

# Path-preserving Anonymization for Inter-domain Routing Policies

Xiaozhe Shao, Hossein Pishro-Nik, and Lixin Gao

Department of Electrical and Computer Engineering, University of Massachusetts Amherst  
{xiaozechao, pishro, lgao}@engin.umass.edu

**Abstract.** The Internet consists of tens of thousands of autonomous systems (ASes), which exchange routing information through policy-based routing protocols, such as Border Gateway Protocol (BGP). BGP allows network administrators/operators to configure routing policies independently. The routing result is derived from complex interactions controlled by the configured policies. Thus, the routing policies of all networks are essential to understand which network a packet uses to reach a destination or whether there is a path to reach a destination. However, the routing policies of a network are proprietary information. The majority of networks in the Internet do not publicly share their routing policies. In this paper, we anonymize the routing policy with the aim to preserve the utility to understand the routing system. To this end, we perturb both the routing policies and the AS-level topology. We propose a series of anonymization schemes that perturb the AS-level topology to derive a  $k$ -anonymity graph. We further perturb routing policies to preserve the best path. We evaluate the anonymization schemes against sophisticated de-anonymization attacks exploiting rich structural information of AS-level topologies, such as neighborhood and reference distance. The experimental result shows that the de-anonymization attack exploiting rich structural information can identify less than 1% of nodes.

**Keywords:** Graph Anonymization · K-anonymity · Network Verification

## 1 Introduction

The Internet consists of tens of thousands of autonomous systems (ASes), each of which belongs to an organization such as an Internet Service Provider (ISP), a university, a company, or an Internet Exchange Point (IXP). The interconnections in the Internet have a hierarchical structure, where tier-1 ISPs provide settlement-free services among each other, and provide transit services to regional ISPs or stub networks, and regional ISPs provide transit services to even smaller ISPs or stub networks. BGP is a policy-based protocol used for routing among ASes.

The route taken by a packet to traverse in the Internet is determined by the route selection process, and a result of complex interactions controlled by the configured routing policies of all networks. Knowing routing policies of all networks allow us to answer a broad set of what-if questions [1, 2]. For example, if a network wants the inbound traffic from Google to go through a specific provider, then before changing its routing configuration the operator wants to verify that the configuration indeed leads to the desired outcome.

Routing policies are considered to be proprietary by their networks. Each network administrator/operator configures its routing policy independently. The routing policy of a network reveals how the network exchanges routes with neighbors. Thus, the commercial agreement (*e.g.*, provider-customer relationships) between neighboring ASes can be inferred through their routing policies. Further, the routing policy of a network reveals that which neighbor is more preferred than others. Therefore, routing policies should be anonymized before sharing.

A trivial method of anonymizing routing policies is to anonymize IDs of all networks, *e.g.* AS Numbers, in the routing policies. Namely, each original AS Number is transformed into a unique ID. However, anonymizing IDs alone is not sufficient. The number of neighbors of a network can be inferred from the anonymized routing policy of the network, even if all network IDs contained in the routing policy are anonymized.

In the Internet, tier-1 ISPs have a large number of neighbors (customers and peers) and the number of neighbors follow power-law distribution [3]. The ASes with a large number of neighbors are few and their degrees (*e.g.* the number of neighbors) are usually unique. The AS-level topology can be derived from the publicly available routing tables of RouteViews [4] and RIPE RIS [5]. The degrees in the AS-level topology can be used to identify the AS. Even if the derived AS-level topology does not exactly reflect the AS-level topology in reality, the ranking of these large ISPs in terms of degree will be same. Thus, these Tier-1 ISPs can still be identified. Identifying a few tier-1 networks can lead to identifying tier-2 networks by what tier-1 neighbors they have. Progressively, more and more ASes might be identified when their higher-tier neighbors have been identified. Therefore, anonymizing the AS-level topology is the key in anonymizing routing policies.

Anonymizing graphs have been studied at the area of online social networks [6–14]. These schemes perturb the graph through adding and/or removing nodes and/or edges in social networks, so that the adversaries can not identify users from the social network. These schemes focus on preserving social network structural properties. However, these schemes do not maintain the best paths derived from the routing policies. The applications exploiting the inter-domain routing policy are very sensitive to the best paths derived from the routing policy. Therefore, these schemes can not be directly applied to the anonymization of routing policies.

In this paper, we propose a path-preserving anonymization scheme for sharing inter-domain routing policies. We anonymize the routing policies so that the best paths can be maintained after anonymization. To do that, we first perturb AS-level topology by adding fake networks and connection between fake networks and between fake and real networks. We then establish routing policies for these newly-added links.

To perturb the AS-level topology, one anonymization scheme is adding neighborhood relationships between pairs of networks. However, simply adding fake edges between networks is not sufficient to maintain the best path utility. The reason is that the local preferences on the fake edge between two networks can not be properly assigned. When network  $u$  is connected to network  $v$  through a fake edge, to avoid the fake edge being included in the best paths, all real edges should be preferred over the fake edge. Thus, network  $u$  and network  $v$  should give each other the lowest local preference. In practice, a network assigns the lowest local preference to its provider and a pair of neighboring ASes can not be the provider of each other. Therefore, the fake edges between a pair of ASes, such as network  $u$  and network  $v$ , can be identified by using the local preferences assigned to the fake edge.

We propose to add fake networks into the AS-level topology and then add the connection between fake networks and real networks to perturb the AS-level topology. In the topology perturbation, we first consider to avoid being identified through the number of neighbors. We propose a deterministic anonymization scheme which derives a  $k$ -anonymity graph. As a result, the probability of identifying a network by the degree information is at most  $\frac{1}{k}$ . However, the resulting  $k$ -anonymity graph is deterministic. We further propose the probabilistic anonymization scheme to generate the  $k$ -anonymity graph with higher randomness. To anonymize the AS-level topology, the probabilistic anonymization scheme might significantly enlarge the graph with fake networks. Adding too many fake networks might make the scheme infeasible in practice. Finally, we propose the hybrid anonymization scheme to take advantage of the strengths of both the deterministic anonymization scheme and the probabilistic anonymization scheme.

After fake networks and fake edges are added, we propose to generate the associated routing policies for those fake networks and edges. To avoid de-anonymization, in the anonymized routing policies, the local preferences assigned to real neighboring ASes and fake neighboring ASes have the same distribution.

We evaluate the proposed anonymization schemes against de-anonymization attacks exploiting rich structural information, such as neighborhood and reference distance. The experimental result shows that the structure-based de-anonymization attacks can identify less than 1% networks from the anonymized AS-level topology. In addition, we use Closeness Centrality as an example metrics to illustrate how hard to identify a network through the structural information. We also show that the forged routing policies of networks added by the anonymization schemes will not be used to identify the added networks from the original networks.

The rest of this paper is organized as follow: In Section 2, we illustrate background and formally define the problem. In Section 3, we discuss the preliminaries necessary for preserving utility to construct the anonymization schemes. To perturb the AS-level topology, in Section 4 we propose the deterministic anonymization scheme, the probabilistic anonymization scheme and the hybrid anonymization scheme respectively. In Section 5, we propose to perturb the routing policies through a local preference assignment scheme. We describe the experiment in Section 6, related work in Section 7, and make the conclusion of the paper in Section 8.

## 2 Background and Problem Formulation

In this section, we formulate the problem to be addressed in the paper.

### 2.1 Routing Policies and AS-level Topology

BGP is the de-facto standard for inter-domain routing in the Internet. BGP allows network operators to configure the routing policy of their own networks. The routing policies of those networks together determine the best routes from source to destination networks. In BGP, a network assigns a local preference to each of its neighbors respectively and prefers the routes from a neighbor network with a larger local preference.

Listing 1.1: A snippet of the routing policy

---

```

router bgp 45000
  neighbor 192.168.1.2 remote-as 40000
  neighbor 192.168.1.2 route-map SET-LOCAL-PREF-CUST in
  neighbor 192.168.1.3 remote-as 50000
  neighbor 192.168.1.3 route-map SET-LOCAL-PREF-PEER in
!
route-map SET-LOCAL-PREF-CUST
  set local-preference 200
route-map SET-LOCAL-PREF-PEER
  set local-preference 100

```

---

In Listing 1.1, we illustrate a snippet of the routing policy for AS 45000. The routing policy of AS 45000 indicates the neighborhood relationships of AS 45000. According to the routing policy, AS 45000 has two neighboring ASes, AS 40000 and AS 50000. This snippet also indicates that AS 45000 sets up the local preference 200 for its customer, AS 40000, and the local preference 100 for its peer, AS 50000. As a result, AS 45000 prefers routes from AS 40000 over that from AS 50000.

The routing policy of a network indicates the neighboring ASes of the network. The routing policies of all networks together implies the AS-level topology of the Internet. Let us denote the AS-level topology by a graph,  $G = \{V, E, R\}$ , where  $V$  is the set of nodes,  $E$  is the set of edges and  $R$  is the set of all routing policies. Each node represents an AS in the Internet. The edge between two nodes represents the link between two networks.

## 2.2 Path-preserving Anonymization

To protect the routing policies, we can obfuscate the key information in the routing policies. One of the most trivial one is pseudonymization. We can replace AS numbers in the routing policies by a pseudo ID, so that the routing policies can not be identified through the AS numbers directly. Let us denote the pseudonymization as a network mapping function,  $\pi$ , that maps the original AS number to the anonymized AS number (the pseudo ID).

Anonymizing IDs only is not enough for the anonymization of routing policies. The neighboring relationship with the real preference and the real neighbor ID implied by the routing policies can be used for identification as we state in the introduction. Beyond pseudonymization, we might also change the neighboring ASes and local preferences for these neighboring ASes in a routing policy. As for the routing policy in Listing 1.1, we might switch the local preferences of AS 40000 and AS 50000. However, the obfuscation might not preserve the best paths derived by the routing policies. For example, after switching the local preferences in the routing policy of Listing 1.1, AS 45000 prefers routes from AS 50000 over that from AS 40000.

Answering what-if questions by exploiting the inter-domain routing policy relies on the best paths derived from the routing policies. If the best paths of the anonymized routing policy is different from the best paths of the original routing policy, then the anonymized routing policy might not be useful to answer what-if questions. In order to preserve the utility of the routing policy, a best path derived by the original routing policy should be a best path derived by the anonymized routing policy. We say that the anonymization scheme preserves the *best path utility*, if, for an AS-level topology  $G = \{V, E, R\}$ ,  $\forall i, j \in V$ ,  $\pi(Path(i, j)) = Path(\pi(i), \pi(j))$ , where  $Path(i, j)$  indicates the best path from network  $i$  to network  $j$  and  $\pi(Path(i, j))$  maps all original AS numbers of ASes in the best path to their anonymized AS numbers.

In this paper, we propose the anonymization schemes that preserve the best path utility. To do so, all networks and the neighboring relationships are maintained after the routing policy anonymization. Namely, the anonymization schemes should not remove nodes or edges from the AS-level topology. Instead, the anonymization schemes can add new nodes and new edges into the AS-level topology. In addition, the local preferences of neighboring ASes can be changed but the ranking of these local preferences should be maintained. Thus, after the anonymization, a network ranks the neighboring ASes in the same way as before the anonymization.

## 2.3 Threat Model

After the anonymization, the routing policies of all networks are publicly available. The adversary can access the anonymized routing policies of all networks and try to figure out the routing policy for each network. Given an AS number, the adversary needs to identify the routing policy of this AS and discovers how this AS ranks its neighboring ASes from the routing policy. To preserve the best path utility, how an AS ranks its neighboring ASes is maintained in

the anonymized routing policy. Therefore, the key of this de-anonymization procedure is to infer the network mapping function,  $\pi$ .

The AS-level topology implied by the anonymized routing policies can be used to identify the routing policy of each network. The routing tables of RouteViews [4] and RIPE RIS [5] are publicly available. The adversary can get the AS-level topology of the Internet by using the routing tables. To preserve the best path utility, all networks and their connections are maintained in the anonymized AS-level topology. We consider that the adversary tries to use the AS-level topology of the Internet to identify the networks in the anonymized AS-level topology and figure out the routing policy of the networks.

### 3 Anonymization Scheme Overview

To anonymize the routing policies, we first consider how to perturb the AS-level topology implied by the routing policies, so that the networks are hard to be identified through the number of neighbors. Then, we modify the routing policies of these networks accordingly, so that the networks are hard to be identified through the local preferences in the anonymized routing policies.

#### 3.1 Anonymizing Graph Structure

To perturb the AS-level topology, we consider to generate *k-anonymity* graphs. A *k-anonymity* graph is a graph that, for any node, there are at least  $k - 1$  other nodes sharing the same degree with the node. If only the degree information is used by the adversary, the probability that the adversary can correctly identify a real node from the *k-anonymity* graph is at most  $\frac{1}{k}$ .

An anonymized graph that is not *k-anonymity* can also guarantee that the adversary can not correctly identify a real node with the probability higher than  $\frac{1}{k}$ . Following this idea, we extend the definition of *k-anonymity* graph as follows. Given an anonymized graph derived from the original graph, the anonymized graph is a *k-anonymity* graph, if the probability that the adversary can correctly identify a real node is no more than  $\frac{1}{k}$  when the original graph size goes to infinity. Formally,  $\forall u \in V$  and  $\forall v \in V'$ ,

$$\lim_{|V| \rightarrow \infty} P(P(\pi(u) = v | \mathbf{Y} = \mathbf{y}) \leq \frac{1}{k}) = 1, \quad (1)$$

where  $V$  and  $V'$  are the sets of nodes in the original graph and the anonymized graph respectively and  $\mathbf{y}$  is the degree sequence of nodes in the anonymized graph. In Section 4, we propose three graph anonymization schemes to generate *k-anonymity* graph.

Although *k-anonymity* with a small  $k$  is vulnerable to de-anonymization attacks, the schemes proposed in this paper can select a large  $k$ , such as  $k$  equals the number of real nodes. The reason is that the schemes in this paper preserve the best path utility. Selecting a large  $k$  does not impact the utility of the anonymized graph. In the experiment, we will show that selecting a large  $k$  can achieve pretty good anonymization performance.

To preserve the best path utility, all nodes and edges in the original graph should be maintained in the anonymized graph. Therefore, the anonymization schemes should not remove nodes or edges. In this paper, we propose to add fake nodes into the graph. Then, fake nodes and real nodes are connected by a set of fake edges to perturb the graph structure.

#### 3.2 Local Preference Assignment for Fake Edges

After the graph anonymization, we need to modify all routing policies for newly-added nodes and edges. For each fake node added into the topology, we need to create a new routing policy for the node. For each fake edge added into the topology, the local preferences should be properly assigned to the edge. In Section 5, we propose a local preference assignment scheme to determine local preferences for these fake edges, so that the fake nodes can not be identified through their local preferences.

## 4 Node Anonymization

In this section, we propose three anonymization schemes to perturb the AS-level topology. In Section 4.1, we first propose the deterministic anonymization scheme which generates a *k-anonymity* graph. However, the resulting graph is determined by the Internet topology and lacks of randomness. To increase the randomness of the resulting graph, we propose the probabilistic anonymization scheme in Section 4.2. The probabilistic anonymization scheme might significantly enlarge the graph with fake networks. Finally, in Section 4.3, we propose the hybrid anonymization scheme that takes the advantages of both the deterministic and probabilistic anonymization schemes.

#### 4.1 Deterministic Anonymization Scheme

In the deterministic anonymization scheme, we partition all nodes of the original graph into several groups, where each group contains at least  $k$  nodes. Then, we add fake edges between these real nodes and fake nodes to make these real nodes in the same group have the same degree. Finally, fake edges are added between fake nodes to make all fake nodes have the same degree.

**4.1.1 Node Grouping** Given a graph with  $n$  nodes, we partition all real nodes into  $\lfloor \frac{n}{k} \rfloor$  groups, where each group at least has  $k$  nodes. To do so, we sort the nodes of the original graph in the descending order of their degrees. For the first  $\lfloor \frac{n}{k} \rfloor - 1$  group, the group  $i$  is composed of the  $(i * k - k + 1)$ -th node to the  $(i * k)$ -th node. The last group is composed of the  $(\lfloor \frac{n}{k} \rfloor * k - k + 1)$ -th node to the  $n$ -th node.

For each group, we select a target degree. All real nodes in the group are supposed to reach the target degree of this group after the anonymization. To increase the degree of real nodes, we add fake nodes and fake edges. We also select a target degree for all fake nodes. When every node reaches their target degrees, we get a  $k$ -anonymity graph. Please refer to Appendix A.1 for how we select the target degrees and the number of fake nodes added into the anonymized graph.

**4.1.2 Graph Construction** Given the number of fake nodes and the target degrees, we add fake edges to generate  $k$ -anonymity graph. We define *degree gap* of a node as the target degree of this node minus its degree. We use  $g_i$  to denote the degree gap of node  $i$ . The degree of fake nodes in the original graph are zero. Through adding fake edges, we will make the degree gaps of all nodes be zero. To do so, we select one node with the largest degree gap and make the node reach its target degree in each step. We show the details as follows.

*Increasing real node degree* To increase the degree of real nodes, we always select the real node that has the largest degree gap. Suppose that we select node  $i$  with the degree gap  $g_i$ . Then, we select  $g_i$  fake nodes that have the largest degree gap and connect these  $g_i$  fake nodes with node  $i$  with fake nodes.

*Increasing fake node degree* To increase the degree of fake nodes, we always select the fake node that has the largest degree gap. Suppose that we select node  $i$  with the degree gap  $g_i$ . Then, besides node  $i$ , we select  $g_i$  fake nodes that have the largest degree gap and connect these  $g_i$  fake nodes with node  $i$  with fake nodes.

#### 4.1.3 Privacy Guarantee

**Theorem 1.** *The deterministic anonymization scheme can generate a  $k$ -anonymity graph.*

Please refer to Appendix A.1 for the proof of this Theorem.

#### 4.2 Probabilistic Anonymization Scheme

The deterministic scheme lacks of uncertainty. The regularity of the anonymized graph generated by the deterministic scheme can be used for de-anonymization. In order to increase the uncertainty of the anonymization scheme, we propose probabilistic anonymization scheme. In probabilistic scheme, we add randomness during the graph anonymization so that the resulting graph is not deterministic.

**4.2.1 Graph Anonymization** The basic idea of the graph anonymization is to connect a real node and a fake node with the probability,  $p$ . For each real node, the number of its fake neighbors is a random variable that follows binomial distribution with the average as  $m * p$  and the variance as  $m * p(1 - p)$ . When the difference between the degrees of two real nodes in the original graph is insignificant comparing to  $m * p(1 - p)$ , these two real nodes can not be identified through their degrees.

Given a graph  $G$ , we refer to a group with  $k$  nodes in  $G$  as  $k$ -node group. In a graph  $G$ , we define  $k$ -node degree delta of a  $k$ -node group as the difference between the maximal and the minimal degrees of nodes in the group. Then, for a specific node, there are a number of  $k$ -node groups that contain this node. We define *minimal  $k$ -node degree delta* of a node as the minimal  $k$ -node degree delta of  $k$ -node groups that contain the node. We use  $\delta_{(i,k)}$  to denote the minimal  $k$ -node degree delta of node  $i$ . We refer to that  $k$ -node group as minimal  $k$ -node group of this node.

We construct the  $k$ -anonymity graph as follows.

*Adding Fake Nodes* We add  $m = \max(\delta_{(max,k)}^2 \ln n, n)$  fake nodes into the graph, where  $\delta_{(max,k)}$  is the maximal  $\delta_{(i,k)}$  for  $\forall i \in [1, n]$ ,  $n$  is the number of nodes in  $G$  and  $d_{max}$  is the maximal degree in the graph  $G$ . We will show that  $m$  is sufficient in Theorem 2.

*Connecting Real Nodes with Fake Nodes* For each real node, we add a fake edge between this real node and each fake node with a probability  $p$ , where  $0 < p < 1$ . On average, each real node will connect to  $m * p$  fake nodes.

*Connecting Fake Nodes with Fake Nodes* Then, we connect fake nodes with fake nodes. The goal is to make the expected degree of fake nodes be the expected degree of the real node with the highest degree in the original graph. To do that, each fake node connects to another fake node with the probability  $q = \frac{(m-n)p + d_{max}}{m-1}$ .

#### 4.2.2 Probabilistic Privacy Guarantee

**Theorem 2.** *The probabilistic scheme generates a  $k$ -anonymity graph.*

We provide the proof of this Theorem in Appendix A.2. In the following, we just illustrate the intuitive idea that the anonymization method can generate a  $k$ -anonymity graph. In the anonymization procedure, a number of fake edges are added to each real and fake node. With a higher probability (the probability goes to one when the graph size goes to infinity), the  $k$  nodes in the minimal  $k$ -node group of each node will have very similar degrees in the anonymized graph. Because, the minimal  $k$ -node degree delta will be asymptotically smaller than the variance for the number of the fake edges connected to each real node. Apparently, it is more likely that the node with higher degree in original graph has the higher degree in the anonymized graph. However, when the degree of a node in the anonymized graph is dominated by the fake edges added in the anonymization procedure, the degree delta of the nodes in the original graph are insignificant. Then, the  $k$  nodes in each group have similar degrees. The probability of successfully identifying nodes by their degrees is  $1/k$ .

### 4.3 Hybrid Anonymization Scheme

In the probabilistic anonymization scheme, the number of fake nodes needed to perturb the graph rises rapidly with the increasing of the degree difference. Since the degree of networks in AS-level topology follow the power-law distribution, the degree difference between Tier-1 networks is large. Then, the number of fake nodes for anonymization is asymptotically larger than the number of real nodes.

In order to reduce the number of added faked nodes, we propose the hybrid anonymization scheme. We will first use the deterministic anonymization scheme to reduce the degree difference in the graph. After that we derive probabilistic  $k$ -anonymity graph by using the probabilistic anonymization scheme.

**4.3.1 Graph Construction** Hybrid anonymization scheme guarantees to generate a  $k$ -anonymity graph. The key idea is to reduce the degree difference through the technique of the deterministic anonymization scheme and then apply the probabilistic anonymization scheme. We describe the two phases as follows.

*Reducing Degree Delta through Deterministic Anonymization Scheme* Given a graph with  $n$  nodes, to reduce the degree delta, we first partition the nodes into  $\lfloor \frac{n}{k} \rfloor$  groups, where each group at least has  $k$  nodes. In this phase, our aim is to make sure that the degree delta of each group is less than a target degree delta,  $\delta_{target}$ . We select  $\delta_{target} = \lfloor \sqrt{n/\ln n} \rfloor$ . When  $\delta_{max,k}$  is less than  $\delta_{target}$ , in a graph with  $n$  nodes, according to Theorem 2,  $n$  fake nodes are enough to generate a  $k$ -anonymity graph in the next phase. The detailed steps of this phase is similar to the deterministic anonymization scheme. We provide the detailed steps of this phase in Appendix A.3 of the technical report <sup>1</sup>.

*Perturbing Node Degrees through Probabilistic Anonymization Scheme* In this phase, we randomly increase the node degrees by using the same technique of the probabilistic anonymization method. We treat both the real nodes and the fake nodes added by the last phase as the real nodes. Then, we can use the probabilistic anonymization scheme to anonymize the graph.

<sup>1</sup> Please access the technical report through the following anonymous link which does not expose the author identities. <https://drive.google.com/file/d/1TwYZO8Y2kCJ3y6G2-k0D35XF.PiJed7X/view?usp=sharing>. (The underscore, “\_”, between F and P in the URL might disappear when it is copy-pasted into the web browser’s address bar.)

### 4.3.2 Algorithm Analysis

**Theorem 3.** *The hybrid anonymization scheme generates a  $k$ -anonymity graph.*

*Proof.* When the degree delta of each group is reduced into  $\lfloor \sqrt{n/\ln n} \rfloor$ , the minimal  $k$ -node degree delta of each node is equal to or less than  $\lfloor \sqrt{n/\ln n} \rfloor$ . According to Theorem 2, we can prove this theorem.

## 5 Routing Policy Anonymization

In this section, we propose a local preference assignment scheme for fake edges. We aim to guarantee that the real nodes and fake nodes have the similar local preferences, so that fake nodes can not be identified by the local preferences. We consider both *outgoing preference* and *incoming preference*. The outgoing preference of a node is the average of local preferences that are assigned to its neighbors by this node. The incoming preference of a node is the average of local preferences that are assigned to this node by its neighbors. The proposed scheme generates local preferences so that the distribution of the outgoing preferences and the incoming preferences for real nodes are same with that for fake nodes.

The best paths are determined by the ranking of the local preferences instead of the exact value of the local preferences. In this paper, we use the ranking number of these local preferences to represent the original values of these local preferences. For example, if a node assigns local preferences to its four neighbors as 100, 200, 200 and 400, then the associated ranking numbers are 3, 2, 2 and 1 respectively.

### 5.1 Accommodating the Distribution of Outgoing Preference

To adjust the outgoing preferences of a node, we need to determine ranking number for neighbors. Namely, how many neighbors have the ranking number 1, how many neighbors have the ranking number 2 and so on. We refer to the set of neighbors with the ranking number  $i$  as the  $i$ -th *ranking set* of the node. The union of all ranking sets of a node should be the neighbor set of the node. Note, in this step, we determine how many fake edges will be assigned with a ranking number instead of assigning the ranking number to specific fake edges. In the following, we illustrate how we determine the size of the ranking sets for each node.

In each step, we randomly select a node and determine the size of the ranking sets. For each node, we initiate the procedure from determining the size of first ranking set and finish the procedure until the total size of ranking sets equals the number of neighbors. To preserve the best path utility, a real node will not change the local preference for its real neighbors and will not prefer fake neighbors over real neighbors. Namely, the top ranking sets of a real node are already determined by the original routing policy. For these real nodes, we select the size of ranking sets for their fake neighboring relationships only.

To determine the size of  $i$ -th ranking set for a node, we investigate the size of  $i$ -th ranking sets of the other nodes. If the size of  $i$ -th ranking sets of the other nodes are all zero, then we just pick the size for  $i$ -th ranking set of this node. Therefore, the total size of ranking sets of this node equals to the node degree. Otherwise, we select an existing size for  $i$ -th ranking set of this node. More specifically, the probability of selecting the size of  $i$ -th ranking set as  $s$  is as follows.

$$p(s, i) = \begin{cases} 0 & \forall k, N(k, i) = 0 \\ \frac{1/N(s, i)}{\sum_{(k|N(k, i) \neq 0)} 1/N(k, i)} & N(s, i) \neq 0. \end{cases} \quad (2)$$

where  $s$  is the size of a ranking set and  $N(s, i)$  is the number of  $i$ -th ranking sets with the size  $s$ .

### 5.2 Accommodating the Distribution of Incoming Preference

Given the size of all ranking sets, we select specific neighbors for each ranking set. The goal is to make the incoming preference of real nodes have the same distribution with that of fake nodes. To preserve the best path utility, real nodes always prefer their real neighbors over fake neighbors. To avoid fake nodes having large incoming preference, we make fake nodes prefer fake neighbors over real neighbors.

The maximal degree of real nodes,  $d_{max}$ , can be used by the adversary to remove fake edges. For a node  $i$ , only if node  $j$  is one of its top  $d_{max}$  neighbors in the ascending order of local preferences, then the edge  $(i, j)$  could be a real edge. Otherwise, the edge  $(i, j)$  has to be a fake edge. To avoid an edge  $(i, j)$  from being identified as a fake edge, node  $i$  and node  $j$  should treat each other as top  $d_{max}$  neighbors. To do so, before assigning the local preference value to each fake edge, for each node, we select a set of neighbors as its top  $d_{max}$  neighbors. We refer to those neighbors of the node

as *high-ranking neighbors*. Apparently, for a node, there are at most  $d_{max}$  high-ranking neighbors. The rest neighbors are *low-ranking neighbors*. Only the edges between the node and its low-ranking neighbors are identified as fake edges.

In Section 5.2.1, we first describe how a node classifies all neighbors as high-ranking neighbors and low-ranking neighbors, so that less fake edges are identified. Then, in Section 5.2.2, we describe how to assign local preference based on the neighbor classification.

**5.2.1 Neighbor Classification** We start with real nodes to determine their high-ranking neighbors. To preserve the best path utility, a real node has to select all its real neighbors as its high-ranking neighbors. Then, we randomly select a specific number of fake neighbors as its high-ranking neighbors so that the number of high-ranking neighbors of a real node is  $d_{max}$ . Next step, we process fake nodes one by one. For a fake node  $i$ , we select  $d_{max}$  high-ranking neighbors in the following order.

1. The real neighbors of node  $i$ .
2. The fake neighbors that select node  $i$  as their high-ranking neighbors.
3. The fake neighbors that have not been processed.
4. The fake neighbors that select node  $i$  as their low-ranking neighbors.

When  $d_{max}$  high-ranking neighbors are selected, the rest of neighbors will be low-ranking neighbors of node  $i$ .

**5.2.2 Assigning Local Preferences** For real nodes, we determine the local preference for fake neighbors only, since all real neighbors already occupy the first few ranking sets. To satisfy the definition of high-ranking neighbors, the high-ranking neighbors will be preferred over the low-ranking neighbors. Therefore, to start with filling up the top ranking sets, we first randomly select the high-ranking fake neighbors. Finally, we fill up the rest of the ranking sets through randomly selecting low-ranking fake neighbors. For fake nodes, we determine the local preference for both real and fake neighbors. We will first fill up the top ranking sets by using the high-ranking fake neighbors. Then, the following ranking sets are filled up by the high-ranking real neighbors. Finally, the low-ranking neighbors are used to fill up the rest of the ranking sets in a random order.

## 6 Experiment

### 6.1 Datasets

To evaluate the proposed anonymization schemes, we anonymize a series of graphs through these schemes. Since Internet topologies are power-law graphs, we synthesize power-law graphs with various sizes for evaluation. In order to generate power-law graphs that have the same properties with Internet topologies, we exploit the graph generation algorithm in [15]. We also use the Internet AS-level topologies in CAIDA AS Relationship Database [16] to evaluate the performance of these schemes on Internet topologies.

### 6.2 The Size of Anonymized AS-level Topology

Table 1: Number of fake nodes needed to anonymize Internet topologies.

Year	1998	2003	2008	2013	2018
Number of Nodes	3,638	15,320	28,411	44,326	60,874
Maximal Degree	770	2,487	2,874	3,912	7,296
Average Degree	3.7	4.5	5.6	6.7	9.9
Hybrid $n$ -anonymity	10,682	45,424	84,095	131,376	181,224
Hybrid $\sqrt{n}$ -anonymity	5,094	20,114	33,835	51,702	74,784
Hybrid $\ln n$ -anonymity	4,944	19,090	32,655	49,588	71,174
Deterministic $n$ -anonymity	3,621	15,294	28,357	44,250	60,792
Deterministic $\sqrt{n}$ -anonymity	750	2,436	2,765	3,753	7,030
Deterministic $\ln n$ -anonymity	675	1,924	2,174	2,695	5,224



The time complexity of answering what-if questions increases as the size of the graph grows [17, 18]. As a result, adding fake nodes enlarges the cost of using the anonymized routing policies. In this experiment, we consider the number of fake nodes added by the anonymization schemes.

We use Internet topologies to evaluate the number of fake nodes added by the proposed anonymization schemes. We evaluate these schemes when  $k = n$ ,  $k = \ln n$  and  $k = \sqrt{n}$ . We illustrate the number of fake nodes needed to anonymize Internet topologies of various sizes in Table 1. As Table 1 shows, the hybrid and deterministic schemes need a moderate number of fake nodes to anonymize the Internet topologies.

### 6.3 Potential Structure-based De-anonymization

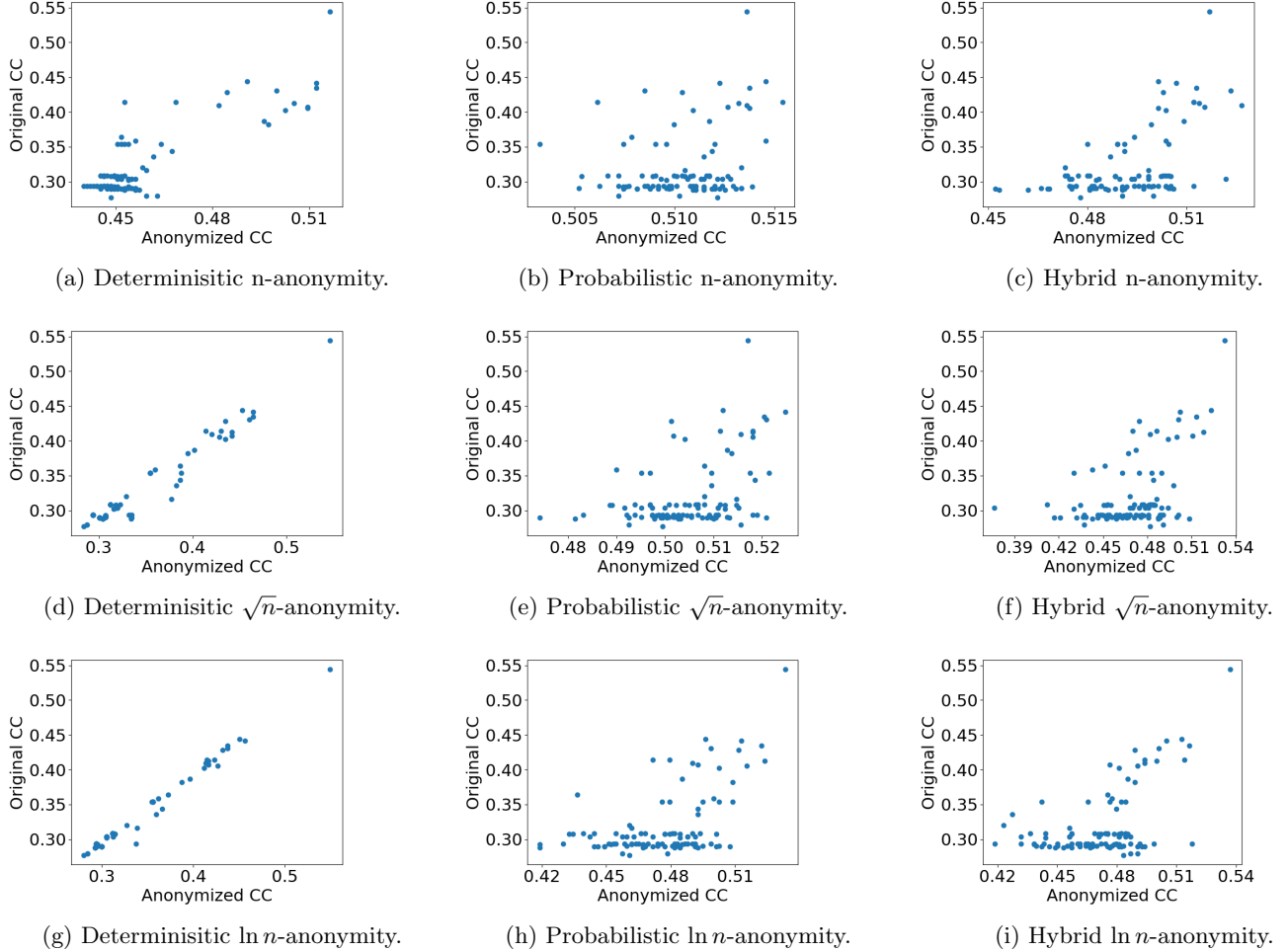


Fig. 1: Distribution of Real Node Closeness Centrality.

All three anonymization schemes of this paper can change the structural feature of a node. To quantify the perturbation, we investigate Closeness Centrality (CC) [19] which measures how long it takes to spread information from a node to all the other nodes sequentially and is defined as the reciprocal of the average shortest path distance to the node over all the other nodes. More specifically, the closeness centrality of node  $u$  is as follow.

$$C(u) = \frac{n-1}{\sum_{v=1}^{n-1} d(v, u)}, \quad (3)$$

where  $d(v, u)$  is the shortest-path distance between node  $v$  and node  $u$ , and  $n$  is the number of nodes that can reach node  $u$ .

We illustrate the distribution of closeness centrality of real nodes before and after anonymization in Figure 1. In the experiment, we generate a power-law graph with size 100 and anonymize the graph with deterministic, probabilistic and hybrid schemes respectively. In those figures, each point represents the CC of a node before and after anonymization, where the y-axis, “Original CC”, is the CC in the original graph and the x-axis, “Anonymized CC”, is the CC in the anonymized graph. According to the experimental results, both probabilistic and hybrid schemes change the original closeness centrality of real nodes significantly. By contrast, as Figure 1(a), Figure 1(d) and Figure 1(g) show, the original CC and the anonymized CC of real nodes have a linear relationship. The linear relationship become stronger when  $k$  is smaller. Since we assume that the adversaries have the structure of the original graph, the strong relationship between the original CC and the anonymized CC of a real node can be potentially used for de-anonymization.

Table 2: Pearson correlation coefficient between Closeness Centrality in the original graph and the anonymized graph.

$k$	Deterministic	Hybrid	Probabilistic
$n$	0.95	0.55	0.26
$\sqrt{n}$	0.97	0.47	0.41
$\ln n$	0.99	0.52	0.53

To quantify the relationship between the original CC and the anonymized CC, we investigate Pearson correlation coefficient of them and illustrate the result in Table 2. As shown in Table 2, the hybrid and probabilistic schemes significantly reduce the correlation between the original CC and the anonymized CC.

#### 6.4 Robustness against De-anonymization Algorithm

There are a number of de-anonymization algorithms which exploit structural information to identify nodes from the anonymized graph. Most of them are seed-based. Those algorithms rely on a seed graph which already maps a subset of real nodes in the original graph to the nodes in the anonymized graph. However, in the context of de-anonymizing network graph, it is hard to get a set of seed nodes. Therefore, in the experiment, we consider the seedless de-anonymization algorithm to identify the real nodes in the anonymized graph. We exploit the state-of-the-art seedless de-anonymization algorithm [20] to identify the nodes from the anonymized graphs.

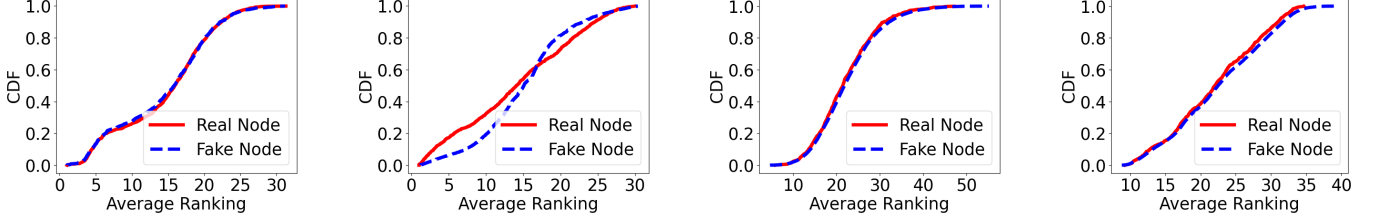
Table 3: Percentage of real nodes that are identified by seedless de-anonymization algorithm [20].

Number of Nodes	100	400	700	1,000
Hybrid $n$ -anonymity	0%	0%	0%	0%
Hybrid $\sqrt{n}$ -anonymity	0.2%	0.1%	0.03%	0%
Hybrid $\ln n$ -anonymity	1.2%	0.15%	0%	0.1%
Deterministic $n$ -anonymity	0%	0%	0%	0.02%
Deterministic $\sqrt{n}$ -anonymity	2.4%	1.05%	0.94%	0.92%
Deterministic $\ln n$ -anonymity	7.2%	5.1%	7.31%	7.28%

Table 3 illustrates the number of real nodes identified through de-anonymization algorithm. In this experiment, we synthesize power-law graphs of various size and perform the anonymization and de-anonymization algorithms. For each size, we generate five graphs and illustrate the average number of real nodes that are correctly identified in Table 3. Although the anonymization schemes aim to perturb the degree of each nodes, to some extent, they can defend the de-anonymization algorithm using structural information. When  $k$  is  $n$ ,  $\sqrt{n}$  or  $\ln n$ , more than 90% real nodes can not be correctly identified.

#### 6.5 Potential De-anonymization Using Local Preference

The local preference of nodes in the anonymized graph might be used to separate real and fake nodes. In this experiment, we investigate the distribution of outgoing and incoming preferences of real nodes and fake nodes. When the distribution for real nodes are similar with that for fake nodes, the local preference is hard to be used to identify real nodes.



(a) Outgoing Preference (DAS). (b) Incoming Preference (DAS). (c) Outgoing Preference (HAS). (d) Incoming Preference (HAS).

Fig. 2: Distribution of local preferences with Deterministic Anonymization Scheme (DAS) and Hybrid Anonymization Scheme (HAS).



Fig. 3: Distribution of local preferences with Hybrid Anonymization Scheme after fake node removal.

We illustrate CDF of outgoing and incoming after deterministic and hybrid anonymization schemes in Figure 2 and Figure 3 respectively. As Figure 2 shows, the distributions for real and fake nodes are very similar. The adversary might or might not use the local preference and the maximal degree to identify fake edges. Namely, for each AS, only the  $d_{max}$  top-ranking links can be real links while the other links have to be fake. These fake edges can be removed. Thus, in the experiments, we consider the scenarios where the adversary uses or does not use fake edge removal. As Figure 3 shows, the distributions for real and fake nodes are still very similar after removing fake edges.

Fake edges in the anonymized graph can be removed through observing the outgoing local preference of each node. Even after removing fake edges, the seedless de-anonymization attack can not identify real nodes from the anonymized graphs. The result is similar to Table 3. Due to the space limit, we do not list the result.

## 7 Related Work

### 7.1 Internet-Wide Network Verification

A number of formal methods have been exploited to verify properties for inter-domain routing. To study the safety property of BGP systems, Satisfiable Module Theories (SMT) [21] and Rewriting Logic [22–24] are used to verify the convergence conditions. A model checking tool is used to search possible attraction attacks on the Internet [25]. A policy-aware model [2] is proposed for routing verification at Internet scale. Those verification schemes need both the Internet topologies and the routing policies of networks for the verification. To preserve the confidentiality of routing policies, the anonymization schemes proposed in this paper can generate anonymized graphs for those verification schemes.

### 7.2 Graph Anonymization

A series of graph anonymization mechanisms have been proposed at the area of online social networks [6–14]. One large class of anonymization mechanisms exploits  $k$ -anonymity, where  $k$  nodes in the anonymized graph can not be separated from each other. The first approach in this direction is proposed in [6], where each node shares the same degrees with at least  $k - 1$  other nodes. At the same time, Zhou and Pei [7] propose the concept of  $k$ -neighborhood anonymity to defend against neighborhood attacks. To defend against more sophisticated structural attacks,  $k$ -automorphism [8] and  $k$ -isomorphism [9] are proposed. Our schemes follow  $k$ -anonymity as well as preserve the best path utility for routing verification.

Another class of anonymization mechanisms use the Differential Privacy technique to provide strong privacy guarantee [10–14]. Those mechanisms generate a synthetic graph which maintains structural similarity to the original graph. However, the best paths of the synthetic graph will be different from that of the original graph, since there is not an injective function mapping from nodes of the original graph to nodes of the synthetic graph.

To preserve edge privacy, Mittal et al. proposed a Random Walk based anonymization technique [26]. Although the preserving edge privacy can avoid local preferences from being revealed, the best paths in the anonymized graph are changed and the graph utility can not be preserved.

### 7.3 Secure Multi-Party Computation

Recently, to preserve policy privacy for networks, Secure Multi-Party Computation (SMPC) methods are proposed for policy-compliant routes computation [27], [28] and [29]. SMPC methods maybe applicable to part of verification problems which focus on the property of the best routes derived from a specific routing policy. Namely, the best routes are derived by SMPC methods and then verified against network properties. However, it is costly to verify the routes derived from a range of routing policies, since the routes need to be derived for each routing policy. It will be infeasible, when the verification is for numerous routing policies. In contrast, the anonymization methods proposed in this paper allow the routing verification tools to be applied to the anonymized graph directly. The overhead is the additional cost to perform routing verification on an enlarged graph.

## 8 Conclusion

In this paper, we anonymize the inter-domain routing policies. We design the anonymization schemes to guarantee the policy privacy of networks without jeopardizing the best path utility of the routing policies. To perturb the Internet topology, we propose deterministic, probabilistic and hybrid anonymization schemes. These three schemes guarantee that degree-based de-anonymization can not be identified network from the anonymized graph with a probability higher than  $\frac{1}{k}$ . Even if the de-anonymization attack using structural information can not identify more than 1% of networks. After topology anonymization, we anonymize the routing policies accordingly. To do so, we propose a local preference assignment scheme. The resulting local preferences in the anonymized graph preserve the best path utility and can not be used to identify networks.

## 9 Acknowledgement

This work was supported in part by NSF grants CNS- 1900866 and CCF-1918187, and the Radcliffe fellowship. We thank Catherine Greenhill for many helpful discussions.

## References

1. K. Weitz, D. Woos, E. Torlak, M. D. Ernst, A. Krishnamurthy, and Z. Tatlock, “Scalable verification of border gateway protocol configurations with an smt solver,” in Proceedings of the 2016 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2016, (New York, NY, USA), pp. 765–780, Association for Computing Machinery, 2016.
2. X. Shao and L. Gao, “Verifying policy-based routing at internet scale,” in IEEE INFOCOM 2020 - IEEE Conference on Computer Communications, pp. 2293–2302, 2020.
3. G. Siganos, M. Faloutsos, P. Faloutsos, and C. Faloutsos, “Power laws and the as-level internet topology,” IEEE/ACM Trans. Netw., vol. 11, p. 514524, Aug. 2003.
4. University of Oregon, “Route Views.” [www.routeviews.org](http://www.routeviews.org), 2021.
5. RIPE NCC, “Routing information service.” <https://www.ripe.net/analyse/internet-measurements/routing-information-service-ris>.
6. K. Liu and E. Terzi, “Towards identity anonymization on graphs,” in Proceedings of the 2008 ACM SIGMOD international conference on Management of data, pp. 93–106, 2008.
7. B. Zhou and J. Pei, “Preserving privacy in social networks against neighborhood attacks,” in 2008 IEEE 24th International Conference on Data Engineering, pp. 506–515, IEEE, 2008.
8. L. Zou, L. Chen, and M. T. Özsu, “K-automorphism: A general framework for privacy preserving network publication,” Proceedings of the VLDB Endowment, vol. 2, no. 1, pp. 946–957, 2009.
9. J. Cheng, A. W.-c. Fu, and J. Liu, “K-isomorphism: privacy preserving network publication against structural attacks,” in Proceedings of the 2010 ACM SIGMOD International Conference on Management of data, pp. 459–470, 2010.

10. A. Sala, X. Zhao, C. Wilson, H. Zheng, and B. Y. Zhao, "Sharing graphs using differentially private graph models," in *Proceedings of the 2011 ACM SIGCOMM conference on Internet measurement conference*, pp. 81–98, 2011.
11. D. Proserpio, S. Goldberg, and F. McSherry, "A workflow for differentially-private graph synthesis," in *Proceedings of the 2012 ACM workshop on Workshop on online social networks*, pp. 13–18, 2012.
12. C. Dwork, F. McSherry, K. Nissim, and A. Smith, "Calibrating noise to sensitivity in private data analysis," in *Theory of cryptography conference*, pp. 265–284, Springer, 2006.
13. Y. Wang and X. Wu, "Preserving differential privacy in degree-correlation based graph generation," *Transactions on data privacy*, vol. 6, no. 2, p. 127, 2013.
14. Q. Xiao, R. Chen, and K.-L. Tan, "Differentially private network data release via structural inference," in *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 911–920, 2014.
15. A. Elmokashfi, A. Kvalbein, and C. Dovrolis, "On the scalability of bgp: The role of topology growth," *IEEE Journal on Selected Areas in Communications*, vol. 28, no. 8, pp. 1250–1261, 2010.
16. CAIDA, "The caida as relationships dataset, <1998-2020>," Oct 2019.
17. A. Abhashkumar, A. Gember-Jacobson, and A. Akella, "Tiramisu: Fast multilayer network verification," in *17th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 20)*, pp. 201–219, 2020.
18. R. Beckett, A. Gupta, R. Mahajan, and D. Walker, "Control plane compression," in *Proceedings of the 2018 Conference of the ACM Special Interest Group on Data Communication, SIGCOMM 18*, (New York, NY, USA), pp. 476–489, Association for Computing Machinery, 2018.
19. L. C. Freeman, "Centrality in social networks conceptual clarification," *Social networks*, vol. 1, no. 3, pp. 215–239, 1978.
20. S. Ji, W. Li, M. Srivatsa, and R. Beyah, "Structural data de-anonymization: Quantification, practice, and implications," in *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security, CCS '14*, (New York, NY, USA), pp. 1040–1053, Association for Computing Machinery, 2014.
21. A. Wang, L. Jia, W. Zhou, Y. Ren, B. T. Loo, J. Rexford, V. Nigam, A. Scedrov, and C. Talcott, "Fsr: Formal analysis and implementation toolkit for safe interdomain routing," *IEEE/ACM Trans. Netw.*, vol. 20, pp. 1814–1827, Dec. 2012.
22. A. Wang, A. J. T. Gurney, X. Han, J. Cao, B. T. Loo, C. Talcott, and A. Scedrov, "A reduction-based approach towards scaling up formal analysis of internet configurations," in *IEEE INFOCOM 2014 - IEEE Conference on Computer Communications*, pp. 637–645, April 2014.
23. A. Wang, C. Talcott, A. J. T. Gurney, B. T. Loo, and A. Scedrov, "Reduction-based formal analysis of bgp instances," in *Proceedings of the 18th International Conference on Tools and Algorithms for the Construction and Analysis of Systems, TACAS'12*, (Berlin, Heidelberg), pp. 283–298, Springer-Verlag, 2012.
24. A. Wang, C. Talcott, L. Jia, B. T. Loo, and A. Scedrov, "Analyzing bgp instances in maude," in *Proceedings of the Joint 13th IFIP WG 6.1 and 30th IFIP WG 6.1 International Conference on Formal Techniques for Distributed Systems, FMOODS'11/FORTE'11*, (Berlin, Heidelberg), pp. 334–348, Springer-Verlag, 2011.
25. A. Sosnovich, O. Grumberg, and G. Nakibly, "Analyzing internet routing security using model checking," in *Proceedings of the 20th International Conference on Logic for Programming, Artificial Intelligence, and Reasoning - Volume 9450, LPAR-20 2015*, (New York, NY, USA), pp. 112–129, Springer-Verlag New York, Inc., 2015.
26. P. Mittal, C. Papamanthou, and D. Song, "Preserving link privacy in social network based systems," *arXiv preprint arXiv:1208.6189*, 2012.
27. G. Asharov, D. Demmler, M. Schapira, T. Schneider, G. Segev, S. Shenker, and M. Zohner, "Privacy-preserving interdomain routing at internet scale," *Proceedings on Privacy Enhancing Technologies*, vol. 2017, no. 3, pp. 147–167, 2017.
28. Q. Chen, C. Qian, and S. Zhong, "Privacy-preserving cross-domain routing optimization - a cryptographic approach," in *2015 IEEE 23rd International Conference on Network Protocols (ICNP)*, pp. 356–365, 2015.
29. D. Gupta, A. Segal, A. Panda, G. Segev, M. Schapira, J. Feigenbaum, J. Rexford, and S. Shenker, "A new approach to interdomain routing based on secure multi-party computation," in *Proceedings of the 11th ACM Workshop on Hot Topics in Networks*, pp. 37–42, 2012.
30. B. D. McKay and F. Skerman, "Degree sequences of random digraphs and bipartite graphs," *arXiv preprint arXiv:1302.2446*, 2013.

## A Appendix

### A.1 Determining Target Degrees in Deterministic Anonymization Schemes

In the deterministic anonymization scheme, we need to determine the target degree of each group and the number of fake nodes needed. In the following, we first describe how to select the target degrees. Then, we will show that the deterministic anonymization scheme generates  $k$ -anonymity graph.

We denote  $d_{target}^i$  as the target degree of the group  $i$ . When fake nodes and real nodes have the same target degrees, in the anonymized graph, fake nodes can not be identified through node degrees. Therefore, we select the target degree of the fake nodes as the target degree of one of the real node groups. When we select the target degree of group  $x$  as the target degree of fake nodes, the target degree of fake nodes is  $d_{target}^x$ . In this case, we denote the number of fake nodes added into the graph as  $m_x$ .

To generate  $k$ -anonymity graph, the target degree and the number of fake nodes have to satisfy a series of constraints as follows.

*The number of fake nodes are adequate* To guarantee that all real nodes can reach their target degrees, the number of fake nodes should be large enough. Since we does not consider multiple edges and self-loops in the graph, the number of fake nodes should be larger than the maximal degree gap of real nodes. We denote  $d_{max}^i$  and  $d_{min}^i$  as the maximal degree and minimal degree in the group  $i$ . Namely,

$$m_x \geq \max\{g_i | i \in V\} \quad (4)$$

In addition, the total degree gap of real nodes should not be more than the total degree gap of fake nodes.

$$m_x * d_{target}^x \geq \sum_{i \in \{RealNodes\}} g_i \quad (5)$$

*All nodes are able to reach their target degrees* In the anonymization procedure, we add fake edges but not remove real edges. Therefore, the target degree of a node can not be less than the original degree of the node. Then, the target degree of each group should not be less than the max degree in the group.

$$d_{target}^i \geq d_{max}^i \quad (6)$$

When all nodes reach their target degrees, the total degrees in the anonymized graph should be even. Namely,  $|E| * 2 + \sum_{i \in \{RealNodes\}} g_i + m_x * d_{target}^x$  is even. That is, we need to select  $m_x$  or  $d_{target}^x$  to make  $m_x * d_{target}^x - \sum_{i \in \{RealNodes\}} g_i$  be even.

Given these constraints, we can minimize  $m_x$ , when  $x \in [1, \lfloor \frac{n}{k} \rfloor]$ . To minimize the anonymization cost, we can pick the number of fake node as  $m = \min(m_x)$ .

**Theorem 1.** *The deterministic anonymization scheme can generate a  $k$ -anonymity graph.*

*Proof.* During each phase of the graph construction in Section 4.1, all constraints hold. Therefore, this theorem is proved.

## A.2 Proof for Probabilistic Anonymization Scheme

In this section, we will prove Theorem 2.

**A.2.1 Model of Probabilistic Anonymization Scheme** Assume there are  $n$  nodes in the original graph, with the index from 1 to  $n$ . The degree of  $i$ th real node is  $d_i$ . We add  $m$  fake nodes into the graph. The  $m$  fake nodes have the index from  $n + 1$  to  $n + m$ . Namely, the degree of  $i$ th fake node is  $d_{n+i} = 0$ . In the anonymization scheme, node ID are removed and each node is assigned a new ID. Without loss of generality, we assume that the new IDs are still in the range  $[1, n + m]$ .

We define a random variable  $R_i$  as a *degree noise* for node  $i$  which is the fake edges that are added in the anonymization scheme. The degree of  $i$ th node is increased from  $d_i$  to  $Z_i = d_i + R_i$ , where  $Z_i$  is the degree of a node after anonymization. We define all degrees in the original graph as a vector  $\mathbf{D} = [d_1, d_2, \dots, d_{n+m}]$ .

We model the anonymization as a random permutation  $\pi$  on the set of  $n + m$  nodes. Each node,  $i$ , is assigned the pseudonym  $\pi(i)$  in the anonymized graph.  $\mathbf{Y}$  is the anonymized version of  $\mathbf{Z}$ . That is  $Y_{\pi(i)} = Z_i$  and  $Y_i = Z_{\pi^{-1}(i)}$ . Thus,

$$\mathbf{Y} = [Z_{\pi^{-1}(1)}, Z_{\pi^{-1}(2)}, \dots, Z_{\pi^{-1}(n+m)}] = [Y_1, Y_2, \dots, Y_{n+m}] \quad (7)$$

**A.2.2 Preliminary** We first show that when the number of real nodes,  $n$ , goes to infinity,  $\forall i \in V$ , the probability that the degree noise,  $R_i$ , is around  $mp$  goes to 1. Thus, in the anonymized graph, the degrees of all nodes are in a narrow range. As a result, it is hard to identify a real node through the anonymized degree.

The degree noise of real nodes is around  $mp$ .

**Lemma 1.** *For  $i \in [1, n]$ ,*

$$P(|R_i - mp| \leq \epsilon) \geq 1 - 2e^{-\frac{\epsilon^2}{3mp}} \quad (8)$$

*Proof.* Since  $R_i$  follows Binomial distribution, the Lemma can be proved through applying Chernoff Bound.

The degree noise of fake nodes is also around  $mp + d_{max}$ .

**Lemma 2.** *For  $i \in [n + 1, n + m]$ ,  $\zeta > 0$  and  $\eta > 0$ ,*

$$P(|R_i - mp - d_{max}| \leq \zeta + \eta) \geq 1 - 2e^{-\frac{\zeta^2}{3np}} - 2e^{-\frac{\eta^2}{3((m-n)p + d_{max})}} \quad (9)$$

*Proof.*  $R_i$  is the sum of two random variables following Binomial distribution. Namely,  $R_i = R_{i,r} + R_{i,f}$ . The first random variable,  $R_{i,r}$ , is the number of fake edges connecting with real nodes while the second random variable,  $R_{i,f}$ , is the number of fake edges connecting with fake nodes. Then, by applying union bound and Chernoff bound, we have the following derivation.

$$\begin{aligned}
& P(|R_i - mp - d_{max}| \leq \zeta_i + \eta_i) \geq P(|R_{i,r} - np| \leq \zeta_i \cap |R_{i,f} - (m-n)p - d_{max}| \leq \eta_i) \\
& = 1 - P(|R_{i,r} - np| \geq \zeta_i \cup |R_{i,f} - (m-n)p - d_{max}| \geq \eta_i) \\
& \geq 1 - P(|R_{i,r} - np| \geq \zeta_i) - P(|R_{i,f} - (m-n)p - d_{max}| \geq \eta_i) \quad (\text{Union bound}) \\
& \geq 1 - 2e^{-\frac{\zeta_i^2}{3np}} - 2e^{-\frac{\eta_i^2}{3((m-n)p + d_{max})}} \quad (\text{Chernoff bound})
\end{aligned} \tag{10}$$

With high probability, the probabilistic scheme generates an anonymized graph, in which the degrees of all nodes are in the range,  $[mp + d_{min} - (1 + \alpha)\sqrt{3mp \ln n}, mp + d_{max} + (1 + \alpha)\sqrt{3mp \ln n}]$ , where  $\alpha > 0$ .

**Lemma 3.** Given any  $\alpha > 0$ ,

$$\lim_{n \rightarrow \infty} P(\cap_{i=1}^{n+m} (mp + d_{min} - (1 + \alpha)\sqrt{3mp \ln n} \leq Y_i \leq mp + d_{max} + (1 + \alpha)\sqrt{3mp \ln n}) = 1 \tag{11}$$

*Proof.* We first consider a narrow range in terms of each  $Y_i$ . For real nodes, we consider the range as  $[mp - (1 + \alpha)\sqrt{3mp \ln n} + d_{min}, mp + (1 + \alpha)\sqrt{3mp \ln n} + d_{max}]$ . For fake nodes, we consider the range as  $[mp + d_{max} - (1 + \alpha)(\sqrt{3np \ln n} + \sqrt{((m-n)p + d_{max}) \ln n}), mp + d_{max} + (1 + \alpha)(\sqrt{3np \ln n} + \sqrt{((m-n)p + d_{max}) \ln n})]$ . Then, we use Lemma 1 and Lemma 2 to prove this Lemma.

For simplicity of expression, we define  $\hat{\epsilon} = (1 + \alpha)\sqrt{3mp \ln n}$ ,  $\hat{\zeta} = (1 + \alpha)\sqrt{3np \ln m}$  and  $\hat{\eta} = (1 + \alpha)\sqrt{3((m-n)p + d_{max}) \ln m}$ . The derivation is as follows.

$$\begin{aligned}
& P(\cap_{i=1}^{n+m} (mp - (1 + \alpha)\sqrt{3mp \ln n} \leq Y_i \leq mp + d_{max} + (1 + \alpha)\sqrt{3mp \ln n}) \\
& \geq P(\cap_{i=1}^n (|R_i - mp| \leq \hat{\epsilon}) \cap_{j=n+1}^{n+m} (|R_j - mp - d_{max}| \leq \hat{\zeta} + \hat{\eta})) \\
& = 1 - P(\cup_{i=1}^n (|R_i - mp| \geq \hat{\epsilon}) \cup_{j=n+1}^{n+m} (|R_j - mp - d_{max}| \geq \hat{\zeta} + \hat{\eta})) \\
& \geq 1 - \sum_{i=1}^n P(|R_i - mp| \geq \hat{\epsilon}) - \sum_{j=n+1}^{n+m} P(|R_j - mp - d_{max}| \geq \hat{\zeta} + \hat{\eta}) \quad (\text{Union bound}) \\
& \geq 1 - 2ne^{-\frac{\hat{\epsilon}^2}{3mp}} - 2me^{-\frac{\hat{\zeta}^2}{3np}} - 2me^{-\frac{\hat{\eta}^2}{3(m-n)p + 3d_{max}}} \geq 1 - \frac{2}{n^\alpha} - \frac{2}{m^\alpha} - \frac{2}{m^\alpha}
\end{aligned} \tag{12}$$

Since,  $m \geq n$ , the lemma is proved.

In the anonymized graph generated by the probabilistic scheme, there is a bipartite graph where all real nodes in one side and all fake nodes in the other side. Whether an edge exists in the bipartite graph is determined by a Bernoulli trial with probability  $p$ . Given the degree sequence of real nodes,  $S = (s_1, \dots, s_n)$ , and the degree sequence of fake nodes,  $T = (t_1, \dots, t_m)$ , in the bipartite graph, we can derive the probability that the probabilistic scheme generates the bipartite graph as follows

**Lemma 4.** If  $\sum_{i=1}^n s_i = \sum_{j=1}^m t_j = K$  and  $p q m n \rightarrow \infty$ , then

$$P(\mathbf{S} = S, \mathbf{T} = T) = (2 + O((pqmn)^{-1})) p^{2K} q^{2mn-2K} \sqrt{\pi pqmn} \prod_{i=1}^n \binom{m}{s_i} \prod_{j=1}^m \binom{n}{t_j} \tag{13}$$

*Proof.* This Lemma is proved in Lemma 4 of [30].

**Lemma 5.** Given  $\tilde{\pi} \in \{\pi | \pi(u) = o\}$  and  $\hat{\pi} \in \{\pi | \pi(u) = v\}$ , where  $u \in \{1, \dots, n\} \setminus \{v\}$ . If  $w = \hat{\pi}^{-1}(o) \leq n$  and  $\text{dif}(\hat{\pi}, \tilde{\pi}) = 2$ , then

$$\lim_{n \rightarrow \infty} \frac{P(\tilde{\pi} \cap \mathbf{Y} = \mathbf{y})}{P(\hat{\pi} \cap \mathbf{Y} = \mathbf{y})} = 1 \tag{14}$$

*Proof.* Given a permutation  $\pi$ , we refer to the vector of real node degree after anonymization as  $s(\mathbf{y}, \pi)$ . Specifically,  $s_i = y_{\pi(i)} - d_i$  for  $i \in [1, n]$ . Accordingly, we refer to the vector of fake node degree after anonymization as  $r(\mathbf{y}, \pi)$ , where  $r_i = y_{\pi(i+n)}$  for  $i \in [1, m]$ .

Each edge of a fake node can be classified into one of two groups: edges connecting to fake nodes and edges connecting to real nodes. Then, for each fake node, we define *real neighbor degree* as the number of edges connecting to real nodes while *fake neighbor degree* as the number of edges connecting to fake nodes. Then,  $r(\mathbf{y}, \pi)$  can be the sum of two vectors,  $u(\mathbf{y}, \pi)$  and  $t(\mathbf{y}, \pi)$ , where  $u(\mathbf{y}, \pi)$  refers to the vector of fake neighbor degrees of all fake nodes and  $t(\mathbf{y}, \pi)$  refers to the vector of real neighbor degrees of all fake nodes.

We use  $K = \sum s(\mathbf{y}, \pi) = \sum t(\mathbf{y}, \pi)$  to represent the number of fake edges between real nodes and fake nodes.

$$\frac{P(\mathbf{Y}=\mathbf{y}|\hat{\pi})P(\hat{\pi})}{P(\mathbf{Y}=\mathbf{y}|\tilde{\pi})P(\tilde{\pi})} = \frac{\sum_{t+u=y} P(s(\mathbf{y}, \hat{\pi}), t(\mathbf{y}, \hat{\pi})|\hat{\pi}) * P(u(\mathbf{y}, \hat{\pi})|\hat{\pi})}{\sum_{t+u=y} P(s(\mathbf{y}, \tilde{\pi}), t(\mathbf{y}, \tilde{\pi})|\tilde{\pi}) * P(u(\mathbf{y}, \tilde{\pi})|\tilde{\pi})} \quad (15)$$

$P(s(\mathbf{y}, \hat{\pi}), t(\mathbf{y}, \hat{\pi})|\hat{\pi})$  is the probability to generate a bipartite graph with one set containing real nodes and the other set containing fake nodes, where  $s(\mathbf{y}, \pi)$  represents the fake degree sequence of the real nodes given the permutation  $\pi$  and  $t(\mathbf{y}, \pi)$  represents the real degree sequence of the fake nodes given the permutation  $\pi$ .  $P(u(\mathbf{y}, \hat{\pi})|\hat{\pi})$  is the probability to generate a graph with fake nodes only, where  $u(\mathbf{y}, \pi)$  represents the fake degree sequence of the fake nodes given the permutation  $\pi$ .

Since the difference of  $\hat{\pi}$  and  $\tilde{\pi}$  is two and the switch is between two real nodes, we have  $t(\mathbf{y}, \hat{\pi}) = t(\mathbf{y}, \tilde{\pi})$ . Accordingly, the number of items in the denominator and the numerator of Equation 15 are same. If the ratio of a pair of items in the denominator and the numerator uniformly converges, then we can prove the limit of Equation 15 is one with  $n \rightarrow \infty$ . For simplicity of expression, we define a function  $F(t)$  to represent that ratio as follows. Then, we use Lemma 4 to calculate  $F(t)$ . For a real degree sequence of fake nodes,  $t(\hat{\pi})$ , where  $0 \leq t_i \leq n$  and  $u_i + t_i = y_{\hat{\pi}^{-1}(i+n)}$ .

$$F(t(\hat{\pi}), t(\tilde{\pi})) = \frac{P(s(\mathbf{y}, \hat{\pi}), t(\mathbf{y}, \hat{\pi})|\hat{\pi}) * P(u(\mathbf{y}, \hat{\pi})|\hat{\pi})}{P(s(\mathbf{y}, \tilde{\pi}), t(\mathbf{y}, \tilde{\pi})|\tilde{\pi}) * P(u(\mathbf{y}, \tilde{\pi})|\tilde{\pi})} = \frac{(2+O((pqmn)^{-1})) \binom{m}{y_o-d_u} * \binom{m}{y_v-d_w}}{(2+O((pqmn)^{-1})) \binom{m}{y_o-d_w} * \binom{m}{y_v-d_u}} \quad (16)$$

Note,  $\forall t(\pi)$ ,  $F(t(\hat{\pi}), t(\tilde{\pi}))$  is in the same order and it is not a function of  $t(\pi)$ . That is, for any  $t(\pi)$ , if  $F(t(\hat{\pi}), t(\tilde{\pi}))$  converge, then  $F(t(\hat{\pi}), t(\tilde{\pi}))$  converges uniformly.

$$\begin{aligned} \lim_{n \rightarrow \infty} F(t(\hat{\pi}), t(\tilde{\pi})) &= \lim_{n \rightarrow \infty} \frac{(2+O((pqmn)^{-1})) \binom{m}{y_o-d_u} * \binom{m}{y_v-d_w}}{(2+O((pqmn)^{-1})) \binom{m}{y_o-d_w} * \binom{m}{y_v-d_u}} = \lim_{n \rightarrow \infty} \frac{\binom{m}{y_o-d_u} * \binom{m}{y_v-d_w}}{\binom{m}{y_o-d_w} * \binom{m}{y_v-d_u}} \\ &= \lim_{n \rightarrow \infty} \frac{\prod_{i=y_o-d_u-1}^{y_o-d_u-1} (m-i) \prod_{j=y_o}^{y_v+1} (j-d_w)}{\prod_{i=y_o-d_w-1}^{y_o-d_w-1} (m-i) \prod_{j=y_o}^{y_v+1} (j-d_u)} = \lim_{n \rightarrow \infty} \frac{\prod_{i=y_v-d_w-1}^{y_v-d_w-1} (m-i) \prod_{j=d_w}^{d_u-1} (y_o-j)}{\prod_{i=y_o-d_u-1}^{y_o-d_u-1} (m-i) \prod_{j=d_w}^{d_u-1} (y_v-j)} \end{aligned} \quad (17)$$

Then, we can get the upper and the lower bound of the limit of  $F(t(\hat{\pi}), t(\tilde{\pi}))$ .

$$\lim_{n \rightarrow \infty} \left( \frac{(m-y_v+d_w+1)(y_o-d_u+1)}{(m-y_o+d_u)(y_v-d_w)} \right)^{d_u-d_w} \leq \lim_{n \rightarrow \infty} F(t(\hat{\pi}), t(\tilde{\pi})) \leq \lim_{n \rightarrow \infty} \left( \frac{(m-y_v+d_u)(y_o-d_w)}{(m-y_o+d_w+1)(y_v-d_u+1)} \right)^{d_u-d_w} \quad (18)$$

Then, we can get the following bounds.

$$\begin{aligned} &\lim_{n \rightarrow \infty} (d_u-d_w)^{\frac{y_o-y_v-d_u+d_w+1}{m-y_o+d_u}} * \lim_{n \rightarrow \infty} (d_u-d_w)^{\frac{y_o-y_v-d_u+d_w+1}{y_o-d_u}} \\ &\leq \lim_{n \rightarrow \infty} F(t(\hat{\pi}), t(\tilde{\pi})) \\ &\leq \lim_{n \rightarrow \infty} (d_u-d_w)^{\frac{y_o-y_v+d_u-d_w-1}{m-y_o+d_w+1}} * \lim_{n \rightarrow \infty} (d_u-d_w)^{\frac{y_o-y_v+d_u-d_w-1}{y_v-d_u+1}} \end{aligned} \quad (19)$$

Note,  $d_u - d_w = O(\delta_{(max,k)})$ , since node  $u$  and node  $w$  should be in the same group. According to Lemma 3, we know that  $y_o - y_v = O(d_{max} - d_{min} + \sqrt{mp \ln n})$ . Since  $m \geq \delta_{(max,k)}^2 \ln n$ , Equation 14 is true. This Lemma is proved.

### A.2.3 Proof of the Theorem

**Theorem 2.** *The probabilistic scheme generates a  $k$ -anonymity graph.*

*Proof.* We define  $\Pi(u, v) = \{\pi | \pi(u) = v\}$  as a set of permutations which map node  $u$  in the original graph to node  $v$  in the anonymized graph.

$$P(\pi(u) = v | \mathbf{Y} = \mathbf{y}) = \frac{\sum_{\hat{\pi} \in \Pi(u, v)} P(\hat{\pi} \cap \mathbf{Y} = \mathbf{y})}{\sum_{r=1}^{n+m} \sum_{\hat{\pi} \in \Pi(u, r)} P(\hat{\pi} \cap \mathbf{Y} = \mathbf{y})} \quad (20)$$

The size of  $\Pi(u, v)$  is  $|\Pi(u, v)| = (n + m - 1)!$ . For any  $\pi \in \Pi(u, v)$ , there are  $(m + n - 1)$  permutations, each of them,  $\pi' \in \Pi(u, o)$ , where  $o \in \{1, \dots, n + m\} \setminus \{v\}$ ,  $\pi'(\pi^{-1}(o)) = v$  and  $\text{diff}(\pi, \pi') = 2$ . Then, for each  $\pi \in \Pi(u, v)$ , we define that  $(m + n - 1)$  permutations as a *one-swap set* of  $\pi$ . According to the one-swap set, we can rewrite Equation 20 as the following.

$$P(\pi(u) = v | \mathbf{Y} = \mathbf{y}) = \frac{\sum_{\hat{\pi} \in \Pi(u, v)} P(\hat{\pi} \cap \mathbf{Y} = \mathbf{y})}{\sum_{r=1}^{n+m} \sum_{\hat{\pi} \in \Pi(u, r)} P(\hat{\pi} \cap \mathbf{Y} = \mathbf{y})} = \frac{\sum_{\hat{\pi} \in \Pi(u, v)} P(\hat{\pi} \cap \mathbf{Y} = \mathbf{y})}{\sum_{\hat{\pi} \in \Pi(u, v)} \sum_{\pi' \in S(\hat{\pi}) \cup \pi} P(\pi' \cap \mathbf{Y} = \mathbf{y})} \quad (21)$$

In addition,  $n - 1$  permutations in  $S(\pi)$  has the property:  $\pi'^{-1}(v) \leq n$ , namely the two different mapping of  $\pi$  and  $\pi'$  happen between two real nodes,  $u$  and  $\pi'^{-1}(v)$ .

According to Lemma 5, we know that the ratio is one if  $m$  and  $n$  go to infinity and the swap happens between two real nodes that the degree delta,  $\Delta_d = o(\sqrt{d_{max}} + n^{\frac{2}{5}})$ . Note, for each node, there is a minimal  $k$ -node group that contains  $k$  nodes with the degree delta smaller than  $\Delta_d$ . When the swap is between a real node and a fake node, the ratio is positive. Therefore,  $\lim_{n \rightarrow \infty} \frac{\sum_{\pi' \in S(\hat{\pi}) \cup \pi} P(\pi' \cap \mathbf{Y} = \mathbf{y})}{P(\hat{\pi} \cap \mathbf{Y} = \mathbf{y})} \geq k$ . Finally, this theorem is proved.