiality. When we compared our algorithm with (α, k) -anonymity, we set $\alpha = 1/m$ and the k value is the same as that use in our algorithm. We evaluated the algorithms in terms of four measurements: *execution time*, *relative error ratio*, *information loss* of QID attributes and *distortion* of sensitive attribute. The distortion of sensitive attribute is calculated by the information loss formula in Definition 6. We give it a different name for the ease of reference. By default, the weighting of each attribute used in the evaluation of information loss is equal to $1/|QID|$, where $|QID|$ is the QID size. For each measurement, we conducted the experiments 100 times and took the average.

We have implemented two different versions of Algorithm MARK: (A) one generalized table is generated and (B) two tables are generated (see the last paragraph in Section 6). For Case (A), we may generalize the QID attributes of the data and distort the sensitive attribute of the data. Thus, we measured these by information loss and distortion, respectively. For Case (B), since the resulting tables do not generalize QID, there is no information loss for QID. The distortion of the sensitive attribute is the same as in Case (A). Hence in the evaluation of information loss and distortion, we only report the results for Case (A).

For case (B) with the generation of two ungeneralized tables, T_{QID} and T_{sen} , as in [25], we measure the error by the *relative error ratio* in answering a aggregate query. We adopt both the form of the aggregate query and the parameters of the *query dimensionality* qd and the *expected query selectivity* s from [25]. For each evaluation in the case of two anonymized tables, we performed 10,000 queries and then reported the average relative error ratio. By default, we set $s = 0.05$ and qd to be the QID size.

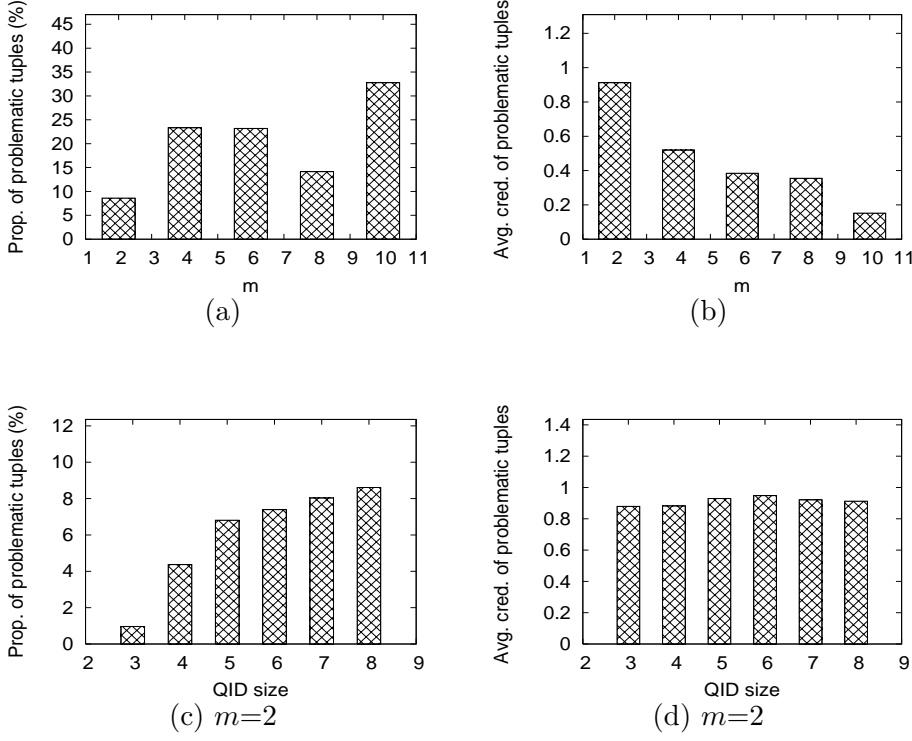


Figure 8: Proportion of problematic tuples and average credibility of problematic tuples against m and QID size

We conducted the experiments by varying the following factors: (1) the QID size, (2) m , (3) k , (4) query dimensionality qd (in the case of two anonymized tables), and (5) selectivity s (in the case of two anonymized tables).

7.2.1 The single table approach

The results for the single table case are shown in Figure 9 and Figure 10. One important observation is that the results are little affected by the values of k which varies from 2 to 10 to 20, this is true for the execution time, the relative error, the information loss and for the distortion. This is important since k is a user parameter and the results indicate that the performance is robust against different choices of the value of k .

A second interesting observation is that the information loss of (α, k) -A is greater than m -conf in some cases. This seems surprising since m -conf has to fend off minimality attack while (α, k) -A does not. The explanation is that in some cases, more generalization is required in (α, k) -A to satisfy l -diversity. However, the first step of m -conf only considers k -anonymity and not l -diversity. Thus, the generalization in m -conf is less compared to (α, k) -A, leading to less information loss. For compensation, the last two steps of m -conf ensure l -diversity and incur distortion, while (α, k) -A has no such steps.

The execution times of the two algorithms are similar because the first step of m -conf occupies over 98% of the execution time on average and the first step is similar to (α, k) -A.

In Figure 9(a), the execution time increases with the QID size, since greater QID size results in more QID-EC's. When k is larger, the execution time is smaller, this is because the number of QID-EC's will be smaller.

Figures 9(b) and (d) show that the average relative error and the distortion of the algorithms increase with the QID size. This is because the number of QID-EC's increases and the average size of each equivalence class decreases. For m -conf, the probability that a QID-EC violates l -diversity (after the k -anonymization step) will be higher. Thus, there is a higher chance for the distortion and higher average relative error. When k is larger, the average relative error of the two algorithms increases. This is because

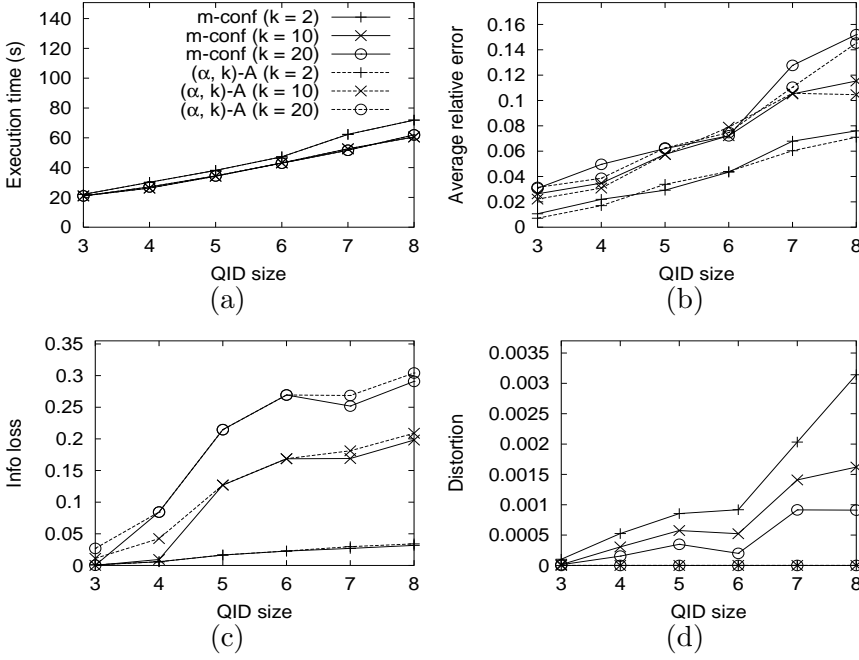


Figure 9: Performance vs QID size ($m = 2$)

the QID attribute will be generalized more, giving rise to more querying errors. If k is larger, the QID-EC size increases, the chance that a QID-EC violates l -diversity is smaller, so the distortion will be less.

In Figure 9(c), when the QID size increases, the information loss of the QID attributes increases since the probability that the tuples in the original table have different QID values is larger. Thus, there is a higher chance for QID generalization leading to more information loss. Similarly, when k is larger, the information loss is larger.

7.2.2 The two tables approach

Our next set of experiments analyze the performance of the two table approach under various conditions.

Effect of k : Figure 11 shows the experimental results when k is varied. The trends are similar to the single table case, and can be explained similarly.

Effect of Query Dimensionality qd : For $m = 2$, Figure 12(a) shows the average relative error increases when the query dimensionality increases. As the query will match fewer tuples, fewer tuples in an equivalence class will match the query, resulting in more relative error. If k is larger, the average relative error is larger because we generalize more data with larger k . Similar trends can also be observed when $m = 10$.

Effect of Selectivity s : In Figure 12(c), the average relative error decreases when s increases. This is because, if s is larger, more tuples will be matched with a given query, and more tuples in an equivalence class is matched with a given query. Similarly, when k is larger, there is more generalization, and the average relative error is larger. We observe similar trends when $m=10$. Similarly, the average relative error is larger when $m=10$.

In conclusion, we find that our algorithm creates very little overhead and pays a very minimal price in information loss in the exchange for m -confidentiality.

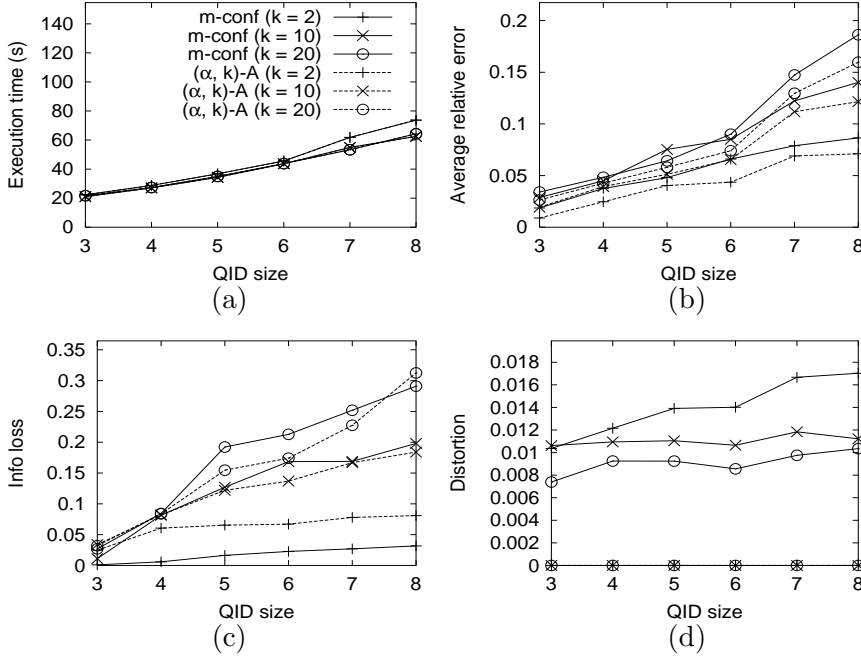


Figure 10: Performance vs QID size ($m = 10$)

8 Conclusion

In existing privacy preservation methods for data publishing, minimality in information loss is an underlying principle. In this paper, we show how this can be used by an adversary to launch an attack on the published data. We call this a minimality attack. We propose the m -confidentiality model which deals with attack by minimality. We also propose an algorithm which generates an m -confidential data set. We conducted experiments to show that our proposed algorithm requires little overhead both in terms of execution time and information loss. For future work we are interested to determine any other kinds of attack that can be related to the nature of the anonymization process.

References

- [1] G. Aggarwal, T. Feder, K. Kenthapadi, R. Motwani, R. Panigrahy, D. Thomas, and A. Zhu. Anonymizing tables. In *ICDT*, pages 246–258, 2005.
- [2] R. Agrawal and R. Srikant. Privacy-preserving data mining. In *Proc. of the ACM SIGMOD Conference on Management of Data*, pages 439–450. ACM Press, May 2000.
- [3] E. Keogh, C. Blake, and C. J. Merz. UCI repository of machine learning databases, <http://www.ics.uci.edu/~mllearn/MLRepository.html>, 1998.
- [4] D. Brumley and D. Boneh. Remote timing attacks are practical. In *USENIX Security Symposium*, 2003.
- [5] Alexandre Evfimievski, Ramakrishnan Srikant, and Johannes Gehrke Rakesh Agrawal. Privacy preserving mining of association rules. In *KDD*, 2002.
- [6] U. M. Fayyad and K. B. Irani. Multi-interval discretization of continuous-valued attributes for classification learning. In *Proceedings of the Thirteenth International Joint Conference on Artificial Intelligence (IJCAI-93)*, pages 1022–1027, San Francisco, 1993. Morgan Kaufmann.
- [7] B. C. M. Fung, K. Wang, and P. S. Yu. Top-down specialization for information and privacy preservation. In *ICDE*, pages 205–216, 2005.
- [8] M.R. Garey and D.S. Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. Freeman, USA, 1979.
- [9] D. Kifer and J. Gehrke. Injecting utility into anonymized datasets. In *SIGMOD*, 2006.
- [10] Paul C. Kocher. Timing attacks on implementations of Diffie-Hellman RSA, DSS, and other systems. In *CRYPTO*, pages 104–113, 1996.

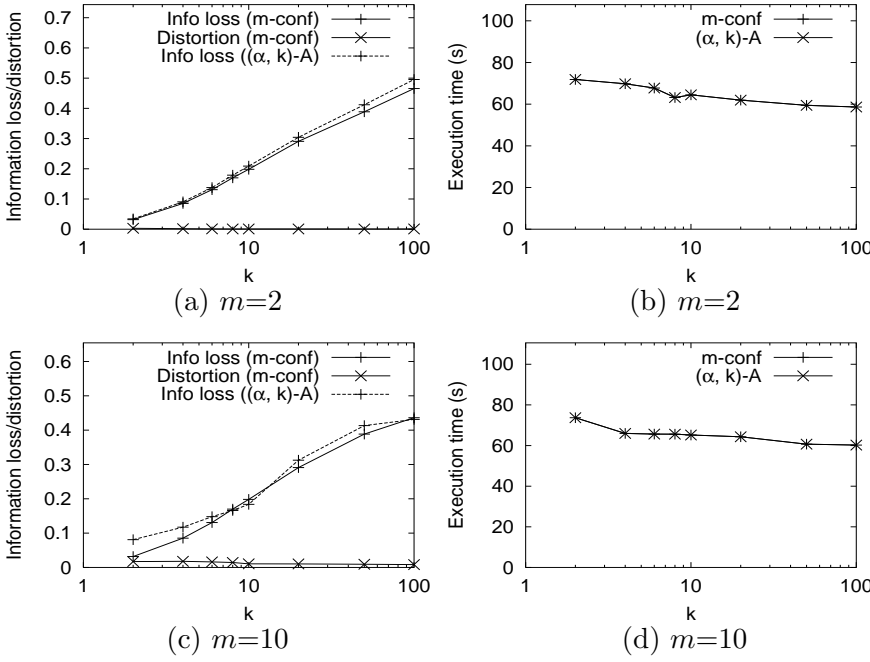


Figure 11: Two Tables case : effect of varying m and k

- [11] K. LeFevre, D. DeWitt, , and R. Ramakrishnan. Multidimensional k -anonymity. In *M. Technical Report 1521, University of Wisconsin*, 2005.
- [12] K. LeFevre, D. J. DeWitt, and R. Ramakrishnan. Incognito: Efficient full-domain k -anonymity. In *SIGMOD Conference*, pages 49–60, 2005.
- [13] N. Li and T. Li. t -closeness: Privacy beyond k -anonymity and l -diversity. In *ICDE*, 2007.
- [14] A. Machanavajjhala, J. Gehrke, and D. Kifer. l -diversity: privacy beyond k -anonymity. In *ICDE06*, 2006.
- [15] D. J. Martin, D. Kifer, A. Machanavajjhala, and J. Gehrke. Worst-case background knowledge for privacy-preserving data publishing. In *ICDE*, 2007.
- [16] A. Meyerson and R. Williams. On the complexity of optimal k -anonymity. In *PODS*, pages 223–228, 2004.
- [17] P. Samarati and L. Sweeney. Protecting privacy when disclosing information: k -anonymity and its enforcement through generalization and suppression, unpublished manuscript. In *unpublished*, 1998.
- [18] L. Sweeney. Weaving technology and policy together to maintain confidentiality. *Journal of Law, Medicine and Ethics*, 25(2-3):98–110, 1997.
- [19] L. Sweeney. Achieving k -anonymity privacy protection using generalization and suppression. *International journal on uncertainty, Fuzziness and knowledge based systems*, 10(5):571 – 588, 2002.
- [20] L. Sweeney. k -anonymity: a model for protecting privacy. *International journal on uncertainty, Fuzziness and knowledge based systems*, 10(5):557 – 570, 2002.
- [21] K. Wang and B. Fung. Anonymizing sequential releases. In *SIGKDD*, 2006.
- [22] K. Wang, B. C. M. Fung, and P. S. Yu. Handicapping attacker’s confidence: An alternative to k -anonymization. In *Knowledge and Information Systems: An International Journal*, 2006.
- [23] K. Wang, P. S. Yu, and S. Chakraborty. Bottom-up generalization: A data mining solution to privacy protection. In *ICDM*, pages 249–256, 2004.
- [24] R.C.W. Wong, J. Li, A. Fu, and K. Wang. (α , k)-anonymity: An enhanced k -anonymity model for privacy-preserving data publishing. In *SIGKDD*, 2006.
- [25] X. Xiao and Y. Tao. Anatomy: Simple and effective privacy preservation. In *Proceedings of the 32th Conference on Very Large Data Base (VLDB06)*, 2006.
- [26] X. Xiao and Y. Tao. Personalized privacy preservation. In *SIGMOD*, 2006.
- [27] Q. Zhang, N. Koudas, D. Srivastava, and T. Yu. Aggregate query answering on anonimized tables. In *ICDE*, 2007.

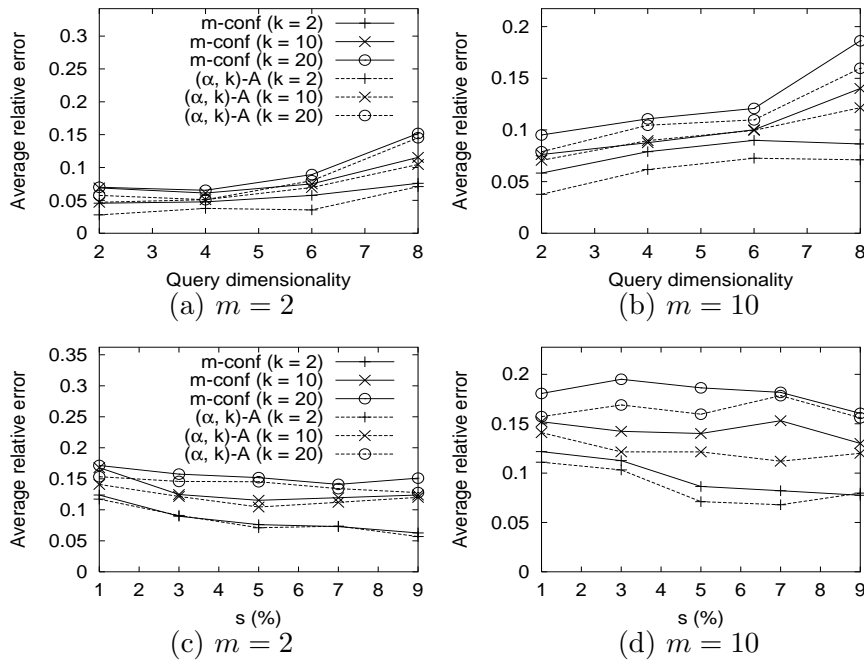


Figure 12: Two Tables Case : effects of varying query dimensionality and selectivity