

OpenWGL: Open-World Graph Learning

Man Wu ^{*}, Shirui Pan [†], Xingquan Zhu ^{*}

^{*}Dept. of Computer & Electrical Engineering and Computer Science, Florida Atlantic University, Boca Raton, USA

[†]Faculty of Information Technology, Monash University, Melbourne, Australia

mwu2019@fau.edu, shirui.pan@monash.edu, xqzhu2@fau.edu

Abstract—In traditional graph learning tasks, such as classification, learning is carried out in a *closed-world* where the number of classes and their training samples provided to help train models, and the learning goal is to correctly classify unlabeled nodes into classes already known in reality, due to limited labeling capability and dynamic of networks, some nodes in the networks may not belong to existing/seen classes, and therefore cannot be correctly classified by closed-world learning algorithms. In this paper, we propose a new *open-world* graph learning paradigm, where the goal is to not only classify nodes belonging to seen classes into correct groups, but also classify nodes not belonging to seen classes to an unseen class. The essential challenge of open-world graph learning is that (1) unseen class has no samples, and may exist in an arbitrary form different from existing seen classes; and (2) both graph feature learning and node prediction should differentiate whether a node may belong to an existing/seen class or an unseen class. To tackle the challenges, we propose an uncertain node representation learning approach, using constrained variational graph autoencoder networks, where the label loss and class uncertainty loss constraints are used to ensure that the node representation learning are sensitive to unseen class. As a result, node embedding features are denoted by distributions, instead of deterministic feature vectors. By using a sampling process to generate multiple versions of feature vectors, we are able to test the certainty of a node belonging to seen classes, and automatically determine a threshold to reject nodes not belonging to seen classes as unseen class nodes. Experiments on real-world networks demonstrate the algorithm performance, comparing to baselines. Case studies and ablation analysis also show the rationale of our design for open-world graph learning.

Keywords—Graph neural network, open-world learning, node classification

I. INTRODUCTION

Networks/Graphs are convenient tools to model interactions and interdependencies between large-scale data. Graph learning, such as node classification, attempts to categorize nodes of graphs into several groups. Such learning tasks are fundamental, but challenging, and have received continuous attention in the research field. Many research efforts have been made to develop reliable and efficient algorithms for different types of node classification tasks. However, existing methods mainly carry out the learning in a “closed-world” setting, where nodes in the test data must belong to one or multiple classes already seen in the training set. As a result, when a new/unseen class node appears in the test set, classifiers cannot detect the new/unseen class and will erroneously classify the node as seen/learned classes in the training data.

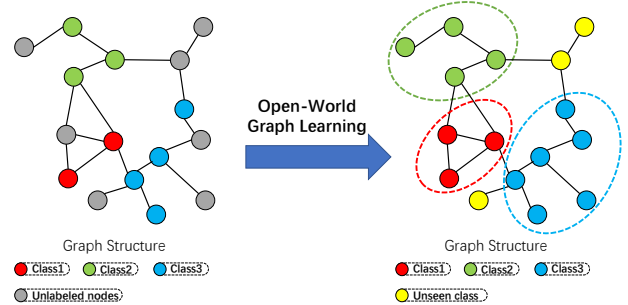


Fig. 1: An example of open-world learning for network node classification. Given a graph with some labeled nodes and unlabeled nodes (left panel), open-world graph learning aims to learn a classifier to classify unlabeled nodes belonging to seen class into its own class, and also detects unlabeled nodes not belonging to any seen class as unseen class nodes.

In reality, new trends emerge constantly and a model that cannot detect these new/unseen trends can hardly work well in a dynamic environment. This problem/phenomenon is referred to as the *open-world classification* or *open classification* problem [1]. The new “open-world” learning (OWL) [2]–[4] paradigm is to not only recognize objects belonging to the classes already seen/learned before, but also detects new class samples which are previously unseen.

Several approaches, such as one-class SVM [5], can be adjusted to address open-world learning by treating all seen classes as the positive class, but they cannot differentiate instances in seen classes and often have poor performance to find unseen class, because no negative data is used. To date, open-world learning has already attracted many interests from Natural Language Processing (NLP) [1] [2] and computer vision fields [6] [7] [8]. In NLP, Shu et al. [1] proposed the solution to open-world learning by setting thresholds before the sigmoid function to reject unseen classes. In computer vision, Scheirer et al. [6] studied the problem of recognizing unseen images that are not in the training data by reducing the half-space of a binary SVM classifier with a positive region. However, to the best of our knowledge, the open-world classification problem has not been previously investigated in graph structure data and graph learning tasks.

Given a graph consisting of seen class and unseen class nodes, the objective of open-world graph learning is to build a classifier to classify nodes belonging to seen classes into correct categories, and also detect nodes not belonging to any

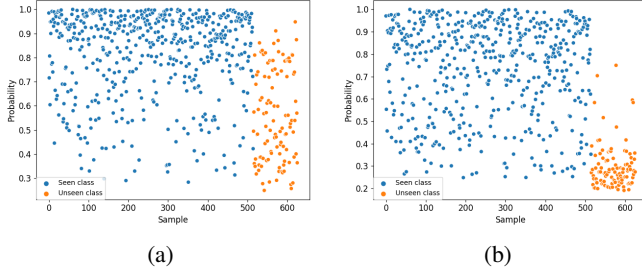


Fig. 2: A visualization of classification probability on seen (blue) and unseen (orange) class test instances. The x -axis denotes the index of test instances (first 500 instances belong to seen classes and the last 100 instances belong to unseen class). The y -axis denotes the maximum probability output of each instance through the softmax classifier. (a) denotes the classification probabilities using only label loss, and (b) denotes the classification probabilities combining both label loss and class uncertainty loss.

seen class as unseen class. An example of open-world learning for graph node classification is illustrated in Fig. 1.

Currently, existing solutions to open-world learning are mainly focused on documents or images, and cannot be directly applied to graph structured data and graph learning tasks because they cannot model graph structural information, which is the core of node classification.

The challenge of graph learning is that graphs have node content and structure information. Furthermore, existing solutions to node classification task are built on the closed-world assumption, in which the classes appeared in the testing data must have appeared in training. For example, the basic idea of graph convolutional networks (GCNs) is to develop a convolutional layer to exploit the graph structure information and use a classification loss function to guide the classification task. However, they directly use softmax as the final output layer, which does not have the rejection capability to unseen class nodes because the prediction probability of each class is normalized across all training/seen classes. In addition, in representation learning level, most existing graph learning methods employ feature engineering or deep learning to extract feature vectors. However, these models can only generate deterministic mappings to capture latent features of nodes. A major limitation of them is their inability to represent uncertainty caused by incomplete or finite available data.

In this paper, we propose to study open-world learning for graph data. Considering the complicated graph data structure and the node classification task, we summarize the main challenges as follows,

- *Challenge 1*: How to design an end to end framework for open-world graph learning in graphs where unseen class has no labeled samples, and may exist in an arbitrary form different from seen classes. Existing graph neural networks (GNNs) are typical built based on closed-world assumption and cannot detect unseen class.
- *Challenge 2*: How can we model the uncertainty of node

representations and promote robustness in graphs. Many existing GNN-based approaches only generate deterministic mappings to capture latent features of nodes.

To overcome the above challenges, we propose a novel open-world graph learning paradigm (OpenWGL) for node classification tasks. For *Challenge 1*, we employ two loss constraints (a label loss and a class uncertainty loss) to ensure that the node representation learning is sensitive to unseen class and assist in our model to differentiate whether a node belongs to an existing/seen class or an unseen class. We visualize a testing dataset in our experiment in Fig. 2, which can illustrate the effectiveness of our method. In Fig. 2(a), we only use the label loss (the cross-entropy loss), which has a good performance on existing/seen class nodes, but unseen class nodes cannot be differentiated and will be classified to seen classes randomly. In Fig. 2(b), we introduce a class uncertainty loss constraint, which can reduce the probability of unseen class nodes being classified as seen class, and therefore help detect unseen class nodes without impact on the classification of nodes in seen classes. For *Challenge 2*, instead of learning deterministic node feature vector, We utilize a graph variational autoencoder module to learn a latent distribution to represent each node. During the classification phase, a novel sampling process is used to generate multiple versions of feature vectors to test the certainty of a node belonging to seen classes, and automatically determine a threshold to reject nodes not belonging to seen classes as unseen class nodes.

Our contributions can be summarized as follows:

- We formulate a new *open-world learning* problem for graph data, and present a novel deep learning model OpenWGL as a solution.
- We propose an uncertain node representation learning approach, by using label loss and class uncertainty loss to constrain variational graph autoencoder to learn node representation sensitive to unseen class.
- We propose to use sampling process to test the certainty of a node belonging to seen classes, and automatically determine a threshold to reject nodes not belonging to seen classes as unseen class nodes.
- Experiments on benchmark graph datasets demonstrate that our approach outperforms the baseline methods.

II. RELATED WORK

A. Open-World Learning

Open-World Learning aims to recognize the classes the learner has seen/learned before and also detect new class it has never seen before. There are some early explorations of open-world learning. Scholkopf et al. [5] employ the one-class SVM as the classifier, which shows poor performance since no negative data is used. Fei and Liu [9] propose a Center-Based Similarity (CBS) space learning method, which first computes a center for each class and converts each document to a vector of similarities to the center. Fei et al. [3] then extend their work by adding the capability of incrementally or cumulatively learning new classes.

Recently, open-world learning has been studied in Natural Language Processing [1] [2] and computer vision (where it is called open-set recognition) [6] [7] [8]. In NLP, Shu et al. [1] propose the deep learning solution to open-world learning by setting thresholds before the sigmoid function to reject unseen classes. Xu et al. [2] propose a new open-world learning model based on meta-learning, which allows new classes to be added or deleted with no need for model re-training. In computer vision, Scheirer et al. [6] study the problem of recognizing unseen images that are not in the training data by reducing the half-space of a binary SVM classifier with a positive region. In [7] and [8], Scheirer et al. utilize the probability threshold to detect new classes, while their models are weak because of lacking prior knowledge.

B. Emerging Class and Outlier Detection

Our research is also related to emerging/new class detection in supervised learning, such as stream data mining [10], [11] and multi-instance learning [12], and outlier detection [13].

In supervised learning, instances are assumed to belong to at least one of the predefined classes, and a classifier is trained to learn discriminative patterns to separate samples into known classes. When a class is unknown or unavailable at the time of training a classifier, in the test stage, an ideal classifier is expected to be able to detect the emerging/new class [14]. A common solution of detecting new class samples is to use a decision threshold to give a confidence score [15]–[17], including multilayer neural network [18] to increase the threshold, and samples with low confidence below threshold are recognized as the new class. Unfortunately, as we have shown in Fig. 2, simply increasing the threshold will make existing class samples being misclassified.

Outlier detection, on the other hand, aims to detect data instances which abnormally deviate from the underlying data [19]. Some distance-based outlier detection methods such as One-class SVM have been proposed, in which the normal data domain is obtained by finding a hyper-sphere enclosing the normal data samples [5] [20]. A recent method [13] proposes to detect outliers from data stream, but new class detection by outliers is not addressed.

C. Graph Neural Networks

Graph Neural Networks (GNNs), introduced in [21] and [22] as a generalization of recursive neural networks to directly deal with a more general class of graphs, *e.g.* cyclic, directed and undirected graphs, are a powerful tool for machine learning on graphs. GNNs have attracted attention all around the world, which are designed to use deep learning architectures on graph-structured data [23] [24] [25]. Many solutions are proposed to generalize well-established neural network models that work on regular grid structure to deal with graphs with arbitrary structures [26] [27] [28].

Among these methods, the most classic model is GCN, which is a deep convolutional learning paradigm for graph-structured data integrating local node features and graph topology structure in convolutional layers [29]. GraphSAGE

[30] is a variant of GCN which designs different aggregation methods for feature extraction. Although GCNs have shown great performance in graph-structured data for semi-supervised learning tasks such as node classification, the variational graph autoencoder (VGAE) [31] extends it to unsupervised scenarios. Specifically, VGAE integrates GCN into the variational encoder framework [32] by using a graph convolutional encoder and a simple inner product decoder.

To the best of our knowledge, the open-world learning problem has not been previously investigated in graph structure data and graph learning tasks. We are the first to study the *open-world graph learning* and propose an novel uncertain node representation learning approach, based on a variant of GCN (*i.e.*, variational graph autoencoder networks) to differentiate whether a node belongs to an existing (seen) class or an unseen class.

III. PROBLEM DEFINITION AND OVERALL FRAMEWORK

A. Problem Statement

Node Classification on Graphs: In this paper, we focus on node classification on graphs. A graph is represented as $G = (V, E, X, Y)$, where $V = \{v_i\}_{i=1, \dots, N}$ is a vertex set representing nodes in a graph, and $e_{i,j} = (v_i, v_j) \in E$ is an edge indicating the relationship between two nodes. The topological structure of a graph G can be represented by an adjacency matrix A , where $A_{i,j} = 1$ if $(v_i, v_j) \in E$; otherwise $A_{i,j} = 0$. $x_i \in X$ indicates content features associated with each node v_i . $Y \in \mathbb{R}^{N \times C}$ is a label matrix of G , where N is the number of nodes in G and C is the number of node categories (classes) already known/seen. If a node $v_i \in V$ is associated with label l , $Y_{(i)}^l = 1$; otherwise, $Y_{(i)}^l = 0$.

Open-World Graph Learning: Given a graph $G = (V, E, X, Y)$, $X = X_{train} \cup X_{test}$, where X_{train} denotes training data (labeled nodes) and X_{test} denotes testing nodes (unlabeled nodes). Assume $X_{test} = S \cup U$, where S are the set of nodes belonging to seen classes already appeared in X_{train} and U are the set of nodes not belonging to any seen class (*i.e.* unseen class nodes). Open-World Learning on Graphs **aims** to learn a $(C + 1)$ -class classifier model, $f(X_{test}) \mapsto Y$, ($Y \in \{1, \dots, C, rejection\}$) to classify each test node S to one of the training/seen classes in Y and reject U to indicate that it does not belong to any training/seen class (*i.e.*, it belongs to the unseen class).

B. Overall Framework

Our framework for open-world graph learning, as shown in Fig. 3, mainly consists of following two components:

- **Node Uncertainty Representation Learning.** Most GCN models generate deterministic mappings to capture latent features of nodes. A major limitation of these models is their inability to represent uncertainty caused by incomplete or finite available data. In order to learn a better representation of each node, we employ a Variational Graph Autoencoder Network to obtain a latent distribution of each node, which enables to represent uncertainty and promote robustness.

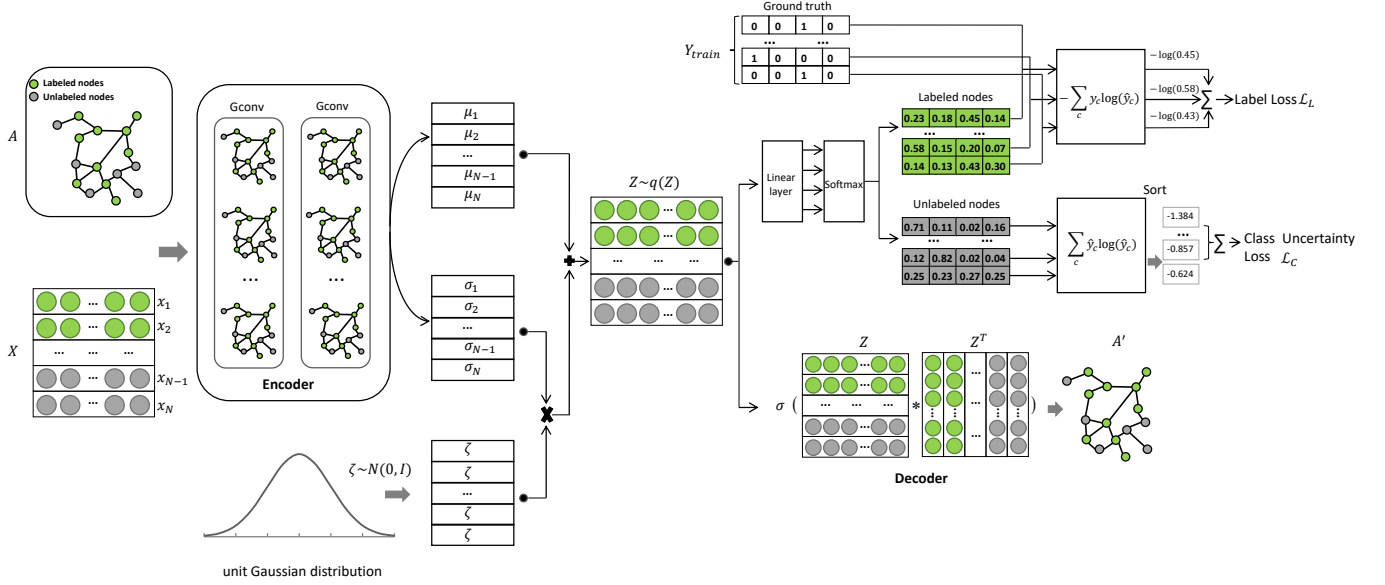


Fig. 3: The overall architecture of the proposed Open-World Graph Learning (OpenWGL) model for node classification task. The input consists of a graph with labeled and unlabeled nodes. The learning objective of OpenWGL, defined in Eq. (8), is constrained by (1) the KL divergence loss and network reconstruction loss (Eq. (7)), and (2) label loss (Eq. (9)), and class uncertainty loss (Eq. (10)). As a result, OpenWGL can learn uncertain node representation sensitive to the class labels and unseen class. More details are given in Section IV.

- **Open-World Classifier Learning.** In order to classify seen class nodes to their own groups and detect unseen class nodes, we introduce two constraints, label loss and class uncertainty loss, to differentiate whether a node belongs to an existing class or an unseen class.

Open-World Classification & Rejection. To perform inference during the testing phase (i.e., perform classification or rejection of an example), we propose a novel sampling process to generate multiple versions of feature vectors to test the certainty of a node belonging to seen classes and automatically determine a threshold to reject nodes not belonging to seen classes as unseen nodes. Our inference framework is given in Fig. 4 with detailed discussion given in Section IV. C.

IV. METHODOLOGY

A. Node Uncertainty Representation Learning

In order to encode latent feature information of each node and obtain an effective representation of uncertainty, we employ Variational Graph Autoencoder network (VGAE) to generate a latent distribution based on extracted node features. This allows our method to leverage uncertainty for robust representation learning.

Graph Encoder Model: Given a graph $G = (X, A)$, in order to represent both node content X and graph structure A in a unified framework, our approach firstly utilizes a two-layer GCN method proposed by [29]. Given the input feature matrix X and adjacency matrix A , the first GCN layer generates a lower-dimensional feature matrix, which is defined as follows:

$$Z^{(1)} = \text{GCN}(X, A) = \text{ReLU} \left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{\frac{1}{2}} X W^{(1)} \right) \quad (1)$$

where $\tilde{A} = A + I_n$ is the adjacency matrix with self-loops ($I_n \in \mathbb{R}^{n \times n}$ is the identity matrix), and $\tilde{D}_{i,i} = \sum_j \tilde{A}_{i,j}$. Accordingly, $\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{\frac{1}{2}}$ is the normalized adjacency matrix. $W^{(1)}$ is the trainable parameters of the network, and $\text{ReLU}(\cdot)$ denotes the activation function.

For the second layer GCN model, instead of generating a deterministic representation, we assume that the output Z is continuous and follows a multivariate Gaussian distribution. Hence, we follow an inference model proposed by [31]:

$$q(Z|X, A) = \prod_{i=1}^N q(z_i|X, A), \quad (2)$$

$$q(z_i|X, A) = \mathcal{N}(z_i | \mu_i, \text{diag}(\sigma_i^2)) \quad (3)$$

Here, $\mu = \text{GCN}_\mu(X, A) = \text{ReLU} \left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{\frac{1}{2}} Z^{(1)} W^{(2)} \right)$ is the matrix of mean vectors μ_i ; σ is the standard variance matrix of the distribution, $\log \sigma = \text{GCN}_\sigma(X, A) = \text{ReLU} \left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{\frac{1}{2}} Z^{(1)} W'^{(2)} \right)$. Then we can calculate Z using a parameterization trick:

$$Z = \mu + \sigma \cdot \zeta, \zeta \sim \mathcal{N}(\mathbf{0}, I) \quad (4)$$

where $\mathbf{0}$ is a vector of zeros and I is the identity matrix. By making use of the latent variable Z , our model is able to capture complex noisy patterns in the data.

Graph Decoder Model: After we get the latent variable Z , we use a decoder model to reconstruct the graph structure A to better learn the relationship between two nodes. Here, the

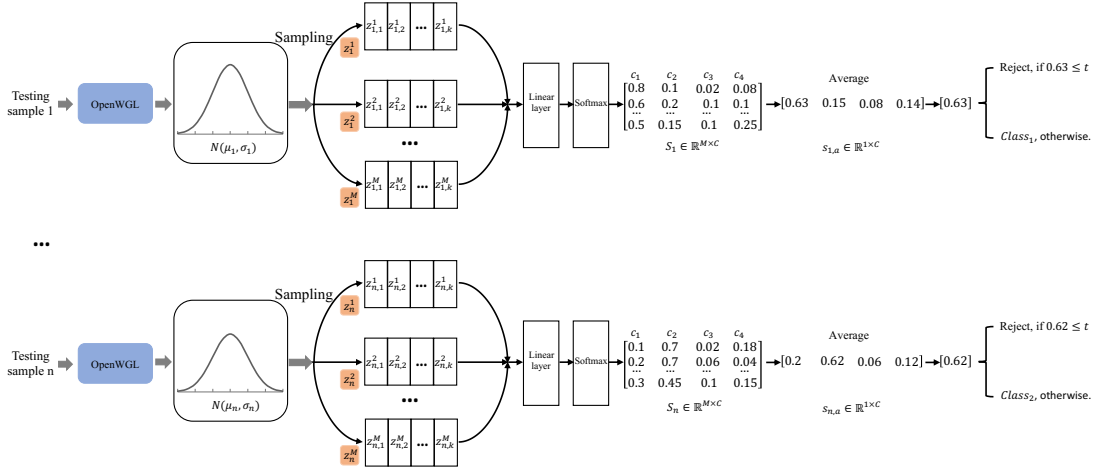


Fig. 4: The classification and rejection process (assuming seen class set has 4 classes). For nodes in the testing set, node uncertainty representation learning generates M different versions of feature vectors for each node by a sampling process. The M different representations are fed into a softmax layer to obtain M probability outputs S_i . The probabilities of each class are averaged to obtain a vector $s_{i,a}$, and the largest average is denoted by $\max(s_{i,a})$. Finally, Eq.(11) is used to decide whether a node belongs to the seen or unseen classes.

graph decoder model is defined by a generative model [31]:

$$p(A|Z) = \prod_{i=1}^N \prod_{j=1}^N p(A_{i,j}|z_i, z_j), \quad (5)$$

$$p(A_{ij} = 1|z_i, z_j) = \sigma(z_i^T z_j), \quad (6)$$

where A_{ij} are the elements of A and $\sigma(\cdot)$ denotes the logistic sigmoid function.

Optimization: To better learn class discriminative node representations, we optimize the variational graph autoencoder module via two losses as follows:

$$\mathcal{L}_{VGAE} = \mathbb{E}_{q(Z|X,A)}[\log p(A|Z)] - KL[q(Z|X,A)||p(Z)] \quad (7)$$

where the first term is the reconstruction loss between the input adjacent matrix and the reconstructed adjacent matrix. The second term $KL[q(Z|X,A)||p(Z)]$ is the Kullback-Leibler divergence between $q(Z|X,A)$ and $p(Z)$, here $p(Z) = \mathcal{N}(\mathbf{0}, \mathbf{I})$.

B. Open-World Classifier Learning

After the variational graph autoencoder network, we obtain the uncertainty embeddings for each node through Eq. (4), which consists of two parts: uncertainty embeddings for labeled/training nodes $Z_{labeled}$ and uncertainty embeddings for unlabeled/test nodes $Z_{unlabeled}$. To better learn an accurate classifier for classifying both seen and unseen nodes in testing data, our proposed model consists of a cooperative module, a label loss as well as a class uncertainty loss working together to differentiate whether a node belongs to an existing class or an unseen class. The overall objective function is as follows:

$$\mathcal{L}_{OpenWGL} = \gamma_1 \mathcal{L}_L + \gamma_2 \mathcal{L}_C + \mathcal{L}_{VGAE} \quad (8)$$

The γ_1, γ_2 are the balance parameters. The \mathcal{L}_{VGAE} is the loss function of the variational graph autoencoder network

mentioned above. The \mathcal{L}_L and \mathcal{L}_C represent the label loss and the class uncertainty loss, respectively. The details are introduced as follows. **Label Loss:** The label loss \mathcal{L}_L is to minimize the cross-entropy loss for the labeled data:

$$\mathcal{L}_L(f_s(Z_{labeled}), Y) = -\frac{1}{N_l} \sum_{i=1}^{N_l} \sum_{c=1}^C y_{i,c} \log(\hat{y}_{i,c}) \quad (9)$$

where $f_s(\cdot)$ is a softmax layer consisting of a full-connected layer with corresponding activation function which can transform $Z_{unlabeled}$ into probabilities that sum to one. N_l is the number of labeled nodes. C denotes the number of seen classes, and $y_{i,c}$ denotes the groundtruth of the i -th node in the labeled data, $\hat{y}_{i,c}$ is the classification prediction score for the i -th labeled node v_i in the c class, respectively.

Class Uncertainty Loss: Since we do not have the class information in the test data and there exists a considerable number of unseen nodes, we need to find a way to differentiate the seen class and unseen class. Unlike the label loss \mathcal{L}_L , which can utilize the abundant training data and have a good performance on the seen class by the cross-entropy loss, the class uncertainty loss is proposed to balance the classification output for each node and have superior effects on the unseen nodes. In our paper, an entropy loss is placed as the class uncertainty loss and our goal is to **maximize** this entropy loss to make the normalized output of each node balanced. The formula is as follows:

$$\mathcal{L}_C(f_s(Z_{unlabeled})) = \frac{1}{N_u} \sum_{i=1}^{N_u} \sum_{c=1}^C \hat{y}_{i,c} \log(\hat{y}_{i,c}) \quad (10)$$

where N_u is the number of unlabeled nodes. $\hat{y}_{i,c}$ is the classification prediction score for the i -th unlabeled node v_i in the c class. Note that we do not put a negative sign in front of the formula as usual because we need to maximize the entropy

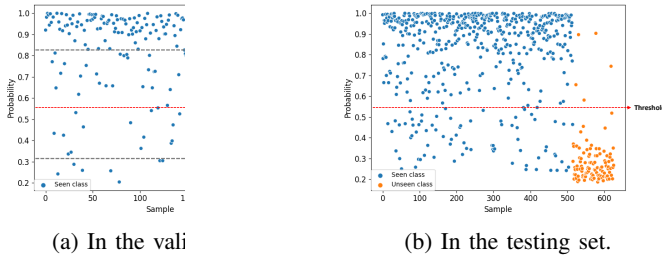


Fig. 5: A visualization of determining the threshold using a validation set (only contains seen class instances). (a) determining the threshold using validation set. (b) applying determined threshold to the test set (contain both seen class and unseen class instances).

loss. In addition, we will not use all the unlabeled data to max the entropy loss. We first sort all the unlabeled data output probability values (choosing the maximum probability for each node) after the softmax layer, and then discard the largest 10% (nodes with large probability values are easily classified into seen classes since their output is discriminative) and the smallest 10% nodes (nodes with small probability means that the node’s output is balanced over each seen class which can be easily detected as the unseen class). Finally the remaining nodes are utilized to maximize their entropy.

The training for label loss and class uncertainty loss is like a adversarial process. On the one hand, we want the label loss to influence the classifier to make the output of each node more discriminative and classify each seen node into the correct class via minimizing Eq. (9). On the other hand, we would like that the class uncertainty loss can make the output of each node more balanced to assist in detecting the unseen class through maximizing the entropy loss.

\mathcal{L}_L , \mathcal{L}_C and \mathcal{L}_{VGAE} are jointly optimized via our objective function in Eq. (8), and all parameters are optimized using the standard backpropagation algorithms.

C. Open-World Classification & Rejection

After performing the node uncertainty representation learning, we obtain a distribution (i.e. the Gaussian distribution) of the node embeddings. Thus we generate M different versions of feature vectors (z_i^1, \dots, z_i^M) for each node v_i from this distribution via Eq. (4) called a reparametrization trick. Then we feed these M different representations into the softmax layer to turn them into probabilities over C classes respectively (each z_i^m can obtain an output vector $s_i^m \in \mathbb{R}^{1 \times C}$).

After this process, for each node we concatenate these M outputs and obtain a sampling matrix $S_i \in \mathbb{R}^{M \times C}$. In S_i , each column denotes M different probabilities of a specific class and we average these probabilities for each class to obtain a vector $s_{i,a} \in \mathbb{R}^{1 \times C}$. For the vector $s_{i,a}$ with C different probabilities, we choose the largest one $\max(s_{i,a})$. To recognize whether each node v_i is the seen or unseen classes for testing data, we have:

$$\hat{y} = \begin{cases} \text{Rejection,} & \text{if } \max_{c \in C} p(c|x_i) \leq t \\ \arg \max_{c \in C} p(c|x_i), & \text{otherwise.} \end{cases} \quad (11)$$

where $p(c|x_i)$ is obtained from the softmax layer output of $f_s(\cdot)$. If none of existing seen classes probability $p(c|x_i)$ value is above the threshold t , we reject x_i as a sample from the unseen class; otherwise, its predicted class is the one with the highest probability. The prediction process of each testing sample is illustrated in Fig. 4.

In open-world graph learning, a key problem is the determining of the threshold t . In our paper, we propose a selection approach to automatically determine a threshold to reject nodes not belonging to seen classes. Specifically, we use the validation set to do the threshold selection. Similarly, for nodes in the validation set, we perform the node uncertainty representation learning, and conduct the same sampling process and choose the largest probability. Then we average these chosen largest probabilities of all the nodes and obtain avg_seen . Because unseen class instances are assumed not appearing in the training set (including the validation set), we choose 10% nodes with the largest entropy as the “expected unseen class nodes”, and their average probability is denoted by avg_E_unseen . The final threshold is calculated by averaging the probabilities: $t = (avg_seen + avg_E_unseen)/2$. Fig. 5(a) shows an example of the determining process in the validation set. We use this determined threshold to classify seen and unseen nodes in the test set in Fig. 5(b).

Therefore, by using a sampling process to generate multiple versions of feature vectors, we are able to test the confidence of a node belonging to seen classes, and automatically determine a threshold to reject nodes not belonging to seen classes as unseen class nodes.

D. Algorithm Description

Our algorithm is illustrated in Algorithm 1. Given a graph $G = (V, E, X, Y)$, our goal is to obtain the node representations and classify the seen nodes and detect the unseen nodes, respectively. Firstly, we employ a variational graph autoencoder network to model the uncertainty of each node (Step 2-10). Here, the output Z is a distribution and we optimize the network through the KL loss and the reconstruction loss (Step 12). Then we propose two loss constraints \mathcal{L}_L and \mathcal{L}_C to make our model capable of classifying seen and unseen classes (Step 13-14). Finally, by jointly considering the label loss, class uncertainty loss and the VGAE loss, our model can better differentiate whether a node belongs to a seen class or an unseen class and capture the uncertainty representations for open-world graph learning.

V. EXPERIMENTS

A. Experimental Setup

Benchmark Datasets We employ three widely used citation network datasets (Cora, Citeseer, DBLP) for node classification [33] [34]. The details of the experimental datasets are reported in Table I.

Test Settings and Evaluation Metrics For each dataset, we hold out some classes as the unseen class for testing and the remaining classes as the seen classes. we randomly sample

Algorithm 1: Open-World Graph Learning

Data: $G = (V, E, X, Y)$: a Graph with links and features;
 $X = X_{train} \cup X_{test}$, $X_{test} = S \cup U$: S are the seen classes appeared in X_{train} and U are the unseen classes; C : the number of seen classes.
Result: $f(X_{test}) \mapsto Y$, $Y \in \{1, \dots, C, rejection\}$.
1: **while** not convergence **do**
2: // Graph Encoder Model
3: For the first layer:
4: $Z^{(1)} \leftarrow ReLU\left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{\frac{1}{2}} X W^{(1)}\right)$
5: For the second layer:
6: $\mu \leftarrow ReLU\left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{\frac{1}{2}} Z^{(1)} W^{(2)}\right)$
7: $\log \sigma \leftarrow ReLU\left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{\frac{1}{2}} Z^{(1)} W'^{(2)}\right)$
8: $Z \leftarrow \mu + \sigma \cdot \zeta, \zeta \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
9: // Graph Decoder Model
10: $p(A_{ij} = 1 | z_i, z_j) \leftarrow \sigma(z_i^T z_j)$
11: // Compute Loss
12: $\mathcal{L}_{VGAE} \leftarrow$ Obtain the variational graph autoencoder loss using Eq. (7)
13: $\mathcal{L}_L \leftarrow$ Obtain the label loss using Eq.(9)
14: $\mathcal{L}_C \leftarrow$ Obtain the class uncertainty loss using Eq.(10)
15: Back-propagate loss gradient using Eq.(8)
16: $[W^{(1)}, W^{(2)}, W'^{(2)}, f_s(\cdot)] \leftarrow$ Update weights
17: **if** early stopping condition satisfied **then**
18: Terminate

TABLE I: Statistics of the experimental datasets.

Dataset	# of Nodes	# of Edges	# of Features	# of Labels
Cora	2,708	5,429	1,433	7
Citeseer	3,312	4,732	3,703	6
DBLP	60,744	52,890	1,587	4

70% of nodes for training, 10% for validation and 20% for testing. Note that, the nodes of unseen class only appear in the testing set. We use the validation set to determine the threshold for rejecting the unseen class. Like the traditional semi-supervised node classification, for each dataset, we feed the whole graph into our model. We vary the number of unseen classes to verify the performance of our model at different unseen class proportion. We use the Macro F1 score and Accuracy for evaluation [1].

Baselines We employ following methods as baselines.

- GCN [29]: GCN is a deep convolutional network for graph-structured data, which directly uses softmax as the final output layer. GCN does not have the rejection capability to the unseen class.
- GCN_Sigmoid: In GCN_Sigmoid, we use multiple 1-vs-rest of sigmoids rather than softmax as the final output layer of the GCN model, which also does not have the rejection capability to the unseen class.
- GCN_Sigmoid_Thre: Based on GCN_Sigmoid, we use the default probability threshold of $t_i = 0.5$ for classification of each class i , which means if all predicted probabilities are less than the threshold 0.5, we will reject it as the unseen class. Otherwise, its predicted class is the one

TABLE II: Experimental results on Cora with different number of unseen classes $|U|$.

Methods	$ U = 1$		$ U = 3$	
	Accuracy	Macro F1	Accuracy	Macro F1
GCN	0.726	0.683	0.345	0.463
GCN_Sigmoid	0.728	0.681	0.338	0.463
GCN_Sigmoid_Thre	0.782	0.786	0.593	0.664
MLP_DOC	0.455	0.452	0.670	0.493
GCN_DOC	0.753	0.769	0.729	0.735
OpenWGL	0.833	0.835	0.775	0.752

TABLE III: Experimental results on Citeseer with different number of unseen classes $|U|$.

Methods	$ U = 1$		$ U = 3$	
	Accuracy	Macro F1	Accuracy	Macro F1
GCN	0.445	0.477	0.263	0.320
GCN_Sigmoid	0.443	0.472	0.258	0.318
GCN_Sigmoid_Thre	0.670	0.609	0.683	0.621
MLP_DOC	0.455	0.433	0.745	0.564
GCN_DOC	0.687	0.613	0.758	0.679
OpenWGL	0.700	0.654	0.766	0.698

with the highest probability.

- MLP_DOC: DOC [1] is the state-of-the-art open-world classification method for text classification. We use a two-layer perceptron to obtain the node representation.
- GCN_DOC: We utilize the rich node relationships and combine the GCN with DOC to compare with our model. In DOC, it uses multiple 1-vs-rest of sigmoids rather than softmax as the final output layer and defines a automatic threshold setting mechanism.

All deep learning algorithms are implemented using Tensorflow and are trained with Adam optimizer. We follow the evaluation protocol in open-world learning [1] [2] and evaluate all approaches through grid search on the hyperparameter space and report the best results of each approach. We feed the whole graph into our model when training. For all baseline methods, we use the same set of parameter configurations unless otherwise specified. For each deep approach, we use a fixed learning rate $1e^{-3}$. For each method, the GCNs contain two hidden layers ($L = 2$) with structure as $32 - 16$. The balance parameters γ_1, γ_2 are set to 1, 0.8, respectively. The dropout rate for each GCN layer is set to 0.3. M is set to 100.

TABLE IV: Experimental results on DBLP with different number of unseen classes $|U|$.

Methods	$ U = 1$		$ U = 2$	
	Accuracy	Macro F1	Accuracy	Macro F1
GCN	0.662	0.562	0.285	0.323
GCN_Sigmoid	0.662	0.562	0.290	0.323
GCN_Sigmoid_Thre	0.657	0.650	0.282	0.326
MLP_DOC	0.643	0.630	0.480	0.477
GCN_DOC	0.657	0.658	0.503	0.506
OpenWGL	0.688	0.689	0.653	0.642

TABLE V: The Macro F1 score and Accuracy on three datasets for closed-world settings (without unseen classes).

Dataset ($ U = 0$)		GCN	OpenWGL
Cora	Accuracy	0.863	0.854
	Macro F1	0.848	0.829
Citeseer	Accuracy	0.774	0.779
	Macro F1	0.752	0.745
DBLP	Accuracy	0.806	0.809
	Macro F1	0.754	0.751

B. Open-world Graph Learning Classification Results

Table II, Table III and Table IV list the Macro F1 score and Accuracy of different methods on open-world node classification task. From the results, we have following observations:

- (1) The GCN and GCN_Sigmoid obtain the worst performance among these baselines in all datasets since they do not have the rejection capability to the unseen class. Therefore, all the unseen nodes will be misclassified and their performance become worse with the number of unseen nodes increases.
- (2) GCN_Sigmoid_Thre and GCN_DOC have better performances than GCN and GCN_Sigmoid, which shows that the threshold can improve the performance of detecting the unseen nodes. In addition, with the number of unseen nodes increases, GCN_Sigmoid_Thre and GCN_DOC become more competitive.
- (3) GCN_DOC has better performance than GCN_Sigmoid_Thre in most cases, confirming that the threshold is not a fixed value and it varies with different datasets and the ratio of unseen class. DOC's automatic threshold setting mechanism can effectively improve the classification results of unseen class.
- (4) The proposed Open-World Graph Learning model consistently outperforms all baselines on three datasets with different numbers of unseen classes. It demonstrates that the proposed constrained graph variational encoder network can better differentiate whether a node belongs to a seen class or an unseen class and capture the uncertainty representation of each node by jointly considering the label loss, class uncertainty loss and the node uncertainty representation learning as a unified learning framework.
- (5) We also report closed-world learning setting results (without unseen class) in Table V. The results show that when networks do not have unseen class, OpenWGL has comparable performance as GCN, confirming its effectiveness and generalization for node classification.

C. Ablation Analysis of OpenWGL Components

Because OpenWGL contains two key constraints, in this subsection, we compare variants of OpenWGL with respect to the following aspects to demonstrate: (1) the effect of the class uncertainty loss, and (2) the impact of the VGAE module (KL loss and reconstruction loss).

The following OpenWGL variants are designed for comparison.

TABLE VI: The Macro F1 score and Accuracy between OpenWGL variants on Cora.

Methods	$ U = 1$		$ U = 3$	
	Accuracy	Macro F1	Accuracy	Macro F1
OpenWGL \neg_C	0.782	0.787	0.700	0.665
OpenWGL \neg_V	0.824	0.829	0.785	0.705
OpenWGL	0.833	0.835	0.775	0.752

TABLE VII: The Macro F1 score and Accuracy between OpenWGL variants on Citeseer.

Methods	$ U = 1$		$ U = 3$	
	Accuracy	Macro F1	Accuracy	Macro F1
OpenWGL \neg_C	0.676	0.645	0.759	0.692
OpenWGL \neg_V	0.691	0.648	0.760	0.683
OpenWGL	0.700	0.654	0.766	0.698

TABLE VIII: The Macro F1 score and Accuracy between OpenWGL variants on DBLP.

Methods	$ U = 1$		$ U = 2$	
	Accuracy	Macro F1	Accuracy	Macro F1
OpenWGL \neg_C	0.675	0.672	0.650	0.631
OpenWGL \neg_V	0.687	0.671	0.651	0.635
OpenWGL	0.688	0.689	0.653	0.642

- OpenWGL \neg_C : A variant of OpenWGL with only the class uncertainty loss being removed.
- OpenWGL \neg_V : A variant of OpenWGL with the KL loss and reconstruction loss being removed.

Tables VI, VII & VIII report the ablation study results.

1) *the effect of the class uncertainty loss*: In order to show the superiority of the class uncertainty loss, we design a variant model OpenWGL \neg_C . As mentioned before, the class uncertainty loss is a constraint on the unlabeled nodes. The ablation study results show the performances of the node classification task on both datasets are improved when the class uncertainty loss is used, indicating its effectiveness of detecting unseen nodes.

2) *the impact of the VGAE module (KL loss and reconstruction loss)*: In order to verify the impact of the VGAE module which can model the uncertainty node representations, we compare OpenWGL model and OpenWGL \neg_V . From the results, we can easily observe the OpenWGL model performs significantly better than OpenWGL \neg_V . This confirms that the usage of KL loss can model the uncertainty to better capture the latent representation of each node, and reconstruction loss can preserve node relationships which will assist in the node representation.

D. Parameter Analysis

Impact of the feature dimensions of node output embeddings Z : As mentioned in the method section, the output of node embeddings is represented as Z . OpenWGL uses 2-layer GCNs with structure as 32 – 16, and feature dimensions d of node output embeddings is 16. We vary d from 4 to 64

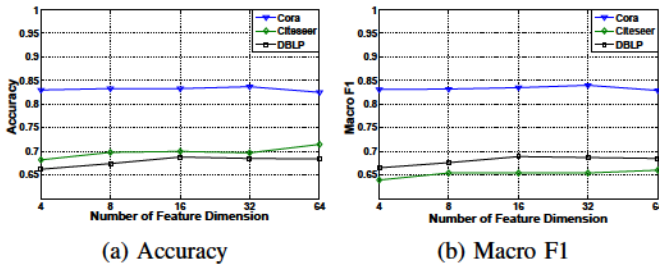


Fig. 6: Impact of feature dimensions of node output embeddings for the accuracy and Macro F1 score on three datasets.

and report the results on three datasets, respectively in Fig. 6. On Citeseer and DBLP datasets, as d increases from 4 to 64, the performance grows gradually to reach a plateau. The performance of Cora dataset is stable with d increasing from 4 to 32 and has a slight decrease at 64. Therefore, only slight differences can be observed with different d values. The increase of d , from 4 to 64, does not necessarily result in performance improvements. The results show that with sufficient feature dimensions ($d \geq 16$), OpenWGL is stable with the increasing number of feature dimensions.

E. Case Study

1) *Visualization of the OpenWGL sampling results:* In order to verify the effectiveness of the sampling process of our model, we randomly choose two testing nodes from Cora dataset for seen and unseen classes (we choose one class as unseen, i.e. $|U| = 1$), respectively. After performing the node uncertainty representation learning, we obtain a distribution of the node embeddings. Then we generate 100 different versions of feature vectors for each node from this distribution and feed them into the softmax layer to turn them into probabilities over 6 classes, respectively. Therefore, after this process, for each node we obtain a 6×100 sampling matrix. In the sampling matrix, each column denotes 100 different probabilities of a specific class. We visualize the sampling matrices of these four nodes through histogram charts with seen and unseen classes in Figs. 7(a) and (b). In Fig. 7, each row represents one node and in each row, there are six subfigures indicating the 100 different probabilities of each class, respectively. From Fig. 7, we can observe that the sampling process have superior performance in differentiating the seen classes and the unseen class, and it is very helpful for determining the threshold. For example, as shown in the first row in Fig. 7(a), only in class 2, most of the 100 different probabilities are distributed on far right side of the histogram (i.e., large probability), while all the other classes (0,1,3,4,5) are distributed on the far left side (i.e., small probability). Thus, through the softmax layer, we can classify this node to class 2 and the ground truth is also class 2. However, if we just use a deterministic feature vector instead of this sampling method, this node may not be classified to the class 2, since class 2 also has cases with small probability values. Similarly, for the unseen nodes as shown in Fig. 7(b), in each seen class, most of the probability values are concentrated on the left side of the histogram (i.e., small probability), so we can easily detect them and classify them

into unseen class. However, If we only obtain one probability output and do not have the sampling process, the unseen node might be misclassified randomly.

2) *The Confusion Matrix:* In order to verify the effectiveness of OpenWGL in differentiating seen class nodes vs. unseen class nodes, Fig. 8 reports the confusion matrix of OpenWGL on Cora network, where “-1” denotes unseen class. The results show that OpenWGL correctly identifies 87% of unseen class nodes and also remains a high accuracy in classifying seen class nodes.

VI. CONCLUSIONS

In this paper, we studied a new open-world graph learning problem. We argued that traditional graph learning tasks are based on the closed-world setting (i.e., classifying nodes into classes already known, so they cannot correctly classify new/unseen class nodes). In the paper, we advocated an open-world graph learning paradigm which not only classifies nodes belonging to seen classes into correct groups, but also classifies nodes not belonging to existing classes to an unseen class. To achieve the goal, we proposed a open-world graph learning (OpenWGL) framework with two major components: (1) node uncertainty representation learning, and (2) open-world classifier learning. The former learns a distribution for each node embedding via a graph variational autoencoder to capture the uncertainty, and the latter minimizes the label loss and class uncertainty loss simultaneously to distinguish seen and unseen class nodes, using automatically determined threshold. Experiments showed that when unseen class presents in test data, OpenWGL significantly outperforms baseline in classifying both seen and unseen class nodes. When networks do not have unseen class nodes (only contain nodes from seen classes), OpenWGL has a comparable performance as the baseline.

ACKNOWLEDGMENT

This research is supported by the U.S. National Science Foundation (NSF) through Grant Nos. IIS-1763452, CNS-1828181, and IIS-2027339.

REFERENCES

- [1] L. Shu, H. Xu, and B. Liu, “Doc: Deep open classification of text documents,” *arXiv preprint arXiv:1709.08716*, 2017.
- [2] H. Xu, B. Liu, L. Shu, and P. Yu, “Open-world learning and application to product classification,” in *Proc. of WWW Conf.*, 2019, pp. 3413–3419.
- [3] G. Fei, S. Wang, and B. Liu, “Learning cumulatively to become more knowledgeable,” in *Proc. of KDD*, 2016, pp. 1565–1574.
- [4] Z. Chen and B. Liu, “Lifelong machine learning,” *Synthesis Lectures on Artificial Intel. and Machine Learning*, vol. 12, no. 3, pp. 1–207, 2018.
- [5] B. Schölkopf, J. C. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson, “Estimating the support of a high-dimensional distribution,” *Neural computation*, vol. 13, no. 7, pp. 1443–1471, 2001.
- [6] W. J. Scheirer, A. de Rezende Rocha, A. Sapkota, and T. E. Boult, “Toward open set recognition,” *IEEE Trans. on pattern analysis and machine intelligence*, vol. 35, no. 7, pp. 1757–1772, 2012.
- [7] W. J. Scheirer, L. P. Jain, and T. E. Boult, “Probability models for open set recognition,” *IEEE Trans. on pattern analysis and machine intelligence*, vol. 36, no. 11, pp. 2317–2324, 2014.
- [8] L. P. Jain, W. J. Scheirer, and T. E. Boult, “Multi-class open set recognition using probability of inclusion,” in *ECCV*. Springer, 2014, pp. 393–409.
- [9] G. Fei and B. Liu, “Social media text classification under negative covariate shift,” in *Proc. of EMNLP*, 2015, pp. 2347–2356.

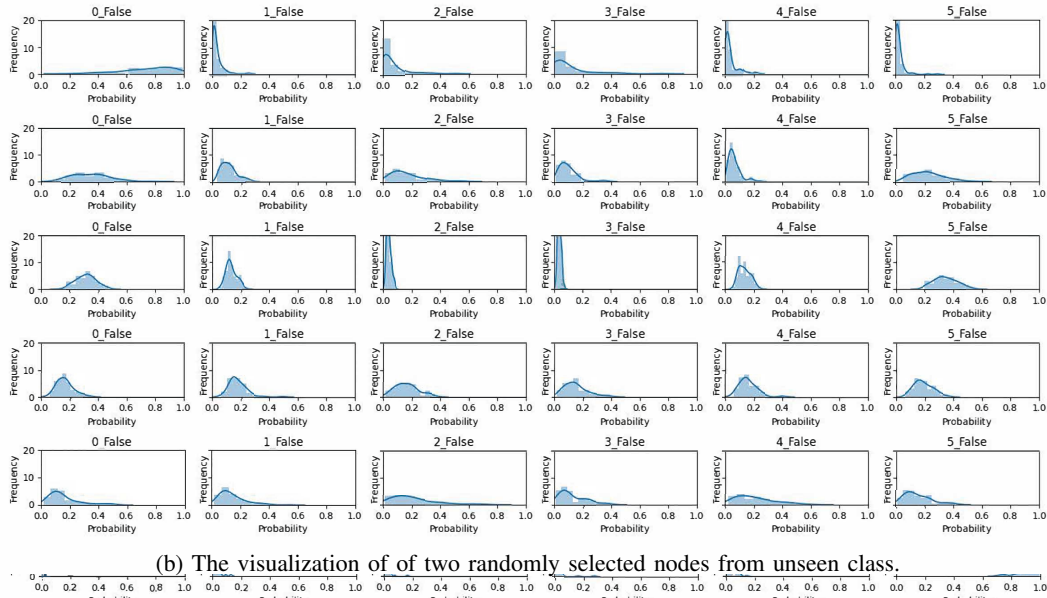


Fig. 7: A case study of the OpenWGL sampling results with two randomly selected nodes from seen classes and unseen class on Cora, respectively. Each row denotes one node, “True” denotes that the node belongs to this class, and “False” means that the node does not belong to this class. (a) two nodes randomly selected from seen classes, and (b) two nodes randomly selected from unseen class. The x -axis denotes the probability output of each node through the softmax classifier, and the y -axis denotes the frequency appearing in each class.

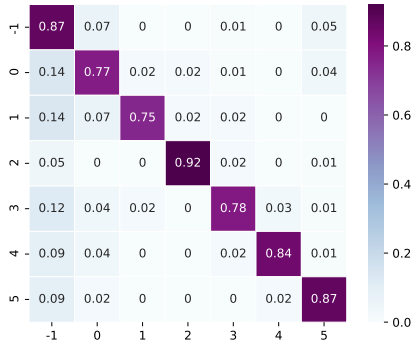


Fig. 8: The confusion matrix of OpenWGL on Cora. “-1” denotes the unseen class and “0,1,2,3,4,5” are seen classes. The (i, j) value of the matrix shows that the percentage value of the i -th class is classified to the j -th category.

- [10] Y. Gao, S. Chandra, Y. Li, L. Kan, and B. Thuraisingham, “Saccos: A semi-supervised framework for emerging class detection and concept drift adaption over data streams,” *IEEE Trans. Knowl. & Data Eng.*, 2020.
- [11] X.-Q. Cai, P. Zhao, K.-M. Ting, X. Mu, and Y. Jiang, “Nearest neighbor ensembles: An effective method for difficult problems in streaming classification with emerging new classes,” in *ICDM*, 2019.
- [12] X.-S. Wei, H.-J. Y. Ye, X. Wu, J. Wu, C. Shen, and Z.-H. Zhou, “Multiple instance learning with emerging novel class,” *IEEE transactions knowledge and data engineering*, 2019.
- [13] G. Na, D. K. Kim, and H. Yu, “Dilof: Effective and memory efficient local outlier detection in data streams,” in *Proc. of KDD*, 2019.
- [14] C. H. Park and H. Shim, “On detecting an emerging class,” in *IEEE Intl. Conf. on Granular Computing (GRC 2007)*. IEEE, 2007, pp. 265–265.
- [15] B. Zadrozny and C. Elkan, “Transforming classifier scores into accurate multiclass probability estimates,” in *Proc. of SIGKDD*, 2002.
- [16] M. Li and I. K. Sethi, “Confidence-based classifier design,” *Pattern Recognition*, vol. 39, no. 7, pp. 1230–1240, 2006.
- [17] K. Proedrou, I. Nouruddinov, V. Vovk, and A. Gammerman, “Transductive confidence machines for pattern recognition,” in *European Conference on Machine Learning*. Springer, 2002, pp. 381–390.
- [18] W. Soares-Filho, J. Seixas, and L. P. Caloba, “Enlarging neural class

detection capacity in passive sonar systems,” in *2002 IEEE International Symposium on Circuits and Systems. Proceedings (Cat. No. 02CH37353)*, vol. 3. IEEE, 2002, pp. III–III.

- [19] E. M. Knorr and R. T. Ng, “Finding intensional knowledge of distance-based outliers,” in *Vldb*, vol. 99, 1999, pp. 211–222.
- [20] E. J. Spinosa and A. Carvalho, “Support vector machines for novel class detection in bioinformatics,” *Genet Mol Res*, vol. 4:3, pp. 608–15, 2005.
- [21] M. Gori, G. Monfardini, and F. Scarselli, “A new model for learning in graph domains,” in *IJCNN*, vol. 2. IEEE, 2005, pp. 729–734.
- [22] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, “The graph neural network model,” *IEEE Trans. on Neural Networks*, vol. 20, no. 1, pp. 61–80, 2008.
- [23] M. Wu, S. Pan, C. Zhou, X. Chang, and X. Zhu, “Unsupervised domain adaptive graph convolutional networks,” in *WWW ’20: The Web Conference, April 20-24, 2020*, 2020, pp. 1457–1467.
- [24] M. Wu, S. Pan, X. Zhu, C. Zhou, and L. Pan, “Domain-adversarial graph neural networks for text classification,” in *IEEE International Conference on Data Mining, ICDM*, 2019, pp. 648–657.
- [25] S. Zhu, L. Zhou, S. Pan, C. Zhou, G. Yan, and B. Wang, “GSSNN: Graph smoothing splines neural networks,” in *AAAI*, 2020, pp. 7007–7014.
- [26] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and P. S. Yu, “A comprehensive survey on graph neural networks,” *TNNLS*, 2020.
- [27] S. Pan, R. Hu, S.-f. Fung, G. Long, J. Jiang, and C. Zhang, “Learning graph embedding with adversarial training methods,” *IEEE Transactions on Cybernetics*, 2019.
- [28] M. Wu, S. Pan, L. Du, I. W. Tsang, X. Zhu, and B. Du, “Long-short distance aggregation networks for positive unlabeled graph learning,” in *Proceedings of CIKM*, 2019, pp. 2157–2160.
- [29] T. N. Kipf and M. Welling, “Semi-supervised classification with graph convolutional networks,” *arXiv preprint arXiv:1609.02907*, 2016.
- [30] W. Hamilton, Z. Ying, and J. Leskovec, “Inductive representation learning on large graphs,” in *Advances in Neural Information Processing Systems*, 2017, pp. 1024–1034.
- [31] T. N. Kipf and M. Welling, “Variational graph auto-encoders,” *arXiv preprint arXiv:1611.07308*, 2016.
- [32] D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” *arXiv preprint arXiv:1312.6114*, 2013.
- [33] C. Yang, Z. Liu, D. Zhao, M. Sun, and E. Y. Chang, “Network representation learning with rich text information,” in *Proc. of IJCAI*, 2015, pp. 2111–2117.
- [34] S. Pan, J. Wu, X. Zhu, C. Zhang, and Y. Wang, “Tri-party deep network representation,” in *Proc. of IJCAI*, 2016, pp. 1895–1901.