

# Benchmarking Fraud Detectors on Private Graph Data

Alexander Goldberg  
Carnegie Mellon University  
Pittsburgh, PA, USA  
akgoldbe@andrew.cmu.edu

Nihar Shah  
Carnegie Mellon University  
Pittsburgh, PA, USA  
nihars@andrew.cmu.edu

Giulia Fanti  
Carnegie Mellon University  
Pittsburgh, PA, USA  
gfanti@andrew.cmu.edu

Steven Wu  
Carnegie Mellon University  
Pittsburgh, PA, USA  
zhiweiw@andrew.cmu.edu

## Abstract

Fraud poses a pervasive challenge across diverse domains. Currently, many types of fraud are managed in part by automated detection algorithms that operate over graph-structured data. We consider a scenario where a data holder wishes to outsource the development of fraud detection algorithms to third parties (e.g., vendors or researchers) on a private graph-structured dataset. The third parties submit their fraud detectors to the data holder, who evaluates these algorithms on the private dataset. The evaluation results are then communicated publicly. We study how to evaluate and release these benchmarking results while satisfying a formal *differential privacy* (DP) guarantee. DP evaluation of fraud detection algorithms over graph data has not been explicitly studied in the literature, so we empirically evaluate two classes of solutions: subsample-and-aggregate and DP synthetic graph data. We demonstrate through extensive empirical experiments that current approaches fail to provide utility when formally guaranteeing privacy. Our results indicate that the error arising from differential privacy trades off between bias from distorting graph structure and variance from adding random noise. Current methods lie on different points along this bias-variance trade-off, but more complex methods tend to require high-variance noise addition, undermining utility. Hence, the most competitive baseline is surprisingly simple: it models the underlying graph as a stochastic block model, estimates the parameters with DP, and then evaluates fraud detectors on a synthetic graph drawn from this learned generative model.

## Keywords

differential privacy, graph analysis, model evaluation, synthetic data, fraud detection

## 1 Introduction

Fraud constitutes a pernicious and widespread problem across numerous domains, manifesting as fake product reviews, fraudulent payments, and the resale of stolen goods, among other harms [6]. The scale of fraud losses is driven in part by the difficulty of detecting fraud: today, the problem is primarily handled by automated

detectors with high false positive and false negative rates. Although many organizations dedicate entire teams to fraud detection, other organizations partially (or entirely) outsource the development of fraud detection mechanisms to third parties, such as vendors of fraud detection software and/or third-party researchers [7, 25, 39]. However, effective outsourcing requires enterprises to share internal fraud data, which can be challenging or impossible due to privacy regulations (e.g., GDPR) and/or the risk of leaking trade secrets through shared datasets. As a result, the lack of publicly shareable data has limited research progress on detection of fraudulent behaviors in privacy-sensitive domains. For example, in scientific peer review, there is a lack of data on reviewer-paper assignments. This unavailability limits researchers’ ability to evaluate the efficacy of potential solutions to the problem of detecting rings of colluding reviewers [28, 41].

In this work, we explore a paradigm for outsourcing fraud detection in which the data does not leave an organization’s boundaries. Instead, third-parties submit fraud detection algorithms based on existing techniques—including domain knowledge, public sources of data, and synthetic data—which are evaluated and ranked by the data holder. These third parties may be motivated by financial and/or reputational rewards for the winning algorithms and developers, based on a leaderboard—a model that has seen significant success in the machine learning community [19]. We study this setting under two key constraints that (together) have not been explored in the literature:

- (1) *Private algorithm evaluation*: We observe that if the accuracy of a fraud detector is released directly, it can leak sensitive information about the underlying test data (Section 2). We therefore consider methods for evaluating algorithms, and releasing their results, under a differential privacy (DP) constraint [9].
- (2) *Graph-structured data*: Many prominent fraud domains, such as financial fraud or product review fraud, have graph-structured datasets. We focus on fraud detection algorithms (and privacy solutions) that can be applied to graph-structured data.

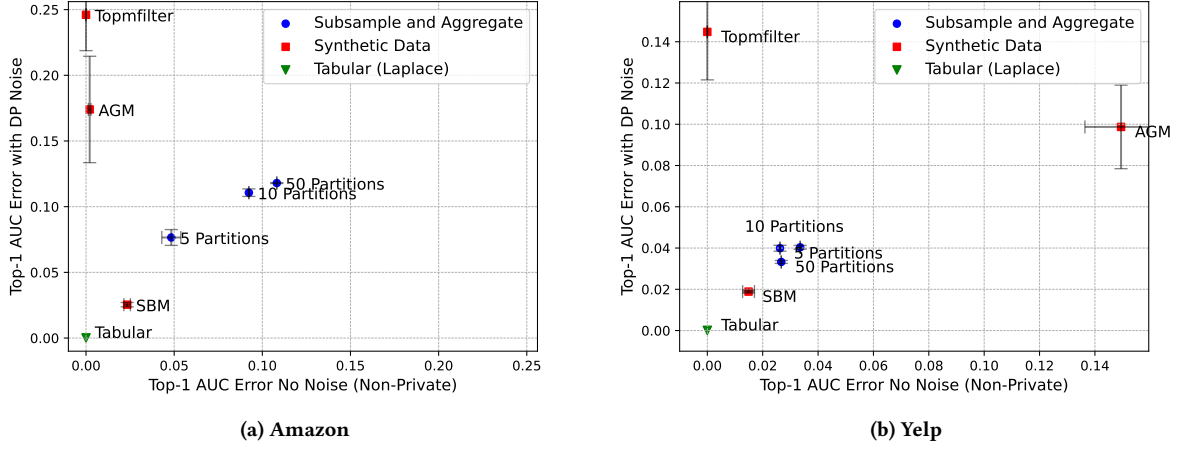
Existing techniques for privately evaluating fraud-detection algorithms cannot be easily applied to graph-structured data (Section 2.2). The main challenge is that existing fraud detection algorithms for graph-structured data rely on queries over the graph with a high global sensitivity, meaning that the query result can change significantly if even a single node’s neighbors are altered in the graph (Definition 2.2). When such high-sensitivity algorithms are combined with existing DP mechanisms, large amounts of noise

This work is licensed under the Creative Commons Attribution 4.0 International License. To view a copy of this license visit <https://creativecommons.org/licenses/by/4.0/> or send a letter to Creative Commons, PO Box 1866, Mountain View, CA 94042, USA.



KDD’25 1(1), 1–16

© 2025 Copyright held by the owner/author(s).



**Figure 1: Comparison of DP benchmarking methods for releasing the best AUC score among 10 fraud detectors with privacy budget of  $\epsilon = 5.0$ . The horizontal axis captures error due to inductive bias (i.e. the underlying graph model, without DP noise); the vertical axis captures error including DP noise. More complex synthetic data methods (Topmfilter and AGM) can model the data well on some datasets without privacy, but suffer from high variance due to noise addition required to satisfy DP, thus undermining utility. Subsample-and-aggregate tends to distort graph structure extensively, even before adding random noise to outputs. All current methods to satisfy DP on graph data incur large utility cost compared to tabular data. Error bars show the standard error of the mean across 10 simulations of each method.**

are required (scaling proportionally to the global sensitivity), thus destroying the utility of the query.

The goal of this work is to instantiate and benchmark different classes of techniques for evaluating fraud detection algorithms over graph-structured data under a differential privacy constraint. Specifically, we evaluate two categories of techniques for dealing with high-sensitivity queries: (1) *Subsample-and-aggregate* partitions the dataset into multiple non-overlapping datasets, then evaluates the fraud detection algorithm over each partition. The average accuracy over the partitions is low-sensitivity, and can be released with less noise than without partitioning. (2) *Synthetic graph data* involves generating a DP synthetic copy of the true graph; then, fraud detection algorithms can be evaluated on this synthetic graph.

In this work, we provide an extensive empirical evaluation of these two categories of approaches for privately evaluating fraud detectors on graph data. We outline the challenges associated with using them, recommendations for when to use which algorithm, and open problems. Our primary contributions are:

- (1) We formulate the problem of *differentially private benchmarking of fraud detectors on private graph data*. We identify challenges unique to graph data that make the model evaluation problem more difficult on graph data than on tabular (non-graph) data. We provide a quantitative and qualitative head-to-head comparison between two frameworks for this problem—subsample-and-aggregate and synthetic graph algorithms.
- (2) Across methods, **we observe a severe trade-off between bias introduced by distorting the graph and noise required to compensate for computing high sensitivity statistics on the graph**. This result is captured in Figure 1, which shows the error in privately benchmarking the best AUC score among

a set of 10 fraud detectors on two fraud datasets. We plot the error of each DP benchmarking method without noise added (inductive bias) on the x-axis against error after adding noise to ensure differential privacy on the y-axis. Among synthetic data methods, more complex methods (TopmFilter and AGM) have lower inductive bias, but much higher noise addition to preserve privacy than the simpler SBM. Subsample-and-aggregate tends to distort graph structure extensively, even before adding random noise to outputs. All current methods to satisfy DP on graph data incur large utility cost compared to tabular data, where it is possible to simultaneously achieve low noise and unbiased estimates.

- (3) To explain these results, we conduct detailed ablations on both subsample-and aggregate and synthetic data methods. While these methods introduce inductive bias in very different ways, we find that both exhibit a similar trade-off—the less we bias our graph representation, the more noise we must add to satisfy DP. We propose a variant of subsample-and-aggregate for fraud graphs where benign and fraudulent vertices are sub-sampled at different rates, which can improve utility at the cost of a weaker privacy guarantee (Section 4). Our extensive empirical analyses suggest that it is helpful to use far fewer partitions than we might expect from applications of subsample-and-aggregate to tabular data (under 50 partitions for  $\epsilon = 0.5$ , compared to over 100 partitions in prior work [31, 34]).

## 2 Problem Formulation

We consider a *benchmarking server*, which has a private graph  $G$  consisting of a set of known fraudulent vertices  $V_1$  (of size  $n_1$ ) and a set of benign vertices  $V_0$  (of size  $n_0$ ).

The benchmarking server's overall goal is to *evaluate* one or more fraud detection algorithms, provided by third parties, and *communicate* the result back to the algorithm designers. We assume the benchmarking server receives a fraud detection algorithm  $\mathcal{A}$ . The fraud detection algorithm takes as input a vertex  $v$  and the entire graph  $G$  and outputs  $\mathcal{A}(G, v)$  which is a numerical score where a higher score indicates a higher likelihood of fraud. For example, the fraud detection algorithm could score a vertex by its degree. The benchmarking server returns an *accuracy statistic* for the fraud detection algorithm on graph  $G$ . Concretely, we consider the *AUC score*, which is defined as:

$$f_{\text{AUC}}(\mathcal{A}, G) = \frac{1}{n_1 n_0} \sum_{v_0 \in V_0} \sum_{v_1 \in V_1} \mathbb{1}[\mathcal{A}(G, v_1) > \mathcal{A}(G, v_0)].$$

The AUC score represents the probability that a randomly chosen fraudulent vertex is scored higher than a randomly chosen benign vertex. It is a commonly used accuracy statistic for class-imbalanced binary classification problems like fraud detection [12].

We consider three different operating modes for the benchmarking server: (1) *one-shot*: the benchmarking server releases the AUC score for a single submitted fraud detector, (2) *full leaderboard*: the benchmarking server returns the AUC score for of a set of submitted fraud detectors, and (3) *top-1 release*: the benchmarking server releases the best-performing fraud detector among a set of submitted algorithms.

*An example attack.* We next show how a malicious actor can compromise the privacy of any individual in the system, even with only a single benchmarking query. Consider a bad actor who wishes to ask whether there exists an edge between two specific vertices in the graph. The adversary needs three capabilities. (1) *An accurate fraud detector*: for example, a known algorithm from the literature which does better than random ( $\text{AUC} > 0.5$ ). (2) *An inaccurate fraud detector*: for example, scoring vertices at random (expected  $\text{AUC} = 0.5$ ). (3) *The ability to identify vertices in  $G$* : this depends on what information the private server gives to the fraud detection algorithm. In many cases,  $G$  may include extensive metadata per vertex, which makes it easy to identify vertices. Even without metadata, there are many de-anonymization attacks leveraging only the graph structure (see [17] for a survey), which enable an adversary to identify vertices.

The adversary can identify the relevant pair of vertices, and then run the accurate fraud detector if an edge exists between the vertices or the inaccurate fraud detector otherwise. If the adversary observes a high AUC score, they learn that an edge exists, while if they observe a low AUC score they learn that the edge does not exist. In effect, the benchmarking server allows a malicious actor to answer any binary query on  $G$  by encoding the answer to their query as either a high-accuracy or low-accuracy fraud detector.

## 2.1 Incorporating Differential Privacy

Motivated by this privacy risk, we propose incorporating differential privacy (DP) [9] into the fraud benchmarking server. Specifically, we want the server to guarantee a relaxation of DP, which promises differential privacy only for benign vertices:

**Definition 2.1** (Protected differential privacy [23]). Two graphs  $G, G'$  are neighboring if:

- (a)  $G$  and  $G'$  share the same partitions of fraudulent and non-fraudulent vertices  $V_1$  and  $V_0$ .
- (b)  $G$  can be obtained from  $G'$  by rewiring the edges of one *benign* vertex and/or changing that vertices' metadata.

Let  $f$  denote the benchmarking server that given a graph and fraud detector outputs an estimate of the AUC score. The server  $f$  satisfies  $\epsilon$ -protected differential privacy if for any two neighboring graphs  $G, G'$ , any fraud detection algorithm  $\mathcal{A}$  and any possible set of outputs  $\mathcal{O}$ :

$$\Pr[f(G, \mathcal{A}) \in \mathcal{O}] \leq e^\epsilon \Pr[f(G', \mathcal{A}) \in \mathcal{O}].$$

Standard DP allows  $G$  and  $G'$  to differ in the data of *any* vertex in the graph, not just a benign vertex. Protected DP is a relaxation of standard DP in that any graphs that are neighbors per the definition of protected DP are also neighbors under standard DP. We will refer to protected differential privacy as DP for brevity throughout.

We primarily adopt this relaxed notion of privacy to improve utility. In many real-world graphs, the rate of fraud is low. Hence, requiring that the released accuracy statistic does not change much if we change the connections of these fraudulent vertices makes it difficult to release high-fidelity benchmarks. Still, we believe this relaxation is useful. In fraud detection, it is natural to hold different privacy expectations for fraudulent participants (many of which may even be fake [15]) compared to legitimate ones.

In the definition of neighboring graphs, we adopt the strong notion of *node* differential privacy, which protects all of the edges of any single benign vertex. Many prior works employ a weaker notion of *edge* differential privacy [4, 18, 21, 32], which defines neighboring graphs as graphs that differ in a single edge. We note that protected DP inherits the *composition property* of standard DP:

**THEOREM 2.1** (COMPOSITION [9]). *For any two fraud detectors  $\mathcal{A}_1$  and  $\mathcal{A}_2$ , if releasing  $f(\mathcal{A}_1, G)$  satisfies  $\epsilon_1$ -protected DP and releasing  $f(\mathcal{A}_2, G)$  satisfies  $\epsilon_2$ -protected DP, then releasing both results on graph  $G$ ,  $(f(\mathcal{A}_1, G), f(\mathcal{A}_2, G))$  satisfies  $(\epsilon_1 + \epsilon_2)$ -DP.*

This property is helpful in moving from one-shot release of fraud detectors to releasing a leaderboard of many fraud detectors.

## 2.2 Challenges of Graph Data

Even evaluating a single fraud detector on graph data proves challenging under DP constraints. To understand why, let us compare our setting to evaluating a fraud detector on tabular data. Evaluating a single fraud detector can be seen as a problem of releasing a (noisy) query result. A simple mechanism that solves the query release problem adds random noise with variance scaled to the "sensitivity" of this query, which is defined as follows.

**Definition 2.2** (Global Sensitivity). For a query  $f : \mathcal{X} \rightarrow \mathbb{R}^d$ , define its *global sensitivity*

$$\Delta_f = \max_{G, G' \text{ neighbors}} \|f(G) - f(G')\|_1$$

as the worst-case change in  $f$  across any two neighboring graphs.

Then, a canonical mechanism, termed the Laplace Mechanism, scales noise to the global sensitivity:

**Definition 2.3** (Laplace Mechanism [9]). On any input  $G$  the *Laplace Mechanism* with privacy parameter  $\epsilon$  releases

$$\tilde{f}(G) = f(G) + \text{Laplace}(\Delta_f/\epsilon).$$

The Laplace Mechanism satisfies  $\epsilon$ -DP.

In the tabular setting, model evaluation is a low sensitivity query and therefore can be released by directly applying the Laplace mechanism. Consider a simple case where fraud detector  $\mathcal{A}$  is a fitted logistic regression model (the weights of the model are fixed). Changing any row of a tabular dataset only changes the features of that row and hence changes at most a single fraud prediction score. Therefore, when evaluating the fixed model on tabular data, the AUC score can only change by  $\frac{1}{n_0}$ . The Laplace mechanism can then release the true AUC score of the fraud detector plus Laplace noise with variance  $\frac{2}{(\epsilon n_0)^2}$ .

In contrast, consider evaluating the logistic regression model on a graph where features of each vertex include graph statistics like the degree of each vertex. Because features of each vertex depend on other vertices, changing any one vertex can change features of all other vertices in the graph. In the worst-case, changing a vertex changes fraud prediction scores for all other vertices in the graph, so the AUC score has a large global sensitivity of 1. As this is the largest possible value AUC can take, the Laplace mechanism must add so much noise that the entire signal is lost.

In this paper, we focus on addressing this challenge of high sensitivity of model evaluation on graph data. In cases where queries of a dataset have large worst-case sensitivity there are three classes of solutions in the DP literature:

- (1) (*Subsample-and-aggregate*) Force low sensitivity of the AUC score by applying “subsample-and-aggregate.”
- (2) (*Synthetic data*) Generate DP synthetic data that captures some structure of the private graph and run fraud benchmarking on this private graph data.
- (3) (*Calibrate noise to “local sensitivity.”*) Estimate (an upper bound) on how sensitive  $f_{\text{AUC}}$  is on the specific graph and fraud detection algorithm  $\mathcal{A}$  and calibrate noise to this sensitivity, which may be much lower than the worst case global sensitivity. This approach includes mechanisms like Propose-Test-Release, Smooth Sensitivity, and the Inverse Sensitivity Mechanism [10, 33] as well as recent work on privatizing black-box scripts run on private data [24].

In this work, we give instantiations of subsample-and-aggregate and synthetic data generation algorithms tailored to the benchmarking server setting and run extensive empirical evaluations to understand opportunities and shortcomings. We do not evaluate local sensitivity based methods [33], because these approaches are computationally infeasible in our setting as they would require enumerating every possible neighboring graphs and evaluating fraud detectors on these graphs to estimate a bound on local sensitivity.

### 3 Related Work

To our knowledge, this work is the first to consider the problem of model evaluation on graph data under differential privacy constraints. For tabular (non-graph) data, there are two lines of work that consider DP model evaluation. One line of work [2, 38, 45] proposes a framework of “verification servers” wherein analysts

fit a model of data (e.g., a linear regression model) on a synthetic dataset and then employ a “verification server” which holds non-synthetic data to perform quality checks that their model is useful like goodness-of-fit tests. While the system design of our work is similar, these works focus on tabular data rather than graph data, which poses specific challenges as we detail in Section 2.2.

A recent line of work in DP machine learning (starting with [30] and extended in [5, 36]), looks at a closely related problem of model selection under differential privacy constraints. These works focus on choosing the (nearly) optimal model in minimizing loss among a large set of models without paying for privacy loss that grows with the number of models. In our work, we observe that on graphs (even ignoring model training) the seemingly straightforward step of one-shot model evaluation is difficult under differential privacy constraints.

A number of works have considered the problem of running arbitrary queries on private data. Subsample-and-aggregate, first proposed in [33], is one popular method for reducing the sensitivity of a query. Practical instantiations of subsample-and-aggregate have been used in popular frameworks for data analysis as in GUPT [31] and for training ML models as in PATE [35]. In our work, we focus on model evaluation rather than training, as the model evaluation task is quite challenging in the graph setting. In contrast, to prior uses of subsample-and-aggregate, we propose up-sampling fraudulent entities satisfying a relaxed notion of privacy and improving utility. We then perform extensive empirical evaluation to understand how the subsample-and-aggregate framework compares to synthetic data generation algorithms for this problem. A recent work [24] also considers the problem of running arbitrary code on a private dataset and gives a new mechanism called TAHOE that is competitive with subsample-and-aggregate in some tabular data settings, but TAHOE is computationally expensive and cannot run efficiently on graph data.

There is a long and rich line of work in differentially private analysis of graph data. We discuss the literature on generating synthetic graphs in more detail in Section 5 where we detail our choice of synthetic graph algorithms to benchmark. These synthetic graph algorithms require estimating statistics of the graph (like degree distribution or number of triangles) under node differential privacy. This introduces new challenges as the prior works on synthetic data generation use the weaker notion of edge-DP to estimate statistics. In our work, we generically transform edge-DP estimation into (reasonably accurate) node-DP estimation using the idea of smoothly projecting a graph to the space of limited-degree graphs from [3, 22]. There may be additional improvements in applying existing synthetic data generation methods under the node-DP privacy regime by applying more tailored estimation procedures for specific graph statistics.

### 4 Subsample-and-Aggregate

The first approach we consider in privatizing fraud benchmarking is the subsample-and-aggregate framework [33]. Recall from Section 2.2 that a key challenge of releasing a DP estimate of the AUC score of a fraud detector on a graph is that this query has global sensitivity of 1, equal to the range of the AUC score. Subsample-and-aggregate forces low sensitivity of the query by first partitioning

the dataset into  $k$  disjoint sets and then estimating AUC on each partition.

Our algorithm follows the template described above for benign vertices, that is, we partition the benign vertices into  $k$  disjoint sets of equal size. However, in fraud graphs, there are often very few fraudulent vertices. For example, in the Elliptic Bitcoin financial fraud dataset [42] there are only 11 fraudsters out of over 6,000 vertices. Partitioning these fraud vertices into a reasonable number of partitions to achieve low sensitivity (say  $k \geq 5$ ) would destroy any structure of the sub-graph of fraud vertices.

We therefore modify typical subsample-and-aggregate for the fraudulent vertices by allowing *duplication* of fraudsters across partitions. For each partition, we sample a subset of fraudulent vertices, where the rate of sub-sampling is controlled by a parameter  $\rho$ . We term this instance of subsample-and-aggregate as Partition, Duplicate, and Aggregate (PDA), described in Algorithm 1. Note that taking  $\rho = 1$  results in duplicating all fraud vertices in each partition, while taking  $\rho = \frac{1}{k}$  is similar to typical subsample-and-aggregate, but with the difference that fraudulent vertices may be sampled into multiple partitions.

---

**Algorithm 1** Partition, Duplicate, and Aggregate
 

---

**Parameters:** privacy parameter  $\epsilon > 0$ , number of partitions  $k$ , fraud sub-sampling rate  $\rho$ .

**Inputs:** fraud detector  $\mathcal{A}$ , accuracy statistic  $f$  with global sensitivity  $\Delta$ , fraud vertices  $V_1$ , benign vertices  $V_0$ , graph  $G$  on vertex set  $V_0 \cup V_1$ .

- Randomly partition non-fraud nodes  $V_0$  into  $k$  equally size sets  $V_0^{(1)}, \dots, V_0^{(k)}$ .
  - Randomly sample  $k$  sets of fraud nodes  $V_1^{(1)}, \dots, V_1^{(k)}$  where each  $V_1^{(i)}$  is sampled independently uniformly from all sub-sets of  $V_1$  of size  $\rho \cdot |V_1|$ .
  - Let  $G_1, \dots, G_k$  be sub-graphs of  $G$  on vertices  $(V_0^{(1)} \cup V_1^{(1)}), \dots, (V_0^{(k)} \cup V_1^{(k)})$ .
  - Release  $Z + \frac{1}{k} \sum_{i=1}^k f(\mathcal{A}, G_i)$  where  $Z \sim \text{Laplace}(\Delta/(k\epsilon))$ .
- 

It is straightforward to prove that Algorithm 1 guarantees differential privacy:

**PROPOSITION 4.1.** *For any choice of sub-sampling rate  $\rho \in (0, 1)$ , number of partitions  $k > 1$  and privacy parameter  $\epsilon > 0$ , Algorithm 1 guarantees  $\epsilon$ -Protected Differential Privacy (Definition 2.1).*

The proof follows from a standard proof of privacy for subsample-and-aggregate: changing the data of any benign vertex impacts at most 1 of the  $k$  partitions between any two neighboring graphs, and the accuracy score on this partition can change by at most 1 since  $f$  has global sensitivity (Definition 2.2) of 1. Hence, the mean across partitions has global sensitivity of  $\frac{1}{k}$  and privacy follows from the Laplace mechanism (Definition 2.3). We note that in practice, the subsample-and-aggregate framework does not introduce substantial computational overhead. We further discuss the computational costs of subsample-and-aggregate in Appendix B.2. We further develop intuition around the error of partition, duplicate, and aggregate as a function of number of partitions and subsample-rate via a simple example in Appendix B.1.

## 4.1 Running Multiple Benchmarks

Algorithm 1 provides a method for one-shot release of the AUC score of a single fraud detector. In order to apply it to full leaderboard release, we can invoke composition (Theorem 2.1) and subdivide the privacy budget among many fraud detectors. For example, if we have 10 algorithms to benchmark, we run each with privacy budget of  $\epsilon/10$  per fraud detector. As it is harder to provide good utility for smaller  $\epsilon$ , we expect our accuracy of estimation to degrade in the number of detectors benchmarked.

In many real-world settings it is useful to release only the best or the top- $m$  fraud detectors, for example when running a competition. In the case of top-1 release, we can use the *Report Noisy Arg Max* mechanism [11]. This mechanism adds Laplace noise to any (finite) number of queries as per the Laplace mechanism, but then only releases the name of the query with the largest (noisy) value. Rather than paying composition cost that grows in the number of queries, this procedure is  $\epsilon$ -DP. In our case, then, we can apply Algorithm 1 to arbitrarily many fraud detectors and then at the end only publicly release the name of the detector with the highest noisy AUC score. This guarantees  $\epsilon$ -DP when each run of Algorithm 1 is run using privacy parameter  $\epsilon$ . While we do not experiment with releasing the top- $m$  fraud detectors, recent work [37] shows that releasing the top- $m$  fraud detectors ranked by noisy AUC (among a larger set of fraud detectors) only incurs total privacy loss of  $m\epsilon$ .

## 5 Synthetic Data Generation

In this section, we describe our choice of synthetic graph data generation algorithms to benchmark; the surveys [14, 27] provide a useful overview of such algorithms. Many methods do not handle *labeled* vertices. Such methods cannot be applied to our problem, as synthetic data for fraud detection benchmarking needs to differentiate between fraudulent and benign vertices. Additionally, most existing work focuses on satisfying the weaker notion of edge-level DP, while we wish to satisfy node-level DP. Therefore, we focus on the following 3 methods that all handle labelled vertices and are amenable to transformation into a node-level DP algorithm:

- (1) *Stochastic block model (SBM)*: Estimate a stochastic block model with two communities (fraud and non-fraud) and sample a graph based on the SBM parameters. For a fixed number of benign and fraud vertices, the stochastic block model has three parameters  $p_1, p_0, p_{01}$ . Each edge in the graph is sampled independently at random with probability  $p_1$  if both of its endpoints are fraudulent,  $p_0$  if both are benign, and  $p_{01}$  if one is fraudulent and the other is benign.
- (2) *Attributed social graph (ASG)*: [18]: Estimate the connection probabilities with and between fraud and non-fraud vertices (as in the SBM), but additionally estimate number of triangles in the graph and the degree sequence of the graph. Then, sample a graph that matches these noisy statistics. We run two versions of this method, with and without the triangle statistic.
- (3) *Top- $m$ -filter* [32]: Directly perturb the adjacency matrix of the graph. In particular, flip each edge in the graph and then perform a filtering step to remove edges to match a noisy estimate of total number of edges.

In general, synthetic data methods first compute graph statistics under differential privacy, which provide a succinct representation

of the graph, and then generate the synthetic graph based on these (noisy) statistics. More expressive graph models may better represent the graph structure, but tend to require the estimation of noisier sufficient statistics due to differential privacy. We choose methods that lie along this spectrum of model complexity.

Other popular synthetic data methods in the literature use exponential random graph models (ERGMs) and more recently graph neural networks (GNNs) [44, 47]. These methods are either computationally intractable for large graphs in the case of ERGMs or would require too much noise addition to sufficient statistics under node-level privacy in the case of GNNs. We discuss these methods in more detail in Appendix C.

*Guaranteeing Node-Level Differential Privacy.* The algorithms we consider were designed to provide edge-level differential privacy. In privately computing sufficient statistics of the graph, these algorithms add Laplace noise proportional to the worst-case sensitivity of a statistic to the change of a single edge in a graph. In order to guarantee node-level privacy in this noise addition step, we use the idea of projecting the graph to the space of graphs with bounded maximum degree from [3, 22] and then adding noise proportional to this “restricted sensitivity.” For a given graph  $G$ , choice of truncation threshold  $T$ , and graph statistic  $g$  the full workflow is:

- (1) (Naive truncation). Truncate graph  $G$  by removing all vertices with degree above  $D$ .
- (2) Estimate the “smooth sensitivity”  $S$  of the naive truncation operation per [22].<sup>1</sup>
- (3) Add Laplace noise with scale proportional to  $S \cdot RS_T(g)$  where  $RS_T(g)$  represents the “restricted sensitivity” of  $g$  on graphs of max degree  $D$ , that is the maximum change in  $g$  between any two node-adjacent graphs of max degree  $T$ .

We summarize the framework for node-private synthetic data release in Algorithm 2. Since the max degree and average degree of the fraud graphs used (see Table 1) tends to be much smaller than the number of vertices in the graph, the restricted sensitivity tends to be much lower than the global sensitivity.

Note that using this method with Laplace noise actually guarantees the relaxation of  $(\epsilon, \delta)$ -differential privacy due to the use of “smooth sensitivity” [33]. We fix  $\delta$  to  $10^{-8}$  for all experiments on synthetic data methods. Additionally, to provide a fair comparison against our subsample and aggregate method which relaxes privacy for fraudulent vertices, we compute statistics that rely only on the fraudulent nodes without noise.

## 6 Experimental Setup

In this section, we describe the datasets, fraud detectors, and metrics used throughout our experiments.

### 6.1 Datasets

We test methods for fraud benchmarking on 4 datasets representing a variety of domains and graph structures. All graphs are undirected unipartite graphs. In *Yelp* [8] and *Amazon* [8] each vertex represents

<sup>1</sup>From [22], Proposition 6.1 we can compute the smooth sensitivity  $S_{trunc}^\beta(G, T)$  of the truncation operation as follows. Let  $N_t(G, T)$  denote the number of benign vertices with degree in range  $[T - t, T + t + 1]$  and  $C_t(G, T) = 1 + t + N_t(G, T)$ . Then,  $S_{trunc}^\beta(G, T) = \max_{t \geq 0} e^{-\beta t} C_t(G, T)$ .

---

### Algorithm 2 Framework for Node-Private Synthetic Data Release

---

**Parameters:** privacy parameters  $\epsilon > 0, \delta \in (0, 1)$ , degree threshold  $T$

- Inputs:** fraud vertices  $V_1$ , benign vertices  $V_0$ , graph  $G$  on vertex set  $V_0 \cup V_1$ , vector of sufficient statistics to compute  $g(G)$  with restricted sensitivity  $\Delta_T$ .
- Remove all benign vertices from  $G$  with degree greater than  $T$ .
  - Compute the  $\beta$ -smooth sensitivity  $S_{trunc}^\beta(G, T)$  of the truncation operation on  $G$ , where  $\beta = -\frac{2\epsilon}{\log(1/2\delta)}$ .
  - Release  $\tilde{g}(G) = g(G) + Z$  where  $Z \sim \text{Laplace}(2S_{trunc}^\beta(G, T) \cdot \Delta_T / \epsilon)$ .
  - Sample output synthetic graph  $\tilde{G}$  based on  $\tilde{g}(G)$ .
- 

a reviewer with edges denoting common reviews on the product/s/restaurants and fraudulent reviewers represent spammers and low rated reviewers respectively. *Peer Review* consists of paper reviewers at a computer science conference with edges denoting mutual bids on each other’s papers [43]. Following [15] we inject a clique of 22 fraudulent reviewers with edge density of 0.8 among these reviewers into the graph, which corresponds to the smallest injected clique that was possible to detect in prior work. Finally, in *Elliptic* [42] each vertex in the graph represents a transaction from the Bitcoin blockchain, an edge represents a flow of Bitcoins between one transaction and the other, and fraudulent nodes are illicit transactions. We take a single time-step from the entire Elliptic graph. We give summary statistics of these datasets in Table 1 in Appendix A.

We run analyses of subsample-and-aggregate and synthetic data algorithms on validation datasets to understand settings of hyperparameters before comparing these methods against each other. We use four validation datasets. For *Yelp*, we use random split of the vertices with 11k vertices in the test set and 11k in the validation set. For *Elliptic*, we use different disjoint time periods for validation and test. For *Amazon* and *Peer Review*, there are not standard train-test splits used in past work. We therefore use the entire graph for evaluation, and generate validation graphs to set hyperparameters by estimating parameters of a stochastic block model (SBM) and sampling from this model.

### 6.2 Fraud Detectors

We benchmark 10 fraud detectors that capture a variety of graph structures used to detect fraud. While many state of the art algorithms incorporate a learning component where they learn parameters from data, we focus on simple approaches that do not require learning. Since even one-shot model evaluation is challenging on graph data (Section 2.2), this allows us to focus on the evaluation part of the problem.

Specifically we evaluate the following fraud detectors:

- (*Negative*) *Degree* [16]: rank by the degree of each vertex (in ascending or descending order).
- (*Negative*) *Clustering Coefficient*: rank by the clustering coefficient of each vertex (in ascending or descending order), inspired by [1].



- **SVD Error [16]:** take the singular value decomposition of the adjacency matrix obtain a low rank approximation (for specified rank  $r$ ). Then, rank each vertex by reconstruction error (aggregating over edges by taking either the sum or the max over edges). We use  $r = 10$  for the sum aggregation and  $r = 50$  for the max aggregation, chosen to maximize average AUC across all datasets in a grid search over choices of  $r$ .
- **Community Detection:** run Leiden community detection [40] to place each vertex in a cluster and then rank by the size of the cluster (with larger clusters less likely to be fraudulent).
- **Aggregations:** take weighted averages of the (normalized) scores or maximum scores obtained from subsets of the prior methods.

These algorithms give a wide range of AUC scores on each dataset. For example, on Yelp, Neg Degree performs the best with an AUC score of 0.69 and SVD Error (Sum) performs poorly with an AUC score of 0.34. In contrast, on Peer Review, SVD Error (Sum) performs the best with an AUC score of 0.88, while Neg Degree has very bad performance with AUC of 0.12.

### 6.3 Measuring Utility

We consider three metrics to compare utility across methods. Each metric corresponds to one of the release modes for the benchmarking server: one-shot, full leaderboard and top-1 release. Let  $\{\mathcal{A}_i\}_{i=1}^m$  denote a set of  $m$  fraud detectors to benchmark on graph  $G$ ,  $f_{\text{AUC}}(\mathcal{A}_i, G)$  denote the true AUC score for fraud detector  $i$  on  $G$  and  $\tilde{f}_{\text{AUC}}(\mathcal{A}_i, G)$  denote the noisy DP estimate of the AUC score. For the one-shot release, where we wish to release AUC for a single fraud detector, we calculate  $L1$  error:  $|f_{\text{AUC}}(\mathcal{A}_i, G) - \tilde{f}_{\text{AUC}}(\mathcal{A}_i, G)|$ .

Finally, when evaluating top-1 release of only the best fraud detector among a set of fraud detectors we measure utility by the distance between the true AUC of the true best fraud detector (computed without any privacy) and the true AUC of the released best fraud detector. That is, we define top-1 error as:

$$f_{\text{AUC}}(\mathcal{A}_{\text{top}}, G) - f_{\text{AUC}}(\mathcal{A}_{\text{top}'}, G) \text{ where } \text{top} = \sigma^{-1}(1), \text{top}' = \tilde{\sigma}^{-1}(1).$$

For the full leaderboard setting we use the weighted Kendall-Tau distance between rankings [26]. We discuss this metric and provide results in this setting in Appendix D.2.

In addition to using the AUC score as an accuracy metric, we evaluate the F1 score of fraud detectors, another popular measure of fraud detector accuracy. We find similar results to that of AUC score, but F1 score tends to be even more difficult to release accurately. We present these additional results in Appendix D.

## 7 Experimental Results

In this section we detail the results of our head-to-head comparison of Subsample-and-Aggregate and Synthetic Graph Generation (7.1). We then describe specific experiments to better understand trade-offs between distorting graph structure and adding noise to preserve privacy for Subsample-and-Aggregate (7.2) and Synthetic Graph Generation (7.3) respectively.

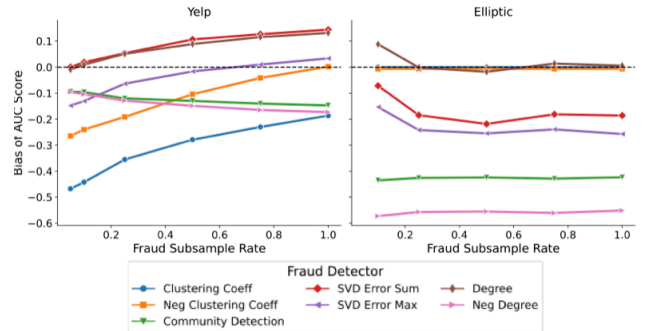
### 7.1 Comparison of Algorithms

We benchmark subsample-and-aggregate against synthetic data algorithms for the concrete task of releasing the best fraud detectors among a set of fraud detectors. We choose parameters of

subsample and aggregate (number of partitions and sub-sampling rate) based on the best parameters for each dataset in releasing a ranking of all fraud detectors on the validation dataset. This follows prior work [31], which assumes that there exist public datasets on which one could reasonably set hyperparameters of subsample and aggregate.

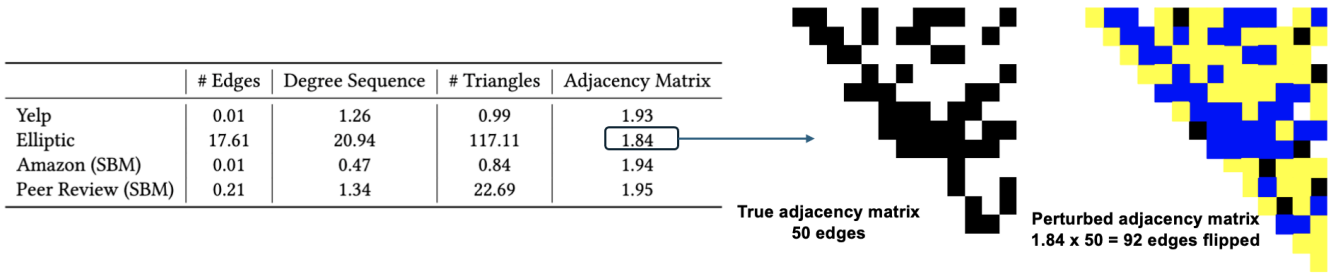
In Figure 1, we show results of the DP benchmarking methods for top-1 release. In order to decompose the error into inductive bias due from how a given method distorts graph structure and error arising from addition of privacy-preserving random noise, we plot each method without privacy-preserving noise on the x-axis and with noise needed to preserve privacy on the y-axis. Specifically, for Non-Private subsample-and-aggregate we only apply the graph partitioning and do not add Laplace noise to the AUC score. For non-private synthetic data methods we compute sufficient statistics for each method without any noise addition and then generate a graph using those sufficient statistics. Among synthetic data methods, Topmfilter has no error without privacy as it releases the full adjacency matrix, while SBM and AGM introduce error even without privacy. However, after adding noise needed for privacy SBM performs the best among synthetic data methods. Subsample-and-aggregate distorts graph structure extensively, even with only 5 partitions, resulting in high error without privacy. We provide additional results for larger privacy budget of  $\epsilon = 5.0$  and other datasets in Appendix D.1, but the general trends are similar.

### 7.2 Subsample-and-Aggregate



**Figure 2: Bias to the AUC score introduced by subsample and aggregate for each fraud detector varying the fraud subsample rate ( $\rho$ ) while fixing number of partitions  $k = 20$ . Subsample-and-aggregate introduces extensive bias to all fraud detectors, with the sign and magnitude of the bias varying widely across fraud detectors.**

In experiments on four validation datasets, we seek to understand how the parameters of the algorithm—number of partitions  $k$  and rate of sub-sampling fraud in each partition  $\rho$ —impact the bias,



**Figure 3: Normalized mean absolute error (MAE) introduced to each of the synthetic graph sufficient statistics. Fixed  $\epsilon = 5.0$  and degree cutoff of 1 times the graph’s max degree. It is possible to estimate SBM parameters accurately, while other parameters have large noise addition. We give an example of what relative error of 1.84 means for the adjacency matrix on the right, where an adjacency matrix with 50 edges has 92 edges flipped, 42 removed (yellow) and 50 added (blue.)**

variance from Laplace noise addition and overall distortion of fraud detector rankings.

On each dataset we run Algorithm 1 for 10 trials for each choice of parameters  $k$ ,  $\rho$  and  $\epsilon$ . We show per-dataset results on the Yelp and Elliptic validation datasets in this section. We provide additional results on the validation datasets for Amazon and Peer Review generated by modeling each of the real graphs using SBMs in Appendix D.3.

In general, we find that partitioning the graph into random sub-graphs introduces significant bias to estimates of graph statistics. The sign and magnitude of this bias can differ widely across fraud detectors. In Figure 2, we show bias per fraud detector fixing the number of partitions at  $k = 20$  and varying the fraud sub-sampling rate. We find that it is not always possible to achieve zero bias for a given fraud detector for a given number of partitions  $k = 20$ . For instance, the clustering coefficient detector has negative bias on the Yelp dataset at all values of  $\rho$ . This makes sense as removing benign vertices from the graph may change the distribution of fraud detection scores for benign vertices such that it is not possible to recover a similar distribution at any rate of sub-sampling fraudulent vertices. We additionally find, as expected, that the magnitude of bias increases with the number of partitions ( $k$ ) although the sign of the bias differs across fraud detectors. We plot bias as a function of number of partitions in Figure 14 of Appendix D.3. These results explain the poor performance of subsample-and-aggregate in Figure 1, as subsampling tends to distort graph structure extensively, biasing different fraud detectors in different ways thereby undermining the utility of the ranking of fraud detectors.

### 7.3 Synthetic Graph Generation

In our experiments we aim to isolate error introduced due to choice of graph model and noisy estimation of sufficient statistics. For each synthetic data generation algorithm, we generate 10 synthetic data sets. For each synthetic graph method, we subdivide the privacy budget evenly between the different parameters to estimate. We note that it may be possible to better distribute privacy budget between different statistics, which is an interesting area for future investigation. We test degree truncation thresholds as a function of the max degree of each graph, so 1.0 is a threshold exactly equal to

the maximum degree benign vertex in a graph while 0.5 removes all nodes with degree  $> 0.5$  times the max degree. In this section, we report results with threshold of 1 and give additional results for 0.5 in Appendix D.4.

We find that outside of the SBM, which has edge counts as sufficient statistics, it is necessary to introduce large distortion to the sufficient statistics of each graph model in order to preserve privacy, as shown in Figure 3. On Yelp, Amazon, and Peer Review it is possible to estimate the edge count sufficient statistics for the SBM with high accuracy at  $\epsilon = 5.0$ , perturbing the edge count by 1% the total number of edges on Yelp and Amazon. Elliptic is an extremely sparse graph (0.04%), so we introduce much larger relative error. For degree sequence and number of triangles, the amount of error is one to two orders of magnitude larger, with error generally at least 50% of the value of the original statistic. Unsurprisingly, the adjacency matrix cannot be accurately estimated under node-DP via direct noise addition. We highlight the amount of noise addition needed to preserve privacy in a simple example of relative error of 1.84 on a 15x15 adjacency matrix, shown in Figure 3. Even after aggressively truncating high-degree nodes, the addition of DP noise results in flipping the same number of edges as were originally in the adjacency matrix. This large distortion of sufficient statistics explains the poor accuracy of AGM and TopMfilter (which require estimates of the degree sequence + triangles and the adjacency matrix respectively.)

## 8 Discussion

In this work we define the novel problem of privately benchmarking fraud detectors on graph-structured data. We benchmark two popular frameworks from the DP literature, subsample-and-aggregate and synthetic data generation on this problem. We characterize a trade-off for each method between error arising from bias due to distorting graph structure and error arising from privacy-preserving noise addition. Generally, our results suggest the need to develop methods that trade-off more effectively between graph distortion and noise addition. There are a number of open directions in moving towards the goal of deploying synthetic graph data methods:



- (1) *Model / hyper-parameter selection under privacy constraints*: our experiments suggest that choice of hyperparameters (e.g., number of partitions in subsample-and-aggregate) and more generally choice of method can have a large impact on utility raising important problem of how to choose the model and hyper-parameters privately.
- (2) *General vs. tailored methods of synthetic graph generation*: there are not existing DP synthetic graph algorithms specifically tailored to fraud detection. In our experiments, we find that existing methods introduced significant bias even without noisy sufficient statistics, suggesting that these models do not capture relevant structure of graphs needed to model fraudulent behavior.
- (3) *Modeling synthetic graph meta-data*: existing synthetic graph methods tend to focus on modeling graph structure. This precludes practical application of these methods for settings with rich meta-data, which characterizes many real-world fraud graphs. Additionally, we hypothesize that modelling graph meta-data can lead to more effective DP synthetic graph generation methods as meta-data may be easier to model under DP constraints.

## Acknowledgments

We thank Patrick Grother and Craig Watson of NIST for constructive discussion and comments. A. Goldberg and G. Fanti acknowledge the Air Force Office of Scientific Research under award number FA9550-21-1-0090, NSF grant CNS-2148359, the Bill & Melinda Gates Foundation, Intel, and the Sloan Foundation for their generous support. A. Goldberg and N. Shah acknowledge the support of NSF grant 2200410 and ONR grant N000142212181. S. Wu acknowledges the support of NSF grant 2232693.

## References

- [1] Leman Akoglu, Mary McGlohon, and Christos Faloutsos. 2010. Oddball: Spotting anomalies in weighted graphs. In *Advances in Knowledge Discovery and Data Mining: 14th Pacific-Asia Conference, PAKDD 2010, Hyderabad, India, June 21-24, 2010. Proceedings. Part II 14*. Springer, 410–421.
- [2] Andrés F Barrientos, Aaron R Williams, Joshua Snoke, and CM Bowen. 2021. *Differentially Private Methods for Validation Servers*. Technical Report. Urban Institute research report.
- [3] Jeremiah Blocki, Avrim Blum, Anupam Datta, and Or Sheffet. 2013. Differentially private data analysis of social networks via restricted sensitivity. In *Proceedings of the 4th Conference on Innovations in Theoretical Computer Science* (Berkeley, California, USA) (ITCS '13). Association for Computing Machinery, New York, NY, USA, 87–96. <https://doi.org/10.1145/2422436.2422449>
- [4] Xihui Chen, Sjouke Mauw, and Yuniior Ramirez-Cruz. 2019. Publishing Community-Preserving Attributed Social Graphs with a Differential Privacy Guarantee. *Proceedings on Privacy Enhancing Technologies* 2020 (2019), 131 – 152. <https://api.semanticscholar.org/CorpusID:202540124>
- [5] Edith Cohen, Xin Lyu, Jelani Nelson, Tam'as Sarl'os, and Uri Stemmer. 2022. Generalized Private Selection and Testing with High Confidence. *ArXiv abs/2211.12063* (2022). <https://api.semanticscholar.org/CorpusID:253761282>
- [6] Federal Trade Commission. 2024. As Nationwide Fraud Losses Top \$10 Billion in 2023, FTC Steps Up Efforts to Protect the Public. <https://www.ftc.gov/news-events/news/press-releases/2024/02/nationwide-fraud-losses-top-10-billion-2023-ftc-steps-efforts-protect-public>. (Accessed on 02/29/2024).
- [7] Datavisor. 2024. Datavisor: AI Powered Fraud Platform for Enterprise. <https://www.datavisor.com/> (Accessed on 02/29/2024).
- [8] Yingdong Dou, Zhiwei Liu, Li Sun, Yutong Deng, Hao Peng, and Philip S. Yu. 2020. Enhancing Graph Neural Network-based Fraud Detectors against Camouflaged Fraudsters. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management* (Virtual Event, Ireland) (CIKM '20). Association for Computing Machinery, New York, NY, USA, 315–324. <https://doi.org/10.1145/3340531.3411903>
- [9] Cynthia Dwork. 2006. Differential Privacy. In *Automata, Languages and Programming*, Michele Bugliesi, Bart Preneel, Vladimiro Sassone, and Ingo Wegener (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 1–12.
- [10] Cynthia Dwork and Jing Lei. 2009. Differential privacy and robust statistics. In *Proceedings of the Forty-First Annual ACM Symposium on Theory of Computing* (Bethesda, MD, USA) (STOC '09). Association for Computing Machinery, New York, NY, USA, 371–380. <https://doi.org/10.1145/1536414.1536466>
- [11] Cynthia Dwork and Aaron Roth. 2014. The Algorithmic Foundations of Differential Privacy. *Found. Trends Theor. Comput. Sci.* 9, 3–4 (aug 2014), 211–407. <https://doi.org/10.1561/04000000042>
- [12] Prince Grover, Julia Xu, Justin Tittelfitz, Anqi Cheng, Zheng Li, Jakub Zablocki, Jianbo Liu, and Hao Zhou. 2023. Fraud Dataset Benchmark and Applications. arXiv:2208.14417 [cs.LG]
- [13] Paul W. Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. 1983. Stochastic blockmodels: First steps. *Social Networks* 5, 2 (1983), 109–137. <https://www.sciencedirect.com/science/article/pii/0378873383900217>
- [14] Y. Hu, F. Wu, Q. Li, Y. Long, G. Garrido, C. Ge, B. Ding, D. Forsyth, B. Li, and D. Song. 2024. SoK: Privacy-Preserving Data Synthesis. In *2024 IEEE Symposium on Security and Privacy (SP)*. IEEE Computer Society, Los Alamitos, CA, USA, 2–2. <https://doi.org/10.1109/SP54263.2024.00002>
- [15] Steven Jecmen, Nihar B. Shah, Fei Fang, and Leman Akoglu. 2024. On the Detection of Reviewer-Author Collusion Rings From Paper Bidding. arXiv:2402.07860 [cs.SI]
- [16] Steven Jecmen, Minji Yoon, Vincent Conitzer, Nihar B. Shah, and Fei Fang. 2022. A Dataset on Malicious Paper Bidding in Peer Review. <https://doi.org/10.48550/ARXIV.2207.02303>
- [17] Shouling Ji, Prateek Mittal, and Raheem Beyah. 2017. Graph Data Anonymization, De-Anonymization Attacks, and De-Anonymizability Quantification: A Survey. *IEEE Communications Surveys & Tutorials* 19, 2 (2017), 1305–1326. <https://doi.org/10.1109/COMST.2016.2633620>
- [18] Zach Jorgensen, Ting Yu, and Graham Cormode. 2016. Publishing Attributed Social Graphs with Formal Privacy Guarantees. In *Proceedings of the 2016 International Conference on Management of Data* (San Francisco, California, USA) (SIGMOD '16). Association for Computing Machinery, New York, NY, USA, 107–122. <https://doi.org/10.1145/2882903.2915215>
- [19] Kaggle. 2024. Kaggle: Your Machine Learning and Data Science Community. <https://www.kaggle.com/> (Accessed on 02/29/2024).
- [20] Vishesh Karwa, Pavel N. Krivitsky, and Aleksandra B. Slavković. 2017. Sharing social network data: differentially private estimation of exponential family random-graph models. *Journal of the Royal Statistical Society. Series C (Applied Statistics)* 66, 3 (2017), 481–500. <http://www.jstor.org/stable/44682587>
- [21] Vishesh Karwa, Sofya Raskhodnikova, Adam Smith, and Grigory Yaroslavtsev. 2014. Private Analysis of Graph Structure. *ACM Trans. Database Syst.* 39, 3, Article 22 (oct 2014), 33 pages. <https://doi.org/10.1145/2611523>
- [22] Shiva Prasad Kasiviswanathan, Kobbi Nissim, Sofya Raskhodnikova, and Adam Smith. 2013. Analyzing Graphs with Node Differential Privacy. In *Theory of Cryptography*, Amit Sahai (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 457–476.
- [23] Michael Kearns, Aaron Roth, Zhiwei Wu, and Grigory Yaroslavtsev. 2016. Private algorithms for the protected in social network search. *Proceedings of the National Academy of Sciences* 113 (01 2016), 201510612. <https://doi.org/10.1073/pnas.1510612113>
- [24] Nitin Kohli and Paul Laskowski. 2023. Differential Privacy for Black-Box Statistical Analyses. *Proceedings on Privacy Enhancing Technologies* 3 (2023), 418–431.
- [25] Kount. 2024. Kount: Fraud Detection and Chargeback Management Solutions. <https://kount.com/> (Accessed on 02/29/2024).
- [26] Ravi Kumar and Sergei Vassilvitskii. 2010. Generalized distances between rankings. In *Proceedings of the 19th International Conference on World Wide Web* (Raleigh, North Carolina, USA) (WWW '10). Association for Computing Machinery, New York, NY, USA, 571–580. <https://doi.org/10.1145/1772690.1772749>
- [27] Yang D. Li, Michaela F. Purcell, Thierry Rakotoarivelo, David B. Smith, Thilina Ranbaduge, and Kee Siong Ng. 2021. Private Graph Data Release: A Survey. *Comput. Surveys* 55 (2021), 1 – 39.
- [28] Michael I. Littman. 2021. Collusion rings threaten the integrity of computer science research. *Commun. ACM* 64, 6 (2021), 43–44.
- [29] Fang Liu, Evercita Eugenio, Ick Hoon Jin, and Claire Bowen. 2020. Differentially Private Generation of Social Networks via Exponential Random Graph Models. In *2020 IEEE 44th Annual Computers, Software, and Applications Conference (COMPSAC)*. 1695–1700. <https://doi.org/10.1109/COMPSAC48688.2020.00-11>
- [30] Jingcheng Liu and Kunal Talwar. 2018. Private Selection from Private Candidates. arXiv:1811.07971 [cs.DS]
- [31] Prashanth Mohan, Abhradeep Thakurta, Elaine Shi, Dawn Song, and David Culler. 2012. GUPt: privacy preserving data analysis made easy. In *Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data*. 349–360.
- [32] Hiep H. Nguyen, Abdessamad Imine, and Michaël Rusinowitch. 2015. Differentially Private Publication of Social Graphs at Linear Cost. In *Proceedings of the 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis*

- and Mining 2015 (Paris, France) (ASONAM '15). Association for Computing Machinery, New York, NY, USA, 596–599. <https://doi.org/10.1145/2808797.2809385>
- [33] Kobbi Nissim, Sofya Raskhodnikova, and Adam Smith. 2007. Smooth Sensitivity and Sampling in Private Data Analysis. In *Proceedings of the Thirty-Ninth Annual ACM Symposium on Theory of Computing* (San Diego, California, USA) (STOC '07). Association for Computing Machinery, New York, NY, USA, 75–84. <https://doi.org/10.1145/1250790.1250803>
- [34] Nicolas Papernot, Martin Abadi, Úlfar Erlingsson, Ian Goodfellow, and Kunal Talwar. 2017. Semi-supervised Knowledge Transfer for Deep Learning from Private Training Data. *arXiv:1610.05755* [stat.ML]
- [35] Nicolas Papernot, Shuang Song, Ilya Mironov, Ananth Raghunathan, Kunal Talwar, and Úlfar Erlingsson. 2018. Scalable Private Learning with PATE. *arXiv:1802.08908* [stat.ML]
- [36] Nicolas Papernot and Thomas Steinke. 2021. Hyperparameter Tuning with Renyi Differential Privacy. *ArXiv abs/2110.03620* (2021). <https://api.semanticscholar.org/CorpusID:238419564>
- [37] Gang Qiao, Weijie Su, and Li Zhang. 2021. Oneshot Differentially Private Top-k Selection. In *Proceedings of the 38th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, 8672–8681. <https://proceedings.mlr.press/v139/qiao21b.html>
- [38] Jerome P Reiter, Anna Oganian, and Alan F Karr. 2009. Verification servers: Enabling analysts to assess the quality of inferences from public use data. *Computational Statistics & Data Analysis* 53, 4 (2009), 1475–1482.
- [39] Riskified. 2024. Riskified: Fraud Prevention & Chargeback Fraud Protection. <https://www.riskified.com/> (Accessed on 02/29/2024).
- [40] Vincent A Traag, Ludo Waltman, and Nees Jan Van Eck. 2019. From Louvain to Leiden: guaranteeing well-connected communities. *Scientific reports* 9, 1 (2019), 5233.
- [41] T. N. Vijaykumar. 2020. Potential Organized Fraud in ACM/IEEE Computer Architecture Conferences. <https://medium.com/@tnvijayk/potential-organized-fraud-in-acm-ieee-computer-architecture-conferences-ccd61169370d>.
- [42] Mark Weber, Giacomo Domeniconi, Jie Chen, Daniel Karl I. Weidele, Claudio Bellei, Tom Robinson, and Charles E. Leiserson. 2019. Anti-Money Laundering in Bitcoin: Experimenting with Graph Convolutional Networks for Financial Forensics. *arXiv:1908.02591* [cs.SI]
- [43] Ruihan Wu, Chuan Guo, Felix Wu, Rahul Kidambi, Laurens van der Maaten, and Kilian Q. Weinberger. 2021. Making Paper Reviewing Robust to Bid Manipulation Attacks. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18–24 July 2021, Virtual Event (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, Virtual Event, 11240–11250. <http://proceedings.mlr.press/v139/wu21b.html>
- [44] Minji Yoon, Yue Wu, John Palowitch, Bryan Perozzi, and Russ Salakhutdinov. 2023. Graph generative model for benchmarking graph neural networks. In *Proceedings of the 40th International Conference on Machine Learning (ICML '23)*. JMLR.org, Honolulu, Hawaii, USA, Article 1680, 24 pages.
- [45] Haoyang Yu and Jerome P Reiter. 2018. Differentially Private Verification of Regression Predictions from Synthetic Data. *Trans. Data Priv.* 11, 3 (2018), 279–297.
- [46] Quan Yuan, Zhikun Zhang, Linkang Du, Min Chen, Peng Cheng, and Mingyang Sun. 2023. PrivGraph: Differentially Private Graph Data Publication by Exploiting Community Information. In *32nd USENIX Security Symposium (USENIX Security 23)*. USENIX Association, Anaheim, CA, 3241–3258. <https://www.usenix.org/conference/usenixsecurity23/presentation/yuan-quan>
- [47] Kiarash Zahrimia, Yaochen Hu, Mark Coates, and Oliver Schulte. 2024. Neural Graph Generation from Graph Statistics. *Advances in Neural Information Processing Systems* 36 (2024).

## A Summary of Datasets

	Yelp	Amazon	Peer Review	Elliptic
Vertices	11,473	11,944	2,483	6,621
Edge Density (%)	0.41	6.17	0.77	0.04
Num Fraud	1,657	821	22	11
Max Degree	236	6,991	255	47
Mean Degree	47.45	736.50	19.12	2.51

Table 1: Graph test datasets

## B Discussion of Subsample-and-Aggregate

### B.1 Controlling Bias and Variance

Algorithm 1 introduces error to the outputted AUC in three ways. First, sub-sampling the graph may introduce bias to the AUC score estimated on each partition; that is,  $E[f_{\text{AUC}}(G_i) - f_{\text{AUC}}(G)] \neq 0$ . Second, the algorithm adds Laplace noise to the released statistic, with variance proportional to the inverse of the number of partitions  $k$ . Finally, estimating on sub-samples of the data may increase the variance of the estimate.

*Tabular data.* In the special case of benchmarking a fraud detector using tabular data (e.g., only using vertex metadata not graph structure) with full duplication of fraudulent vertices ( $\rho = 1$ ), Algorithm 1 introduces error only from noise addition. In particular, changing the value of any one row in a dataset does not change properties of other rows of the dataset. Therefore, on tabular data, partitioning does not change the fraud scores of individual data-points compared to running the fraud detector on the entire dataset. Then from the definition of the AUC score, the AUC score of sub-partition  $i$  is

$$f_{\text{AUC}}(\mathcal{A}, G_i) = \frac{k}{n_1 n_0} \sum_{v_0 \in V_0^{(i)}} \sum_{v_1 \in V_1} \mathbb{1}[\mathcal{A}(G, v_1) > \mathcal{A}(G, v_0)]$$

where  $V_1^{(i)} = V_1$  since we duplicate all fraud vertices in each partition. Then, the mean over all partitions is:

$$\begin{aligned} & \frac{1}{k} \sum_{i=1}^k f_{\text{AUC}}(\mathcal{A}, G_i) \\ &= \frac{1}{n_1 n_0} \sum_{v_0 \in V_0} \sum_{v_1 \in V_1} \mathbb{1}[\mathcal{A}(G, v_1) > \mathcal{A}(G, v_0)] = f_{\text{AUC}}(\mathcal{A}, G) \end{aligned}$$

so the mean AUC over partitions exactly recovers the AUC score evaluated on the whole graph. In fact, taking  $k = n_0$  we add Laplace noise with scale proportional to  $\frac{1}{n_0}$  in the aggregation step, exactly recovering the Laplace mechanism for a query with global sensitivity of  $\frac{1}{n_0}$ .

*Graph data.* For graph data, partitioning can introduce bias to the estimate of AUC of each partition. The magnitude and direction of the bias may depend on the combination of graph and fraud detector under evaluation. We show this by way of a stylized example. Consider a fraud detector  $\mathcal{A}$  which scores each vertex in the graph by its degree. Let graph  $G$  be a random sample from a very simple stochastic block model (SBM) [13]. The SBM is defined as follows: fix a number of fraudulent vertices  $n_1$  and benign vertices  $n_0$ . Then for each pair of fraudulent vertices in the graph, sample an edge between the two independently at random with probability  $p_1$ . For each pair of benign vertices, sample an edge between the two vertices i.i.d. with probability  $p_0$ . Letting the random variable  $D$  be the difference in degree between a random fraudulent vertex and a benign vertex, we are concerned with the expected AUC score of a graph sampled from the SBM model, which is exactly  $E[\mathbb{1}[D > 0]] = \Pr[D > 0]$ .

The expected difference between degree of a fraud vertex and degree of a benign vertex is

$$E[D] = (n_1 - 1)p_1 - (n_0 - 1)p_0$$

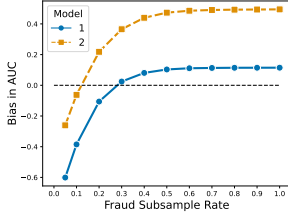
while its standard deviation is given by

$$SD[D] = \sqrt{(n_1 - 1)p_1(1 - p_1) + (n_0 - 1)p_0(1 - p_0)}$$

For large  $n_1p_1$  and  $n_0p_0$ ,  $D$  can be approximated by a Normal distribution so we can estimate the expected AUC score as

$$\Pr[D > 0] \approx \Pr\left[Z > -\frac{E[D]}{SD[D]}\right]$$

where  $Z$  is a standard normal random variable. Suppose that  $p_1 > p_0$  and  $E[D] > 0$ , so ranking by a degree is a good estimator in that  $\Pr[D > 0] > 0.5$ . Now, consider what happens to expected AUC score of a sub-partition, where fraud and benign vertices are sub-sampled at the same rate ( $\rho = \frac{1}{k}$ ). In this case,  $E[D]$  decreases by a factor of roughly  $\frac{1}{k}$  while  $SD[D]$  decreases by a factor of  $\frac{1}{\sqrt{k}}$  so  $-\frac{E[D]}{SD[D]}$  gets larger, negatively biasing the AUC score downwards. In contrast, consider setting  $\rho = 1$  so all fraud vertices are duplicated. In this case,  $E[D]$  actually increases since we have only down-sampled benign vertices, while the variance decreases, so  $-\frac{E[D]}{SD[D]}$  decreases putting upward bias on the AUC score.



**Figure 4: Comparison of the bias in sub-graph AUC score of the “rank by degree” detector as a function of fraud subsampling rate  $\rho$  fixing  $k = 20$ . Results are for two simulated SBM models on 1,000 benign and 100 fraud vertices. The bias-minimizing choice of parameter is between  $1/k$  and 1, but quite different on the two datasets.**

In theory, then, we would like to choose a sub-sampling rate somewhere between  $\frac{1}{k}$  and 1 to minimize this bias. Unfortunately, it is unclear how to set this sub-sampling rate in general. For example, in Figure 4, we show simulations for simple SBM models on 1,000 benign and 100 fraud vertices where Model 1 has  $n_0p_0 = 5$ ,  $n_1p_1 = 10$ , while Model 2 has  $n_0p_0 = 9$ ,  $n_1p_1 = 10$ . In both cases, we can observe that taking  $\rho = 1.0$  leads to positive bias while  $\rho = 1/k$  leads to negative bias. However, the two differ in optimal choice of subsample rate. For Model 1, the best subsample rate for 0 bias is roughly 0.3, which gives large positive bias on Model 2. Meanwhile, for Model 2, the best subsample rate for 0 bias is around 0.1, which gives significant negative bias on Model 2. It is unclear how to choose the error-minimizing parameters as privately estimating the error in AUC requires estimating the AUC, the original estimation problem. In prior work on subsample-and-aggregate building a system named GUPT, the authors advocate for choosing parameters

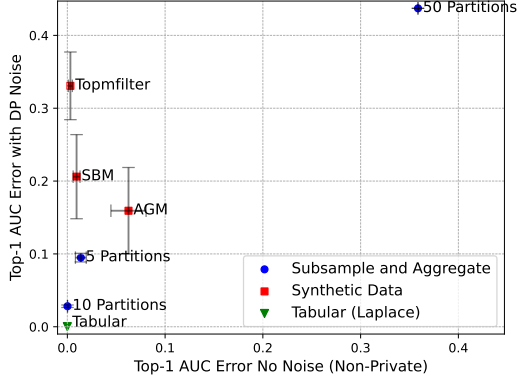
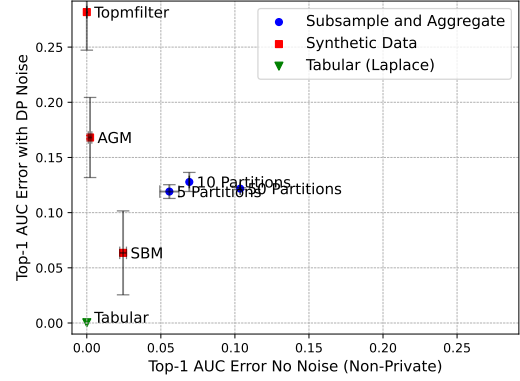
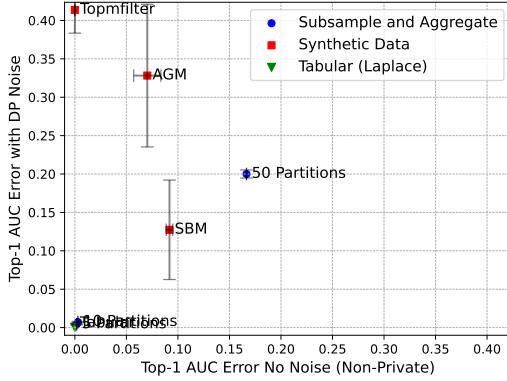
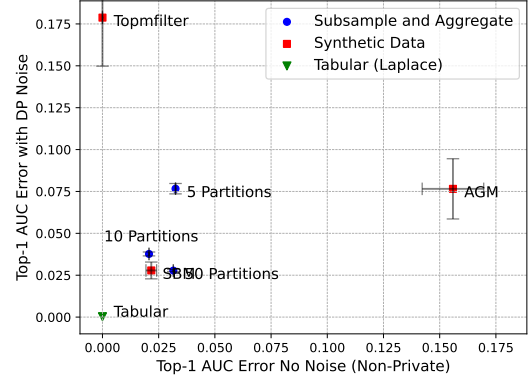
of subsample-and-aggregate based on older (now) public data that is similar to the private dataset [31]. However, such data may be difficult to find in the fraud setting. In our work, we empirically evaluate error as a function of  $k$ ,  $\rho$ , and  $\epsilon$  on validation datasets in Section 7.2 and then use the best choice of parameters on test datasets for comparison against synthetic data methods.

## B.2 Runtime

We observe in experiments that subsample-and-aggregate often significantly speeds up evaluation of fraud detectors. In this appendix, we give some theoretical intuition for why this may be. Algorithm 1 requires running the same fraud detector on  $k$  partitions of the dataset, each of size roughly  $\frac{1}{k}$  the original number of vertices in the dataset (in fact, slightly larger due to duplication of fraud vertices). In many cases, the runtime of a fraud detector actually decreases by a factor of more than  $\frac{1}{k}$  per partition. This can happen for two reasons. First, many graph algorithms are polynomial in the number of vertices in the graph (for example, an algorithm that cubes the adjacency matrix to compute number of triangles per vertex). Hence, the partitioning gives a polynomial  $\frac{1}{k}$  improvement in total computational cost. Second, partitioning can only decrease the total number of edges across all partitions since no edge can be duplicated in multiple partitions. Therefore, any algorithm with runtime dependent on number of edges is faster when run on all the partitioned graphs rather than the original graph. In addition, Algorithm 1 can be parallelized by running the fraud detector on each partition separately, which may make it easier for the benchmarking server to efficiently execute submitted fraud detectors in practice.

## C Other Synthetic Graph Algorithms

We briefly discuss other popular synthetic graph models considered in the literature. One approach uses exponential random graph models (ERGMs) to model vertex-labelled graph data [20, 29]. These methods are difficult to scale to graphs of more than a few hundred vertices, and prior empirical evaluations are limited to graphs of this size. Hence, they are not applicable to the types of fraud graphs we consider which are larger by an order of magnitude. Recent work [44, 47] has considered using Graph Neural Networks (GNNs) to generate synthetic graph data from graph statistics. They find that directly using DP-SGD (stochastic gradient descent) to train the GNN leads to poor utility. However, it is possible to obtain useful synthetic data by computing vectors of vertex-level graph statistics (like histograms of triangles and 2-paths) under edge-level DP. Estimating these sub-graph histogram statistics under *node-level* DP requires much larger noise addition. For instance, even assuming that a graph has no vertices of degree greater than  $T$ , the triangle histogram has sensitivity of greater than  $T^2$ , while we would expect most vertices to participate in far fewer than  $T$  triangles. Because we would need to add much more noise, than the experiments from this work (which just add noise scaled to  $\frac{1}{\epsilon}$  to the statistics) we do not focus on these methods in this work. Finally, some methods incorporate a community detection step [4, 46] that first clusters vertices of the graph and then estimates connection probability parameters between these clusters to incorporate into

Figure 5: Top-1 AUC,  $\epsilon = 5.0$  on Elliptic DataFigure 7: Top-1 AUC,  $\epsilon = 2.0$  on Amazon DataFigure 6: Top-1 AUC,  $\epsilon = 5.0$  on Peer Review DataFigure 8: Top-1 AUC error,  $\epsilon = 2.0$  on Yelp Data

the graph generative model. It is unclear how to make this clustering step satisfy node-level DP with reasonable utility.

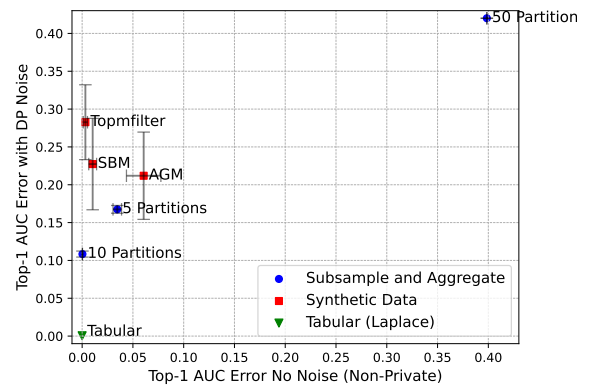
## D Additional Results

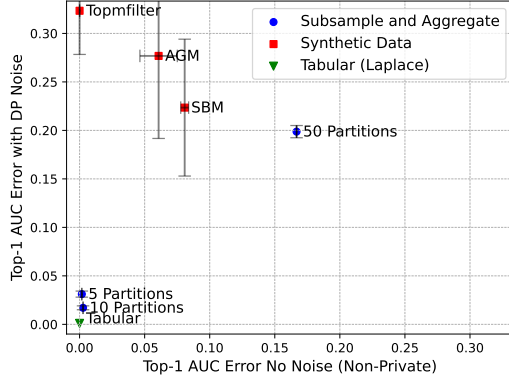
In this section we present additional results of our empirical experiments.

### D.1 Comparison of Algorithms, Top-1 Release

In this section, we provide additional results for the comparison of all algorithms at top-1 release as measured by the error in AUC score of the released best fraud detector vs. the actual best fraud detector. First, in Figures 5 and 6, we show the same plot as Figure 1 for the Elliptic and Peer Review datasets. We generally observe a similar trend of increasingly complex methods performing worse after private noise addition than simpler methods. Interestingly, on Elliptic, subsample and aggregate works well with  $\epsilon = 5.0$ . We note that Elliptic is much sparser than the other graphs, so it may admit different effective algorithms.

In Figures 7, 8, 9, and 10, we give results for the stricter privacy budget of  $\epsilon = 2.0$ .

Figure 9: Top-1 AUC error,  $\epsilon = 2.0$  on Elliptic Data


 Figure 10: Top-1 AUC error,  $\epsilon = 2.0$  on Peer Review Data

## D.2 Comparison of Algorithms, Full Leaderboard

For the full leaderboard, to capture distance between the true ranking of fraud detectors and the privacy-preserving noisy ranking, we use a similarity-weighted Kendall-Tau distance [26], which counts the number of inversions between two rankings, weighted by the difference in true AUC scores of the swap. Precisely, let  $\sigma(i)$  denote the rank of fraud detector  $\mathcal{A}_i$  in the true AUC leaderboard and  $\tilde{\sigma}(i)$  denote the rank of fraud detector  $\mathcal{A}_i$  in the noisy AUC leaderboard. Then, the similarity weighted Kendall-Tau distance is given by:

$$\sum_{(i,j):\sigma(i)<\sigma(j)} \mathbb{1}[\tilde{\sigma}(i) > \tilde{\sigma}(j)] (f_{\text{AUC}}(\mathcal{A}_i, G) - f_{\text{AUC}}(\mathcal{A}_j, G)).$$

As a baseline value for the Kendall-Tau similarity on our set of 10 fraud detectors on each dataset, we can compute the expected distance between the true leaderboard and a random permutation of the fraud detectors for each dataset. This yields values in the range of 5 to 8 for each dataset (which we show as baselines in our results section). For further validation of the metric, we consider the distance between rankings on validation and test sets for the Yelp and Elliptic datasets. We find that the distance from test to validation is 0.003 and 0.021 respectively reflecting that test and validation sets reliably produce similar leaderboards.

In Figures 11 and 12, we show the Kendall-Tau distance on each dataset for the best choice of parameters with privacy budgets of  $\epsilon = 5.0$  and  $\epsilon = 2.0$  respectively. Subsample-and-aggregate is generally less competitive at full leaderboard release than top-1 release, because it requires splitting the privacy budget across all of the 10 fraud detectors benchmarked. In contrast, synthetic data methods only use up privacy budget once to generate the synthetic data and then can benchmark any number of fraud detectors with no further privacy loss.

## D.3 Subsample-and-Aggregate

In this section, we give additional results for subsample and aggregate. First, in Figure 13, we show the bias to different fraud detectors as a function of fraud subsampling rate for the Amazon and Peer Review datasets (as in Figure 2 in the main text.) Then, in Figure 14,

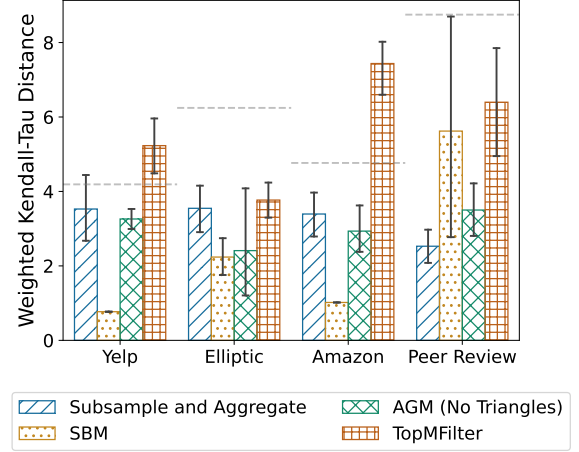


Figure 11: Head-to-head comparison of DP benchmarking methods for  $\epsilon = 5$  overall. Dashed lines show expected Kendall-Tau distance of a random permutation. Error bars show standard errors over 10 trials. SBM and subsample-and-aggregate are the most competitive approaches, though neither uniformly outperforms the other.

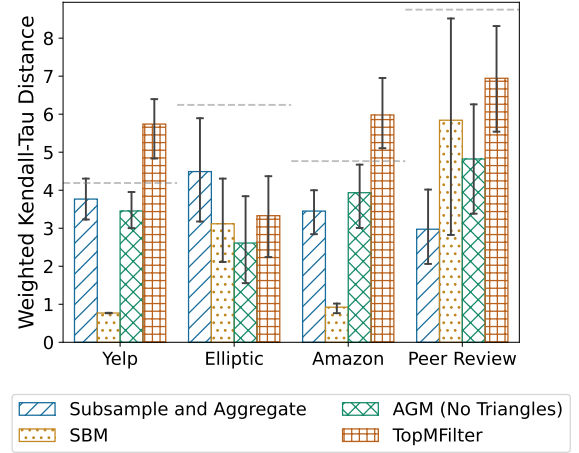
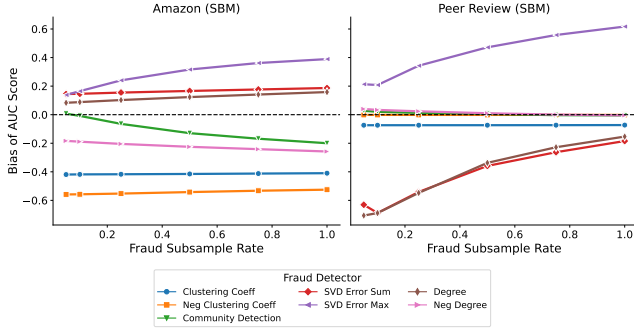
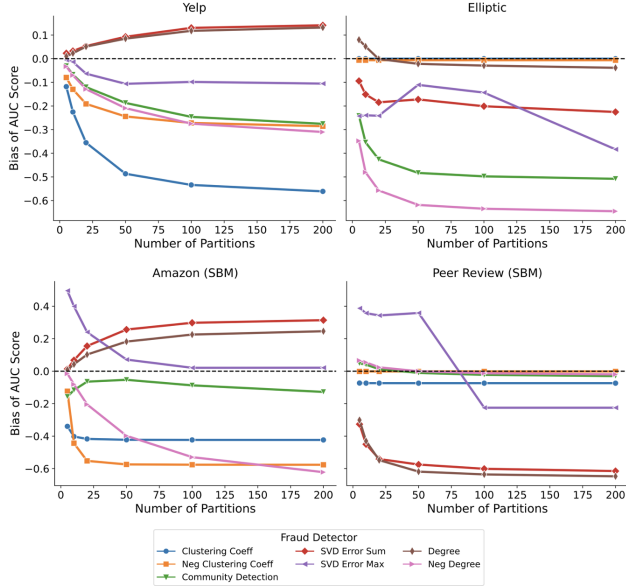


Figure 12: Head-to-head comparison of DP benchmarking methods for  $\epsilon = 2$  overall. Dashed lines show expected Kendall-Tau distance of a random permutation. Error bars show standard errors over 10 trials. SBM and subsample-and-aggregate are the most competitive approaches, though neither uniformly outperforms the other.

we show the bias of each fraud detector as a function of the number of partitions in subsample-and-aggregate. As expected, bias for each fraud detector increases in magnitude with more partitions (hence more graph distortion), but the sign and magnitude differ across different fraud detectors.



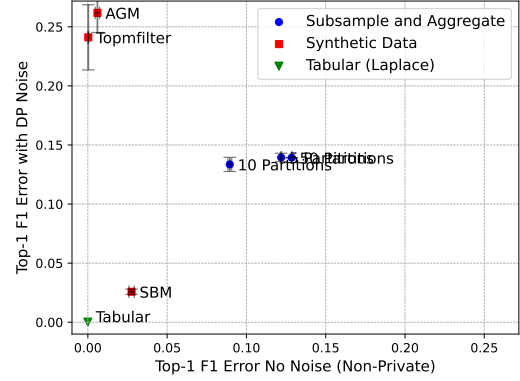
**Figure 13: Bias to the AUC score introduced by subsample and aggregate for each fraud detector varying the fraud subsample rate ( $\rho$ ) while fixing number of partitions  $k = 20$  for additional datasets.**



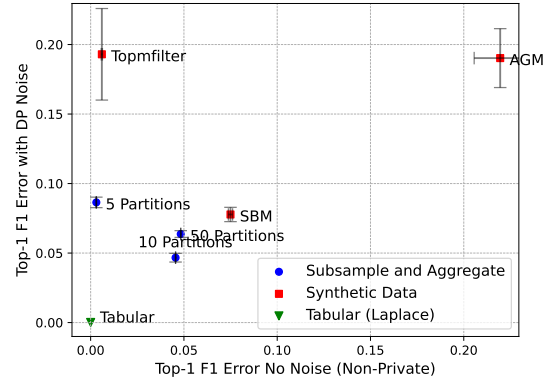
**Figure 14: Bias to the AUC score introduced by subsample and aggregate for each fraud detector varying the number of partitions ( $k$ ) while fixing fraud sub-sampling rate of  $\rho = 0.5$ . (Absolute) bias increases with  $k$ , but magnitude and sign of the bias varies per fraud detector and dataset.**

#### D.4 Synthetic Data

In this section, we provide additional results for synthetic data methods. In Table 2, we show the error to sufficient statistics using a more aggressive degree truncation threshold of 0.5 times the max degree, compared to 1.0 times the max degree in Figure 3 in the main text. Truncating more aggressively generally increases the error, except on Elliptic, where it decreases the error due to Elliptic being a highly sparse graph.



**Figure 15: Top-1 F1 Score,  $\epsilon = 5.0$  on Amazon Data**



**Figure 16: Top-1 F1 Score,  $\epsilon = 5.0$  on Yelp Data**

#### D.5 F1 Score

In this section, we give additional results using the F1 Score to benchmark fraud detectors instead of the AUC score. The F1 Score is the harmonic mean of the precision and recall of a classifier and has range  $[0, 1]$ . As it depends on a threshold chosen to convert a fraud detection score into a fraud/benign label, for each fraud detector we compute the F1 score as the best F1 score across all possible thresholds. We show results for F1 score, analogous to Figure 1 in the main text, with  $\epsilon = 5.0$  for all datasets in Figures 15 and 16.



	# Edges	Degree Sequence	# Triangles	Adjacency Matrix
Yelp	0.45	0.78	1.00	1.50
Elliptic	17.61	25.64	25.13	1.91
Amazon (SBM)	0.99	0.66	0.96	1.00
Peer Review (SBM)	0.21	0.53	1.94	1.66

**Table 2: Normalized mean absolute error (MAE) introduced to each of the synthetic graph sufficient statistics. Fixed  $\varepsilon = 5.0$  and degree cutoff of 0.5 the graph’s max degree.**

