

THEORY

On the Complexity of Optimal k -Anonymity: A New Proof Based on Graph Coloring

YAVUZ CANBAY^{ID}

Department of Computer Engineering, Kahramanmaraş Sütçü İmam University, 46040 Kahramanmaraş, Türkiye

e-mail: yavuzcanbay@ksu.edu.tr

ABSTRACT Privacy is a complex balancing problem between risks and utility of data. k -anonymity, a fundamental model for preserving privacy, guarantees that an item cannot be differentiated from at least $k-1$ other items. Due to the k -anonymity is a hard problem, which means obtaining an optimal solution within a reasonable time is not possible, researchers endeavor to create near-optimal solutions. There are some researches in the literature demonstrating the NP-Hardness of achieving k -anonymity. The problems of k -dimensional perfect matching, edge partition into triangles, minimum vertex covering, and maximum k -dimensional matching with k -occurrences are some examples of NP-Complete problems commonly used for reduction to prove the NP-Hardness of k -anonymity. However, previous proofs use large alphabet size and suppress more cells causing less utility. This study presents a significant contribution by providing a new proof for the NP-Hardness of k -anonymity and enhances both the alphabet size and the number of suppressed cells. The proof is achieved by using a reduction from the graph coloring problem, which is being provided for the first time.

INDEX TERMS Graph coloring, k -anonymity, np-hardness, proof.

I. INTRODUCTION

The digitalization of society enables the handling of an increasing amount of data pertaining to the real world, machines, individuals, and so on. Currently, numerous institutions or parties, referred as data curators, gather and retain data from various individuals and entities such as clients, patients, users, firms, and institutions. The primary objectives of these acts include fulfilling their goals, enhancing services such as customer modeling, identifying behavioral patterns, diagnosing diseases, formulating plans, establishing regulations, constructing decision-making procedures, etc. In some cases, it is necessary to publish or share the data in order to maximize the benefits. Through this approach, one can achieve outcomes that have a direct and positive impact on respondents at all levels, ranging from the individual to the country as a whole [1], [2], [3]. On the other hand, one of the most crucial concerns with data publishing is privacy.

The concept of privacy was initially presented by Warren and Brandeis in 1890 [4] and defined as the “right to be

let alone.” Today, protection of this right is ensured by legal measures, as it is considered both essential and hot topic. Preserving the privacy of individuals, who may be encounter with numerous cyber-attacks, can be defined as the right to maintain one’s individuality in both physical and digital domains, wherein individuals establish their personal boundaries. The boundaries exhibit variability throughout different cultures, countries, religions, and even among individuals [5]. In addition, it may be helpful to consider certain definitions for data privacy in order to better grasp the concept of privacy in the context of data.

In the literature, some works provide definitions for data privacy, such as “informational self-determination” [6] and “the appropriate use of responders’ information and the ability to decide what information of a responder goes where” [7]. Other definitions can be provided such as, the selective control of data owner about the borders of data sharing, including with whom, for what purpose, and to what extent, and the right to be data. In the light of the growing accumulation of personalized data, protecting data privacy has emerged as a critical necessity and an essential prerequisite for conducting data analysis [8], [9]. Any attempt to direct

The associate editor coordinating the review of this manuscript and approving it for publication was Zijian Zhang^{ID}.

publishing of raw data may violate the privacy of responders. Hence, it is vital to employ strategies that eradicate breaches of privacy.

Anonymization and cryptography are the primary methods used for protecting the privacy of data [5]. Cryptography employs cryptographic keys to encrypt data, whereas anonymization masks the identity of respondents. Due to the lack of usefulness of encrypted data, encryption is not a recommended way for data publishing. Anonymization is the method of grouping data records in such a way that each member of a group cannot be differentiated from others based on certain features. This strategy is often favored for maintaining data privacy while yet ensuring the utility of data [10].

Several privacy-preserving approaches have been established to resist privacy-focused disclosure attacks which can be listed as record linkage, attribute linkage and probability attacks. k-anonymity, l-diversity and t-closeness are well-known and important models for protecting privacy. k-anonymity is a method that addresses the issue of record linkage attack by guaranteeing that a record cannot be differentiated from at least k-1 other records [11]. l-diversity ensures the presence of diverse sensitive information within equivalence classes, effectively mitigating record linkage and attribute linkage attacks [12]. t-closeness addresses the issues of attribute linkage and probability attacks. It ensures a balance between the distribution of sensitive data both inside each equivalence class and over the entire table [13]. A comprehensive list and their explanations can be found in [14]. While previous studies related to this topic acknowledge that each record is associated with a distinct individual, subsequent researches have begun to acknowledge the possibility that one individual may have many records [15], [16], [17].

While k-anonymity offers a solution to privacy concerns, the complexity of k-anonymity is an additional matter that needs to be handled. Previous researches have shown that brute force methods for achieving k-anonymity exhibit an exponential relationship between input size and the number of possible solutions. Hence, the literature emphasizes that achieving k-anonymity is a computationally challenging task and so it classified as NP-Hard, necessitating the use of near-optimal methods [18], [19].

This study proposes a new proof on the NP-Hardness of k-anonymity by demonstrating a reduction from the graph coloring problem. The paper is organized as follows. Section II provides some briefs about the previous proofs. In Section III, a novel proof utilizing graph coloring was introduced to establish the NP-Hardness of k-anonymity. Finally, the conclusion and discussions were provided in Section IV.

II. PREVIOUS PROOFS ON THE HARDNESS OF K-ANONYMITY

Several studies in the literature specifically address the NP-Hardness of k-anonymity.

Meyerson and Williams [18] employed a reduction from k-dimensional perfect matching problem to examine the complexity of k-anonymity. They stated that if there is no limitation on the size of the alphabet, achieving k-anonymity becomes computationally difficult and falls under the category of NP-Hard problem for $k \geq 3$. Additionally, they mentioned that the maximum number of suppressed cells is also $n(m-1)$.

However, Aggarwal et al. [19], [20] decreased the size of the alphabet to 2, while the number of suppressed cells were remained as $n(m-1) \lceil \log_2(n(m-1)/3) \rceil$ (in [19], the number of suppressed cells is determined as $9m$ where m indicates the number of triangles. Nevertheless, in order to do a comparison between the number of suppressed cells in all proofs, we generalized and assumed $9m$ as $n(m-1)$, where n represents the number of rows and m indicates the number of columns for $n(m-1)$. Therefore, we accepted that $9mt$ equals to $n(m-1) \lceil \log_2(n(m-1)/3) \rceil$). They employed a reduction technique to transform the edge partition problem into triangles problem.

In a similar manner, Sun et al. [21] reported an alphabet size of 2, while the number of repressed cells was determined as nm (in [21], the number of suppressed cell is presented as $48m$. However, in order to do a comparison between the number of suppressed cells in all proofs, we generalized and assumed $48m$ as nm , where n indicates the number of rows and m indicates the number of columns. Consequently, we accepted that $48m$ equals to nm). In their proof, they utilized edge partitioning into 4-cliques and claimed that (p, α) -sensitive k-anonymity is NP-Hard.

In addition, Bonizzoni et al. [22], specifically examined two restricted instances of the k-anonymity problem. The researchers demonstrated that achieving 3-anonymity is APX-Hard under the constraint of a binary alphabet. Furthermore, they established that achieving 4-anonymity remains APX-Hard even when the number of rows has a length of 8. The problem of minimum vertex cover on a cubic graph was also employed.

In a different study, Blocki and Williams [23] presented a proof that utilizes a reduction from the problem of maximum 3-dimensional matching with 3 occurrences. The number of attributes in each record was limited to 27, and it was demonstrated that with this constraint, achieving 3-anonymity is a computationally hard problem known as Max SNP-Hard.

Further, the proof provided by Scott et al. [24] indicates that the process of anonymizing k-attribute is also NP-Hard, for $k \geq 2$. They employed a reduction from c-Hitting Set problem.

Finally, LeFevre et al. [25] employed a reduction from partition to demonstrate the NP-Hardness of optimal k-anonymous multidimensional partitioning. They used discernibility metric to approximate optimal solution.

Based on the aforementioned proofs, we conducted detailed reviews of two studies considering the hardness of k-anonymity. The reviews are provided below.

A. THE PROOF VIA K-DIMENSIONAL PERFECT MATCHING

In [18], the authors employed a reduction from k-dimensional perfect matching to prove that k-anonymity is NP-Hard for $k \geq 3$. To understand the proof better, some notations used in the related paper were given and described in Table 1.

TABLE 1. A descriptive representation of the notations used in [18].

Notation	Description
m	dimension
Σ	alphabet of attribute values
V	table, $V \subseteq \Sigma^m$
t	suppressor function
v	record in $V = \{v_1, \dots, v_n\}$
v'	suppressed record in $V' = \{v'_1, \dots, v'_n\}$
V'	anonymized table
i	row indices $\{i_1, \dots, i_n\}$
j	column indices $\{j_1, \dots, j_m\}$
n	number of records
k	minimum cardinality in an equivalence class
*	symbol used in suppression
l	maximum number of suppressed vector coordinate in V'
H	hypergraph
U	set of vertices
E	set of edges
S	subset of edge set
u	each vertex in $U = \{u_1, \dots, u_n\}$
e	each edges in $E = \{e_1, \dots, e_n\}$
$j(i)$	indices of unique hyper edge containing u_i

suppressor function: let t be a map from V to $(\Sigma \cup \{*\})^m$. For all $v \in V$ and $j = 1, \dots, m$; if $t(v)[j] \in \{v[j], *\}$ then t is a suppressor function.

k-anonymizer: let t be a suppressor function on $V = \{v_1, \dots, v_n\} \subseteq \Sigma^m$. $t(V)$ is k-anonymous if for all $v_i \in V$, there exists $k-1$ indices $i_1, i_2, \dots, i_{k-1} \in \{1, \dots, n\}$ ensuring that $t(v_{i_1}) = t(v_{i_2}) = \dots = t(v_{i_{k-1}}) = t(v_i)$. Therefore, t is called as a k-anonymizer on V .

k-anonymity as a decision problem: for a given table V and an integer number $l \in \mathbb{N}$, is there a suppressor t which makes V k-anonymous and suppresses maximum l coordinates?

It was claimed that if the alphabet size is unlimited, optimal k-anonymity is a hard problem for $k \geq 3$.

Theorem: k-anonymity is NP-Hard for $k \geq 3$ even $|\Sigma| \geq |V|$.

A reduction from k-dimensional perfect matching: let $H = (U, E)$ be a k-hypergraph, n and m be the number of vertices and the number of edges, respectively. In this case, in n/k hyperedges, is there a subset $S \in E$, such that each vertex of U is covered by one hyperedge of S ?

Assume H is a k-dimensional simple hypergraph, $U = \{u_1, \dots, u_n\}$ and $E = \{e_1, \dots, e_m\}$ are vertices and edges of H , respectively and finally $\Sigma = \{0, 1, \dots, n\}$ is the alphabet. Table V can be constructed as below and for each

u_i ; m dimensional vector $v_i \in \Sigma^m$ can be defined as;

$$v_i[j] := \begin{cases} 0 & \text{if } u_i \in e_j, \\ 1 & \text{otherwise.} \end{cases}$$

Suppose V includes the series of v_i such as $V := \{v_1, \dots, v_n\}$. Assume that t suppresses minimum number of vector coordinates and ensures k-anonymity. It was claimed that the number of coordinates suppressed by t is maximum $n(m-1)$ if there exists a k-dimensional perfect matching in H .

This claim was proved for $k = 3$. Firstly, it was accepted that there exists a perfect 3-dimensional matching M in H . For $i = 1, \dots, n$, let $j(i)$ be the indices of $e_{j(i)}$ which is the unique hyperedge in M containing vertex u_i .

Suppressor t is defined as follows;

$$t(v_i)[j'] := \begin{cases} 0 & \text{if } j' = j(i), \\ * & \text{otherwise.} \end{cases}$$

Because of u_i is on the hyperedge $e_{j(i)}$, the following states occur by definition, $v_i[j(i)] = 0$ and all other coordinates are *. Therefore, t is a suppressor on V .

Consider any three nodes $u_i, u_{i'}, u_{i''}$ on $e_{j(i)}$ and each node has identical anonymized vectors such as $t(v_i) = t(v_{i'}) = t(v_{i''})$. Therefore, $t(V)$ has three identical vectors, which shows that $t(V)$ is 3-anonymous.

Since every $t(v) \in t(V)$ has at most one non-* coordinate, the value of the solution is $n(m-1)$. Therefore, the optimum solution of 3-anonymity includes at most $n(m-1)$ number of *'s in these vectors.

B. THE PROOF VIA EDGE PARTITION INTO TRIANGLES

In [19] and [20], it was shown that k-anonymity is NP-Hard for $k \geq 3$. The authors employed edge partition into triangles reduction. To understand the proof in details, some notations used in the papers were presented in Table 2.

Theorem: k-anonymity is NP-Hard even for $\Sigma = \{0, 1, 2\}$.

Edge partition into triangles reduction: for a given graph $G = (V, E)$ with $|E| = 3m$ for any integer m , can the edges of G be partitioned into m triangle whose edges are disjoint?

The proof starts with the reduction of edge partition into triangles and 4-stars before explaining the reduction of edge partition into triangles.

The proof has two phases. In the first phase, it is shown that for a graph $G = (V, E)$ with $|E| = 3m$, if and only if G can be partitioned into triangle and 4-stars, the optimal 3-anonymity solution for T is $9m$. In the second phase, it is shown that if and only if G can be partitioned into m disjoint-triangles, the 3-anonymity solution is maximum $9m \lceil \log_2(3m) \rceil$ for table T' .

Edge partition into triangles and 4-stars: for a given graph $G = (V, E)$ with $|E| = 3m$ and $|V| = n$, a table T with $3m$ row and n attribute is created. The row $r_{a,b}$ corresponding to edge (a, b) has 1 in positions corresponding to attributes a and b , and 0 otherwise. Assume that graph G can be partitioned into m disjoint-triangles and 4-stars. Let a, b and c be the

TABLE 2. A descriptive representation of the notations used in [19], [20].

Notation	Description
m	number of disjoint triangles
Σ	alphabet of attribute values
T	preliminary table for vertex-edge relationship
T'	replication of T for 4-stars
$r_{a,b}$	the row corresponding to edge (a, b)
a, b, c	vertices of triangles
d	central vertices
*	symbol used in suppression
t	number of blocks
n	number of columns
G	complete graph
V	set of vertexes
E	set of edges
v	common vertex
e	edges in $E = \{e_1, \dots, e_n\}$
$conf_i$	Configuration, $i = \{0,1\}$
i	indices of blocks
k	minimum number of cardinality in an equivalence class

vertices of triangle. In the rows $r_{a,b}$, $r_{b,c}$ and $r_{a,c}$, by suppressing the positions of a , b and c , three identical rows, each containing 3 *s and 0, are obtained. Now consider a 4-star has the vertices a, b, c, d and edges (a, d) , (b, d) , (c, d) . In the rows $r_{a,d}$, $r_{b,d}$, $r_{c,d}$, by suppressing corresponding positions of a, b and c , three identical rows each containing 3 *s, with a single 1 and 0 anywhere else. Hence, for every triangle and 4-stars, 3 identical generalized records are obtained by suppressing 9 cells. In conclusion, table T is 3-anonym with cost $9m$. Figure 1 shows an illustration about this phase.

Edge partition into triangles: Assume $t = 1 + \lceil \log_2(3m) \rceil$. Let T' is a table whose each row has t blocks and n columns. For an arbitrary order of edges in E , the rank of an edge $e = (a, b)$ can be expressed in a binary form such as b_1, \dots, b_t . In tuples corresponding to edge e , each block is 0 except a and b . Any block can be in two configurations based on the values of a and b . $conf_0$: 1 in position a and 2 in position b , $conf_1$: 2 in position a and 1 in position b . i^{th} block in corresponding row of e , has $conf_{b_i}$. For example, the edges in Figure 1 are ranked from 1 to 6 and T' is presented in Figure 2. It can be understood that if and only if E is partitioned into m disjoint-triangles, optimal 3-anonymity cost is at most $9mt$ for T' .

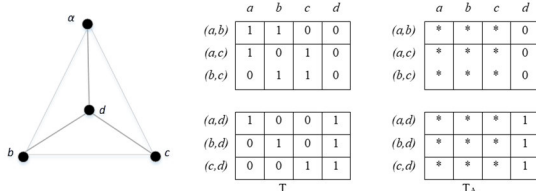


FIGURE 1. Anonymizing of triangle and 4-stars obtained from graph.

001	(a,b)	1	2	0	0	1	2	0	0	2	1	0	0
010	(a,c)	1	0	2	0	2	0	1	0	1	0	2	0
011	(b,c)	0	1	2	0	0	2	1	0	0	2	1	0
100	(a,d)	2	0	0	1	1	0	0	2	1	0	0	2
101	(b,d)	0	2	0	1	0	1	0	2	0	2	0	1
110	(c,d)	0	0	2	1	0	0	2	1	0	0	1	2
001	(a,b)	*	*	*	0	*	*	*	0	*	*	*	0
010	(a,c)	*	*	*	0	*	*	*	0	*	*	*	0
011	(b,c)	*	*	*	0	*	*	*	0	*	*	*	0
100	(a,d)	*	*	*	1	*	*	*	*	*	*	*	*
101	(b,d)	*	*	*	1	*	*	*	*	*	*	*	*
110	(c,d)	*	*	*	1	*	*	*	*	*	*	*	*

FIGURE 2. Anonymized table obtained from graph in Figure 1.

III. A NEW PROOF BASED ON GRAPH COLORING

A new proof based on graph coloring was presented in this section. In the literature, there are some studies using different types of NP-Complete problems for reductions. In this study, graph 3-coloring problem was preferred to use for the reduction. The aim is to examine whether graph coloring problem can be used for a reduction to prove the NP-Hardness of k-anonymity and whether it enables us to improve both the alphabet size and the number of suppressed cells. In addition, graph coloring based representation also presents simplicity for a better understanding.

Garey and Johnson [26], [27] proved that graph 3-coloring with no vertex degree exceeding 4 is NP-Complete. In our study, we borrowed and adopted this idea and then investigated the availability of graph 3-coloring problem to prove the NP-Hardness of k-anonymity.

A. REVISIT OF THE PROOF OF GAREY AND JOHNSON

Garey and Johnson [26], [27] proposed to restrict the maximum vertex degree of graph will be colored. If the maximum vertex degree is restricted with a small enough degree, many graph coloring problem is solved in polynomial time. Hence, finding the most powerful constraint of vertex degree that will keep the problem in NP-Complete is very important. Table 3 presents the complexity classes of problems that a subproblem can be belong to any of them with a degree constraint.

TABLE 3. Classification of subproblems based on complexity classes that any d vertex degree limited graph can be belong to.

Problem	P ($D \leq$)	NP-Complete ($D \geq$)
Exact Cover	2	3
Hamilton Cycle	2	3
Graph 3-colorability	3	4
Feedback vertex set	2	3

In this context, maximum vertex degree in graph 3-coloring problem is presented as 4. This implies that subproblem is still in NP-Complete class even the constraint of vertex degree of graph 3-coloring problem is determined as 4. Hence, in the proof, maximum vertex degree is limited with 4.

In order to prove the results of degree limited NP-Completeness, vertex substitute approach is used. Vertex substitution is defined as substituting a vertex with a subgraph that meets some certain criteria.

Theorem: Graph 3-colorability with no vertex degree exceeding 4 is an NP-Complete problem.

Proof: Assume $G = (V, E)$ is an arbitrary graph of a general problem, $G' = (V', E')$ is a restricted instance of G that no vertex of G' have the degree exceeding 4. If and only if G is 3-colorable, G' is also 3-colorable.

In this vertex substitution approach, a graph with eight vertex is considered. In Figure 3.a, graph H_3 has three outlets with labels 1, 2 and 3. Let x be the number of outlets, for $x \geq 4$, H_x which is a vertex substitution with x number of outlets is formed with adjoining a copy of H_3 to substitution H_{x-1} . H_5 , which is used to prove the NP-Completeness of graph 3-colorability with degree restriction, is presented in Figure 3.b. In these graphs, all outlets have the same color.

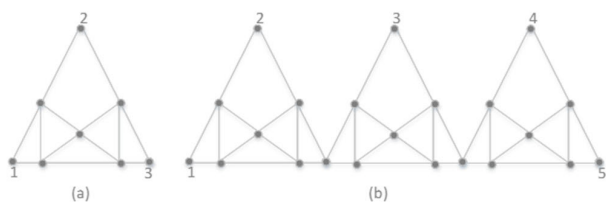


FIGURE 3. Graph H_3 (a) and vertex substitution H_5 (b).

As indicated in Figure 3, for $x \geq 3$, the following situations appear;

1. H_x has $7(x - 2) + 1$ number of outlets, including x number of labelled outlets.
2. H_x has no vertex whose degree exceed 4.
3. The degree of each outlet of H_x is 2.
4. H_x is 3-colorable, but not 2-colorable. In every different 3-coloring way, outlets of H_x have the same color.

Assume an arbitrary graph G , composing from s number of vertices v_1, v_2, \dots, v_s and containing vertexes with degree exceeding 4, composed from graph array as shown below;

$$G = G_0, G_1, \dots, G_s = G'$$

Each $G_l, 1 \leq l \leq s$, is constructed from G_{l-1} . Let d be the degree of v_l in G_{l-1} and $\{v_1, v_l\}, \{v_2, v_l\}, \dots, \{v_3, v_l\}$ be the edges containing v_l . To form G_l, v_l is deleted from G_{l-1} and is replaced with a copy of H_d . Each edge $\{v_m, v_l\}$ is replaced with an edge joining outlet m and v_m . In the new construction, for $0 \leq x \leq s$, the number of vertexes of G_x exceeding 4 is $s - x$, if and only if graph G is 3-colorable, then G_x is 3-colorable. Therefore, $G' = G_s$ is obtained.

The overall approach is illustrated in Figure 4. In graph G , if there exist any vertex whose degree exist 4, by replacing those vertices with vertices in G', G becomes 3-colorable. Because of the graph, which belongs to a general problem and was showed in Figure 4.a, has vertexes with degree exceeding 4, it is replaced with the substitution in Figure 4.b. Thus, graph G becomes 3-colorable. As a result, graph 3-colorability problem with 4-degree restriction is NP-Complete problem.

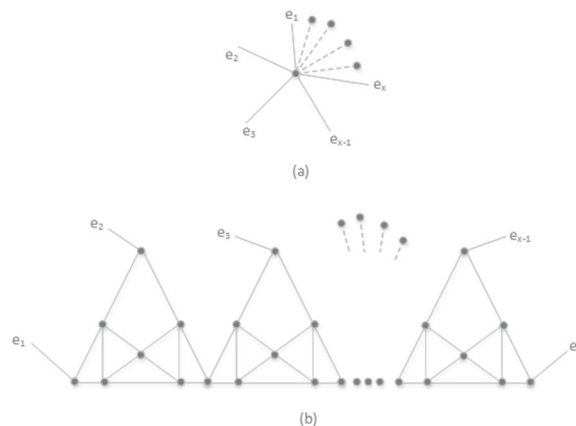


FIGURE 4. An example of vertex with degree exceeding 4 (a) and structure of vertex substitution (b).

B. THE PROPOSED PROOF BASED ON GRAPH COLORING

The NP-Hardness of k-anonymity for $k \geq 3$ was proved using a reduction from degree-limited graph 3-coloring. Some notations used in the proof was presented in Table 4.

Definition: Let t be a map from R to $(\sum \cup \{*\})^m$. For all $r \in R$ and $j = 1, \dots, m$; if $t(r)[j] \in \{r[j], *\}$ then t is a suppressor function.

Each vector $r \in R$ has corresponding $t(r) = r'$ in anonymized table $R' \subseteq (\sum \cup \{*\})^m$. In addition, the coordinates of r' are the same as with the coordinates of r .

In order to work on vector sets in R, t can be extended as follows. $t(R)$ is accepted as a multiple set when one or more vectors in r is mapped to the same suppressed vectors. For instance, $r_1 \neq r_2 \neq r_3 \in R$ but when t is applied, $t(r_1) = t(r_2) = t(r_3)$ is obtained.

Definition: Let t be a suppressor on $R = \{r_1, \dots, r_n\} \subseteq \sum^m$. $t(R)$ is k-anonymous if for all $r_i \in R$, there exist $k - 1$ number of indices $i_1, i_2, \dots, i_{k-1} \in \{1, \dots, n\}$ providing $t(r_{i_1}) = t(r_{i_2}) = \dots = t(r_{i_{k-1}}) = t(r_i)$. Therefore, t can be called as a k-anonymizer on R .

k-anonymity: for a given table R and an integer $l \in N$, is there a suppressor t which makes table R k-anonymous and suppresses maximum l number of coordinates?

Theorem: k-anonymity is NP-Hard for $k \geq 3$ even $\sum = \{0, 1\}$.

A reduction from degree-limited graph 3-coloring: given a graph $G = (V, E)$, is G 3-colorable such that for every edge $e_{i,j} \in E$, the color c_i of v_i is different from the color c_j of v_j and there exist no vertex degree exceeding 4.

We investigated this situation for $|V| = 6w$, for any integer w . Assume G is a simple graph, $V = \{v_1, \dots, v_n\}$ and $E = \{e_{1,2}, \dots, e_{q,z}\}$ are vertices and edges of G , respectively and $\sum = \{0, 1\}$ is the alphabet. A table R can be created as follows. For every $e_{i,j}$, m-dimensional $r_i \in \sum^m$ vector can be defined as;

$$r_i[j] := \begin{cases} 1 & \text{if } e_{i,j} \in E, \\ 0 & \text{otherwise.} \end{cases}$$

TABLE 4. A descriptive representation of the notations used in the proposed proof.

Notation	Definition
m	dimension
Σ	alphabet of attribute values
R	table (adjacency matrix), $R \subseteq \Sigma^m$
t	suppressor function
r	record in $R = \{r_1, \dots, r_n\}$
r'	suppressed record in $R' = \{r'_1, \dots, r'_n\}$
R'	anonymized version of R
i	row indices $\{i_1, \dots, i_n\}$
j	column indices $\{j_1, \dots, j_m\}$
n	number of records
k	minimum cardinality in an equivalence class
*	symbol used in suppression
l	maximum number of suppressed vector coordinates in R'
G	graph
V	set of vertices
E	set of edges
C	set of colors
v_i	vertices in $V = \{v_1, \dots, v_n\}$
$e_{i,j}$	edges in $E = \{e_{1,2}, \dots, e_{q,z}\}$
$M_{x,y,z}$	candidate matrix of corresponding coordinates of x, y, z
x, y, z	some indices
P	color sharing table
P'	anonymized version of P
c_i	colors in $C = \{c_1, \dots, c_r\}$

Set $R := \{r_1, \dots, r_n\}$. The relations between each vertices can be obtained through R . A color sharing table P can be created as follows. Firstly, one-complement of R is taken and then principal diagonal elements are assigned to zeros. Hereby, P facilitates us to determine which vertices cannot share the same color with v_i . P can be obtained as follows;

$$P = \text{diag}(\bar{R})_{\text{zeros}}$$

Assume that t suppresses minimum number of vector coordinates and provides k-anonymity. In our proof, it is claimed that if and only if G is 3-colorable, the number of suppressed vector coordinates by t is at most $n(m - 3)$.

This claim was proved for $k = 3$. Firstly, it is assumed that G is 3-colorable and has $6w$ number of vertices for any integer w , and V can be partitioned into some disjoint groups each containing triple vertices with colors c_1, c_2, c_3 .

For $0 \leq x, y, z \leq n$, if any x, y, z coordinates of any 3 vertices are zeros, then these parts are left as they are, otherwise they are replaced with *s. In other words, if $M_{x,y,z}$ is a non-overlapping zeros matrix, it saves the original form, but if it is not then it is suppressed with *s. In this case, candidate matrix $M_{x,y,z}$ and suppressor t on $M_{x,y,z}$ can be defined as;

$$M_{x,y,z} = P_{x,y,z} \begin{bmatrix} r_{x,x} & r_{x,y} & r_{x,z} \\ r_{y,x} & r_{y,y} & r_{y,z} \\ r_{z,x} & r_{z,y} & r_{z,z} \end{bmatrix}$$

$$t(M_{x,y,z}) := \begin{cases} 0 & \text{if } M_{x,y,z} = 0, \\ * & \text{otherwise.} \end{cases}$$

Because of each v_i meets the condition of 3-colorability, the following states occur by definition, each candidate matrix $M_{x,y,z} = 0$ and all other coordinates is *s. Hence, t is a suppressor on P .

Consider graph G given in Figure 5. Assume three vertices v_1, v_2, v_3 of graph G is colored with three different colors c_1, c_2, c_3 . If P is anonymized, we have the same anonymized vectors $t(r_1) = t(r_2) = t(r_3)$. Therefore, $t(P)$ contains three identical vectors with respect to v_1, v_2, v_3 , and this situation shows that $t(P)$ is 3-anonymous. $t(P)$ includes 9 non-*s in each triple anonymized vector. Therefore, an optimal 3-anonymous solution has at most $n(m - 3)$ number of *s.

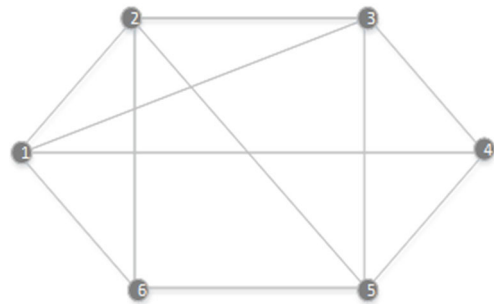


FIGURE 5. A sample graph G .

With regard to Figure 5, Table R can be created as presented in Table 5.

TABLE 5. Table R (or adjacency matrix).

	1	2	3	4	5	6
1	0	1	1	1	0	1
2	1	0	1	0	1	1
3	1	1	0	1	1	0
4	1	0	1	0	1	0
5	0	1	1	1	0	1
6	1	1	0	0	1	1

If we take one complement of R and then change the values of principal diagonal elements with zeros, we obtain color sharing table P as presented in Table 6.

Table 6 guides us to obtain the following statements. In Figure 5, v_3, v_4 and v_5 cannot share the same color and

TABLE 6. Color sharing table P.

	1	2	3	4	5	6
1	0	0	0	0	1	0
2	0	0	0	1	0	0
3	0	0	0	0	0	1
4	0	1	0	0	0	1
5	1	0	0	0	0	0
6	0	0	1	1	0	0

each of these vertices is colored with one out of three different colors. Similarly, v_1, v_2 and v_6 have the same condition, and also there may be many other different selections. Table 6 shows that G can be divided into two groups and each group includes exactly three elements with three different colors. For this example, group one contains 1, 2 and 6 while group two includes 3, 4 and 5. A possible 3-coloring of G was presented in Figure 6.

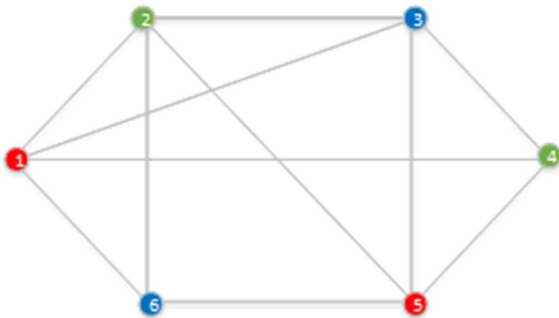


FIGURE 6. A possible 3-coloring of G.

Table 7 indicates the anonymized version of P . Therefore, as it was stated in the proof, maximum number of *s is $n(m - 3)$.

TABLE 7. Anonym table P'.

	1	2	3	4	5	6
1	0	0	*	*	*	0
2	0	0	*	*	*	0
3	*	*	0	0	0	*
4	*	*	0	0	0	*
5	*	*	0	0	0	*
6	0	0	*	*	*	0

Within the context of existing assumptions, the number of suppressed cells of 3-anonymization of G is at most $n(m - 3)$. Since the previous proofs present $n(m - 1)$, $n(m - 1) \lceil \log_2(n(m - 1)/3) \rceil$ and nm number of suppressed cells, respectively, our proof reduces it to $n(m - 3)$. Hence, the average information losses for each result are obtained as follows;

$$\frac{n(m - 3)}{nm} < \frac{n(m - 1)}{nm} < \frac{n(m - 1) \lceil \log_2(n(m - 1)/3) \rceil}{nm} < \frac{nm}{nm}$$

It can be clearly seen that our proof presents the minimum information loss and alphabet size in comparison with other previous proofs and this result may be a good reason for employing graph coloring approach for reductions.

In Table 8, we listed a number of studies available in the literature on the hardness of k-anonymity. We tabularized these studies based on some criteria such as reduction methods, alphabet sizes and average information loss. The results show that our proof provides acceptable outcomes and better results.

TABLE 8. Comparison table for the proposed proof and the previous studies (N/D* = No certain value is defined).

Paper	Reduction Method	Alphabet Size	Average Information Loss
Meyerson and Williams [18]	k-dimensional perfect matching	n	$\frac{n(m - 1)}{nm}$
Aggarwal et al. [19, 20]	edge partition into triangles	3	$\frac{n(m - 1) \lceil \log_2(n(m - 1)/3) \rceil}{nm}$
Bonizzoni et al. [22]	minimum vertex cover on cubic graph	2	N/D*
Scott et al. [24]	maximum 3-dimensional matching with 3 occurrences	N/D*	N/D*
Sun et al. [21]	edge partition into 4-cliques	2	$\frac{nm}{nm}$
Chen et al. [28]	vertex cover	N/D*	N/D*
Our proof	graph 3-coloring	2	$\frac{n(m - 3)}{nm}$

IV. CONCLUSION AND DISCUSSION

As introduced previously, this paper focuses on the computational complexity of k-anonymity and introduces a new approximation approach. Since k-anonymity is an NP-Hard problem and optimal solutions cannot be achieved in a reasonable time, near-optimal solutions are always required. To prove the NP-Hardness of k-anonymity, especially graph problems are employed for reduction frequently.

To the best of our knowledge, this article proved the NP-Hardness of k-anonymity using a reduction from degree-limited graph 3-coloring for the first time. We also improved both the alphabet size and the average information loss in comparison with some previous proofs which were listed in Table 8. The results showed that reductions utilize graph coloring presents better results than the others. However, in the future, other NP-Complete problems can be examined in terms of whether they present better results.

REFERENCES

[1] A. S. M. T. Hasan and Q. Jiang, "A general framework for privacy preserving sequential data publishing," in *Proc. 31st Int. Conf. Adv. Inf. Netw. Appl. Workshops (WAINA)*, Mar. 2017, pp. 519–524.

- [2] M. M. Almasi, T. R. Siddiqui, N. Mohammed, and H. Hemmati, "The risk-utility tradeoff for data privacy models," in *Proc. 8th IFIP Int. Conf. New Technol., Mobility Secur. (NTMS)*, Nov. 2016, pp. 1–5.
- [3] X. Chen and V. Huang, "Privacy preserving data publishing for recommender system," in *Proc. IEEE 36th Annu. Comput. Softw. Appl. Conf. Workshops*, Jul. 2012, pp. 128–133.
- [4] S. D. Warren and L. D. Brandeis, "The right to privacy," *Harvard Law Rev.*, vol. 1, pp. 193–220, Jul. 1890.
- [5] W. Fang, X. Z. Wen, Y. Zheng, and M. Zhou, "A survey of big data security and privacy preserving," *IETE Tech. Rev.*, vol. 34, no. 5, pp. 544–560, Sep. 2017.
- [6] M. Chibba and A. Cavoukian, "Privacy, consumer trust and big data: Privacy by design and the 3 C's," in *ITU Kaleidoscope, Trust in the Information Society*. Barcelona, Spain: IEEE Press, Dec. 2015, doi: 10.1109/Kaleidoscope.2015.7383624.
- [7] P. Jain, M. Gyanchandani, and N. Khare, "Big data privacy: A technological perspective and review," *J. Big Data*, vol. 3, no. 1, p. 25, Dec. 2016.
- [8] J. J. V. Nayahi and V. Kavitha, "Privacy and utility preserving data clustering for data anonymization and distribution on Hadoop," *Future Gener. Comput. Syst.*, vol. 74, pp. 393–408, Sep. 2017.
- [9] Q. Tang, Y. Wu, S. Liao, and X. Wang, "Utility-based k -Anonymization," in *Proc. 6th Int. Conf. Networked Comput. Adv. Inf. Manage.*, Aug. 2010, pp. 318–323.
- [10] B. C. Fung, K. Wang, A. W. Fu, and S. Y. Philip, *Introduction to Privacy-Preserving Data Publishing: Concepts and Techniques*. Boca Raton, FL, USA: CRC Press, 2010.
- [11] L. Sweeney, " K -Anonymity: A model for protecting privacy," *Int. J. Uncertainty, Fuzziness Knowl.-Based Syst.*, vol. 10, no. 5, pp. 557–570, Oct. 2002.
- [12] A. Machanavajjhala, J. Gehrke, D. Kifer, and M. Venkatasubramanian, " L -Diversity: Privacy beyond k -anonymity," in *Proc. Int. Conf. Data Eng.*, 2006, pp. 24–35.
- [13] N. Li, T. Li, and S. Venkatasubramanian, " T -closeness: Privacy beyond k -Anonymity and l -Diversity," in *Proc. IEEE 23rd Int. Conf. Data Eng.*, Apr. 2007, pp. 106–115.
- [14] S. A. Abdelhameed, S. M. Moussa, and M. E. Khalifa, "Privacy-preserving tabular data publishing: A comprehensive evaluation from web to cloud," *Comput. Secur.*, vol. 72, pp. 74–95, Jan. 2018.
- [15] Q. Gong, J. Luo, M. Yang, W. Ni, and X.-B. Li, "Anonymizing 1: M microdata with high utility," *Knowl.-Based Syst.*, vol. 115, pp. 15–26, Jan. 2017.
- [16] Y. Tao, Y. Tong, S. Tan, S. Tang, and D. Yang, "Protecting the publishing identity in multiple tuples," in *Proc. Annu. Conf. Data Appl. Secur. Privacy*, 2008, pp. 205–218.
- [17] G. Poulis, G. Loukides, A. Gkoulalas-Divanis, and S. Skiadopoulos, "Anonymizing data with relational and transaction attributes," in *Proc. Joint Eur. Conf. Mach. Learn. Knowl. Discovery Databases*, 2013, pp. 353–369.
- [18] A. Meyerson and R. Williams, "On the complexity of optimal k -anonymity," in *Proc. 23rd ACM SIGMOD-SIGACT-SIGART Symp. Princ. Database Syst.*, Jun. 2004, pp. 1–20.
- [19] G. Aggarwal, T. Feder, K. Kenthapadi, R. Motwani, R. Panigrahy, D. Thomas, and A. Zhu, "Approximation algorithms for k -Anonymity," *J. Privacy Technol.*, vol. 1, pp. 1–18, Aug. 2005.
- [20] G. Aggarwal, T. Feder, K. Kenthapadi, R. Motwani, R. Panigrahy, D. Thomas, and A. Zhu, "Anonymizing tables," in *Proc. Int. Conf. Database Theory*, 2005, pp. 246–258.
- [21] X. Sun, L. Sun, and H. Wang, "Extended k -Anonymity models against sensitive attribute disclosure," *Comput. Commun.*, vol. 34, no. 4, pp. 526–535, Apr. 2011.
- [22] P. Bonizzoni, G. Della Vedova, and R. Dondi, "Anonymizing binary and small tables is hard to approximate," *J. Combinat. Optim.*, vol. 22, no. 1, pp. 97–119, Jul. 2011.
- [23] J. Blocki and R. Williams, "Resolving the complexity of some data privacy problems," in *Proc. Int. Colloq. Automata, Lang., Program.*, 2010, pp. 393–404.
- [24] A. Scott, V. Srinivasan, and U. Stege, " K -Attribute-Anonymity is hard even for $k=2$," *Inf. Process. Lett.*, vol. 115, no. 2, pp. 368–370, Feb. 2015.
- [25] K. LeFevre, D. J. DeWitt, and R. Ramakrishnan, "Mondrian multidimensional k -Anonymity," in *Proc. 22nd Int. Conf. Data Eng. (ICDE)*, 2006, pp. 1–12.
- [26] M. R. Garey and D. S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Completeness*. San Francisco, CA, USA: Freeman, 1979. [Online]. Available: <https://dl.acm.org/doi/10.5555/578533>
- [27] M. R. Garey, D. S. Johnson, and L. Stockmeyer, "Some simplified NP-complete graph problems," *Theor. Comput. Sci.*, vol. 1, no. 3, pp. 237–267, Feb. 1976.
- [28] R. Chen, B. C. M. Fung, N. Mohammed, B. C. Desai, and K. Wang, "Privacy-preserving trajectory data publishing by local suppression," *Inf. Sci.*, vol. 231, pp. 83–97, May 2013.



YAVUZ CANBAY received the Ph.D. degree from the Department of Computer Engineering, Gazi University, in 2019. Currently, he is an Assistant Professor with Kahramanmaraş Sütçü İmam University. He is also the Director of the Data Vision Laboratory (DEVLAB). His main research interests include data privacy, information security, and artificial intelligence.

• • •