

Accurately Modeling Biased Random Walks on Weighted Graphs Using *Node2vec+*

Renming Liu,¹ Matthew Hirn,^{1,2,3,*} Arjun Krishnan^{1,4,*}

¹Department of Computational Mathematics, Science & Engineering, ²Department of Mathematics, ³Center for Quantum Computing, Science & Engineering, ⁴Department of Biochemistry and Molecular Biology, Michigan State University, East Lansing, MI 48824, USA

liurenmi@msu.edu, mhirn@msu.edu, arjun@msu.edu (*Corresponding authors)

Abstract

Node embedding is a powerful approach for representing the structural role of each node in a graph. *Node2vec* is a widely used method for node embedding that works by exploring the local neighborhoods via biased random walks on the graph. However, *node2vec* does not consider edge weights when computing walk biases. This intrinsic limitation prevents *node2vec* from leveraging all the information in weighted graphs and, in turn, limits its application to many real-world networks that are weighted and dense. Here, we naturally extend *node2vec* to *node2vec+* in a way that accounts for edge weights when calculating walk biases, but which reduces to *node2vec* in the cases of unweighted graphs or unbiased walks. We empirically show that *node2vec+* is more robust to additive noise than *node2vec* in weighted graphs using two synthetic datasets. We also demonstrate that *node2vec+* significantly outperforms *node2vec* on a commonly benchmarked multi-label dataset (Wikipedia). Furthermore, we test *node2vec+* against GCN and GraphSAGE using various challenging gene classification tasks on two protein-protein interaction networks. Despite some clear advantages of GCN and GraphSAGE, they show comparable performance with *node2vec+*. Finally, *node2vec+* can be used as a general approach for generating biased random walks, benefiting all existing methods built on top of *node2vec*. *Node2vec+* is implemented as part of PecanPy, which is available at <https://github.com/krishnanlab/PecanPy>.

Introduction

Graphs and networks naturally appear in many real-world datasets including social networks and biological networks. The graph structure provides insightful information about the role of each node in the graph, such as protein function in a protein-protein interaction network (Liu et al. 2020). However, it is typically impractical to access all neighborhood information of the graph due to the massive scale of many real-world networks with millions of nodes (Hu et al. 2021). Node embeddings are one type of solution that aim to resolve this issue by representing all the nodes in a graph in a low dimensional vector space, where some mutual relations between nodes are preserved (Cui et al. 2018; Hamilton, Ying, and Leskovec 2017).

Node2vec is a second-order random walk based embedding method (Grover and Leskovec 2016). It is widely used for unsupervised node embedding for various tasks, particularly in computational biology (Nelson et al. 2019), such as gene function prediction (Liu et al. 2020; Ata et al. 2018) and essential protein prediction (Wang et al. 2021a; Zeng et al. 2021), due to its superior performance than other matrix factorization and neural network based methods (Yue et al. 2019). Some recent works built on top of *node2vec* aim to adapt *node2vec* to more specific types of networks (Wang et al. 2021b; Valentini et al. 2021), generalize *node2vec* to higher dimension (Hacker 2021), augment *node2vec* with additional downstream processing (Chattopadhyay and Ganguly 2020; Hu et al. 2020), or even study *node2vec* theoretically (Grohe 2020; Davison and Austern 2021; Qiu et al. 2018). Nevertheless, none of these follow-up works account for the fact that *node2vec* is less effective for weighted graphs, where the edge weights reflect the (potentially noisy) similarities between pairs of nodes. This failing is due to the inability of *node2vec* to differentiate between small and large edges connecting the previous vertex with a potential next vertex in the random walk, which subsequently causes less accurate modeling of the intended walk bias. In some extreme cases where the weighted graphs are fully connected, such as integrative biological networks (Greene et al. 2015), *node2vec* completely loses its ability to devise a biased random walk.

Meanwhile, another line of recent works on graph neural networks (GNNs) have shown remarkable performance in prediction tasks that involve graph structure, including node classification (Bronstein et al. 2021; Zhou et al. 2021; Zhang et al. 2021; Wu et al. 2021; Xia et al. 2021). Although GNNs and embedding methods like *node2vec* are related in that they both aim at projecting nodes in the graph to a feature space, two main differences set them apart. First, GNNs typically require labeled data while embedding methods do not. This label dependency makes the embeddings generated by a GNN tied to the quality of the labels, which in some cases, like in biological networks, are noisy and scarce. Moreover, in a biological network, the node labels, such as gene functions and disease associations, are typically skewed, with many negatives but few positives (Ata et al. 2018). This data imbalance issue poses another challenge to properly train a GNN. Second, GNNs require node features as input to

train, which are not always available. In the absence of given node features, one needs to generate them and often GNN algorithms in this case use trivial node features such as the constant feature or node degree. These two differences give node embedding methods a unique place in node classification, apart from the GNN methods.

Here, we propose an improved version of *node2vec* that is more effective for weighted graphs by taking into account the edge weight connecting the previous vertex and the potential next vertex. The proposed method, called *node2vec+*, is a natural extension of *node2vec*; when the input graph is unweighted, the resulting embeddings of *node2vec+* and *node2vec* are equivalent in expectation. Moreover, when the bias parameters are set to neutral, *node2vec+* recovers a first-order random walk, just as *node2vec* does. In order to demonstrate the utility of *node2vec+*, we empirically show that 1) *node2vec+* is more robust to additive noise than *node2vec* using two synthetic datasets; and 2) *node2vec+* can achieve equivalent or better performance in multi-label node classification tasks using gene labels in biological networks compared to *node2vec* and GNNs including GCN (Kipf and Welling 2016) and GraphSAGE (Hamilton, Ying, and Leskovec 2017).

Method

We start by briefly reviewing the *node2vec* method. Then we illustrate that *node2vec* is less effective for weighted graphs due to its inability to identify *out* edges. Finally, we present a natural extension of *node2vec* that resolves this issue.

Node2vec overview

In the setting of node embeddings, we are interested in finding a mapping $f : V \rightarrow \mathbb{R}^d$ that maps each node $v \in V$ to a d -dimensional vector, so that the mutual proximity between pairs of nodes in the graph is preserved. In particular, a random walk based approach aims to maximize the probability of reconstructing the neighborhoods for any node in the graph, based on some sampling strategy S . Formally, given a graph $G = (V, E)$ (the analysis generalizes to directed and/or weighted graphs), we want to maximize the log probability of reconstructing the sampled neighborhood $\mathcal{N}_S(u)$ for each $u \in V$:

$$\max_f \sum_{u \in V} \log \mathbb{P}(\mathcal{N}_S(u) | f(u)) \quad (1)$$

Under the conditional independence assumption, and the parameterization of the probabilities as the softmax normalized inner products (Grover and Leskovec 2016; Mikolov et al. 2013a), the objective function above simplifies to:

$$\max_f \sum_{u \in V} \left(\sum_{v \in \mathcal{N}_S(u)} \langle f(v), f(u) \rangle - \log Z_u \right) \quad (2)$$

In practice, the partition function Z_u is approximated by negative sampling (Mikolov et al. 2013b) to save computational time. Given any sampling strategy S , equation (2) can find the corresponding embedding f , which is achieved in practice by feeding the random walks generated to the skipgram with negative sampling (Mikolov et al. 2013a).

Node2vec devises a second order random walk as the sampling strategy. Unlike a first order random walk (Perozzi, Al-Rfou, and Skiena 2014), where the transition probability $\mathbb{P}(c_i | c_{i-1})$ depends only on the current vertex c_{i-1} , a second order random walk depends also on the previous vertex c_{i-2} , with transition probability $\mathbb{P}(c_i | c_{i-1}, c_{i-2})$. It does so by applying a bias factor $\alpha_{pq}(c_{i-2}, c_i)$ to the edge $(c_i, c_{i-1}) \in E$ that connects the current vertex and a potential next vertex. This bias factor is a function that depends on the relation between the previous vertex and the potential next vertex, and is parameterized by the *return* parameter p , and the *in-out* parameter q . For the ease of notation, in the following, we denote $x = c_i$ as the potential next vertex, $v = c_{i-1}$ as the current vertex, and $t = c_{i-2}$ as the previous vertex (see Figure 1). In this way, the random walk can be generated based on the following transition probabilities:

$$\mathbb{P}(x|v, t) = \begin{cases} \frac{\alpha_{pq}(t,x)w(v,x)}{\sum_{x' \in \mathcal{N}(v)} \alpha_{pq}(t,x')w(v,x')} & \text{if } (v, x) \in E \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

where the bias factor is defined as:

$$\alpha_{pq}(t, x) = \begin{cases} \frac{1}{p} & \text{if } t = x \\ 1 & \text{if } t \neq x \text{ and } (t, x) \in E \\ \frac{1}{q} & \text{if } t \neq x \text{ and } (t, x) \notin E \end{cases} \quad (4)$$

According to this bias factor, *node2vec* differentiates three types of edges: 1) the *return* edge, where the potential next vertex is the previous vertex (Figure 1a); 2) the *out* edge, where the potential next vertex is *not* connected to the previous vertex (Figure 1b); and 3) the *in* edge, where the potential next vertex is connected to the previous vertex (Figure 1c). Note that the first order (or unbiased) random walk can be seen as a special case of the second order random walk where both the *return* parameter and the *in-out* parameter are set to neutral ($p = 1, q = 1$).

We now turn our attention to weighted networks, where the edge weights are not necessarily zeros or ones. Consider the case where t and x are connected, but with a small weight (Figure 1d), i.e. $(t, x) \in E$ and $0 < w(t, x) \ll 1$. According to the definition of the bias factor, no matter how small $w(t, x)$ is, (v, x) would always be considered as an *in* edge. Since in this case t and x are barely connected, (v, x) should in fact be considered as an *out* edge. In the extreme case of a fully connected weighted graph, where $(u, v) \in E$ for all $u, v \in V$, *node2vec* completely loses its ability to identify *out* edges.

Thus, *node2vec* is less effective for weighted networks due to its inability to identify potential *out* edges where the terminal vertex x is loosely connected to a previous vertex t . Next, we propose an extension of *node2vec* that resolves this issue, by taking into account of the edge weight $w(t, x)$ in the bias factor.

Node2vec+

The main idea of extending *node2vec* is to identify potential out edges $(v, x) \in E$ coming from node t , where x is loosely connected with t . Intuitively, we can determine the “looseness” of (v, x) based on some threshold edge value.

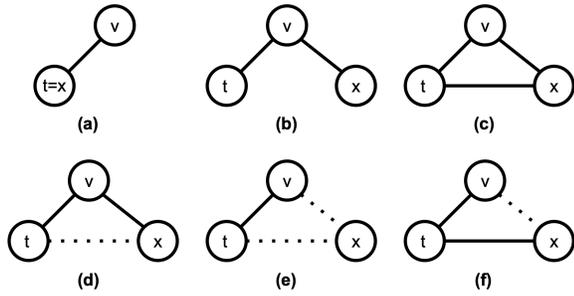


Figure 1: Illustration of different settings of *return* and *in-out* edges. The solid and dotted lines represent edges with large and small edge weights, respectively.

However, given that the distribution of edge weights of any given node in the graph is not known *a priori*, it is hard to come up with a reasonable threshold value for all networks. Instead, we define the “looseness” of (v, x) by comparing it with the average edge weights for that node v . Formally, a node $v \in V$ is said to be *loosely connected* to $u \in V$ if $w(u, v) < \tilde{d}(u)$, where $\tilde{d}(u) = \frac{\sum_{v' \in \mathcal{N}(u)} w(u, v')}{|\mathcal{N}(u)|}$ is the average edge weight of u . For the ease of notation, we also consider v being *loosely connected* to u if $(u, v) \notin E$, as a slight abuse of notation. Finally, an (directed or undirected) edge (u, v) is *loose* if v is loosely connected to u , and otherwise it is *tight*.

Based on the definition of looseness of edges, and assuming $t \neq x$, there are four types of (v, x) edges (see Figure 1, (c-f)), depending on the looseness of (v, x) and (x, t) . Following *node2vec*, we categorize these edge types into *in* and *out* edges. Furthermore, to prevent amplification of noisy connections, we added one more edge type called the *noisy* edge, which is always suppressed.

Out edge (1b, 1d) As a direct generalization to *node2vec*, we consider (v, x) to be an *out* edge if (v, x) is tight and (x, t) is loose. The *in-out* parameter q then modifies the out edge to differentiate “inward” and “outward” nodes, and subsequently leads to Breadth First Search or Depth First Search like searching strategies (Grover and Leskovec 2016). Unlike *node2vec*, however, we further parameterize this bias factor based on $w(x, t)$, and for simplicity, we choose to use linear interpolation. Specifically, for an *out* edge (v, x) , the bias factor is computed as $\alpha_{pq}(t, v, x) = \frac{1}{q} + (1 - \frac{1}{q}) \frac{w(x, t)}{\tilde{d}(x)}$. Thus the amount of modification to the *out* edge depends on the level of looseness of (x, t) . When $w(x, t) = 0$, or equivalently $(x, t) \notin E$, the bias factor for (v, x) is $\frac{1}{q}$, same as that defined in *node2vec*.

Noisy edge (1e) We consider (v, x) to be a *noisy* edge if both (v, x) and (x, t) are loose. The *noisy* edges are not very informative and thus should be suppressed in all cases regardless of the setting of q to prevent amplification of noise. Thus, the bias factor for a *noisy* edge is set to be $\min\{1, \frac{1}{q}\}$.

In edge (1c, 1f) Finally, we consider (v, x) to be an *in* edge if (x, t) is tight, regardless of $w(v, x)$. The corresponding

bias factor is set to neutral as in *node2vec*.

Combining the above, the bias factor for *node2vec+* is defined as follows:

$$\alpha_{pq}(t, v, x) = \begin{cases} \frac{1}{p} & \text{if } t = x \\ 1 & \text{if } w(x, t) \geq \tilde{d}(x) \\ \min\{1, \frac{1}{q}\} & \text{if } w(x, t) < \tilde{d}(x) \\ & \text{and } w(v, x) < \tilde{d}(v) \\ \frac{1}{q} + (1 - \frac{1}{q}) \frac{w(x, t)}{\tilde{d}(x)} & \text{if } w(x, t) < \tilde{d}(x) \\ & \text{and } w(v, x) \geq \tilde{d}(v) \end{cases} \quad (5)$$

Note that the last case in equation (5) includes cases where $(x, t) \notin E$. Based on the biased random walk searching strategy using this bias factor, the embedding can be generated accordingly using (2). One can verify, by checking equation (5), that this is indeed a natural extension of *node2vec* in the sense that

- For an unweighted graph, the *node2vec+* is equivalent to *node2vec*.
- When p and q are set to 1, *node2vec+* recovers a first order random walk, same as *node2vec* does.

Finally, by design, *node2vec+* is able to identify potential *out* edges that would have been obliterated by *node2vec*. *Node2vec+* is implemented as part of PecanPy (Liu and Krishnan 2021).

Experiments

Synthetic datasets

We first demonstrate the ability of *node2vec+* to identify potential *out* edges in weighted graphs using a barbell graph and the hierarchical cluster graphs.

Barbell graph A barbell graph, denoted as B , is constructed by connecting two complete graphs of size 20 with a common bridge node (Figure 2a). All edges in B are weighted 1. There are three types of nodes in B , 1) the bridge node; 2) the peripheral nodes that connect the two modules with the bridge node; 3) the interior nodes of the two modules. By changing the *in-out* parameter q , *node2vec* could put the peripheral nodes closer to the bridge node or interior nodes in the embedding space.

When q is large, *node2vec* suppresses the *out* edges, e.g., an edge connecting a peripheral node to the bridge node, coming from an interior node. Consequently, the biased random walks are restricted to the network modules. In this case, the transition from the peripheral nodes to the bridge node becomes less likely compared to a first-order random walk, thus pushing the embeddings between the bridge node and the peripheral nodes away from each other. Conversely, when q is small, the transition between the peripheral nodes and the bridge node is encouraged. In this case, the embeddings of the bridge node and the peripheral nodes are pulled together. To see this, we run *node2vec* with fixed $p = 1$, and three different settings of $q = [1, 100, 0.01]$. Indeed, for $q = 100$, *node2vec* tightly clusters interior nodes and pushed the bridge node away from the peripheral nodes, and

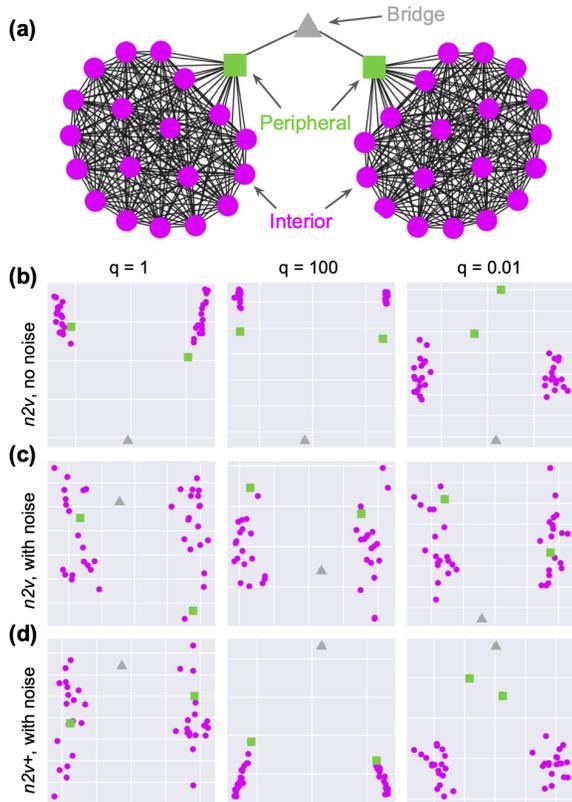


Figure 2: Barbell graph. (a) Illustration of the barbell graph, and three different types of nodes. (b)-(d) Embeddings of the barbell graph (or the noisy version of it), with three different settings of $q = [1, 100, 0.01]$, and the color coding is the same as that in (a).

for $q = 0.01$, the peripheral nodes are pushed away from the interior nodes (Figure 2b).

Next, the barbell graph is perturbed by adding loose edges with edge weights of 0.1, making the graph fully connected. This perturbed barbell graph is denoted \tilde{B} . As expected, in this case, *node2vec* can not make any difference in the embedding space by changing q (Figure 2c), since none of the edges are identified as an *out* edge. On the other hand, *node2vec+* still picks out potential *out* edges and thus produces the desired outcome (Figure 2d). Note that both *node2vec* and *node2vec+* have similar results for \tilde{B} when $q = 1$. This confirms that *node2vec+* and *node2vec* are equivalent when p and q are set to neutral, corresponding to embedding with unbiased random walks. Finally, when using non-neutral settings of q , *node2vec+* is able to suppress some noisy edges, resulting in less scattered embeddings of the interior nodes (Figure 2d).

Hierarchical CLUSTER graph We use a modified version of the CLUSTER dataset (Dwivedi et al. 2020) to further demonstrate the advantage of the *node2vec+* due to identifying potential *out* edges. Specifically, the hierarchical cluster graph K3L2 contains $L = 2$ levels (3 including the root level) of clusters, and each parent cluster is asso-

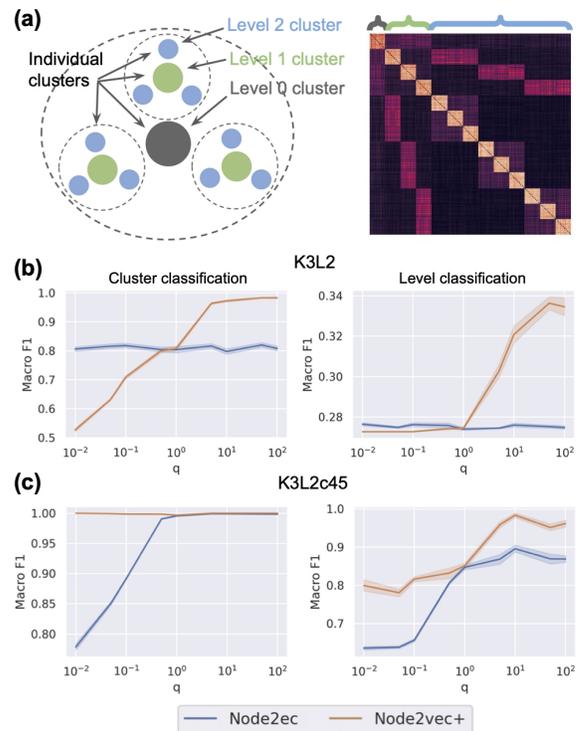


Figure 3: Hierarchical CLUSTER graph classification task. (a) Illustrations of the K3L2 hierarchical clusters. Left: top-down view of the clusters. Right: adjacency matrix of K3L2; colored brackets indicate the corresponding cluster levels of the nodes. (b) Classification evaluation on K3L2. (c) Classification evaluation on K3L2c45

ciated with $K = 3$ children clusters (Figure 3a). There are 30 nodes in each cluster, resulting in a total of 390 nodes. To generate the hierarchical cluster graph, we first generate data points via a Gaussian process in a latent space so that the Euclidean distance between two points from two sibling clusters is about twice ($\sqrt{2}$ to be precise) the expected Euclidean distance from one of the two points to a point in the parent cluster, which is set to be 1. The noisiness of the clusters is controlled by the parameter σ , which is set to 0.1 by default. These data points are then turned into a fully connected weighted graph using a RBF kernel (see supplement).

Similar to CLUSTER, the classification task is to classify each node in the hierarchical cluster graph into the corresponding cluster (*cluster classification*). On the other hand, we can also classify the nodes based on the level of their corresponding clusters (*level classification*). We perform stratified split to separate nodes into 10% training and 90% testing. We use the multinomial logistic regression model with l2 regularization for prediction, and evaluate the performance by macro f1 score. This evaluation process, including the embedding generation, is repeated ten times. As shown in Figure 3b, the performance of *node2vec* is not affected by the q parameter because the graph is fully connected. Meanwhile, *node2vec+* achieves significantly better performance than *node2vec* for large q settings for both classification

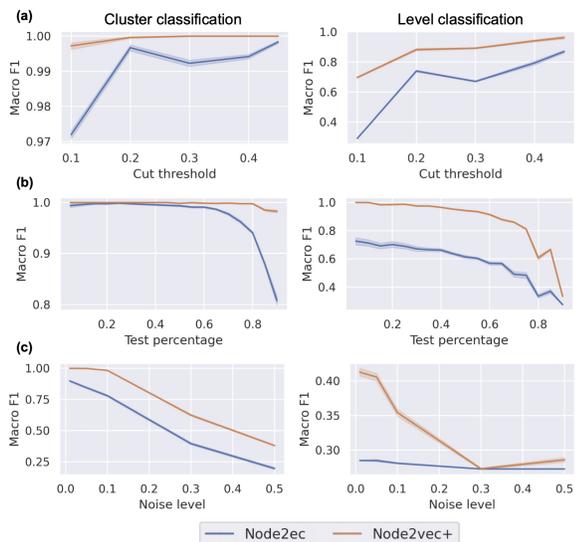


Figure 4: Fine-grained performance evaluation on K3L2. (a) Varying cut threshold. (b) Varying percentage of testing. (c) Varying noise level.

tasks. This is expected since *node2vec+* identifies potential *out* edges and in the setting of large q , *out* walks are discouraged, leading to random walks that are localized in tightly clustered network modules. And as illustrated in the previous section, this localization of the random walk causes the nodes in a network module to be oriented closer together in the embedding space, which is particularly helpful for distinguishing different clusters in the embedding space. Nevertheless, we note that the performance of *node2vec+* for *level classification* on K3L2 is sub-optimal, despite being significantly better than *node2vec*. Similar results are observed on a couple different hierarchical cluster graphs K3L3, K5L1, and K5L2 (see supplement).

However, it is not surprising that *node2vec+* does better than *node2vec* because K3L2 is fully connected. To make a more fair comparison, we maximally sparsify K3L2 using a global edge threshold value of 0.45, which preserves the connectdness of the graph (see supplement). The resulting sparsified graph is denoted as K3L2c45. In this case, *node2vec* shows significant improvement for large q compared to $q = 1$, resulting in perfect prediction for the *cluster classification* task (Figure 3c). Meanwhile, *node2vec+* achieves perfect *cluster classification* prediction too and, more importantly, also achieves near perfect *level classification* performance. This drastic improvement in *level classification* for *node2vec+* as a result of sparsification suggests that even *node2vec+* cannot perfectly resolve the problem with fully connected weighted graphs, and further improvement can be made.

We then compare the methods in more fine-grained settings by varying the percentage of testing data, the noise level, and the cut threshold. Since both *node2vec* and *node2vec+* favor large q in this dataset, we fix $p = 1$ and $q = 100$ throughout the analyses on K3L2. Overall, *node2vec*

and *node2vec+* both admit similar trends that 1) the prediction performance increases as the cut ratio increases, 2) the prediction performance decreases when the testing size increases, leaving a smaller training size, and when the noise increases (Figure 4). Notice that reducing the training size (larger testing size) and including more edges with small edge weights (smaller cut threshold) have less detrimental effects on *node2vec+* than *node2vec*, particularly for the *cluster classification* task. Finally, *node2vec+* consistently outperforms (or at least is comparable to) *node2vec* across all different settings and tasks, highlighting the benefits of identifying potential out edges.

Real world datasets

We further assess the quality of node embeddings generated by *node2vec+* using node classification tasks on real world datasets. Two commonly used datasets, BlogCatalog and Wikipedia, are tested. In addition, two popular protein-protein interaction (PPI) networks STRING (Szkarczyk et al. 2021) and GIANT-TN (Greene et al. 2015) are also tested using a variety of challenging gene classification tasks. Following our previous approach for making a more fair comparison on GIANT-TN, we also include a sparse GIANT-TN-c01 network, by applying a global edge threshold of 0.01 to GIANT-TN, which preserves the connectivity of the network (see supplement). More detailed information for each network can be found in Table 1.

Baseline Methods In our evaluation on BlogCatalog and Wikipedia, the main comparisons are between *node2vec* and *node2vec+*, as before. The key parameters including embedding dimension, window-size, walk-length, and number of walks per node are set to 128, 10, 80, and 10, respectively, by default. We exclude the comparison against some popular node embedding methods like DeepWalk (Perozzi, Al-Rfou, and Skiena 2014) and LINE (Tang et al. 2015), as they were shown to be inferior to *node2vec* (Grover and Leskovec 2016). We also exclude some neural network based methods like GraRep (Cao, Lu, and Xu 2015) and deepNF (Gligoric-jevic, Barot, and Bonneau 2018), since they can not scale efficiently to large and dense networks.

For gene classification on STRING and GIANT-TN, we include the comparison against two popular GNNs, GCN (Kipf and Welling 2016) and GraphSAGE (Hamilton, Ying, and Leskovec 2017), which have shown exceptional performance in many node classification tasks. We use full-batch GraphSAGE with mean pooling aggregation following the Open Graph Benchmark (Hu et al. 2021). To match the dimension for *node2vec* embeddings for a fair comparison,

Network	Weighted	$ V $	Edge density
BlogCatalog	No	10,312	6.28E-03
Wikipedia	Yes	4,777	8.10E-03
GIANT-TN	Yes	25,825	1.00E+00
GIANT-TN-c01	Yes	25,689	1.18E-01
STRING	Yes	17,352	2.42E-02

Table 1: Network information

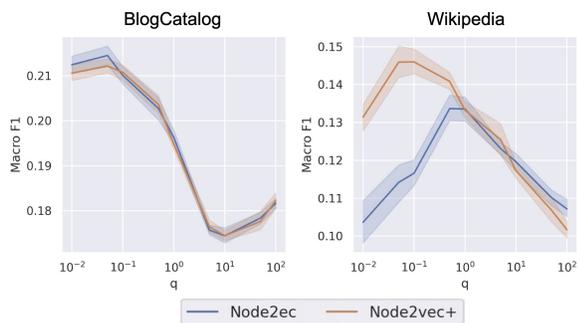


Figure 5: Node classification performance for BlogCatalog and Wikipedia

we use one hidden layer models for both GCN and GraphSAGE, with dimensions 128 for GCN and 64 for GraphSAGE. Since the PPIs here do not come with node features, we use constant feature for GCN, and degree feature for GraphSAGE.

As mentioned earlier, one of the challenges for training GNNs on datasets like the gene classification tasks here is the data imbalance issue, with many negative examples but only a few positive examples. Specifically, when using the GIANT-TN network, the ratios of negatives to positives are on average 30, 19, and 33 for GOBP, KEGGBP, and DisGeNet, respectively. To mitigate this data imbalance issue, we up-weight the loss for positive examples using the corresponding negative to positive ratio. We observe that by doing so, the model converges significantly faster than if we did not include the scaling ratios (see supplement). Additionally, we tuned the learning rate via grid search from 10^{-5} to 10^{-1} . The optimal learning rates that result in decent convergence rate without diverging are 0.01 and 0.0005 for GCN and GraphSAGE, respectively (see supplement). Finally, we set the maximum epochs to 100,000 so that the models are decently converged (see supplement).

Experiment setup We follow *node2vec* and randomly split the datasets into 50% training and 50% testing for BlogCatalog and Wikipedia. The networks are then embedded using varying values of q from 0.01 up to 100. In this evaluation, we fix the p parameter at neutral since the biased factor for p was not modified in *node2vec+* algorithm and should have the same effect as in *node2vec*. Then, a one vs rest logistic regression model with l2 regularization is trained on the embedding. The final test performance is reported as the macro f1 score. This evaluation process, including the embedding generation step, is repeated 10 times.

Meanwhile, for the PPIs, we test three types of gene classification tasks including two gene function predictions (KEGGBP, GOBP) and one disease gene prediction (DisGeNet), which are obtained as in (Liu et al. 2020). We call each type of these gene classification tasks a gene set collection. Similar as before, we train a one vs rest logistic regression with l2 regularization using the embeddings from *node2vec* and *node2vec+*, while for GNNs, each gene set collection is used to train the models for a specific com-

bination of network and gene set collection. Thus, we also note that the comparisons between GNNs and embeddings are slightly unfair since the one vs rest logistic regression model for the embeddings has access to only one task, while a GNN model has access to labels from all the tasks in a gene set collection, effectively increasing the number of training examples for the GNN model.

We use two types of validation schemes for evaluation, 5-fold cross-validation and study-bias holdout, following (Liu et al. 2020). The cross-validation evaluations for GNNs are excluded due to extra computational time. In the study-bias holdout evaluation, we use 60% of the most well-studied genes, according to the number of PubMed publications associated with the genes in each dataset, for training, 20% of the least studied genes for testing, and the rest for validation. We filter out gene sets with less than ten positives in any splits for each PPI network. Thus, despite being a more realistic and rigorous evaluation scheme, the study-bias holdout is more restrictive than 5-fold cross-validation, leaving out a lot of gene sets that are scarce (Table 2). The validation set was used for tuning the p, q combinations via grid search for each gene set. This optimally tuned p, q combination is then used for the final testing, reported as $\log_2 \frac{\text{auPRC}}{\text{prior}}$, which is more appropriate for evaluating performance on imbalanced data (Liu et al. 2020). Because of this tuning step, we consider the embedding approach being semi-supervised.

Experiment results For the two commonly used datasets, *node2vec+* is either comparable or better than *node2vec* (Figure 5). In particular, the identical performance on BlogCatalog confirms that *node2vec+* reduces to *node2vec* for unweighted graphs. On the other hand, for Wikipedia, *node2vec+* achieves the best result when q is around 0.1, and significantly outperforms *node2vec*, whose performance peaks at around $q = 1$. This performance difference is attributed to the ability of *node2vec+* to identify potential *out* edges, and consequently lead to more faithful *out*-biased walks.

For gene classification in the 5-fold cross-validation setting (Figure 6, top row), *node2vec+* significantly outperforms *node2vec* (Wilcoxon test p-value < 0.01) when using GIANT-TN for all three gene set collections. When using sparser networks like GIANT-TN-c01 and STRING, *node2vec+* still achieves similar performance as *node2vec* in most cases. We note that although sparsifying GIANT-TN to GIANT-TN-c01 improves the overall performance compared to GIANT-TN, and at the same time, reduces the performance gap between *node2vec* and *node2vec+*, finding the optimal sparsification is not a trivial task. For example, one can take various sparsification approaches, including

	GOBO	KEGGBP	DisGeNet
GIANT-TN	94 (41)	65 (28)	245 (102)
STRING	85 (35)	58 (26)	234 (97)

Table 2: Number of gene sets for each combinations of network and gene set collection in the 5-fold cross validation (study-bias holdout) setting.

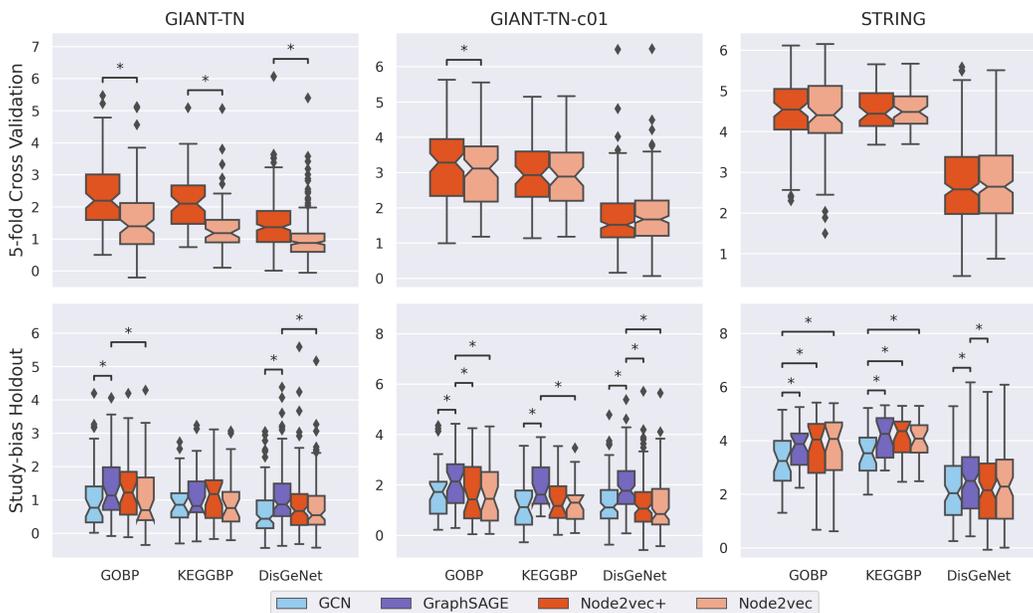


Figure 6: Gene classification tasks using protein-protein interaction networks. Starred (*) pairs indicate that the performance between two methods are significantly different (Wilcoxon p-value < 0.01).

the global edge thresholding used here, a more complicated node-specific edge thresholding, even spectral sparsification (Spielman and Teng 2010), and not to mention that each of these approaches come with parameters that need to be tuned. Our evaluation results here indicate that *node2vec+* does no worse than *node2vec*, and in some cases, when the weighted edges are noisy, and an appropriate graph sparsification is not trivial to be applied, *node2vec+* can perform significantly better than *node2vec*.

In the study-bias holdout setting (Figure 6, bottom row), *node2vec+* and *node2vec* are comparable with no significant difference. This is likely because most gene sets that distinguish the two methods, e.g., the ones that are scarcer and are harder to predict, are removed from the gene set collection during the necessary filtering step for the study-bias holdout. Meanwhile, *node2vec+* performs similarly to GCN, and in some cases, such as GOBP and KEGGBP using STRING, *node2vec+* significantly outperforms GCN. GraphSAGE, on the other hand, performs much better than GCN. Yet, except for GIANT-TN-c01, *node2vec+* still holds up well against GraphSAGE. Overall, despite having a substantial advantage of more training examples and being fully supervised, GCN and GraphSAGE did not entirely outperform the semi-supervised approach *node2vec+*.

Discussion and conclusion

In this paper, we proposed an improved version of the second-order random walk in *node2vec* for weighted graphs by considering the edge weights, which effectively identifies potential out edges. Consequently, the corresponding embeddings are improved whenever the in-out walks are appropriate for the task. The proposed modification *node2vec+* is a natural extension of *node2vec* for weighted graphs. We

empirically confirmed, by the evaluation on BlogCataog, that *node2vec+* reduces to *node2vec* for unweighted graphs. We note that *node2vec+* also serves as a general procedure for biased random walks and thus can be easily adapted to other methods that are built on top of *node2vec* such as KG2Vec (Wang et al. 2021b) and Het-Node2vec (Valentini et al. 2021).

We illustrated the ability of *node2vec+* to identify potential out edges on weighted graphs, as opposed to *node2vec*, using a synthetic barbell graph with uniformly added noise. Furthermore, using synthetic hierarchical cluster graphs and real-world datasets like Wikipedia, we showed that in many cases, when the graphs are weighted, *node2vec+* outperforms *node2vec*. Finally, we evaluated *node2vec+* against GCN and GraphSAGE using a variety of challenging gene classification tasks. We effectively managed the data imbalance issue for GNNs by applying scaling ratios of negatives to positives in the loss function. In our setup, although GNN methods have some clear advantages over the *node2vec+*, i.e., having more training examples and being fully supervised, the performance of the *node2vec+* is still comparable to the GNNs on the GIANT-TN and STRING networks.

Although *node2vec+* does not outperform GraphSAGE, we emphasize that our main contribution was improving the biased random walk strategy. In fact, as a future direction, we would like to explore the unsupervised GraphSAGE, which is trained similarly to the skipgram with negative sampling using the first-order random walks on the graph (Equation (2)), by replacing the first-order random walks with the *node2vec+* biased random walks. Using *node2vec+* as the biased random walk strategy is a promising strategy for further improving unsupervised GraphSAGE because 1) *node2vec* outperforms DeepWalk because of more flex-

ible searching strategies, and 2) *node2vec*+ identifies potential out edges, thus providing the intended biased searching strategies on weighted graphs.

References

- Ata, S. K.; Ou-Yang, L.; Fang, Y.; Kwok, C.-K.; Wu, M.; and Li, X.-L. 2018. Integrating node embeddings and biological annotations for genes to predict disease-gene associations. *BMC Systems Biology*, 12(9): 138.
- Bronstein, M. M.; Bruna, J.; Cohen, T.; and Veličković, P. 2021. Geometric Deep Learning: Grids, Groups, Graphs, Geodesics, and Gauges. *arXiv:2104.13478 [cs, stat]*. ArXiv: 2104.13478.
- Cao, S.; Lu, W.; and Xu, Q. 2015. GraRep: Learning Graph Representations with Global Structural Information. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management, CIKM '15*, 891–900. New York, NY, USA: Association for Computing Machinery. ISBN 978-1-4503-3794-6.
- Chattopadhyay, S.; and Ganguly, D. 2020. Community Structure aware Embedding of Nodes in a Network. *arXiv:2006.15313 [physics]*. ArXiv: 2006.15313.
- Cui, P.; Wang, X.; Pei, J.; and Zhu, W. 2018. A Survey on Network Embedding. *IEEE Transactions on Knowledge and Data Engineering*, 1–1.
- Davison, A.; and Austern, M. 2021. Asymptotics of Network Embeddings Learned via Subsampling. *arXiv:2107.02363 [cs, math, stat]*. ArXiv: 2107.02363.
- Dwivedi, V. P.; Joshi, C. K.; Laurent, T.; Bengio, Y.; and Bresson, X. 2020. Benchmarking Graph Neural Networks. *arXiv:2003.00982 [cs, stat]*. ArXiv: 2003.00982.
- Gligorijević, V.; Barot, M.; and Bonneau, R. 2018. deepNF: deep network fusion for protein function prediction. *Bioinformatics (Oxford, England)*, 34(22): 3873–3881.
- Greene, C. S.; Krishnan, A.; Wong, A. K.; Ricciotti, E.; Zelaya, R. A.; Himmelstein, D. S.; Zhang, R.; Hartmann, B. M.; Zaslavsky, E.; Sealfon, S. C.; Chasman, D. I.; FitzGerald, G. A.; Dolinski, K.; Grosser, T.; and Troyanskaya, O. G. 2015. Understanding multicellular function and disease with human tissue-specific networks. *Nature genetics*, 47(6): 569–576.
- Grohe, M. 2020. word2vec, node2vec, graph2vec, X2vec: Towards a Theory of Vector Embeddings of Structured Data. *PODS*.
- Grover, A.; and Leskovec, J. 2016. Node2Vec: Scalable Feature Learning for Networks. In *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16*, 855–864. New York, NY, USA: ACM. ISBN 978-1-4503-4232-2. Event-place: San Francisco, California, USA.
- Hacker, C. 2021. k-simplex2vec: a simplicial extension of node2vec. *arXiv:2010.05636 [cs, math]*. ArXiv: 2010.05636.
- Hamilton, W. L.; Ying, R.; and Leskovec, J. 2017. Inductive Representation Learning on Large Graphs. *arXiv:1706.02216 [cs, stat]*. ArXiv: 1706.02216.
- Hu, F.; Liu, J.; Li, L.; and Liang, J. 2020. Community detection in complex networks using Node2vec with spectral clustering. *Physica A: Statistical Mechanics and its Applications*, 545: 123633.
- Hu, W.; Fey, M.; Zitnik, M.; Dong, Y.; Ren, H.; Liu, B.; Catasta, M.; and Leskovec, J. 2021. Open Graph Benchmark: Datasets for Machine Learning on Graphs. *arXiv:2005.00687 [cs, stat]*. ArXiv: 2005.00687.
- Kipf, T. N.; and Welling, M. 2016. Semi-Supervised Classification with Graph Convolutional Networks. *arXiv:1609.02907 [cs, stat]*. ArXiv: 1609.02907.
- Liu, R.; and Krishnan, A. 2021. PecanPy: a fast, efficient and parallelized Python implementation of node2vec. *Bioinformatics*, (btab202).
- Liu, R.; Mancuso, C. A.; Yannakopoulos, A.; Johnson, K. A.; and Krishnan, A. 2020. Supervised learning is an accurate method for network-based gene classification. *Bioinformatics*, 36(11): 3457–3465. [eprint: https://academic.oup.com/bioinformatics/article-pdf/36/11/3457/33329234/btaa150.pdf](https://academic.oup.com/bioinformatics/article-pdf/36/11/3457/33329234/btaa150.pdf).
- Mikolov, T.; Chen, K.; Corrado, G.; and Dean, J. 2013a. Efficient Estimation of Word Representations in Vector Space. *arXiv:1301.3781 [cs]*. ArXiv: 1301.3781.
- Mikolov, T.; Sutskever, I.; Chen, K.; Corrado, G.; and Dean, J. 2013b. Distributed Representations of Words and Phrases and their Compositionality. *arXiv:1310.4546 [cs, stat]*. ArXiv: 1310.4546.
- Nelson, W.; Zitnik, M.; Wang, B.; Leskovec, J.; Goldenberg, A.; and Sharan, R. 2019. To Embed or Not: Network Embedding as a Paradigm in Computational Biology. *Frontiers in Genetics*, 10: 381.
- Perozzi, B.; Al-Rfou, R.; and Skiena, S. 2014. DeepWalk: Online Learning of Social Representations. *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining - KDD '14*, 701–710. ArXiv: 1403.6652.
- Qiu, J.; Dong, Y.; Ma, H.; Li, J.; Wang, K.; and Tang, J. 2018. Network Embedding as Matrix Factorization: Unifying DeepWalk, LINE, PTE, and node2vec. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, WSDM '18*, 459–467. New York, NY, USA: Association for Computing Machinery. ISBN 978-1-4503-5581-0.
- Spielman, D. A.; and Teng, S.-H. 2010. Spectral Sparsification of Graphs. *arXiv:0808.4134 [cs]*. ArXiv: 0808.4134.
- Szklarczyk, D.; Gable, A. L.; Nastou, K. C.; Lyon, D.; Kirsch, R.; Pyysalo, S.; Doncheva, N. T.; Legeay, M.; Fang, T.; Bork, P.; Jensen, L. J.; and von Mering, C. 2021. The STRING database in 2021: customizable protein–protein networks, and functional characterization of user-uploaded gene/measurement sets. *Nucleic Acids Research*, 49(D1): D605–D612.
- Tang, J.; Qu, M.; Wang, M.; Zhang, M.; Yan, J.; and Mei, Q. 2015. LINE: Large-scale Information Network Embedding. In *Proceedings of the 24th International Conference on World Wide Web, WWW '15*, 1067–1077. Republic and

Canton of Geneva, CHE: International World Wide Web Conferences Steering Committee. ISBN 978-1-4503-3469-3.

Valentini, G.; Casiraghi, E.; Cappelletti, L.; Ravanmehr, V.; Fontana, T.; Reese, J.; and Robinson, P. 2021. Het-node2vec: second order random walk sampling for heterogeneous multigraphs embedding. *arXiv:2101.01425 [physics]*. ArXiv: 2101.01425.

Wang, N.; Zeng, M.; Li, Y.; Wu, F.-x.; and Li, M. 2021a. Essential Protein Prediction Based on node2vec and XGBoost. *Journal of Computational Biology*, 28(7): 687–700. Publisher: Mary Ann Liebert, Inc., publishers.

Wang, Y.; Dong, L.; Jiang, X.; Ma, X.; Li, Y.; and Zhang, H. 2021b. KG2Vec: A node2vec-based vectorization model for knowledge graph. *PLOS ONE*, 16(3): e0248552. Publisher: Public Library of Science.

Wu, Z.; Pan, S.; Chen, F.; Long, G.; Zhang, C.; and Yu, P. S. 2021. A Comprehensive Survey on Graph Neural Networks. *IEEE Transactions on Neural Networks and Learning Systems*, 32(1): 4–24. ArXiv: 1901.00596.

Xia, F.; Sun, K.; Yu, S.; Aziz, A.; Wan, L.; Pan, S.; and Liu, H. 2021. Graph Learning: A Survey. *IEEE Transactions on Artificial Intelligence*, 1(01): 1–1. Publisher: IEEE Computer Society.

Yue, X.; Wang, Z.; Huang, J.; Parthasarathy, S.; Moosavinasab, S.; Huang, Y.; Lin, S. M.; Zhang, W.; Zhang, P.; and Sun, H. 2019. Graph Embedding on Biomedical Networks: Methods, Applications, and Evaluations. *arXiv:1906.05017 [cs]*. ArXiv: 1906.05017.

Zeng, M.; Li, M.; Fei, Z.; Wu, F.-X.; Li, Y.; Pan, Y.; and Wang, J. 2021. A Deep Learning Framework for Identifying Essential Proteins by Integrating Multiple Types of Biological Information. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 18(1): 296–305. Conference Name: IEEE/ACM Transactions on Computational Biology and Bioinformatics.

Zhang, X.-M.; Liang, L.; Liu, L.; and Tang, M.-J. 2021. Graph Neural Networks and Their Current Applications in Bioinformatics. *Frontiers in Genetics*, 12: 1073.

Zhou, J.; Cui, G.; Hu, S.; Zhang, Z.; Yang, C.; Liu, Z.; Wang, L.; Li, C.; and Sun, M. 2021. Graph Neural Networks: A Review of Methods and Applications. *arXiv:1812.08434 [cs, stat]*. ArXiv: 1812.08434.

Supplementary Materials for Accurately Modeling Biased Random walks on Weighted Graphs Using *Node2vec+*

Renming Liu,¹ Matthew Hirn,^{1,2,3,*} Arjun Krishnan^{1,4,*}

¹Department of Computational Mathematics, Science & Engineering, ²Department of Mathematics, ³Center for Quantum Computing, Science & Engineering, ⁴Department of Biochemistry and Molecular Biology, Michigan State University, East Lansing, MI 48824, USA

liurenmi@msu.edu, mhirn@msu.edu, arjun@msu.edu (*Corresponding authors)

A Hierarchical cluster graphs

A.1 Construction details

The hierarchical cluster graphs are constructed by first sampling in a representation space constructed based on the corresponding tree structure and then applying an RBF kernel.

Tree construction We first construct the cluster centroids using a tree structure. In particular, a *perfect binary tree* is a binary tree where all nodes except for the leaf nodes have two children, and all leaf nodes have the same level. This definition can be generalized to *perfect K-trees*, in which all the interior nodes have K number of nodes, for K greater than or equal to one. We denote $T_{K,L}$ as the *perfect K-tree* with maximum level L . Figure A.1 shows the example of $T_{2,2}$, a *perfect binary tree* (or *2-tree*) with a maximum level of two.

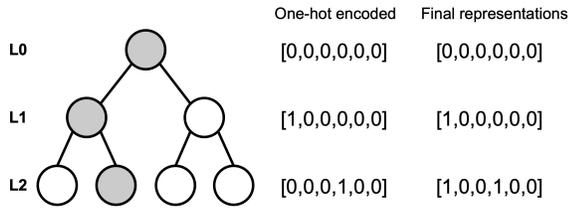


Figure A.1: A perfect binary tree with maximum level of two, with example one-hot encoded and final representations of the grey nodes.

Representing nodes in the tree A straightforward solution to represent the nodes in $T_{K,L}$ is one-hot encoding. For a more compact representation, we leave out the indicator for the root node, and represent it as all zeros. Thus, the dimension of the indicator array is equal to the total number of nodes in $T_{K,L}$, excluding the root node, which is the following

$$|V(T_{K,L})| - 1 = \sum_{l=1}^L K^l \quad (\text{A.1})$$

However, if we use one-hot encoding, then all nodes are equally distanced in the Euclidean space. Instead, we com-

bine the one-hot encoded representations of all the ancestor nodes as the final representation for each node, denoted as μ_i , $i = 1, \dots, |V(T_{K,L})|$. In this way, all sibling nodes are equally distanced, with $\sqrt{2}$ times the distance from the parent node. Figure A.1 shows the example of both the one-hot encoded and the final representations of the grey nodes. Notice the difference between the final representation and the one-hot encoded representation of the grey leaf node.

Hierarchical cluster We then draw data points x from a Gaussian distribution around each node in the $T_{K,L}$ tree:

$$x \sim \mathcal{N}(\mu_i, \sigma), \quad i = 1, \dots, |V(T_{K,L})| \quad (\text{A.2})$$

In the case of K3L2, the data points are drawn using $T_{3,2}$. The parameter σ , which controls the noisiness of the sampled data points, is set to 0.01 by default. Throughout the study, we fix the number of data points per node in the tree to 30. Finally, we turn the sampled data points into a fully connected weighted graph using the RBF kernel.

A.2 Sparsifying K3L2

We apply global edge threshold to K3L2 by removing all edges below a certain edge threshold value. Sweeping through $[0.01, 0.9]$, we found the maximum global edge threshold that preserves the connectivity of the graph at around 0.45 (Figure A.2). Notice that by doing so, the edge density drastically reduces to around 0.1.

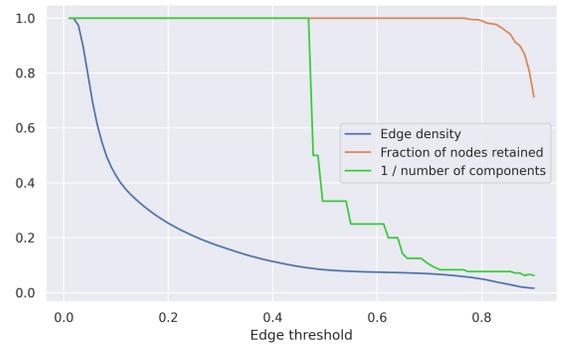


Figure A.2: Sparsifying K3L2 via global edge thresholding.

A.3 Other hierarchical cluster graph results

Similar to the K3L2 graph, we evaluate the cluster and level classification tasks using three other hierarchical cluster graphs, including K3L3, K5L1, and K5L2.

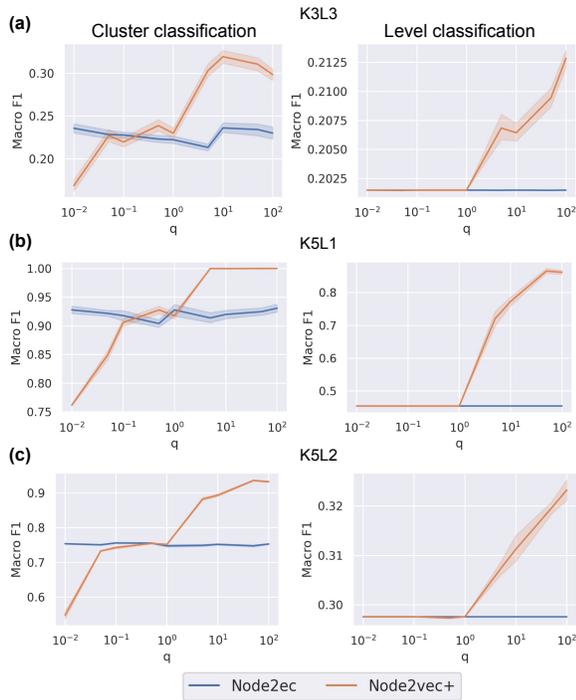


Figure A.3: Classification evaluation results for other hierarchical cluster graphs.

B Gene classification datasets

B.1 Sparsifying GIANT-TN

Using a similar approach for sparsifying the K3L2 graph, we maximally sparsify the GIANT-TN network. The maximal global edge threshold, in this case, is 0.01.

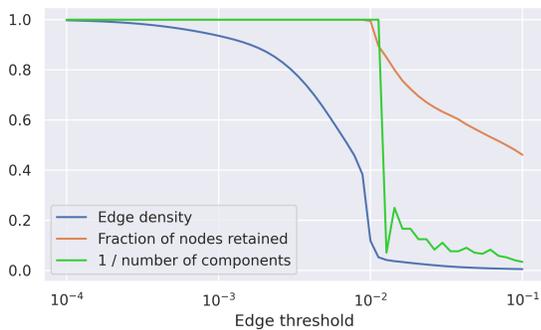


Figure B.1: Sparsifying GIANT-TN via global edge thresholding.

B.2 GIANT-TN fine-grained sparsification

In addition to the maximally sparsified GIANT-TN, we also tested two less restrictive edge thresholds, 0.001 and 0.005, resulting in GIANT-TN-c001 and GIANT-TN-c005, respectively. We observe that in these cases, node2vec+ still achieves equal or better performance than node2vec, using both the 5-fold cross validation schemes (Figure B.2) and the study-bias holdout (Figure B.3).

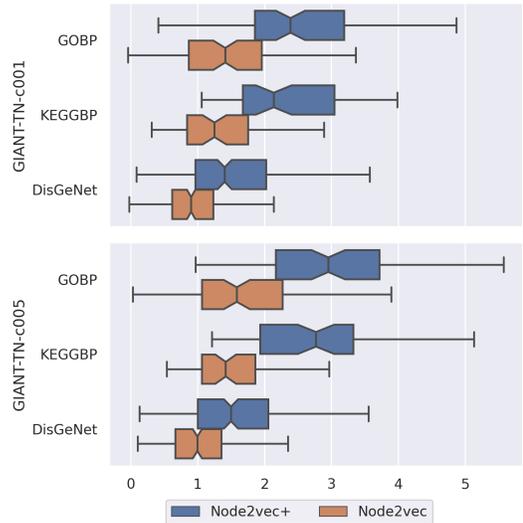


Figure B.2: Gene classification evaluation on fine-grained GIANT-TN sparsifications using 5-fold cross validation.

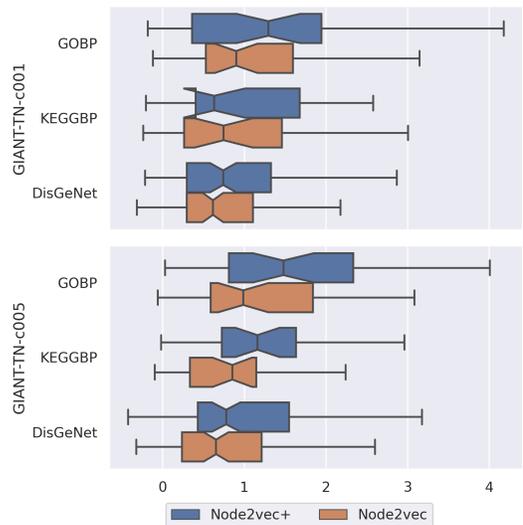


Figure B.3: Gene classification evaluation on fine-grained GIANT-TN sparsifications using study-bias holdout.

B.3 Mitigating data imbalance issue for GNN

We apply a scaling factor as the negative-to-positive ratio for each task to the loss function to mitigate the data imbalance

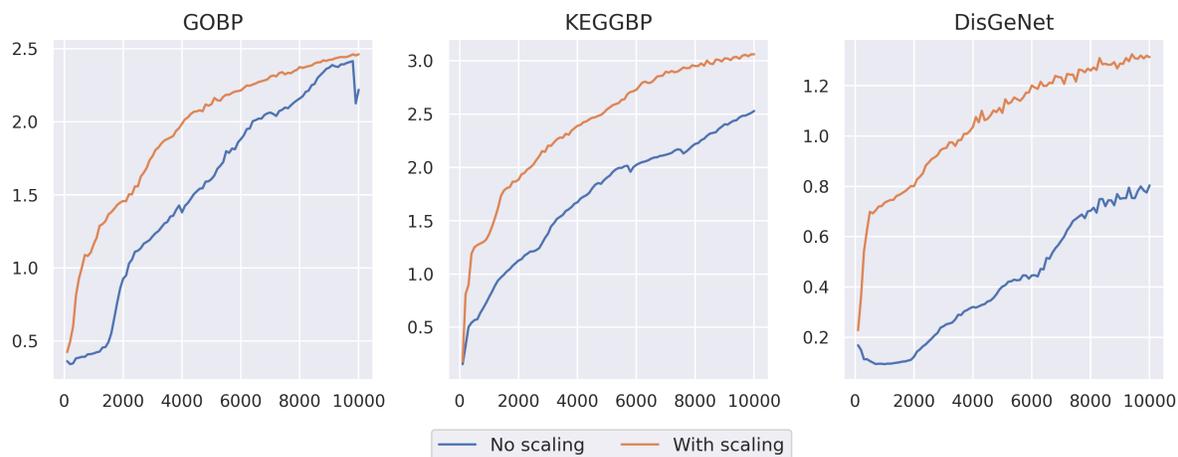


Figure B.4: Testing scores for GCN using the STRING network with and without scaling factors.

issue. The test results for GCN using the STRING network show that applying the scaling factors speeds up the convergence of the GCN (Figure B.4).

B.4 Tuning GNN learning rates

See Figure B.5 for tuning GCN learning rate and Figure B.6 for GraphSAGE.

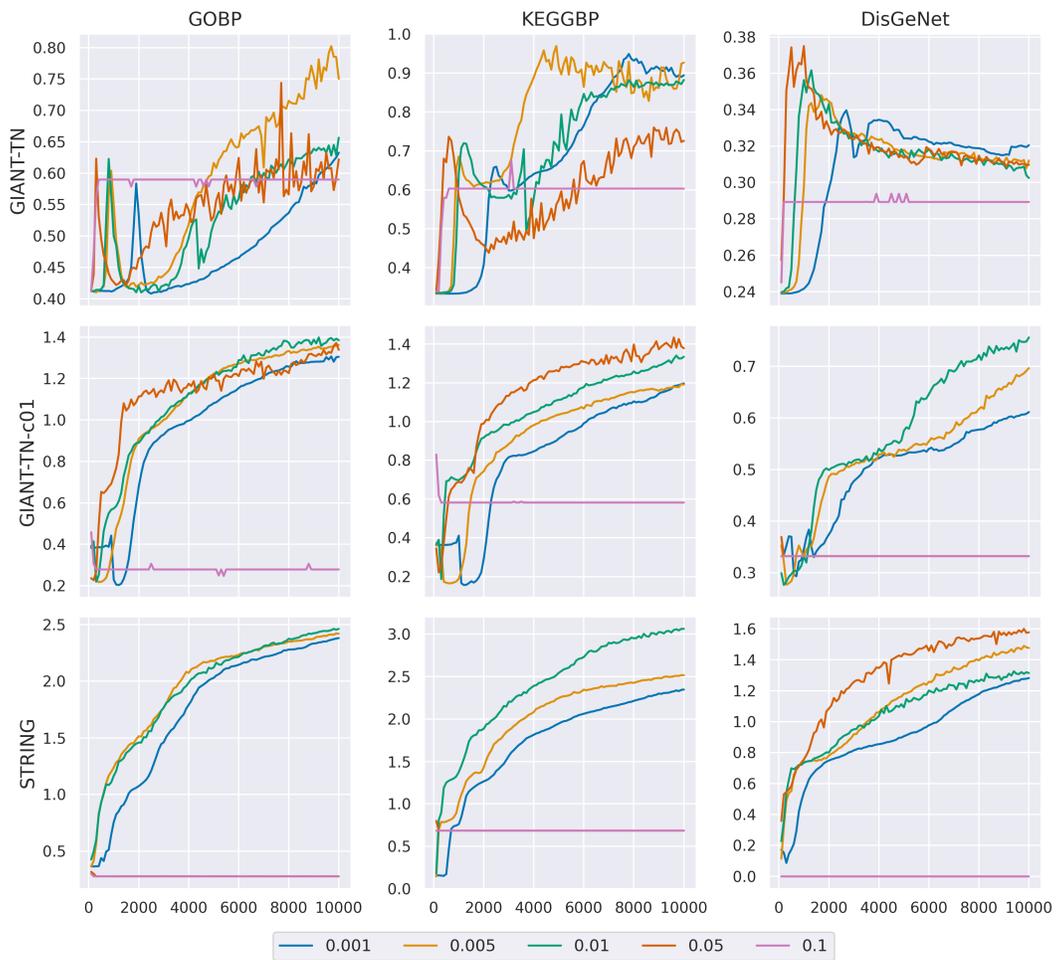


Figure B.5: Validation scores for GCN using different learning rates.

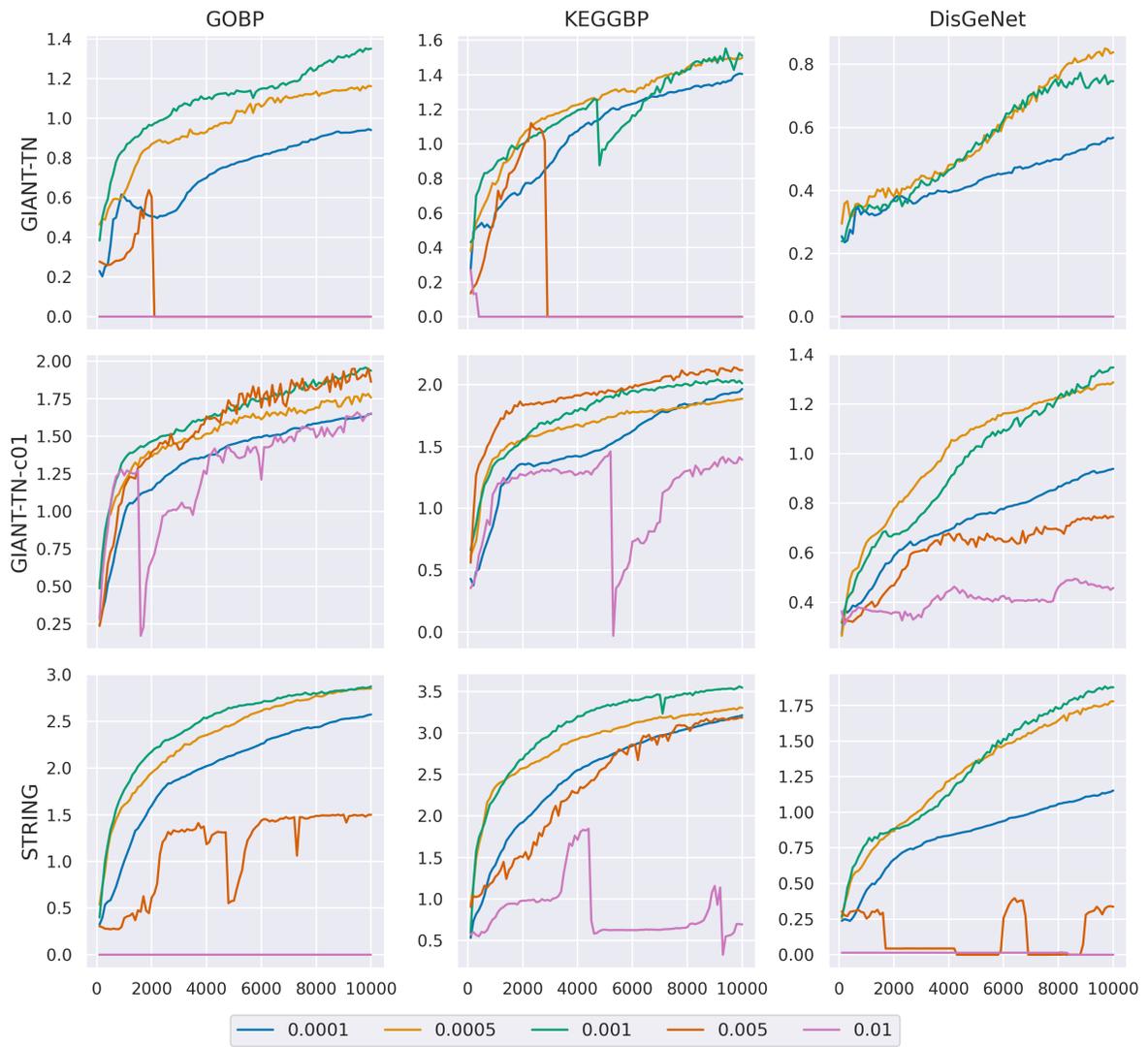


Figure B.6: Validation scores for GraphSAGE using different learning rates.