# A Appendix

## A.1 About *DGraph*

During the data collection stage of *DGraph* construction, the behavior data collected by *Finvolution Group* and the type of data have been approved by the user. The following is a detailed description of our agreement (***"User Privacy Protection Policy"***[7], Article 1):

- *"When you start the Paipaidai (Finvolution) loan service, you need to perform real-name verification. We will collect your name, mobile phone number, ID number, ID photo..."*
- *"When you apply for Paipaidai (Finvolution) to evaluate the loan amount, you need to provide your personal information for credit extension. You need to provide the following necessary information: ...emergency contact information..."*

From an ethical point of view, we follow the GDPR data minimization principle, and all data used by *DGraph* is necessary for the platform's anti-fraud algorithm.

It is worth noting that *DGraph* strictly follows the ***"Personal Information Protection Law of the People's Republic of China"***[8] as all the raw data of *DGraph* is collected within China. Meanwhile, we also followed General Data Protection Regulation (GDPR). The data we disclose has equipped with a strict encryption algorithm to ensure that the data is disclosed in an anonymous way. Anonymized data is defined by the "Personal Information Protection Law of the People's Republic of China" at:

- Article 4, *"Personal information refers to various information related to identified or identifiable natural persons and recorded electronically or in other ways, excluding anonymized information."*
- Article 73, *"The meanings of the following terms in this Law: (3) De-identification refers to the process in which personal information is processed so that it cannot identify a specific natural person without the aid of additional information. (4) Anonymization refers to the process in which personal information cannot identify a specific natural person and cannot be recovered after processing."*

In terms of the specific implementation, we anonymize the user by deleting the personal identification and randomizing the user order. The user ID thus cannot be traced back. Moreover, since the user features are not unique, users cannot be identified through the data set, which ensures the anonymity of *DGraph*. So, the concerns of the GDPR, such as the correction and deletion requirements from users, will not affect *DGraph* (as no one can trace or recognize anyone's data in *DGraph*, and *DGraph* does not tie to individuals' right). Even so, *Finvolution Group* still will process data strictly in accordance with the user's data rights.

Besides, we will continue to track the use of *DGraph* more closely. Currently, before downloading the data, the user must provide his or her name and email address, and confirm that he or she has read and agreed with the user license, which ensures the non-commercial, unethical, unfair research, or causing negative social impact usage of *DGraph*. The license of *DGraph* can be found at: https://dgraph.xinye.com/clause.

## A.2 Explanations of the observation

The results in observation are consistent with common sense overall. Fraudsters' primary purpose is to defraud from the platform, which motivates them to exhibit a variety of abnormal traits. We can explain the results of the observation (Sec. 3.2) from this point, as follows:

- Fig. 2 (a). A lower average out-degree indicates that fraudsters tend to fill fewer emergency contacts in general. This phenomenon coincides with their purpose because filling more emergency contacts are helpless in defrauding money.
- Fig. 2 (b). The emergency contact (EC) relationship is a kind of social connection and the literature suggests that users with social connections are more similar. However, as fraudsters could provide the platform with a false list of emergency contacts to avoid being caught, they may not have social connections with the emergency contacts they filled. As a result, according to [25], the average feature similarity of fraudsters' out-edges is lower than that of normal users.

---

[7]User Privacy Protection Policy of Finvolution Group.https://loancontract.ppdai.com/latest/agency/privacy_policy.html

[8]Personal Information Protection Law of the People's Republic of China. https://www.npc.gov.cn/npc/c30834/202108/a8c4e3672c74491a80b53a172bb753fe.shtml

- Fig. 2 (c). Since their purpose is to borrow money as soon as possible, fraudsters will not carefully fill out optional items in their personal profile, which will cause their node features to have a large percentage of missing values.
- Fig. 2 (d). This result suggests that some fraudsters may fill in multiple emergency contacts, but they often fill in their emergency contact within a short time. This is in line with their purpose. Adding more emergency contact are helpless for defrauding money.

Note that the background node is a new concept proposed by this paper and they are a unique component of *DGraph*. This paper only makes a preliminary exploration of them (in Sec. 4.3 and Sec. 5.3) to show their properties (values). We anticipate that subsequent research could focus on the background node and offer novel discoveries.

### A.3 Experiment details

#### A.3.1 Data splitting

We randomly divide the nodes of *DGraph* into training/validation/test sets, with a split of 70/15/15, respectively. We fix this split and provide it in the public dataset.

#### A.3.2 Methods

**Baseline methods.** We select MLPs as the baseline methods. Its' input is the node feature.

**General graph models.** We evaluate 4 general graph models on *DGraph*, which are Node2Vec, GCN, SAGE, and TGAT. Specifically, Node2Vec only utilizes graph structure information. GCN and SAGE can utilize both structure information and node features. TGAT is a general dynamic GNNs. It can handle dynamic edges by a time encoder. Since node time is required by TGAT, we define node time as the min time of a node's out-edges. It can approximately represent the earliest activation time of a user on Finvolution, because the user must fill in emergency contact as soon as he or she starts the first loan.

**Supervised GAD methods.** We evaluate four anomaly detection methods, which are DevNet, CARE-GNN, PC-GNN and AMNet. All of them have special components to handle the extreme imbalance of samples. Among them, DevNet is similar to MLPs, in which input is only node features. Other methods are GNNs-based methods, which can both utilize structure information and node features. Table 5 summarize the difference of the above methods.

**Unsupervised methods.** We evaluate 7 unsupervised anomaly detection methods, They are: SCAN, MLPAE, GCNAE, Radar, DOMINANT, GUIDE, and OCGNN. And we implement them by PyGOD.

#### A.3.3 Setup

We optimize each model's hyper-parameters based on their AUC performance on the validation set. For all experiments, the number of epochs is set to 1000 except for Node2Vec, where the model is pre-trained for 600 epochs to get the nodes embedding that is further used to train MLPs for 1000 epochs to classify the nodes. To evaluate the models, we repeat all the experiments for five runs and take the average performance. For anomaly detection methods, we use source code provided by their authors and modify hyper-parameters in accordance with their instructions. Since the imbalanced class, we search the class weight of loss function range from [1:1,1:25,1:50,1:100] for general methods, excluding the search of general hyper-parameter settings (such as hidden size). We report the AUC and AP for each model on the test set. Our experiments are conducted in Python3 on a Dell PowerEdge T640 with 48 CPU cores and 1 Tesla P100 GPU. More details can see https://github.com/hxttkl/DGraph_Experiments.

### A.4 Ablation Studies

To further demonstrate the effects of network structure (emergency contact) in experiments, we supply an ablation study on different graph structures. They are:

- **Only Self-loops**. In this structure, there are not any connections between different nodes. We add self-loops for each node. In this setting, GCNs is somehow like MLPs.

Table 5: Summary of supervised methods. ✓denotes a method have a particular component to handle a specific factor.

| Method | Structure | Neighbor | Dynamics | Direction | Anomaly | MV | BN |
|---|---|---|---|---|---|---|---|
| MLP | | | | | | | |
| Node2Vec | ✓ | | | | | | |
| GCN | ✓ | ✓ | | | | | |
| SAGE | ✓ | ✓ | | | | | |
| TGAT[1] | ✓ | ✓ | ✓ | | | | |
| DevNet[2] | | | | | ✓ | | |
| CARE-GNN[3] | ✓ | ✓ | | | ✓ | | |
| PC-GNN [4] | ✓ | ✓ | | | ✓ | | |
| AMNet[5] | ✓ | ✓ | | | ✓ | | |

[1]`https://github.com/StatsDLMathsRecomSys/`
`Inductive-representation-learning-on-temporal-graphs`
[2]`https://github.com/GuansongPang/deviation-network`
[3]`https://github.com/YingtongDou/CARE-GNN`
[4]`https://github.com/PonderLY/PC-GNN`
[5]`https://github.com/zjunet/AMNet`

Table 6: Comparison of different network structures of *DGraph* based on GCN.

| Structure | Validation | | Test | |
|---|---|---|---|---|
| | AUC | AP | AUC | AP |
| Only Self-loops | $0.716_{\pm0.001}$ | $0.026_{\pm0.000}$ | $0.722_{\pm0.002}$ | $0.027_{\pm0.000}$ |
| Random Network | $0.668_{\pm0.001}$ | $0.021_{\pm0.000}$ | $0.666_{\pm0.001}$ | $0.022_{\pm0.000}$ |
| KNN Network | $0.719_{\pm0.000}$ | $0.027_{\pm0.000}$ | $0.726_{\pm0.001}$ | $0.028_{\pm0.001}$ |
| **Emergency contact** | $\mathbf{0.746_{\pm0.001}}$ | $\mathbf{0.035_{\pm0.000}}$ | $\mathbf{0.751_{\pm0.002}}$ | $\mathbf{0.037_{\pm0.000}}$ |

- **Random Networks**. In this structure, each node is randomly connected with others. For comparison, we limit the edge number to as same as the original graph structure.
- **K-Nearest Neighbor (KNN) Networks**. In this structure, similar nodes trends to have edges. Considering the node number, we first randomly generate 100,000,000 edges, the score of each edge is the cosine similarity of its end node features. Then we preserve the 4,300,999 (as same as the original structure) top-scoring edges and filter others.

Then, we conduct experiment for these different network based on GCN. The result is shown in Table 6. Based on emergency contact network, GCN achieves the based performance comparing with other network structure. This result demonstrates again that emergency contact have high correlation with fraud behaviors. How to model emergency contact structure is a key factor on *DGraph* for GAD methods.