

Towards Graph Analytic and Privacy Protection

Dongqing Xiao

A Proposal for a PhD Dissertation in Computer Science

Worcester Polytechnic Institute, Worcester, MA

Dec 2016

Committee Members:

Dr. Mohamed Y. Eltabakh, Assistant Professor, Worcester Polytechnic Institute.

Dr. Elke A. Rundensteiner, Professor, Worcester Polytechnic Institute. Advisor.

Dr. Xiangnan Kong, Assistant Professor, Worcester Polytechnic Institute.

Dr. Yuanyuan Tian, Researcher, IBM Almaden, External member

Abstract

Talking about the focus of this researchh

Each component with its item

Experimental studies including performance evaluation and user studies will be conducted on various real datasets to evaluate both the effectiveness and efficiency of the proposed approaches.

Contents

1	Introduction	2
1.1	Motivation	2
1.2	State-Of-the-Art	2
1.3	Challenges	3
2	Preliminaries	4
2.1	Distributed Triangle Listing	4
2.1.1	Triangle Listing Problem	4
2.1.2	Sequential Triangle Listing	4
2.1.3	MapReduce Overview	6
2.1.4	Triangle Listing in MapReduce	7
2.2	Privacy Protection on Uncertain Graphs	8
2.2.1	Uncertain Graph	8
2.2.2	Attack Model	8
2.2.3	Privacy Criteria	9
2.2.4	Reliability-based Utility Criteria	11
3	Distributed Triangle Listing Algorithm	13
3.1	Overview	13
3.2	Bermuda Edge-Centric Node++	14
3.3	Bermuda Vertex-Centric Node++	18
3.3.1	Message Sharing Management	21
3.3.2	Discussion	24
4	Degree Anonymization over Uncertain Graphs	26
4.1	Naive Approach: Anonymization via Representative Instance	26
4.2	Reliability-Preserving anonymization on uncertain graphs	28

4.2.1	Anonymization Procedure	28
4.2.2	Reliability-oriented Edge Selection	30
4.2.3	Anonymity-oriented Edge Perturbing	36
Bibliography		39

Acknowledgments

I express my sincere thanks to my advisor Prof. Mohamed Eltabakh. My thank you goes to the members of my Ph.D. committee, Prof. Elke Rundensteiner, Prof. Xiangnan Kong and Dr. Yuanyuan Tian. My thank you also goes to all other previous and current DSRG members for their useful discussion and feedback.

Chapter 1

Introduction

1.1 Motivation

Network data has become increasingly important in recent years of the greater importance of various application domains such as Web, social network, biological networks, and communication networks. The semantics and the interpretation of the nodes and the links may vary with application domain, *e.g.*, in a social network node can represent individuals and their connection captures friendship, while in a protein-to-protein network (PPI) nodes are proteins and connections refer to their interactions. These graph data carries valuable information and are analyzed in various ways to exact knowledge. For example, social networks provide significant insight about psychological behavior of individuals or PPI network are widely studied to elucidate mechanism of diseases.

Before such graph data can be released for XX .

1.2 State-Of-the-Art

XX. A few recent efforts have been proposed to support distributed XX. Some utilize different distributed computing paradigms that, unlike MapReduce, either suffer from a bottleneck by requiring a central master node to split and broadcast data to each slave node XX, or they allow all nodes unrestricted exchange of data with each other XXX. In contrast, the MapReduce infrastructure neither assumes a central node nor does it allow data exchange among the mappers (nor among the reducers) to enable easier distribution of tasks and higher scalability. Therefore none of these techniques is applicable on the MapReduce infrastructure. To support the distributed outlier detection on MapReduce,

partitioning approaches that make sure each reducer can detect outliers independently from other reducers have to be designed.

In the literature Map-Reduced based approaches have been proposed for XXX techniques such as XXX. Although they investigate the key concepts in distributed systems such as *load balancing* and *efficient partitioning* that determine the efficiency of distributed mining algorithms, their methods cannot be applied in our outlier detection area as shown below.

For example XXX.

In recent years researchers started to look at the problem of XXX []. Specifically [] proposed solutions for XX. [] improves upon this solution [] by now XX. All solutions leverage the overlap of sliding windows and thus avoid XX.

However, these existing techniques [] didn't explore the optimization opportunities enabled by the *critical observation* below. That is, they didn't exploit the fact that outliers by nature XX.

Furthermore, all of the above approaches focus on handling XX with. The simultaneous execution of XX.

In the broader area of XX, the main focus of previous work has been on XX. Their methods include XX. However the key problem we aim to address in this dissertation is different from the more general purpose optimization effort required by the traditional SQL query sharing. XXX.

1.3 Challenges

Scaling the outlier detection techniques to big data is challenging.

Second, designing an effective partitioning strategy for XX approach is challenging.

Third, it is challenging to XXX. First, due to the algorithmic complexity of mining techniques [?], processing each outlier request from scratch over big datasets each time when it is submitted clearly cannot satisfy the response time requirement of interactive system. On the other hand XXX.

Chapter 2

Preliminaries

In my dissertation, I will study triangle listing problem for web-scale graphs and privacy protection issue in the context of uncertain graphs. For the later topics, we focus on resisting the degree-based node re-identification. Here, we give the definitions of uncertain graph and attack model, privacy criteria related with uncertain graph anonymization problem. We also give the formal formulation of proposed problems.

2.1 Distributed Triangle Listing

2.1.1 Triangle Listing Problem

Suppose we have a simple undirected graph $G(V, E)$, where V is the set of vertices (nodes), and E is the set of edges. Let $n = |V|$ and $m = |E|$. Let $N_v = \{u | (u, v) \in E\}$ denote the set of *adjacent nodes* of node v , and $d_v = |N_v|$ denote the degree of node v . We assume that G is stored in the most popular format for graph data, *i.e.*, the adjacency list representation. Given any three distinct vertices $u, v, w \in V$, they form a triangle \triangle_{uvw} , *iif* $(u, v), (u, w), (v, w) \in E$. We define the set of all triangles that involve node v as $\triangle(v) = \{\triangle_{uvw} | (v, u), (v, w), (u, w) \in E\}$. Similarly, we define $\triangle(G) = \bigcup_{v \in V} \triangle(v)$ as the set of all triangles in G .

Problem 1 *Triangle Listing Problem:* *Given a large-scale distributed graph $G(V, E)$, our goal is to report all triangles in G , *i.e.*, $\triangle(G)$, in a highly distributed way.*

2.1.2 Sequential Triangle Listing

Algorithm 1 NodeIterator++

Preprocessing step

```
1: for all  $(u, v) \in E$  do  
2:   if  $u \succ v$ , store  $u$  in  $N_v^H$   
3:   else store  $v$  in  $N_u^H$   
4: end for
```

Triangle Listing

```
5:  $\Delta(G) \leftarrow \emptyset$   
6: for all  $v \in V$  do  
7:   for all  $u \in N_v^H$  do  
8:     for all  $w \in N_v^H \cap N_u^H$  do  
9:        $\Delta(G) \leftarrow \Delta(G) \cup \{\Delta_{vuw}\}$   
10:    end for  
11:  end for  
12: end for
```

Sequental triangle listing algorithms have been extensively studied. Here, we present a sequential triangle listing algorithm which is widely used as the basis of parallel approaches [18, 2, 14]. In this work, we also use it as the basis of our distributed approach.

A naive algorithm for listing triangles is as follows. For each node $v \in V$, find the set of edges among its neighbors, i.e., pairs of neighbors that complete a triangle with node v . Given this simple method, each triangle (u, v, w) is listed six times—all six permutations of u, v and w . Several other algorithms have been proposed to improve on and eliminate the redundancy of this basic method, e.g., [17, 5].

One of the algorithms, known as *NodeIterator++* [17], uses a total ordering over the nodes to avoid duplicate listing of the same triangle. By following a specific ordering, it guarantees that each triangle is counted only once among the six permutations. Moreover, the *NodeIterator++* algorithm adopts an interesting node ordering based on the nodes' degrees, with ties broken by node IDs, as defined below:

$$u \succ v \iff d_u > d_v \text{ or } (d_u = d_v \text{ and } u > v)$$

This degree-based ordering improves the running time by reducing the diversity of the effective degree \hat{d}_v . The running time of *NodeIterator++* algorithm is $O(m^{3/2})$. A comprehensive analysis can be found in [17].

The standard *NodeIterator++* algorithm performs the degree-based ordering comparison during the final phase, i.e., the triangle listing phase. The work in [2] and [18] further

improves on that by performing the comparison $u \succ v$ for each edge $(u, v) \in E$ in the preprocessing step (Lines 1-3, Algorithm 1). For each node v and edge (u, v) , node u is stored in the effective list of v (N_v^H) if and only if $u \succ v$, and hence $N_v^H = \{u : u \succ v \text{ and } (u, v) \in E\}$. The preprocessing step cuts the storage and memory requirement by half since each edge is stored only once. After the preprocessing step, the effective degree of nodes in G is $O(\sqrt{m})$ [17]. Its correctness proof can be found in [2]. The modified NodeIterator++ algorithm is presented in Algorithm 1.

2.1.3 MapReduce Overview

MapReduce is a popular distributed programming framework for processing large datasets [10]. MapReduce, and its open-source implementation Hadoop [19], have been used for many important graph mining tasks [18, 14]. In this paper, our algorithms are designed and analyzed in the MapReduce framework.

Computation Model. An analytical job in MapReduce executes in two rigid phases, called the *map* and *reduce* phases. Each phase consumes/produces records in the form of *key-value* pairs—We will use the keywords *pair*, *record*, or *message* interchangeably to refer to these key-value pairs. A pair is denoted as $\langle k; val \rangle$, where k is the key and val is the value. The *map* phase takes one key-value pair as input at a time, and produces zero or more output pairs. The *reduce* phase receives multiple key-listOfValues pairs and produces zero or more output pairs. Between the two phases, there is an implicit phase, called *shuffling/sorting*, in which the mappers’ output pairs are shuffled and sorted to group the pairs of the same key together as input for reducers.

Our proposed solution will leverage and extend some of the basic functionality of MapReduce, which are:

- **Key Partitioning:** Mappers employ a *key partitioning function* over their outputs to partition and route the records across the reducers. By default, it is a hash-based function, but can be replaced by any other user-defined logic.
- **Multi-Key Reducers:** Typically, the number of distinct keys in an application is much larger than the number of reducers in the system. This implies that a single reducer will sequentially process multiple keys—along with their associated groups of values—in the same reduce instance. Moreover, the processing order is defined by *key sorting function* used in shuffling/sorting phase. By default, a single reduce

Algorithm 2 MR-Baseline

Map: Input: $\langle v; N_v^H \rangle$
1: emit $\langle v; (v, N_v^H) \rangle$
2: **for all** $u \in N_v^H$ **do**
3: emit $\langle u; (v, N_v^H) \rangle$
4: **end for**

Reduce: Input: $[\langle u; (v, N_v^H) \rangle]$
5: initiate N_u^H
6: **for all** $\langle u; (v, N_v^H) \rangle$ **do**
7: **for all** $w \in N_u^H \cap N_v^H$ **do**
8: emit \triangle_{vuw}
9: **end for**
10: **end for**

instance processes each of its input groups in total isolation from the other groups with no sharing or communication.

2.1.4 Triangle Listing in MapReduce

Both [18] and [2] use the NodeIterator++ algorithm as the basis of their distributed algorithms. [18] identifies the triangles by checking the existence of pivot edges, while [2] uses set intersection of effective adjacency list (Line 7, Algorithm 1). In this section, we present the MapReduce version of the NodeIterator++ algorithm similar to the one presented in [2], referred to as *MR-Baseline* (Algorithm 2).

The general approach is the same as in the NodeIterator++ algorithm. In the map phase, each node v needs to emit two types of messages. The first type is used for the initiation its own effective adjacency list in the reduce side, referred to as a *core message* (Line 1, Algorithm 2). The second type is used for identifying triangles, referred to as *pivot messages* (Lines 2-3, Algorithm 2). All pivot messages from v to its effective adjacent nodes are identical. In the reduce phase, each node u will receive a core message from itself, and a pivot message from adjacent nodes with the lower degree. Then, each node identifies the triangles by performing a set intersection operation (Lines 5-6, Algorithm 2).

We omit the code of the pre-processing procedure since its implementation is straightforward in MapReduce. In addition, we will exclude the pre-processing cost for any further consideration since it is typically dominated by the actual running time of the triangle

listing algorithm, plus it is the same overhead for all algorithms.

2.2 Privacy Protection on Uncertain Graphs

2.2.1 Uncertain Graph

Let $\mathcal{G} = (V, E, p)$ be an uncertain graph, where V is the set of nodes, E is the set of edges, and function $p : E \rightarrow [0, 1]$ assigns a probability of existence to each edge, denoted as $p(e)$. We consider the edge probabilities are independent, and we assume *possible world* semantics, consistent with existing literature [16, 1, 11]. Specifically, the *possible world* semantics interprets \mathcal{G} as a set of possible deterministic graphs $W(\mathcal{G})$, each deterministic graph $G \in W(\mathcal{G})$ including all vertices of \mathcal{G} and a subset of edges $E_G \subset E$. The probability of observing any possible world $G = (V, E_G) \in W(\mathcal{G})$ is

$$Pr[G] = \prod_{e \in E_G} p(e) \prod_{e \in E \setminus E_G} (1 - p(e))$$

We assume that the uncertain graph is simple, *i.e.*, the uncertain graph is undirected containing no self-loops.

2.2.2 Attack Model

Our uncertain graph anonymization problem is related to the conventional graph anonymization problem, which has been extensively studied. As Hay *et al.* pointed out, the structural information such as node degree, 1-neighborhood can be utilized to breach user privacy. For example, [] consider the case that , while [] consider the case that XX. Different anonymization techniques are designed to resist different kinds of privacy attack.

The first step to anonymization is to know what external information of a graph may be acquired by an adversary. In our work, we assume that the information of victim node degree is easy to be collected by the adversary, consistent to the literature []. In fact, degree-based de-anonymization is one of the most serious privacy attacks in the context of deterministic graphs. The knowledge could also be understood as some *assertion* of the individual, which could be evaluated to *true* or *false* based on the topology structure of the network.

In the context of deterministic graphs, such knowledge can be expressed as one assertion with certainty such as Ana has 3 neighbors. In the context of uncertain graphs, such

knowledge can be a little different. In an uncertain graph, for a given node v , its degree d_v is defined in a probabilistic way. For example, node a in Figure ??, its node degree may be 3 or 2 in different possible worlds.

It is difficult to model the specific degree-related knowledge used by the adversaries by advance. In our work, we first consider the adversary can and only acquire the expected node degree of victim nodes. For example, the adversary knows that the expected value of the number of Ana's neighbors should 3 as shown in Figure ?. Later, we consider the adversary can acquire more comprehensive knowledge of node degree of victim nodes such as the distribution function. For example, the adversary may know the distribution of the number of Ana's neighbors as $(3 :)$, $(2 :)$, $(1 :)$, $(0 :)$. In this work, we will present the corresponding solution to resist ever-discussed privacy attack in the context of uncertain graphs.

The second step to anonymization is to know how the adversary links the external knowledge such as node degree to a node in the perturbed graph. In the deterministic graph, degree-based node re-identification is straightforward. Let \mathcal{G} be a graph and $\hat{\mathcal{G}}$ be the anonymized version of the given graph. The adversary can get a number of `match` vertices in $\hat{\mathcal{G}}$, *i.e.*, the nodes U_c with the `match` degree. When $\hat{\mathcal{G}}$ is a deterministic one, for a given node in $\hat{\mathcal{G}}$, the assertion is evaluated either *true* or *false* with certainty. When $\hat{\mathcal{G}}$ is an uncertain one, the assertion is evaluated based on the graph structure combined with possible world semantics and is defined a probabilistic way. For example, it is said that the probability that the node a was observed with degree 3 is 0.504. Following this path, we consider the matching attack is performed between a random variable, namely, find the equal random variable. Namely, for a given node u in $\hat{\mathcal{G}}$, it is defined as the probability that the random variable d_u is equal to the adversary knowledge d_v , denotes as $P(d_u = d_v)$. It can be computed as

$$P(d_u = d_v) = \sum_{i \in Z} P(d_u = i)P(d_v = i) \quad (2.1)$$

2.2.3 Privacy Criteria

As ever discussed, we want to publish the anonymized version of uncertain graphs where node identity is obfuscated enough. We adopt (k, ϵ) -*obf* model, proposed in [7], to quantify the anonymity level achieved by an anonymized graph (uncertain graph). Its formal definition is as follow:

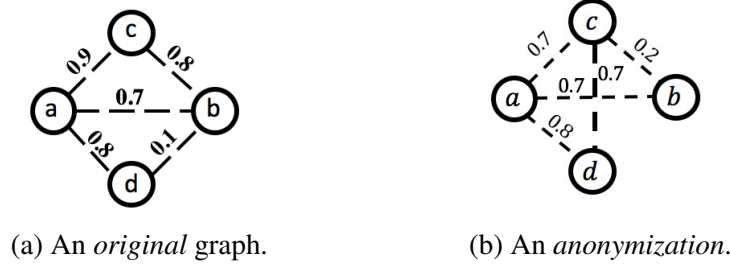


Figure 2.1: Uncertain graph anonymization.

The probability that a node v with degree value ω

$$Y_v(\omega) := \frac{X_v(\omega)}{S(\omega)}$$

Normalization

The probability that a node v is the image of the target node has degree value ω

$X_v(\omega)$	deg = 0	deg = 1	deg = 2	deg = 3
A	0.006	0.092	0.398	0.504
B	0.054	0.348	0.542	0.056
C	0.020	0.260	0.720	0.000
D	0.180	0.740	0.080	0.000
$S(\omega)$	0.260	1.440	1.740	0.560

The probability that a node v is the image of the target node has degree value ω

$Y_v(\omega)$	deg = 0	deg = 1	deg = 2	deg = 3
A	0.023	0.064	0.229	0.900
B	0.208	0.242	0.311	0.100
C	0.077	0.180	0.414	0.000
D	0.692	0.514	0.046	0.000
$H(\omega)$	1.248	1.688	1.742	0.469

Figure 2.2: The degree uncertainty for each node (top) and normalized values for each degree (bottom)

Definition 1 ((k, ϵ) -obf [7]) Let P be a vertex property, $k \geq 1$ be a desired level of anonymity, and $\epsilon > 0$ be a tolerance parameter. The uncertain graph $\tilde{\mathcal{G}}$ is said to k -obfuscate a given vertex $v \in V_{\tilde{\mathcal{G}}}$ with respect to P if the entropy of the distribution $Y_{P(v)}$ over the vertices of $\tilde{\mathcal{G}}$ is greater than or equals to $\log_2 k$:

$$H(Y_{P(v)}) \geq \log_2 k.$$

The uncertain graph \mathcal{G} is (k, ϵ) -obf with respect to property P if it k -obfuscates at least $(1 - \epsilon)|V|$ vertices in $V_{\mathcal{G}}$. P can be any node properties.

Figure 2.2 gives an example of how to compute degree entropy for the uncertain graph in Figure 2.1(a). Each row in the top side is the degree distribution for the corresponding node. For example, node a has degree 0 with probability $(1 - 0.9)(1 - 0.8)(1 - 0.7) = 0.006$. The last row in the top side is the sum of probability over all the nodes. For

example, $S(0) = 0.006 + 0.054 + 0.020 + 0.180 = 0.26$. The bottom one normalizes values in each column to get distribution $Y_{P(v)}$. For example, $Y_0(a) = \frac{X_0(a)}{S(0)} = 0.023$. The entropy $H(Y_d(v))$ for each degree value ω is shown in the last bottom row. Given $k = 3$, $\log_2 3 = 1.58$, node b, c with expected degree 2 and d with true degree 1 satisfies k -obf while a with expected degree 3 fails, thus $\epsilon = 0.25$.

2.2.4 Reliability-based Utility Criteria

Inspired by the importance of reliability, we suggest to quantify the utility loss, anonymization cost, as the reliability difference between the anonymized result $\tilde{\mathcal{G}}$ and the original one \mathcal{G} , reliability discrepancy. We proceed to explain its formal definition.

In the context of uncertain graphs, *reliability* generalizes the concept of connectivity by capturing the probability that two given (sets of) nodes are reachable over all possible worlds. We focus on two-terminal reliability.

Definition 2 Two-Terminal Reliability [9] Given an uncertain graph \mathcal{G} , and two distinct nodes u and v in V , the reliability for two nodes (u, v) is defined as follows:

$$R_{u,v}(\mathcal{G}) = \sum_{G \subseteq W(\mathcal{G})} \mathcal{I}_G(u, v) \Pr[G]$$

where $\mathcal{I}_G(u, v)$ is 1 when u and v are contained in a connected component in G , and 0 otherwise.

Definition 3 Two-Terminal Reliability Discrepancy The reliability discrepancy of two distinct nodes (u, v) is defined as the absolute difference of corresponding reliability values in the original graph \mathcal{G} and the anonymized result $\tilde{\mathcal{G}}$. Its formal definition is as follows:

$$RD_{u,v}(\tilde{\mathcal{G}}) = |R_{u,v}(\mathcal{G}) - R_{u,v}(\tilde{\mathcal{G}})|$$

Thus, we come up the definition of reliability discrepancy of the overall anonymized result against the true graph.

Definition 4 Reliability Discrepancy Compared to the input uncertain graph $\mathcal{G} = (V, E, p)$, the reliability discrepancy of one anonymized instance $\tilde{\mathcal{G}} = (V, E, \tilde{p})$ is the sum of RD over all vertex pair (u, v) . The formal definition is as follows:

$$\Delta(\tilde{\mathcal{G}}) = \sum_{(u,v)} RD_{u,v}(\tilde{\mathcal{G}})$$

Example 1 Given an input uncertain graph \mathcal{G} in Figure 2.1(a) and its anonymization $\tilde{\mathcal{G}}$ in Figure 2.1(b) and suppose one wants to compute $RD_{a,b}(\tilde{\mathcal{G}})$. Following the definition, the reliability of node pairs (a, b) for \mathcal{G} in Figure 2.1(a) equals 0.92272. While, the reliability of node pairs in Figure 2.1(b) equals 0.75208. Thus, the reliability discrepancy of node pair $RD_{a,b}(\tilde{\mathcal{G}})$ is $|0.92272 - 0.75208| = 0.17064$.

We are ready to formulate the problem, reliability preserving anonymization in the context of uncertain graphs.

Problem 2 Reliability Preserving Anonymization over Uncertain graphs Given an uncertain graph $\mathcal{G} = (V, E, p)$ and anonymization parameters k, ϵ , the objective is to find a possible (k, ϵ) -obfuscated graph $\tilde{\mathcal{G}} = (V, E, \tilde{p})$ with minimal reliability discrepancy $\Delta(\tilde{\mathcal{G}})$, as

$$\underset{\tilde{\mathcal{G}}}{\operatorname{argmin}} \quad \Delta(\tilde{\mathcal{G}})$$

Subject to $\tilde{\mathcal{G}}$ is (k, ϵ) -obj

Chapter 3

Distributed Triangle Listing Algorithm

In this chapter, I will briefly review the distributed triangle listing techniques I developed for MapReduce Framework. These techniques have been submitted as a conference paper [1].

In this work, we present Bermuda method which effectively reduces the size of the intermediate data via redundancy elimination and sharing of messages whenever possible. As a result, Bermuda achieves orders-of-magnitudes of speedup and enables processing larger graphs that other techniques fail to process under the same resources. Bermuda exploits the locality of processing, i.e., in which reduce instance each graph vertex will be processed, to avoid the redundancy of generating messages from mappers to reducers. Bermuda also proposes novel message sharing techniques within each reduce instance to increase the usability of the received messages.

3.1 Overview

With the MR-Baseline Algorithm, one node needs to send the same pivot message to multiple nodes residing in the same reducer. Therefore, the network traffic can be reduced by sending only one message to the destination reducer, and either have main-memory cache or distributing the message to actual graph nodes within each reducer. Although this strategy seems quite simple, and other systems such as GPS and X-Pregel [4, ?] have implemented it, the trick lies on how to efficiently perform the caching and sharing. In this section, we propose new effective caching strategies to maximize the sharing benefit while encountering little overhead. We also present novel theoretical analysis for the proposed techniques.

In the frameworks of GPS and X-Pregel, adjacency lists of high degree nodes are used for identifying distinct destination reducer and distributing the message to target nodes in the reduce side. This method requires extensive memory and computations for message sharing. In contrast, in Bermuda, each node uses the *universal key partition function* to group its destination nodes. Thus, each node would only send the same pivot message to each reduce instance only once. At the same time, reduce instances will adopt different message-sharing strategies to guarantee the correctness of algorithm. As a result, Bermuda achieves a trade off between reducing the network communication—which is known to be a big bottleneck for map-reduce jobs—and increasing the processing cost and memory utilization.

3.2 Bermuda Edge-Centric Node++

A straightforward (and intuitive) approach for sharing the pivot messages within each reduce instance is to organize either the pivot or core messages in main-memory for efficient random access. We propose the Bermuda Edge-Centric Node++ (Bermuda-EC) algorithm, which is based on the observation that for a given input graph, it is common to have the number of core messages smaller than the number of pivot messages. Therefore, the main idea of Bermuda-EC algorithm is to first read the core messages, cache them in memory, and then stream the pivot messages, and on-the-fly intersect the pivot messages with the needed core messages (See Figure 3.1). The MapReduce code of the Bermuda-EC algorithm is presented in Algorithm 3.

In order to avoid pivot message redundancy, a universal key partitioning function is utilized by mappers. The corresponding modification in the map side is as follows. First, each node v employs a universal key partitioning function $h()$ to group its destination nodes (Line 3, Algorithm 3). This grouping captures the graph nodes that will be processed by the same reduce instance. Then, each node v sends a pivot message including the information of N_v^H to each non-empty group (Lines 4-6, Algorithm 3). Following this strategy, each reduce instance receives each pivot message exactly once even if it will be referenced multiple times.

Moreover, we use tags to distinguish core and pivot messages, which are not listed in the algorithm for simplicity. Combined with the MapReduce internal sorting function, Bermuda-EC guarantees that all core messages are received by the reduce function before any of the pivot messages as illustrated in Figure 3.1. Therefore, it becomes feasible to

Algorithm 3 Bermuda-EC

Map: Input: $\langle v; N_v^H \rangle$
Let $h(\cdot)$ be a key partitioning function into $[0, k-1]$

- 1: $j \leftarrow h(v)$
- 2: emit $\langle j; (v, N_v^H) \rangle$
- 3: *Group* the set of nodes in N_v^H by $h(\cdot)$
- 4: **for all** $i \in [0, k-1]$ **do**
- 5: **if** $gp_i \neq \emptyset$ **then**
- 6: emit $\langle i; (v, N_v^H) \rangle$
- 7: **end if**
- 8: **end for**

Reduce: Input: $[\langle i; (v, N_v^H) \rangle]$

- 9: initiate all the core nodes' N_u^H in main memory
- 10: **for all** pivot message $\langle i; (v, N_v^H) \rangle$ **do**
- 11: **for all** $u \in N_v^H$ and $h(u) = i$ **do**
- 12: **for all** $w \in N_v^H \cap N_u^H$ **do**
- 13: emit \triangle_{vuw}
- 14: **end for**
- 15: **end for**
- 16: **end for**

cache only the core messages in memory, and then perform the intersection as the pivot messages are received.

The corresponding modification in the reduce side is as follows. For a given reduce instance R_i , it first reads all the core message into main-memory (Line 7, Algorithm 3). Then, it iterates over all pivot message. Each pivot message is intersected with the cached core messages for identifying the triangles. As presented in the MR-Baseline algorithm (Algorithm 2), each pivot message (v, N_v^H) needs to be processed in reduce instance R_i only for nodes $u : u \in N_v^H$ where $h(u) = i$. Interestingly, this information is encoded within the pivot message. Thus, each pivot message is processed for all its requested core nodes once received (Lines 9-11, Algorithm 3).

Analysis of Bermuda Techniques

Extending the analysis in Section 2.1.4, we demonstrate that Bermuda-EC achieves improvement over MR-Baseline w.r.t both space usage and execution efficiency. Furthermore, we discuss the effect of the number of reducers k on the algorithm performance.

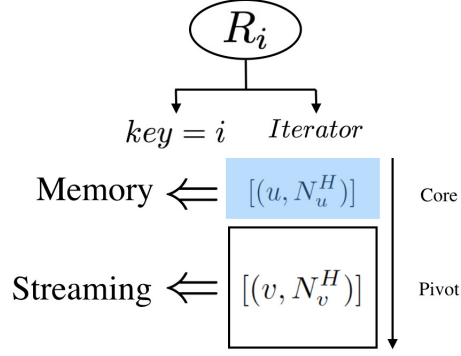


Figure 3.1: Bermuda-EC Execution.

Theorem 1 *For a given number of reducers k , we have:*

- *The expected total size of the map output is $O(km)$.*
- *The expected size of core messages to any reduce instance is $O(m/k)$.*

Proof 1 *As shown in Algorithm 3, the size of the map output generated by node v is at most $k * \hat{d}_v$. Thus, the total size of the map output T is as follows:*

$$T < \sum_{v \in V} k \hat{d}_v = k \sum_{v \in V} \hat{d}_v = km$$

*For the second bound, observe that a random edge is present in a reduce instance R_i and represented as a core message with probability $1/k$. By following the Linearity of Expectation, the expected number of the core messages to any reduce instance is $O(m * \frac{1}{k})$.*

Space Usage. Theorem 1 shows that when $k \ll \sqrt{m}$ (the usual case for massive graphs), then the total size of the map output generated by Bermuda-EC algorithm is significantly less than that generated by the MR-Baseline algorithm. In other words, Bermuda -EC is able to handle even larger graphs with limited compute clusters. **Execution Time.** A positive consequence of having a smaller intermediate result is that it requires less time for generating and shuffling/sorting the data. Moreover, the imbalance of the map outputs is also reduced significantly by limiting the replication factor of the pivot messages up to k . The next theorem shows the approximate variance of the number of the intermediate result from mappers. When $k < E(x)$, it implies smaller variance among the mappers than that of the MR-Baseline algorithm. Together, Bermuda -EC achieves better performance and scales to larger graphs compared to the MR-Baseline algorithm.

Theorem 2 For a given graph $G(V, E)$, let a random variable x denotes the effective degree of any node in G and the variance of x is denotes as $\text{Var}(x)$. Then the expectation of x ($E(x)$) equals the average degree and computed as $E(x) = \frac{m}{n}$. For typical graphs, $\text{Var}(x) \neq 0$ and $E(x) \neq 0$ always hold. Since each mapper starts with approximately the same input size (say receives c graph nodes), the variance of the map output's size under the Bermuda-EC Algorithm is $O(2ck^2\text{Var}(x))$, where k represents the number of reducers.

Proof 2 Assume the number of reducers is k . Given a graph node v , where its effective degree $\hat{d}_v = x$. Let random variable $y(x)$ be the number of distinct reducers processing the effective neighbors of v , and thus $y(x) \leq k$. Then, the size of the map output generated by a single node u would be xy , denoted as $g(x)$ (Lines 3-4, Algorithm 3). Thus, the total size of the map output generated by c nodes in a single mapper $T(X) = \sum_{i=1}^c g(x_i)$. Since x_1, x_2, \dots, x_c are independent and identically distributed random variables, then $\text{Var}(T(x)) = c * \text{Var}(g(x))$. The approximate variance of $g(x)$ is as follows

$$\begin{aligned} \text{Var}(xy) &= E(x^2y^2) - E(xy)^2 \\ &< E(x^2y^2) \\ &< k^2 E(x^2) \\ &< k^2 (E(x)^2 + \text{Var}(x)) \\ &< 2k^2 \text{Var}(x) \end{aligned}$$

As presented in [17], $E(x^2) \approx \frac{m^{\frac{3}{2}}}{n}$ and $E(x) = \frac{m}{n}$. Thus $\frac{E(x^2)}{E(x)^2} \approx \frac{n}{\sqrt{m}}$. In many real graphs where $n^2 > m$ it implies $\frac{n}{\sqrt{m}} > \sqrt{m} > 2$. It implies $E(x^2) > 2E(x)^2$, thus $\text{Var}(x) = E(x^2) - E(x)^2 > E(x)^2$.

We now study in more details the effect of parameter k (the number of reducers) on the space and time complexity for the Bermuda-EC algorithm.

Effect on Space Usage. The reducers number k trades off the memory used by a single reduce instance and the size of the intermediate data generated during the MapReduce job. The memory used by a single reducer should not exceed the available memory of a single machine, i.e., $O(m/k)$ should be sub-linear to the size of main memory in a single machine. In addition, the total space used by the intermediate data must also remain bounded, i.e., $O(km)$ should be no larger than the total storage. Given a cluster of machines, these two constraints define the bounds of k for a given input graph $G(V, E)$.

Effect on Execution Time. The reducers number k trades off the reduce computation time and the time for shuffling and sorting. As the parallelization degree k increases, it reduces the computational time in the reduce phase. At the same time, the size of the intermediate data, *i.e.*, $O(km)$ increases significantly as k increases (notice that m is very large), and thus the communication cost becomes a bottleneck in the job’s execution. Moreover, the increasing variance among mappers $O(2ck^2Var(x))$ implies a more significant straggler problem which slows down the execution progress.

In general, Bermuda-EC algorithm favors the smaller setting of k for higher efficiency while subjects to memory bound that the expected size of core message $O(m/k)$ should not exceed the available memory of a single reduce instance.

Unfortunately, for processing web-scale graphs such as *ClueWeb* with more than 80 billion edges (and total size of approximately 700GBs)—which as we will show the state-of-art techniques cannot actually process—the number of reducers needed for Bermuda-EC for acceptable performance is in the order of 100s. Although, this number is very reasonable for most mid-size clusters, the intermediate results $O(km)$ will be huge, which leads to significant network congestion.

Disk-Based Bermuda-EC: A generalization to the proposed Bermuda-EC algorithm that guarantees no failure even under the case where the core messages cannot fit in a reducer’s memory is the *Disk-Based Bermuda-EC* variation. The idea is straightforward and relies on the usage of the local disk of each reducer. The main idea is as follows: (1) Partition the core messages such that each partition fits into main memory, and (2) Buffer a group of pivot messages, and then iterate over the core messages one partition at a time, and for each partition, identify the triangles as in the standard Bermuda-EC algorithm. Obviously, such method trades off between disk I/O (pivot message scanning) and main-memory requirement. For a setting of reduce number k , the expected size of core messages in a single reduce instance is $O(m/k)$, thus the expected number of rounds is $O(\frac{m}{kM})$ where M represents the size of available main-memory for single reducer. The expected size of pivot message reaches $O(m)$. Therefore, the total disk I/O reaches $O(\frac{m^2}{kM})$. In the case of massive graph, it implies longer time.

3.3 Bermuda Vertex-Centric Node++

As discussed in Section 3.2, the Bermuda-EC algorithm assumes that the core messages can fit in the memory of a single reducer. However, it is not always guaranteed to be the

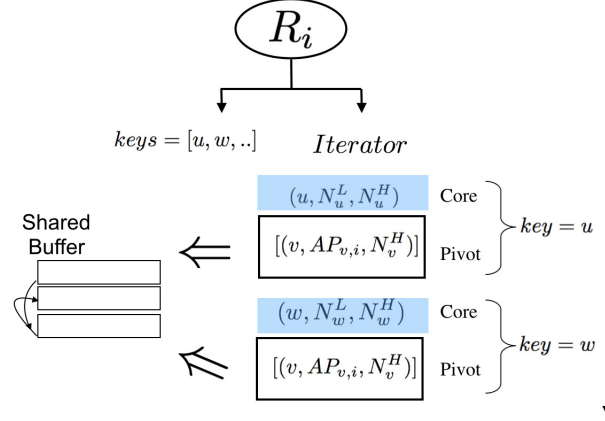


Figure 3.2: Bermuda-VC Execution.

case, especially in web-scale graphs.

One crucial observation is that the access pattern of the pivot messages can be learned and leveraged for better re-usability. In MapReduce, a single reduce instance processes many keys (graph nodes) in a specific sequential order. This order is defined based on the key comparator function. For example, let $h()$ be the key partitioning function and $l()$ be key comparator function within the MapReduce framework, then $h(u) = h(w) = i$ and $l(u, w) < 0$ implies that the reduce instance R_i is responsible for the computations over nodes u , w , and also the computations of node u precede that of node w . Given these known functions, the relative order among the keys in the same reduce instance becomes known, and the access pattern of the pivot message can be predicted. The knowledge of the access pattern of the pivot messages holds a great promise for proving better caching and better memory utilization.

Inspired by these facts, we propose the Bermuda-VC algorithm which supports random access over the pivot messages by caching them in main-memory while streaming in the core messages. More specifically, Bermuda-VC will reverse the assumption of Bermuda-EC, where we now try to make the pivot messages arrive first to reducers, get them cached and organized in memory, and then the core messages are received and processed against the pivot messages. Although the size of the pivot messages is usually larger than that of the core messages, their access pattern is more predictable which will enable better caching strategies as we will present in this section. The Bermuda-VC algorithm is presented in Algorithm 4.

The Bermuda-VC algorithm uses a shared buffer for caching the pivot messages. And

Algorithm 4 Bermuda-VC

Map: Input: $(\langle v; (N_v^L, N_v^H) \rangle)$
Let $h(\cdot)$ be a key partitioning function into $[0, k-1]$
Let $l(\cdot)$ be a key comparator function

- 1: emit $\langle v; (v, N_v^L, N_v^H) \rangle$
- 2: *Group* the set of nodes in N_v^H by $h(\cdot)$
- 3: **for all** $i \in [0, k-1]$ **do**
- 4: **if** $gp_i \neq \emptyset$ **then**
- 5: $gp_i \leftarrow \text{sort}(gp_i)_{\text{based on } l(\cdot)}$
- 6: $u \leftarrow gp_i.\text{first}$
- 7: $AP_{v,i} \leftarrow \text{accessPattern}(gp_i)$
- 8: emit $\langle u; (v, AP_{v,i}, N_v^H) \rangle$
- 9: **end if**
- 10: **end for**

Reduce: Input: $[\langle u; (v, AP_{v,i}, N_v^H) \rangle]$

- 11: initiate the core node u , N_u^L, N_u^H in main memory
- 12: **for all** pivot message $\langle u; (v, AP_{v,i}, N_v^H) \rangle$ **do**
- 13: **for all** $w \in N_v^H \cap N_u^H$ **do**
- 14: emit Δ_{vuw}
- 15: **end for**
- 16: Put $(v, AP_{v,i}, N_v^H)$ into shared buffer
- 17: $N_u^L \leftarrow N_u^L - \{v\}$
- 18: **end for**
- 19: **for all** $r \in N_u^L$ **do**
- 20: Fetch $(r, AP_{r,i}, N_r^H)$ from shared buffer
- 21: **for all** $w \in N_r^H \cap N_u^H$ **do**
- 22: emit Δ_{ruw}
- 23: **end for**
- 24: **end for**

then, for the reduce-side computations over a core node u , the reducer compares u 's core message with all related pivot messages—some are associated with u 's core message, while the rest should be residing in the shared buffer. Bermuda-VC algorithm applies the same scheme to avoid generating redundant pivot messages. It utilizes a universal key partitioning function to group effective neighbors N_v^H of each node v . In order to guarantee the availability of the pivot messages, a universal key comparator function is utilized to sort the destination nodes in each group (Line 5, Algorithm 4). As a result, destination nodes are sorted based on their processing order. The first node in group gp_i indicates the earliest request of a pivot message. Hence, each node v sends a pivot

message to the first node of each non-empty group by emitting key value pairs where key equals the first node ID (Lines 6-8, Algorithm 4).

Combined with the sorting phase of the MapReduce framework, Bermuda-VC guarantees the availability of all needed pivot messages of any node u when u 's core message is received by a reducer, i.e., the needed pivot messages are either associated with u itself or associated with another nodes processed before u .

The reducers' processing mechanism is similar to that of the MR-Baseline algorithm. Each node u reads its core message for initiating N_u^H and N_v^L (Line 9), and then it iterates over every pivot message associated with key u against its effective adjacency list N_v^H to enumerate the triangles (Lines 10-12). As discussed before, not all expected pivot messages are carried with key u . The rest of the related pivot messages reside in the shared buffer. Here, N_v^L is used for fetching the rest of these pivot messages (Line 14, Algorithm 4), and enumerating the triangles (Lines 15-18, Algorithm 4). Moreover, the new coming pivot messages associated with node u are pushed into the shared buffer for further access by other nodes (Line 13). Figure 3.2 illustrates the reduce-side processing flow of Bermuda-VC. In the following sections, we will discuss in more details the management of the pivot messages in the shared buffer.

3.3.1 Message Sharing Management

Note that, each reduce instance uses its shared buffer to organize sharing pivot messages in memory. Message sharing via shared buffer is effective because the reduce computation program access the same pivot message over and over. By keeping as much of pivot message as possible in the shared buffer, it avoids access the slower medium-disk. In general, there are two types of operations over the shared buffer inside a reduce instance, which are: “*Put*” for adding new incoming pivot messages into the shared buffer (Line 13), and “*Get*” for retrieving the needed pivot messages (Lines 15-18). For massive graphs, the main memory may not hold all the pivot messages. This problem is similar to the classical caching problem studied in [12, 15], where a *reuse-distance* factor is used to estimate the distances between consecutive references of a given cached element, and based on that effective replacement policies can be deployed. We adopt the same idea in Bermuda-VC.

Interestingly, in addition to the reuse distance, all access patterns of each pivot message can be easily estimated in our context. The access pattern AP of a pivot message is defined as the sequence of graph nodes (keys) that will reference this message. In

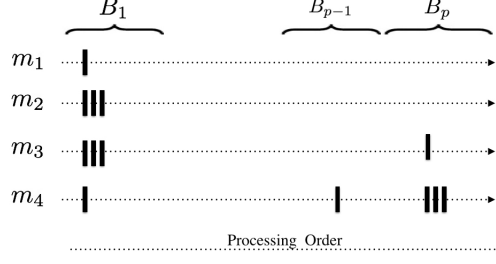


Figure 3.3: Various access Patterns for Pivot Messages.

particular, the access pattern of a pivot message from node v to reduce instance R_i can be computed based on the sorted effective nodes gp_i received by R_i . Several interesting metrics can be derived from this access pattern. For example, the first node in gp_i indicates the occurrence of the first reference, the size of gp_i equals the cumulative reference frequency. Such access pattern information is encoded within each pivot message (Lines 7-8, Algorithm 4). With the availability of this access pattern, effective message sharing strategies can be deployed under limited memory.

As an illustrative example, Figure 3.3 depicts different access patterns for four pivot messages $\{m_1, m_2, m_3, m_4\}$. The black bars indicate requests to the corresponding pivot message, while the gaps represent the re-use distances (which are idle periods for this message). Pivot messages may exhibit entirely different access patterns, e.g., pivot message m_1 is referenced only once, while others are utilized more than once, and some pivot messages are used in dense consecutive pattern in a short interval, e.g., m_2 and m_3 . Inspired by these observations, we propose two heuristic-based replacement policies, namely *usage-based tracking*, and *bucket-based tracking*. They trade off the tracking overhead with memory hits as will be described next.

•**Usage-Based Tracking** Given a pivot message originated from node v , the total use frequency is limited to \sqrt{m} , referring to the number of its effective neighbors, which is much smaller than the expected number of nodes processed in a single reducer, which is estimated to n/k . This implies that each pivot message may become useless (and can be discard) as a reducer progresses, and it is always desirable to detect the earliest time at which a pivot message can be discarded to maximize the memory’s utilization.

The main idea of the usage-based tracking is to use a usage counter per pivot message in the shared buffer. And then, the tracking is performed as follows. Each *Put* operation sets the counter as the total use frequency. And, only the pivot messages whose usage counter is larger than zero are added to the shared buffer. Each *Get* operation decre-

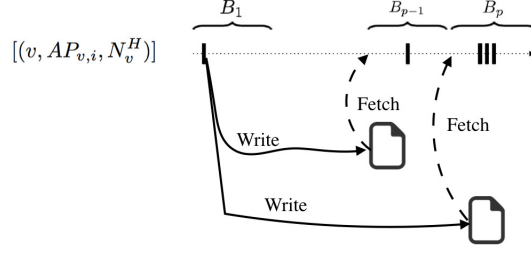


Figure 3.4: Read and Write Back-up files

ments the counter of the target pivot message by one. Once the counter reached zero, the corresponding pivot message is evicted from the shared buffer.

The usage-based scheme may fall short in optimizing sparse and scattered access patterns. For example, as shown in Figure 3.3, the reuse distance of message m_4 is large. Therefore, the usage-based tracking strategy has to keep m_4 in the shared buffer although it will not be referenced for a long time. What's worse, such scattered access is common in massive graphs. Therefore, pivot messages may unnecessarily overwhelm the available memory of each single reduce instance.

Bucket-Based Tracking We introduce the *Bucket-based* tracking strategy to optimize message sharing over scattered access patterns. The main idea is to manage the access patterns of each pivot message at a smaller granularity, called a bucket. The processing sequence of keys/nodes is sliced into buckets as illustrated in Figure 3.3. In this work, we use the range partitioning method for balancing workload among buckets. Correspondingly, the usage counter of one pivot message is defined per bucket, i.e., each message will have an array of usage counters of a length equals to the number of its buckets. For example, the usage count of m_4 in the first bucket is 1 while in the second bucket is 0. Therefore, for a pivot message that will remain idle (with no reference) for a long time, its counter array will have a long sequence of adjacent zeros. Such access pattern information can be computed in the map function, encoded in access pattern (Line 7 in Algorithm 4), and passed to the reduce side.

The corresponding modification of the *Put* operation is as follows. Each new pivot message will be pushed into the shared buffer (in memory) and backed up by local files (in disk) based on its access pattern. Figure 3.4 illustrates this procedure. For the arrival of a pivot message with the access pattern $[1, 0, \dots, 1, 3]$, the reduce instance actively adds this message into back-up files for buckets B_{p-1} (next-to-last bucket) and B_p (last bucket).

And then, at the end of each bucket processing and before the start of processing the next bucket, all pivot messages in the shared buffer are discarded, and a new set of pivot messages is fetched from the corresponding back-up file into memory (See Figure 3.4).

The Bucket-based tracking strategy provides better memory utilization since it prevents the long retention of unnecessary pivot messages. In addition, usage-based tracking can be applied to each bucket to combine both benefits, which is referred to as the *bucket-usage* tracking strategy.

3.3.2 Discussion

In this section, we show the benefits of the Bermuda-VC algorithm over the Bermuda-EC algorithm. Furthermore, we discuss the effect of parameter p , which is the number of buckets, on the performance.

Under the same settings of the number of reducers k , the Bermuda-VC algorithm generates more intermediate message and takes longer execution time. Firstly, the Bermuda-VC algorithm generates the same number of pivot messages while generating more core messages (*i.e.*, additional N_v^L for reference in the reduce side). Thus, the total size of the extra N_v^L core message is $\sum_{v \in V} N_v^L = m$. Such noticeable size of extra core messages requires additional time for generating and shuffling. Moreover, an additional computational overhead (Lines 13-14) is required for the message sharing management.

However, because of the proposed sharing strategies, the Bermuda-VC algorithm can work under smaller settings for k —which are settings under which the Bermuda-EC algorithm will probably fail. In this case, the benefits brought by having a smaller k will exceed the corresponding cost. In such cases, Bermuda-VC algorithm will outperform Bermuda-EC algorithm.

Moreover, compared to the disk-based Bermuda-EC algorithm, the Bermuda-VC algorithm has a relatively smaller disk I/O cost because the predictability of the access pattern of the pivot messages, which enable purging them early, while that is not applicable to the core messages. Notice that, for any given reduce instance, the expected usage count of pivot message from u is d_u^H/k . Thus, the expected usage count for any pivot message is $E(d_u^H)/k$, equals m/nk . Therefore, the total disk I/O with pivot messages is at most m^2/nk , smaller than disk I/O cost of Bermuda-EC algorithm m^2/Mk , where M stands for the size of the available memory in a single machine.

Effect of the number of buckets p : At a high level, p trades off the space used by the shared buffer with the I/O cost for writing and reading the back-up files. Bermuda-VC

algorithm favors smaller settings of p in the capacity of main-memory. As p decreases, the expected number of reading and writing decreases, however the total size of the pivot messages in the shared buffer may exceed the capacity of the main-memory. For a setting of p , the expected size of the pivot messages for any bucket is $O(m/kp)$. Therefore, a visible solution for $O(m/kp) \leq M$ is $\geq O(m/kM)$. Here, p is set as $O(m/kM)$ where m is the size of the input graph, and M is the size of the available memory in a single machine.

Chapter 4

Degree Anonymization over Uncertain Graphs

In this chapter, I will briefly review the degree anonymization techniques I developed for uncertain graphs. These techniques have been submitted as a conference paper [].

In this work, we study the novel problem of anonymization in the context of uncertain graphs. We extend the existing framework to work over uncertain graph by integrating uncertainty into anonymization process. In particular, we introduce a new *reliability-based* utility metric suitable for uncertain graphs, in contrast to the existing metrics which are all geared towards deterministic graphs. We present an efficient approach which achieves the desired level of anonymity at a slight cost of reliability by perturbing edge uncertainties judiciously. To this purpose, we develop two uncertainty-aware heuristics based on the uncertain graph theory, which is *reliability-oriented* edge selection and *anonymity-oriented* edge perturbing. We show that the incorporation of uncertainty is necessary and beneficial. We experimentally evaluate the proposed approach using different real-world datasets and study the behavior of the algorithms under the different heuristics. The results demonstrate the effectiveness of our approach.

4.1 Naive Approach: Anonymization via Representative Instance

A naive approach of anonymizing an uncertain graph is to first somehow transform it to a deterministic graph, then perform anonymization processing over the deterministic

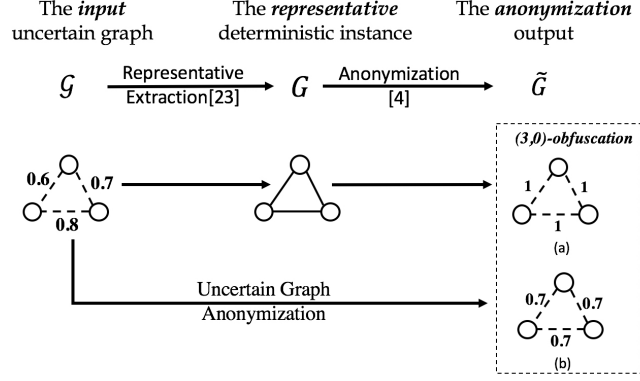


Figure 4.1: Illustration of anonymizing an uncertain graph through its representative deterministic instance and its drawback.

one. Fortunately, an increasing research effort was dedicated to the topic of extracting representative deterministic graphs from an uncertain graph [13]. Parchas *et al.* [13] ever introduced algorithms for extracting deterministic representative graph which captures key properties of the input uncertain graph. Now, it becomes realizable to anonymize an uncertain graph in two steps as shown in Figure 4.1. We first extract one deterministic representative instance G from the input uncertain graph \mathcal{G} . Then, we anonymize the extracted deterministic graph G , and output this result as the anonymized result of the original uncertain graph \mathcal{G} (referred as Rep-An).

The Rep-An approach is attracting since it does not require any new anonymization techniques specific designed for uncertain graphs. When the extracted representative deterministic graph G is close enough to the input uncertain graph \mathcal{G} in terms of graph properties, its anonymized result is expected to be a good anonymization of the input uncertain one. However, there is a non-negotiable difference between the input uncertain graph \mathcal{G} and its deterministic representative instance G , as exemplified in Figure 4.1. The anonymized result of G which is structurally similar to itself instead of the input uncertain graph, consequently, may be far different from the optimal solution, as exemplified in Figure 4.1. Therefore, we believe that for many applications, the Rep-An approach, introducing a high level of noise in such fashion, do reduce the overall graph utility. In experiment section, we will further illustrate this phenomenon over real-world datasets.

Algorithm 5 Anonymization over uncertain graphs

Input: Uncertain graph \mathcal{G} , adversary knowledge ak , obfuscation level k , tolerance level ϵ , size multiplier c and white noise level q

Output: The anonymized result $\tilde{\mathcal{G}}_f$

- 1: Initiation of a lower bound σ_l and an upper bound σ_u
 - 2: **while** (not terminated) **do**
 - 3: Search of (k, ϵ) -obf using standard deviation σ
 $\langle \tilde{\epsilon}, \tilde{\mathcal{G}} \rangle \leftarrow \text{GenerateObfuscation}(\mathcal{G}, \sigma_u, ak, k, \epsilon, c, q)$
 - 4: Update the anonymized result if necessary
 - 5: Update the lower bound or the upper bound
 - 6: Update the size multiplier c if necessary
 - 7: **end while**
 - 8: return the anonymized result $\tilde{\mathcal{G}}_f$
-

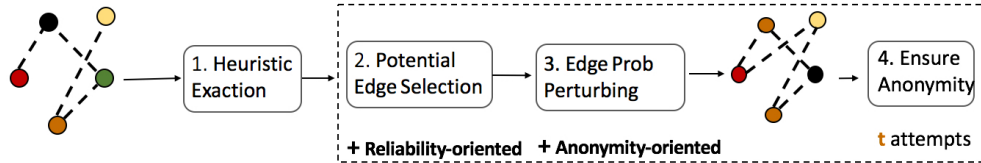


Figure 4.2: The pipeline of function `GenerateObfuscation`

4.2 Reliability-Preserving anonymization on uncertain graphs

In this section, we first give a brief review of the state-of-art framework which anonymizes deterministic graphs by injecting uncertainty to certain edges. Then, we discuss how to extend this framework to work over uncertain graphs by explicitly incorporating edge uncertainty.

4.2.1 Anonymization Procedure

Boldi *et al.* proposed a seminal approach of anonymizing *deterministic graphs* which injects uncertainty in the existence of the edges of the graph and publishes the resulting *uncertain graph*. Their method injects uncertainty to edges in a deterministic graph as follows: for each existing sampled edge e , it is assigned a probability $1 - r_e$, and for each non-existing sampled edge e , it is assigned a probability $r_e \leftarrow R_{\sigma_e}$. By this way, it converts the input graph into an uncertain one. In particular, the perturbation variable r_e

is generated by a truncated $[0, 1]$ normal distribution, $r_e \leftarrow R_{\sigma_e}$;

$$R_{\sigma}(r) := \begin{cases} \frac{\Phi_{0,\sigma}(r)}{\int_0^1 \Phi_{0,\sigma}(x)dx} & r \in [0, 1] \\ 0 & \text{otherwise} \end{cases}$$

where $\Phi_{0,\sigma}$ is the density function of a Gaussian distribution. As the standard deviation σ of the normal distribution decreases, a greater mass of R_{σ} will concentrate near $r = 0$ and the amount of injected uncertainty will be smaller. Namely, smaller values of σ contribute towards better maintaining the characteristics of the original graph. Targeting at high utility, they formulated the graph anonymization problem as the minimization of σ and computed the minimal amount of uncertainty via a binary search on the value of the uncertainty parameter σ . The overall anonymization procedure is illustrated in Algorithm 5.

The search flow of Algorithm 5 is determined by the function `GenerateObfuscation`, whose pipeline is shown as Figure 4.2. The function `GenerateObfuscation` aims to find a (k, ϵ) -obfuscation of the original *graph* using a given standard deviation parameter σ . A key feature of this function is the utilization of heuristics for judicious edge selection and edge prob perturbing. The function first computes heuristics values based on node properties in the input graph. After that, the search of (k, ϵ) -*obf* instance starts in a randomized way: t attempts are performed. Each attempt performs the following steps:

- Selecting a subset of edges E_c
- Perturbing existence probabilities of selected edges
- Checking the resulting uncertain graph whether satisfies the anonymity requirement

If the algorithm finds a (k, ϵ) -obfuscated graph in one of its t trials, it returns the obfuscated graph with minimal ϵ . Otherwise, the algorithm indicates the failure by returning $\tilde{\epsilon} = 1$.

In this work, we extend this broad framework by incorporating its key function `GenerateObfuscation` with uncertainty. The specific method proposed in [6], has two drawbacks for anonymizing uncertain graphs. First, their method does not consider the structural relevance of edges in the critical edge selection step, which leads to unnecessary structural distortion. Second, its interior edge perturbing mechanism assumes the existence of edges is known with certainty, thus fails to handle uncertain graphs where the

existence of edges is probabilistic. We will present corresponding solutions (reliability-oriented edge selection and anonymity-oriented edge perturbing).

4.2.2 Reliability-oriented Edge Selection

The most important step is the selection of potential edges, which impacts further anonymization effort significantly. Finding the optimal set of edges E_c that balances privacy and utility is a typical combinational optimization problem. The problem is computationally expensive since there are exponential many combinations to be considered.

To alleviate the combinational intractability, kind of heuristics have been utilized in the context of deterministic graphs. These heuristics can be classified into two main categories (1) Anonymity-oriented ones that suggest implementing larger perturbation to less anonymized nodes [20] (2) Utility-oriented ones that suggest implementing smaller perturbation to more influential edges/nodes [8, 21].

In this work, we present a sampling-based approach which combines both heuristics together. First, we adopt the idea of *uniqueness score* for quantifying the anonymity level of nodes. Second, we introduce a novel edge relevance metric, *reliability relevance* (RR), for quantifying the impact of edge modification to the overall uncertain graph. In contrast to the existing metrics such as edge betweenness [8], which are all defined in deterministic graphs, reliability relevance geared towards uncertain graphs. It allows us to select a subset of edges subjected to perturbation with less impact on reliability (reliability-oriented edge selection).

Uniqueness Score

Uniqueness score is proposed to measure how typical the node is among all the nodes in terms of its property value [6]. More formally, the uniqueness score is defined as follows.

Definition 5 Uniqueness Score [6] Let $P : V \rightarrow \Omega_P$ be a property on the set of nodes V of the graph \mathcal{G} , let d be a distance function on Ω_P , and let $\theta > 0$ be a parameter. Then the θ – commonness of the property values $\omega \in \Omega_P$ is $C_\theta(\omega) := \sum_{u \in V} \Phi_{0,\theta}(d(\omega, P(v)))$, while the θ -uniqueness of $\omega \in \Omega_P$ is $U_\theta := \frac{1}{C_\theta(\omega)}$.

Note that, the commonness of the property value ω captures how typical the value ω is among all the nodes in the graph via a weighted average function, where the weight decays exponentially as the distance d . As above defined, uniqueness scores depend on the parameter θ , which determines the decay rate of average weights. In this work, we set

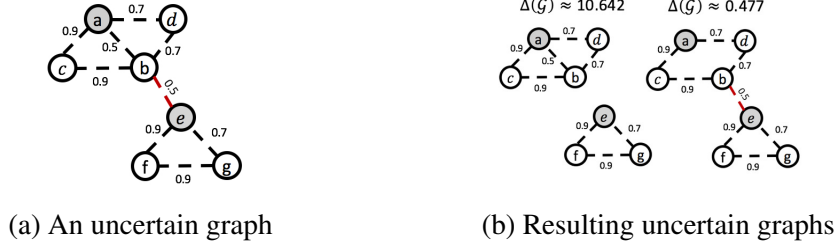


Figure 4.3: Illustration of varying structural distortions by removing different uncertain edges (b, e) and (a, b) over an uncertain graph.

$\theta = \sigma_{\mathcal{G}}$, where $\sigma_{\mathcal{G}}$ represents how frequently the property value spread in the uncertain graph \mathcal{G} . A comprehensive discussion can be found in the literature [6].

Reliability Relevance

As exemplified in Figure 4.3, each edge modification will have *varying* impact to the graph structure in the context of uncertain graphs. From the perspective of *uniqueness*, nodes a and e are identical since adjacent edges have identical uncertainties. Accordingly, anonymity-oriented heuristics would select and perturb edges (a, b) and (b, e) without bias. As shown in Figure 4.3, the deletion of the uncertain edge (b, e) , the only link connecting two *reliable* clusters, clearly incurs much larger structure distortion than the deletion of the edge (a, b) . This observation calls for an effective measure of edge influence in the context of uncertain graphs.

Inspired by the significance of *reliability*, we provide a compact analytic form of reliability deviation caused by individual uncertain edge modification in a fine-grained way, namely changing its edge probability in any granularity. Following the path, we introduce edge reliability relevance for edges' influence estimation, which can be quantified by the sum of reliability discrepancy of all the node pairs. Besides, we provide an algorithm for accessing edge *reliability relevance* (RR) efficiently.

Reliability Relevance Definition Lets start with analyzing the impact of single uncertain edge alteration to the connectivity of a specific node pair (Two-terminal Reliability).

Definition 6 Two-terminal Reliability Relevance Given an uncertain graph \mathcal{G} , and two nodes u, v , we consider the probability that u is reachable to v , $R_{u,v}(\mathcal{G})$, as defined in Definition 2, as a multivariate function involving all the edge probabilities. Thus, the partial derivative of $R_{u,v}(\mathcal{G})$ with respect to an individual variable $p(e)$, existence probability

of an uncertain edge, denotes as $RR_{u,v}(e)$. It represents the sensitivity to the change of reliability $R_{u,v}$ which is determined by $p(e)$ with the others held constant. Its definitions is as follows:

$$RR_{u,v}(e) = \frac{\partial R_{u,v}(\mathcal{G})}{\partial p(e)}$$

Lemma 1 Factorization Lemma Given an uncertain graph \mathcal{G} , the reliability of the node pair (u, v) $R_{u,v}(\mathcal{G})$ can factorized via one uncertain edge e as follows:

$$R_{u,v}(\mathcal{G}) = p(e)R_{u,v}(\mathcal{G}_e) + (1 - p(e))R_{u,v}(\mathcal{G}_{\bar{e}})$$

where uncertain graph \mathcal{G}_e is identical with \mathcal{G} except its edge e is an **certain** edge. Similarly, the uncertain graph $\mathcal{G}_{\bar{e}}$ is identical with \mathcal{G} except its edge e is an **certain non** edge.

According to the factorization lemma, the partial derivative $RR_{u,v}(e)$ can be rewritten as

$$RR_{u,v}(e) = R_{u,v}(\mathcal{G}_e) - R_{u,v}(\mathcal{G}_{\bar{e}})$$

The equation indicates the sensitivity of reliability is constant and it does not depend on the probability value of the uncertain edge e . Another important point to remember is that because all the connected pairs remain connected after the addition of an edge, $R_{u,v}(\mathcal{G}_e) - R_{u,v}(\mathcal{G}_{\bar{e}}) \geq 0$.

For a given uncertain edge e , the derivatives $RR_{u,v}(e)$ can be arranged in a $|V| \times |V|$ matrix, where each element corresponds to one node pair (u, v) and it gives the partial derivative for the corresponding reliability $R_{u,v}$. As ever discussed, all the element are equal to or greater than zero. In this work, we define $RR(e)$ to be reliability relevance of an edge e which can be quantified by the sum of all the $RR_{u,v}(e)$, as

$$\begin{aligned} RR(e) &= \sum_{u,v} |RR_{u,v}(e)| \\ &= \sum_{u,v} |R_{u,v}(\mathcal{G}_e) - R_{u,v}(\mathcal{G}_{\bar{e}})| \\ &= \sum_{u,v} R_{u,v}(\mathcal{G}_e) - \sum_{u,v} R_{u,v}(\mathcal{G}_{\bar{e}}) \end{aligned}$$

Note that, $RR(e)$ equals the difference of the expected number of connected pairs between uncertain graphs \mathcal{G}_e and $\mathcal{G}_{\bar{e}}$ by explicit incorporation of edge uncertainty. In the context of edge relevance, reliability relevance can be seen as generalization of cut-edges, which quantifies the impact of partial edge deletion or addition on the connectivity in the uncertain graph. The higher reliability relevance score of an edge, the bigger impact of

Algorithm 6 Reliability Relevance Evaluation

Input: $\mathcal{G} = (V, E, p)$, K is the number of sampled graphs; usually $K = 1000$

Output: RR Reliability relevance of edges in \mathcal{G}

```
1:  $CC_e \leftarrow 0, CC_{\bar{e}} \leftarrow 0$ 
2: for  $i=1$  to  $K$  do
3:    $G \leftarrow$  A deterministic sampled instance
4:    $Ind(G)$  is edge existence of sampled graph  $\mathcal{G}$ 
5:    $cc(G) \leftarrow$  the number of connected pairs of  $G$ 
6:    $CC_{e+} = Ind(G) \cdot cc(G), CC_{\bar{e}+} = (1 - Ind(G)) \cdot cc(G)$ 
7: end for
8:  $RR = CC_e/p - CC_{\bar{e}}/1 - p$ 
```

edge perturbation over the overall graph. On this basis, we define the reliability centrality of node $v, RC(v)$ to be the overall influence of graph reliability which can be quantified by the weighted sum of reliability relevance of all adjacent edges.

$$RC(u) = \sum_e p(e)RR(e)$$

Reliability Relevance Evaluation Note that, the evaluation of $RR(e)$ need computing all the $R_{u,v}$ over two uncertain graphs \mathcal{G}_e and $\mathcal{G}_{\bar{e}}$. The two-terminal reliability detection problem is a #P-complete problem [3]. Thus, the exact computation is not practical for large uncertain graphs. A basic approach is to get reliability approximation by sampling: for each uncertain edge e 1) we first sample K possible graphs G_1, \dots, G_K of \mathcal{G}_e according to edge probability p , and 2) we then compute the number difference of connected pairs in each sample graph after edge e is added. The basic sampling method can be rather computationally expensive. The total running time is $\Theta(|E| \cdot K \cdot \alpha(|V|)|E|)$ assuming incremental algorithm is used for computing connected components. It is impractical for large uncertain graphs with huge size of edges.

Here, we introduce a much faster approach for computing $RR(e)$ as shown in Algorithm 6. The basic idea is to partition all sampled graphs into groups according to existing edges so that the evaluation of the number of connected pairs can be reused. The complexity of reliability relevance evaluation of all edges is $\Theta(K \cdot \alpha(|V|)|E|)$.

Reliability-oriented Edge Selection Procedure

Algorithm 7 GenerateObfuscation

Input: Uncertain graph $\mathcal{G} = (V, E, p)$, ak, k, ϵ, c, q ,
and standard deviation σ

Output: A pair $\langle \tilde{\epsilon}, \tilde{\mathcal{G}} \rangle$

```
1: compute the uniqueness  $U(v)$  for all the nodes
2: compute the centrality  $RC(v)$  for all the nodes
3:  $Q(v) \leftarrow U(v) \cdot RC(v)$ 
4:  $H \leftarrow$  the set of  $\lceil \frac{\epsilon}{2} |V| \rceil$  with largest  $Q(v)$  {Excluding}
5: Normalized  $RC(v)$ ;  $Q(v) \leftarrow U(v) \cdot 1 - RC(v)$ ;
6:  $\tilde{\epsilon} \leftarrow 1$ 
7: for  $t$  times do
8:    $E_C \leftarrow E$  {Reliability-oriented Edge Selection}
9:   repeat
10:    randomly pick a vertex  $u \in V \setminus H$  according to  $Q$ 
11:    randomly pick a vertex  $v \in V \setminus H$  according to  $Q$ 
12:    draw  $w$  uniformly at random from  $[0, 1]$ 
13:    if  $(u, v) \in E$  then
14:      if  $w > p(e)$  then
15:         $E_C \leftarrow E_C \setminus \{(u, v)\}$ 
16:      end if
17:    else
18:       $E_C \leftarrow E_C \cup \{(u, v)\}$ 
19:    end if
20:  until  $E_C = c|E|$ 
21:  for all  $e \in E_C$  do
22:    compute  $\sigma(e)$  {Edge Probability Perturbation}
23:    draw  $w$  uniformly at random from  $[0, 1]$ 
24:    if  $w < q$  then
25:       $r_e \leftarrow U(0, 1)$ 
26:    else
27:       $r_e \leftarrow R_{\sigma(e)}$ 
28:    end if
29:     $\tilde{p}(e) \leftarrow p(e) + 2(0.5 - p(e)) \cdot r_e$ 
30:  end for
31:   $\hat{\epsilon} \leftarrow checkAnonymity(\tilde{\mathcal{G}})$  {Ensure Anonymity}
32:  Update result  $\langle \tilde{\epsilon}, \tilde{\mathcal{G}} \rangle$  if  $\hat{\epsilon} < \tilde{\epsilon}$ 
33: end for
34: return  $\langle \tilde{\epsilon}, \tilde{\mathcal{G}} \rangle$ 
```

Now, we are ready to introduce our function `GenerateObfuscation` aims at finding a (k, ϵ) -obf instance for an input *uncertain graph* \mathcal{G} using a given standard deviation

parameter σ as shown in Algorithm 7. It resembles the method proposed in [6] with explicit incorporation of edge uncertainty in edge selecting and perturbing step. Another important contribution of our work is the utilization of reliability relevance (RR) for capturing structural relevance of edges in the uncertain graph.

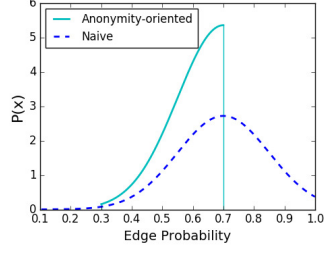
First, it computes the uniqueness level and reliability centrality for each node $v \in \mathcal{G}$. Intuitively, the more unique a node is, the harder it is to be obfuscated; the more *important* an edge is, the bigger utility loss it incurs. Such heuristic information is important for our privacy-preserving and utility-preserving purpose. In order to use the “uncertainty budget” in the most efficient way, the algorithm performs the following steps as shown in Algorithm 7.

(Line 4: Excluding) Since it is allowed not to obfuscate $\epsilon|V|$ of the nodes, the algorithm selects the set H of $\frac{\epsilon}{2}|V|$ nodes with largest uniqueness scores or reliability centrality scores, and exclude them from the subsequent obfuscation efforts. In later steps, the algorithm will inject uncertainty only to edges that are not adjacent to any of vertices in H . The key feature of our method is that it strictly preserve the 1-neighborhoods of influential nodes.

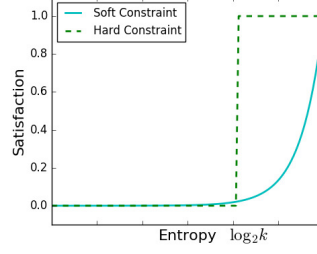
(Line 8-18: Potential Edge Selection) The set of nodes not in H will need to be anonymized. To anonymize more unique vertices, higher uncertainty is necessary. Thus, edges needed to be sampled with higher probability if they are adjacent to unique nodes. In order to better preserve graph structure, edges needed to be sampled with smaller probability if they are adjacent to more influential nodes. In order to handle this sample process, our algorithm assigns a probability $Q(v)$ to every $v \in V$, which depends on the uniqueness level $U(v)$ and reliability centrality $RC(v)$ of v .

Each attempt begins by selecting a subset of edges E_c , which will be subjected to edge probability perturbation. Then set E_c , whose target size $\|E_c\| = c \|E\|$, is initialized to be E . Then, the algorithm randomly selects two distinct vertices u and v , according to probability distribution Q . The pair of vertices (u, v) is removed from E_c with the probability $p(e)$ if it is an edge in the original graph, or added to E_c otherwise. The process is repeated until E_c reaches the require size. In typical uncertain graphs, the number of non-edges is significantly larger than the number of uncertain edges, the loop ends very quickly, for small values of c . And, the resulting set E_c includes most of edges in E .

(Line 20) Next, we redistribute the perturbation budget among all selected edges $e \in E_c$ in proportion to their intermediate representations $Q(v)$. Specially, we define for each



(a) Anonymity-oriented Edge Perturbing



(b) Relaxing k -obfuscation constraint

$e = (u, v) \in E_c$, its uncertainty level,

$$Q(e) := \frac{Q(u) + Q(v)}{2}$$

and then set its edge perturbation parameter

$$\sigma(e) = \sigma|E_c| \cdot \frac{Q(e)}{\sum_{e \in E_c} Q(e)}$$

so that the average of $\sigma(e)$ over all $e \in E_C$ equals σ .

4.2.3 Anonymity-oriented Edge Perturbing

As ever discussed, the existing uncertainty injecting scheme is designed for deterministic graphs and can not be used to handle uncertain graphs directly. Given the uncertainty level $\sigma(e)$, a natural strategy is to consider the partial edge addition and deletion in a random way as shown in Figure 4.4(a). Through the analysis of the impact of a single edge probability alteration, we give a anonymity-oriented perturbing heuristics which is able to constraint the potential range of edge prob perturbing (referred to C). For each selected edge $e \in E_c$ and the assigned uncertainty level $\sigma(e)$, we alter its existence probability as

$$\tilde{p}(e) := p(e) + (1 - 2p(e)) \cdot r_e$$

Where the random perturbation r_e is generated according to $\sigma(e)$ as Eq. 4.2.1. Namely, for an edge with the probability p , we only consider potential edge probability \tilde{p} in the limited range that more likely contributes to higher graph anonymity. Clearly, previous scheme defined in deterministic graph becomes a special case of our approach. We proceed to elaborate the rationality and the benefit of this anonymity-oriented edge perturbing scheme by treating it as a constraint satisfaction problem.

Basics

Let us consider k -obfuscate a given node v as a single constraint c_v for the target anonymization graph. Then, according to the definition 1, the satisfaction of the constraint c_v is defined as

$$c_v := \begin{cases} 1 & H(Y_{P(v)}) \geq \log_2 k \\ 0 & \text{otherwise} \end{cases}$$

Namely, given the considered degree-based re-identification scenario, we lower bound the entropy of degree distribution over the anonymized graph by $\log_2 k$. Accordingly, whether the anonymized graph k -obfuscates all the nodes can be expressed as degree of joint satisfaction.

$$\begin{aligned} \mathcal{C}(\mathcal{G}) &= \prod_{v \in V} c_v \\ &= \prod_{\omega} \underbrace{c \dots c}_{s(\omega)} \end{aligned}$$

The uncertain graph is said to k -obfuscate all the nodes if and only if the $\mathcal{C}(\mathcal{G})$ equals 1.

Re-visiting graph anonymity

As shown in Figure 4.4(b), the original satisfaction function is not everywhere differentiable. To simplify the anonymity analysis, we approximate the individual constraint c_v to a fuzzy relation in which the satisfaction degree of a constraint is defined as a continuous and differentiable function, going from fully violated to fully satisfied. A natural candidate for soft satisfaction function is,

$$C_v = e^{H(Y_{P(v)}) - \log_2 |V|}$$

According to this approximation, the satisfaction function of k -obfuscate all the nodes can be rewritten as

$$\mathcal{C}(\mathcal{G}) \propto \prod_{\omega} \underbrace{e^{H(Y(\omega))} \dots e^{H(Y(\omega))}}_{s(\omega)}$$

While the original satisfaction is a binary value, the approximated satisfaction value lies in the continuous range $[0, 1]$. Note that, the higher satisfaction score indicates the higher level of anonymity achieved. Taking logarithm for both side, we get the concise formula as

$$\log \mathcal{C}(\mathcal{G}) = \sum_{\omega} s(\omega) \cdot H(Y(\omega))$$

Regarding the property P , it equals the weighted sum of entropy over all possible values ω , where $s(\omega)$ is the expected number of vertices with property value ω over all possible worlds. Targeting at high anonymity, we wish to increase the weighted sum of entropy.

Greedy search of graph anonymity

The remaining issue is to connect the single edge probability alteration with the objective. With respects to node degree, a graph can be represented as one matrix as shown in figure 4.4(b). The weighted sum of entropy is related with graph coding. Here, we utilize entropy encoding, especially Huffman coding. From the coding perspective, we have two different angles to perform graph coding: row or column of degree matrix. The encoded length should be equal to each other, as shown in the following equation.

$$\sum_v H(v) + n \log n = \sum_{\omega} s(\omega) \cdot H(Y(\omega)) + H(s(\omega))$$

It indicates that higher global anonymity of a uncertain graph can be achieved by increasing the entropy of individual nodes when the global distribution $H(\omega)$ keeps constant.

Note that, the probability change over a specific edge only affects degree distributions of its connected nodes. To further simply the problem, we assume that the impact of edge perturbing is independent to each other. Recall that, we suggest implementing perturbation to less anonymized nodes. In the case of degree obfuscation, they are nodes with high degree. For a node v with high degree, the probability distribution of its degree d_v may be approximated as the normal distribution as implied by the Central Limit Theorem. Hence, its induced entropy $H(v)$ can be approximated as $\frac{1}{2} \ln(2\pi e \sigma^2)$.

Targeting at maximizing the graph anonymity, we take steps proportion to the positive gradient with respect to the choice of individual edge probabilities.

$$\frac{\partial \sum_{\omega} s(\omega) H(Y(\omega))}{\partial p(e)} \propto (1 - 2p(e))$$

Therefore,

$$\tilde{p}(e) := p(e) + (1 - 2p(e)) \cdot r_e$$

Namely, our approach simulates one iteration of batch gradient ascent method for finding the local maximum of weighted entropy sum or the anonymity level.

Bibliography

- [1] E. Adar and C. Re. Managing uncertainty in social networks. *IEEE Data Eng. Bull.*, 2007.
- [2] S. Arifuzzaman, M. Khan, and M. Marathe. Patric: A parallel algorithm for counting triangles in massive networks. *CIKM*, pages 529–538, 2013.
- [3] M. O. Ball. Computational complexity of network reliability analysis: An overview. *IEEE Transactions on Reliability*, 1986.
- [4] N. Bao and T. Suzumura. Towards highly scalable pregel-based graph processing platform with x10. *WWW*, pages 501–508, 2013.
- [5] L. Becchetti, P. Boldi, C. Castillo, and A. Gionis. Efficient semi-streaming algorithms for local triangle counting in massive graphs. *KDD*, pages 16–24, 2008.
- [6] P. Boldi, F. Bonchi, A. Gionis, and T. Tassa. Injecting uncertainty in graphs for identity obfuscation. *SIGMOD*, 2012.
- [7] F. Bonchi, A. Gionis, and T. Tassa. Identity obfuscation in graphs through the information theoretic lens. *ICDE*, 2014.
- [8] J. Casas-Roma. Privacy-preserving on graphs using randomization and edge-relevance. *Modeling Decisions for Artificial Intelligence*, 2015.
- [9] Colbourn and Colbourn. The combinatorics of network reliability. 1987.
- [10] J. Dean and S. Ghemawat. Mapreduce: Simplified data processing on large clusters. *Commun. ACM*, 2008.
- [11] D. Kempe, J. Kleinberg, and É. Tardos. Maximizing the spread of influence through a social network. 2003.

- [12] G. Keramidas and P. Petoumenos. Cache replacement based on reuse-distance prediction. *ICCD*, 2007.
- [13] Parchas, Gullo, Papadias, and Bonchi. The pursuit of a good possible world: extracting representative instances of uncertain graphs. *SIGMOD*, 2014.
- [14] H. Park, F. Silvestri, U. Kang, and R. Pagh. MapReduce triangle enumeration with guarantees. *CIKM*, pages 1739–1748, 2014.
- [15] P. Petoumenos and G. Keramidas. Instruction-based reuse-distance prediction for effective cache management. *MSP*, pages 60–68, 2009.
- [16] M. Potamias, F. Bonchi, A. Gionis, and G. Kollios. K-nearest neighbors in uncertain graphs. *VLDB*, 2010.
- [17] T. Schank. Algorithmic aspects of triangle-based network analysis. *Phd in computer science*, 2007.
- [18] S. Suri and S. Vassilvitskii. Counting triangles and the curse of the last reducer. *KDD*, pages 607–614, 2011.
- [19] T. White. Hadoop: The definitive guide. 2010.
- [20] X. Ying, K. Pan, X. Wu, and L. Guo. Comparisons of randomization and k-degree anonymization schemes for privacy preserving social network publishing. *SNA-KDD*, 2009.
- [21] X. Ying and X. Wu. Randomizing social networks: a spectrum preserving approach. pages 739–750, 2008.