

Why Waldo befriended the dummy? k -Anonymization of social networks with pseudo-nodes

Sean Chester · Bruce M. Kapron · Ganesh Ramesh ·
Gautam Srivastava · Alex Thomo · S. Venkatesh

Received: 5 December 2011 / Revised: 28 June 2012 / Accepted: 4 September 2012 / Published online: 26 September 2012
© Springer-Verlag 2012

Abstract For a graph-based representation of a social network, the identity of participants can be uniquely determined if an adversary has background structural knowledge about the graph. We focus on degree-based attacks, wherein the adversary knows the degrees of particular target vertices and we aim to protect the anonymity of participants through k -anonymization, which ensures that every participant is equivalent to at least $k - 1$ other participants with respect to degree. We introduce a natural and novel approach of introducing “dummy” participants into the network and linking them to each other and to real

participants in order to achieve this anonymity. The advantage of our approach lies in the nature of the results that we derive. We show that if participants have labels associated with them, the problem of anonymizing a subset of participants is NP-Complete. On the other hand, in the absence of labels, we give an $\mathcal{O}(nk)$ algorithm to optimally k -anonymize a subset of participants or to near-optimally k -anonymize all real and all dummy participants. For degree-based-attacks, such theoretical guarantees are novel.

Keywords Privacy · k -Anonymization · Social networks · Complexity · Dynamic programming

A preliminary, short version (Chester et al. 2011) of this paper appeared at ADBIS 2011.

S. Chester (✉) · B. M. Kapron · G. Srivastava · A. Thomo ·
S. Venkatesh
Department of Computer Science, University of Victoria,
PO Box 3055, STN CSC, Victoria, BC V8W 3P6, Canada
e-mail: schester@uvic.ca

B. M. Kapron
e-mail: bmkapron@cs.uvic.ca

G. Srivastava
e-mail: gsrivast@uvic.ca

A. Thomo
e-mail: thomo@cs.uvic.ca

S. Venkatesh
e-mail: venkat@cs.uvic.ca

G. Ramesh
Yahoo! Inc., 4401 Great America Parkway,
Santa Clara, CA 95054, USA
e-mail: rganesh@yahoo-inc.com

1 Introduction

The advent of the social web has brought about an explosive increase in the pervasiveness of large-scale social networks. Embedded within these networks is a wealth of information of multidisciplinary interest for the industrious analyst. However, unlike many other data, this information is about *people*; so, it is only ethical to respect their right to privacy. As a data owner, one must respect that some user groups expect a level of privacy and so that privacy must be provided to that subgroup if the data is to be shared with external organizations, if not the entire network. The study by Ferri et al. (2012), for example, reveals that although some user groups are less concerned by data owners sharing data about them, up to 90 % of members in other groups disagree with the principle.

If the social network is “sanitized” prior to release, the privacy of individuals within it can be protected while satisfying the needs of analysts. We assume the network is represented as a vertex-labelled graph $\mathcal{G} = (V, E, \Sigma, \ell)$, with vertices representing participants and edges representing the

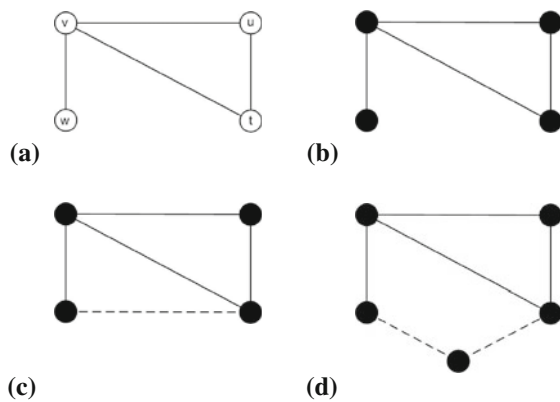


Fig. 1 An example of 2-degree-anonymity by means both of edge- (as in literature) and vertex- (as in this paper) addition. **a** An example unlabelled social network graph with identifiers for each vertex. **b** The same unlabelled social network graph with the identifiers blacked out. Vertices v and w can still be uniquely identified because of their degree. **c** An optimal k -degree-anonymization of the example graph by means of edge addition. No longer can vertices v and w be uniquely identified by degree. **d** An optimal k -degree-anonymization of the example graph by means of vertex addition. As before, vertices v and w cannot be uniquely identified by degree, but the method by which this property was achieved is different.

relationships between them.¹ As shown by Backstrom et al. (2007), removing the identifiers from the graph is insufficient to guarantee privacy: an adversary can still reveal the identity of individuals by exploiting background structural knowledge that he has about the network. Figure 1a and b illustrates how degree, the structural property we study, can be sufficient to uniquely identify vertices: v is the only person with exactly three relationships, so stripping his identifier as in Fig. 1b does not provide him any privacy against an adversary who knows his degree.

We pursue the well-known alternative of k -anonymizing the social network data prior to release. If an adversary has some knowledge of a structural property \mathcal{P} of an arbitrary individual u who he seeks to identify, one can conceal the identity of u —and, indeed, everyone else—by establishing that every individual in the network is identical to at least $k - 1$ others with respect to \mathcal{P} . In this way, the knowledge of the adversary renders at best a $1/k$ chance of uniquely identifying u . Stated informally, the k - \mathcal{P} -anonymization problem is thus: given a vertex-labelled graph \mathcal{G} and a structural property \mathcal{P} , transform \mathcal{G} into a k - \mathcal{P} -anonymous graph \mathcal{G}' such that \mathcal{G}' is as close as possible to \mathcal{G} . (Fig 1c, d is the example of k - \mathcal{P} -anonymous graphs for $\mathcal{P} = \text{degree}$.)

Certainly, there are many possible structural properties \mathcal{P} of which the adversary may have some knowledge.

Many of these have been studied in literature (e.g., Liu and Terzi 2008; Zhou and Pei 2011; Thompson and Yao 2009; Wu et al. 2010) and appropriate k -anonymization techniques have been developed and shown to work well in experimental settings. Nonetheless, there are some pertinent, fundamental questions for any choice of \mathcal{P} :²

- Is it possible to produce an optimal k - \mathcal{P} -anonymous graph in polynomial time?
- If not, can theoretical guarantees, such as bounds on approximation ratios, be derived for non-optimal algorithms?
- Can better theoretical guarantees be determined for special classes of graphs?

For any choice of adversarial knowledge, \mathcal{P} , these appear to be quite challenging problems. Consider the seemingly simple adversary of study in this work who knows only the degree of his target. Achieving k -degree anonymity to thwart him was studied by Lui and Terzi (2008), who derived an experimentally effective algorithm. Meanwhile, the works of Chester et al. (2012b) and of Zhou and Pei (2011) offer hardness results of k - \mathcal{P} -anonymization for other choices of \mathcal{P} . Still, none of these three most closely related works answer the questions above for k -degree-anonymization on arbitrary graphs, implying that the problems are still quite non-trivial for this adversary.

Thus, here, we introduce a new approach to k - \mathcal{P} -anonymization, changing both the way one goes about constructing \mathcal{G}' and, correspondingly, the optimisation condition. In previous work, the problem formulation has been to transform $\mathcal{G} = (V, E)$ into a k -anonymous graph $\mathcal{G}' = (V, E \cup E')$; that is to say, only edges are added to \mathcal{G} in order to construct \mathcal{G}' . (see Fig. 1c.) The optimisation condition, then, is to minimise $|E'|$. We embark here on the very natural study of adding new vertices as well, transforming $\mathcal{G} = (V, E)$ to $\mathcal{G}' = (V \cup V', E \cup E')$, and minimising $|V'|$.³ (see Fig. 1d.) For this optimisation problem, we require that the new edges must have a new vertex as an endpoint ($E' \subseteq V' \times (V \cup V')$); otherwise, the complete graph $K_{|V|}$ offers an optimal solution with $|V'| = 0$, which is clearly not what we want. This alternative formulation presents notable utility. Ultimately, the intent of releasing data is to facilitate analysis, and the analysis is conducted at the aggregate level. Introducing few new vertices with similar characteristics to those already in the network could quite accurately preserve the aggregate characteristics of the network. Some network characteristics, such as the

¹ We define a vertex-labelled graph as the four-tuple (V, E, Σ, ℓ) , where V is a vertex set, $E \subseteq V \times V$ is a set of undirected edges, Σ is a set of sensitive labels, and $\ell : V \rightarrow \Sigma$ is a labelling function that assigns a label to each vertex. We discuss in the paper two types of labels, sensitive and identifying. By Σ , we refer to the former, assuming the latter is stripped from the graph.

² We mention specific cases in which these questions have been answered in our discussion of related work in Sect. 6 Even in these cases, however, not all three questions have been fully addressed.

³ Precise formulations of the problem appear in Sect. 2 for unlabelled graphs and in Sect. 5 for labelled graphs.

number of large cliques, are exactly preserved as a consequence of our constraint on E' .

Perhaps more importantly, however, is that by changing the problem formulation, we are able to address the aforementioned three problems for a degree-based attack. In particular, we prove hardness for k -anonymization against degree-based attacks on arbitrary, vertex-labelled graphs using this vertex-addition model (Theorem 2). By fixing the size of the label set to one (i.e., considering unlabelled graphs), we provide an efficient and optimal algorithm for subset anonymization (Corollary 2) and an efficient, near-optimal algorithm for complete anonymization (Theorem 1). Furthermore, for graphs with certain properties that are likely to arise in social networks, our optimality guarantee is improved to being within one vertex of optimal (Corollary 1).

1.1 Our contributions

We introduce a vertex addition approach for k - \mathcal{P} -anonymization. For when \mathcal{P} is degree, we offer the following results:

- We introduce for unlabelled graphs an $\mathcal{O}(k * |V|)$ k -degree-anonymization algorithm based on dynamic programming and prove that, on any arbitrary graph, the minimisation of $|V'|$ is optimal within an additive factor of k . For a special class of graphs that is likely to include social networks, the algorithm is optimal within 1 for reasonable values of k (Sect. 3);
- We conduct an empirical evaluation of our algorithm on several well-known network datasets, demonstrating that it quite largely preserves the utility of the original graph with respect to standard structural parameters like clustering coefficient, average path length and connectivity, even as k approaches percentages of $|V|$ that are quite high for the context of graph anonymization (Sect. 4);
- We demonstrate that for the more general case of labelled graphs ($|\Sigma| > 1$), k -degree-anonymization with a predetermined number of vertex additions (the decision version of the problem) is NP-Complete, by giving a reduction from a known NP-Hard table anonymization problem (Sect. 5).

2 Subset anonymity and related concepts

We begin by introducing the concepts under study in this paper. At a high level, subset anonymity acknowledges that a certain subset of members in a social network may be highly identifiable as a group and, simultaneously, need an assurance of individual anonymity. For example, consider

a social network constructed from a company email corpus where each node of the network corresponds to an email address and each link between two nodes implies an exchange of emails between them. The internal email addresses are likely to have much higher degree and be more sensitive than the external email addresses. The objective of subset anonymity is to ensure that members who belong to an identifiable subset of a social network, in this case the internal email addresses, cannot be distinguished from each other with certainty greater than $1/k$. This way, although they are identifiable as a group, they are not identifiable individually.

We model a social network as an unlabeled, undirected graph, $\mathcal{G} = (V, E)$ ⁴ which contains an *anonymizing subset* $X \subseteq V$ of vertices which need to be anonymized. The objective is to produce a similar graph $\mathcal{G}' = (V \cup V', E \cup E')$ in which X is anonymous. This differs from previous work, not only in the focus on subset anonymity, but also in the permission to introduce new, “fake” vertices (V').

An important consideration is what it means to be *anonymous* in a graph. Because graphs embed copious structural information, anonymity depends on what background structural information an attacker might have. As we detail in Sect. 6, various definitions have been introduced which assume an attacker has progressively more background structural information. Since no work, however, considers the implications of expanding the vertex set V , we embark on this *foundational* study in which we assume the simplest adversarial knowledge, that the attacker knows how many connections his target has, to contrast our conclusions with those derived from previous works which have not permitted an expansion of the vertex set (i.e., require $V' = \emptyset$).

To state the problem formally, we first need to introduce a few definitions. We begin by explaining k -degree-anonymity, the focus of this paper:

Definition 1 The *degree* of a vertex v in a graph $\mathcal{G} = (V, E)$ is the number of neighbours it has, $|\{u \in V : (u, v) \in E\}|$.

For example, the uppermost vertex in Fig. 2b has a degree of 5, because it is connected to 5 other vertices (neighbours). The rightmost node in the same graph has degree 1 because it has only one incident edge. This is the information that we assume an attacker possesses.

Definition 2 A *degree sequence* of a set of vertices V is the sequence (d_1, \dots, d_n) composed by sorting in descending order the degrees of every node in V .

⁴ For simplicity in this section, we regard a graph as a 2-tuple. We note that equivalently, for consistency, we could express an unlabelled graph as $\mathcal{G} = (V, E, \Sigma, \ell)$ where $\exists \sigma \in \Sigma : \forall v \in V, \ell(v) = \sigma$. However, the simpler notation simplifies the exposition.

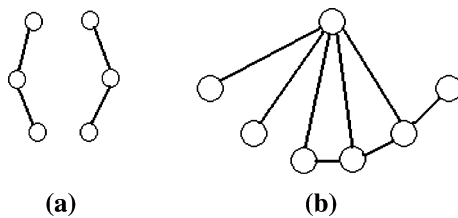


Fig. 2 The two small example graphs used to illustrate our anonymization procedure

Again referring to Fig. 2b, the degree sequence is (5, 3, 3, 2, 1, 1, 1), the degrees of each of the seven vertices sorted in descending order.

Definition 3 A k -partitioning of a degree sequence is a partitioning of the degree sequence into disjoint partitions such that every degree appears in exactly one partition and every partition contains at least k elements.

Two possible 3-partitioning of the degree sequence in our running example are ((5, 3, 3), (2, 1, 1, 1)) and ((5, 3, 3, 2), (1, 1, 1)). There are other possible 3-partitioning of this degree sequence that have non-contiguous partitions, such as ((5, 3, 2, 1), (3, 1, 1)), but we will indicate later in Proposition 1 that these are never better choices for our algorithm.

Definition 4 A degree sequence is k -anonymous if partitioning it into the sets of distinct elements induces a k -partitioning.

That is to say, a degree sequence is k -anonymous if every degree appears at least k times in the sequence. Then, a subset of vertices can be defined as k -degree-anonymous depending on their degree sequence:

Definition 5 A subset of vertices $X \subseteq V$ is said to be k -degree-anonymous in a graph $\mathcal{G}' = (V, E)$ iff the degree sequence of X is k -anonymous. Similarly, a graph is said to be k -degree-anonymous if its entire vertex set is k -degree-anonymous.

This provides sufficient material to define the problem of study here:

Problem Definition Given an input graph $\mathcal{G} = (V, E)$ and an anonymizing subset X , construct a graph $\mathcal{G}' = (V \cup V', E \cup E')$, $E' \cap (V \times V) = \emptyset$, such that X and $X \cup V'$ are both k -degree-anonymous in \mathcal{G}' and $|V'|$ is minimised.

Note that X must be k -degree-anonymous in \mathcal{G}' because it may be readily identifiable. This corresponds to a scenario when X is a small percentage of V and V' is indistinguishable from $V \setminus X$.⁵ In other scenarios, say when users

opt in to X as conceived by Yuan et al. (2010), X could be a large percentage of V and not readily identifiable so we would need that $X \cup V'$ is k -degree-anonymous to hide the “fake” vertices among X rather than among $V \setminus X$. Thus the need for anonymity within both subsets.

As a last note, we introduce a couple definitions that are useful for describing the algorithm in the next section.

Definition 6 The max deficiency of a k -partitioning of a degree sequence is the largest difference between highest and smallest degree within any partition.

For example, taking the two 3-partitionings in our running example, the max deficiency of ((5, 3, 3), (2, 1, 1, 1)) is $\max(5 - 3, 2 - 1) = 2$ and of ((5, 3, 3, 2), (1, 1, 1)) is $\max(5 - 2, 1 - 1) = 3$. Finally,

Definition 7 The total deficiency of a k -partitioning of a degree sequence is the sum over every d_i in the sequence of the difference between d_i and the largest degree in the same partition.

Using the same example we find that the total deficiency of ((5, 3, 3), (2, 1, 1, 1)) is:

$$((5 - 5) + (5 - 3) + (5 - 3) + (2 - 2) + (2 - 1) + (2 - 1) + (2 - 1)) = 7$$

and of ((5, 3, 3, 2), (1, 1, 1)) is:

$$((5 - 5) + (5 - 3) + (5 - 3) + (5 - 2) + (1 - 1) + (1 - 1) + (1 - 1)) = 7.$$

3 An efficient algorithm to near-optimally k -degree-anonymize unlabelled graphs

We present here immediately our most encouraging result, that for the special case of unlabelled graphs, k -degree-anonymization can be solved very near-optimally in linear time. We produce a k -degree-anonymous graph $\mathcal{G}' = (V \cup V', E \cup E')$ from an original. In \mathcal{G}' , we require that all the original vertices, V , are k -degree-anonymous—and in Corollary 2, we will relax this constraint to an input subset $X \subseteq V$. We also require that the new vertices are concealed as well so that they cannot be readily identified and removed from the graph in order to recover \mathcal{G} (i.e., $V \cup V'$ is k -degree-anonymous in \mathcal{G}'). As mentioned in the introduction, we seek to minimise $|V'|$, while maintaining the constraint that $E \subseteq V' \times (V \cup V')$. We will prove the following theorem at the end of this section:

Theorem 1 Our algorithm produces a k -degree-anonymous graph \mathcal{G}' containing the input graph \mathcal{G} as an induced subgraph, using $\mathcal{O}(nk)$ time and $\mathcal{O}(n)$ space. The number of new vertices added is optimal up to an additive factor of k .

⁵ Considering the Enron email corpus on which we experiment in Sect. 4.1, $|V| > 65,000$, but only 151 vertices correspond to internal email addresses.

At a high level, the algorithm proceeds in three stages. At first, we design a precise recursion to group the vertices of V by target degree (the degree they will have in \mathcal{G}'). The recursion establishes a grouping such that the max deficiency, a parameter in determining with how many nodes V must be augmented, is minimised. We evaluate the recursion using dynamic programming with $\mathcal{O}(nk)$ execution cost.

The second stage is to determine precisely how many vertices with which we wish to augment V in order to guarantee that we can k -anonymise all of V' . This number is a function of k and max deficiency, the parameter arising out of and minimised by the recursion evaluated in Stage 1.

Finally, we introduce a particular means of adding new edges, each of which has at least one endpoint in V' , with the objective of satisfying all the target degrees established during the recursion of Stage 1 and k -anonymizing the new vertices added during Stage 2. A critical property of our specific approach is that the edges are added in such a manner as to guarantee tractability of the problem of k -anonymizing the new vertices, a problem which may be hard in the general case.

As we describe the three stages of the algorithm in the following subsections, we illustrate their execution by 3-anonymizing the graphs in Fig. 2.

3.1 Stage 1: determining each vertex's target degree

Our algorithm first proceeds by identifying which vertices should have the same degree and what degree that should be. Similar to Liu and Terzi (2008), we construct a recursion on the degree sequence of \mathcal{G} to compute these groups and degrees. Contrary to their work, however, we minimise max deficiency rather than total deficiency, because, as we prove later in Lemma 2, max deficiency tightly lower bounds the number of new vertices that must be added in order to k -anonymize the degree sequence of \mathcal{G} . Thus, we define a k -partitioning that minimises max deficiency to be optimal.

To produce an optimal k -partitioning of the degree sequence of \mathcal{G} , we offer an incremental algorithm which operates from left (position 1) to right (position n) on the degree sequence, maintaining the optimal (i.e., with minimum max deficiency) k -partitioning of those values in the degree sequence seen so far. The ability to do this with $\mathcal{O}(nk)$ cost is a consequence of the following propositions:

Proposition 1 *The max deficiency of a partition containing a highest degree of d_i and a smallest degree of d_j will be less or equal to the max deficiency of any partition containing d_i and any d_{j+c} or containing d_j and any d_{i-c} , $\forall c \in \mathbb{N}$.*

Proposition 2 *For any partition (d_i, \dots, d_l) , its max deficiency is greater or equal to that of the partitions $(d_i, \dots, d_j), (d_{j+1}, \dots, d_l)$, for $i < j < l$. That is to say, it never produces a higher max deficiency when one splits a partition.*

Both propositions follow from the facts that the degree sequences are sorted and max deficiency (i.e., difference) is transitive. Importantly, they allow us to construct a recursion, and an incremental, dynamic programming algorithm to evaluate that recursion, because they imply that there are only k ways to produce an optimal k -partitioning of the first x elements if the best possible k -partitionings are known for the first i elements, $\forall i < x$. From Proposition 1, it is clear that the x th element should be added to the right of the first $(x - 1)$ elements. From Proposition 2, it is clear that if there is an optimal split point for the rightmost partition that is farther than $2k - 1$ positions left of x , then it can be split into sub-partitions such that the rightmost partition begins at some other position at least as far right as $x - 2k + 1$. The rightmost position must also have at least k elements; therefore, the rightmost partition of the optimal k -partitioning on the first x elements must begin between positions $x - 2k + 1$ and $x - k$, inclusive.

Our algorithm evaluates the recursion “bottom-up”, constructing an optimal k -partitioning of the degree sequence by incrementally adding the next x 'th rightmost degree and determining the best possible k -partitioning. If there are fewer than $2k$ degrees in the sequence, there is not any choice but to group them all together, because at least $2k$ elements are required to make two partitions of size $\geq k$. When, on the other hand, there are at least $2k$ degrees, then we evaluate the cost of splitting off a rightmost partition at any of the k positions between $x - 2k + 1$ and $x - k$, inclusive, and choose the rightmost of the cheapest among them. We invoke the recursion by recognising that the max deficiency incurred by creating a rightmost partition that starts at some position i is exactly the larger of the best possible k -partitioning up to $i - 1$ and the degree differences between the i 'th and x 'th degrees in the sequence.

The following recursion evaluates the cost of splitting and thus constructs an optimal k -partitioning. In the statement of the recursion below, the Δ function computes the max deficiency of a particular partition; the Start function keeps track of where x 's partition starts, should x be the rightmost degree in it; and the Cost function computes the overall cost (max deficiency) of the best possible partitioning up to the x 'th element. The Start function allows us to retrace the best possible k -partitioning up to any x 'th position: the rightmost partition is given by $[\text{Start}(x), x]$; the next rightmost position is given by $[\text{Start}(\text{Start}(x) - 1), \text{Start}(x) - 1]$; &c. Finally, note that

we assume argmin returns the maximal point at which a function is minimised.

Let

$$\text{Cost_Split} = \min_{i \in [\max(k, x-2k+1), x-k]} \times (\max(\text{Cost}(1, i-1), \Delta(i, x))),$$

$$\text{Pos_Split} = \text{argmin}_{i \in [\max(k, x-2k+1), x-k]} \times (\max(\text{Cost}(1, i-1), \Delta(i, x))).$$

Then,

$$\Delta(x, y) = d_x - d_y,$$

$$\text{Cost}(1, x) = \Delta(1, x), \text{ if } x < 2k,$$

$$\text{Cost}(1, x) = \text{Cost_Split}, \text{ if } x \geq 2k,$$

$$\text{Start}(x) = 1, \text{ if } x < 2k,$$

$$\text{Start}(x) = \text{Pos_Split}, \text{ if } x \geq 2k.$$

For the graph of Fig. 2b, the degree sequence is (5, 3, 3, 2, 1, 1, 1). The optimal k -partitioning of this degree sequence is ((5, 3, 3), (2, 1, 1, 1)), which we arrive at by evaluating the recursion, as tabulated in Table 1.

Two parameters important for the next stages of our algorithm arise out of the k -partitioning of the degree sequence, namely max deficiency and total deficiency . As we show later in Lemma 2, the max deficiency is a lower bound on $|V'|$ in an optimal solution. We necessarily must bound the value of total deficiency , because this allows us to upper bound the number of edges added by—and thus execution cost of—our entire algorithm.

Lemma 1 *The total deficiency of an optimal k -partitioning of a sorted degree sequence is upper-bounded by $(n-1)(2k-1)$.*

Proof First, note that from Proposition 2 that any optimal partitioning with a partition sized $2k$ or greater can be split into two partitions such that the max deficiency is not increased. Furthermore, doing so will decrease the total deficiency unless the contribution of the partition is already zero. Therefore, no optimal partition should contain $2k$ or more elements. \square

If $[\text{front}(d_i), \text{end}(d_i)]$ denotes the partition containing d_i , $[\text{front}(p_j), \text{end}(p_j)]$ denotes the partition p_j , and $|\mathcal{P}|$

denotes the number of partitions, then the total deficiency of an optimal k -partitioning is given by:

$$\begin{aligned} & \sum_{d_i} \text{front}(d_i) - d_i \\ & \leq \sum_{d_i} \text{front}(d_i) - \text{end}(d_i) \\ & \leq (2k-1) \sum_{p_j} \text{front}(p_j) - \text{end}(p_j) \\ & \leq (2k-1)(\text{front}(p_1) - \text{end}(p_{|\mathcal{P}|})) \\ & \leq (2k-1)(n-1), \end{aligned}$$

where the second inequality follows because the deficiency contributed by a particular degree is certainly no more than the max deficiency contributed by its entire partition, the third inequality follows because $\text{front}(p_{j+1}) \leq \text{end}(p_j)$, and the fourth inequality follows because the maximum degree in a simple graph is $n-1$.

3.2 Stage 2: determining m , the number of new nodes

The next step in anonymizing a graph is to determine precisely how many vertices should be added. Ideally, we would add exactly max deficiency vertices, because this is a lower bound on how many *must* be added, as shown in Lemma 2:

Lemma 2 *To k -degree-anonymize a set of vertices $X \subseteq V$ by means of vertex addition, one must add at least as many new vertices as the max deficiency of an optimal k -partitioning of the degree sequence of X .*

Proof First note that any graph \mathcal{G}' in which X is k -degree-anonymous corresponds to some k -partitioning of the original degree sequence of X . In order to satisfy every target degree arising from the k -partitioning, some vertex will require as many new edges as the max deficiency of that k -partitioning. In addition, because edges can only be added from each $v \in X$ to new vertices and because we investigate only simple graphs (i.e., those which do not contain multiple edges between the same source and destination nodes), clearly these new edges must connect to max deficiency new vertices in order to satisfy that particular vertex's requirement. \square

From among all k -partitionings of the degree sequence of X , corresponding to all graphs \mathcal{G}' in which X is k -degree-anonymous, the optimal choice minimises max deficiency , so also minimises how many new vertices are required in \mathcal{G}' .

However, for security, we anonymize the new vertices as well (ensure that $V \cup V'$ is k -degree-anonymous in \mathcal{G}'), and this may require adding more than max deficiency nodes. Our approach to adding edges we describe in more

Table 1 The values of the recursion for the 3-degree-anonymization of the example graph from Fig. 2b

Pos	1	2	3	4	5	6	7
Deg Seq	5	3	3	2	1	1	1
Cost (1, x)	0	2	2	3	4	2	2
Start (x)	1	1	1	1	1	4	4

detail in the next subsection, but there are a couple important things to note here because they influence this stage. For readability, let $md = \max \text{ deficiency}$ and let $td = \text{total deficiency}$. Then, once every target degree is achieved through edge additions in Stage 3, $td(\bmod md)$ of the new vertices will have some degree, call it d , and the other $md - td(\bmod md)$ will have a degree of $d - 1$. If these both appear in the target degrees of the k -partitioning, the entire graph \mathcal{G} is k -anonymous. Figure 3 illustrates an example of this, where the graph on the left is anonymized by adding one additional vertex which coincidentally then has a degree $d = 1$ that is already present on three other vertices. On the other hand, if d or $d - 1$ does not appear in the target degrees, the new vertices need to be explicitly anonymized as well. To accomplish this, we create a k -anonymous group of the new vertices by introducing intra-new-vertex edges to establish that they all have the same degree. This requires that: (1) there are k new vertices with which to form a k -anonymous group; and (2) there are an odd number of new vertices (we explain why in the next subsection).

Both these conditions are satisfied when we add exactly $m = (1 + \max(md, k))(\bmod 2) + \max(md, k)$

new vertices to the graph to create a new vertex set V' of size $n + m$.

3.3 Stage 3: the edge insertion

The final stage of our anonymization algorithm is to add new edges to the graph (restricting the addition to those with an endpoint in $V' \setminus V$, as per the problem definition), in order to meet the target degrees established in Stage 1. This must be done carefully because not any arbitrary approach is guaranteed to succeed: the task of anonymizing a sub-graph of m new vertices added in Stage 2 is not generally trivial. Hence the motivation behind our *cycling* approach: it is designed such that it will always produces a scenario in which the anonymization of the new vertices is of linear cost.

First, we order the m additional vertices (arbitrarily). Let $def(v_i)$ be the discrepancy between the i 'th degree in the



Fig. 3 An example for which the additional vertex has the same degree as a 3-degree-anonymous group of original vertices. **a** The example graph from Fig. 2a. **b** A 3-degree-anonymization of **a** in which the new vertex ends up with the same degree as three vertices from the original graph, so does not need to be explicitly anonymized

degree sequence and the largest degree within the same partition (the *deficiency* of the i 'th vertex). We then connect the first $def(v_1)$ of the m additional vertices to the vertex corresponding to v_1 , the next $def(v_2)$ additional vertices to the vertex corresponding to v_2 , and so on until all m additional vertices have an edge. This process ends at some v_j .

We continue with another iteration, this time starting at v_j and with subsequent iterations until we have satisfied the deficiency of every node in the original graph. This is illustrated for the example graph of Fig. 2b in the first five steps of Fig. 4, reading left to right, then top to bottom. Seven edges must be added and they are done so cyclically from left to right.

Because of the nature of this *cycling* procedure, always adding an edge to an additional vertex that has not yet been visited on that particular iteration, we can guarantee that the degrees of the additional vertices are all within one. In fact, if exactly d iterations are required in order to anonymize the original graph, then (as we hinted in Stage 2) $td(\bmod m)$ of the m additional vertices will have degree d and the remaining $m - td(\bmod m)$ will have degree $d - 1$. Because we serviced each original vertex in turn, we can be certain to not accidentally introduce the same edge twice. Therefore, we have successfully k -anonymized every vertex that was in the original graph. Quite importantly, we have done so in a manner that guarantees that the remaining edge-addition-based problem is efficient and successful because there is an odd number of vertices and the degrees are all within one.

The last detail is to ensure that all the new vertices are themselves k -anonymous. As a first recourse, if d and $d - 1$ are both present as target degrees from the recursion in Stage 1, then the additional vertices will already belong to some k -anonymous group. Figure 3 demonstrates that 3-anonymizing the example graph of Fig. 2a is such a case, where the additional vertex fits nicely into the 3-anonymous group of degree 1 vertices.

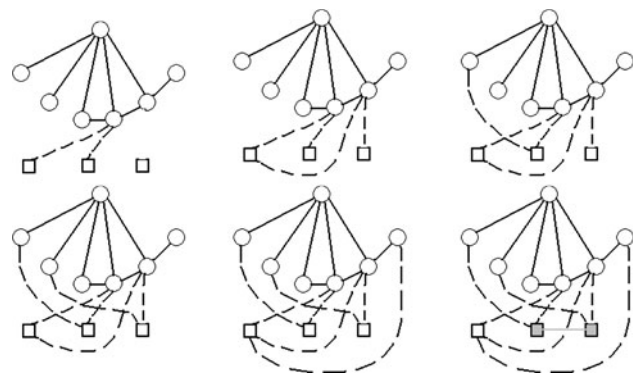


Fig. 4 3-Anonymizing the example graph from Fig. 2b with three additional vertices

In the event that either d or $d - 1$ is not present in the anonymized degree sequence, we explicitly anonymize the new vertices. For the $m - td(\bmod m)$ vertices with degree $d - 1$, we randomly pair them and add an edge between each pair. If $m - td(\bmod m)$ is even and $m \geq k$, we know this is sufficient to guarantee all new vertices have the same degree (namely d) as at least $k - 1$ other vertices.

If, instead, $m - td(\bmod m)$ is odd, then this pairing will leave out one last vertex, call it r . Because of our diligence in selecting m , we know that $m - 1$ is even (and at least k), so we can add an edge from r to each of two other additional vertices so that all three have degree $d + 1$. The remaining $m - 3$ vertices with degree d can then all be paired off again (since $m - 3$ is even) and all additional vertices will be anonymized with degree $d + 1$. Some care must be taken to not accidentally re-add an edge between additional nodes, but this is really quite trivial because of how we proceeded with the preceding edge addition. Figure 5 illustrates this scenario. The degrees (4 and 3) of the two additional nodes are irreconcilable and $md - 1$ is odd, so we instead our algorithm added the extra vertex as in Fig. 4.

3.4 Algorithmic analysis

From the algorithm described in this section, we prove the following theorem:

Theorem 1 *Our algorithm produces a k -degree-anonymous graph \mathcal{G}' containing the input graph \mathcal{G} as an induced subgraph, using $\mathcal{O}(nk)$ time and $\mathcal{O}(n)$ space. The number of new vertices added is optimal up to an additive factor of k .*

Proof First, we note that the degree sequence partitioning requires $\mathcal{O}(n)$ space and $\mathcal{O}(nk)$ time. If this recursion is evaluated bottom-up (i.e., from $x = 1$ to $x = n$) and the results are memoised after each iteration, then the running time is linear in nk because the degree sequence has length n and for each element of the degree sequence there is an iteration that will cost at most a comparison to $2k$, lookups of Cost, and computations of Δ for each $i \in [k, 2k)$. The memory cost of the memoization is $3n$ units of memory: an array of size n each for storing the results of the Cost and

Start calculations and an additional array of length n in which the original degree sequence is kept. \square

Second, the number of edges added is bounded by $td + m$, since td edges are added between original and additional vertices; the subsequent anonymization of additional vertices never changes an additional vertex's degree by more than $(d + 1) - (d - 1) = 2$; and at least one additional vertex already has a degree of d . We introduce at most m additional vertices, at most either $md + 1$ or $k + 1$. k -Degree-anonymizing \mathcal{V} implies, from Lemma 2, a lower bound on the optimal number of additional vertices is $md \geq 1$, noting also that any graph with a md of 0 is already n -anonymous. So, we add at most $md + 1 - md = 1$ more vertex than optimal or we add at most $k + 1 - 1 = k$ more vertices than optimal.

Since m is bounded by n (because neither md nor k can be larger than n) and the td is bounded by $(n - 1)(2k - 1)$, the addition of these edges requires $\mathcal{O}(nk)$ time. Since this requires constant memory, the overall space usage is $\mathcal{O}(n)$.

Note that we assume in this proof that the sorted degree sequence can be provided based on the representation of the graph. If this is not the case, an $\mathcal{O}(n \log n)$ preprocessing step to compute the sorted degree sequence subsumes our algorithm.

Corollary 1 *For input graphs in which $\exists i: (d_{i-1} - d_i) \geq k$ and $(d_i - d_{i+k-1}) \geq k$, our algorithm is optimal up to an additive factor of 1.*

Proof If $\exists i: (d_{i-1} - d_i) \geq k$ and $(d_i - d_{i+k-1}) \geq k$, then the max deficiency of the degree sequence is necessarily at least k because d_i must be grouped either with d_{i-1} or d_{i+k-1} by Proposition 1. Consequently, the scenario in which the algorithm adds some number of vertices other than max deficiency or max deficiency + 1 does not arise. \square

This corollary is particularly interesting with respect to social network graphs, the degrees of which tend to follow a power law distribution 7.

3.5 On the security of \mathcal{G}'

Perhaps an attacker's best chance at deconstructing our anonymization comes from cycling through all possible choices of the new equivalence classes (k -anonymous groups), removes them from the graph and then runs our algorithm to check if the output matches the original anonymized graph. We note here, that the last step where the adversary has to check if the output graph is the "same" as the initial anonymized graph requires an algorithm for checking graph isomorphism. No efficient (polynomial time) algorithm is known for this problem. In addition, there are examples in which an equivalence class

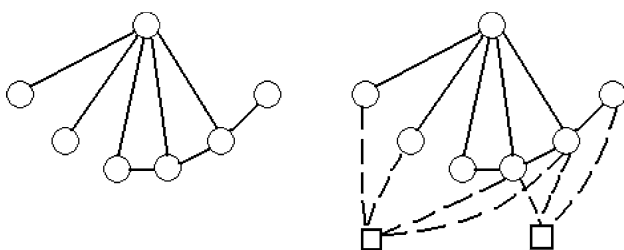


Fig. 5 An example in which our algorithm will not optimally 3-degree-anonymise Fig. 2b with only two additional vertices because of the difficulty resolving the anonymity of the new vertices

of additional vertices is merged with an equivalence class of original vertices because they have the same degree at the end of the anonymization procedure. For such graphs, this attack will not be able to extract the additional vertices only.

4 Experimental evaluation of the utility of \mathcal{G}'

Our algorithm, and indeed a vertex-addition approach in general, offers many advantages in terms of theoretical guarantees. In Sect. 3.4, for example, we proved asymptotic performance and near-optimality.

If we consider again the motivation for anonymizing the network in the first place, it is because ultimately we want to release it for analysis. In the case of an unlabelled graph, this is necessarily structural analysis. So, preserving the utility of the released data by not especially disrupting the structural characteristics of the graph is important. Thus more theoretical guarantees that our vertex-addition model offers: the number of cliques of size >3 is exactly maintained between \mathcal{G} and \mathcal{G}' . More broadly, for any monotone property of a graph, our approach can only possibly add false positives, never generating false negatives.

However, there are interesting non-monotone properties, too. In particular, for *clustering coefficient*, *average path length*, and the *hop plot* of a graph, the effect of any anonymization procedure is more difficult to predict and input-dependent, but the properties are focal points of other anonymization and social network analysis papers. In this section we present an experimental study in two parts of to what extent the application of our algorithm in the previous section disrupts the utility of the network with respect to these less predictable properties.

First, we investigate the scalability of the algorithm with respect to how well graph characteristics are preserved with increasingly demanding test scenarios (Sect. 4.1). Then, in Sect. 4.2 we investigate how well graph characteristics are preserved in the context of alternative approaches.

4.1 Scalability tests

4.1.1 Metrics and setup

Before describing the results of our experimental evaluation, we first detail our experimental setup. In particular, we describe the choice of datasets and metrics, and the implementation and machine details.

Datasets Four datasets from diverse domains form the subjects of our empirical study. We use one large dataset that represents an email communication network in a company (namely, Enron) (Leskovec et al. 2005). In addition, we select three other datasets that we expect to exhibit near

Table 2 Structural properties of datasets for our empirical study

Graph	Nodes	Edges	APL	CC
Enron	36,692	183,831	3.39	0.09
Net science	1,589	2,742	5.76	0.69
Prefuse	129	159	3.16	0.07
Football	115	613	2.51	0.61

worst-case behaviour. The *Net Science* graph (Newman 2006) has a large discrepancy in degrees and thus incurs substantial max deficiency and total deficiency during the anonymization process. The *Prefuse* (Heer 2005) and *Football* (Girvan and Newman 2002) graphs are small, so the effect of adding nodes and/or edges has a larger percentage effect on the properties of the graphs. The properties of these datasets are shown in Table 2. Note that Net Science, in particular, has been remarked to be a strong test set for network analysis by Leskovec et al. (2008). The Enron dataset was obtained from the Stanford SNAP repository,⁶ and the Net Science and Football datasets were obtained from Mark Newman's repository.⁷

Metrics We measure the distortion introduced by the algorithm via some metrics which are commonly studied properties in the social network literature according to the survey by Chakrabarti and Faloutsos (2006). The three metrics we study are defined below.

1. **Clustering coefficient (CC)** (Barrat and Weigt 2000): Informally, clustering coefficient measures the percentage of paths of length 2 which are also triangles. This metric in some sense measures *triadic closure* of graphs—social networks are known to have significant triadic closure (friends of a person are also likely to know each other). More formally, for all ordered triples $u, v, w \in V$,

$$CC = \frac{||\{u, v, w \in V : (u, v) \in E \wedge (u, w) \in E \wedge (v, w) \in E\}||}{||\{u, v, w \in V : (u, v) \in E \wedge (u, w) \in E\}||}$$

2. **Average path length (APL)**: This metric is a measure of the expected path length in the graph between any two randomly chosen connected vertices. This metric is highly relevant as it is directly related to the *six degrees of separation* that is known to exist between randomly chosen people in social networks, since the study of Milgram (1967). Define a predicate $C(u, v)$ to be *true* if u and v are connected in the graph and *false* if they are not connected. Define $CP = \{(u, v) : C(u, v) = \text{true}\}$ to be the set of all the pairs of vertices that are connected. We define the average path length to be:

⁶ <http://snap.stanford.edu/data/>.

⁷ <http://www-personal.umich.edu/mejn/netdata/>.

$$\text{APL} = \frac{\sum_{(u,v) \in \text{CP}} \text{PathLength}(u,v)}{|\text{CP}|}.$$

We assume that $\text{PathLength}(u,u) = 0$ for all $u \in V$.

3. Hop Plot (Faloutsos et al. 1999): The connectivity of a graph can be graphically modeled using the *hop plot*. The hop plot studies reachability for each path length k . For a given value of k , the hop plot displays, summed over all the vertices, the number of nodes reachable from that vertex using paths of length at most k . The maximum value for any value of k is n^2 where n is the number of vertices in the graph. The smallest value of k for which the maximum value of n^2 is reached is the *diameter* of the social network, the path length using which any two nodes in the graph can reach each other. Changing or distorting the connectivity of a graph drastically would change the shape of their hop plots. This is the main motivation behind studying these plots.

Experimental setup A java implementation of the algorithm was used to measure the distortion based on the metrics for the five chosen datasets defined earlier in this section. The resulting graphs were manually verified to be k -anonymous. All experiments were performed on a quad-core Intel Xeon 5140 2.33 GHz processor with 4 MB of L2 cache and 6 GB of RAM.

Experimental studies in literature typically study small values of k (close to 3 and rarely even 100). We vary k as a fraction of n for our experiments, while still maintaining that $k \ll d$. For large and midsize datasets, we vary k from $k = 0.25$ up to 2 % of the number of nodes in the graph. For the two small datasets (with over 100 nodes), this is too refined, so we vary k from $k = 1$ up to 5 % of the number of nodes in the original graphs. For the large and midsize datasets, 2 % of the number of nodes would still translate to k values of 80 to 720 which is a substantial number for the context.

4.1.2 Results and discussion

We first discuss the results of the experimentation on clustering coefficient. Notice first the plots for the Enron (Fig. 6) and Football (Fig. 7) datasets, which are examples of very good performance. The horizontal axis represents the value of k as a function of the number of nodes in the dataset. The vertical axis represents the clustering coefficient. Zero distortion would correspond to the solid line, which represents the characteristics of the anonymized data, exactly overlaying the dotted line, which represents the value in the original dataset. The performance of our algorithm on the Football dataset and for larger values of k on the Enron dataset is excellent because it very nearly achieves this perfect overlay.

Our performance with respect to clustering coefficient on the other two datasets, Net Science and Prefuse, is depicted in Figs. 8 and 9, respectively. That the opposite effect was observed on each dataset as k increased demonstrates the unpredictability of the effect of anonymization on clustering coefficient. In this case, the results match intuition, because of the properties of the dataset. In particular, Net Science has a high clustering coefficient originally; on the other hand, Prefuse features a couple nodes with especially high degree relative to the other nodes, so each newly added vertex is connected to a large percentage of original nodes. This is necessary in order to balance out the very high *td:md* ratio at higher values of k .

Next, we discuss the results of the measurements on average path length. The plots that we have created are to be interpreted in the same manner as those just seen for

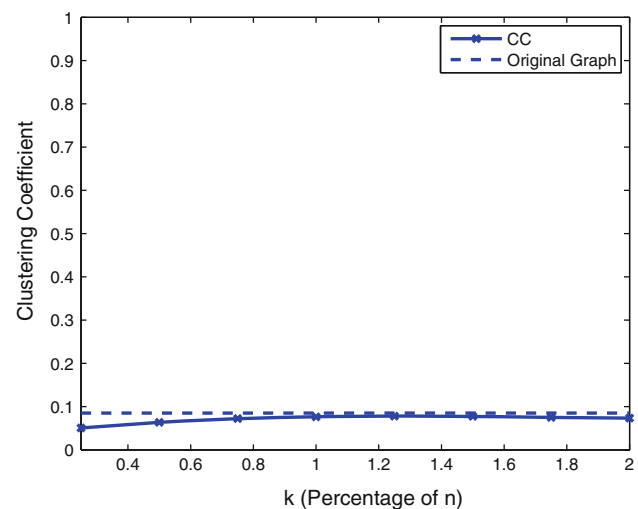


Fig. 6 Distortion of CC of Enron dataset as a function of k

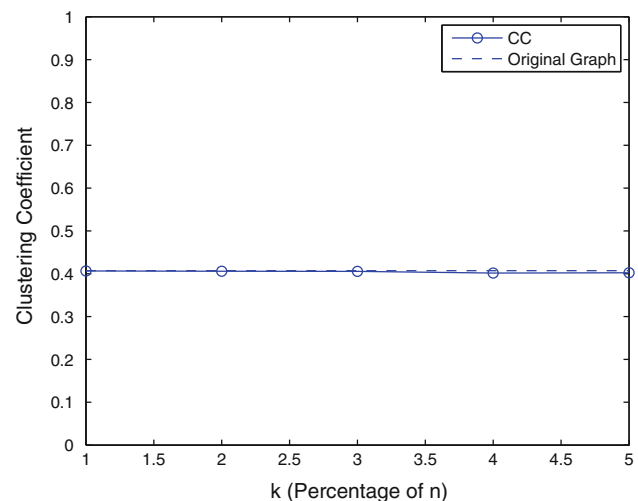


Fig. 7 Distortion of CC of Football dataset as a function of k

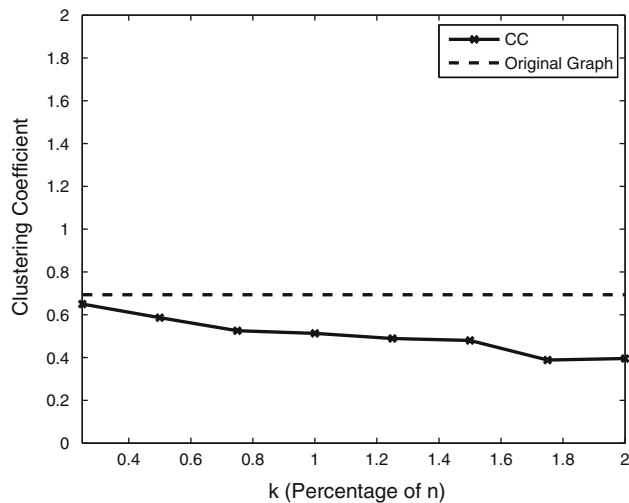


Fig. 8 Distortion of CC of Net Science dataset as a function of k

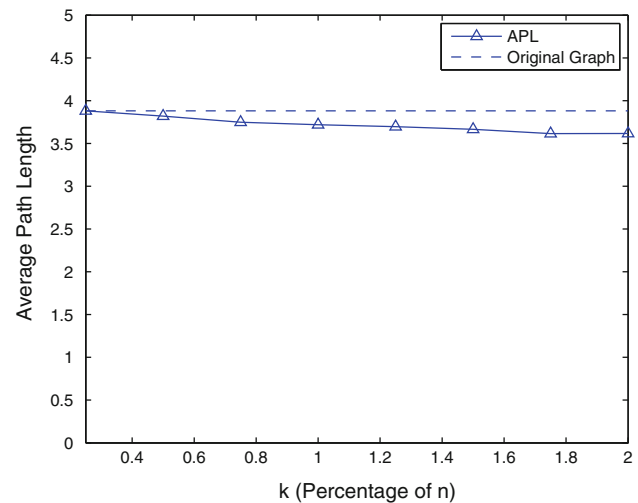


Fig. 10 Distortion of APL of Enron dataset as a function of k

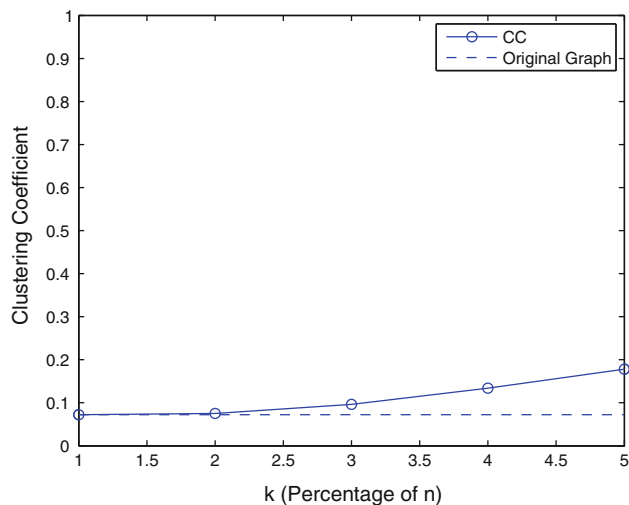


Fig. 9 Distortion of CC of Prefuse dataset as a function of k

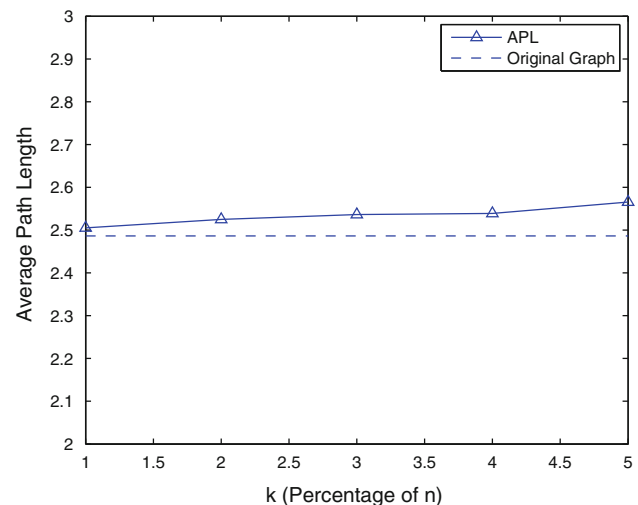


Fig. 11 Distortion of APL of Football dataset as a function of k

clustering coefficient. That is to say, perfect performance occurs when the solid line, representing the anonymized graphs, overlays the dotted line, representing the value in the original graph. We again start by showing our performance on the Enron (Fig. 10) and the Football (Fig. 11) datasets, to demonstrate especially good results.

It is worth noting again the unpredictability of these metrics. The average path length on the Football dataset, unlike any of the others, rises in the anonymized graph. We note that the only especially large difference between the original graph and anonymized graph occurs on the Net Science dataset (Fig. 12), which has a relatively high original average path length compared to the other datasets. Although there is a notable change on the Prefuse dataset (Fig. 13) for $k = 2\%$, it decreases at a very slow rate for all subsequent values $k > 2\%$.

The final set of experiments that we ran on the structural distortion of our algorithm for non-monotone properties was in measuring the hop plot of the graphs. The structure of these plots is different, but the interpretation is similar. For each value of k , a separate line is depicted. In addition, one line depicts the hop plot for the original graph. Because the number of vertices changes, the lines cannot possibly overlay each other. Instead, good results are indicated by the similarity of the *shape* of the curves. Of the results on the Enron (Fig. 14), Football (Fig. 15), Net Science (Fig. 16), and Prefuse (Fig. 17) datasets, all but Fig. 16 illustrate especially good performance, where the shape of the curve is consistent despite its increasing max value.

Notice that the hop plot lines for the Football dataset nearly *do* overlay each other: this is because the number of additional nodes m required to anonymize the dataset is quite small. The

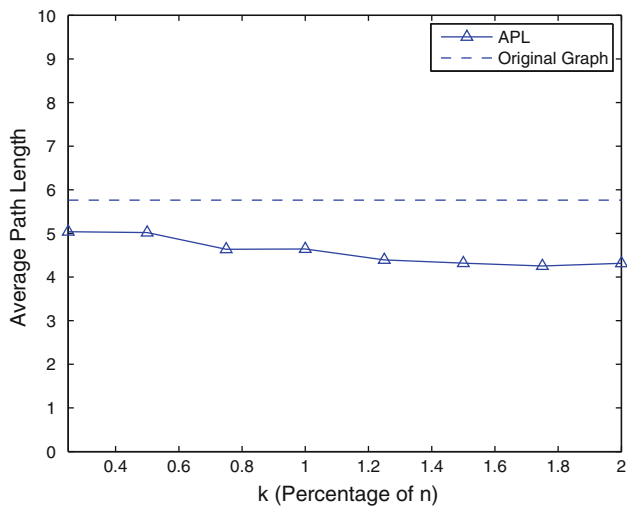


Fig. 12 Distortion of APL of Net Science dataset as a function of k

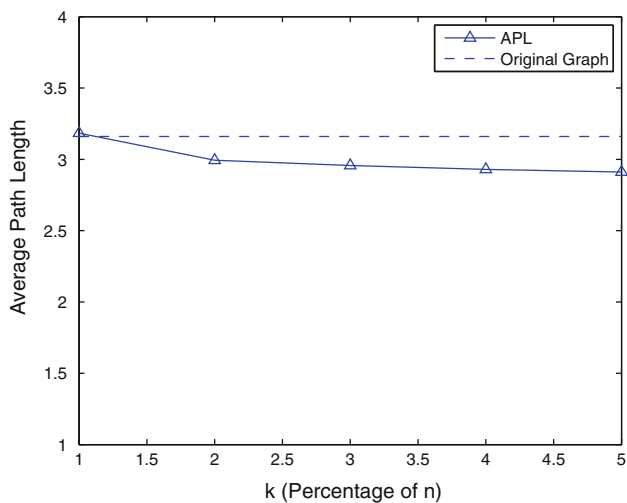


Fig. 13 Distortion of APL of Prefuse dataset as a function of k

Net Science dataset (Fig. 16) has a hop plot that varies more substantially after anonymization. Again, however, most of this change occurs at a low value of k ; for subsequent, higher values, not much more change is witnessed.

A last comment from our scalability evaluation is with respect to execution time. Indeed, we proved good asymptotic performance, but, nonetheless, absolute running time is always a curious measure. For the largest dataset (Enron), the running time for the actual anonymization took a little over 1 min (70 s) for all the values of k in the plots, reflecting the proven efficiency of the algorithm. The running time reported is the average of five independent runs of the algorithm for each value of k . The evaluation time was dominated by the naive computation of the metrics on the original and distorted graphs. These took about 20–30 min each.

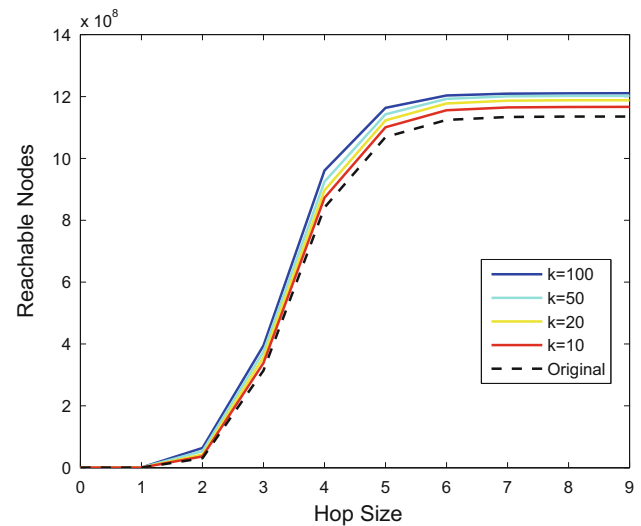


Fig. 14 Distortion to the hop plot of Enron dataset for several values of k

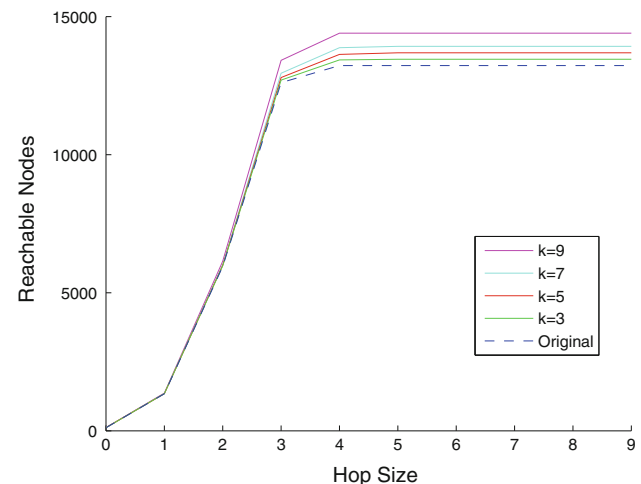


Fig. 15 Distortion to the hop plot of Football dataset for several values of k

The times on smaller graphs were much lower and had the same trend where the computation of the metrics dominated the running time.

4.2 Comparability tests

Although our work here is quite different in nature, we still believe it is important to present a comparison to the nearest research, namely the state-of-the-art edge-based k -degree-anonymization algorithm (Liu and Terzi 2008). The purpose of these experiments is to evaluate whether the vertex-based approach that we introduce is empirically competitive in terms of what is known about the edge-based approach.

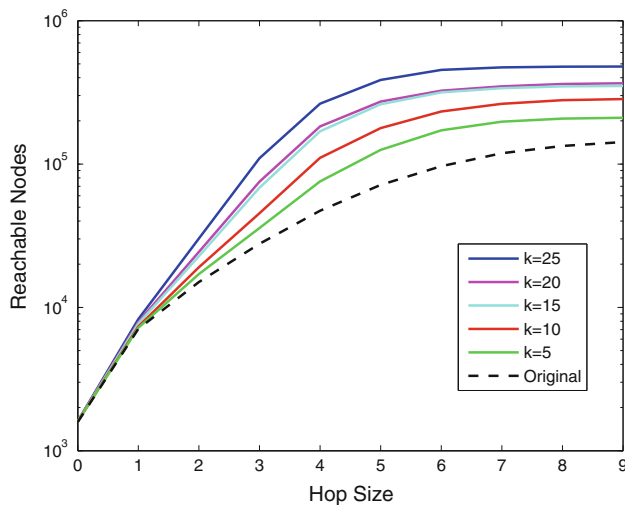


Fig. 16 Distortion to the hop plot of Net Science dataset for several values of k

4.2.1 Methodology

An experiment was run earlier by Ying et al. (2009) to compare the k -degree-anonymization of Lui and Terzi (2008) to an approach that randomly adds and removes edges to provide anonymity through obscurity. Our methodology for comparing to the algorithm of Lui and Terzi is to run a series of experiments identical those of Ying et al. and compare the performance of our algorithm to these published values. Thus, as before, we anonymize the entire vertex set, V .

Datasets We conduct the experiments on the same *polblogs* dataset (Adamic and Glance 2005) used by Ying et al., which we obtained from the Carnegie Mellon CASOS repository.⁸ The dataset consists of 1,222 vertices, each corresponding to a political blog b_i , and 16,714 edges (b_i, b_j) indicating that a hyperlink existed from b_i to b_j and/or from b_j to b_i . In other words, the dataset is an undirected citation network.

Metrics Ying et al. evaluate four metrics that were presented in the comprehensive review of network characteristics by Costa et al. (2007). We evaluate three of these; the modularity measure, Q , on the other hand, does not apply in our case because we assume vertices are all of the same type (i.e., unlabeled), and it evaluates the extent to which the network is modularized by vertex type.

1. Transitivity measure (C) (Barrat and Weigt 2000): The transitivity measure used by Ying et al. is exactly the *clustering coefficient* that we introduced in Sect. 4.1. We will continue to refer to this as the clustering coefficient of the graph.

⁸ http://www.casos.cs.cmu.edu/computational_tools/datasets/external/polblogs/index11.php.

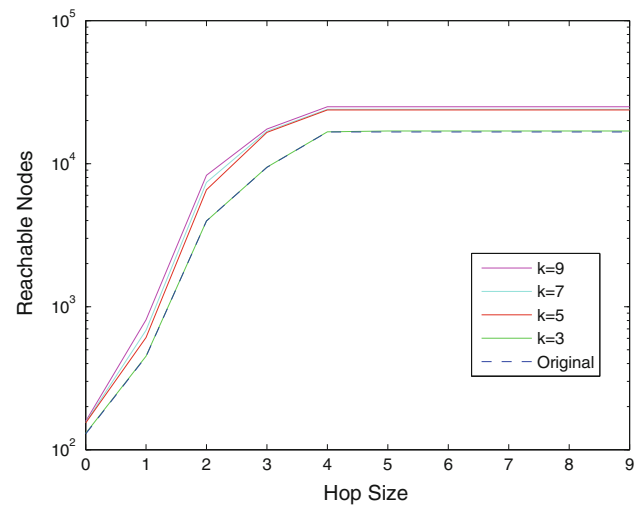


Fig. 17 Distortion to the hop plot of Prefuse dataset for several values of k

2. Harmonic Mean of the shortest path (h) (Latora and Marchiori 1987): The *harmonic mean* is an evaluation of connectivity, similar to the *average path length* that we used in Sect. 4.1. Let d_{ij} be the distance of the shortest path from vertex i to vertex j , or ∞ if they not connected. Then,

$$\frac{1}{h} = \frac{1}{|V|(|V| - 1)} \sum_{i \neq j} \frac{1}{d_{ij}}.$$

3. Subgraph centrality (SC) (Estrada and Rodriguez-Velazquez 2005): The *subgraph centrality* is an evaluation of how many (not necessarily simple) cycles emanate, on average, from each vertex, where each cycle is weighted by the reciprocal of the factorial of its length. Let W_i^l be the number of walks⁹ of length l that start and end at vertex i . The subgraph centrality is the following convergent infinite summation, which converged to the precision we report in these experiments at a length of $l \approx 120$:

$$SC = \frac{1}{|V|} \sum_{i=1}^{|V|} \sum_{l=2}^{\infty} \frac{W_i^l}{l!}.$$

Experimental setup The experiments are conducted identically to those in our scalability tests (Sect. 4.1), except that we instead vary k through integers $k \in [2, 10]$, for precise compatibility with the study by Ying et al. We report all values to the same precision as offered by Ying et al.

4.2.2 Results

The results for the comparability study are disclosed in Table 3 below. Each row of the table corresponds to an

⁹ Recall that a walk is any sequence of adjacent edges, including those which revisit edges and/or vertices.

Table 3 Values of *clustering coefficient* (CC), *harmonic mean* (h), and *subgraph centrality* (SC), respectively, for the *polblogs* dataset after applying *k*-degree-anonymization by means of vertex addition (V), as in this paper or strictly edge addition (E), as in Liu and Terzi (2008)

	CC-V	CC-E ^a	h-V	h-E ^a	SC-V ($\times 10^{29}$)	SC-E ^a ($\times 10^{29}$)
$k = 1$	0.226	0.226	2.51	2.51	1.21	1.21
$k = 2$	0.219	0.225	2.51	2.50	1.30	2.73
$k = 3$	0.215	0.223	2.49	2.48	1.41	1.87
$k = 4$	0.207	0.224	2.48	2.49	2.16	3.61
$k = 5$	0.205	0.221	2.48	2.48	2.88	3.40
$k = 6$	0.200	0.222	2.46	2.47	2.66	1.45
$k = 7$	0.226	0.220	2.46	2.46	5.55	6.94
$k = 8$	0.190	0.219	2.45	2.46	5.37	6.25
$k = 9$	0.185	0.221	2.44	2.49	11.0	4.46
$k = 10$	0.183	0.221	2.43	2.46	8.25	4.04

The row $k = 1$ represents the original characteristics of the graph (i.e., only a guarantee of 1-anonymity). Note that the scale is changed for the two columns reporting *subgraph centrality*

^a Values for edge-addition (E) taken from Ying et al. (2009)

experiment with a different anonymity threshold. Each column indicates the measured value for the corresponding graph characteristic. Each characteristic is reported twice, once corresponding to an edge-addition algorithm and indicated by the suffix *-E* and once corresponding to our node-addition algorithm and indicated by the suffix *-V*. Perfect performance in a row would be indicated by achieving exactly the same characteristics as in the original graph (the $k = 1$ row); deviation in either direction is undesirable.

4.2.3 Discussion

The results of the comparability experiments are very encouraging. Our algorithm does not rely on any randomization, yet performs as well as the edge-addition algorithm on the first two attributes. For *clustering coefficient* our algorithm is within 0.018 of the edge-additions values, on average, with a standard deviation across points of 0.015. This is very close, considering the precision of the experiments is reported only to within 0.001. For *harmonic mean*, these numbers are even closer: $\mu = 0.01$ and $\sigma = 0.02$ with an experimental precision of 0.01.

These values are closer even than they seem because of the randomization that is involved in the edge-addition algorithm. Whereas our algorithm will always produce the same values for these characteristics, any values reported for the algorithm of Lui and Terzi are subject to randomness and could be either higher or lower on any subsequent run of the experiments. This suggests that either approach to *k*-anonymization can produce anonymous graphs that have similar structural characteristics.

Results for the final metric under consideration, *subgraph centrality*, are somewhat inconclusive. Our algorithm outperforms the edge-addition on all but one case prior to

$k = 9$. This is reasonable to expect, because adding edges to existent nodes is more likely to create new, short loops than is adding edges to “fake” nodes and the addition of new nodes by our algorithm increases the denominator to help stunt the growth of the numerator. For the highest values of k , 9 and 10, the edge-addition has an unintuitive 40 % leap in performance in contrast to our predictable degradation. Consequently, it is difficult to extrapolate what might occur with this metric in new experiments.

5 A hardness result for subset anonymization on labelled graphs

Until this point, we have been considering the special case of unlabelled graphs. For example, we stripped the political lean labels from the *polblogs* dataset during our experiments. In this section, we increase $|\Sigma|$ and study the tractability of *k*-anonymization against degree-based attacks. With $|\Sigma| > 1$, the nature of a degree-based attack changes, because the adversary may know the number of outbound edges from a target v to vertices of each label. So, a stronger anonymity notion is required. (See Fig. 18.) Also, we recognise that often it is not the case that *every* vertex need be anonymous (for example, see the work of Yuan et al. (2010)), so consider *subset anonymization*. We begin by formally defining the problem for the labelled setting. Note that for $|\Sigma| = 1$, this reduces to the problem introduced in Sect. 3.

The analogue in the labelled setting of a degree sequence is a label sequence:

Definition 8 (*Label sequence*) For $v \in V$, we say that $S_v = (l_1, l_2, \dots, l_m)$ is a label sequence for v if it corresponds to some ordering of the labels of v and the vertices

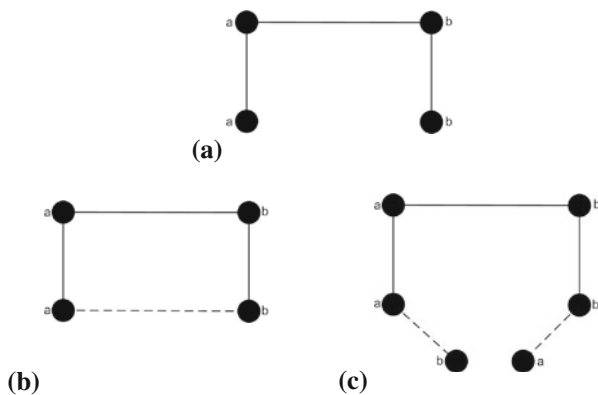


Fig. 18 An example of 2-sequence-anonymity. Notice that 2-degree-anonymity **(a)** is insufficient to guarantee anonymity in a vertex-labelled setting. **b** The example vertex-labelled graph made to be 2-sequence-anonymous by means of edge addition. **c** The example vertex-labelled graph made to be 2-sequence-anonymous by means of vertex addition

that are adjacent to v . We will consider label sequences of vertices to be equivalent up to reordering.

The four label sequences in the example of Fig. 18a, anticlockwise from the bottom left, are $((v, v), (v, v, u), (v, u, u), (u, u))$. Then, k -anonymity for labelled graphs relates to the uniqueness of label sequences:

Definition 9 Given a vertex-labelled graph $\mathcal{G} = (V, E, \Sigma, \ell)$, a subset $X \subseteq V$ of vertices is k -sequence-anonymous in \mathcal{G} if for every vertex $v \in X$, there are at least $k - 1$ other vertices in X whose label sequence is the same as the label sequence of v .

We denote by $u \equiv v$ that vertices u and v have the same label sequence. Clearly, \equiv is an equivalence relation and hence induces a partition X/\equiv of X into equivalence classes. If X is k -anonymous in \mathcal{G} , every equivalence class is of size at least k . Then, the problem of k -sequence anonymization is defined as:

k -Labelled subgraph anonymization

Input: A vertex-labelled graph $\mathcal{G} = (V, E, \Sigma, \ell)$, a set $X \subseteq V$ of vertices, and an integer t .

Question: Is there a vertex-labelled graph $\mathcal{G}' = (V \cup V', E \cup E', \Sigma \cup \Sigma', \ell \cup \ell')$ such that $|V'| \leq t$, $E' \subseteq (V \times V') \cup (V' \times V)$, $\Sigma'_{|V} = \Sigma$, $\ell'_{|V} = \ell$ and X is k -anonymous in \mathcal{G}' ?

That is, can we k -sequence-anonymize X by adding at most t new labelled vertices? New edges must have at least one endpoint in V' .

Theorem 2 k -Labelled Subgraph Anonymization is NP-Complete for $k \geq 3$.

Proof To show that this problem is NP-Hard, we build a reduction from the following anonymization problem on tables that was shown to be NP-Hard by Meyerson and Williams (2004):

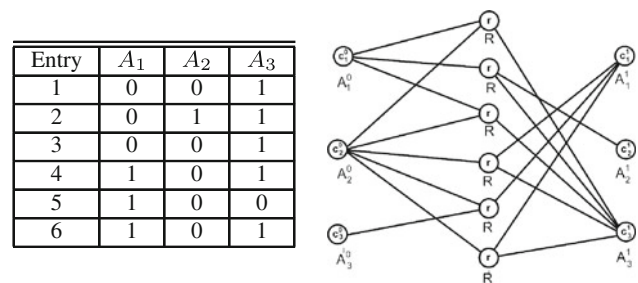


Fig. 19 An example table T and its transformation into a graph \mathcal{G}_T

k -Attribute-anonymity

Input: A table T with n rows and l columns (also called attributes) with entries over $\{0, 1\}$ and integers t and $k \geq 3$.

Question: Can the rows of T be k -anonymized by suppressing at most t attributes of T ? An attribute is said to be suppressed if all its entries (0 and 1) are replaced by $*$.

Reduction: Our reduction is described as follows:

Given a Table T , let $T_{(m,j)} \in \{0, 1\}$ denote the value of attribute j in row m . Then, the vertex-labelled graph \mathcal{G}_T corresponding to T is constructed as follows:

- $V_T = \{r_1, r_2, \dots, r_n\} \cup \{c_j^i | 1 \leq j \leq l, i \in \{0, 1\}\}$.
- Let $E_T = \{(r_m, c_j^i) | T_{(m,j)} = i\}$ where $1 \leq m \leq n$, $1 \leq j \leq l$ and $i \in \{0, 1\}$.
- $\Sigma = \{R\} \cup \{A_j^i | 1 \leq j \leq l, i \in \{0, 1\}\}$.
- $\ell(r_i) = R$ and $\ell(c_j^i) = A_j^i$.
- Finally, remove all isolated vertices from \mathcal{G}_T .

In other words, we encode a binary table as a graph in which a row vertex r_m with label R is connected to a column vertex c_j^0 (c_j^1) with label A_j^0 (A_j^1) if the (m, j) th entry is 0 (1). Figure 19 illustrates this reduction.

Let $X = \{r_1, \dots, r_n\}$ denote the set of row vertices of \mathcal{G}_T . We will show that T can be k -anonymized by suppressing at most t attributes if and only if we can k -anonymize X by adding at most $2t$ new labelled vertices.

Let \mathcal{G}'_T be any graph obtained from \mathcal{G}_T such that X is k -anonymous in \mathcal{G}'_T . We will now show that we may assume without loss of generality, \mathcal{G}'_T satisfies a set of properties. If it does not, we convert it into one that satisfies the properties without increasing the number of new labelled vertices added in \mathcal{G}'_T .

Our first lemma shows that the anonymization procedure does not introduce vertices with new labels not in Σ_T . Furthermore, none of the new vertices added to \mathcal{G}_T will be labelled by R . \square

Lemma 3 $\Sigma'_T = \Sigma_T$ and $\ell'_T(v) \neq R$ for any $v \in V'_T$.

Proof Suppose there is at least one vertex with a new label l , $l \notin \Sigma_T$ in \mathcal{G}'_T . Let Y be any equivalence class of X/\equiv . Every vertex $y \in Y$ has the same number of labels l in its label sequence in \mathcal{G}'_T . Therefore, the label sequences of

all vertices in Y will remain the same after all the new vertices with label l are removed from \mathcal{G}'_T . Since this is true of any equivalence class Y , we can assume that $\Sigma'_T = \Sigma_T$. To show that $\ell'_T(v) \neq R$ for any $v \in V'_T$, note that the label R does not appear in the label sequence of any vertex of Y in \mathcal{G}_T . In addition, R must appear the same number of times in the label sequence of any vertex of Y in \mathcal{G}'_T . Therefore, by the same reasoning as above, all new nodes with the label R in \mathcal{G}'_T can be removed and Y will remain k -anonymous. \square

Lemma 4 *For every $i, i \in \{0, 1\}$, and every $j, j \in \{1, 2, \dots, l\}$, the label A_j^i appears at most once in the label sequence of a vertex from X in the graph \mathcal{G}'_T .*

Proof Let Y be any equivalence class of X . All vertices in Y have the same label sequence in \mathcal{G}'_T . If a label A_j^i appears in the label sequence of a vertex y of Y in \mathcal{G}_T , we can assume that no edge between y and a new vertex labelled A_j^i is added in \mathcal{G}'_T . If not, this label will appear more than once in the label sequence of y and hence in the label sequence of every other vertex of Y . Also, the number of occurrences of this label in every label sequence will be the same. Since at most one of these will be due to an edge of \mathcal{G}_T , for each vertex v in Y , we can remove all but one edge in $E_T \cup E'_T$ between v and a node labelled A_j^i in \mathcal{G}'_T and still preserve anonymity of vertices in Y . \square

Lemma 5 *For every $i, i \in \{0, 1\}$, and every $j, j \in \{1, 2, \dots, l\}$, at most one new vertex with the label A_j^i is added in the graph \mathcal{G}'_T .*

Proof Suppose not. We merge all the new nodes with the label A_j^i into a single node. This only reduces the number of new vertices added with label A_j^i . In Lemma 4, we have shown that every label A_j^i appears at most once in the label sequence of every vertex in X . Therefore, the merge operation will not create multiple edges and all the label sequences of the vertices of X will remain the same as before after this modification. \square

Lemma 6 *All new vertices appear in pairs. That is, if $A_j^0 \in \ell(V'_T)$, then $A_j^1 \in \ell(V'_T)$ and vice versa.*

Proof Suppose a new vertex with the label A_j^0 is added. Let us consider all the vertices in X to which this vertex is adjacent in \mathcal{G}'_T . Then one of those vertices v must be in an equivalence class Y that had a vertex v' with an edge to vertex with label A_j^0 in \mathcal{G}_T . If there is no such Y , it implies that all the vertices in X adjacent to this new vertex had the label A_j^1 in their label sequence before anonymization. Therefore, the new vertex with the label A_j^0 can be removed and all equivalence classes of X will remain k -anonymous. Now note that v has both A_j^0 and A_j^1 in its label sequence in

\mathcal{G}'_T . Hence, the anonymization procedure must create a new vertex labelled A_j^1 and add an edge between v' and this new vertex to preserve anonymity. \square

We now give a proof of correctness of our reduction.

Lemma 7 *T can be k -anonymized by suppressing t attributes if and only if the subset X of vertices in \mathcal{G}_T can be k -anonymized by adding $2t$ new vertices.*

Proof (only if:) Suppose T can be k -anonymized by suppressing t attributes. For every column A_j that is suppressed, we add two new vertices A_j^0 and A_j^1 and add an edge from every row vertex labelled R to one of these new vertices that was not in its label sequence originally. At the end of this procedure, if a set of rows of T formed a k -anonymous group, the corresponding set of row vertices in X form a k -anonymous group. To do this, we added $2t$ new vertices. \square

(if:) As shown in Lemma 6, we can assume that the anonymization procedure adds new vertices in pairs. If the procedure adds a new pair (A_j^0, A_j^1) , then suppress the column A_j in T . Note that if $2t$ new vertices are added in \mathcal{G}'_T , t attributes are suppressed in T . At the end of this procedure, it is easy to see that if a set of row vertices in X form a k -anonymous group in \mathcal{G}'_T , then the corresponding set of rows of T are k -anonymous.

We now remark that this problem is also in NP . At the first glance, this is not clear as there is no natural bound on the integer t in terms of the size of the graph \mathcal{G} . However, based on Lemmas 3 and 4, it is clear that we can assume without loss of generality $t \leq |\Sigma|$. Therefore, a membership in this language can be shown by a list of at most $|\Sigma|$ new vertices and a list of new edges. Therefore, this problem is NP -Complete. This completes the proof of the theorem.

As a last note on the subject of tractability, we remark that our algorithm given earlier implies tractability (in fact, describes an $\mathcal{O}(nk)$ algorithm) for the analogous subset anonymization problem on unlabelled graphs. In the case of subset anonymization, the new vertices need not be also anonymized. The problem is interesting from a practical perspective, because it is non-obvious for an adversary to distinguish between vertices of V' and of $V \setminus X$.

Corollary 2 *k -Labelled Subgraph Anonymization is in P if $|\Sigma| = 1$, for all k .*

Proof This follows from the algorithm in Sect. 3, because anonymization of a labelled subgraph with an alphabet size of 1 is equivalent to anonymization of an unlabelled graph simply by removing all the labels, performing the unlabelled anonymization, and then relabelling every vertex. Because $V' \cap X = \emptyset$, the third step of the algorithm presented in Sect. 3 is not required. Consequently, exactly

max deficiency new vertices are sufficient to anonymize X , because the additional vertices were required for the anonymization of V' . \square

6 Related work

The notion of k -anonymity was introduced by Sweeney (2002) within the context of relational databases, with the insight that tabular microdata could be published without compromising privacy if every tuple of a relation was made to look identical to at least $k - 1$ other tuples with respect to identifying and quasi-identifying attributes. Subsequent to this, work by Meyerson and Williams (2004) and by Agarwal et al. (2005) demonstrated hardness for k -anonymity of tables, even if the quasi-identifiers come from a small-sized alphabet. Outside the relational model, the privacy of statistical tables was studied towards inference control (Domingo-Ferrer 2002; González 2002; Robertson and Ethier 2002).

Later, the desire to publish microdata for—and, consequently the concern of privacy with respect to—other types of data emerged, particularly for the growing body of social network graphs. The knowledge transfer from the relational database community to the study of social networks was pioneered by Backstrom et al. (2007), Zheleva and Getoor (2007), and Hay et al. (2008), who formalised the first notions of attacks against social network data.

This sparked a series of works on developing algorithms for anonymity against progressively stronger adversaries, originating with the study of Liu and Terzi (2008) on degree-based attacks in unlabelled graphs and of Zhou and Pei (2011) on neighbourhood attacks in vertex-labelled graphs. Both of these papers offered an experimentally effective algorithm, while the latter also offered a proof of NP -Hardness. The empirical merits of this approach was verified by Ying et al. (2009), who demonstrated that deliberate k -anonymization can preserve structural characteristics of graphs much better than adding random noise does. Subsequently, Thompson and Yao (2009) studied i -hop degree-based attacks. That is, an adversary's prior knowledge includes the degree of the target and the degree of its neighbours within i hops. They develop an inter-cluster matching method for anonymizing graphs against 1-hop attacks through edge addition and deletion. Cheng et al. (2010), in their work on k -isomorphism, form k pairwise isomorphic subgraphs to achieve protection against two specific classes of attacks. Wu et al. (2010) proposed the k -symmetry model, wherein for any vertex v , there exist at least $k - 1$ other vertices to which v can be mapped using an automorphism of the underlying graph.

Each of the above adversarial models assumed that an adversary's objective was so-called *identity disclosure*: the

identification of a vertex. Another attack is *attribute disclosure*, which seeks not necessarily to identify a vertex, but to reveal sensitive labels of the vertex. In the context of relational databases, this sort of attack was formalised by Machanavajjhala et al. (2007) with the introduction of l -diversity as an anonymity measure, which requires that at least l attribute values appear in each equivalence class. Li et al. (2007) then introduced t -closeness, which required instead that the *distribution* of attribute values in each equivalence class was within t of the entire dataset. More recently, Chester and Srivastava adopted these ideas for social networks with the introduction of α -proximity (Chester and Srivastava 2011).

The need for these progressively stronger adversarial models stems from the difficulty in deriving an analogous notion for equivalence of tuples when adopting the k -anonymization model from relational databases. Despite all these results, our three questions from the introduction of hardness, approximation guarantees, and class-specific analysis are unaddressed for the majority of these models. Our work here, and, indeed, our introduction of a vertex-addition approach, is meant to deepen understanding of the known adversarial models in the same manner that the work of Agarwal et al. (2005) and of Meyerson and Williams (2004) did in the relational database setting.

Subset anonymity is a relatively new consideration. Studies by Yuan et al. (2010) and by Ferri et al. (2012) reveal that privacy concerns and needs are varied across and among different user groups. This has prompted a couple recent works (Chester et al. 2012a, b) that consider the more general k -subset-anonymity.

Another set of related works is as follows. König (1936) showed that, given a graph \mathcal{G} with maximum degree d , it is always possible to convert \mathcal{G} into a d -regular graph H such that \mathcal{G} is an induced subgraph of H . In a subsequent paper, Erdős and Kelly (1967) strengthened the result of König by giving an efficient algorithm to determine the minimum number of new vertices that must be added to \mathcal{G} to obtain such a graph H . This results were extended by the more recent work of Akiyama et al. (1983) and Bodlaender et al. (2000). This result can be viewed as an optimal n -degree-anonymization of \mathcal{G} . Our approach in this paper can be viewed as a generalisation of these results with two relaxations. First, as discussed by Yuan et al. (2010), we may require that only a subset of nodes be anonymized. For example, different people have different expectations of privacy and some may not be concerned if they can be reidentified in the anonymous graph. Second, we require k -anonymity of the graph for an arbitrary k (which we typically assume to be some reasonably small value much less than d).

Our work here and the above references on k -anonymization, encompass at a high level “data anonymization”

methods. These methods first transform the data and then release them. We note lastly that there is another general family of methods for achieving data privacy, which does not release data, but, rather, only the output of an analysis computation. The released output is such that it is very difficult to infer from it any information about an individual input datum. The differentially-private methods (e.g., Dwork 2006; McSherry and Mironov 2009) belong in this family, which has a much different objective than anonymization.

7 Conclusion and future work

In this paper, we focused on the problem of k -degree-anonymizing a subset of a graph \mathcal{G} to protect against a degree-based attack. We established theoretical results that hold for any input graph by introducing a natural alternative to the formulation of the problem present in literature. Specifically, we introduced a vertex-addition approach wherein, given $\mathcal{G} = (V, E, \Sigma, \ell)$, the new graph $\mathcal{G}' = (V \cup V', E \cup E', \Sigma \cup \Sigma', \ell \cup \ell')$ is constructed such that $E' \subseteq V' \times (V' \cup V)$. The optimisation constraint is to minimise $|V'|$, which ensures that \mathcal{G}' is near to \mathcal{G} .

Within this setting, we established that k -degree-anonymizing a given subset X of V for $k > 3$ is *NP*-Complete, unless $|\Sigma| = 1$. For $|\Sigma| = 1$, we gave an efficient, exact algorithm for subset anonymization of X and an efficient, near-optimal algorithm to *also* ensure the k -degree-anonymity of all $X \cup V'$. Existing techniques in literature provide no characterisation of tractability nor approximation guarantees.

Furthermore, we demonstrated empirically that the resultant graph \mathcal{G}' of our algorithm is quite near to the input graph \mathcal{G} with respect to five structural characteristics of graphs. We contrasted the output graphs for our algorithm with those produced by the state-of-the-art algorithm which does not allow a vertex set expansion, achieving highly comparable performance. This suggests that the vertex-addition approach performs well in real-life in addition to the strong theoretical claims that we have proven.

The success here implies that there are interesting research directions with which to extend this work. We focused on foundational work with a simpler adversarial model to establish the credibility of adding “fake” vertices. Clearly, extending the algorithmic work to more challenging adversaries would be well worthwhile since structural graph characteristics that may equally well preserved by doing so. Our hardness result dictates that an appropriate pursuit for the case of vertex-labeled graphs would be to design approximate or heuristic algorithms. Finally, another research direction would be to enhance the algorithm by adding edges among $V' \times V'$ during the third

step of our algorithm in a utility-oriented approach such as that introduced by Wang et al. (2011).

References

- Adamic L, Glance N (2005) The political blogosphere and the 2004 u.s. election: divided they blog. In: Proceedings of WWW 2005 workshop on the weblogging ecosystem
- Aggarwal G, Feder T, Kenthapadi K, Motwani R, Panigrahy R, Thomas D, Zhu A (2005) Anonymizing tables. In: Proceedings of international conference on database theory (ICDT), pp 246–258
- Akiyama J, Era H, Harary F (1983) Regular graphs containing a given graph. *Am Math Month* 83:15–17
- Backstrom L, Dwork C, Kleinberg JM (2007) Wherefore art thou r3579x?: anonymized social networks, hidden patterns, and structural steganography. In: Proceedings of conference on world wide web (WWW), pp 181–190
- Barrat A, Weigt M (2000) On the properties of small-world network models. *Eur Phys J B* 13(3):547–560
- Bodlaender HL, Tan RB, van Leeuwen J (2000) Finding a delta-regular supergraph of minimum order. *Tech Rep UU-CS-2000-29*, Dept of Computer Science, Utrecht University, Utrecht
- Chakrabarti, D., Faloutsos, C (2006) Graph mining: laws, generators, and algorithms. *ACM Comput Surv* 38(1):2. doi:[10.1145/1132952.1132954](https://doi.org/10.1145/1132952.1132954)
- Cheng J, Fu AWC, Liu J (2010) K-isomorphism: privacy preserving network publication against structural attacks. In: Proceedings of ACM Special Interest Group on Management of Data (SIGMOD), pp 459–470
- Chester S, Srivastava G (2011) Social network privacy for attribute disclosure attacks. In: Proceedings of advances in social networks analysis and mining (ASONAM)
- Chester S, Kapron B, Ramesh G, Srivastava G, Thomo A, Venkatesh S (2011) k -anonymization of social networks by vertex addition. In: Proceedings of advances in databases and information systems (ADBIS)
- Chester S, Gaertner J, Stege U, Venkatesh S (2012a) Anonymizing subsets of social networks with degree constrained subgraphs. In: Proceedings of advances in social networks analysis and mining (ASONAM)
- Chester S, Kapron B, Srivastava G, Venkatesh S (2012b) Complexity of social network anonymization. *Soc Netw Anal Min*. doi:[10.1007/s13278-012-0059-7](https://doi.org/10.1007/s13278-012-0059-7)
- Costa LdF, Rodrigues FA, Travieso G, Villas Boas PR (2007) Characterization of complex networks: a survey of measurements. *Adv Phys* 56:167–242
- Domingo-Ferrer J (ed) (2002) Inference Control in statistical databases, from theory to practice. In: *Lecture Notes in Computer Science*, vol 2316. Springer, Berlin
- Dwork C (2006) Differential privacy. In: *ICALP*. Springer, Berlin, pp 1–12
- Erdős P, Kelly P (1967) The minimal regular graph containing a given graph. *Am Math Month* 70:1074–1075
- Estrada E, Rodriguez-Velazquez JA (2005) Spectral measures of bipartivity in complex networks. *Phys Rev E* 72(4):046105. doi:[10.1103/PhysRevE.72.046105](https://doi.org/10.1103/PhysRevE.72.046105)
- Faloutsos M, Faloutsos P, Faloutsos C (1999) On power-law relationships of the internet topology. *SIGCOMM Comput Commun Rev* 29(4):251–262. doi:[10.1145/316194.316229](https://doi.org/10.1145/316194.316229)
- Ferri F, Grifoni P, Guzzo T (2012) New forms of social and professional digital relationships: the case of facebook. *Soc Netw Anal Min* 2(2):121–137

- Girvan M, Newman MEJ (2002) Community structure in social and biological networks. *Proc Natl Acad Sci USA* 99:7821–7826
- González JJS (2002) Extending cell suppression to protect tabular data against several attackers. In: *Inference Control in Statistical Databases*, pp 34–58
- Hay M, Miklau G, Jensen D, Towsley DF, Weis P (2008) Resisting structural re-identification in anonymized social networks. *Proc Very Large Datab* 1(1):102–114
- Heer J (2005) Prefuse: a toolkit for interactive information visualization. In: *CHI 05: Proceedings of the SIGCHI conference on human factors in computing systems*. ACM Press, New York, pp 421–430
- König D (1936) Akademische verlagsgesellschaft. Leipzig
- Latora V, Marchiori M (2001) Efficient behavior of small-world networks. *Phys Rev Lett* 87. doi:[10.1103/PhysRevLett.87.198701](https://doi.org/10.1103/PhysRevLett.87.198701)
- Leskovec J, Kleinberg J, Faloutsos C (2005) Graphs over time: Densification laws, shrinking diameters and possible explanations. In: *Proceedings of international conference on knowledge discovery and data mining (KDD)*
- Leskovec J, Lang KJ, Dasgupta A, Mahoney MW (2008) Statistical properties of community structure in large social and information networks. In: *Proceedings of conference on world wide web (WWW)*, pp 695–704
- Li N, Li T, Venkatasubramanian S (2007) *t*-closeness: privacy beyond *k*-anonymity and *l*-diversity. In: *Proceedings of IEEE 23rd international conference on data engineering (ICDE07)*
- Liu K, Terzi E (2008) Towards identity anonymization on graphs. In: *Proceedings of ACM Special Interest Group on Management of Data (SIGMOD)*, pp 93–106
- Machanavajjhala A, Kifer D, Gehrke J, Venkatasubramanian M (2007) *L*-diversity: Privacy beyond *k*-anonymity. *ACM Trans. Knowl. Discov. Data* 1(1). doi:[10.1145/1217299.1217302](https://doi.org/10.1145/1217299.1217302)
- McSherry F, Mironov I (2009) Differentially private recommender systems: building privacy into the netflix prize contenders. In: *Proceedings of international conference on knowledge discovery and data mining (KDD)*, pp 627–636
- Meyerson A, Williams R (2004) On the complexity of optimal *k*-anonymity. In: *Principles of database systems*, pp 223–228
- Milgram S (1967) The small world problem. *Psychol Today* 2:60–67
- Newman MEJ (2006) Finding community structure in networks using the eigenvectors of matrices. *Phys Rev E* 74(3). doi:[10.1103/PhysRevE.74.036104](https://doi.org/10.1103/PhysRevE.74.036104)
- Robertson DA, Ethier R (2002) Cell suppression: experience and theory. In: *Inference control in statistical databases*, pp 8–20
- Sweeney L (2002) *k*-anonymity: A model for protecting privacy. *Int J Uncertainty Fuzziness Knowl Based Syst* 10(5):557–570
- Thompson B, Yao D (2009) The union-split algorithm and cluster-based anonymization of social networks. In: *Proceedings of ACM symposium on information, computer and communications security (ASIACCS)*, pp 218–227
- Wang Y, Xie L, Zheng B, Lee KCK (2011) Utility-oriented *k*-anonymization on social networks. In: *Proceedings of the 16th international conference on Database systems for advanced applications, vol Part I, DASFAA'11*. Springer, Berlin, pp 78–92
- Wu W, Xiao Y, Wang W, He Z, Wang Z (2010) *k*-symmetry model for identity anonymization in social networks. In: *Proceedings of international conference on extending database technology (EDBT)*, pp 111–122
- Ying X, Pan K, Wu X, Guo L (2009) Comparisons of randomization and *k*-degree anonymization schemes for privacy preserving social network publishing. In: *Proceedings of 3rd workshop on social network mining and analysis (SNA-KDD)*. ACM, New York, pp 10:1–10:10
- Yuan M, Chen L, Yu PS (2010) Personalized privacy protection in social networks. *Proc Very Large Datab* 4(2):141–150
- Zheleva E, Getoor L (2007) Preserving the privacy of sensitive relationships in graph data. In: *Proceedings of privacy, security, and trust in KDD (PinKDD)*, pp 153–171
- Zhou B, Pei J (2011) The *k*-anonymity and *l*-diversity approaches for privacy preservation in social networks against neighborhood attacks. *Knowledge Information Systems* 28(1):47–77