



**Titre:** Information theoretic-based privacy risk evaluation for data  
Title: anonymization

**Auteurs:** Anis Bkakria, Frédéric Cuppens, Nora Boulahia Cuppens, & Aimilia  
Authors: Tasidou

**Date:** 2021

**Type:** Article de revue / Article

**Référence:** Bkakria, A., Cuppens, F., Boulahia Cuppens, N., & Tasidou, A. (2021). Information  
Citation: theoretic-based privacy risk evaluation for data anonymization. Journal of  
Surveillance, Security and Safety, 2, 83-102.  
<https://doi.org/10.20517/jsss.2020.20>

 **Document en libre accès dans PolyPublie**  
Open Access document in PolyPublie

**URL de PolyPublie:**  
PolyPublie URL: <https://publications.polymtl.ca/9466/>

**Version:** Version officielle de l'éditeur / Published version  
Révisé par les pairs / Refereed

**Conditions d'utilisation:**  
Terms of Use: CC BY

 **Document publié chez l'éditeur officiel**  
Document issued by the official publisher

**Titre de la revue:** Journal of Surveillance, Security and Safety (vol. 2)  
Journal Title:

**Maison d'édition:** OAE Publishing Inc  
Publisher:

**URL officiel:** <https://doi.org/10.20517/jsss.2020.20>  
Official URL:

**Mention légale:** © The Author(s) 2020. Open Access This article is licensed under a Creative Commons  
Legal notice: Attribution 4.0 International License (<https://creativecommons.org/licenses/by/4.0/>),  
which permits unrestricted use, sharing, adaptation, distribution and reproduction in  
any medium or format, for any purpose, even commercially, as long as you give  
appropriate credit to the original author(s) and the source, provide a link to the Creative  
Commons license, and indicate if changes were made.

Original Article

Open Access



# Information theoretic-based privacy risk evaluation for data anonymization

Anis Bkakria<sup>1</sup>, Frédéric Cuppens<sup>1,2</sup>, Nora Cuppens<sup>1,2</sup>, Aimilia Tasidou<sup>1</sup>

<sup>1</sup>IMT Atlantique, Rennes 35510, France.

<sup>2</sup>Polytechnique Montréal, Montréal H3T 1J4, Canada.

**Correspondence to:** Dr. Anis Bkakria, SRCD department, IMT Atlantique, 2 rue de la châtaigneraie, Rennes 35510, France. E-mail: anis.bkakria@imt-atlantique.fr

**How to cite this article:** Bkakria A, Cuppens F, Cuppens N, Tasidou A. Information theoretic-based privacy risk evaluation for data anonymization. *J Surveill Secur Saf* 2021;2:83-102. <https://dx.doi.org/10.20517/jsss.2020.20>

**Received:** 1 Jun 2020 **First Decision:** 25 Aug 2020 **Revised:** 2 Sep 2020 **Accepted:** 2 Sep 2020 **Published:** 29 Aug 2021

**Academic Editor:** Xiaofeng Chen **Copy Editor:** Xi-Jun Chen **Production Editor:** Xi-Jun Chen

## Abstract

**Aim:** Data anonymization aims to enable data publishing without compromising the individuals' privacy. The re-identification and sensitive information inference risks of a dataset are important factors in the decision-making process for the techniques and the parameters of the anonymization process. If correctly assessed, measuring the re-identification and inference risks can help optimize the balance between protection and utility of the dataset, as too aggressive anonymization can render the data useless, while publishing data with a high risk of de-anonymization is troublesome.

**Methods:** In this paper, a new information theoretic-based privacy metric (ITPR) for assessing both the re-identification risk and sensitive information inference risk of datasets is proposed. We compare the proposed metric with existing information theoretic metrics and their ability to assess risk for various cases of dataset characteristics.

**Results:** We show that ITPR is the only metric that can effectively quantify both re-identification and sensitive information inference risks. We provide several experiments to illustrate the effectiveness of ITPR.

**Conclusion:** Unlike existing information theoretic-based privacy metrics, the ITPR metric we propose in this paper is, to the best of our knowledge, the first information theoretic-based privacy metric that allows correctly assessing both re-identification and sensitive information inference risks.



© The Author(s) 2020. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, sharing, adaptation, distribution and reproduction in any medium or format, for any purpose, even commercially, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.



**Keywords:** Data anonymization, identification risk, disclosure risk, information theoretic-based privacy metrics.

## 1 INTRODUCTION

Digital services today rely on the availability and processing of often sensitive data, while data analytics brings about important benefits both for the service providers and the individuals. However, the protection of sensitive data during this extensive processing has become a growing concern for Internet users, as the widespread use of IoT services, mobile devices, and location services lead to constant monitoring and vast amounts of sensitive data being gathered and stored.

There is a trade-off among data protection, data availability and data utility that needs to be tackled to achieve services which ensure privacy protection and produce usable results at the same time. This need becomes more imminent today, both due to the individuals' sensitization to data protection matters, which can lead to their lack of cooperation if they do not trust the service, and the data protection legislation being put in action, which holds the data handler responsible for any data breaches<sup>[1]</sup>.

Another important issue concerns privacy-preserving data publishing<sup>[2]</sup>. Anonymization techniques are being used to sanitize datasets prior to data publishing, so that similar data processing results are produced, while preserving data privacy. However, oftentimes the data characteristics are such that re-identification is easier than expected. De-anonymization techniques have advanced significantly and the abundance of external auxiliary information concerning individuals can lead to the de-anonymization of seemingly safe datasets<sup>[3]</sup>.

In this paper, we study and compare existing information theoretic-based privacy metrics to show that, in several cases, they cannot correctly assess both re-identification risk and sensitive information inference risk. We then propose, to the best of our knowledge, the first information theoretic-based privacy metric that can correctly assess both re-identification and sensitive information inference risks.

**Outline** The paper is organized as follows. In Section 2, we describe the problem statement and the motivation for this work. In Section 3, we present the background regarding anonymization risk metrics, as well as the related work. The proposed information theoretic-based privacy risk metric (ITPR) is presented in Section 4. In Section 5, we describe the produced experimental results, as well as compared them to the related work. The conclusions are presented in Section 6.

## 2 PROBLEM STATEMENT

Motivated by the desire to provide statistical information (sum, count, average, maximum, minimum, *etc.*) without disclosing sensitive information about individuals, a lot of research has been conducted in the area of statistical databases<sup>[4-8]</sup>. Fellegi and Sunter<sup>[7]</sup> proposed a Bayesian approach to estimate the advantage of an adversary to disclose identities in continuous data. Unfortunately, metrics that use the adversary's success probability to estimate the identification risk are known to be useful with averaged population of individuals. An individual may still suffer a high identification risk even when the adversary's success probability is low. In addition, the proposed metric is not useful for measuring the inference risk of sensitive attribute values. In<sup>[9]</sup>, Fellegi proposed an approach for checking residual disclosure. In contrast to our metric, the proposed approach is not generic since it considers only disclosure caused by the publication of counts or aggregate query results. Bethlehem *et al.*<sup>[5]</sup> presented a metric based on the Poisson–Gamma model for estimating the risk of identification using sample data. Unfortunately, the proposed model is not useful when it comes to measuring inference risk. One of the first research works that studied both the re-identification risk and the harm that can be caused by re-identification was conducted by Lambert<sup>[8]</sup>. The metric proposed to measure the harm

can be adapted to be used for measuring inference risk. In particular, the author introduced the definition of true and false disclosure and shows that harm is difficult to measure because false re-identifications often lead to false inferences which may have different harm. In<sup>[6]</sup>, Duncan and Lambert used Bayesian decision theory to propose a metric for assessing potential compromises on individual confidentiality due to the publication of aggregated statistical query results. They proposed a method for disclosure limiting where the disclosure is modeled as the difference between the prior and the posterior beliefs of the adversary about the data contents before and after the release of statistical query results.

Most of the privacy metrics proposed in the area of statistical databases (e.g.,<sup>[6,8,9]</sup>) address the problem of measuring the disclosure based on queried information. In addition, they are useful for measuring the disclosure risk that results from the matching of the results of specific statistical queries with other available information. They are often measured based on the information that the adversary can gather from these available data sources about the individuals in the released data, which we believe to be very difficult to determine in real-world situations. Instead, the privacy metric we propose in this paper, as well as the ones we consider for comparison, aim to measure the disclosure risk of individual-specific data on which users can produce summaries according to their own needs. In particular, the considered metrics make no assumption about the entity.

Privacy-preserving data publishing<sup>[2,10]</sup> enables the utilization of the data collected by organizations and service providers, without compromising the dataset participants' privacy. The goal is to release microdata from the dataset for processing, while protecting private information. To achieve that, the initial dataset is anonymized. Metrics such as k-anonymity<sup>[11]</sup>, l-diversity<sup>[12]</sup>, and t-closeness<sup>[13]</sup> have been proposed to assess the quality of the anonymization process, as well as the disclosure risk, setting thresholds for the anonymized dataset characteristics to be considered safe.

In<sup>[14]</sup>, Delanaux *et al.* proposed a declarative framework for privacy preserving linked data publishing in which privacy and utility policies are specified as SPARQL queries. The proposed framework allows determining the sequence of anonymization operations that should be used to satisfy the specified policy in the context of linked data. In the same context, Grau *et al.*<sup>[15]</sup> proposed a framework for ensuring that anonymized RDF graphs can be published on the Semantic Web with provable privacy guarantees.

It is common that the data publisher does not know in advance the entities that will access the released data and the operations that will be performed on the released data. Therefore, the data publisher needs a method to assess the characteristics of the dataset and the resulting disclosure risk, in order to make informed choices on what to include in the released dataset.

Tailoring the anonymization techniques and parameter setting to the dataset characteristics is important, so that the maximum possible data utility is preserved, while the data remain protected. If the anonymization process affects the data more than needed, then data utility diminishes, while lighter anonymization can lead to unintended data disclosure. Therefore, the challenge in the anonymization process is achieving a balance between protection and utility<sup>[16]</sup>.

While anonymized data utility measurement depends heavily on the type of analysis to be performed over the data, the measurement of the privacy protection level depends mainly on two factors: (1) the re-identification risk that measures the risk for any individual to be identified in the anonymized dataset; and (2) the inference risk that measures the risk for any sensitive information in the dataset to be linked to a specific individual. The level of privacy protection that can be ensured depends considerably on the distribution of the values of the different attributes that compose the dataset to be anonymized. Thus, depending on the considered dataset, we often end up with balancing the values of the re-identification and sensitive information inference risks to

find an acceptable trade-off between data utility and data privacy. However, in the literature, these risks are often measured using different metrics, e.g.,  $k$ -anonymity<sup>[11]</sup> for re-identification risk and  $l$ -diversity<sup>[12]</sup> and  $t$ -closeness<sup>[13]</sup> for inference risk. The different risk measurement methods and in some cases behaviors of the existing metrics often make combining the value of two risks to find the best balance between minimizing re-identification and sensitive information inference risks hard.

In this paper, we aim to provide a new information theoretic-based privacy metric that is able to assess both re-identification and sensitive information inference risks. We show that information theoretic-based risk metrics proposed in the literature, are mainly average values. As a result, they do not assess effectively the contribution of a single record to the risk value<sup>[17]</sup>. As individual risk values can fluctuate greatly (as we illustrate in Section 3.2), average values are not suitable to represent both the re-identification and the sensitive information inference risks.

### 3 BACKGROUND AND RELATED WORK

In this section, we present the relevant background regarding information theoretic anonymity metrics, as well as related work on re-identification risk metrics. The typical setting of an anonymization process involves data, contained in one or more tables. Each table row represents a data record and the columns of the table are the record attributes. These attributes are categorized into three main categories:

- Identifiers directly identify an individual, such as name, social security number, *etc.*
- Quasi-identifiers (key attributes) can be used in combination to identify an individual, such as age, gender, occupation, postal code, *etc.*
- Sensitive attributes contain sensitive information concerning an individual, such as health data, financial data, *etc.*

To protect individuals from the disclosure of their sensitive data, anonymization techniques can be employed, such as data generalization, suppression, and perturbation, as well as noise addition techniques. De-anonymization attacks on the released data can lead to both identity and attribute disclosure. In the case of identity disclosure, an individual is directly linked to a dataset record and the sensitive information it contains. In the case of attribute disclosure, the individual is associated with an attribute value but not a specific record.

Anonymity is defined as the state of being not identifiable within a set of subjects, called the anonymity set<sup>[18]</sup>. Statistical disclosure control (SDC) methods propose minimum requirements for each attribute in the dataset. To conform with  $k$ -anonymity<sup>[11]</sup>, it is required that all quasi-identifier groups contain at least  $k$  records. As the value of the quasi-identifier can be the same in the whole group,  $k$ -anonymity does not protect against homogeneity attacks<sup>[12]</sup>. For example, let us consider a 4-anonymity table composed of two attribute columns: Age and Disease. If we assume that all individuals having the value “4\*” for Age in the considered table are suffering from “HIV”, then, to perform an homogeneity attack, an adversary only needs to know that an individual present in the table is between 40 and 49 years old to know his/her disease.

To address this issue,  $l$ -diversity<sup>[12]</sup> has been proposed, as it requires each group to contain at least  $l$  distinct values for each quasi-identifier. Fulfilling  $l$ -diversity still fails to protect against skewness attacks<sup>[13]</sup>, which allow sensitive information disclosure when their distribution in the quasi-identifier group is significantly different from the corresponding distribution over the entire dataset. To deal with this issue,  $t$ -closeness<sup>[13]</sup> requires that the distance between the distributions of sensitive attributes in the quasi-identifier group and the whole dataset remains under  $t$ . However, although these methods provide an objective way of assessing and enforcing privacy in the datasets, they do not constitute a uniform risk-assessment metric.

**Table 1. First anonymized dataset ( $k_1 = 3, l_1 = 3, \text{ and } t_1 = 0$ )**

**(a) Initial dataset**

#	Zip Code	Age	Disease
1	35510	21	Asthma
2	35591	42	HIV
3	35593	47	Asthma
4	35210	38	Diabetes
5	35273	32	HIV
6	35517	20	Diabetes
7	35599	49	Diabetes
8	35262	33	Asthma
9	35511	26	HIV

**(b) Anonymized dataset**

#	Zip Code	Age	Disease
Group 1	3551*	2*	Asthma Diabetes HIV
Group 2	3559*	4*	HIV Asthma Diabetes
Group 3	352*	3*	Diabetes HIV Asthma

**Table 2. Second anonymized dataset ( $k_2 = 5, l_2 = 3, \text{ and } t_2 = 0$ )**

**(a) Initial dataset**

#	Zip Code	Age	Disease
1	35510	21	Asthma
2	35591	42	HIV
3	35593	47	Asthma
4	35210	38	Diabetes
5	35273	32	HIV
6	35517	20	Diabetes
7	35599	49	Diabetes
8	35262	33	Asthma
9	35511	26	HIV
10	35212	39	Diabetes
11	35281	32	Diabetes
12	35596	41	Diabetes
13	35592	46	Diabetes
14	35515	23	Diabetes
15	35511	26	Diabetes

**(b) Anonymized dataset**

#	Zip Code	Age	Disease
Group 1	3551*	2*	Asthma Diabetes Diabetes Diabetes HIV
Group 2	3559*	4*	HIV Asthma Diabetes Diabetes Diabetes
Group 3	352*	3*	Diabetes HIV Asthma Diabetes Diabetes

**3.1 The limitations of k-anonymity, l-diversity, and t-closeness models**

As shown in the previous section, k-anonymity was proposed to mitigate identity disclosure, l-diversity was proposed to mitigate homogeneity attacks, and t-closeness was proposed to prevent skewness attack. However, when we analyze carefully the three models, we realize that they are not useful for computing the effective inference risk (disclosure risk) of sensitive attributes information.

To illustrate, let us take the example of the two anonymized datasets in Tables 1 and 2. The anonymized dataset in Table 1b satisfies 3-anonymity, 3-diversity, and 0-closeness while the anonymized dataset in Table 2b satisfies 5-anonymity, 3-diversity, and 0-closeness. Thus, if we limit our analysis to the computed three values for k-anonymity, l-diversity, and t-closeness in these two cases, the overall level of ensured privacy is better in the second dataset since  $k_2 > k_1$  (i.e., the re-identification risk is lower in the second anonymized dataset than the first one),  $l_1 = l_2$ , and  $t_1 = t_2$ . However, if we look carefully at the distribution of the values of the attribute Disease in the two anonymized datasets, we realize that, if an adversary knows that an individual is in Group 1, 2, or 3, he has a bigger probability (0.60) of inferring the individual’s disease in the second anonymized dataset than in the first one (0.33). This example proves that the combination of the k-anonymity, l-diversity, and t-closeness models does not measure the effective disclosure risk but instead the accomplishment of the anonymization process.

To evaluate the performance of an anonymization method and to be able to compare the effectiveness among different methods, we need to define a common evaluation framework that can measure effectively both the re-identification risk and the sensitive attributes inference risk.

### 3.2 Information theoretic risk metrics

Information theory can be applied to the data protection context to evaluate the amount of information carried by a dataset and the possibility that disclosing these data leads to identity or attribute leakage. Information theoretic risk metrics provide the ability to be applied to different anonymity systems<sup>[19]</sup>. Information can be represented as a variable that can contain different values, and an information theoretic risk metric aims at measuring the amount of information leaked from a dataset.

Entropy is a key concept of information theory<sup>[20]</sup> that quantifies the uncertainty of a random variable. Uncertainty enhances privacy, as it hinders an adversary from effectively estimating attribute values<sup>[21]</sup>. In the following paragraphs, we provide definitions for the key concepts in entropy-based anonymity metrics.

We consider  $X$  and  $Y$  to be two random variables, corresponding to two attributes in a dataset.

Similar to uncertainty, information theory can be used to produce metrics that quantify information loss or gain for an adversary.

**Entropy** The entropy of a discrete random variable  $X$  is:

$$H(X) = - \sum_{x \in \mathcal{X}} p(x) \log p(x) \quad (1)$$

where  $p(x)$  is the probability of occurrence for value  $x \in X$ .

**Conditional Entropy** The conditional entropy of a discrete random variable  $X$ , given a discrete random variable  $Y$ , is:

$$\begin{aligned} H(X|Y) &= - \sum_y p(y) \sum_x p(x|y) \log p(x|y) \\ &= - \sum_{x,y} p(x,y) \log p(x|y) \end{aligned} \quad (2)$$

where  $p(x|y)$  is the conditional probability of occurrence for value  $x \in X$ , given the occurrence of  $y \in Y$ . Conditional entropy expresses how much information is needed to describe  $X$ , knowing the value of  $Y$ . The maximum of conditional entropy  $H(X|Y)$  is the entropy  $H(X)$ <sup>[21]</sup>. Therefore, normalized conditional entropy is computed by the following fraction:

$$\frac{H(X|Y)}{H(X)} \quad (3)$$

**Joint Entropy** The joint entropy of two discrete random variables  $X$  and  $Y$  is:

$$H(X, Y) = - \sum_{x \in \mathcal{X}} p(x, y) \log p(x, y) \quad (4)$$

where  $p(x, y)$  is the joint probability of occurrence for the value pair  $(x, y)$ .

### 3.3 Related Work

In this section, we present the related work on information theoretic-based privacy risk metrics. After presenting the metrics and their characteristics, we present some example cases which show that these metrics are unable to assess correctly the re-identification and inference risks of a dataset in certain cases.

### 3.3.1 Discrimination Rate

Discrimination Rate (DR) is an attribute-centric privacy metric, which aims to measure the degree to which attributes are able to refine an anonymity set and to measure their identification capability<sup>[22,23]</sup>. For this purpose, attributes are represented as discrete random variables. Considering two discrete random variables  $X$  and  $Y$ , DR is used to measure the identification capacity of attribute  $Y$  over the set of  $X$ .

$$DR_X(Y) = 1 - \frac{H(X|Y)}{H(X)} \quad (5)$$

DR is bounded on  $[0,1]$ , where 1 means that an identifier reduces the anonymity set to a single individual.

### 3.3.2 Mutual Information

Mutual Information (MI) has been proposed as a metric for the disclosure risk and the utility of a dataset<sup>[24–28]</sup>. The mutual information of two discrete random variables  $X$  and  $Y$  represents the average amount of knowledge about  $X$  gained by knowing  $Y$ , or alternatively the amount of shared information between  $X$  and  $Y$ . Therefore, mutual information is an information gain metric. Intuitively, if  $X$  and  $Y$  are independent, their mutual information is equal to zero. Mutual Information is computed as the difference between entropy  $H(X)$  and conditional entropy  $H(X|Y)$ :

$$\begin{aligned} I(X;Y) &= H(X) - H(X|Y) = H(Y) - H(Y|X) \\ &= H(X) + H(Y) - H(X,Y) \end{aligned} \quad (6)$$

### 3.3.3 Conditional Privacy

Conditional privacy (CP) is a privacy metric proposed in<sup>[29]</sup> for quantifying the fraction of privacy of a sensitive attribute  $X$ , which is lost by revealing another attribute  $Y$ . Conditional privacy can be seen as a mutual information normalization and is formalized as follows:

$$Priv_{CP} = 1 - 2^{-I(X;Y)} \quad (7)$$

where  $I(X;Y)$  represents the mutual information of  $X$  and  $Y$ .

### 3.3.4 Maximum Information Leakage

Maximum Information Leakage (MIL) is a modification of mutual information metric to consider only a single instance of a quasi-identifier attribute  $Y$ . It measures the maximum amount of information about a sensitive attribute  $X$  that can be learned by observing a single instance of  $Y$ <sup>[30]</sup>.

$$Priv_{MIL} = \max_{y \in Y} I(X;Y=y) \quad (8)$$

where  $I(X;Y)$  represents the mutual information of  $X$  and  $Y$ .

### 3.3.5 Entropy $l$ -Diversity

Entropy  $l$ -Diversity (ELD) is proposed as an instantiation of the  $l$ -diversity principle<sup>[12]</sup>. It states that a table is entropy  $l$ -diverse for a sensitive attribute  $X$  if the following condition holds for all quasi-identifier groups  $q$ :

$$-\sum_{x \in X} p(q,x) \cdot \log(p(q,x)) \geq \log(l)$$

where  $p(q,x)$  is the fraction of records in  $q$  that have the value  $x$  from the attribute  $X$ . We note that ELD, as proposed in<sup>[12]</sup>, does not allow measuring the inference risk of a sensitive attribute, but it is used as an entropy-based condition that must be satisfied to achieve the  $l$ -diversity property for a sensitive attribute. We adapt ELD



**Table 3. Example datasets**

Identifier	Cases of <b>Age</b> Attribute					Cases of <b>Disease</b> Attribute		
	Case 1	Case 2	Case 3	Case 4	Case 5	Case 1	Case 2	Case 3
#1	30	30	30	30	30	Diabetes	Diabetes	Diabetes
#2	62	30	30	30	30	Epilepsy	Diabetes	Diabetes
#3	37	30	30	30	30	Asthma	Epilepsy	Diabetes
#4	21	30	47	47	30	Allergies	Depression	Diabetes
#5	19	30	30	30	47	Depression	HIV	Diabetes
#6	47	30	30	47	47	HIV	Heart Disease	HIV
#7	71	30	30	30	47	Heart Disease	Cancer	Asthma
#8	73	30	30	30	47	Cancer	Allergies	Allergies

**Table 4. Information theoretic-based privacy metrics analysis results**

Privacy Risk Type	Considered Dataset	Results				
		DR	MI	CP	MIL	ELD
Re-identification risk	Identifier + Case 1 of Age	1.0	3.0	0,875	3.0	1.0
	Identifier + Case 2 of Age	0.0	0,0	0,0	0,0	0,125
	Identifier + Case 3 of Age	0,18	0,54	0,31	3,0	1,0
	Identifier + Case 4 of Age	0,27	0,81	0,43	2,75	0,5
	Identifier + Case 5 of Age	0,33	1,0	0,5	2,0	0,25
Inference risk	Identifier + Case 5 of Age + Case 1 of Disease	0,33	1,0	0,5	2,0	0,25
	Identifier + Case 5 of Age + Case 2 of Disease	0,36	1,0	0,5	1,0	0,35
	Identifier + Case 5 of Age + Case 3 of Disease	0,35	0,54	0,31	1,0	1,0

to allow inference risk measurement as following: Given a sensitive attribute  $X$  and the set of quasi-identifiers groups  $Q$ , the inference risk of the attribute  $X$  can be measured as:

$$pELD(X) = 2^{\min_{q \in Q} \sum_{x \in X} p(q,x) \cdot \log(p(q,x))}$$

#### Examples - Problem cases

As mentioned above, information theoretic-based privacy metrics do not always succeed in assessing the disclosure risk, since they do not effectively assess the contribution of individual records to the risk value. We present in Table 3 some examples of datasets, where the considered metrics fail to express one or both of the enhanced re-identification and inference risks of the data. As shown in Table 3, we consider a dataset composed of eight records containing values of three attributes: Identifier, Age, and Disease. Five cases are considered for the Age attribute:

- Case 1: All records contain unique attribute values.
- Case 2: All records contain the same attribute value.
- Case 3: Only one record contains a different attribute value.
- Case 4: Only two records contain a different attribute value.
- Case 5: Half of the records contain one and half of the records contain another attribute value.

For the Disease attribute, we consider the first three cases considered above for the attribute Age.

For Cases 1 and 3 of the attribute Age, we expect the highest value for the re-identification risk since some values of the attribute Age uniquely identify an individual. For Case 2, we expect the lowest value for the re-identification risk since the knowledge of the value of the attribute Age gives no additional information about the considered identify. For Case 4, we expect a lower value for the re-identification risk than the one assigned to Case 3 and a higher value than the one assigned to Case 5.

For the inference risk, we expect the risk assigned to Case 1 (Identifier + Case 5 of Age + Case 1 of Disease) to be lower that the ones assigned to Cases 2 (Identifier + Case 5 of Age + Case 2 of Disease) and 3 (Identifier +

Case 5 of Age + Case 3 of Disease). In addition, we expect the risk to be assigned to Case 2 to be lowest that the one assigned to Case 3.

In Table 4, we examine the behavior of the information theoretic-based privacy metrics previously presented in Sections 3.3.1, 3.3.2, 3.3.3, and 3.3.4, depending on the distribution of the values in the considered attributes cases. According to the obtained results, DR, MI, and MIL correctly express the re-identification risk when Cases 1 and 2 of the attribute Age are considered (Rows 2 and 3 of Table 4). However, all three metrics fail to reflect the level of re-identification risk for Cases 3–5 of the attribute Age (Rows 4–6 of Table 4). Out of these previously mentioned three cases, DR, MI, and MIL metrics output lower values to the case that represents a higher re-identification risk (Case 3 of the attribute Age) and higher values to the case that represents a lower re-identification risk (Case 5 of attribute Age). As for CP metric, it correctly reflects the re-identification risk of the dataset instance in which Case 2 of the attribute Age is considered and fails to correctly reflect the re-identification risk for other cases. The HLD metric seems to correctly reflect both re-identification and inference risks for all considered cases. However, we show in Section 5 that HLD fails to measure correctly the inference risk caused by the difference between the distribution of the values of the sensitive attribute in the whole table and their distribution in the different quasi-identifier groups.

When it comes to measuring the inference risk, all considered metrics successfully reflect the re-identification risk of the dataset instance in which Case 5 of the attribute Age and Case 1 of the attribute Disease are considered (Row 6 of Table 4) and fail to correctly reflect the re-identification risk for other cases (Rows 7 and 8 of Table 4).

#### 4 THE NEW INFORMATION THEORETIC-BASED PRIVACY RISK METRIC

To address the lack of ability of existing information theoretic based privacy metrics to effectively assess the contribution of individual records of a dataset to the re-identification risk value and correctly quantify the inference risk that stems from the correlation between a quasi-identifier attribute (e.g., Age) and a sensitive attribute (e.g., Disease), we propose a new information theoretic-based privacy risk metric (ITPR). The value of ITPR can effectively express, on the one side, the probability of the attacker to refine the anonymity set and re-identify a dataset participant, based on the knowledge of an (quasi-identifier) attribute and, on the other side, the probability of an adversary to refine the anonymity set and link an identity to a value of a sensitive attribute.

To develop the formula for the ITPR metric, we follow a similar logic as in the discrimination rate and mutual information metrics, still relying on information theory and entropy calculations. However, instead of using the average value of the attribute values' entropy, we take the maximum value of entropy among attribute values.

To compute the remaining identification information of attribute X, given attribute Y, we compute  $H(X) - H(X|Y)$ . We then divide by  $H(X)$  to normalize the computed value, resulting in the following representation:

$$\max_{y \in \Omega_Y} \left( \left\{ 1 - \frac{\tilde{H}(X|Y = y)}{H(X)} \right\} \right) \quad (9)$$

where  $\Omega_Y$  is the sample space of the discrete random variable Y and  $\tilde{H}(X|Y = y) = p(Y = y) \times H(X|Y = y)$ .

Using this equation, the results produced depend on the number of distinct values for Y; thus, for example, for the case of two distinct values ( $|\Omega_Y| = 2$ , e.g. Case 5 of attribute Age in Table 3), the produced results span between 0.5 and 1.0, in the case of three distinct values ( $|\Omega_Y| = 3$ ), ITPR values span between 0.66 and 1.0, and

so on. To counteract this behavior, we introduce the number of distinct values in attribute  $Y$  as a parameter in the ITPR metric, leading to the following definition.

**Definition 1.** (Simple ITPR) Given two attributes  $X$  and  $Y$  of a dataset, the simple ITPR of attribute  $Y$  relative to attribute  $X$  quantifies the capacity of attribute  $Y$  to refine the set of values of attribute  $X$  and is measured as follows:

$$ITPR_X(Y) = \max_{y \in \Omega_Y} \left( \left\{ 1 - \frac{|\Omega_Y| \times \tilde{H}(X|Y=y)}{H(X)} \right\} \right) \quad (10)$$

where  $|\Omega_Y|$  denotes the number of different values of  $Y$ .

Definition 1 can be generalized to define combined ITPR, which quantifies the ITPR measure related to the combination of the values of several attributes to perform re-identification and/or inference attacks.

**Definition 2.** (Combined ITPR) Given a set of attributes  $X, Y_1, Y_2, \dots, Y_n$  of a dataset  $\mathcal{D}$ , let us denote  $\mathcal{T}$  as the set of  $\langle Y_1, Y_2, \dots, Y_n \rangle$  distinct tuples in  $\mathcal{D}$ . The combined ITPR of attributes  $Y_1, Y_2, \dots, Y_n$  relative to attribute  $X$  quantifies the capacity of attributes  $Y_1, Y_2, \dots, Y_n$  to refine the set of values of attribute  $X$  and is computed according to the following formula:

$$ITPR_X(Y_1, Y_2, \dots, Y_n) = \max_{\langle y_1, y_2, \dots, y_n \rangle \in \mathcal{T}} \left( \left\{ 1 - \frac{|\mathcal{T}| \times \tilde{H}(X|Y_1=y_1, Y_2=y_2, \dots, Y_n=y_n)}{H(X)} \right\} \right)$$

where

$$\tilde{H}(X|Y_1=y_1, Y_2=y_2, \dots, Y_n=y_n) = p \left( \bigcap_{i=1}^n Y_i = y_i \right) \times H(X|Y_1=y_1, Y_2=y_2, \dots, Y_n=y_n)$$

**Proposition 1.** The output of both the simple and combined ITPR is bounded by 0 and 1.

*Proof.* We start by proving that the output of the simple ITPR (Definition 1) is bounded by 0 and 1. The proof is by contradiction. Let us suppose that  $\forall X, \forall Y, ITPR_X(Y) < 0$ . We get

$$\begin{aligned} & \max_{y \in \Omega_Y} \left( \left\{ 1 - \frac{|\Omega_Y| \times \tilde{H}(X|Y=y)}{H(X)} \right\} \right) < 0 \\ \iff & \forall y \in \Omega_Y : 1 - \frac{|\Omega_Y| \times \tilde{H}(X|Y=y)}{H(X)} < 0 \\ \iff & \forall y \in \Omega_Y : |\Omega_Y| \times \tilde{H}(X|Y=y) > H(X) \\ \iff & \forall y \in \Omega_Y : |\Omega_Y| \times p(Y=y) \times H(X|Y=y) > H(X) \\ \iff & |\Omega_Y| \times \sum_{y \in \Omega_Y} p(y) \times H(X|Y=y) > \sum_{y \in \Omega_Y} H(X) \\ \iff & |\Omega_Y| \times \sum_{y \in \Omega_Y} p(y) \times H(X|Y=y) > H(X) \times |\Omega_Y| \\ \iff & \sum_{y \in \Omega_Y} p(y) \times H(X|Y=y) > H(X) \\ \iff & H(X|Y) > H(X) \end{aligned} \quad (11)$$

which cannot be true for any  $X$  and  $Y$ . Therefore, for all  $X$  and  $Y$ ,  $ITPR_X(Y) \geq 0$ .

**Table 5. ITPR results**

#	Privacy Risk Type	Considered Dataset	ITPR Results
1	Re-identification risk	Identifier + Case 1 of Age	1.0
2		Identifier + Case 2 of Age	0.0
3		Identifier + Case 3 of Age	1.0
4		Identifier + Case 4 of Age	0.83
5		Identifier + Case 5 of Age	0.33
6		Identifier + Case 2 of Age + Case 1 of Zip Code (Table 6)	0.6
7		Identifier + Case 2 of Age + Case 1 of Zip Code (Table 6)	0.75
8	Inference risk	Identifier + Case 5 of Age + Case 1 of Disease	0.33
9		Identifier + Case 5 of Age + Case 2 of Disease	0.45
10		Identifier + Case 5 of Age + Case 3 of Disease	1.0

**Table 6. Cases of the attribute Zip Code**

Case 1	Case 2
35000	35000
35000	35000
35000	35000
35510	35510
35510	35510
35510	35510
35510	35200
35510	35200

Now, let us suppose that  $\forall X, \forall Y, ITPR_X(Y) > 1$ . We get

$$\begin{aligned}
 & \max_{y \in \Omega_Y} \left( \left\{ 1 - \frac{|\Omega_Y| \times \tilde{H}(X|Y=y)}{H(X)} \right\} \right) > 1 \\
 \iff & \exists y \in \Omega_Y : 1 - \frac{|\Omega_Y| \times \tilde{H}(X|Y=y)}{H(X)} > 1 \tag{12} \\
 \iff & \exists y \in \Omega_Y : \frac{\tilde{H}(X|Y=y)}{H(X)} < 0
 \end{aligned}$$

which cannot be true for all  $X, Y$ , and  $y \in Y$ . Therefore, for all  $X$  and  $Y, ITPR_X(Y) \leq 1$ . We can use the same approach to prove that for all  $X, Y_1, \dots, Y_n, 0 \leq ITPR_X(Y_1, Y_2, \dots, Y_n) \leq 1$ .  $\square$

As stated in Proposition 1, the returned results for ITPR are normalized between values 0 and 1 and effectively represent, as we illustrate below: (1) the re-identification risk of an (identifier) attribute  $X$ , given a (quasi-identifier) attribute  $Y$ ; and (2) the inference risk caused by a sensitive attribute  $X$  when a (quasi-identifier) attribute  $Y$  is published.

Table 5 illustrates the expressiveness of the ITPR metric, using the same examples as used in Table 4.

As shown in Table 5, compared to the existing information theoretic-based metrics we analyzed and reported in Table 4, ITPR correctly quantifies the re-identification risk for all the considered cases (Rows 1–5 of Table 5). Rows 6 and 7 of Table 5 show that ITPR can effectively measure the re-identification risk when two (quasi-identifier) attributes are combined. Moreover, the results show that the ITPR metric correctly measures the inference risk represented by the attribute Disease (Rows 8–10 of Table 5).

The behavior of the ITPR privacy metric regarding the distribution of considered attribute values is more thoroughly tested and illustrated in the next section.

## 5 EXPERIMENTAL RESULTS

In this section, we provide the experimental results for the ITPR metric, compared to the discrimination rate, mutual information, conditional privacy, maximum information leakage, and entropy l-diversity metrics. To assess the behavior of the functions of these metrics, we first calculated the metric values for a dataset of 10,000 records, containing two distinct  $Y$  attribute values (for example,  $y_1 = Male$  and  $y_2 = Female$ ). We denote by  $\epsilon$  the maximum difference between the number of occurrences of the values of the attribute  $Y$ :

$$\epsilon = \max_{y_1, y_2 \in \Omega_Y} (|occ(y_1) - occ(y_2)|) \quad (13)$$

where  $\Omega_Y$  denotes the set of distinct values of  $Y$  and  $occ()$  is a function that returns the number of occurrences of a value. Obviously, the bigger the value of  $\epsilon$  is, the smaller the number of occurrences of  $y_1$  or  $y_2$  will be, resulting in a high re-identification risk.

The results are illustrated in Figure 1. One can observe that the value of ITPR begins from 1 for the case of  $\epsilon = 9998$  (e.g.,  $|y_1| = 1$ ,  $|y_2| = 9999$ ) and diminishes smoothly while the value  $\epsilon$  decreases (i.e., the sizes of the two value groups ( $|Y_1|$ ,  $|Y_2|$ ) move closer to each other), converging to a very small value ( $7 * 10^{-2}$ ) when the value of  $\epsilon = 0$  (i.e., the two group sizes are equal  $|Y_1| = |Y_2| = 5000$ ).

In the DR case, the metric stays below 0.1 for this dataset, failing to accurately express the re-identification risk for the different cases. In the MI and CP cases, the metric output appears to increase as the number of records of each attribute value move towards being equal, failing also to express that the re-identification risk is higher when a smaller number of records contains one of the attribute values and the majority contains the other. The ELD output increases extremely slowly as the value of  $\epsilon$  increases. In Figure 1b, we observe that the MIL metric output increases as the value of  $\epsilon$  increases, which represents a correct behavior regarding the re-identification risk represented in the different considered cases. Unfortunately, the MIL metric suffers from two drawbacks: (1) the wide range of output values (e.g., between 7 and 13) makes the interpretation of the output of the metric difficult; and (2) the MIL metric does not correctly express the inference risk represented by a sensitive attribute, as illustrated in Figure 3b.

We note that, for all studied metrics, the same behaviors can be observed when several values are considered for the attribute  $Y$ , as described in Figure 2. As the results indicate, ITPR is able to effectively express:

- the lower existence of risk when the attribute values are distributed equally among the dataset records;
- the gradual enhancement of risk, as the number of records containing a certain value decreases; and
- the higher risk value when a certain value appears only in a small number of records in the dataset.

Furthermore, we compared the ability of the considered metrics to assess the inference risk represented by the publication of a sensitive attribute. For this, we consider two attributes Age and Disease in a dataset composed of 10,000 records. For simplicity, we consider the Age to be composed of five different values uniformly distributed over the 10,000 dataset records and the attribute Disease to be composed of instances of 10 different values. Since the inference risk depends mainly on the distribution of the different values of the considered sensitive attribute, we analyze the output of the considered metrics regarding the difference between the most used and the least used values of the attribute Disease that we denote  $\lambda$ .

$$\lambda = \max_{x \in \Omega_A, y_1, y_2 \in \Omega_D} (|occ(y_1|Age = x) - occ(y_2|Age = x)|) \quad (14)$$

where  $\Omega_A$  and  $\Omega_D$  denote the set of distinct values of attributes Age and Disease, respectively, and  $occ(y|Age = x)$  denotes the number of occurrences of  $y$  in the dataset when the value of  $Age = x$ . Note that the inference risk is expected to increase as the value of  $\lambda$  increases since the higher is the value of  $\lambda$  the higher is the number of occurrences of a specific value of a sensitive attribute in an anonymity class.

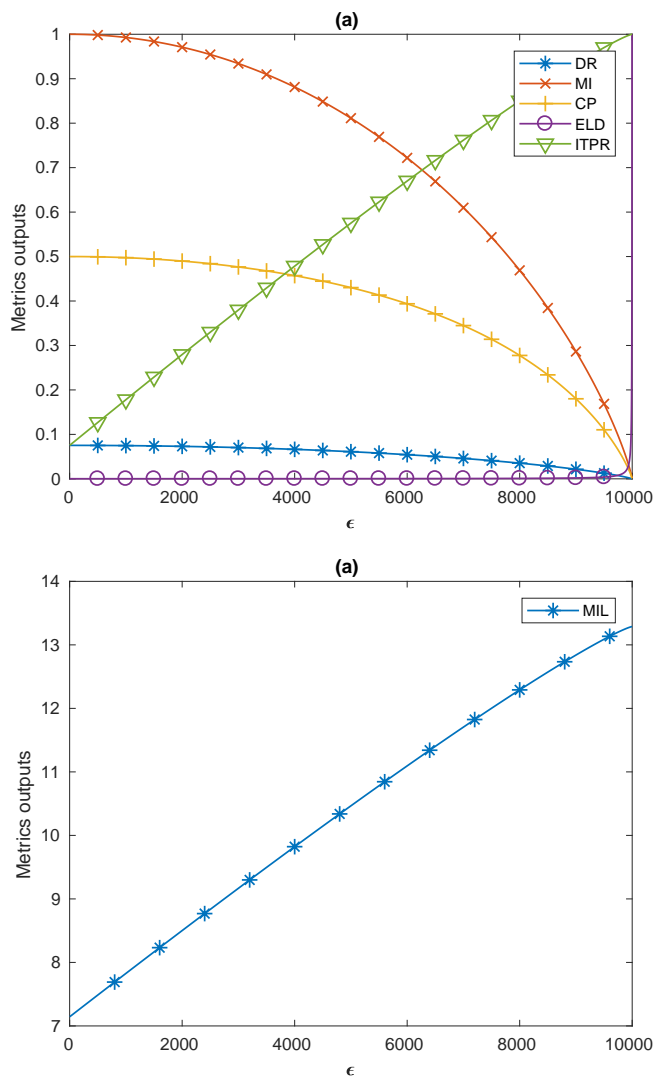
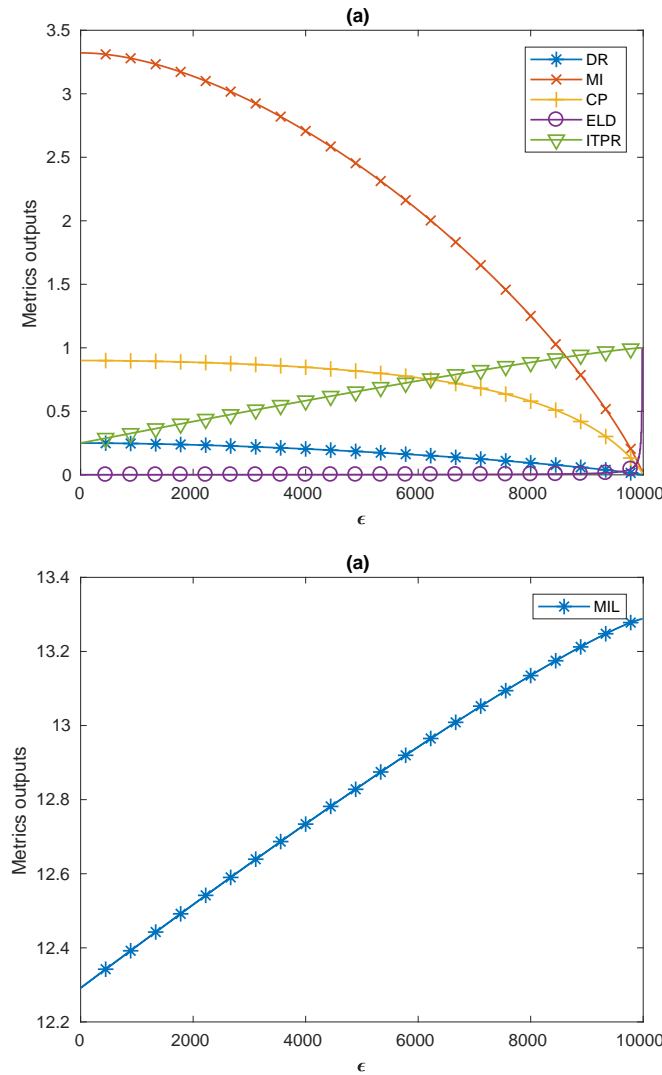


Figure 1. Comparison of ITPR with DR, MI, CP, ELD, and MIL for re-identification risk quantification (two different values for  $\lambda$ ).

Figure 3 illustrates the obtained results. Figure 3a shows that the outputs of ITPR, DR, MI, ELD, and CP increase according to the growth of the value of  $\lambda$ , which is compatible with the behavior we previously expected. For the three metrics DR, MI, and CP, the particular output ranges make the interpretation of the risk difficult since without knowing the output of the metric in the worst case it is hard to evaluate the severity of the output value. The ITPR and ELD metrics do not suffer from this limitation since their output always ranges between 0 and 1. As for the MIL metric, Figure 3b shows that its value decreases according to the growth of the value of  $\lambda$ , failing to effectively assess the inference risk represented by the Disease attribute.

Li et al. [13] showed that the inference risk does not depend only on the distribution of the values of the considered sensitive attribute in the quasi-identifier groups. The variation between the global distribution of the values of the sensitive attribute in the considered dataset and the local distribution of the values of the sensitive attribute in the quasi-identifier groups can drastically impact the inference risk. To illustrate, let us consider a dataset that has only one sensitive attribute, Disease, and is composed of  $10^8$  records. Furthermore, suppose that each record in the dataset is associated with a different individual and that only 1000 records contain “HIV” as a value for the attribute Disease. This means that anyone in the considered dataset has  $10^{-3}\%$  possibility of



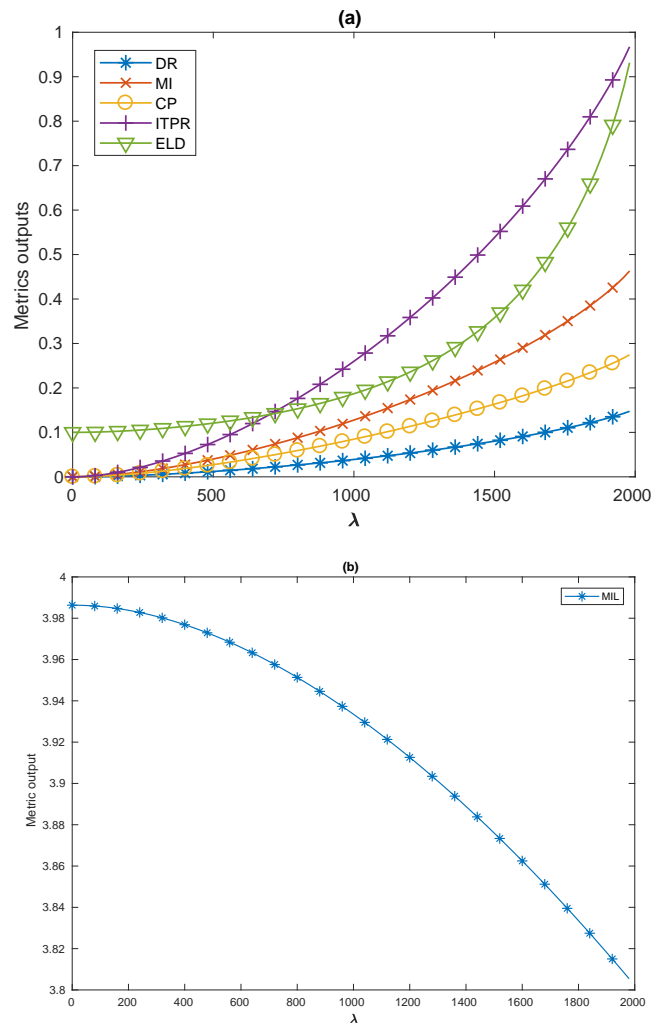
**Figure 2.** Comparison of ITPR with DR, MI, CP, ELD, and MIL for re-identification risk quantification (10 different values for Y).

having “HIV”. Now, let us suppose that one of the quasi-identifiers groups created by the used anonymization mechanism contains five records out of 100 that have “HIV” as a value for the attribute Disease. Clearly, this presents a serious privacy risk, because anyone in the considered quasi-identifiers group would be considered to have 5% possibility of having “HIV”, compared to the  $10^{-3}\%$  of the overall population. Thus, a correct inference risk measurement associated with a sensitive attribute  $X$  should take into consideration the variation  $\vartheta_X$  between the global distribution of the values of  $X$  in the considered dataset and the local distribution of the values  $X$  in the quasi-identifier groups. In fact, the higher this variation is, the higher the inference risk associated to the considered sensitive attribute must be. The variation  $\vartheta_X$  is formalized as follows:

$$\vartheta_X = \max_{q \in Q} (H(X) - H_q(X)) \tag{15}$$

where  $Q$  denotes the set of quasi-identifier groups,  $H(X)$  denotes the entropy of  $X$  in the hole dataset, and  $H_q(X)$  denotes the entropy of  $X$  in the quasi-identifier group  $q$ .

We compare the ability of the considered metrics for assessing the inference risk represented by the sensitive



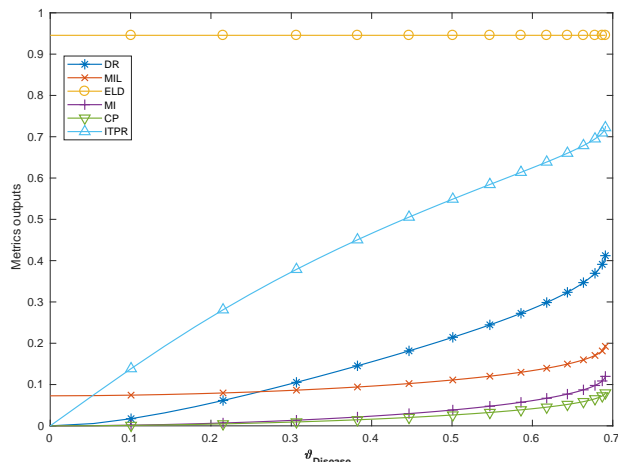
**Figure 3.** Comparison of ITPR with DR, MI, MIL, and ELD for inference risk output regarding  $\lambda$ .

attribute *Disease* regarding the variation  $\vartheta_{Disease}$  between the global distribution of its values in the considered dataset and the local distribution of its values in the quasi-identifier groups. Note that the inference risk of a sensitive attribute  $X$  is expected to increase according to raise of the variation  $\vartheta_X$ . The obtained results are illustrated in Figure 4. We can observe that the ELD metric does not take into consideration the variation  $\vartheta_{Disease}$  since its output is constant in the function of  $\vartheta_{Disease}$ . This experimentally proves that the ELD metric does not measure correctly the inference risk of the attribute Disease.

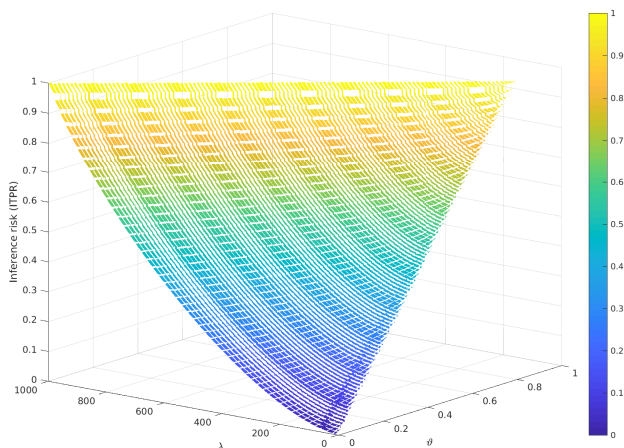
Moreover, we studied the behavior of the ITPR for different values of  $\epsilon(0 - 1200)$  and  $\vartheta(0 - 0.9)$ . In this experiment, we considered a dataset with 10,000 records composed of 10 different values for the attribute Age and 20 different values for the attribute Disease. The result is depicted in Figure 5.

It shows that the ITPR metric output increases smoothly while the value of  $\epsilon$  or the value of  $\vartheta$  increases which represents the expected behavior of a correct inference risk assessment metric. As for the MIL, MI, and CP metrics, their outputs increase extremely slowly as the value of  $\vartheta_{Disease}$  increases. For example, if we suppose that  $\vartheta_{Disease} = 0$  represents the case described above in which anyone in the database has  $10^{-3}\%$  possibility of having HIV,  $\vartheta_{Disease} = 0.4$  can represent the case in which anyone in a specific quasi-identifier group has





**Figure 4.** Comparison of ITPR with DR, MI, MIL, and ELD for inference risk output regarding  $\vartheta$ .



**Figure 5.** Inference risk assessment using ITPR.

12% possibility of having HIV. However, when we examine the variations of the outputs of MIL, MI, and CP between  $\vartheta_{Disease} = 0$  and  $\vartheta_{Disease} = 0.4$ , they increase only from 0.07 to 0.09, from 0 to 0.02, and from 0 to 0.01, respectively. Finally, this experiment shows that our proposed ITPR metric has the best behavior compared to considered metrics regarding the increase of the variation  $\vartheta_{Disease}$ .

Anonymization processes often deal with a large volume of data. As a result, the computation effectiveness of such a metric should be evaluated. For this, we considered a table composed of three attributes: Identifier, Age, and Disease. The attribute Age contains 120 different values while the attribute Disease contains 100 different values. The evaluation was performed on a Spark cluster of four nodes with 100 workers with one core and 1 GB per worker. Figure 6 shows the time needed for the computation of the ITPR metric regarding the number of rows in the considered table.

Finally, we evaluated the computation effectiveness of the ITPR metrics regarding the number  $n$  of the considered quasi-identifier attributes (Definition 2). For this, we considered a table composed of  $10^7$  records and

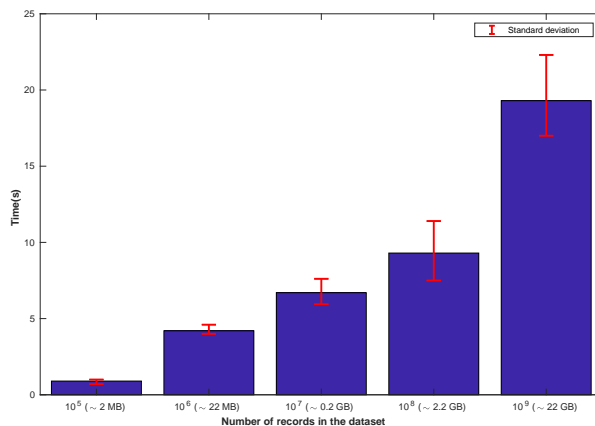


Figure 6. ITPR computation time per number of records.

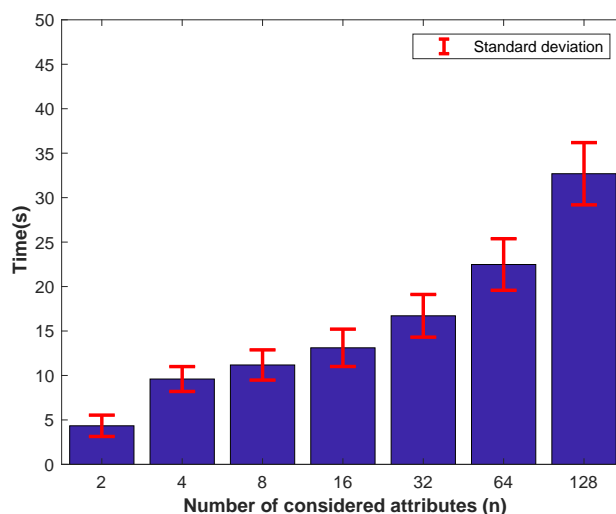


Figure 7. ITPR computation time as function of the number of considered quasi-identifier attributes  $n$ .

130 attributes. The obtained results are reported in Figure 7. It shows that the time required for computing ITPR increases linearly with the number of considered quasi-identifiers, which demonstrates ITPR scalability when dealing with wide flattened tables.

The proposed ITPR metric can be used in both unprocessed and anonymized datasets. In the case of raw datasets, ITPR can assist the data owners with making decisions on which attributes of the dataset are more sensitive (i.e., the ones that have high inference risks) and which ones should be included in the anonymized dataset and with what anonymization parameters.

As with all information theoretic risk metrics, defining the threshold values on what constitutes low, moderate, or high risk, remains an issue. The selection of such thresholds always depends on the characteristics of the dataset and the objectives of the data owner; however, in the case of ITPR, selecting thresholds for the risk values appears to be more straightforward, as the ITPR value increases gradually as the number of records containing a certain value decreases.

**Table 7. Benefits of the ITPR compared to existing metrics**

Privacy Metrics	Re-identification risk	Inference Risk
Discrimination Rate (DR)	X	(✓)
Mutual Information (MI)	X	(✓)
Conditional Privacy (CP)	X	(✓)
Maximum Information Leakage (MIL)	(✓)	X
Entropy L-Diversity (ELD)	(✓)	X
<b>ITPR</b>	✓	✓

(✓): Correct assessment but hard interpretation of the risk. The maximum output of the metric (i.e., the highest risk value) is not always the same, and it depends on the analyzed data. This makes the interpretation of the risk level hard.

Another interesting use of ITPR can be the quality control and cleaning of datasets, as datasets which contain errors in their records or contain a few outlier values will produce high values of re-identification risk, allowing data owners to clean errors in their datasets or decide to suppress outlier values to facilitate the successful dataset anonymization process.

## 6 CONCLUSIONS

In this paper, the ITPR information theoretic-based metric is proposed, for assessing both re-identification and inference risks within datasets. This metric aims at effectively representing the contribution of individual records of a dataset to the re-identification and inference risks value. To achieve that, ITPR takes into account the maximum value of entropy among the dataset attribute values. To facilitate the comparison of risk values among different anonymization processes and between different datasets, the ITPR value is normalized and bounded between 0 and 1. The experimental results show that ITPR succeeds in expressing both re-identification and inference risks. The comparison with existing Information theoretic-based privacy metrics (Table 7) shows that ITPR is the only metric that can effectively assess both re-identification and inference risk.

## 7 DECLARATIONS

### 7.1 Acknowledgments

The authors would like to thank the anonymous reviewers for their careful reading of our manuscript and their useful comments and suggestions.

### 7.2 Authors' contributions

Made substantial contributions to conception and design of the study and performed data analysis and interpretation: Bkakria A., Cuppens F, Cuppens N., Tasidou A.

### 7.3 Availability of data and materials

The data that used to validate the utility of the proposed metric contain real medical records and therefore cannot be published. The implementation of the ITPR metric we are proposing in this paper can be found at <https://github.com/nserser/ITPR>.

### 7.4 Financial support and sponsorship

None.

### 7.5 Conflicts of interest

All authors declared that there are no conflicts of interest.

### 7.6 Ethical approval and consent to participate

Not applicable.

## 7.7 Consent for publication

Not applicable.

## 7.8 Copyright

© The Author(s) 2021.

## REFERENCES

1. Regulation GDP. Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46. Official Journal of the European Union (OJ) 2016;59:294.
2. Fung BCM, Wang K, Chen R, Yu PS. Privacy-preserving Data Publishing: A Survey of Recent Developments. *ACM Comput Surv* 2010;42:14:1–4:53.
3. Narayanan A, Shmatikov V. Robust de-anonymization of large sparse datasets. In: *Security and Privacy, 2008. SP 2008. IEEE Symposium on. IEEE*; 2008. pp. 111–25.
4. Adam NR, Worthmann JC. Security-control methods for statistical databases: a comparative study. *ACM Computing Surveys (CSUR)* 1989;21:515–56.
5. Bethlehem JG, Keller WJ, Pannekoek J. Disclosure control of microdata. *Journal of the American Statistical Association* 1990;85:38–45.
6. Duncan GT, Lambert D. Disclosure-limited data dissemination. *Journal of the American statistical association* 1986;81:10–18.
7. Fellegi IP, Sunter AB. A theory for record linkage. *Journal of the American Statistical Association* 1969;64:1183–210.
8. Lambert D. Measures of disclosure risk and harm. *JOURNAL OF OFFICIAL STATISTICS-STOCKHOLM-* 1993;9:313–13.
9. Fellegi IP. On the question of statistical confidentiality. *Journal of the American Statistical Association* 1972;67:7–18.
10. Livraga G. Privacy in microdata release: Challenges, techniques, and approaches. In: *Data-Driven Policy Impact Evaluation*. Springer; 2019. pp. 67–83.
11. Sweeney L. k-anonymity: A model for protecting privacy. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems* 2002;10:557–70.
12. Machanavajjhala A, Kifer D, Gehrke J, Venkatasubramanian M. L-diversity: Privacy Beyond K-anonymity. *ACM Trans Knowl Discov Data* 2007 Mar;1. Available from: <http://doi.acm.org/10.1145/1217299.1217302>. [DOI: 10.1145/1217299.1217302]
13. Li N, Li T, Venkatasubramanian S. t-Closeness: Privacy Beyond k-Anonymity and l-Diversity. In: *2007 IEEE 23rd International Conference on Data Engineering*; 2007. pp. 106–15. [DOI: 10.1109/ICDE.2007.367856]
14. Delanaux R, Bonifati A, Rousset M, Thion R. Query-Based Linked Data Anonymization. In: *Vrandečić D, Bontcheva K, Suárez-Figueroa MC, Presutti V, Celino I, et al., editors. The Semantic Web - ISWC 2018 - 17th International Semantic Web Conference, Monterey, CA, USA, October 8-12, 2018, Proceedings, Part I. vol. 11136 of Lecture Notes in Computer Science. Springer; 2018. pp. 530–46. Available from: [https://doi.org/10.1007/978-3-030-00671-6\\_31](https://doi.org/10.1007/978-3-030-00671-6_31). [DOI: 10.1007/978-3-030-00671-6\_31]*
15. Grau BC, Kostylev EV. Logical Foundations of Linked Data Anonymisation. *J Artif Intell Res* 2019;64:253–314. Available from: <https://doi.org/10.1613/jair.1.11355>. [DOI: 10.1613/jair.1.11355]
16. Li T, Li N. On the tradeoff between privacy and utility in data publishing. In: *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining. ACM*; 2009. pp. 517–26.
17. Bezzi M. An information theoretic approach for privacy metrics. *Trans Data Privacy* 2010;3:199–215.
18. Pfitzmann A, Köhntopp M. Anonymity, unobservability, and pseudonymity—a proposal for terminology. In: *Designing privacy enhancing technologies*. Springer; 2001. pp. 1–9.
19. Diaz C. Anonymity metrics revisited. In: *Dagstuhl Seminar Proceedings. Schloss Dagstuhl-Leibniz-Zentrum für Informatik*; 2006. .
20. Shannon CE. A mathematical theory of communication. *Bell system technical journal* 1948;27:379–423.
21. Wagner I, Eckhoff D. Technical privacy metrics: a systematic survey. *ACM Computing Surveys (CSUR)* 2018;51:57.
22. Sondeck LP, Laurent M, Frey V. Discrimination rate: an attribute-centric metric to measure privacy. *Annales des Télécommunications* 2017;72:755–66. Available from: <https://doi.org/10.1007/s12243-017-0581-8>. [DOI: 10.1007/s12243-017-0581-8]
23. Sondeck LP, Laurent M, Frey V. The Semantic Discrimination Rate Metric for Privacy Measurements which Questions the Benefit of t-closeness over l-diversity. In: *Proceedings of the 14th International Joint Conference on e-Business and Telecommunications (ICETE 2017) - Volume 4: SECRYPT, Madrid, Spain, July 24-26, 2017.*; 2017. pp. 285–94. Available from: <https://doi.org/10.5220/0006418002850294>. [DOI: 10.5220/0006418002850294]
24. Bezzi M, De Capitani di Vimercati S, Foresti S, Livraga G, Samarati P, et al. Modeling and preventing inferences from sensitive value distributions in data release 1. *Journal of Computer Security* 2012;20:393–436.
25. Cover TM, Thomas JA. *Elements of information theory*. John Wiley & Sons; 2012.
26. Diaz M, Wang H, Calmon FP, Sankar L. On the Robustness of Information-Theoretic Privacy Measures and Mechanisms. *IEEE Transactions on Information Theory* 2020;66:1949–78.
27. Lendasse A. Practical Estimation of Mutual Information on Non-Euclidean Spaces. In: *Machine Learning and Knowledge Extraction: First IFIP TC 5, WG 8.4, 8.9, 12.9 International Cross-Domain Conference, CD-MAKE 2017, Reggio, Italy, August 29–September 1, 2017, Proceedings. vol. 10410. Springer; 2017. p. 123.*
28. Liao J, Kosut O, Sankar L, du Pin Calmon F. Tunable Measures for Information Leakage and Applications to Privacy-Utility Tradeoffs.

- IEEE Transactions on Information Theory 2019;65:8043–66.
29. Agrawal D, Aggarwal CC. On the design and quantification of privacy preserving data mining algorithms. In: Proceedings of the twentieth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems. ACM; 2001. pp. 247–55.
  30. du Pin Calmon F, Fawaz N. Privacy against statistical inference. In: 2012 50th Annual Allerton Conference on Communication, Control, and Computing (Allerton). IEEE; 2012. pp. 1401–8.