

Graph-Fraudster: Adversarial Attacks on Graph Neural Network Based Vertical Federated Learning

Jinyin Chen, Guohan Huang, Haibin Zheng, Shanqing Yu, Wenrong Jiang, Chen Cui

Abstract—Graph neural network (GNN) has achieved great success on graph representation learning. Challenged by large scale private data collected from user-side, GNN may not be able to reflect the excellent performance, without rich features and complete adjacent relationships. Addressing the problem, vertical federated learning (VFL) is proposed to implement local data protection through training a global model collaboratively. Consequently, for graph-structured data, it is a natural idea to construct a GNN based VFL framework, denoted as GVFL. However, GNN has been proved vulnerable to adversarial attacks. Whether the vulnerability will be brought into the GVFL has not been studied. This is the first study of adversarial attacks on GVFL. A novel adversarial attack method is proposed, named Graph-Fraudster. It generates adversarial perturbations based on the noise-added global node embeddings via the privacy leakage and the gradient of pairwise node. Specifically, first, Graph-Fraudster steals the global node embeddings and sets up a shadow model of the server for the attack generator. Second, noise is added into node embeddings to confuse the shadow model. At last, the gradient of pairwise node is used to generate attacks with the guidance of noise-added node embeddings. Extensive experiments on five benchmark datasets demonstrate that Graph-Fraudster achieves the state-of-the-art attack performance compared with baselines in different GNN based GVFLs. Furthermore, Graph-Fraudster can remain a threat to GVFL even if two possible defense mechanisms are applied. Additionally, some suggestions are put forward for the future work to improve the robustness of GVFL. The code and datasets can be downloaded at <https://github.com/hgh0545/Graph-Fraudster>.

Index Terms—Vertical federated learning, graph neural network, adversarial attack, privacy leakage, defense.

I. INTRODUCTION

FEDERATED learning (FL) [1], [2], a privacy-preserving and distributed learning paradigm, establishes machine learning models based on distributed datasets across multiple parties/devices. It aims to protect client’s local data and mitigate privacy leakage in the context of legal restrictions, user side privacy protection and commercial competition. Benefiting from the collaborative training of the server model and privacy protection of the local raw data, FL may be widely applied to various industries, e.g., mobile service [3], healthcare [4], finance [5], face annotation [6], [7]. According to the data distribution characteristic, FL can be roughly categorized into horizontal FL (HFL), vertical FL (VFL) and federated transfer learning (FTL). HFL is suitable for the clients sharing datasets of same feature space but different examples, while VFL is designed for the clients sharing datasets of same examples but different feature spaces. FTL is proposed for the client’s dataset differs neither in feature space nor in examples. Among the three FLs, most studies are focused on HFL [8]–[12], while the other two still need to be further explored.

Vertical data distribution is typical in real-world applications, for instance, financial data (e.g., transaction records, income) is often vertically partitioned and owned by different financial institutions (e.g., banks or lending platforms). The banks attempt to avoid lending to users with low credit ratings. Thus, a reliable evaluation agency is necessary to assess the same users among financial parties while sharing different features. To avoid raw data sharing between banks, a VFL framework is a good option for such a practical scenario. As described, some financial data can be constructed as graphs. Consequently, graph neural network (GNN) based VFL, denoted as GVFL, is suitable for this scenario.

In such a practical scenario, there could be some potential threats in GVFL. Fig. 1 shows the threat model of the adversarial attack on GVFL. To be evaluated as high credit rating users, some may hide their true behaviors by adding some extra transaction records. Even worse, some dishonest lending platforms may tamper their data to help the low credit rating users by evading the detection of the evaluation agency. These maliciously crafted operations will cause the bank to issue loans to low-credit users.

We further analyze the possibility of the adversarial attack on GVFL, and find out that the threats could come from the innate deficiency in GVFL. Specifically, each participant in GVFL uploads the intermediate information, i.e., the node embeddings extracted by the local GNN model instead of the raw data. The adversarial perturbations (e.g., fake relationship, extra transaction record) on the local raw data will influence the updated node embeddings, thus posing a threat to the server model by making a wrong decision. Besides, privacy leakage of embeddings of GVFL will provide conditions for adversarial attacks, e.g., the leakage of the embeddings will help the malicious client to establish a precise shadow model of the server. These threats of adversarial attacks on GVFL have not been studied yet. Motivated by the practical application of GVFL and its potential vulnerability towards adversarial attacks, a novel attack method is proposed in this paper as the first work of security research of GVFL.

The contributions of our work are summarized as follows:

- To the best of our knowledge, this is the first work of proposing and formulation the adversarial attack on GVFL, by revealing its vulnerability towards a crisis of distrust in practical applications. Due to privacy leakage and data bias of the local model, GVFL unintentionally provides the conditions for successful adversarial attacks.
- A novel adversarial attack method is proposed against GVFL, denoted as Graph-Fraudster. Aiming at confusing the server model in GVFL, Graph-Fraudster generates

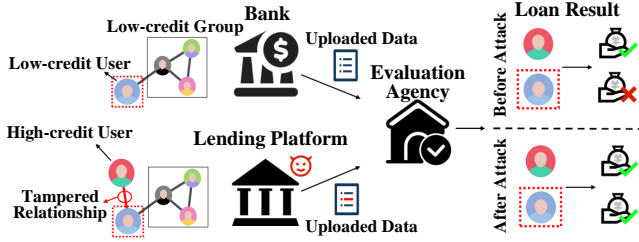


Fig. 1. An illustration of the adversarial attack on GVFL. Small perturbations (e.g., fake relationship, extra transaction record) of the graph-structured data could lead to misclassification of the evaluation agency based on GVFL.

adversarial perturbations based on the noise-added global node embeddings via GVFL’s privacy leakage and the gradient of pairwise node.

- Extensive experiments are conducted on five real-world graph-structured datasets for node classification with different GNN structured GVFLs. Graph-Fraudster is validated that it performs significantly better than three advanced adversarial attacks transferable to GVFL from centralized framework on seven metrics, achieving the state-of-the-art performance.
- Furthermore, we propose two possible defenses against the attack. Graph-Fraudster can still pose a threat to the defensive GVFL. Additionally, according to the observations of the experiments, some suggestions are provided to improve the robustness of GVFL, including privacy preserving, bias removal and defense/detection deployment. The code and datasets are available in <https://github.com/hgh0545/Graph-Fraudster>.

The rest of the paper is organized as follows. Related works on GNN based federated learning, inference attacks on federated learning and adversarial attacks on GNN models are reviewed in Section II. Section III defines the adversarial attack on GNN model and GVFL, sets up threat model, and then introduces Graph-Fraudster in details. In Section IV, extensive experiments are presented to demonstrate the performance of Graph-Fraudster. At last, the paper is concluded, and the future work is enumerated in Section V.

II. RELATED WORK

In line with the focus of the work, we briefly summarize the existing works of GNN based federated learning, inference attacks on federated learning and adversarial attacks on GNN.

A. GNN based Federated Learning

Though FL has shown the superiority in various domains, it has not been widely applied to graph-structured data. According to the investigation, most of the research focus on horizontal scenarios. For instance, considering data privacy preservation in non-independent identically and distribution (Non-IID), Zheng et al. [13] used clients to implement message communication, and leveraged Bayesian optimization for hyper-parameters fine-tuning. They further proposed a novel learning paradigm, named ASFGNN. Inspired by meta

learning, GraphFL [14] is proposed to address the node classification under Non-IID. These proposed FLs are supported by machine learning models designed for graph mining. In the aspect of GNN as the local model for FL, He et al. [15] presented FedGraphNN, an open-source FL system based on GNN models. But it is not suitable for VFL. As for GNN based vertical federated learning, the research is still in its fancy. Zhou et al. [16] gave the first vertical learning paradigm for privacy-preserving GNN models for node classification, by splitting the graph into two parts for different clients. Secure multi-party computation is adopted as well to ensure data privacy and efficiency. To ensure privacy, Ni et al. [17] adopted additively homomorphic encryption, and proposed a federated GCN learning paradigm for the privacy-preserving node classification task, named FedVGCN. Besides the framework study of GNN based FL, it is also applied to real-world applications. For instance, Wu et al. [18] proposed a GNN structured FL for recommendation to accomplish server model training and user-side privacy.

B. Inference Attack on Federated Learning

Since the inference attack on GNN based federated learning has not been studied yet, we summarize some inference attacks on federated learning. Numerous attacks on HFL have been proposed. For instance, in [19], generative adversarial network is used for inference attack on HFL for the first time, which generates the same distribution as the training data. Nasr et al. [20] exploited the vulnerability of the stochastic gradient descent (SGD), and designed a membership attack in white-box setting. Besides, deep leakage from gradients (DLG) [21] uses public shared gradients to obtain the local training data without extra knowledge of the data. To sum up, most of them greatly rely on the gradient which exchanged during training process. In vertical scenarios, Luo et al. [22] formulated the problem of feature inference attack in model’s inference process for the first time, and they proposed two attack methods on the logistic regression (LR) and decision tree (DT) models, named equality solving attack (ESA) and path restriction attack (PRA), respectively. Further, they proposed a general attack method named generative regression network attack to handle more complex models. Zhang et al. [23] devoted to gathering the man-made poison data and the corresponding intermediate outputs in the model’s training process. Then they built the training set for the decoder to recover other private data of the victim. Besides, Li et al. [24] explored whether the labels can be uncovered by sharing gradient in the VFL setting, and proposed a label inference attack method based on the gradient norm.

C. Adversarial Attacks on GNN Models

Existing technologies may have multiple vulnerabilities that enable hackers or criminals to manipulate data for malicious purposes [25], and the adversarial examples may jeopardize the operation of the reality systems such as the Internet of Things [26]. Extensive studies have proven that GNN models are vulnerable to imperceptible adversarial perturbations that affect the performance of downstream applications. For node

classification, Zügner et al. [27] proposed the first adversarial attack on GNN models, denoted as NETTACK, by generating attacks iteratively according to score functions. Reinforcement learning is applied to generate perturbations in [28] as well. Aiming at fooling GNN models by attacking embeddings of the target node, Chen et al. [29] proposed a fast gradient attack method (FGA) via the maximal absolute edge gradient. Analogously, IG-FGSM and IG-JSMA are introduced in [30]. To deal with the large-scale graph, Li et al. [31] proposed a multi-stage attack framework SGA, which relies on a much smaller subgraph centered at the target node. Another modification strategy is proposed by inserting fake nodes [32] instead of adding/deleting edges in the adversarial examples. Considering a black-box scenario, GF-Attack [33] is constructed by the graph filter and feature matrix without accessing any knowledge of the target classifiers. For graph classification, Wu et al. [34] used rewiring to hide the adversarial edges, which preserves some important properties of the graph such as the number of nodes, edges and total degrees of the graph. For community detection, genetic algorithm is adopted in Q-Attack [35] to fail the detection method. To sum up, the existing graph adversarial attacks mainly degrade the performance of GNN models in various tasks by adding or deleting key edges.

III. METHODOLOGY

First, we formalize the problem of adversarial attack on GVFL, and give its definition. Then, the threat model of adversarial attack on GVFL is elaborated. Afterward, we introduce the Graph-Fraudster in aspect of framework, implementation details and algorithm rationality. For convenience, the definitions of symbols used in this paper are listed in the TABLE I.

TABLE I
THE DEFINITIONS OF SYMBOLS.

Symbol	Definition
$G = (V, E)$	the original graph with sets of nodes and edges
A, \hat{A}	the adjacency matrix / adversarial adjacency matrix of G
X, \hat{X}	the feature matrix / adversarial feature matrix of nodes
$f_{\theta}(\cdot)$	the GNN model with parameter θ
K	the number of participants
h_m	the node embeddings of the malicious participant
h_{global}	the global node embeddings
V_L	the set of nodes with labels
T	the set of target nodes
ρ	the activation function
Y, Y'	the real / predicted label list
$ F $	the number of class for nodes in the G
N	the number of nodes in the graph G
P	the probabilities with the real global node embeddings
\tilde{P}	the probabilities with the fake global node embeddings
d	the dimensions of local node embeddings
h_{fake}	the fake global node embeddings
$\mathcal{S}(\cdot)$	the server model in GVFL
$\tilde{\mathcal{S}}(\cdot)$	the shadow server model
C_d^0, R_d^0	the weights and bias of the simplified $\mathcal{S}(\cdot)$
$B_S, B_{\tilde{S}}$	the boundary of \mathcal{S} and $\tilde{\mathcal{S}}$

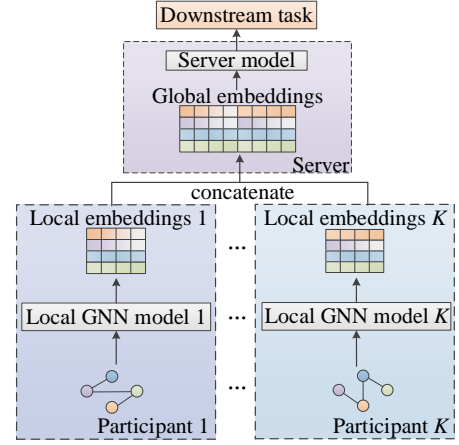


Fig. 2. Forward propagation of GVFL. The server model concatenates the local embeddings as global embeddings, and completes the downstream task.

A. Problem Definition

Definition 1: Adversarial attack on GNN model. For node classification, the attacker is aiming to mislead the target GNN model to output the expected label of the adversarial example, which is carefully crafted with imperceptible perturbations (e.g., rewiring edges, fake nodes) on the benign one. As expected, fed by the adversarial example, the GNN model will extract low quality node embeddings, leading to the wrong prediction. Specifically, $G = (V, E)$ represents a graph and $f_{\theta}(\cdot)$ represents a GNN model. Denote A and X as the adjacency matrix of G and node features, respectively. The adversarial attack on GNN model can be formalized as,

$$\begin{aligned} \max \sum_{t \in T} L_{atk}(\text{softmax}(f_{\theta^*}(\hat{A}, \hat{X}, v_t)), y_t) \\ \text{s.t. } \theta^* = \arg \min_{\theta} \sum_{v_l \in V_L} L_{train}(f_{\theta}(A, X, v_l), y_l) \end{aligned} \quad (1)$$

where \hat{A} represents the perturbed adjacency, \hat{X} expresses the perturbed node features. v_t represents the target node and y_t is the ground truth of v_t . The parameters of GNN model θ are trained with labeled node sets V_L , resulting θ^* . Note that, the attacker can manipulate edges and node features respectively, or both. The research shows that modifying edges is more effective than modifying node features [30]. Thus, only adding/deleting edges is considered as perturbations in this work.

Definition 2: Adversarial attack on GVFL. In GVFL, each client trains a local GNN model with its private data, and updates the embedding for the aggregation of the server. The attacker's goal is to disrupt the local embeddings, causing the server model to produce wrong predictions. Fig. 2 is a schematic diagram of the forward propagation for GVFL.

In this paper, concatenating the local node embeddings is used as the combination strategy:

$$h_{global} = h_1 || h_2 || \dots || h_K \quad (2)$$

where K denotes the number of participants, and $||$ is the concatenating operator.

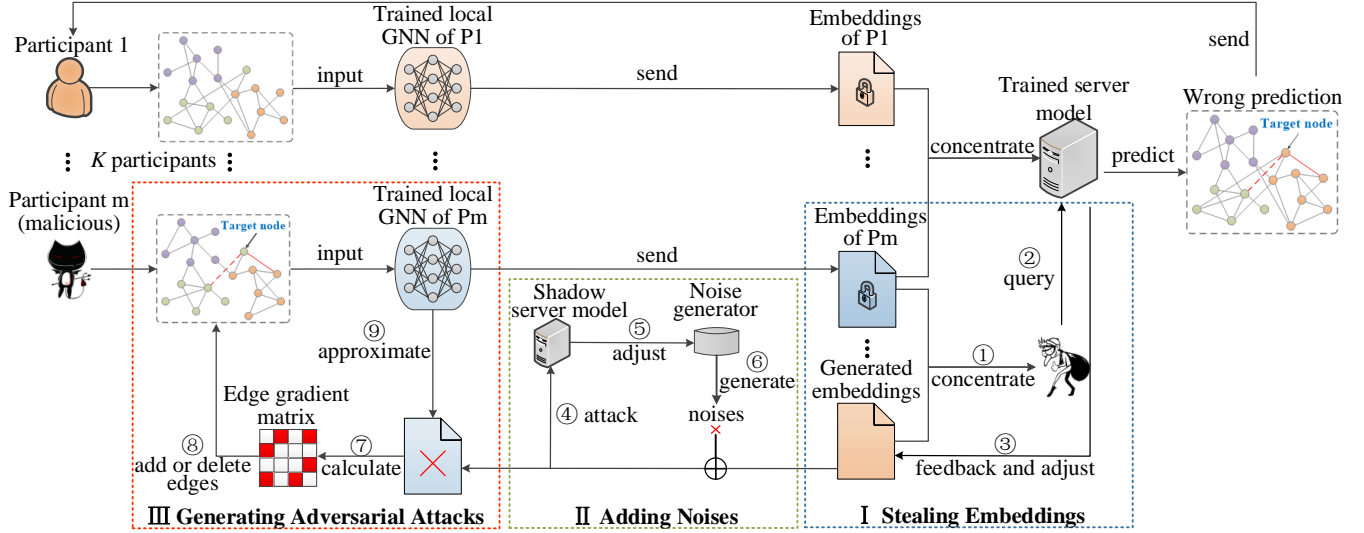


Fig. 3. The framework of Graph-Fraudster. One participant is randomly selected as the malicious participant from all. The adversarial attack is conducted through three stages: (1) stealing embeddings: a global node embeddings is generated, which adjusted by querying result of trained server model; (2) adding noises: noises are added into the generated embeddings to attack the shadow model of the server; (3) generating adversarial attacks: guided by noise-added node embeddings of malicious participant, adversarial attacks are generated by using gradient of pairwise node. As a result, the wrong prediction of the target node will be produced by the trained server model.

For a node classification task, the server model \mathcal{S} is a classifier to make prediction by:

$$Y' = \text{softmax}(W_1 \cdot \rho(\dots \rho(W_0 \cdot h_{\text{global}}))) \quad (3)$$

where $\{W_0 \dots W_1\}$ are weight matrices of the classifier. $\rho(\cdot)$ represents activation function such as ReLU.

Further, the GVFL uses the cross-entropy error over all labeled examples, which can be learned by gradient descent:

$$L_{\text{train}} = - \sum_{l=1}^{|V_L|} \sum_{n=1}^{|F|} Y_{ln} \ln(Y'_{ln}) \quad (4)$$

where V_L is the set of nodes with labels, Y_{ln} is the ground truth, and $|F|$ is the number of node classes in the graph G .

Then, the adversarial attack on GVFL is defined as:

$$\hat{Y}'_t = \text{softmax}(W_1 \cdot \rho(\dots \rho(W_0^{v_t} \cdot (h_1^{v_t} || h_2^{v_t} || \dots || \hat{h}_m^{v_t} || h_K^{v_t})))) \quad (5)$$

s.t. $\hat{h}_m^{v_t} = f_{\theta^*}(A, \hat{X}, v_t)$

where $\hat{h}_m^{v_t}$ is the perturbed node embeddings of the target node v_t which is uploaded by a malicious participant. It means a malicious participant can add perturbations into the graph by manipulating A or X . Then, the local GNN model may be confused, which will produce low quality or even targeted node embeddings and upload the perturbed node embeddings to the server model. Therefore, the server model will make wrong decision.

The target loss function L_{atk} of target node v_t can be defined as:

$$L_{\text{atk}} = - \sum_{n=1}^{|F|} Y_{tn} \ln(Y'_{tn}) \quad (6)$$

which measures the difference between the prediction and the ground truth. It is easy to find the prediction of the model is worse with larger value of L_{atk} .

B. Threat Model

Scenario. Taking more general scenarios into consideration, GVFL is applied to multiple participants training a central server collaboratively for node classification. When the server model is well-trained, the prediction of the server model will be shared by all participants. Besides, the server model is semi-honest, i.e., the server model serves to be queried by the malicious participant without revealing the inner parameters.

Knowledge of the malicious participant. Assume that the malicious participant knows the structure of the server, and only can access to its own data. In other word, the malicious participant can attack the server model by manipulating its own data merely. In addition, the trained server model will be served as an API for the malicious participant, and it only returns probability for each query. Thus, more details, such as the parameters of the server model and the gradient information, are unable to be obtained by the malicious participant.

The goal of the malicious participant. For the malicious participant, the goal is to attack the server model by uploading perturbed embeddings. So that the other benign participants, which rely on the server model, receive the unexpected prediction for the target example.

C. Graph-Fraudster

As described in Section III-B, the attacker can manipulate the own edges or node features in the testing to fool the global GVFL. During the aggregate process, modifying the edges affects all dimensions of features belonging to the target node, while modifying features only affects a few features. Thus, Graph-Fraudster is designed to attack the server model by manipulating edges of the local graph. It is conducted through three stages: (1) stealing node embeddings of other participants; (2) adding noises into embeddings via constructing

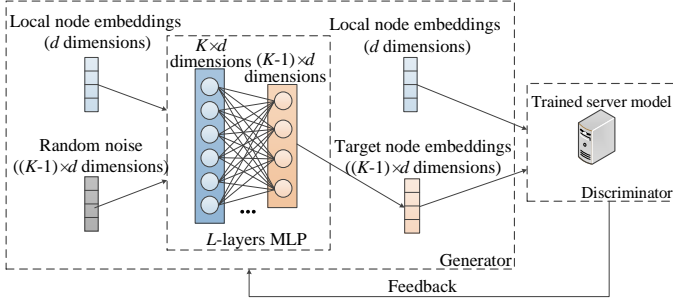


Fig. 4. Stealing node embeddings via GRN.

the shadow server model; (3) generating adversarial attacks to mislead the global GVFL. The framework of Graph-Fraudster is shown as Fig. 3. The details and rationality of Graph-Fraudster is described in this subsection.

1) *Embeddings Stealing Strategy*: Similar to [22], Graph-Fraudster uses generative regression network (GRN) to steal embeddings uploaded by other participants. First, as shown in Fig. 4, a $(K-1) \times d$ random noise vector is generated, where K is the number of participants in GVFL, and d is dimensions of local node embeddings. Then, a L -layers multilayer perceptron (MLP) is used to generate the target node embeddings, i.e., the embeddings uploaded by other participants. In this paper, a three-layer MLP with ReLU is applied as the generator. The dimension of an input is $K \times d$, the dimension of an output is $(K-1) \times d$, and the hidden units are 512 and 256, respectively. Then, the local node embeddings of the malicious participant and the target node embeddings are concatenated, represented as h_{fake} , as input of the discriminator. The classification probabilities of an adversarial example can be obtained by querying the server model. Note that the parameters of the trained server model are fixed. The target of stealing privacy can be converted to a regression problem. Therefore, mean square error (MSE) is adopted for the training process, which can be formulated as:

$$L_{GRN} = \frac{1}{N \times |F|} \sum_{i=1}^N \sum_{j=1}^{|F|} (\tilde{P}_{i,j} - P_{i,j})^2 \quad (7)$$

where N is the number of nodes, $\tilde{P} = \mathcal{S}(h_{fake})$ is the probabilities returned when node embeddings of an adversarial example are fed. P is the probabilities returned at the end of GVFL's training process.

To recognize the GRN in details, we simplify the model structure of GRN with a single-layer generator model and a discriminator model \mathcal{S} with fixed parameters. Thus, for the real node embeddings h_{real} , its ground-truth prediction output of the discriminator model can be represented as

$$P = \mathcal{S}(h_{real}) = \text{softmax}(C_d^0 h_{real} + R_d^0) \quad (8)$$

Due to the parameters of \mathcal{S} are fixed, C_d^0 is used to represent the weights of \mathcal{S} . R_d^0 is bias of \mathcal{S} . For a fake node embeddings

h_{fake} , the output of \mathcal{S} is

$$\begin{aligned} P_1 &= \mathcal{S}(h_{fake}) \\ &= \mathcal{S}(g(n_{in}) || h_m) \\ &= \text{softmax}(((n_{in} || h_m)w_0 + b_0) || h_m) C_d^0 + R_d^0 \quad (9) \\ \text{s.t. } g(n_{in}) &= (n_{in} || h_m)w_0 + b_0 \end{aligned}$$

where n_{in} is a $(K-1) \times d$ dimensions noise. h_m is local node embeddings of the malicious participant, $g(\cdot)$ is the generator, w_0 is the weights of the generator model and b_0 is bias. GRN back-propagates the loss and updates the parameters of the MLP model.

2) *Noise Addition via Stolen Embeddings*: In order to degrade the performance of the server model, it's significant to produce the adversarial embeddings crossing the decision boundary to mislead the server model. An effective and efficient method is using the gradient information of the server model to produce adversarial examples, which has been widely used to construct attack on GNN. However, in GVFL, the attackers cannot get gradient information but only the returned probabilities. Therefore, it is necessary to establish a shadow model in local to simulate the server model. By stealing the global node embeddings, the shadow model $\tilde{\mathcal{S}}$ can be successfully established based on the necessary knowledge, i.e., training data h_{fake} , model structure and target probabilities P . According to the setting in Section III-B, the model structure and P are known. To the end, the shadow model will output the similar probability P to the server's when fed by the same examples. MSE is applied as the target as follows:

$$L_{shadow} = \frac{1}{N \times |F|} \sum_{i=1}^N \sum_{j=1}^{|F|} (\tilde{\mathcal{S}}(h_{fake}) - P_{i,j})^2 \quad (10)$$

where $\tilde{\mathcal{S}}(h_{fake})$ is the output of $\tilde{\mathcal{S}}$.

Then, fast gradient sign method (FGSM) [36] is adopted as noise generator by using the gradient information of the shadow model $\tilde{\mathcal{S}}$.

For a target node v_t , its noise-added embeddings $\bar{h}_{fake}^{v_t}$ could be represented as:

$$\bar{h}_{fake}^{v_t} = h_{fake}^{v_t} + \epsilon \cdot \text{sign}\left(\frac{\partial L_{atk}}{\partial h_{fake}^{v_t}}\right) \quad (11)$$

where $\text{sign}(\cdot)$ is gradient's direction, $\epsilon \in [0, 1]$ is noise scale.

In particular, the reason for using MSE as the target to establish the shadow model is as follows. Suppose using a mapping equation to represent the server model:

$$\hat{k}_{\mathcal{S}}(h) = \arg \max_k \mathcal{S}(h) \quad (12)$$

where h is the node embeddings and k is k -th class of \mathcal{S} . The classification boundary function can be express as:

$$\mathcal{S}_k(h) - \hat{k}_{\mathcal{S}}(h) = 0 \quad (13)$$

By applying the first-order Taylor expansion, the approximate classification boundary function of \mathcal{S} is denoted as $\tilde{\mathcal{S}}$:

$$\begin{aligned} B_{\mathcal{S}} : \mathcal{S}_k(h_{real}) + h_{real} \nabla \mathcal{S}_k(h_{real}) \\ - \hat{k}_{\mathcal{S}}(h_{real}) - h_{real} \nabla \hat{k}_{\mathcal{S}}(h_{real}) = 0 \end{aligned} \quad (14)$$

$$\begin{aligned} B_{\tilde{S}} : \tilde{S}_k(h_{fake}) + h_{fake} \nabla \tilde{S}_k(h_{fake}) \\ - \hat{k}_{\tilde{S}}(h_{fake}) - h_{fake} \nabla \hat{k}_{\tilde{S}}(h_{fake}) = 0 \end{aligned} \quad (15)$$

The MSE of the shadow model is aiming to approximate the classification boundary of \tilde{S} and S . Thus, the target can be converted to:

$$\begin{aligned} \arg \min_{\tilde{S}} MSE(\tilde{S}(h_{fake}), S(h_{real})) \\ \Leftrightarrow \arg \min_{\tilde{S}} |B_{\tilde{S}} - B_S| \\ \Leftrightarrow \arg \min_{\tilde{S}} |\tilde{S}(h_{fake}) - S(h_{real})| \\ + |(h_{fake} - h_{real})(\nabla \tilde{S}(h_{fake}) - \nabla S(h_{real}))| \end{aligned} \quad (16)$$

Then, assuming $\nabla \tilde{S}(h_{fake}) \approx \nabla S(h_{real})$, and the noise generated by FGSM against the shadow model can also attack the server model in GVFL successfully since the approximate direction of the gradient is obtained.

3) *Adversarial Attack via Noise-added Embeddings*: Based on the above steps, the noise-added embeddings $\bar{h}_{fake}^{v_t}$ are obtained. Next, the part of the noise-added embeddings $\bar{h}_m^{v_t}$ that belongs to the malicious participant is extracted from $\bar{h}_{fake}^{v_t}$. Assuming that $\bar{h}_{fake}^{v_t}$ can confuse the shadow model successfully, the rest of the work is modifying some suitable edges in the original graph to make the node embeddings approximate to $\bar{h}_m^{v_t}$. Thus, Graph-Fraudster devotes to making the target node's embeddings and $\bar{h}_m^{v_t}$ more similar by adding perturbations into A .

$$\arg \min_{\hat{A}} MSE(f_{\theta^*}(\hat{A}, X, v_t), \bar{h}_m^{v_t}) \quad (17)$$

where $\|\hat{A} - A\|_0$ is twice the budget Δ , i.e., $\|\hat{A} - A\|_0 \leq 2\Delta$ due to the symmetry of the adjacency matrix.

To speed up the process of searching the suitable adversarial edges, the gradient of pairwise node is applied, followed the work of Chen et al. [29]. The edge gradient matrix is computed as follows:

$$\begin{aligned} g_{u,v} = \frac{\partial L_t}{\partial A_{u,v}} \\ s.t. \quad L_t = \frac{1}{d} \sum_{i=1}^d ([f_{\theta^*}(\hat{A}, X, v_t)]_i - [\bar{h}_{fake}^{v_t}]_i)^2 \end{aligned} \quad (18)$$

where (u, v) is any pairwise node in the graph G , $[\cdot]_i$ denotes the i -th element of node embeddings and d is the dimensions of local node embeddings. Considering the symmetry of the adjacency matrix of the undirected graphs, the gradient matrix is symmetric as:

$$g_{sym} = \frac{g + g^T}{2} \quad (19)$$

where T is the transpose symbol. The adversarial edge e_{adv} is selected by:

$$e_{adv} = (u, v) \leftarrow \arg \max_{u,v} -g_{sym} \quad (20)$$

To minimize the difference between the local GNN model's output and the noise-added embeddings rather than maximizing L_{atk} , the gradient value should be inverted. At last, the

adversarial edges will be added into clean adjacency iteratively within budget Δ :

$$\hat{A}_{u,v} = -A_{u,v} + 1 \quad (21)$$

where (u, v) comes from e_{adv} and \hat{A} is symmetrical.

D. Algorithm

The pseudo-code for Graph-Fraudster is given in Algorithm 1.

Algorithm 1: Graph-Fraudster

Input: Original adjacency A , node features X , target node v_t , attack budget Δ , dimensions of local node embeddings d , number of participants K .
Output: The adversarial network \hat{A} .

- 1 Train the GVFL to obtain the server model S and the local GNN models $f_{\theta^*}(\cdot)$.
- 2 Generate a $(K - 1) \times d$ dimensions noise randomly.
- 3 Concatenate generated noise and the local node embeddings of malicious participant h_m as input h_{in} .
- 4 **for** $t = 1$ to T_{GRN} **do**
- 5 $h_{target} \leftarrow MLP(h_{in})$
- 6 Concatenate h_{target} and h_m as h_{fake} for the server model S .
- 7 Minimize the MSE in Equation 7.
- 8 **end**
- 9 Save h_{fake} .
- 10 Initialize the parameter of the shallow server model \tilde{S} .
- 11 **for** $t = 1$ to T_{shadow} **do**
- 12 Minimize the MSE in Equation 10 with h_{fake} .
- 13 **end**
- 14 Generate the noises and add them into h_{fake} as Equation 11 and return $\bar{h}_{fake}^{v_t}$ for target node v_t .
- 15 Initialize $\hat{A}^0 = A$.
- 16 **for** $i = 1$ to Δ **do**
- 17 Calculate the symmetrical edge gradient matrix as Equation 18 and 19.
- 18 Select adversarial edge $e_{adv} = (u, v)$ by Equation 20.
- 19 Add adversarial edge $e_{adv} = (u, v)$ into \hat{A}^{i-1} as Equation 21 and obtain \hat{A}^i .
- 20 **end**
- 21 **Return** the adversarial adjacency matrix \hat{A} .

E. Complexity Analysis

For generation of adversarial attacks, the calculation time of Graph-Fraudster is divided into four parts, including GRN's training time T_{GRN} , shadow model's training time T_{shadow} , noise adding time and adversarial edges selecting time. Thus, the time complexity of Graph-Fraudster is:

$$\mathcal{O}(T_{GRN}) + \mathcal{O}(T_{shadow}) + \mathcal{O}(1) + \mathcal{O}(\Delta) \sim \mathcal{O}(N) \quad (22)$$

where Δ is the attack budget. $\mathcal{O}(T_{GRN})$ and $\mathcal{O}(T_{shadow})$ are the complexity of the training time of GRN and the shadow, depending on maximum training epochs T_{GRN} and T_{shadow} , respectively. $\mathcal{O}(1)$ indicates the complexity of noise adding time because FGSM is a one-step noise generator. $\mathcal{O}(\Delta)$ is the complexity of adversarial edges selecting time, which controlled by the attack budget Δ . $\mathcal{O}(N)$ indicates the complexity of Graph-Fraudster is linear, according to the total number of the above steps N .

The parameters of Graph-Fraudster include GRN’s parameters, shadow model’s parameters. Therefore, the space complexity is:

$$\mathcal{O}(K \times d \times n_0 + n_0 \times n_1 + \dots + (K-1) \times d \times n_{L-1}) + \mathcal{O}(\tilde{S}) \sim \mathcal{O}(M^2) \quad (23)$$

where K is the number of participants, d is the dimensions of local node embeddings, $[n_0, \dots, n_{L-1}]$ is the number of the units in GRN models and \tilde{S} is the shadow model. $\mathcal{O}(M^2)$ indicates that the space complexity depends on the memory occupied by the model’s weight matrix, whose is squared-level.

IV. EXPERIMENTS

In order to comprehensively evaluate the performance of Graph-Fraudster, both dual-participants based GVFL and multi-participant based GVFL are testified in aspect of attack performance, two possible defense strategies against GVFL and parameter sensitivity analysis.

A. Datasets

Five graph-structured datasets are used to evaluate the performance of the proposed method, including Cora [37], Cora_ML [37], Citeseer [37], Pol.Blogs [38] and Pubmed [39]. Their basic statistics are summarized in TABLE II, and the specific information is the same as the reference source.

In GVFL, each dataset will be split for different clients. Assuming there are no overlap edges for each client, then with more clients participant in, the more isolated nodes appear in clients’ graph data. It will be not conducive to collaboratively train the server model. Consequently, two splitting strategies are adopted for different number of clients. Specifically, for dual-participants, both edges and node features of the dataset are divided into two parts on average, in which a small part of isolated nodes will appear. For multi-participant, node features are divided averagely to each participant, and the complete topology of the graph is retained. Due to the randomness of the segmentation, the dataset is divided and tested for 10 times. The average accuracy of the trained server model will be recorded.

TABLE II
THE BASIC STATISTICS OF FIVE NETWORK DATASETS

Datasets	#Nodes	#Edges	#Features	#Classes	#Average Degree
Cora [37]	2708	5429	1433	7	2.00
Cora_ML [37]	2810	7981	2879	7	2.84
Citeseer [37]	3327	4732	3703	6	1.42
Pol.Blogs [38]	1222	16714	/	2	13.68
Pubmed [39]	19717	44325	500	3	2.25

B. Local GNN Model

To demonstrate that Graph-Fraudster is effective in various GNN structures based GVFLs, three GNN models are adopted as local participants.

- **Graph convolutional network (GCN) [40]:** it uses a layer-wise propagation rule based on a first-order approximation of spectral convolutions on graphs. For node

classification, a two-layer GCN is adopted as each local GNN model in GVFL. The node is represented as:

$$h_i^{(l+1)} = \rho \left(\sum_{j \in ne(i)} \frac{1}{\sqrt{\tilde{D}_{ii}\tilde{D}_{jj}}} h_j^{(l)} W^{(l)} \right) \quad (24)$$

where $h^{(l+1)}$ is $(l+1)$ -th layer’s node representation and $W^{(l)}$ is the parameters of l -th layer. Node j belongs to neighbor node set ne of node i . $\tilde{D}_{ii} = \sum_j (A + I_N)$ is the degree matrix of $A + I_N$, and I_N is the identity matrix.

- **Simple graph convolution (SGC) [41]:** it is a linearized version of GCN without the activation function. The aggregation function can be formulated as:

$$h_i^{(l+1)} = \sum_{j \in ne(i)} \frac{1}{\sqrt{\tilde{D}_{ii}\tilde{D}_{jj}}} h_j^{(l)} W^{(l)} \quad (25)$$

SGC usually outperforms GCN for node classification.

- **Robust GCN (RGCN) [42]:** it adopts Gaussian distributions as the hidden representations of nodes, which can absorb the effects of adversarial attack:

$$\begin{aligned} \mu_i^{(l+1)} &= \rho \left(\sum_{j \in ne(i)} \frac{1}{\sqrt{\tilde{D}_{ii}\tilde{D}_{jj}}} (\mu_j^{(l)} \odot \alpha_j^{(l)}) W_\mu^{(l)} \right) \\ \sigma_i^{(l+1)} &= \rho \left(\sum_{j \in ne(i)} \frac{1}{\tilde{D}_{ii}\tilde{D}_{jj}} (\sigma_j^{(l)} \odot \alpha_j^{(l)} \odot \alpha_j^{(l)}) W_\sigma^{(l)} \right) \end{aligned} \quad (26)$$

where μ is the means and σ is the variances of Gaussian distributions. W_μ and W_σ are weight matrix of the means and the variances, respectively. \odot is the element-wise product.

C. Attack Baselines

Since this is the first work of adversarial attack on GVFL, three attack methods in the centralized setting are chosen as baselines. To make the attack methods transferable to GVFL, the probabilities returned by the server model and the local data are used to train a surrogate model, which is the same as the local GNN model. The baselines are briefly described as follows.

- **RND [27]:** we assume that the connection of different types of nodes will affect the prediction result. For a given target node, RND randomly samples a node whose predicted label is unequal to the target node. Then, add an edge between the target node and the sampled node.
- **NETTACK [27]:** it generates the adversarial edges guided by two evaluation functions from selected candidate edges and features. It modifies the highest-scoring edges or features, and updates the adversarial network iteratively to confuse GNNs.
- **FGA [29]:** it constructs the edge gradient network based on original network firstly. Then, it generates the adversarial edges iteratively with the maximal absolute edge gradient til it reaches the attack budget Δ .

TABLE III
THE ACCURACY(\pm STD) OF GVFL BASED ON DIFFERENT GNN MODELS AGAINST FOUR ADVERSARIAL ATTACK METHODS ON MULTIPLE DATASETS (IN DUAL-PARTICIPANTS CASE).

Local Model	Datasets	Centralized	Method				
			Clean	RND	NETTACK	FGA	Graph-Fraudster
GCN	Cora	0.806	0.719 \pm 0.017	0.663 \pm 0.015	0.541 \pm 0.043	0.565 \pm 0.051	0.435\pm0.034
	Cora_ML	0.843	0.809 \pm 0.008	0.766 \pm 0.011	0.600 \pm 0.046	0.654 \pm 0.035	0.480\pm0.032
	Citeseer	0.682	0.629 \pm 0.019	0.562 \pm 0.033	0.507 \pm 0.032	0.496 \pm 0.037	0.376\pm0.028
	Pol.Blogs	0.955	0.923 \pm 0.015	0.835 \pm 0.024	0.775 \pm 0.058	0.745 \pm 0.024	0.713\pm0.032
	Pubmed	0.789	0.718 \pm 0.022	0.645 \pm 0.011	0.381 \pm 0.042	0.418 \pm 0.053	0.344\pm0.031
SGC	Cora	0.808	0.722 \pm 0.015	0.669 \pm 0.020	0.544 \pm 0.052	0.564 \pm 0.062	0.448\pm0.033
	Cora_ML	0.848	0.814 \pm 0.011	0.768 \pm 0.013	0.593 \pm 0.049	0.625 \pm 0.034	0.491\pm0.018
	Citeseer	0.689	0.637 \pm 0.021	0.565 \pm 0.032	0.507 \pm 0.028	0.483 \pm 0.049	0.367\pm0.019
	Pol.Blogs	0.954	0.922 \pm 0.015	0.844 \pm 0.027	0.783 \pm 0.053	0.770 \pm 0.039	0.714\pm0.036
	Pubmed	0.788	0.718 \pm 0.019	0.647 \pm 0.012	0.381 \pm 0.033	0.412 \pm 0.055	0.339\pm0.026
RGCN	Cora	0.808	0.725 \pm 0.009	0.671 \pm 0.012	0.518 \pm 0.026	0.581 \pm 0.036	0.454\pm0.027
	Cora_ML	0.845	0.812 \pm 0.010	0.769 \pm 0.016	0.620 \pm 0.071	0.681 \pm 0.068	0.529\pm0.042
	Citeseer	0.674	0.642 \pm 0.013	0.566 \pm 0.016	0.503 \pm 0.037	0.519 \pm 0.029	0.386\pm0.038
	Pol.Blogs	0.954	0.949 \pm 0.005	0.833 \pm 0.026	0.785 \pm 0.074	0.775 \pm 0.049	0.708\pm0.018
	Pubmed	0.789	0.740 \pm 0.013	0.662 \pm 0.018	0.387 \pm 0.037	0.431 \pm 0.044	0.352\pm0.017

D. Experiment Setup

For each local GNN model, a two-layer GNN model is applied to extract the local node embeddings, whose dimension is set to 16. The number of hidden units is fixed to 32. For GCN and SGC, the activation function is ReLU. ELU and ReLU are applied for means and variances respectively in RGCN. The GVFL is trained for 200 epochs using Adam with learning rate of 0.01. Considering sparsity of the dataset and concealment of attack, the attack budget Δ is fixed to 1. Besides, accuracy of the server model is applied as the main evaluation metric, and the lower the accuracy, the better the attacker’s performance.

Our experimental environment consists of Intel XEON 6240 2.6GHz x 18C (CPU), Tesla V100 32GiB (GPU), 16GiB memory (DDR4-RECC 2666) and Ubuntu 16.04 (OS).

E. Attack Performance

In this subsection, Graph-Fraudster is conducted on five real-world datasets in two main scenes, i.e., dual-participants based GVFL and multi-participant based GVFL.

1) *Attack on Dual-participants Based GVFL*: As described above, in this setting, datasets are divided into two pieces both in edges and node features randomly. Multi-perspective metrics are used to evaluate the performance of the attackers, including the node classification accuracy, precision, recall, F1-Score, mean absolute error (MAE) and Log Loss.

As shown in TABLE III, in order to verify the generality of Graph-Fraudster in different graph segmentation situations, we conduct the experiments 10 times and report the average accuracy with standard deviation. The best attack performance is highlighted in bold. In Fig. 5, other five metrics are shown to evaluate the proposed method variously. Specifically,

precision, recall and F1-Score are applied to measure the performance of GVFL. The lower the attacker’s score, the better the attack performance is. MAE and Log Loss are used to evaluate the degree of confusion of the target server model. Higher value represents better performance of attack. Some observations are concluded in this experiment.

- *Graph-Fraudster achieves the SOTA attack performance compared with baselines.* For instance, the accuracy of Graph-Fraudster degrades the baselines more than 10% on Cora, Cora_ML and Citeseer, which shows the superior performance of Graph-Fraudster. For Pol.Blogs, it can be seen that the performance of GVFL drops less than other datasets under attacks. For example, Graph-Fraudster reduces the GVFL’s performance by about 21%, 20.8%, 24.1%, corresponding to GCN, SGC and RGCN respectively, but gains more than 25% on the other datasets. It is caused by the rich relationships of Pol.Blogs which make attacks more difficult within limited perturbation budget. It’s worth mentioning that the standard deviation shows that Graph-Fraudster is more stable than NETTACK and FGA generally. Fig. 5 reports the results on other five metrics, and Graph-Fraudster gains better attack performance than baselines on each metric.
- *The integrity of graph-structured data affects the performance of the model.* Compared with the performance of centralized GNN models, GVFL’s performance drops. It can be analyzed that the local GNN models cannot capture complete node information while edges are assigned to different participants, which weakens the quality of the local node embeddings. Furthermore, GVFL performs better on the graph with higher degree, e.g., the perfor-

		Cora					Cora ML					Citeseer					Pol.Blogs					Pubmed				
GCN	Clean	0.68	0.73	0.69	0.10	1.00	0.80	0.78	0.79	0.07	0.65	0.64	0.64	0.63	0.15	1.14	0.94	0.94	0.94	0.14	0.21	0.75	0.76	0.75	0.21	0.70
	RND	0.62	0.66	0.63	0.12	1.09	0.74	0.72	0.73	0.09	0.76	0.58	0.59	0.58	0.17	1.21	0.80	0.80	0.80	0.23	0.37	0.67	0.68	0.67	0.26	0.83
	NETTACK	0.54	0.47	0.49	0.15	1.75	0.59	0.55	0.54	0.14	1.79	0.48	0.44	0.45	0.19	1.68	0.76	0.76	0.76	0.25	0.48	0.47	0.38	0.41	0.43	3.98
	FGA	0.55	0.57	0.56	0.14	1.25	0.62	0.59	0.59	0.12	1.19	0.55	0.53	0.53	0.18	1.26	0.71	0.71	0.71	0.29	0.56	0.44	0.38	0.40	0.42	3.60
	Graph-Fraudster	0.44	0.39	0.41	0.18	2.09	0.46	0.43	0.43	0.18	2.19	0.41	0.32	0.33	0.23	2.15	0.69	0.69	0.69	0.30	0.64	0.43	0.29	0.31	0.46	4.73
SGC	Clean	0.70	0.72	0.70	0.10	0.95	0.81	0.78	0.80	0.07	0.65	0.64	0.64	0.64	0.14	1.11	0.94	0.93	0.94	0.14	0.21	0.75	0.76	0.76	0.21	0.69
	RND	0.65	0.68	0.66	0.12	1.00	0.75	0.72	0.73	0.09	0.75	0.59	0.60	0.59	0.16	1.18	0.85	0.85	0.85	0.21	0.33	0.67	0.68	0.68	0.26	0.83
	NETTACK	0.58	0.58	0.57	0.14	1.30	0.60	0.57	0.56	0.14	1.72	0.49	0.46	0.47	0.19	1.57	0.78	0.78	0.78	0.24	0.42	0.46	0.39	0.41	0.43	3.79
	FGA	0.57	0.55	0.55	0.14	1.26	0.63	0.58	0.58	0.12	1.19	0.55	0.53	0.53	0.18	1.23	0.76	0.76	0.76	0.27	0.50	0.43	0.38	0.40	0.42	3.47
	Graph-Fraudster	0.47	0.46	0.45	0.17	1.66	0.48	0.45	0.46	0.17	2.01	0.48	0.39	0.40	0.21	1.86	0.68	0.68	0.68	0.33	0.82	0.41	0.29	0.31	0.46	4.66
RGCN	Clean	0.69	0.73	0.71	0.10	0.96	0.83	0.78	0.80	0.07	0.61	0.63	0.63	0.62	0.15	1.18	0.95	0.95	0.95	0.13	0.19	0.74	0.77	0.75	0.20	0.74
	RND	0.63	0.67	0.64	0.12	1.03	0.79	0.73	0.75	0.09	0.69	0.56	0.55	0.55	0.17	1.30	0.85	0.85	0.85	0.21	0.34	0.66	0.69	0.67	0.26	0.83
	NETTACK	0.57	0.58	0.56	0.14	1.28	0.71	0.62	0.64	0.12	1.07	0.47	0.45	0.43	0.20	1.74	0.78	0.78	0.78	0.24	0.41	0.49	0.49	0.49	0.38	2.59
	FGA	0.63	0.63	0.62	0.13	1.09	0.71	0.68	0.69	0.11	0.88	0.55	0.52	0.51	0.18	1.32	0.76	0.76	0.76	0.25	0.44	0.52	0.56	0.52	0.35	2.08
	Graph-Fraudster	0.49	0.48	0.47	0.16	1.61	0.61	0.50	0.52	0.15	1.59	0.38	0.36	0.36	0.22	1.94	0.69	0.68	0.68	0.32	0.75	0.41	0.36	0.37	0.42	2.86
		Precision	Recall	F1-Score	MAE	Log Loss	Precision	Recall	F1-Score	MAE	Log Loss	Precision	Recall	F1-Score	MAE	Log Loss	Precision	Recall	F1-Score	MAE	Log Loss	Precision	Recall	F1-Score	MAE	Log Loss

Fig. 5. Multi-perspective metrics to measure the performance of attacks. For precision, recall, and F1-Score, the lower the attacker’s score, the better the attack performance is, and for MAE and Log Loss, the higher the attacker’s score, the better the attack performance is.

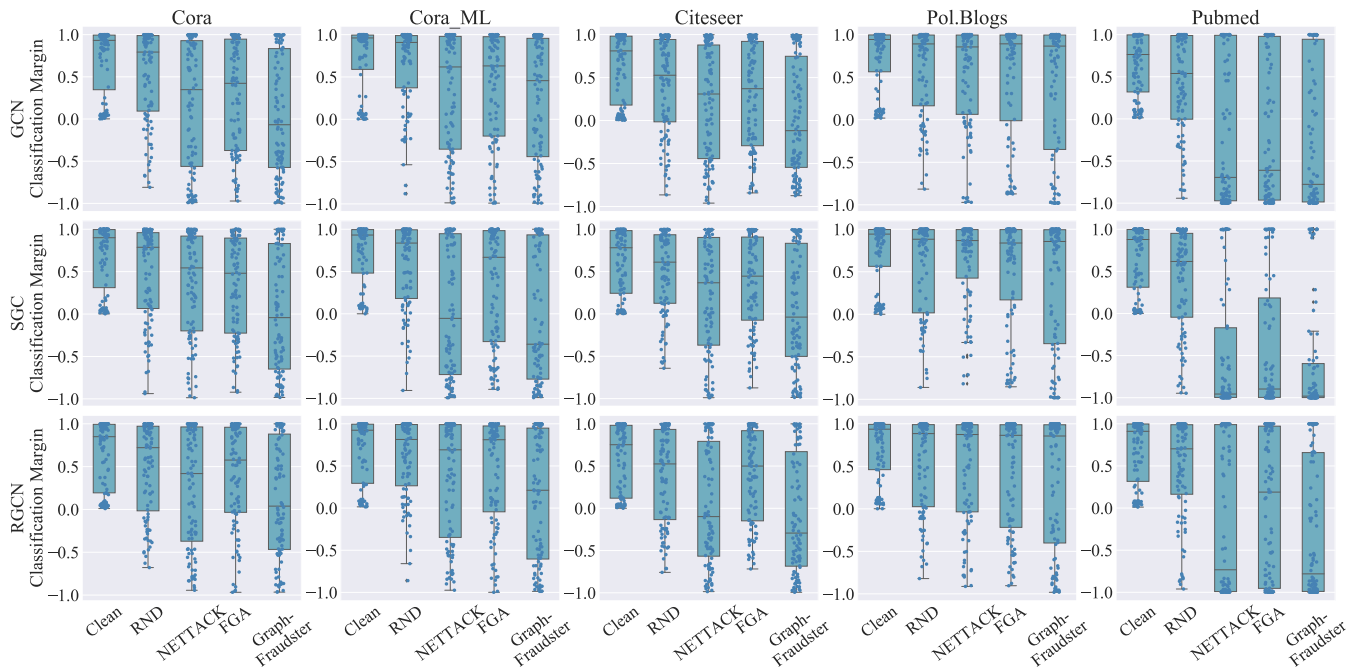


Fig. 6. Classification margin of target nodes predicted by GVFL on five datasets. The lower classification margin, the better attack performance is.

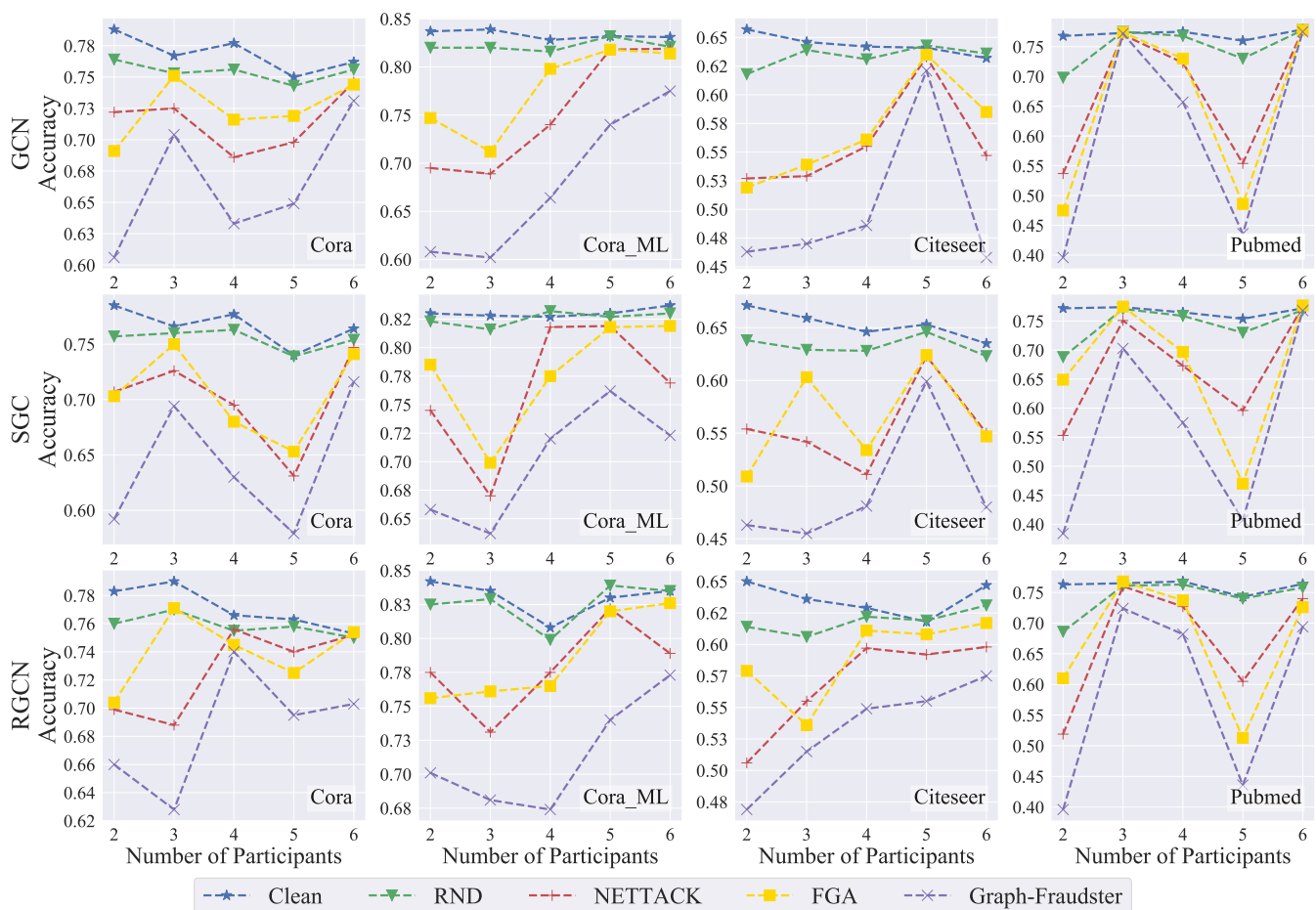


Fig. 7. Attack on multi-participant based GVFL. The x-axis represents the number of participants, and the y-axis represents node classification tasks’ accuracy. Each row represents experiments on different datasets of the same local GNN model.

mance of GVFL based on GCN drops by 8.7% on Cora (the average degree is 2.00) while it only drops by 3.2% on Pol.Blogs (the average degree is 13.68). There are similar results on different GVFLs. The phenomenon can be interpreted as the denser the graph can retain more relational information after data segmentation.

- *Graph-Fraudster outperforms the other baselines on the robust GNN model.* RGCN was proposed as a kind of robust GNN model with defense mechanism. Thus, in most cases, RGCN-based GVFL outperforms other GVFLs without defense mechanism. Despite the performance of attack methods declines, Graph-Fraudster still keeps the state-of-the-art performance. Additionally, benefiting from its randomness, RND always maintains its performance, but its attack performance is poor.

To further compare the performance of attack methods in more details, classification margin is used as a metric to measure the influence of the attack methods on the probabilities of the model’s prediction, following Zügner’s work [27]. We select 100 nodes from the test set, which are correctly classified. These nodes satisfy:

- the 20 nodes with the highest classification margin.
- the 20 nodes with the lowest classification margin but still be classified correctly.

- the 60 nodes are selected randomly.

It can be observed from Fig. 6 that Graph-Fraudster outperforms baselines. More than half of the target nodes’ classification margin below 0, on Cora, Citeseer and Pubmed, corresponding to the higher attack performance of Graph-Fraudster on these datasets. Moreover, more nodes gain low classification margin than NETTACK and FGA, under Graph-Fraudster’s attack, on Cora_ML. Although, for Pol.Blogs, all attack methods cannot achieve the same performance as well as other datasets due to its density. Graph-Fraudster still performs better. On Pubmed, Graph-Fraudster achieves similar performance to NETTACK and FGA, but stays ahead. For RND, only a few nodes obtains low classification margin, which explains its worst attack performance. Additionally, there are still some nodes with higher classification margin, corresponding to nodes that have not been successfully attacked. It may be caused by the single-edge attack, which is unable to make the GVFL misclassify these nodes.

The above experiments suggest that GVFL suffers from adversarial attacks. Although the existing adversarial attacks transferable to GVFL are not powerful enough without the knowledge of GVFL, GVFL still performs as expected. It reveals that GVFL is vulnerable to adversarial attacks. The vulnerability may be caused by the GNN model. The ad-

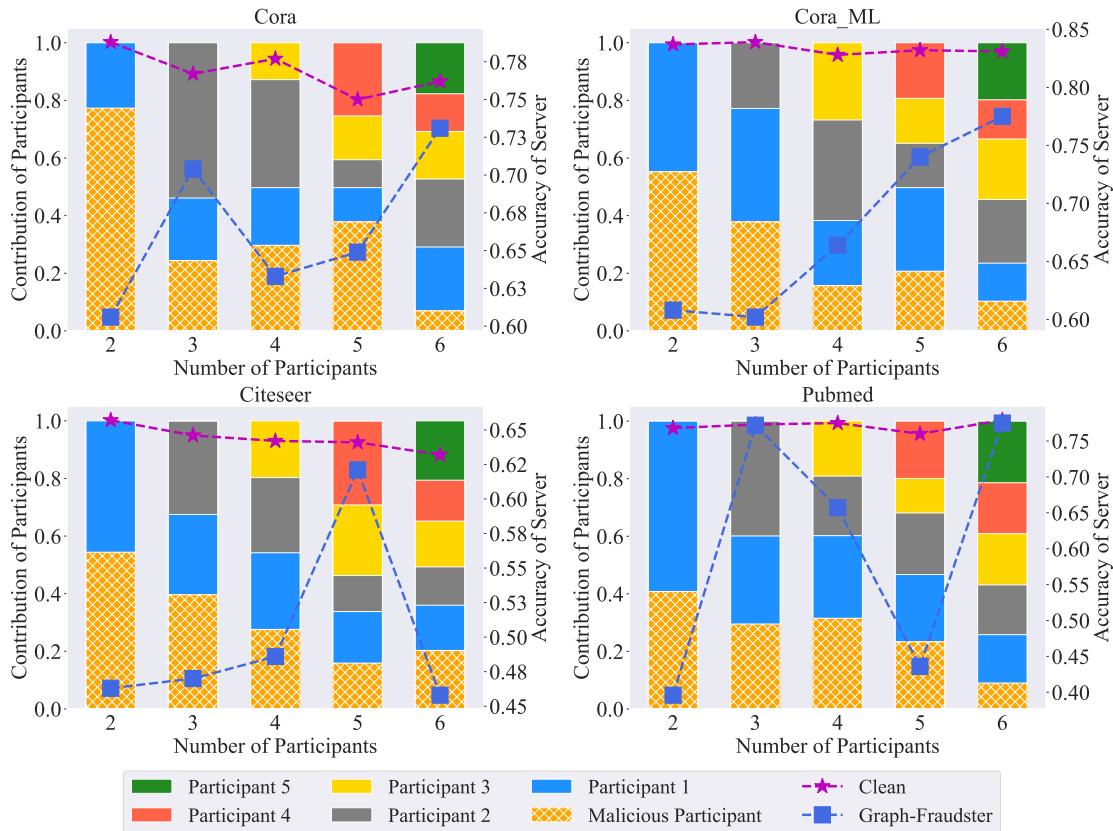


Fig. 8. Relationship between attack performance of Graph-Fraudster and contribution of the malicious participant on GCN-based VFL.

versarial attacks make the GNN model extract the node embedding vector incorrectly and finally confuse the server model. Further, benefiting from the global node embeddings leakage in GVFL, Graph-Fraudster establishes a shadow server model which guides the attack to succeed. More knowledge of GVFL helps the performance of attacks, which explains the superior performance of Graph-Fraudster.

2) *Attack on Multi-participant Based GVFL*: More generally, there are multiple participants in FL. In order to verify the capacity of Graph-Fraudster, the number of participants is set from 2 to 6 in the experiments. To avoid the emergence of numerous isolated nodes, we just distribute node features randomly to each participant without segmentation of edges. Since Pol.Blogs has no node features, the experiments will be conducted on the remaining four datasets. The results are shown in Fig. 7.

First, comparing to edges segmentation, the performance of GVFL is more close to the centralized GNN models in this setting due to the retention of complete relationships between node pairs. Then, as the number of participants increases, the accuracy of GVFL decreases gradually. It is caused by information loss in node features segmentation. Same as edges segmentation setting, Graph-Fraudster maintains optimal attack performance with two participants. In general, the performance of attack methods decreases with the increasing of participants number. Because the server model treats every participant equally, as the number of participants increases, the contribution of the malicious participant's node embeddings

decreases. As a result, the server model is harder to be confused. For example, on Pubmed, all attack methods almost fail with the single-edge attack when the participant number is 6. Interestingly, it not always happens in the experiments. Taking Citeseer as an example, the accuracy of GVFL falls when the number of participants are 6. We give the relationship between attack performance of Graph-Fraudster and contribution of the malicious participant in Fig. 8. Contribution of i -th participant is defined as $C_i = \frac{acc_i}{\sum_j^K acc_{j=1}}$, where K is the number of participants and acc_i is the accuracy of the server model when only participant i 's node embeddings is input. It can be observed that on Cora_ML, when the number of participants increases, the contribution of participants gradually decreases, and the server model is less susceptible to attacks. On Citeseer, when the number of participants is 6, the proportion of malicious participant's contribution increases, causing the server model to be attacked by Graph-Fraudster. We attribute the reason to the segmentation of the dataset. In this setting, the node features are divided randomly, which may cause that not all participants get useful node features. Thus, a small perturbation is powerful enough to fool the server model. On the contrary, the performance drops when the benign participants get some key node features and the slight structural perturbation from malicious participants is not enough to affect the performance of the server model. That is, the model is biased against the data. Besides, the properties of datasets affect attack performance. Taking Citeseer, which is the most

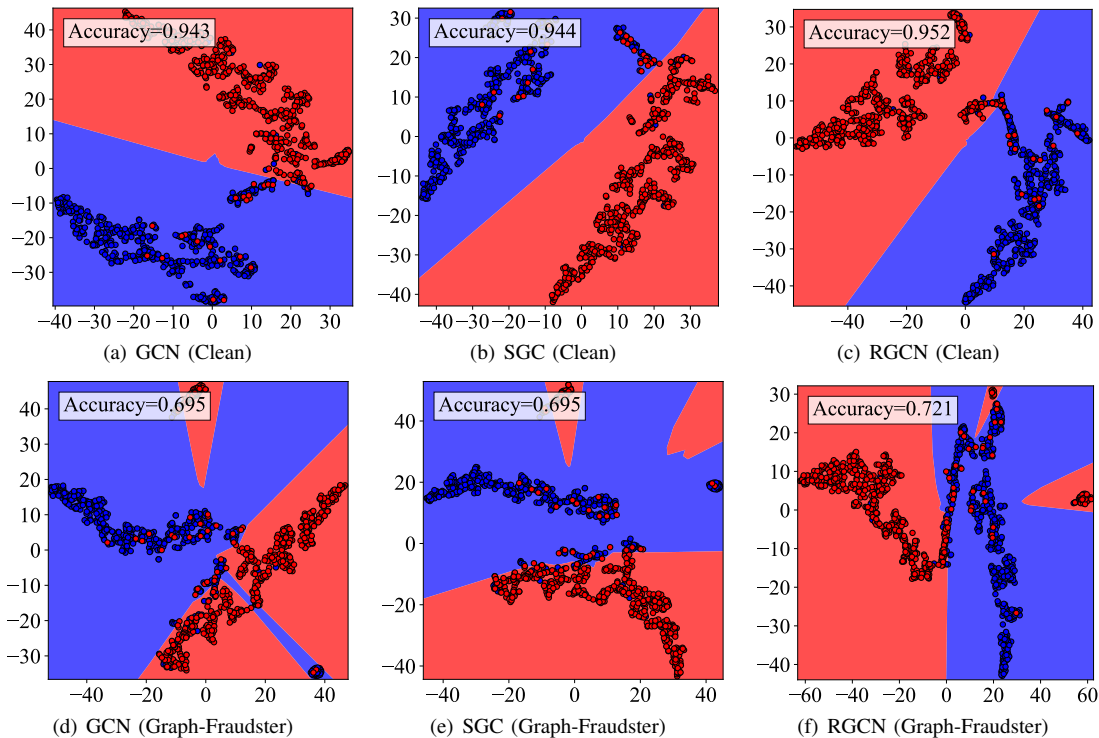


Fig. 9. Decision boundary before (1st row) and after (2nd row) Graph-Fraudster’s attack on Pol.Blogs.

sparse, and it can be seen that Graph-Fraudster maintains good attack performance. Due to its sparseness, it improves the influence of structural attacks on the perturbations of the local node embeddings.

In summary, in the multi-participant case, with the number of participants increases, the performance of attackers generally shows a downward trend. It may be caused by the influence of the server model of each participant decreases while the number of participants increases. Thus, considering adding perturbations into both structure and node features may a better strategy, and it will be considered in our future work. Further, the model’s bias against the data is one of the possible reasons for the success of the adversarial attack. It indicates data bias of the model may lead the server model more rely on high-contributing participants. As a result, it is easier for attacks which initiated by high-contributing participants to cheat the server model. Therefore, setting up a fair evaluation mechanism for participants may help to improve the robustness of GVFL against adversarial attacks.

3) *The Impact of Attacks on Decision Boundary*: To explain the attack results of Graph-Fraudster better, the decision boundaries on Pol.Blogs before and after attack are plotted in Fig. 9. The setting is the same as dual-participants based GVFL. Combining the accuracy of the server model, the decision boundary is regular without perturbations, which corresponds to high accuracy. After attack, as mentioned in the previous experiments, the server model is not easy to be perturbed with high accuracy. However, some nodes break away from their corresponding groups, which makes the decision boundary fracture. In other words, these outliers will be misclassified.

F. Possible Defense

To mitigate the performance of Graph-Fraudster, Differential Privacy (DP) mechanism and Top- k mechanism are introduced as the possible defenses. DP is used to prevent Graph-Fraudster from stealing the node embeddings, and Top- k devotes to filtering out perturbations in the local node embeddings. The two strategies are briefly introduced as follows:

- **DP [43]**: randomized noise is injected into the local node embeddings before uploading, obeying the principle of DP. The Laplacian mechanism is used in this paper with the scale of noise β .
- **Top- k [44]**: before uploading the local node embeddings, we sort the embedding value and preserve the top k values to reduce the influence of some useless information.

The scale of Laplacian noise β is set from 0 to 0.5, and the value of k is set from 8 to 16. Fig. 10 shows the result on GCN-based GVFL. Note that $\beta = 0$ and $k = 16$ present the GVFL without defense mechanism. With the increase of β or the decrease of k , these two defense mechanisms reduce the performance of GVFL to a certain extent. It shows that there should be a trade-off between defense performance and GVFL’s performance. As shown in Fig. 10(a), DP has limited defensive capability against Graph-Fraudster, or even useless, e.g., the performance on Citeseer. It may be that the noise added by DP cannot offset the noise added by Graph-Fraudster well. As for Top- k , as shown in Fig. 10(b), it is more effective than DP, but it is unstable. For example, on Cora, Top- k with $k = 8$ outperforms other value and $k = 12$ is second, which is not monotonous with the change of k . It is caused by the non-concentrated distribution of perturbations in the node embeddings. Specifically, the perturbations do not only exist in

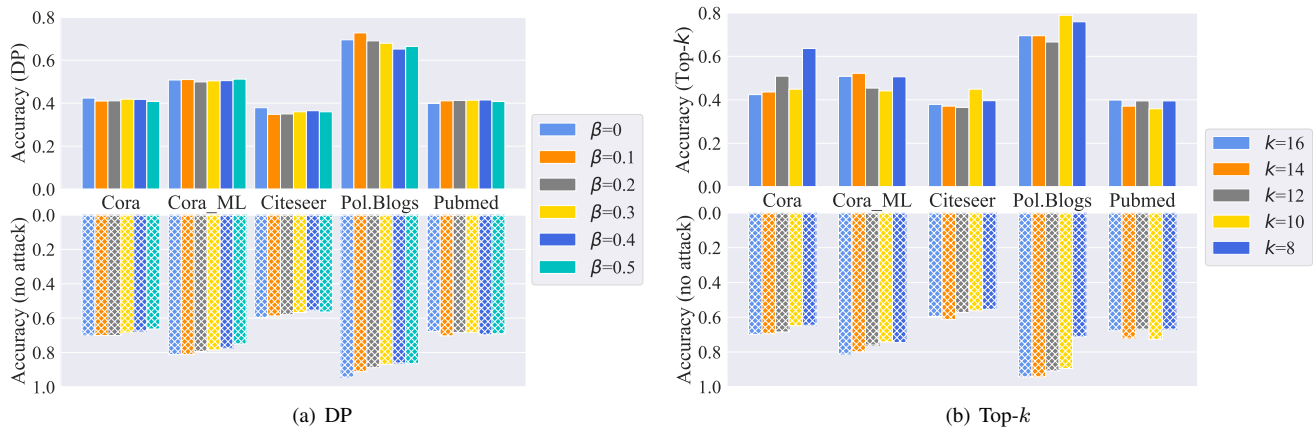


Fig. 10. Possible defense against Graph-Fraudster on GCN-based VFL. The upper part is the attack performance of Graph-Fraudster under the defense mechanisms, and the lower part is the performance of GVFL without attack.

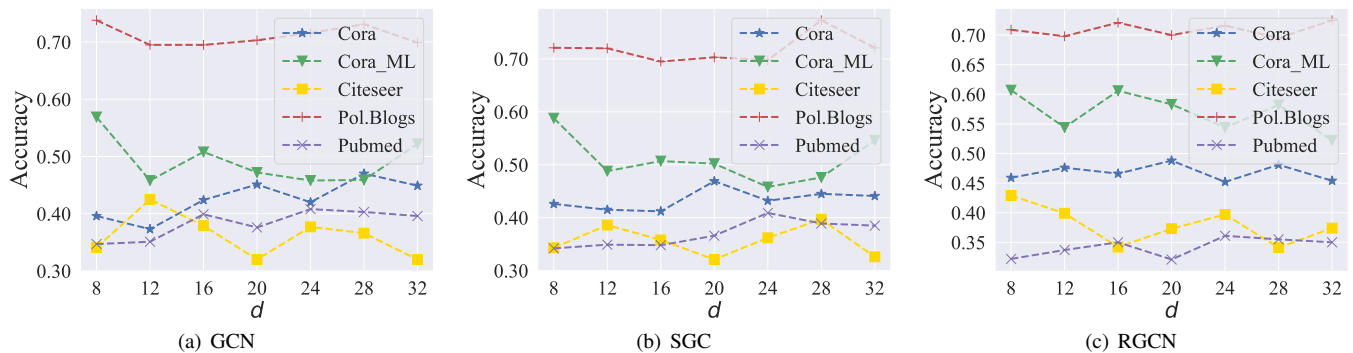


Fig. 11. Performance of Graph-Fraudster with different dimensions of the local node embeddings.

lower values. Moreover, the price of better defensive capability against Graph-Fraudster is the performance degradation of the main task. In summary, Graph-Fraudster can still perform well under two possible defense mechanisms. In future works, filtering perturbations from the local node embeddings is a feasible solution against adversarial attacks in GVFL, and keeping the performance of GVFL while gains defensive capability is also a great challenge.

G. Parameter Analysis

In this subsection, the sensitivity of the dimensions of the local node embeddings d and the scale of noise ϵ for Graph-Fraudster is explored. The value of these parameters is adjusted to see how they affect the performance of Graph-Fraudster. In more detail, d is set from 8 to 32 with a step size of 4, and ϵ is changed from 0 to 0.01.

The performance change of Graph-Fraudster is shown in Fig. 11 and Fig. 12. It can be observed that as the value of d increases, the attack performance of Graph-Fraudster is not significantly affected (less than 10%), which shows stability of Graph-Fraudster. In addition, we take the performance of Graph-Fraudster on GCN-based VFL with different noise scale ϵ as an example. The results are shown in Fig. 12. Similar results are observed in other settings. It can be seen that the accuracy of the server model drops rapidly in a small range of ϵ (from 0 to 0.002). Single-edge attack may be the main reason, because the perturbations are too small to further approximate the node embeddings with large-scale noise. Thus, for single-

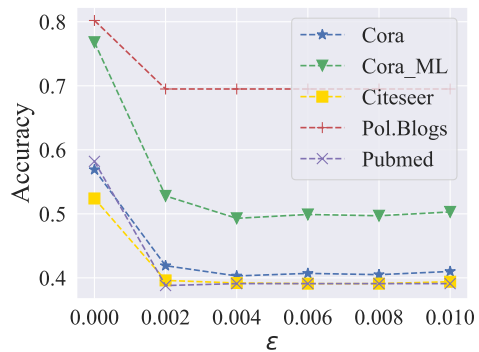


Fig. 12. Performance of Graph-Fraudster with different noise scale ϵ on GCN-based VFL.

edge attack, a small value of ϵ is enough (e.g., 0.004 in this setting) to guide the attack, and choosing an appropriate value of ϵ can boost the performance of Graph-Fraudster.

V. CONCLUSION

Contributions. The first novel adversarial attack method on GVFL, named Graph-Fraudster, is proposed by combining GVFL's privacy leakage and the gradient of pairwise node. Benefiting from the information of GVFL, Graph-Fraudster achieves the state-of-the-art attack performance compared with the baselines. Extensive experiments have testified the effectiveness and efficiency of Graph-Fraudster and further

revealed the vulnerability of GVFL even under defensive circumstances.

Observations. Besides of the adversarial attack problem formulation and Graph-Fraudster attack, several observations are gained according to the comprehensive experiments and analysis. (1) GVFL does suffer from adversarial attacks, especially when some information (e.g., node embeddings, structure of model and gradient of model etc.) is explored by the attackers; (2) graph splitting strategy of GVFL will not only affect the performance of GVFL, but also affect the success rate of adversarial attack, thus a proper splitting strategy should be seriously considered in real-world scenarios; (3) the more a participant contributes to the server, the more likely the perturbations added to the participant is to construct a successful adversarial attack; (4) although DP and Top- K fail to fully defend Graph-Fraudster, the results show that denoising the uploaded node embeddings is still a feasible defensive strategy.

Limitations. Graph-Fraudster devotes to fooling the server model, even with single-edge modifying. However, the small number of perturbations does not mean that the perturbations are imperceptible, and these perturbations may change some properties of the graph significantly. Besides, frequent queries are needed for Graph-Fraudster, which are expensive and detectable.

Future Works. Adversarial attack on GVFL will raise our attention about the robustness of GVFL. There are two sides of suggestions are provided for future works. **Recommendations for vulnerability mining:** (1) For perturbation strategies, both structure and node features should be considered for perturbation, to reduce the impact of the data imbalance of attacks. (2) Since frequent queries on the server model are expensive and abnormal, reducing queries by stealing more information from GVFL is a possible solution. (3) Exploring the concealment of the perturbations for diverse downstream tasks, such as link prediction, graph classification, should be discussed as well. **Suggestions for defenders:** (1) Researchers are advised to focus on the information leakage problem on GVFL, which provides convenience for attackers. Preventing key information of GVFL from privacy inferencing is a direct and efficient method. (2) A fair evaluation agency for participants should be established. Over-reliance on high-contributing participants may bring potential security risks, i.e., the impact of adversarial attacks may be amplified by the data bias of the model. (3) Some defense or detection mechanisms can be deployed in terminals or servers to eliminate the impact of adversarial perturbations.

REFERENCES

- [1] H. B. McMahan, E. Moore, D. Ramage, and B. A. y Arcas, "Federated learning of deep networks using model averaging," *arXiv preprint arXiv:1602.05629*, 2016.
- [2] Q. Yang, Y. Liu, T. Chen, and Y. Tong, "Federated machine learning: Concept and applications," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 10, no. 2, pp. 1–19, 2019.
- [3] T. Yang, G. Andrew, H. Eichner, H. Sun, W. Li, N. Kong, D. Ramage, and F. Beaufays, "Applied federated learning: Improving google keyboard query suggestions," *arXiv preprint arXiv:1812.02903*, 2018.
- [4] S. Ge, F. Wu, C. Wu, T. Qi, Y. Huang, and X. Xie, "Fedner: Medical named entity recognition with federated learning," *arXiv preprint arXiv:2003.09288*, 2020.
- [5] G. Long, Y. Tan, J. Jiang, and C. Zhang, "Federated learning for open banking," in *Federated Learning*. Springer, 2020, pp. 240–254.
- [6] A. Suruliandi, A. Kasthuri, and S. Raja, "Deep feature representation and similarity matrix based noise label refinement method for efficient face annotation," *International Journal of Interactive Multimedia and Artificial Intelligence*, vol. 7, no. 2, pp. 66–78, 2021.
- [7] A. Kasthuri, A. Suruliandi, and S. Raja, "Gabor-oriented local order feature-based deep learning for face annotation," *International Journal of Wavelets, Multiresolution and Information Processing*, vol. 17, no. 05, p. 1950032, 2019.
- [8] J. Konečný, H. B. McMahan, D. Ramage, and P. Richtárik, "Federated optimization: Distributed machine learning for on-device intelligence," *arXiv preprint arXiv:1610.02527*, 2016.
- [9] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial Intelligence and Statistics*. PMLR, 2017, pp. 1273–1282.
- [10] M. Mohri, G. Sivek, and A. T. Suresh, "Agnostic federated learning," in *International Conference on Machine Learning*. PMLR, 2019, pp. 4615–4625.
- [11] M. Yurochkin, M. Agarwal, S. Ghosh, K. Greenewald, N. Hoang, and Y. Khazaeni, "Bayesian nonparametric federated learning of neural networks," in *International Conference on Machine Learning*. PMLR, 2019, pp. 7252–7261.
- [12] V. Mugunthan, A. Peraire-Bueno, and L. Kagal, "Privacyfl: A simulator for privacy-preserving and secure federated learning," in *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, 2020, pp. 3085–3092.
- [13] L. Zheng, J. Zhou, C. Chen, B. Wu, L. Wang, and B. Zhang, "Asfgnn: Automated separated-federated graph neural network," *arXiv e-prints*, pp. arXiv–2011, 2020.
- [14] B. Wang, A. Li, H. Li, and Y. Chen, "Graphfl: A federated learning framework for semi-supervised node classification on graphs," *arXiv preprint arXiv:2012.04187*, 2020.
- [15] C. He, K. Balasubramanian, E. Ceyani, Y. Rong, P. Zhao, J. Huang, M. Annavararam, and S. Avestimehr, "Fedgraphnn: A federated learning system and benchmark for graph neural networks," *arXiv preprint arXiv:2104.07145*, 2021.
- [16] J. Zhou, C. Chen, L. Zheng, X. Zheng, B. Wu, Z. Liu, and L. Wang, "Privacy-preserving graph neural network for node classification," *arXiv preprint arXiv:2005.11903*, 2020.
- [17] X. Ni, X. Xu, L. Lyu, C. Meng, and W. Wang, "A vertical federated learning framework for graph convolutional network," *arXiv preprint arXiv:2106.11593*, 2021.
- [18] C. Wu, F. Wu, Y. Cao, Y. Huang, and X. Xie, "Fedgmn: Federated graph neural network for privacy-preserving recommendation," *arXiv preprint arXiv:2102.04925*, 2021.
- [19] B. Hitaj, G. Ateniese, and F. Perez-Cruz, "Deep models under the gan: information leakage from collaborative deep learning," in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, 2017, pp. 603–618.
- [20] M. Nasr, R. Shokri, and A. Houmansadr, "Comprehensive privacy analysis of deep learning: Passive and active white-box inference attacks against centralized and federated learning," in *2019 IEEE symposium on security and privacy (SP)*. IEEE, 2019, pp. 739–753.
- [21] L. Zhu and S. Han, "Deep leakage from gradients," in *Federated Learning*. Springer, 2020, pp. 17–31.
- [22] X. Luo, Y. Wu, X. Xiao, and B. C. Ooi, "Feature inference attack on model predictions in vertical federated learning," *arXiv preprint arXiv:2010.10152*, 2020.
- [23] S. Zhang, L. Xiang, X. Yu, P. Chu, Y. Chen, C. Cen, and L. Wang, "Privacy-preserving federated learning on partitioned attributes," *arXiv preprint arXiv:2104.14383*, 2021.
- [24] O. Li, J. Sun, X. Yang, W. Gao, H. Zhang, J. Xie, V. Smith, and C. Wang, "Label leakage and protection in two-party split learning," *arXiv preprint arXiv:2102.08504*, 2021.
- [25] C. Iwendí, S. U. Rehman, A. R. Javed, S. Khan, and G. Srivastava, "Sustainable security for the internet of things using artificial intelligence architectures," *ACM Transactions on Internet Technology (TOIT)*, vol. 21, no. 3, pp. 1–22, 2021.
- [26] U. Ahmed, J. C.-W. Lin, and G. Srivastava, "Generative ensemble learning for mitigating adversarial malware detection in iot," in *2021 IEEE 29th International Conference on Network Protocols (ICNP)*. IEEE, 2021, pp. 1–5.
- [27] D. Zügner, A. Akbarnejad, and S. Günnemann, "Adversarial attacks on neural networks for graph data," in *Proceedings of the 24th ACM*

- SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018, pp. 2847–2856.
- [28] H. Dai, H. Li, T. Tian, X. Huang, L. Wang, J. Zhu, and L. Song, “Adversarial attack on graph structured data,” *arXiv preprint arXiv:1806.02371*, 2018.
- [29] J. Chen, Y. Wu, X. Xu, Y. Chen, H. Zheng, and Q. Xuan, “Fast gradient attack on network embedding,” *arXiv preprint arXiv:1809.02797*, 2018.
- [30] H. Wu, C. Wang, Y. Tyshetskiy, A. Docherty, K. Lu, and L. Zhu, “Adversarial examples on graph data: Deep insights into attack and defense,” *arXiv preprint arXiv:1903.01610*, 2019.
- [31] J. Li, T. Xie, C. Liang, F. Xie, X. He, and Z. Zheng, “Adversarial attack on large scale graph,” *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [32] X. Wang, M. Cheng, J. Eaton, C.-J. Hsieh, and F. Wu, “Attack graph convolutional networks by adding fake nodes,” *arXiv preprint arXiv:1810.10751*, 2018.
- [33] H. Chang, Y. Rong, T. Xu, W. Huang, H. Zhang, P. Cui, W. Zhu, and J. Huang, “A restricted black-box adversarial framework towards attacking graph embedding models,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 3389–3396.
- [34] Y. Ma, S. Wang, T. Derr, L. Wu, and J. Tang, “Attacking graph convolutional networks via rewiring,” *arXiv preprint arXiv:1906.03750*, 2019.
- [35] J. Chen, L. Chen, Y. Chen, M. Zhao, S. Yu, Q. Xuan, and X. Yang, “Ga-based q-attack on community detection,” *IEEE Transactions on Computational Social Systems*, vol. 6, no. 3, pp. 491–503, 2019.
- [36] I. J. Goodfellow, J. Shlens, and C. Szegedy, “Explaining and harnessing adversarial examples,” *arXiv preprint arXiv:1412.6572*, 2014.
- [37] A. K. McCallum, K. Nigam, J. Rennie, and K. Seymore, “Automating the construction of internet portals with machine learning,” *Information Retrieval*, vol. 3, no. 2, pp. 127–163, 2000.
- [38] L. A. Adamic and N. Glance, “The political blogosphere and the 2004 us election: divided they blog,” in *Proceedings of the 3rd international workshop on Link discovery*, 2005, pp. 36–43.
- [39] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Galligher, and T. Eliassirad, “Collective classification in network data,” *AI magazine*, vol. 29, no. 3, pp. 93–93, 2008.
- [40] T. N. Kipf and M. Welling, “Semi-supervised classification with graph convolutional networks,” *arXiv preprint arXiv:1609.02907*, 2016.
- [41] F. Wu, A. Souza, T. Zhang, C. Fifty, T. Yu, and K. Weinberger, “Simplifying graph convolutional networks,” in *International conference on machine learning*. PMLR, 2019, pp. 6861–6871.
- [42] D. Zhu, Z. Zhang, P. Cui, and W. Zhu, “Robust graph convolutional networks against adversarial attacks,” in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019, pp. 1399–1407.
- [43] C. Dwork, “Differential privacy: A survey of results,” in *International conference on theory and applications of models of computation*. Springer, 2008, pp. 1–19.
- [44] R. Liu, Y. Cao, M. Yoshikawa, and H. Chen, “Fedsel: Federated sgd under local differential privacy with top-k dimension selection,” in *International Conference on Database Systems for Advanced Applications*. Springer, 2020, pp. 485–501.