

Research Article

Quasi-Identifier Recognition Algorithm for Privacy Preservation of Cloud Data Based on Risk Reidentification

Huda O. Mansour ^{1,2} Maheyzah M. Siraj ² Fuad A. Ghaleb ¹ Faisal Saeed ³
Eman H. Alkhamash ⁴ and Mohd A. Maarof¹

¹Faculty of Engineering, School of Computing, Universiti Teknologi Malaysia (UTM), Johor 81310, Malaysia

²Department of Computer Science, Faculty of Computer Science and Information Technology, University of Kassala, Kassala 31111, Sudan

³College of Computer Science and Engineering, Taibah University, Medina, Saudi Arabia

⁴Department of Computer Science, College of Computers and Information Technology, Taif University, P.O. Box 11099, Taif 21944, Saudi Arabia

Correspondence should be addressed to Fuad A. Ghaleb; abdulgaleel@utm.my

Received 30 April 2021; Revised 26 June 2021; Accepted 9 August 2021; Published 26 August 2021

Academic Editor: Ihsan Ali

Copyright © 2021 Huda O. Mansour et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Cloud computing plays an essential role as a source for outsourcing data to perform mining operations or other data processing, especially for data owners who do not have sufficient resources or experience to execute data mining techniques. However, the privacy of outsourced data is a serious concern. Most data owners are using anonymization-based techniques to prevent identity and attribute disclosures to avoid privacy leakage before outsourced data for mining over the cloud. In addition, data collection and dissemination in a resource-limited network such as sensor cloud require efficient methods to reduce privacy leakage. The main issue that caused identity disclosure is quasi-identifier (QID) linking. But most researchers of anonymization methods ignore the identification of proper QIDs. This reduces the validity of the used anonymization methods and may thus lead to a failure of the anonymity process. This paper introduces a new quasi-identifier recognition algorithm that reduces identity disclosure which resulted from QID linking. The proposed algorithm is comprised of two main stages: (1) attribute classification (or QID recognition) and (2) QID dimension identification. The algorithm works based on the reidentification of risk rate for all attributes and the dimension of QIDs where it determines the proper QIDs and their suitable dimensions. The proposed algorithm was tested on a real dataset. The results demonstrated that the proposed algorithm significantly reduces privacy leakage and maintains the data utility compared to recent related algorithms.

1. Introduction

In the modern information age, many companies are using external sources of data for processing, storing, or obtaining some services such as data mining. Unlimited computational resources, reduced costs, nonburden of maintenance, and nondiligence to learn the skills of proficiency in certain services, all of these were temptations to advance to the modern change. However, there are still security and privacy concerns that hinder the use of the features offered by the cloud [1]. Numerous studies clarified that attackers often reveal the information from third-party services or third-party

clouds [2]. For example, one of the security breaches in October 2014 was a breakthrough for Dropbox. The attackers stole 700 user passwords to obtain cash values of its Bitcoins (BTC). In 2015, a lot of users' information, which exceeds 4 million, such as the user's name, date of birth, address, e-mail, phone number, and other sensitive data, were leaked through the TalkTalk service provider in the UK. In 2016, Time Warner, one of the largest cable television companies in the United States, has announced that about 32 million passwords and e-mail of the users have been stolen via an attacker. In 2017, more than 200 million data of the users containing users' names, phone numbers,

e-mail addresses, home addresses, and other data have been disclosed through the API of McDelivery Company in India [2, 3]. A fresh security violation in Google displayed that any administrator of the server who has access to the secret information can misuse it easily. The worst problem is that administrator of the honest-but-curious server can violate privacy without being discovered [4].

Three kinds of the disclosure can cause privacy leakage, identity disclosure, attribute disclosure, and membership disclosure [5]. In attribute disclosure and identity disclosure, the intruder identifies that the tuple of the target individual is found in the released dataset and he aims to acquire some private/sensitive data about that individual from the released dataset [6]. Serious issues that lead to identity disclosure are quasi-identifier (QID) value linking and the attacker's knowledge background. The QIDs are the dataset attributes that if each of them is considered separately does not distinguish the individual, but when several attributes are combined they can give a distinctive identification of individuals [7]. For example, when looking at the attributes of date of birth, gender, and ZIP code together, one can reidentify the individuals as stated in [8]. Reidentification of the individuals through linking their QIDs leads to what are called linking attacks. Therefore, the careless publication of QIDs will lead to leakage of privacy [9].

One of the popular practices to avoid privacy leakage is anonymization. The anonymization can be performed via several types of transformations, by removing the values, changing the structure, replacing the values by taxonomy, and combining the values. The anonymization-based methods use one or a combination of operations to accomplish an optimum level of concealment [10]. A commonly utilized privacy criterion of anonymization is k -anonymity introduced by Sweeney [8]. The k -anonymization model is aimed at making any record in the released dataset that cannot be distinguished from at least $(k - 1)$ other records [1, 11]. To avoid the linking attacks, k -anonymization can be used. The effective method to determine the real QIDs is the primary issue for privacy-preserving methods based on k -anonymity or other anonymization models seek to prevent QID linking. While most of the current methods neglected this issue or just determine QIDs manually, this reduces the validity of the anonymization method as well as negatively affects the usefulness of anonymous data [9]. This study is aimed at overcoming the identity disclosure resulting from QID linking and reducing the leakage of privacy by proposing a QID recognition (QIR) algorithm based on risk rate reidentification. The proposed algorithm comprises two main stages: (1) attribute classification (or QIDs Recognition) and (2) QID dimension identification. The algorithm works based on the reidentification of risk rate for all attributes and the dimension of QIDs where it determines the proper QIDs and their suitable dimensions. Figure 1 shows the cause-effect diagram of privacy leakage. The dark boxes in Figure 1 explain the privacy leakage causes addressed by the proposed QID recognition (QIR) algorithm in this study. As shown in Figure 1, it is essential to properly identify the QID attributes to overcome the identity disclosure to reduce the leakage of privacy resulting from QID linking. This

paper is made up of 5 sections. Section 2 describes the state of the art of privacy-preserving data mining (PPDM) over the cloud, whereby some of the current methods and algorithms that address the issue of identification QIDs accurately to avoid identity disclosure are presented. A detailed description of the proposed algorithm has been provided in Section 3. Section 4 demonstrates the experimental evaluation, discussion, and comparison with related work. Section 5 concludes this work.

2. Related Work

The research of privacy-preserving outsourced data focuses on anonymization-based methods [12–18], cryptographic-based methods [19–24], hybrid methods [2, 25–27], and methods that seek to improve the data utility [26, 28, 29]. Some recent studies have demonstrated the privacy requirements of incremental datasets [30–32] and multiple sensitive attributes [33–35]. However, most of these studies neglected the issue of identification of the right QIDs, despite its importance in the success of the anonymity process. Few of these studies have attempted to introduce methods so that identification of the QIDs is required in the anonymization process, as presented in the next section.

Huang and others [36] introduce a new method that depends on the hypergraph to find a group of related views and QID set. This method maps the group of related views into a hypergraph and includes all paths available between every two nodes instead of finding the group of related views. The weakness of this method is that the QID group produced may include so many attributes. Further, it has high computational complexity resulting from the process of degeneration of the common graph from the hypergraph.

Omer and Mohamad [37] introduce a new method to select a quasi-identifier (QID) to achieve k -anonymity. Selective and decompose algorithms depend on nominating multiple attributes as a set and then generating power set $P(S)$ for them. Following that, the distinct values of the power set $P(S)$ elements were computed and listed in a table. Finally, the candidate element from the power set is the element with the maximum distinct value. The main problem in this method is selecting the primary nominate set of attributes, where the accuracy of the selection depends on the user experience [9]. Furthermore, it is impractical to generate $P(S)$ if the number of attributes is big (e.g., more than 8).

Y. J. Lee and K. H. Lee [38] examine the factors and the likelihood of an individual reidentified for medical information through inferable QIDs. The QIDs were considered as database variables that enable the reidentification of individuals by linking their QIDs with available external information or a specific individual. They selected five factors to form QID attributes to prevent patient privacy violations. The factors were selected based on their influence on the likelihood of reidentification and the possibility of inferring it from background knowledge. One of the disadvantages of this study is that the QIDs that can be extracted to reidentify patients' records may exceed 5. Besides, the paper focused only on the problem of reidentification of patients' records and avoiding leakage of privacy in the medical

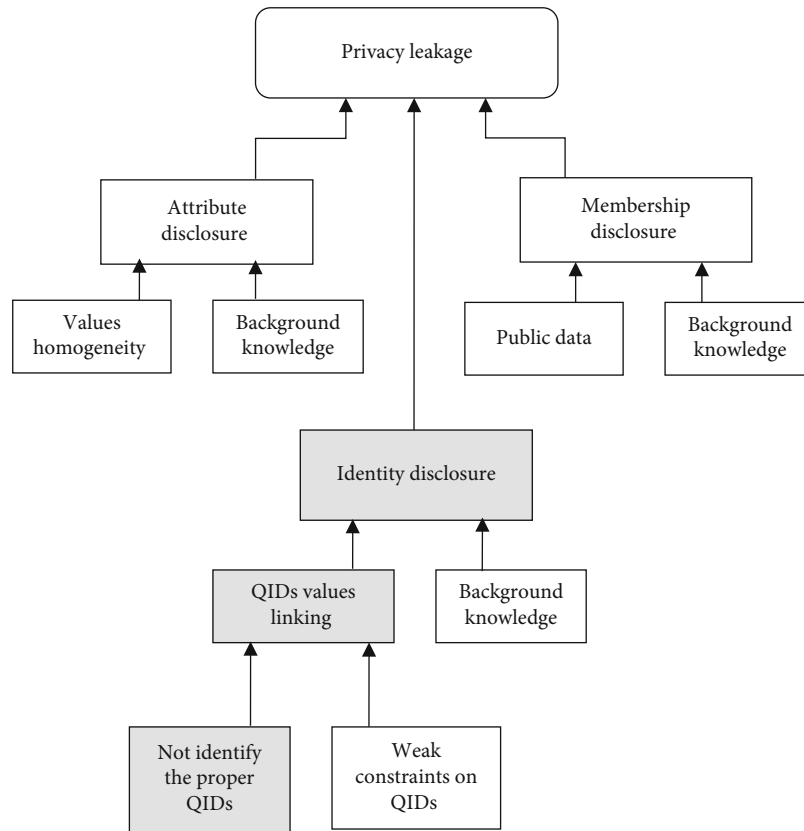


FIGURE 1: Privacy leakage causes addressed by the QIR algorithm.

records, lacking a public method that could be used for general data publishing. Bampoulidis and others [7] assume that some QIDs are more important than others (i.e., in data mining/analysis) and, therefore, should be distorted as little as possible in the anonymization process. They present a tool to address the issue of QIDs by utilizing a local recoding algorithm for k -anonymity. The tool outperforms the ARX (data anonymization tool) in terms of dataset quality. The major problem with this method is that it depends on the user in defining the QID attributes, giving priority to each attribute, as the user relies on his personal experience in determining the QID attributes, which are usually not accurate.

Kaur and Agrawal [10] study the impact of QIDs on the anonymization process. They gave new ways to consider before choosing the quasi-identifiers. The reidentification risks have been examined using different QIDs, diverse parameters, and different sizes of a data sample. The results of their work showed that when making the variance in selecting the QIDs for anonymization operation, note that the risk of reidentification increases when the number of QIDs increases, and it decreases when using QIDs that contain fewer categories. Although it is good to take into account these observations before starting the anonymity process, it should be noted that these observations extracted by the study are not fixed and may change from one dataset to another.

Wong and others [39] do not reveal the complete set of quasi-identifiers (QID) to the data collector before and after

the data anonymization process. They believed that the QIDs can be both sensitive values and identifying values; they allow the respondents/data owners to hide sensitive QID attributes from other parties. The first issue with this method is that the QID attributes that respondents consider them are sensitive which may contain data that are very useful in mining or may adversely affect mining outcomes. The second issue is if respondents submit inaccurate data, there is no guarantee of the usefulness of the results obtained from data analysis.

Sei and others [40] consider that some QIDs are regarded as sensitive QIDs and they propose novel privacy models, namely, $(l1, \dots, lq)$ – diversity and $(t1, \dots, tq)$ – closeness, and a method that can treat sensitive QIDs. Their proposed method comprises two algorithms: anonymization and reconstruction algorithms that can treat sensitive QIDs. Although this method can perform anonymity while preserving the quality of the data, it suffers from the problem of the Wong [39] method; this is because there is no effective method to accurately determine which of the QID attributes is considered sensitive QIDs.

Victor and Lopez [41] offer a (k, n, m) anonymity method for sensitive/private data based on the k -anonymity. The graph algorithms were used to perform QIDs and are moreover been improved by selecting similar QIDs based on the composite and derived attributes. The set of QIDs obtained from the methods in [36, 41] may include too many attributes, which increases the information loss in models based on generalizations like the k -anonymity [9].

3. The Proposed QID Recognition Algorithm

There are two main stages involved in the QID recognition algorithm (QIR) to prevent privacy leakage of outsourced data. *First*, classify the dataset attributes into quasi-identifiers (QIDs), sensitive attributes (SAs), and nonsensitive attributes (NSs). That is, each attribute in the dataset is classified into one of the aforementioned groups (QIDs, SAs, or NSs). In the attributes' classification (QID recognition) stage, the IDs (identifier attributes) are usually removed from the dataset by the data owner. The quasi-identifiers (QIDs) are the attributes that, when linked together, define the individual, for example, age, gender, and ZIP. The sensitive attributes (SAs) are the attributes that explain sensitive/private information about an individual such as medical information, financial records, and location. Meanwhile, the nonsensitive attributes (NSs) are the other attributes in the dataset that do not fall under the previously mentioned categories, as they do not help reidentify the identity of the individual, for example, state and religious attributes. In the basic privacy models (such as k -anonymity [7–9, 11–13, 18, 28], l -diversity [40, 42], and t -closeness [34, 43]), the attributes of a dataset were categorized into two groups: sensitive and nonsensitive. Meanwhile, most of the recent researchers such as in [9, 44–47] divide the dataset attributes into three types: QID, SA, and NS (not including identifiers) directly. Accordingly, the classification of dataset attributes in this study is divided into three types of QID, SA, and NS (not including identifiers) utilizing the same definitional meaning of each category as in the previous work in [9, 44–47].

Second, determine the actual dimension of QIDs that should be used in an anonymization operation that will achieve optimum case. If the set of QIDs contains too many attributes, the loss of information caused by generalization will be exacerbated. Nonetheless, sometimes the minimal set of QID does not imply the most appropriate privacy protection setting because the method does not consider what attributes the adversary could potentially have [37]. Therefore, we need a mechanism that determines the appropriate dimension of the QIDs to avoid these problems. In the QID dimension determining stage, the proposed algorithm performs this task. Figure 2 illustrates the general procedure of the two main phases of the QIR algorithm. The following subsections explain these two stages in more detail.

3.1. QID Recognition Stage. In this stage, the algorithm classifies the attributes depending on the reidentification risk rate for each attribute in the dataset, and then, the risk rate of the attribute is compared to the threshold values of the classification. As shown in Figure 2, the attribute classification stage comprises four main activities. These activities include (1) dataset preprocessing, (2) computing risk rate for all attributes, (3) selecting the classification thresholds, and (4) classifying the attributes according to the selected thresholds.

In the first activity, the dataset is preprocessed which includes filling the missing values, fixing the inconsistencies in the dataset, and data normalization. Then, in the second

activity, the risk rate is computed according to the g -distinct which is adopted in computing the reidentification risk rate [48]. A detailed description of the g -distinct method is presented in the next section. In the third activity, the classification thresholds were selected based on the maximum and minimum risk of reidentification as follows. These thresholds are denoted by β and α in this study; α threshold represents the maximum risk of reidentification of the individual while β represents the minimum risk of reidentification. The threshold values can be determined by the user or the data owner after calculating the reidentification risk for all attributes. Based on percentages of the highest and lowest attribute risk, one can choose the α value to be less than the highest risk value and choose the β value to be less than the lowest risk value. The nature of the data and the degree of importance of each attribute affect the selection of the threshold values. So, these thresholds are adjustable and differ from one dataset to another. For instance, let the dataset (D) contain attributes (A_1, A_2, \dots, A_n) , i.e., $D = A_1, A_2, \dots, A_n$; let $\beta = 0.05\%$ and $\alpha = 30\%$. Let Risk_{A_i} be the reidentification risk of attribute A_i and $\text{Risk}_{A_i} = 35\%$. As $\text{Risk}_{A_i} > \alpha$, then A_i is classified as SA. Suppose Risk_{A_3} and Risk_{A_5} are 23 and 0.01, respectively, then A_3 is classified as QID while A_5 will be classified as NS, respectively. Reidentification risk rate of attribute A_i computes the degree that makes the records distinguished based on this attribute. Finally, the fourth activity includes classifying the attributes according to the selected thresholds using rules represented by *if-else* testaments (see Algorithm 1, lines 27–39). In the following subsection, a detailed description of computing the reidentification risk rate (g -distinct) is presented. More explanation of the QID recognition stage is also presented.

3.2. g -Distinct. The g -distinct is adopted in computing the reidentification risk rate [48]. A person or record in any dataset is said to be unique if he/she or it has a combination of attributes that is not for someone/record else. The person/record is g -distinct if their combination of attributes is matching to $g-1$ or less than other people/records in the dataset [48]. Thus, uniqueness is the base situation of 1-distinct. In general, g -distinct is the total of the number of subgroups with i individuals, which is computed as

$$h_n(g) = \sum_{i=1}^g i \cdot f_n(i), \quad (1)$$

where $f_n(i)$ refers to the expected number of subgroups with i individuals that can be derived from a given aggregated group and g represents the whole number of individuals in a subgroup. That is, g is associated with the g -distinct to represent the number of distinguished individuals in the subgroup. For example, when we say 3-distinct, it means that three individuals have common QID characteristics out of the total number of people g in the subgroup. The sum of all g -distinct of individuals in a specific attribute represents the reidentification risk rate that the attribute potential to cause it. We can compute the general risk of the whole

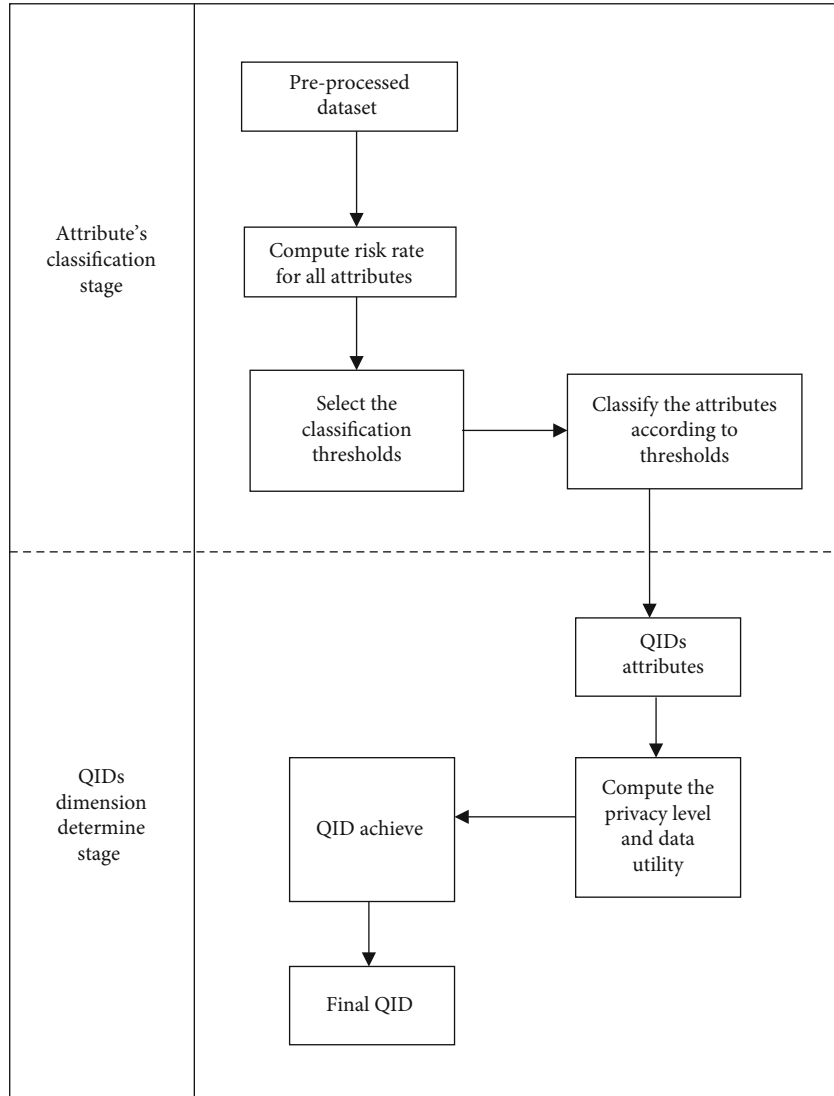


FIGURE 2: The general procedure of the proposed QIR algorithm.

dataset through equation (2) where b is the number of possible subgroups.

$$R_n^j(g) = \left(\frac{j}{n}\right) b^{1-n} (b^n - (b-1)^n). \quad (2)$$

Finally, the attribute classification stage returns the reidentification risk rate for each attribute in the dataset. Based on the resulting reidentification risk rates, the dataset attributes are classified to sensitive and nonsensitive according to the rate of the reidentification risk for each attribute in addition to threshold values β , α . The outcomes of this stage will be input into the QID dimension identification stage to determine the dimension of QIDs that is suitable to achieve optimal privacy requirements. The practical steps of the classification stage are explained by Algorithm 1. Lines 2–16 in Algorithm 1 are to compute the g -distinct for all dataset attributes while lines 18–26 are to calculate the reidentification risk rate based on the attributes' g -distinct. Finally, lines

28–40 addressed the process of attribute classification using the reidentification risk rate of each attribute to produce three categories of attributes: QIDs, SAs, and NSs.

The importance of this stage of the proposed algorithm represented by Algorithm 1 is that it contributes to reducing the attribute disclosure resulting from linking the QID values due to a weakness/failure in defining the QID characteristics correctly. This contribution helps in minimizing the leakage of information and avoiding privacy violations.

3.3. QID Dimension Identification Stage. This stage of the algorithm is aimed at determining the best dimension of QIDs that will achieve optimum cases. The optimum case gives high privacy with a high/reasonable percentage of preserving data quality. In other words, it has high privacy gain (PG) with high/reasonable nonuniform entropy (NUE). Algorithm 2 describes the implementation steps for this stage. The algorithm takes a sample of data with the QID that has the highest reidentification risk rate. Following that, the QIR calculates the PG and NUE base on k -anonymity

```

Input: dataset  $D$ ,  $\beta$ ,  $\alpha$ .
Output: classified dataset.
1: //Compute  $g$ -distinct for all dataset tuples for each attribute.
2:  $Dg_{Attr} \leftarrow g$ -distinct of the attribute (Attr)
3:  $n \leftarrow$  attribute domain
4:  $m \leftarrow$  tuple domain
5:  $Attr \in n$ 
6:  $g \in m$ 
7:  $tv \leftarrow$  attribute value of a specific tuple
8:  $Attr_{Dg}[i][j] = 0$ 
9: For  $i := 1$  to  $n.length$  do
10:   For  $j := 1$  to  $m.length$  do
11:      $Dg_{Attr}[i] = 1 / \int (tv)_j$ 
12:      $Attr_{Dg}[i][j] = Attr_{Dg}[i][j] + Dg_{Attr}(i)$ ;
13:      $j = j + 1$ ;
14:   End
15:    $i = i + 1$ ;
16: end
17: //Compute reidentification risk rate for all dataset attributes.
18:  $Risk_{Attr}[i] = 0$ 
19:  $Risk_{Attr} \leftarrow$  reidentification risk rate of Attr
20: For  $i := 1$  to  $Attr_{Dg}[i].length$  do
21:   For  $j := 1$  to  $m.length$  do
22:      $Risk_{Attr}[i] = Risk_{Attr}[i] + Dg_{Attr}[i][j]$ 
23:      $j = j + 1$ ;
24:   End
25:    $i = i + 1$ ;
26: End
27: //Classify the attributes based on risk rate and threshold values.
28:  $QIDs[] = 0$ 
29:  $SAs[] = 0$ 
30:  $NSs[] = 0$ 
31: For  $i := 1$  to  $Risk_{Attr}[i].length$  do
32:   If ( $Risk_{Attr}[i]$  in range( $\beta$ ))
33:      $QIDs[i] = QIDs[] + Risk_{Attr}[i]$ ;
34:   Else If ( $Risk_{Attr}[i]$  in range( $\alpha$ ))
35:      $SAs[i] = SAs[] + Risk_{Attr}[i]$ ;
36:   Else
37:      $NSs[i] = NSs[] + Risk_{Attr}[i]$ ;
38:    $i = i + 1$ ;
39: end
40: Return( $QIDs[], SAs[], NSs[]$ )

```

ALGORITHM 1: Attribute classification.

through equations (3) and (4). In the next step, the QID number is increased, and PG and NUE are calculated again and so on until all QIDs are finished.

Finally, the algorithm determines the optimum case that gives high privacy with a high/reasonable percentage of preserving data quality. The best QID dimension is the QIDs with the optimum case. Algorithm 2 provides the executive steps of this stage; lines 5–12 implement the anonymization by k -anonymity on a sample of the dataset. It begins with QID that has the highest reidentification risk rate. After that, the algorithm calculates the privacy gain (PG) and nonuniform entropy (NUE) through equations (3) and (4). Then, the QID number is increased; PG and NUE have been calculated repeatedly until all the QIDs are finished. Lastly, in lines 13–15, the algorithm determines the best QID dimen-

sion (QidD) that achieves the optimum case to be involved in the anonymization process.

It was observed in study [9] that in most cases, when the QID dimension is large, the data loss increases. However, when the QID dimension is small, the privacy protection is not applied optimally because one cannot know what the actual QIDs an attacker possesses [37]. Therefore, determining an appropriate QID dimension is important to reduce data loss.

3.4. Performance Measures. Two performance evaluation measures were used in this study: the privacy gain (PG) and the nonuniform entropy (NUE). More explanation and the derivation of these measures are presented in the following subsections.

```

Input: dataset sample  $d$ , QIDs [], privacy parameter  $k$ .
Output: optimal dimension of QIDs.
1: QidD  $\leftarrow$  dimension of QIDs
2: QidD  $\in$  QIDs []
3: Optimal_QidD  $\leftarrow$  Optimal dimension of QIDs
4: QidD[] = 0
5: For  $i := 1$  to QIDs [].lengthdo
6:   QidD[ $i$ ] = QidD[] + QIDs[ $i$ ];
7:   Anonymized_data[ $i$ ] =  $k$ -anonymity( $d$ , QidD[ $i$ ],  $k$ );
8:   PG [  $i$  ] = Privacy_gain(Anonymized_data[ $i$ ]);
9:   NUE [  $i$  ] = Nonuniform_Entropy (Anonymized_data [  $i$  ]);
10:  Difference[ $i$ ] = PG [  $i$  ] - EIL[ $i$ ];
11:   $i = i + 1$ ;
12: end
13: If((PG[] == max)&&(NUE[] == max))
14:   Optimal_QidD[] = QidD[ $i$ ];
15: Return(Optimal_QidD[]).

```

ALGORITHM 2: QID dimension identification.

3.4.1. The Privacy Gain. To evaluate the privacy level for the proposed algorithm, equation (3) and Definition 1 are used as follows.

$$PG = A_{t(\text{gen})} - A_{b(\text{gen})}, \quad (3)$$

where $A_{t(\text{gen})}$ is anonymity after generalization (gen) and $A_{b(\text{gen})}$ is anonymity before generalization [27, 49, 50].

Definition 1. Anonymity quasi-identifier: a quasi-identifier qid is an anonymity quasi-identifier if $|\text{QIG}(\text{qid})| = \min_{\text{qid}' \in \text{QID}} |\text{QIG}(\text{qid}')|$, where $||$ represents the size of a QI group [50].

3.4.2. Nonuniform Entropy. In the context of data deidentification, the nonuniform entropy is to compare the frequencies of attribute values in the transformed dataset according to frequencies in the input dataset; it was originally introduced as a model for measuring the loss of information [51]. When a dataset D is transformed into another dataset D' , nonuniform entropy is defined as

$$\Delta(D, D') = \sum_{x \in D} -\log \left(\frac{f(D, x)}{f(D', x)} \right). \quad (4)$$

4. Experimental Evaluation

In this section, the experimental evaluation of our implementation algorithm will be presented in terms of PG and NUE. In Dataset Setup, we describe the datasets we have used for running the experiments and the experimental environment setup. In Experimental Results, we present the first set of experiments and provide the results from our algorithm. In Performance Benchmark and Discussion, we provide benchmark and discussion results of our algo-

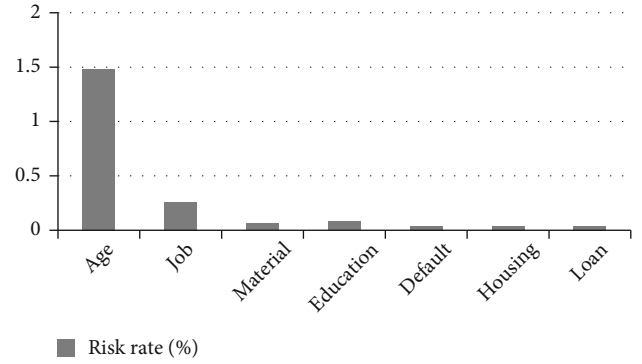


FIGURE 3: Risk rate of the bank dataset attributes.

gorithm against a close recent algorithm introduced by Omer and Mohamad [37].

4.1. Dataset Setup. Two real-life datasets from the University of California–Irvine were used in this study to demonstrate the performance of the proposed algorithms. The first is the bank direct marketing dataset [52]. The bank dataset consists of 17 attributes and 45,211 tuples and does not include any missing values. The dataset attributes are divided into three divisions which are (1) data of bank clients: age, job, marital status, education, default, balance, housing, and loan; in this paper, we will consider these attributes because these attributes are significant for bank clients and reidentification purposes; (2) data related to the last contact of the current campaign; and (3) other attributes like the campaign and days. The second dataset is the adult dataset [53] used as a standard for anonymization algorithm evaluation [7] consisting of 48,842 census records and 15 attributes.

ARX data anonymization software is open source introduced and developed by Prasser et al. [54] for data anonymization; we used it to implement the algorithms as explained in the following sections. The experiments were executed on

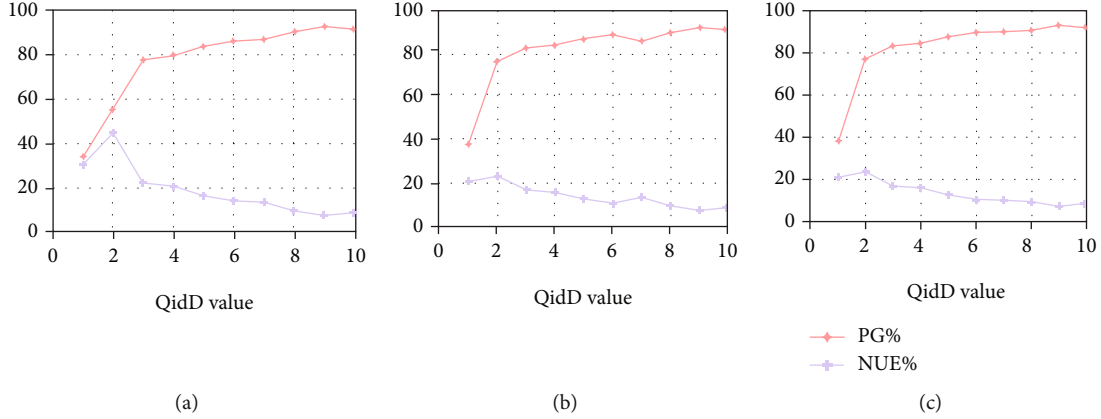
FIGURE 4: (a–c) The best QidD selection for the bank dataset by QIR on different k values.

TABLE 1: Classification of the bank dataset.

Classification	Threshold value $\alpha = 30, \beta = 0$	Attributes
SAs	$R_{risk} > \alpha$	Balance
QIDs	$\beta \leq R_{risk} < \alpha$	Age, job, education, and marital status
NSs	$R_{risk} < \beta$	Default, housing, and loan

TABLE 2: Classification of the adult dataset.

Classification	Threshold value $\alpha = 0.2, \beta = 0.01$	Attributes
SAs	$R_{risk} > \alpha$	Capital gain, capital loss
QIDs	$\beta \leq R_{risk} < \alpha$	Hours-per-week, work-class, age, native-country, education, education-num, occupation, marital-status, relationship, and race
NSs	$R_{risk} < \beta$	Sex, income

TABLE 3: Experimental results for selecting the best QidD in the adult dataset.

QID value	$k = 5$		$k = 15$		$k = 25$	
	PG %	NUE %	PG %	NUE %	PG %	NUE %
1	33.8	30.41	38.35	21.05	38.35	21.05
2	55.34	44.65	76.86	23.13	76.86	23.13
3	77.94	22.05	83.17	16.82	83.17	16.82
4	79.53	20.46	84.39	15.6	84.39	15.6
5	83.62	16.37	87.48	12.51	87.48	12.51
6	85.91	14.08	89.56	10.43	89.56	10.43
7	86.65	13.34	86.65	13.34	89.69	10.3
8	90.51	9.48	90.51	9.48	90.51	9.48
9	92.68	7.31	92.68	7.31	92.68	7.31
10	91.59	8.4	91.59	8.4	91.59	8.4

a machine with an Intel Core i7 2.7 GHz processor with 8 GB RAM, under Windows 10.

4.2. *Experimental Results.* The first experiment is to classify the dataset attributes according to their risk rate. Figures 3

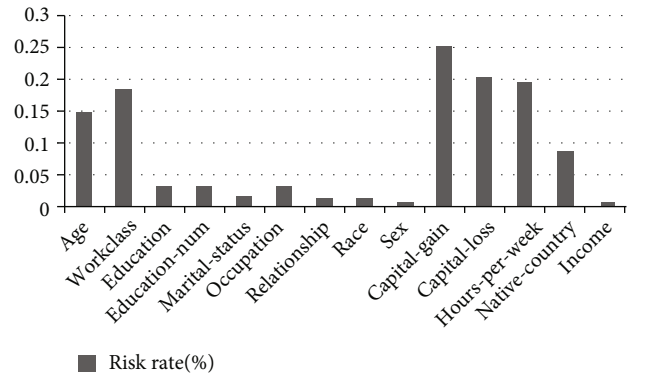


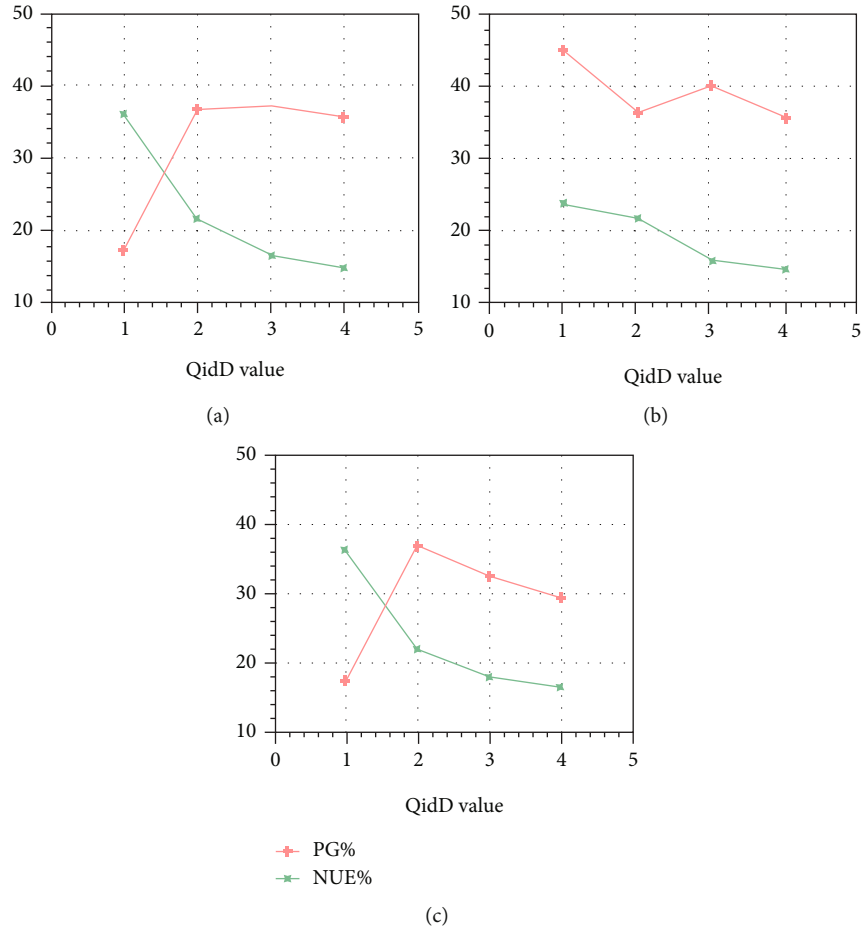
FIGURE 5: Risk rate of the adult dataset attributes.

and 4 illustrate the risk rate for bank attributes and adult attributes, respectively.

For the bank dataset, we identify α and β as $\alpha = 30, \beta = 0$. Table 1 demonstrates bank attribute classification. In the adult dataset, we add $\alpha = 0.2, \beta = 0.01$ to classify the attributes. Table 2 demonstrates the classification of the adult dataset. Because the “balance” attribute has a risk of 52.04 %, which is large compared to other attributes, it is excluded

TABLE 4: Experimental results for selecting the best QidD in the bank dataset.

QidD	QID	$k = 5$		$k = 15$		$k = 25$	
		PG %	NUE %	PG %	NUE %	PG %	NUE %
1	Age	23.89	45.28	36.12	17.27	36.12	17.27
2	Age, job	21.83	36.65	21.83	36.65	21.83	36.65
3	Age, job, marital status	15.83	40.35	16.67	37.18	17.94	32.37
4	Age, job, marital status, education	14.88	35.93	14.88	35.93	16.43	29.26

FIGURE 6: (a–c) The best QidD selection for the bank dataset by QIR on different k values.

from Figure 3 to highlight the difference between the attributes that have relatively small risk values.

After calculating the risk rate of each attribute in the dataset, the attribute is classified according to the selected threshold α and β as was explained in QID Recognition Stage. Tables 1 and 2 show the classification results of the bank dataset and the adult dataset, respectively, according to the selected classification thresholds α and β for each dataset. After the classification stage, the best dimension of QIDs that achieves optimum case should be determined. In the bank dataset, the QID dimension (QidD) is four (QidD = 4) while in the adult dataset QidD is 10 (QidD = 10). For each dataset, the initial value of QID dimension is set to one (QidD = 1) to be used as input into the proposed QID dimension identification algorithm (as explained in Algorithm 2) Identification of QID dimension

begins with the initial value of QidD, and it is incremented until the maximum number of QID dimension. Identification of QID dimension begins also with a sample size equal to 10% of the dataset with k -anonymity of 5, and it is incremented until $k = 25$ for each QidD value (sample size is changeable). Then, the privacy gain (PG) and the nonuniform entropy (NUE) are calculated for each sample and each new QidD until QidD values reach four (QidD = 4) for the bank dataset and QidD = 10 for the adult dataset.

Finally, the proposed algorithm returns the QidD that achieves the optimum case to be as the best dimension will be used in the anonymization process. Table 3 demonstrates the results of finding the best QidD for the adult dataset.

According to Table 3, we observed that QidD = 2 is the optimum case that increases the privacy gain as well as the NUE. Moreover, we can notice that the privacy level also

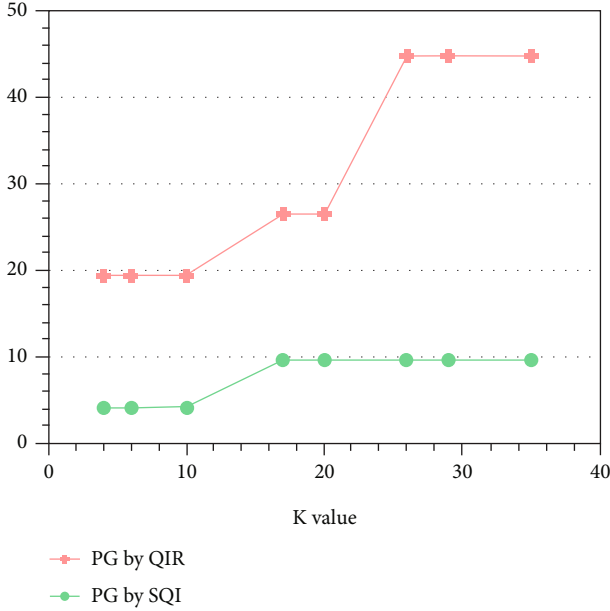


FIGURE 7: PG at several k values. Dataset: adult dataset 48,842 tuples. QidD of QIR = 2 (work class, HPW). QidD of SQI = 1 (age).

increases when QidD value increases. The privacy gain reaches 91.59% when QidD is 10. On the other hand, NUE decreases, and accordingly, the data utility decreases when QidD increases. Figures 4(a)–4(c) demonstrate the selection of the best QidD for the bank dataset by the proposed QIR algorithm on different k -anonymity values, 5, 15, and 25, respectively. In the bank dataset, the proposed algorithm’s selected QID attributes are work-class and hours-per-week (HPW). These two attributes achieve the highest reidentification risk; thus, they must be involved in the anonymization process (see Figure 5).

To determine the best QidD in the bank dataset, track Table 4 and Figures 6(a)–6(c); it is clear that when QidD = 1 the proposed algorithm achieves the optimum case as it gives high privacy in several cases of k values. It can be also observed in Table 4 that NUE drops from 45.28% when $k = 5$ to 17.27% when k increases above 15. It is also noticeable in the bank database that privacy decreases as the value of QidD increases which is normal with the level of privacy provided.

4.3. Performance Benchmark and Discussion. To evaluate the proposed QIR algorithm, we compare it based on k -anonymity against recent similar work SQI algorithm [37]. The comparison was conducted in terms of their privacy gain (PG) and nonuniform Entropy (NUE). Multiple k values and different dataset sizes of the adult dataset will be used. In Figures 7 and 8, the privacy provided by QIR is more than the privacy achieved by SQI, where the improvement average exceeds 23%. Although SQI outperformed the QIR in data utility represented by NUE at $k = 26, 29, 35$, with a privacy rate of 9.57%, this is considered a deficiency because QIR provided data utility higher than that with much higher privacy at $k = 4, 6, 10, 17$, and 20.

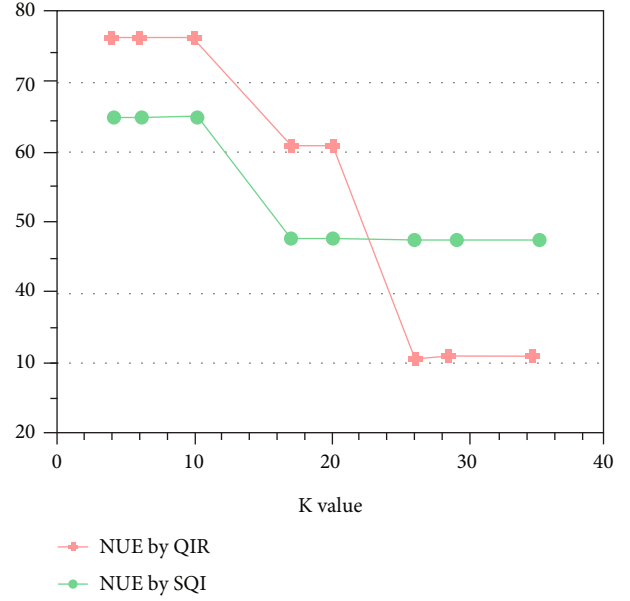


FIGURE 8: NUE at several k values. Dataset: adult dataset 48,842 tuples. QidD of QIR = 2 (work class, HPW). QidD of SQI = 1 (age).

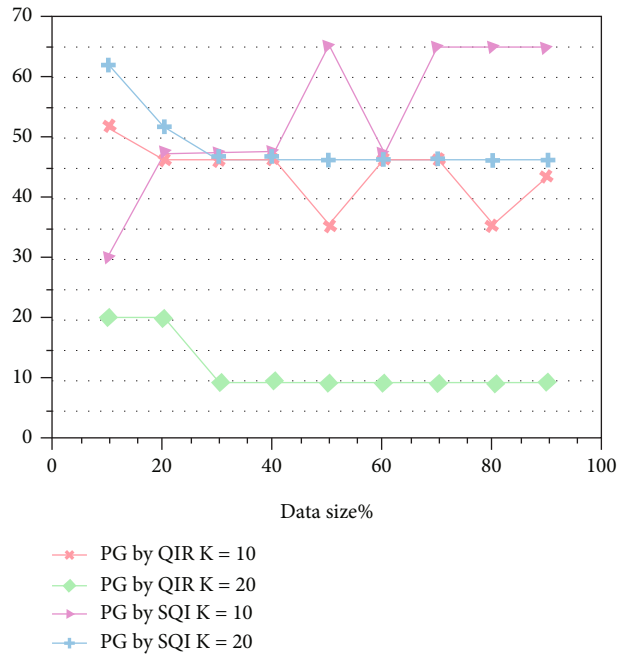


FIGURE 9: PG at several data sizes. Dataset: adult dataset. QidD of QIR = 2 (work class, HPW). QidD of SQI = 1 (age). $k = 10, 20$.

In Figures 9 and 10, it can be observed that at 10% of the dataset and $k = 10$ the privacy achieved by the proposed QIR algorithm is more than double the privacy achieved by the SQI algorithm with slight increases in data utility, that is, the proposed QIR algorithm outperforms the SQI algorithm in terms of preserving privacy and data utility. With data size 20% and $k = 20$, NUE obtained by SQI and QIR is 30.27 and 31.66%, respectively, while the privacy given by SQI is 20.52% and that by QIR is 51.82 which is twice more than that achieved by SQI. Similar results were obtained at

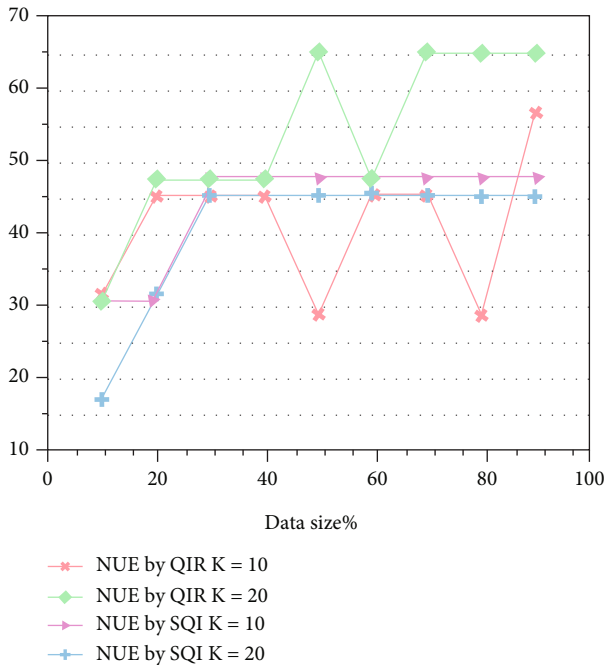


FIGURE 10: NUE at several data sizes. Dataset: adult dataset. QidD of QIR = 2 (work class, HPW). QidD of SQI = 1 (age). $k = 10, 20$.

$k = 20$ and data size = 30% and 90%, respectively. In most cases, when data size increases the privacy decreases, and therefore, the data utility increases.

Generally, for the whole adult data, results of the experiments at $k = 10$ and $k = 20$ show that the average privacy percentage presented by SQI is 10.17% with 48.62% data utility, while the average privacy percentage offered by the proposed QIR is 46.49% with 41.04% data utility. Also, for the whole adult dataset and all k values experimented, the average privacy provided by SQI is 7.51% against 54.13% data utility, while the average privacy percentage achieved by QIR is 30.67% against 55.46% data utility; hence, using QIR for identification of the real QIDs is considered more ideal.

5. Conclusions

Accurate identification of QIDs is an important issue for the success and validity methods of privacy-preserving outsourced data that seek to avoid privacy leakage caused by QID linking. This paper is aimed at classifying dataset attributes before the anonymization process and determining the proper QIDs that should be involved in the anonymity operation. A new algorithm is proposed based on the calculation of the reidentification risk for dataset attributes to classify attributes to SAs, QIDs, and NSs based on prespecified thresholds. In addition to attribute classification, the algorithm determines the actual dimension of QIDs that is required in the anonymization process depending on the amount of privacy provided versus a loss of the quality of the data. The experiment results indicated that the proposed identification algorithm has better performance and is more perfect in terms of privacy provided against data utility when

compared with other works. Although the proposed algorithm is suitable to be used with any method or privacy model concerned with QID attributes, in this paper, we have relied on the k -anonymity model.

Data Availability

All data used in this article are available in the machine learning repository at the University of California, Irvine (UCI): <https://archive.ics.uci.edu/ml/datasets/>.

Conflicts of Interest

The authors declare that there is no conflict of interest regarding the publication of this paper.

Acknowledgments

The authors would like to acknowledge Taif University Researchers Supporting Project (number TURSP-2020/292) Taif University, Taif, Saudi Arabia.

References

- [1] J. Domingo-Ferrer, O. Farràs, J. Ribes-González, and D. Sánchez, "Privacy-preserving cloud computing on sensitive data: a survey of methods, products, and challenges," *Computer Communications*, vol. 140–141, pp. 38–60, 2019.
- [2] S. Aldeen Yousra and S. Mazleena, "A new heuristic anonymization technique for privacy preserved datasets publication on cloud computing," *Journal of Physics: Conference Series*, vol. 1003, p. 012030, 2018.
- [3] C. Bradford, "7 most infamous cloud security breaches - StorageCraft," *storagecraft*, 2020, <https://blog.storagecraft.com/7-infamous-cloud-security-breaches/>.
- [4] B. Chen, P. Cheung, P. Cheung, and Y. Kwok, "Cypherdb: a novel architecture for outsourcing secure database processing," *IEEE Transactions on Cloud Computing*, vol. 6, no. 2, pp. 372–386, 2018.
- [5] B. C. M. Fung, K. Wang, R. Chen, and P. S. Yu, "Privacy-preserving data publishing," *ACM Computing Surveys*, vol. 42, no. 4, pp. 1–53, 2010.
- [6] S. A. Abdelhameed, S. M. Moussa, and M. E. Khalifa, "Privacy-preserving tabular data publishing: a comprehensive evaluation from web to cloud," *Computers & Security*, vol. 72, pp. 74–95, 2018.
- [7] A. Bampoulidis, I. Markopoulos, and M. Lupu, "PrioPrivacy: a local recoding K -anonymity tool for prioritised quasi-identifiers," in *IEEE/WIC/ACM International Conference on Web Intelligence - Companion Volume*, pp. 314–317, ACM: New York, 2019.
- [8] L. Sweeney, "Achieving k -anonymity privacy protection using generalization and suppression," *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, vol. 10, no. 5, pp. 571–588, 2002.
- [9] Y. Yan, W. Wang, X. Hao, and L. Zhang, "Finding quasi-identifiers for k -anonymity model by the set of cut-vertex," *Engineering Letters*, vol. 26, no. 1, 2018.
- [10] G. Kaur and S. Agrawal, "Differential privacy framework," in *Impact of Quasi-identifiers on Anonymization*, vol. 46, Springer, Singapore, 2019.

- [11] D. Wei, K. Natesan Ramamurthy, and K. R. Varshney, "Distribution-preserving k -anonymity," *Statistical Analysis and Data Mining: The ASA Data Science Journal*, vol. 11, no. 6, pp. 253–270, 2018.
- [12] P. R. Bhaladhare and D. C. Jinwala, "Novel approaches for privacy preserving data mining in k -anonymity model," *Journal of Information Science and Engineering*, vol. 32, no. 1, pp. 63–78, 2016.
- [13] M. S. Simi, K. S. Nayaki, and M. S. Elayidom, "An extensive study on data anonymization algorithms based on K -anonymity," *IOP Conference Series: Materials Science and Engineering*, vol. 225, p. 012279, 2017.
- [14] H. Kaur, N. Kumar, and S. Batra, "ClamPP: a cloud-based multi-party privacy preserving classification scheme for distributed applications," *The Journal of Supercomputing*, vol. 75, no. 6, pp. 3046–3075, 2019.
- [15] G. G. Dagher, B. C. M. Fung, N. Mohammed, and J. Clark, "SecDM: privacy-preserving data outsourcing framework with differential privacy," *Knowledge and Information Systems*, vol. 62, no. 5, pp. 1923–1960, 2020.
- [16] A. F. Westin, "Privacy and freedom," *American Sociological Review*, vol. 33, no. 1, p. 173, 1968.
- [17] M. Templ, *Statistical disclosure control for microdata*, Springer International Publishing, Cham, 2017.
- [18] W. Mahanan, W. A. Chaovalitwongse, and J. Natwichai, "Data anonymization: a novel optimal k -anonymity algorithm for identical generalization hierarchy data in IoT," *Service Oriented Computing and Applications*, vol. 14, no. 2, pp. 89–100, 2020.
- [19] S. Mayil, M. Vanitha, C. Science, J. J. College, and T. St, "A survey on privacy preserving data mining techniques for clinical decision support system," *International Research Journal of Engineering and Technology*, vol. 5, no. 5, pp. 6054–6056, 2016.
- [20] N. Uttarwar and M. A. Pradhan, "K-NN data classification technique using semantic search on encrypted relational data base," in *2016 International Conference on Computing Communication Control and automation (ICCUBEA)*, Pune, India, 2017.
- [21] K. El Makkaoui, A. Beni-Hssane, A. Ezzati, and A. El-Ansari, "Fast Cloud-RSA scheme for promoting data confidentiality in the cloud computing," *Procedia Computer Science*, vol. 113, pp. 33–40, 2017.
- [22] W. Wang, L. Chen, and Q. Zhang, "Outsourcing high-dimensional healthcare data to cloud with personalized privacy preservation," *Computer Networks*, vol. 88, pp. 136–148, 2015.
- [23] K. El Makkaoui, A. Beni-Hssane, and A. Ezzati, "Speedy Cloud-RSA homomorphic scheme for preserving data confidentiality in cloud computing," *Journal of Ambient Intelligence and Humanized Computing*, vol. 10, no. 12, pp. 4629–4640, 2019.
- [24] D. Chandravathi and P. V. Lakshmi, "Privacy preserving using extended Euclidean algorithm applied to RSA-homomorphic encryption technique," *VOLUME-8 ISSUE-10, AUGUST 2019, REGULAR ISSUE*, vol. 8, no. 10, pp. 3175–3179, 2019.
- [25] P. Shyja Rose, J. Visumathi, and H. Haripriya, "Research paper on privacy preservation by data anonymization in public cloud for hospital management on big data," *International Journal of Control Theory and Applications*, 2016.
- [26] Y. A. A. S. Aldeen and M. Salleh, "Privacy preserving data utility mining architecture," in *Smart Cities Cybersecurity and Privacy*, pp. 253–268, Elsevier Inc., 2019.
- [27] Y. A. A. S. Aldeen and M. Salleh, "Techniques for privacy preserving data publication in the cloud for smart city applications," in *Smart Cities Cybersecurity and Privacy*, pp. 129–145, Elsevier Inc., 2019.
- [28] Y. A. A. S. Aldeen and M. Salleh, "A hybrid K -anonymity data relocation technique for privacy preserved data mining in cloud computing," *Journal of Internet Computing and Services*, vol. 17, no. 5, pp. 51–58, 2016.
- [29] H. Lee, S. Kim, J. W. Kim, and Y. D. Chung, "Utility-preserving anonymization for health data publishing," *BMC Medical Informatics and Decision Making*, vol. 17, no. 1, p. 104, 2017.
- [30] Y. A. A. S. Aldeen, M. Salleh, and Y. Aljeroudi, "An innovative privacy preserving technique for incremental datasets on cloud computing," *Journal of Biomedical Informatics*, vol. 62, pp. 107–116, 2016.
- [31] S. R. P. Reddy, K. V. S. V. N. Raju, and V. V. Kumari, "Personalized privacy preserving incremental data dissemination through optimal generalization," *International Journal of Engineering & Technology*, vol. 7, no. 2.20, p. 197, 2018.
- [32] R. V. Sudhakar and T. C. M. Rao, "Security aware index based quasi-identifier approach for privacy preservation of data sets for cloud applications," in *Cluster Computing*, pp. 1–11, Springer, 2020.
- [33] S. A. Onashoga, B. A. Bamiro, A. T. Akinwale, and J. A. Oguntuase, "KC-Slice: a dynamic privacy-preserving data publishing technique for multisensitive attributes," *Information Security Journal: A Global Perspective*, vol. 26, no. 3, pp. 121–135, 2017.
- [34] R. Wang, Y. Zhu, T.-S. Chen, and C.-C. Chang, "Privacy-preserving algorithms for multiple sensitive attributes satisfying t -closeness," *Journal of Computer Science and Technology*, vol. 33, no. 6, pp. 1231–1242, 2018.
- [35] S. Sriyjanthi, T. Sethukarasi, and A. Thilagavathy, "Efficient anonymization algorithm for multiple sensitive attributes," *International Journal of Innovative Technology and Exploring Engineering*, vol. 9, no. 1, pp. 4961–4963, 2019.
- [36] L. Huang, J. Song, Q. Lu, X. Liu, and C. Zhang, "Hypergraph-based solution for selecting quasi-identifier," *International Journal of Digital Content Technology and its Applications*, vol. 6, no. 20, pp. 597–606, 2012.
- [37] A. M. Omer and M. M. Bin Mohamad, "Simple and effective method for selecting quasi-identifier," *Journal of Theoretical and Applied Information Technology*, vol. 89, no. 2, pp. 512–517, 2016.
- [38] Y. J. Lee and K. H. Lee, "Re-identification of medical records by optimum quasi-identifiers," in *2017 19th international conference on advanced communication technology (ICACT)*, pp. 428–435, PyeongChang, Korea, 2017.
- [39] K. S. Wong, N. A. Tu, D. M. Bui, S. Y. Ooi, and M. H. Kim, "Privacy-preserving collaborative data anonymization with sensitive quasi-identifiers," in *2019 12th CMI Conference on Cybersecurity and Privacy (CMI)*, Copenhagen, Denmark, 2019.
- [40] Y. Sei, H. Okumura, T. Takenouchi, and A. Ohsuga, "Anonymization of sensitive quasi-identifiers for l -diversity and t -closeness," *IEEE Transactions on Dependable and Secure Computing*, vol. 16, no. 4, pp. 580–593, 2019.
- [41] N. Victor and D. Lopez, "Privacy preserving sensitive data publishing using (k, n, m) anonymity approach," *Journal of communications software and systems*, vol. 16, no. 1, pp. 46–56, 2020.

- [42] A. Machanavajjhala, J. Gehrke, D. Kifer, and M. Venkatasubramanian, "L-diversity: privacy beyond k -anonymity," in *22nd International Conference on Data Engineering (ICDE'06)*, pp. 24–24, Atlanta, GA, USA, 2006.
- [43] N. Li, T. Li, and S. Venkatasubramanian, "t-Closeness: privacy beyond k -anonymity and l -diversity," in *2007 IEEE 23rd International Conference on Data Engineering*, pp. 106–115, Turkey, 2007.
- [44] H. Y. Tran and J. Hu, "Privacy-preserving big data analytics a comprehensive survey," *Journal of Parallel and Distributed Computing*, vol. 134, pp. 207–218, 2019.
- [45] K. Patel and G. B. Jethava, "Privacy preserving techniques for big data: a survey," in *2018 Second International Conference on Inventive Communication and Computational Technologies (ICICCT)*, pp. 194–199, Coimbatore, India, 2018.
- [46] E. E. Brown, "Improving privacy preserving methods to enhance data mining for correlation research," in *Southeast-Con 2017*, pp. 3–6, Concord, NC, USA, 2017.
- [47] X. Jiang, A. D. Sarwate, and L. Ohno-Machado, "Privacy technology to support data sharing for comparative effectiveness research: a systematic review," *Medical Care*, vol. 51, 8 Supplement 3, pp. S58–S65, 2013.
- [48] K. Benitez and B. Malin, "Evaluating re-identification risks with respect to the HIPAA privacy rule," *Journal of the American Medical Informatics Association*, vol. 17, no. 2, pp. 169–177, 2010.
- [49] X. Zhang, C. Liu, S. Nepal, C. Yang, W. Dou, and J. Chen, "Combining top-down and bottom-up: scalable sub-tree anonymization over big data using MapReduce on cloud," in *2013 12th IEEE International Conference on Trust, Security and Privacy in Computing and Communications*, pp. 501–508, Melbourne, VIC, Australia, 2013.
- [50] X. Zhang, C. Liu, S. Nepal, C. Yang, W. Dou, and J. Chen, "A hybrid approach for scalable sub-tree anonymization over big data using MapReduce on cloud," *Journal of Computer and System Sciences*, vol. 80, no. 5, pp. 1008–1020, 2014.
- [51] F. Prasser, R. Bild, and K. A. Kuhn, "A generic method for assessing the quality of de-identified health data," *Studies in Health Technology and Informatics*, vol. 228, pp. 312–316, 2016.
- [52] S. Moro, P. Cortez, and P. Rita, *UCI Machine Learning Repository: Bank Marketing Data Set*, 2014, <https://archive.ics.uci.edu/ml/datasets/Bank+Marketing>.
- [53] R. Kohavi and B. Becker, *Adult Census Income | Kaggle*, 2016, <https://www.kaggle.com/uciml/adult-census-income>.
- [54] F. Prasser, K. A. Kuhn, and J. Eicher, "Flexible data anonymization using ARX—current status and challenges ahead," *Software: Practice and Experience*, vol. 50, no. 7, pp. 1277–1304, 2020.