

Re-identification risk prediction paradigm using incomplete statistical information and recursive hypergeometric distribution

Zhigang Yang (✉ ayzg163@163.com)

Chongqing University of Posts and Telecommunications <https://orcid.org/0000-0002-7268-5390>

Ruyan Wang

Chongqing University of Posts and Telecommunications

Dapeng Wu

Chongqing University of Posts and Telecommunications

Boran Yang

Chongqing University of Posts and Telecommunications

Daizhong Luo

Chongqing University of Arts and Sciences

Article

Keywords: recursive hypergeometric distribution, statistical information, risk prediction paradigm

Posted Date: May 17th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-533949/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

¹ **Re-identification risk prediction paradigm using incomplete
2 statistical information and recursive hypergeometric distribution**

³ **Author Information**

⁴ **Affiliations**

⁵ **Contributions**

⁶ **Corresponding author**

⁷ **Abstract**

⁸ The dataset anonymization has not eliminated the re-identification risk, the evaluation of which remains a
⁹ huge challenge, especially given incomplete statistical information. The re-identification risk of
¹⁰ individuals depends on their tuple frequency. The paper proposes the recursive hypergeometric (RH)
¹¹ distribution to accurately calculate the tuple frequency and leverages the binomial distribution to
¹² approximate the RH distribution and to efficiently predict the re-identification risk of individuals in both
¹³ generated and real-world datasets. The experimental results show that our tuple frequency based re-
¹⁴ identification risk (TFRR) prediction model has a superior performance (average AUC 0.86~0.98) for all
¹⁵ types of datasets. Furthermore, we exploit the value dependence knowledge to rectify the prediction result
¹⁶ for some subsets (average AUC 0.95~0.98). Our research reveals the general rule of the tuple frequency
¹⁷ distribution and enables individuals and regulators to responsively predict the re-identification risk.

¹⁸ **Introduction**

¹⁹ Today we are living in an era of data explosion¹. We have easier access to information services than at
²⁰ any time in history, but we also face unprecedented privacy risks because your service providers are
²¹ extremely likely to know you better than you do^{2,3,4}. Although service providers often allege that they
²² have to collect as much personal data as possible to improve user experience, they fail to properly protect

23 user privacy^{5,6}. In this regard, the government of many countries promulgated privacy protection laws,
24 such as the General Data Protection Regulation⁷ (GDPR) in Europe and Personal Information Security
25 Specification (PISS) in China⁸. PISS emphasizes that all collected personal data should be immediately
26 de-identified and stored separately from their profile data^{9,10}. However, even after the de-identification,
27 anonymized personal data still face re-identification risk and are vulnerable to linkage attacks launched by
28 either honest but curious data collectors or malicious hackers^{11,12}. Therefore, the re-identification risk of
29 individual data not only reflects the privacy risk level of individuals but also supports regulators in
30 formulating privacy protection policies. Beyond this, it is difficult for individuals and regulatory agencies
31 to obtain the complete dataset maintained by service providers, and they can only infer the re-
32 identification risk from the released incomplete statistical information.

33 The re-identification risk of an individual is closely related to her/his tuple frequency. The tuple
34 frequency is defined as the count of a specific data value combination, where a high tuple frequency
35 signifies a low re-identification risk. If an attacker has sufficient background knowledge for the linkage
36 attack, individuals will be re-identified by her/his unique data records with 100% probability. Therefore,
37 the uniqueness of individual data has attracted extensive research attention¹³. According to the 1990 and
38 2000 U.S. census data releases, it takes only three attributes, namely the date of birth, gender, and zip
39 code, to uniquely identify 87% and 63% of the population^{14,15}. Montjoye found that it takes only four
40 spatiotemporal points in trajectory data to uniquely identify 95% of the individuals in the location dataset
41 and 90% in the credit card dataset^{16,17}. By exploiting the uniqueness contained in the sampled data records
42 or statistical characteristics of datasets, a latent attacker can measure the uniqueness of individuals given
43 incomplete statistical information¹⁸ and even recover the original personal data¹⁹. However, using the
44 uniqueness to describe the re-identification risk is sometimes inaccurate, because non-unique data records
45 can still be exploited to re-identify individuals from anonymized datasets with a certain probability²⁰.

46 **Results**

47 The attribute dependence of experimental datasets.

48 Inspired by k -anonymity^{21,22,23}, we propose to leverage k -indistinguishability as an indicator to describe
49 the re-identification risk of individuals. If the tuple frequency of an individual in an anonymized dataset is
50 not less than k , then this individual is k -indistinguishable. If the probability of a specific individual being
51 k -indistinguishable can be derived for $k = 2, 3, \dots$, one can have a relatively more comprehensive
52 understanding of her/his re-identification risk. Unfortunately, given incomplete dataset information, the
53 state-of-the-art privacy risk research cannot determine the probability of an individual being k -
54 indistinguishable when $k > 2$. In light of this, the paper presents how to accurately predict the re-
55 identification risk for a given individual with only the incomplete statistical information of the target
56 dataset. Specifically, given some statistical information, the probability mass function (PMF) of the RH
57 distribution can be used to estimate the frequency of the tuples containing strong dependent attribute
58 pairs. In real-world applications, an approximate distribution of the RH distribution is employed to
59 calculate the tuple frequency in an anonymized dataset for computational efficiency, and to further derive
60 the probability of an individual being k -indistinguishable in the target dataset. Our experiments use
61 random²⁴, demographic²⁵, medical¹³, and educational²⁶ datasets, and the results show that for all involved
62 datasets, the average AUC of our proposed TFRR is 0.86~0.98, suggesting a high prediction accuracy.
63 For datasets containing strongly dependent attribute pairs, the value dependence knowledge is introduced
64 to rectify the prediction results and the average AUC reaches 0.95~0.98. Our research reveals a general
65 rule determining the distribution of the tuple frequency, which is applicable for all random datasets and
66 most real-world datasets and provides a concise yet effective tool for the re-identification risk prediction
67 of anonymized datasets. With the incomplete statistical information of the target dataset, both individuals
68 and regulators can easily use this tool to predict the re-identification risk. Beyond this function, one can
69 even predict the re-identification risk of submitting data to service providers according to their published
70 data formats, statistical information, and privacy protection plans, and accordingly question whether they
71 obey the existing privacy protection laws, so as to foresee and prevent privacy threats.

72 Considering dataset D is a table with columns representing attributes and rows representing data records.
73 Each cell in the table maintains the value of a particular attribute of a particular data record. A tuple is
74 defined as an ordered list drawing one value per attribute, to enumerate all possible cases of data records
75 in D , some of which may not appear in D . From the perspective of probability theory, the frequency of
76 a specific tuple in a target dataset follows the RH distribution (see Methods). However, the dependence
77 between the values in the tuple will affect the tuple frequency distribution. Therefore, we define value
78 dependence (see Methods) to describe the dependence between the value pairs of a tuple, and use the
79 value dependence knowledge of a particular tuple to rectify the prediction results. To grasp a general
80 understanding of the dependence between an attribute pair, we define the attribute dependence and
81 analyze the dependence between each attribute pair in experimental datasets (including random and real-
82 world datasets).

83 The attribute dependence profiles an asymmetric relation between two attributes. The dependence of
84 attribute B on attribute A can be calculated as follows,

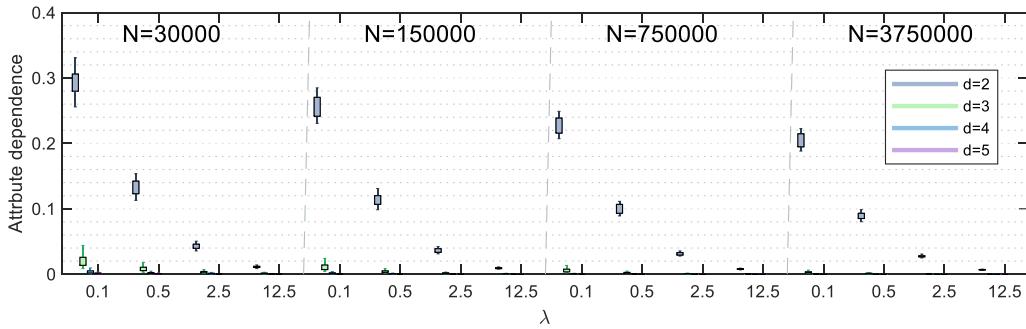
$$85 \quad D(A \Rightarrow B) = \frac{I(A; B)}{H(B)}, \quad (1)$$

86 where $I(A; B)$ is the mutual information of A and B , and $H(B)$ is the entropy of B . (See Methods in
87 Supplementary)

88 If the value of $D(A \Rightarrow B)$ or $D(B \Rightarrow A)$ reaches or exceeds a certain threshold, the attribute pair
89 (A, B) is called as a strongly dependent attribute pair. The threshold is set to 0.5 in this paper, which is
90 very consistent with the subsequent experimental results.

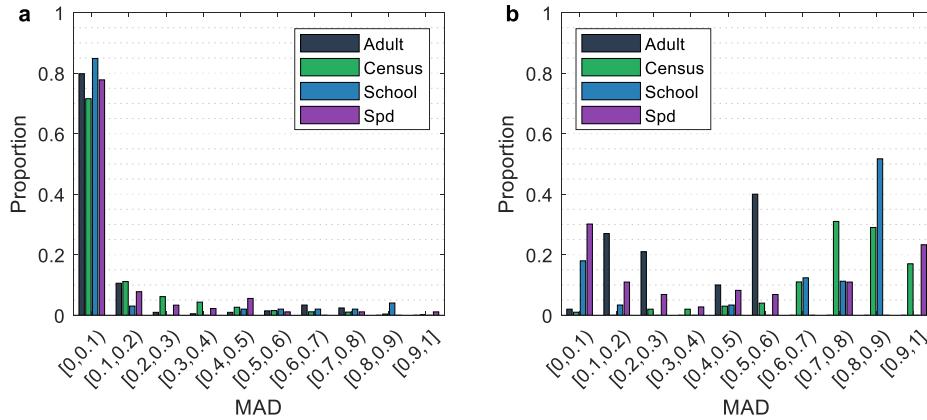
91 We use a parameter set (N, d, λ) to generate random datasets, where N is the size of the dataset, d is
92 the number of attributes, λ is the indistinguishability indicator equal to the ratio of N to the size of
93 sample space \mathcal{X} (see Supplementary Note 1). For parameter combinations from
94 $N \in \{30000, 150000, 750000, 3750000\}$, $d \in \{2, 3, 4, 5\}$, $\lambda \in \{0.1, 0.5, 2.5, 12.5\}$, 64 different

95 parameter sets generate 64×50 random datasets, and the attribute dependence of these datasets are
 96 illustrated in Fig. 1.



97
 98 **Fig.1** The attribute dependence of random datasets. 64 boxes represent the attribute dependence of the
 99 datasets randomly generated by the 64 parameter sets. The grey dotted vertical line divides the figure into
 100 4 sectors, each of which corresponds to a different value of N , the X-axis labels λ , the gray, green,
 101 blue, and purple boxes indicate that $d = 2, 3, 4$, and 5 , respectively. With two of the parameters fixed,
 102 the attribute dependence decreases as the third parameter increases. The decreasing trend with d is the
 103 most obvious, followed by λ and N . When $d = 2$ and $\lambda = 0.1$, the average attribute dependence lies
 104 between 0.2 and 0.3. When $d \geq 3$, the attribute dependence drops below 0.05, which signifies a very low
 105 attribute dependence for random datasets above three dimensions. When $d = 2$ and $\lambda = 0.1$, the high
 106 attribute dependence is due to the small λ and the limited dataset size (see Supplementary Note 2).

107
 108 We use 4 real-world datasets collected from public sources: the populations of the U.S. (ADULT,
 109 Census), medical dataset (SPD), and education dataset (SCHOOL). The attribute dependence of each real-
 110 world dataset is shown in Fig. 2a. From each dataset, we create subsets by randomly selecting attributes
 111 (columns). The resulting 362 subsets cover a large range of attribute dependence values (0.000062-
 112 0.9886), numbers of attributes (2-28), and data records (44842-3985166 individuals).
 113 The distribution of Max Attribute Dependence (MAD) of the 362 subsets is shown in Fig. 2b, where
 114 MAD is the maximum attribute dependence of all attribute pairs in a subset.



115

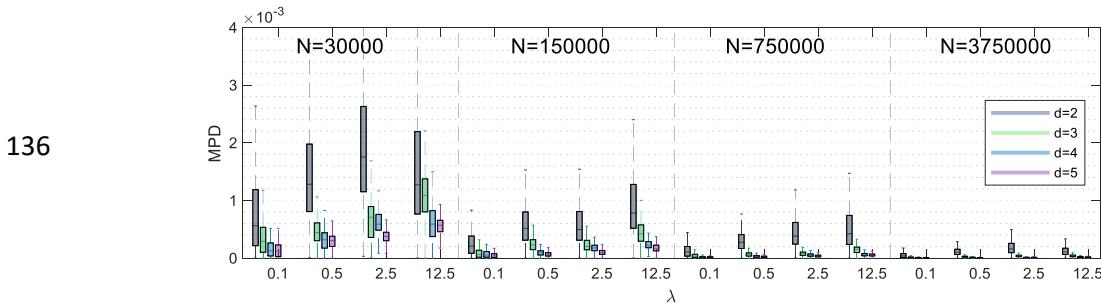
116 **Fig.2** The attribute dependence of real-world datasets. **a** The attribute dependence of real-world datasets. The black, green, blue, and purple columns represent the attribute dependence of Adult, Census, School, and SPD, respectively. 70% of the attribute dependence are less than 0.1 and only 13% of the attribute dependence are greater than 0.3, which means that most of the attribute pairs in these four real-world datasets have a low dependence. **b** The MAD of subsets. 63.3% of the subsets have a $\text{MAD} \geq 0.5$, Adult(40%), Census(92%), School(75.3%), SPD (41.1%), suggesting more than half of the subsets contain at least 1 pair of strongly dependent attributes.

123 The approximate RH distribution.

124 Because of the computational complexity of the RH distribution, we expect to find an approximate
 125 distribution to reduce the computational burden. According to the analysis in Methods, we find that when
 126 $N \rightarrow \infty$, the binomial distribution $B(n_d, P_{d-1})$ can be employed to approximate the RH distribution. To
 127 have a clearer understanding of the difference between them, we randomly select many tuples from the
 128 random datasets and use the PMFs of the two distributions to calculate the occurrence probability of these
 129 tuples. The maximum probability distance (MPD) is used to measure the difference between the binomial
 130 distribution and the RH distribution, defined as,

$$131 \quad \text{MPD} = \max_k (f_B(k) - f_{RH}(k)), \quad (2)$$

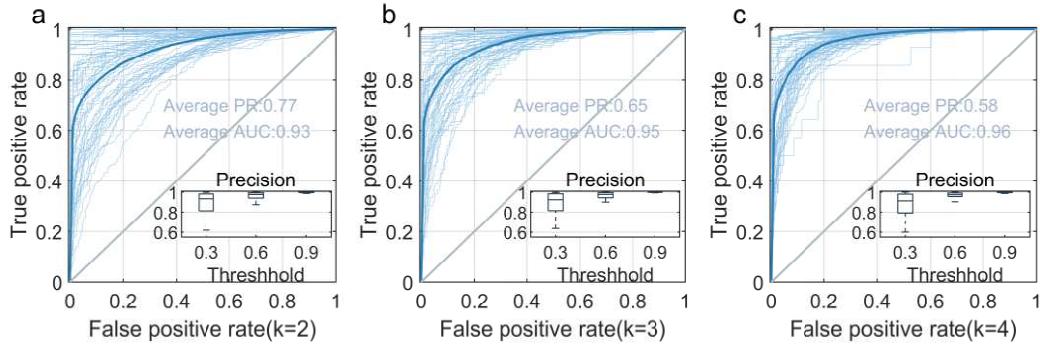
132 where $f_B(k)$ and $f_{rH}(k)$ are the PMFs of $B(n_d, P_{d-1})$ and $rH(N, d, n_1, \dots, n_d)$.
 133 We use the same 64 parameter sets as in the previous experiment to generate random datasets, and
 134 randomly select 1000 tuples from each dataset. The MPD between the binomial distribution and the RH
 135 distribution is shown in Fig. 3.



136
 137 **Fig. 3** The MPD between $B(n_d, P_{d-1})$ and $rH(N, d, n_1, \dots, n_d)$. The grey dotted vertical line divides the
 138 figure into 4 sectors, each of which corresponds to a different value of N , the X-axis labels λ , the gray,
 139 green, blue, and purple boxes indicate that $d = 2, 3, 4$, and 5 , respectively. The MPD decreases as N
 140 or d increases if the other two parameters are fixed. When $N = 30000$, the MPD lies between 0-0.005.
 141 When $N \geq 750000$, the MPD drops to 0-0.001. This signifies that, if $N \rightarrow \infty$, the difference between
 142 $B(n_d, P_{d-1})$ and $rH(N, d, n_1, \dots, n_d)$ is negligibly small and $B(n_d, P_{d-1})$ can be exploited to
 143 approximate $rH(N, d, n_1, \dots, n_d)$ in calculating the tuple frequency.

144 The possibility of an individual being k -indistinguishable in random datasets.
 145 We randomly select the data records of 1000 individuals from 64 random datasets and use Eq. 11 to
 146 estimate the possibility of these individuals being k -indistinguishable (see Methods for details). The
 147 result of binary classification is shown in Fig. 4.

148

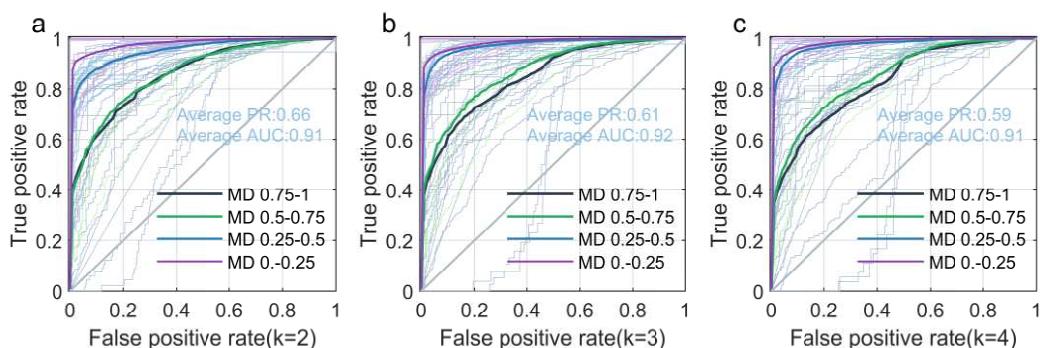


149 **Fig. 4** Receiver operating characteristic (ROC) curves for random datasets (light blue for each dataset and
 150 dark blue for the average). **a**, **b**, and **c** identify the impact of different values of k on prediction
 151 performance. The AUC of our proposed TFRR prediction is over 0.93, showing a high prediction
 152 performance for random datasets. The average positive rate declines sharply as k grows, suggesting
 153 fewer individuals will be k -indistinguishable in random datasets if k increases. When the classification
 154 threshold is set to 0.6, the precision of our proposed TFRR exceeds 0.9.

155 The possibility of an individual being k -indistinguishable in real-world datasets.

156 We randomly select the data records of 1000 individuals from each subset and then calculate the
 157 possibility of these individuals being k -indistinguishable. If all the data samples from a subset are either
 158 positive or negative, no binary prediction will be performed on this subset (See Supplementary Methods).
 159 For readability, we present only the results for subsets from Spd in Fig. 5. Other prediction results are
 160 available in Supplementary.

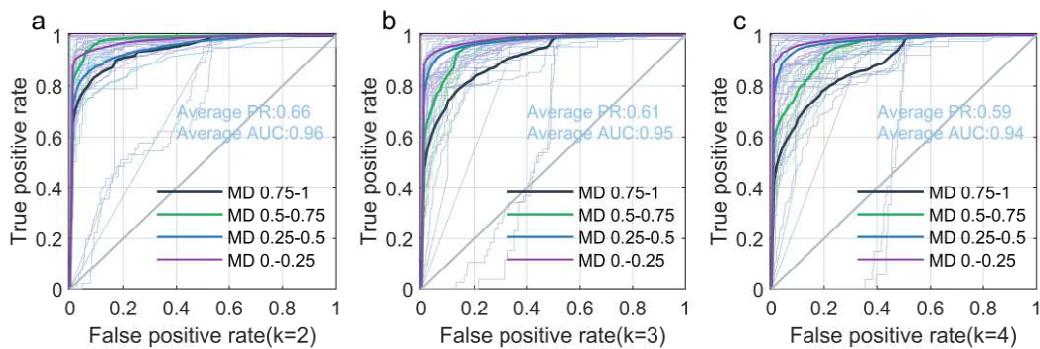
161



162 **Fig. 5** ROC curves for real-world datasets. **a**, **b**, and **c** identify the impact of different values of k on
 163 prediction performance. Gray (black), light green (green), light blue (blue), and light purple (purple)
 164 represent the ROC (average ROC) curves when the MADs are [0.75-1], [0.5-0.75], [0.25-0.5], and [0-
 165 0.25], respectively. For subsets whose $\text{MAD} < 0.5$, the average AUC lies between 0.83~0.87. The average
 166 Positive Rate (PR) and average AUC of each subfigure are **a** average PR=0.66 average AUC=0.91. **b**
 167 average PR=0.61 average AUC=0.92 **c** average PR=0.59 average AUC=0.91. See Fig. 1 in
 168 Supplementary for more information.

169 Exploiting value dependence knowledge to rectify prediction results.

170 For the subsets containing strongly dependent attribute pairs, approximation distribution $B(n_i p_{ij}, p')$
 171 (see Methods for details) is exploited to recalculate the prediction scores of the true positive samples
 172 whose prediction scores are less than 0.5. Only the samples of the subsets from Spd are used to
 173 recalculate, because most of the subsets from the other real-world datasets contain at least two strongly
 174 dependent attribute pairs, meaning that the samples in these datasets may contain two or more strongly
 175 dependent value pairs. But $B(n_i p_{ij}, p')$ can only rectify the negative impact of 1 strongly dependent
 176 value pair on the prediction performance. The improved ROC curves are shown in Fig. 6.



177
 178 **Fig. 6** Rectified ROC curves. This experiment only considers the value dependence knowledge of
 179 attribute pair (ZIP code, county), (ZIP code, hospital ID), and (hospital ID, county). The complex
 180 dependence knowledge involving more than 3 attribute values is ignored (see Methods). The average PR

181 and average AUC of each subfigure are **a** average PR=0.66 average AUC=0.96. **b** average PR=0.61
182 average AUC=0.95 **c** average PR=0.59 average AUC=0.94. The AUC curves of all the tested subsets are
183 rectified by using the value dependence knowledge, except for three subsets with extremely low positive
184 or negative sample ratio (AUC: 0.54~0.80), and the AUC of the other subsets are relatively high (AUC:
185 0.86~1). This shows that the value dependence knowledge can be leveraged to efficiently improve the
186 prediction performance (see Note 3 in Supplementary).

187 **Discussion**

188 The tuple frequency follows the RH distribution. This rule enables the re-identification risk prediction of
189 individuals before even obtaining the complete dataset. Due to the computational complexity of the RH
190 distribution, the binomial approximation can help predict the privacy risk of individuals.

191 Both attribute dependence and value dependence are employed to measure the relation between attributes
192 and values. The dependence between values will break the randomness of the value distribution, and the
193 prediction error will be unacceptable when using the binomial approximation to predict the tuple
194 frequency of strongly dependent value pairs. To estimate the frequency of such tuples, we need to not
195 only know the parameters including N, d, n_1, \dots, n_d , but also obtain the attribute dependence knowledge,
196 which can be obtained from:

- 197 1. Statistical information published by statistical agencies or data service providers^{27,28}.
- 198 2. Data sampled from complete datasets. Since the value proportion in the sample is an unbiased estimate
199 of the value proportion in the subset, if the sample and the size of the subset are obtained, the attribute
200 dependence knowledge of the subset can be derived using statistical inference. This enables the estimation
201 of the frequency distribution of the tuple matching a particular individual using the sampled data from the
202 complete dataset.

203 The knowledge of attribute dependence and value dependence can reveal the internal relation between
204 different attributes and values of data records, which are important indicators for the value distribution of

205 data²⁹. Existing privacy protection methods, such as differential privacy^{30,31,32}, can hide the original data
206 while ensuring their availability by adding random noise regularly. Although we can obtain useful
207 information, such as the tuple frequency distribution and top- k data from the de-identified dataset³³, the
208 value and utilization of the de-identified dataset are substantially reduced due to the impact of the added
209 noise on the attribute dependence and value dependence³⁴. Therefore, we plan to study how to customize
210 the differential privacy budgets and noise generation methods according to the predicted re-identification
211 risk for specific individuals, and how to maximize the preservation of the attribute dependence and value
212 dependence information. In addition, trajectory datasets typically contain temporal and sequential location
213 data with strong attribute dependence³⁵, which makes the binomial approximation ineffective in privacy
214 risk prediction. This also points out a new research direction for future work.

215 **Methods**

216 The d -RH distribution.

217 The d -dimensional recursive hypergeometric (d -RH) distribution describes the frequency of a specific
218 tuple in a d -dimensional dataset, namely the tuple frequency. For a d -dimensional dataset D of size
219 N , \mathcal{X} is the universe of all possible tuples in D . Considering A_1, A_2, \dots, A_d are the value sets of d
220 attributes, $x^i = \{x_1^i, x_2^i, \dots, x_d^i\} (1 \leq i \leq |\mathcal{X}|)$ denotes a possible tuple in D , where $x_j^i \in A_j (1 \leq j \leq d)$ is
221 the j -th value of x^i . Let the frequencies of $x_1^i, x_2^i, \dots, x_d^i$ in corresponding columns (attributes) be
222 n_1, n_2, \dots, n_d . Let the frequency of x^i in D be X_d , such that $X_d \sim_r H(N, d, n_1, \dots, n_d)$, then
223 $_r H(N, d, n_1, \dots, n_d)$ is the d -RH distribution.
224 Considering $X_j (1 \leq j \leq d)$ satisfies the j -RH distribution characterized by N, j, n_1, \dots, n_j .
225 When $j = 1$,

$$P(X_j = k) = \begin{cases} 1 & k = n_1 \\ 0 & \text{else} \end{cases} \quad (3)$$

227 When $j \geq 2$, the PMF of the j -RH distribution can be obtained recursively as follows,

$$\begin{aligned}
 P(X_j = k) &= \sum_{m=k}^{n_{j-1}} P(X_{j-1} = m) P(X_j = k / X_{j-1} = m) \\
 &= \sum_{m=k}^{n_{j-1}} P(X_{j-1} = m) \cdot \frac{C_m^k C_{N-m}^{n_j-k}}{C_N^{n_j}}
 \end{aligned} \tag{4}$$

229 Equation 4 can be interpreted as that, given that sub-tuple

230 $x^i(j-1) = \{x_1^i, \dots, x_{j-1}^i\}$ ($1 \leq i \leq |\mathcal{X}|$, $2 \leq j \leq d$) appears m times in the first $j-1$ columns of D , the
 231 times that $x^i(j)$ appears in the first j columns follow the hypergeometric distribution $H(N, m, n_i)$.

232 Therefore, the PMF of ${}_r H(N, d, n_1, \dots, n_d)$ can be obtained similarly by the recursive calculation of Eq.
 233 4.

234 The PMF of 2-RH distribution $H(N, 2, n_1, n_2)$ is,

$$235 \quad P(X_2 = k) = \frac{C_{n_l}^k C_{N-n_l}^{n_2-k}}{C_{n_2}^{n_2}}. \quad (5)$$

236 The 2-RH distribution $H(N, 2, n_1, n_2)$ and the hypergeometric distribution $H(N, n_1, n_2)$ are identical.

Therefore, we consider that the hypergeometric distribution is only a special case of d -RH distributions with $d = 2$.

²³⁹ The binomial approximation of the d-RH distribution.

240 Since the frequency of the d -th value of x^i is n_d , the frequency of x^i must be less than or equal to
 241 n_d , namely $X_d \leq n_d$. Considering r is a record in D whose d -th value is x_d^i and a n_d round
 242 experiment is designed. In each round, we randomly draw one value per attribute in the first $d-1$

243 columns without replacement and combine them with the fixed x_d^i to obtain r . The count of rounds in
 244 which record r matching tuple x^i is equal to X_d , such that $X_d \sim_r H(N, d, n_1, \dots, n_d)$.

245 Because we draw values without replacement in each round, the probability of record r matching tuple
 246 x^i varies³⁶. But when $N \rightarrow \infty$, this probability stays almost the same. Therefore, the count of rounds
 247 satisfies the binomial distribution. The probability of record r matching tuple x^i is as follows,

$$248 P_{d-1} = n_1 \times \dots \times n_{d-1} / N^{d-1} \quad (6)$$

249 Then the probability of record r matching tuple x^i for k rounds is as follows,

$$250 P(X_d = k) = C_{n_d}^k P_{d-1}^k (1 - P_{d-1})^{n_d - k}. \quad (7)$$

251 In conclusion, when $N \rightarrow \infty$, the binomial distribution $B(n_d, P_{d-1})$ approximates the d -RH distribution
 252 $_r H(N, d, n_1, \dots, n_d)$.

253 Value dependence.

254 There is a dependence between each pair of values in a tuple. For example, in a shopping dataset,
 255 someone who bought a long dress is highly likely to be female. The value long dress determines to a large
 256 extent the value female. The confidence and lift³⁷ in association rule learning are employed as indicators
 257 to describe the dependence between a pair of values.

258 Let $x_1 \in A_1$ and $x_2 \in A_1$ be two values, E_x and E_y be the events that record r contains x_1 and x_2 ,
 259 then the confidence and lift of the association rule that x_1 implies x_2 are as follows.

$$260 \text{conf}(x_1 \Rightarrow x_2) = \frac{|\{r \in D; \{x_1, x_2\} \subseteq r\}|}{|\{r \in D; \{x_1\} \subseteq r\}|} = P(E_{x_2} | E_{x_1}) \quad (8)$$

$$261 \text{lift}(x_1 \Rightarrow x_2) = \frac{|\{r \in D; \{x_1, x_2\} \subseteq r\}|}{|\{r \in D; \{x_1\} \subseteq r\}| \times |\{r \in D; \{x_2\} \subseteq r\}|} = \frac{P(E_{x_1} \cap E_{x_2})}{P(E_{x_1})P(E_{x_2})} \quad (9)$$

262 If $\text{conf}(x_1 \Rightarrow x_2)$ or $\text{conf}(x_2 \Rightarrow x_1)$ reaches or exceeds a certain threshold, the value pair (x_1, x_2) is
263 called as a strongly dependent value pair, and the threshold is set to 0.5 in this paper.

264 The frequency distribution of tuples with strongly dependent value pairs.

265 The physical significance of d -RH distribution can be summarized as follows. Let $x = (x_1, \dots, x_d)$ be a
266 possible tuple of a d -dimensional dataset D with size N , X be the frequency of tuple x in D . Let the
267 frequencies of x_1, \dots, x_d in the corresponding column be n_1, n_2, \dots, n_d . Considering \mathcal{D}^i is the set of all
268 datasets in which the frequency of x is i , and $\mathcal{D} = \mathcal{D}^0 \cup \dots \cup \mathcal{D}^\infty$ is the set of all datasets that may
269 contain x . The d -RH distribution is based on classical probability theory. From the perspective of
270 classical probability theory³⁸, each dataset in \mathcal{D} has the same odds to be D , and the probability of tuple
271 x appearing k times in D can be obtained by $P(X=k) = |\mathcal{D}^k| / |\mathcal{D}|$. When $n_1, n_2, \dots, n_d \ll N$, in
272 most datasets of \mathcal{D} , x does not contain strongly dependent value pairs. Therefore, for the tuples without
273 strongly dependent value pairs, each dataset in \mathcal{D} has almost the same odds to be D and the d -RH
274 distribution is suitable for predicting the frequency of such tuples. For the tuples with strongly dependent
275 value pairs, the occurrence probability of each dataset is unequal, and therefore the d -RH distribution
276 will fail in prediction.

277 Fortunately, introducing the background knowledge of value dependence can improve the prediction
278 result. Assume that all the values of the tuple x are divided into groups A and B, and all values in A are
279 the members of the strongly dependent value pair, whereas all values in B are not the members. The times
280 of group A appearing in D is denoted by m , the number of members in B is s , ordered by i_1, \dots, i_s , then
281 the frequency of tuple x can be denoted by an $(s+1)$ -RH distribution with parameter
282 $N, s+1, n_{i_1}, \dots, n_{i_s}, m$.

283 When the number of values in A is greater than 2, the dependence between members becomes complex.
 284 Therefore, we only consider in this paper the cases that group A has only 2 values, and give the binomial
 285 approximation for these cases. Considering group A only has two values x_i and x_j , where $1 \leq i, j \leq d$
 286 and $\text{conf}(x_i, x_j) = p_{ij} > \text{conf}(x_j, x_i)$, then

$$287 n' = \begin{cases} n_i p_{ij} & d > 2 \\ n_i & d = 2 \end{cases} \quad (10)$$

$$288 p' = \begin{cases} \left(\prod_{k=1, k \neq i, j}^d n_k \right) / N^{d-2} & d > 2 \\ p_{ij} & d = 2 \end{cases} \quad (11)$$

289 The frequency distribution of X in D can be approximated by $B(n', p')$.

290 The probability of a specific individual being k-indistinguishable.

291 Considering $r = (r_1, \dots, r_d)$ is the record of an individual p in the d -dimensional dataset D , X_d is the
 292 frequency of the tuple matching r in D , such that X_d follows the d -RH distribution.

293 Then the probability of p being k-indistinguishable when $k \geq 2$ can be calculated as follows

$$\begin{aligned} 294 & P(p \text{ is } k\text{-indistinguishable}) \\ & = P(X_d \geq k \mid r \text{ in } D) \\ & = P(X_d \geq k) / P(r \text{ in } D) \\ & = \{1 - \sum_{i=0}^{k-1} P(X_d = i)\} / \{1 - P(X_d = 0)\} \end{aligned} \quad (12)$$

295

296 **Data Availability**

297 All data information related to the article is provided in the Supplementary Methods.

298 **Code Availability**

299 All simulations were implemented in Matlab. The source code to reproduce the experiments will be is
300 deposited in Code ocean or Github.

301 **References**

- 302 1. Marx, V. (2013). The big challenges of big data. *Nature*, 498(7453), 255-260.
- 303 2. Blumenstock, J., Cadamuro, G., & On, R. (2015). Predicting poverty and wealth from mobile phone
304 metadata. *Science*, 350(6264), 1073-1076.
- 305 3. Thompson, S. A., & Warzel, C.. How to Track President Trump. *The New York Times*, 12, 20
306 (2019).
- 307 4. Raskar, R., Schunemann, I., Barbar, R., Vilcans, K., Gray, J., Vepakomma, P., ... & Werner, J.
308 (2020). Apps gone rogue: Maintaining personal privacy in an epidemic. *arXiv preprint*
309 *arXiv:2003.08567*.
- 310 5. Kozlowska, I. (2018). Facebook and data privacy in the age of Cambridge Analytica. Seattle, WA:
311 The University of Washington. Retrieved August, 1, 2019.
- 312 6. Ablon, L., & Libicki, M. (2015). Hacker's bazaar: The markets for cybercrime tools and stolen data.
313 *Def. Counsel J.*, 82, 143.
- 314 7. Voigt, P., & Von dem Bussche, A. (2017). The eu general data protection regulation (gdpr). A
315 Practical Guide, 1st Ed., Cham: Springer International Publishing, 10, 3152676.
- 316 8. Wilmer Hale. China Issues New Personal Information Security Specification.3, 24 (2020).
317 <https://www.wilmerhale.com/en/insights/client-alerts/20200324-china-issues-new-personal-information-security-specification>
- 319 9. Shi, M. L., Sacks, S., Chen, Q. H., & Webster, G. (2019). Translation: China's personal information
320 security specification. *New America*. <https://www.newamerica.org/cybersecurity-initiative/digichina/blog/translation-chinas-personal-information-security-specification/>

- 322 10. Hong, Y. et al. Information security technology—Personal information security specification.
323 <https://www.tc260.org.cn/upload/2020-09-18/1600432872689070371.pdf> (2020)
- 324 11. Clauß, S., Kesdogan, D., & Kölsch, T. (2005, November). Privacy enhancing identity management:
325 protection against re-identification and profiling. In Proceedings of the 2005 workshop on Digital
326 identity management (pp. 84-93).
- 327 12. Henriksen-Bulmer, J., & Jeary, S. (2016). Re-identification attacks—A systematic literature review.
328 International Journal of Information Management, 36(6), 1184-1192.
- 329 13. Sánchez, D., Martnez, S. & Domingo-Ferrer, J. Comment on unique in the shopping mall: on the
330 reidentifiability of credit card metadata. Science 351, 1274 (2016).
- 331 14. Sweeney, L. (2000). Uniqueness of simple demographics in the US population. LIDAP-WP4, 2000.
- 332 15. Golle, P. (2006, October). Revisiting the uniqueness of simple demographics in the US population.
333 In Proceedings of the 5th ACM Workshop on Privacy in Electronic Society (pp. 77-80).
- 334 16. De Montjoye, Y. A., Hidalgo, C. A., Verleysen, M., & Blondel, V. D. (2013). Unique in the crowd:
335 The privacy bounds of human mobility. Scientific reports, 3(1), 1-5.
- 336 17. De Montjoye, Y. A., Radaelli, L., & Singh, V. K. (2015). Unique in the shopping mall: On the
337 reidentifiability of credit card metadata. Science, 347(6221), 536-539.
- 338 18. Rocher, L., Hendrickx, J. M., & De Montjoye, Y. A. (2019). Estimating the success of re-
339 identifications in incomplete datasets using generative models. Nature communications, 10(1), 1-9.
- 340 19. Tu, Z., Xu, F., Li, Y., Zhang, P., & Jin, D. (2018). A new privacy breach: User trajectory recovery
341 from aggregated mobility data. IEEE/ACM Transactions on Networking, 26(3), 1446-1459.
- 342 20. Benitez, K., & Malin, B. (2010). Evaluating re-identification risks with respect to the HIPAA privacy
343 rule. Journal of the American Medical Informatics Association, 17(2), 169-177.

- 344 21. Sweeney, L. (2002). k-anonymity: A model for protecting privacy. International Journal of
345 Uncertainty, Fuzziness and Knowledge-Based Systems, 10(05), 557-570.
- 346 22. Sweeney, L. (2002). Achieving k-anonymity privacy protection using generalization and
347 suppression. International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems, 10(05),
348 571-588.
- 349 23. El Emam, K., & Dankar, F. K. (2008). Protecting privacy using k-anonymity. Journal of the
350 American Medical Informatics Association, 15(5), 627-637.
- 351 24. Yang, Z., Wang, R., Luo, D., & Xiong, Y. (2020). Rapid Re-Identification Risk Assessment for
352 Anonymous Data Set in Mobile Multimedia Scene. IEEE Access, 8, 41557-41565.
- 353 25. Asuncion, A., & Newman, D. (2007). UCI machine learning repository.
- 354 26. Romania government. (2014). Primary school enrollment 2014.
- 355 27. Fung, B. C., Wang, K., Chen, R., & Yu, P. S. (2010). Privacy-preserving data publishing: A survey
356 of recent developments. ACM Computing Surveys (Csur), 42(4), 1-53.
- 357 28. Abella, A., Ortiz-de-Urbina-Criado, M., & De-Pablos-Heredero, C. (2019). The process of open data
358 publication and reuse. Journal of the Association for Information Science and Technology, 70(3),
359 296-300.
- 360 29. Zheng, Y., Zhang, L., Xie, X., & Ma, W. Y. (2009, November). Mining correlation between
361 locations using human location history. In Proceedings of the 17th ACM SIGSPATIAL international
362 conference on advances in geographic information systems (pp. 472-475).
- 363 30. Dwork, C., & Roth, A. (2014). The algorithmic foundations of differential privacy. Foundations and
364 Trends in Theoretical Computer Science, 9(3-4), 211-407.

- 365 31. Friedman, A., & Schuster, A. (2010, July). Data mining with differential privacy. In Proceedings of
366 the 16th ACM SIGKDD international conference on Knowledge discovery and data mining (pp. 493-
367 502).
- 368 32. Xiong, J., Ma, R., Chen, L., Tian, Y., Li, Q., Liu, X., & Yao, Z. (2019). A personalized privacy
369 protection framework for mobile crowdsensing in IIoT. IEEE Transactions on Industrial Informatics,
370 16(6), 4231-4241.
- 371 33. Fanti, G., Pihur, V., & Erlingsson, Ú. (2015). Building a RAPPOR with the unknown: Privacy-
372 preserving learning of associations and data dictionaries. arXiv preprint arXiv:1503.01214.
- 373 34. Ren, X., Yu, C. M., Yu, W., Yang, S., Yang, X., McCann, J. A., & Philip, S. Y. (2018). LoPub:
374 High-Dimensional Crowdsourced Data Publication with Local Differential Privacy. IEEE
375 Transactions on Information Forensics and Security, 13(9), 2151-2166.
- 376 35. Yang, Z., Wang, R., Wu, D., & Luo, D. (2020). UTM: A trajectory privacy evaluating model for
377 online health monitoring. Digital Communications and Networks.
- 378 36. Berkopac, A. (2007). HyperQuick algorithm for discrete hypergeometric distribution. Journal of
379 Discrete Algorithms, 5(2), 341-347.
- 380 37. Hornik, K., Grün, B., & Hahsler, M. (2005). arules-A computational environment for mining
381 association rules and frequent item sets. Journal of statistical software, 14(15), 1-25.
- 382 38. Jaynes, E. T. (2003). Probability theory: The logic of science. Cambridge university press.

383 **Acknowledgments (optional)**

384 This work was supported.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Supplementary.pdf](#)