

# Stochastic Block Model-Aware Topological Neural Networks for Graph Link Prediction

Anonymous authors

Paper under double-blind review

## Abstract

Link prediction is an important learning task for graph-structured data and is indispensable to understanding graphs' properties. Recent works focus on designing complicated graph neural networks (GNNs) architectures to explore and capture various pairwise interactions among graph nodes. Most GNNs are based on combining graph structural and node feature information by iterative message-passing schemes. However, despite GNNs revolutionizing the field of graph representation learning, some thorny questions are raised concerning whether GNNs can efficiently learn the edge probabilities based on topological structures (i.e., higher-order interactions) and node features, and provide statistically rigorous uncertainty estimates. In this paper, we tackle these challenges and propose a novel stochastic block model (SBM)-aware topological neural networks, called SBM-TNN, that uses SBMs to infer the latent community structure of nodes from graph structures and uses persistent homology to encode higher-order information. Furthermore, we theoretically study the entrywise bound and asymptotic normality of the estimated edge probability matrix to quantify the uncertainty in statistical inference of the edge probabilities. Our extensive experiments for link prediction on both graphs and knowledge graphs show that SBM-TNN achieves state-of-the-art performance over a set of popular baseline methods.

## 1 Introduction

Graph data are ubiquitous throughout the natural and social sciences, e.g., many real-world objects can be represented by graphs, such as molecules, ecosystems, transportation systems, energy systems, citation networks, and internet networks Sen et al. (2008); Li et al. (2018); Xia et al. (2021); Chen et al. (2023). Tremendous advances in graph analysis have been achieved in recent years, especially in the field of geometric deep learning (GDL) Defferrard et al. (2016); Bronstein et al. (2017); Zhang et al. (2020). In particular, graph neural networks (GNNs) have emerged as effective architectures for various prediction problems, e.g., node classification Kipf & Welling (2017); Veličković et al. (2018a); Hamilton et al. (2017), community detection Chen et al. (2018); Shchur & Günnemann (2019), and graph classification Xu et al. (2018); Ying et al. (2018). Specifically, GNNs are neural network architectures designed to handle graph-structured data. The fundamental idea behind GNNs involves treating the underlying graph as a computation graph and leveraging neural network primitives to generate node embeddings. The key processes involve message passing, propagating, and aggregating node features and graph structural information throughout the graph.

In GNNs, the graph convolutional layer builds upon the observed graph adjacency matrix, also called the connection matrix. The adjacency matrix can be viewed as a noisy version of an edge probability matrix with additive noises. In practice, graphs often contain communities, and thus the probability of an edge between any two nodes can depend on their group assignment, which is popularly modeled through the stochastic block model (SBM) Holland et al. (1983) or its variant the degree-corrected stochastic block model (DCSBM) Karrer & Newman (2011). When a graph indeed has communities, we propose to use the estimated edge probability matrix to replace the adjacency matrix for graph representation learning. We estimate the edge probability matrix using the spectral clustering method von Luxburg (2007), which is computationally fast. Compared to the adjacency matrix, which only has two values 0 and 1, the estimated edge probability matrix contains the probability value of each edge, which can provide more information

about the relationship between nodes and better recover the graph structure. In addition, the estimated probability matrix is proven to have asymptotic normality, which can be used for uncertainty quantification and confidence interval estimation. By contrast, the adjacency matrix can not achieve this goal. We also establish a uniform error bound for the estimated edge probability matrix in terms of the entrywise maximum norm. These theoretical results are novel in the SBM literature. Furthermore, GNNs tend to majorly focus on information propagation among nodes and thus the model capability is limited, i.e., almost fail in learning topological and structural information. However, as recently shown by Wasserman (2018); Hensel et al. (2021); Pun et al. (2022), such topological structures, e.g., connected components and holes might be an important step in graph knowledge discovery. For instance, persistent homology Edelsbrunner (2013); Zomorodian & Carlsson (2004) has been used to study the topological information encoded in the graph Zhao & Wang (2019); Chen et al. (2021); Horn et al. (2021). However, these ideas have never been yet applied in conjunction with knowledge representation learning.

Aiming to solve the above challenges, we turn to the idea of combining learned community information from using an SBM-based model and different types of topological features by using multiple descriptor functions that can generate more expressive node embedding. The SBM is a powerful tool to learn from graph-structured data, as it is designed to model graphs with clear community structures. In GNNs, these communities represent groups of nodes having similar behaviors. By leveraging SBM in our topological neural networks (TNNs), the model can simultaneously learn the latent communities and capture the topological structures to enhance the performance in link prediction and node classification. Moreover, the community-based learning from SBM enables our TNNs to capture local structures to improve prediction accuracy, when dealing with sparse graphs with few connections, which in general is a great challenge for GNNs. SBM, on the other hand, provides a probabilistic framework that can help the model infer relationships between sparse connections as well as providing a better interpretability and understanding of the relationships between nodes. We use the learned community information from SBM together with different types of topological features by using multiple descriptor functions to generate more expressive latent node embedding. The primary contributions of this work can be summarized as follows:

- We propose a Stochastic Block Model-Aware Topological Neural Networks (SBM-TNN), a novel TNN-based model equipped with SBM concepts that captures topological structures, node features, and structure of neighborhood relations. It is the first approach bringing the concepts of topological signature representation learning and stochastic block models to graph learning.
- We further study the important problem, i.e., how to quantify the uncertainty for the estimated edge probability matrix. To achieve this goal, we establish an entrywise error bound and asymptotic normality of the estimated edge probability matrix that can be used to construct asymptotically valid confidence intervals for the edge probabilities, and help quantify the accuracy and uncertainty of the estimated edge probability matrix and provide theoretical guarantee for the follow up procedures of the proposed SBM-TNN model.
- Extensive experiments on benchmark datasets clearly show that SBM-TNN delivers state-of-the-art link prediction and knowledge graph completion tasks with a significant margin.

## 2 Related Work

**Graph Neural Networks.** Recently, Graph Neural Network (GNN) has emerged as a primary tool for node classification, link prediction, and graph forecasting Wu et al. (2020); Zhou et al. (2020). Based on the spectral graph theory, Bruna et al. (2014) introduces a graph-based convolution in the Fourier domain. However, the complexity of this model is very high since all Laplacian eigenvectors are needed. To tackle this problem, ChebNet Defferrard et al. (2016) integrates spectral graph convolution with Chebyshev polynomials. Then, Graph Convolutional Networks (GCNs) Kipf & Welling (2017) simplifies the graph convolution with a localized first-order approximation. SEAL Zhang & Chen (2018) extracts local enclosing subgraphs around the target links and learns a function mapping the subgraph patterns to link existence. Graph2Gauss (G2G) Bojchevski & Günnemann (2018) designs an unsupervised model that handles inductive link prediction by using a deep encoder to embed each node as a Gaussian distribution. Deep Generative Latent Feature Relational Model

(DGLFRM) Mehta et al. (2019) proposes an overlapping stochastic blockmodel for community discovery tasks. In addition, the Hyperbolic Graph Convolutional Neural Networks (HGNC) Chami et al. (2019) leverages both the hyperbolic geometry and GCN framework to learn node representations. Another interesting recent strategy is to use pairwise topological features to find latent representations of the geometrical structure of graph using GCN Yan et al. (2021). A common limitation is that they fail to accurately capture correlated and rich topological properties of graphs and incorporate rich structure and topological information both in local and global domains.

**Stochastic Block Model.** The SBM Holland et al. (1983) is a probabilistic model to generate networks with community structures, where nodes are partitioned into blocks and the probability of edges between nodes depends on their block memberships. The past few decades have seen various methods for recovering community memberships based on the observed network Abbe (2018). Spectral clustering stands out because of its computational tractability. The statistical properties of spectral clustering under SBM or its variant DCSBM have been widely studied. For example, the weak consistency of clustering, i.e., the proportion of misclustered nodes converges to zero as the number of nodes increases, has been investigated by Rohe et al. (2011); Lei & Rinaldo (2015); Joseph & Yu (2016), among others, and the strong consistency, namely, the memberships can be perfectly recovered in large samples, has been established in Su et al. (2019). Moreover, the minimax rate of the estimator for the edge probability matrix (i.e., the population counterpart of the adjacency matrix in terms of matrix norms such as Frobenius or spectral norm has been provided Gao et al. (2015). However, how accurately the spectral clustering method can estimate each entry of the edge probability matrix is unclear. We develop the entrywise bound for the estimated edge probability matrix, which is essential to quantify the uncertainty in statistical inference for the population counterpart of the adjacency matrix in SBMs and DCSBMs. In addition, we establish the asymptotic normality of the estimated edge probability matrix, which turns out to be asymptotic efficient, under SBMs and DCSBMs, so we can construct asymptotically entrywise confidence intervals for the probability matrix. The asymptotic Gaussian behavior for the estimators of the block probability matrices Tang et al. (2022) and the eigenvector matrix Tang & Priebe (2018); Cape et al. (2019); Xie (2024) have been studied under SBMs. However, the asymptotic behavior of the estimator for the edge probability matrix under more general DCSBMs is undeveloped.

### 3 Stochastic Block Model-Aware Topological Neural Networks: Undirected Graph

#### 3.1 Mixed-Up Undirected Graph Construction

To capture the topological information from the graph  $\mathcal{G}$  and node features, we construct a mixed-up graph  $\mathcal{G}_{\mathcal{M}} = (\mathbf{A}_{\mathcal{M}}, \mathbf{X})$  based on original input graph  $\mathcal{G}_{\mathcal{O}} = (\mathbf{A}_{\mathcal{O}}, \mathbf{X})$  and  $k$ -hop graph  $\mathcal{G}_{\mathcal{K}} = (\mathbf{A}_{\mathcal{K}}, \mathbf{X})$  (where  $\mathbf{A}_{\mathcal{M}}$ ,  $\mathbf{A}_{\mathcal{O}}$ , and  $\mathbf{A}_{\mathcal{K}}$  denote adjacency matrices of the mixed-up graph, original graph, and  $k$ -hop graph respectively).  $\mathbf{X} = [\mathbf{x}_1^\top, \mathbf{x}_2^\top, \dots, \mathbf{x}_N^\top] \in \mathbb{R}^{N \times d}$  is the node feature matrix where  $N$  is the number of nodes,  $d$  is the number of features, and  $\mathbf{x}_i$  denotes the node features of the node  $u_i$ .

**Original Graph Representation Learning.** We adopt the Graph Convolutional Layer (GCL) to perform message passing on the original graph  $\mathcal{G}_{\mathcal{O}} = (\mathbf{A}_{\mathcal{O}}, \mathbf{X})$  where  $\mathbf{A}_{\mathcal{O}}$  denotes the adjacency matrix of the original graph. It utilizes the original graph structure of  $\mathcal{G}_{\mathcal{O}}$  with its node feature matrix  $\mathbf{X}$  through the graph convolution operation and a multi-layer perceptron (MLP). Specifically, the designed graph convolution operation proceeds by multiplying the input of each layer with the  $\tau$ -th power of the normalized adjacency matrix. The  $\tau$ -th power operator contains statistics from the  $\tau$ -th step of a random walk on the graph (in this study, we set  $\tau$  to be 2), thus nodes can indirectly receive more information from farther nodes in the graph. Combined with a multi-layer perceptron (MLP), the representation learned at the  $\ell$ -th layer is given by:

$$\mathbf{z}_{\mathcal{G}_{\mathcal{O}}}^{(\ell+1)} = f_{\text{MLP}}(\sigma(\hat{\mathbf{A}}_{\mathcal{O}}^\tau \mathbf{H}_{\mathcal{G}_{\mathcal{O}}}^{(\ell)} \mathbf{W}^{(\ell)})), \quad (1)$$

where  $\hat{\mathbf{A}}_{\mathcal{O}} = \tilde{\mathbf{D}}_{\mathcal{O}}^{-\frac{1}{2}} \tilde{\mathbf{A}}_{\mathcal{O}} \tilde{\mathbf{D}}_{\mathcal{O}}^{\frac{1}{2}}$ ,  $\tilde{\mathbf{A}}_{\mathcal{O}} = \mathbf{A}_{\mathcal{O}} + \mathbf{I}$ , and  $\tilde{\mathbf{D}}$  is the corresponding degree matrix of  $\tilde{\mathbf{A}}$ ,  $\mathbf{H}_{\mathcal{G}_{\mathcal{O}}}^{(0)} = \mathbf{X}$ ,  $f_{\text{MLP}}$  is an MLP which has 2 layers with batch normalization,  $\sigma(\cdot)$  is the non-linear activation function,  $\mathbf{W}^{(\ell)}$  is a trainable weight of  $\ell$ -th layer.

**$\mathcal{K}$ -Nearest Neighbor Graph Representation Learning.** First, in order to capture graph structural information of nodes in topology and feature spaces, we build a  $\mathcal{K}$ -nearest neighbor ( $\mathcal{K}$ NN) graph, i.e.,

$\mathcal{G}_{\mathcal{K}} = (\mathbf{A}_{\mathcal{K}}, \mathbf{X})$ . In our study, we first define the similarity matrix  $\mathbf{S}_{\mathcal{K}} \in \mathbb{R}^{N \times N}$  among  $N$  nodes and we consider three different methods as follows: (i) *Cosine Similarity*: It uses the cosine value of the angle between two vectors to measure the similarity, i.e.,  $\mathbf{S}_{uv} = \frac{\mathbf{x}_u \cdot \mathbf{x}_v}{\|\mathbf{x}_u\| \|\mathbf{x}_v\|}$ ; (ii) *Gaussian Kernel*: It is based on the idea of the heat equation, a partial differential equation that describes how heat propagates over time  $t$ , which can be defined as follows  $\mathbf{S}_{uv} = \exp(-\|\mathbf{x}_u - \mathbf{x}_v\|^2/t)$ ; and (iii) *Node Embedding Similarity*: Let  $\mathbf{H}^{(\ell+1)}$  be the node embedding of  $(\ell)$ -th layer of GNN. For any  $u, v \in \mathcal{V}$ , we can calculate the similarity score  $\mathbf{S}_{uv}$  between nodes  $u$  and  $v$  as (i) Cosine Similarity:  $\mathbf{S}_{uv} = \frac{\mathbf{H}_u^{(\ell+1)} \cdot \mathbf{H}_v^{(\ell+1)}}{\|\mathbf{H}_u^{(\ell+1)}\| \|\mathbf{H}_v^{(\ell+1)}\|}$  or (ii) Gaussian Kernel:  $\mathbf{S}_{uv} = \exp(-\|\mathbf{H}_u^{(\ell+1)} - \mathbf{H}_v^{(\ell+1)}\|^2/t)$  (where  $t$  is a free parameter). Then, the adjacency matrix  $\mathbf{A}_{\mathcal{K}}$  can be obtained by selecting top- $\mathcal{K}$  similar neighboring nodes of each node. Similarly, we can use Eq. 1 to learn the  $(\ell + 1)$ -th layer node embeddings of the above  $\mathcal{K}$ NN graph, which is denoted by  $\mathcal{Z}_{\mathcal{G}_{\mathcal{K}}}^{(\ell+1)}$ .

**Mixup for Graph Construction.** Here we adopt the node-level attention mechanism to learn the hidden connectivity between nodes. Specifically, given a node pair  $(u, v)$ , the importance coefficient between nodes  $u$  and  $v$  can be formulated as (for the simplicity, we omit  $(\ell + 1)$  for  $\mathcal{Z}_{\mathcal{G}_{\mathcal{O}}}^{(\ell+1)}$  and  $\mathcal{Z}_{\mathcal{G}_{\mathcal{K}}}^{(\ell+1)}$ ):

$$\begin{aligned} \mathbf{e}_{uv}^{\mathcal{M}} &= \mathbf{W}_{\mathcal{M}}[\mathcal{Z}_{\mathcal{G}_{\mathcal{O}}}, \mathcal{Z}_{\mathcal{G}_{\mathcal{K}}}], \\ \alpha_{e_{uv}^{\mathcal{M}}} &= \text{Softmax}(\mathbf{e}_{uv}^{\mathcal{M}}) = \frac{\exp(\sigma(\mathbf{W}'_{\mathcal{M}} \mathbf{e}_{uv}^{\mathcal{M}}))}{\sum_{v' \in \mathcal{V}} \exp(\sigma(\mathbf{W}'_{\mathcal{M}} \mathbf{e}_{uv'}^{\mathcal{M}}))}, \end{aligned}$$

where  $[\cdot, \cdot]$  represents the concatenation operation,  $\mathbf{W}_{\mathcal{M}}$  and  $\mathbf{W}'_{\mathcal{M}}$  are trainable parameters,  $\sigma(\cdot)$  denotes the LeakyReLU function with negative input slope as 0.1. After the above calculation, we can get the mixup attention score  $\alpha_{e_{uv}^{\mathcal{M}}}$  which represents the weight of the edge between nodes  $u$  and  $v$ .

### 3.2 Stochastic Block Models for Undirected Graph

We consider two classes of probabilistic models for generating undirected networks with communities. The first is the SBM Holland et al. (1983). The second is the DCSBM Karrer & Newman (2011). Suppose the  $N$  nodes are assigned to  $K$  non-overlapping communities. The  $k$ -th community has  $N_k$  Nodes with  $\sum_{k=1}^K N_k = N$ , and denote  $\pi_k := N_k/N$ . Let  $g_i \in \{1, \dots, K\}$  be the community assignment (i.e., cluster) of node  $i$ . Alternatively, the community assignments can be represented by a membership matrix  $\mathbf{Z} \in \{0, 1\}^{N \times K}$ , where  $\mathbf{Z}_{ik} = 1$  if and only if  $g_i = k$ , and  $\mathbf{Z}_{ij} = 0$  otherwise. Let  $\mathbf{B} \in \mathbb{R}^{K \times K}$  be the block probability matrix.

Given  $\mathbf{B}$  and  $\mathbf{Z}$ , the SBM assume that each entry  $\mathbf{A}_{ij}$  ( $i < j$ ) of  $\mathbf{A}$  is generated independently by  $\mathbf{A}_{ij} \sim \text{Bernoulli}(\mathbf{B}_{g_i g_j})$ . In SBMs, the nodes within each community are stochastic equivalent. To incorporate the node heterogeneity, the more general model DCSBM is considered as follows. Let  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N)^{\top} \in \mathbb{R}^N$  be the node propensity parameters and denote  $\boldsymbol{\Theta} = \text{diag}\{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N\}$ . Given  $\mathbf{B}$ ,  $\mathbf{Z}$  and  $\boldsymbol{\Theta}$ , the DCSBM assume each entry  $\mathbf{A}_{ij}$  ( $i < j$ ) of  $\mathbf{A}$  is generated independently by  $\mathbf{A}_{ij} \sim \text{Bernoulli}(\boldsymbol{\theta}_i \boldsymbol{\theta}_j \mathbf{B}_{g_i g_j})$ . It is then easy to see that

$$\mathbf{P} := \boldsymbol{\Theta} \mathbf{Z} \mathbf{B} \mathbf{Z}^{\top} \boldsymbol{\Theta} \in \mathbb{R}^{N \times N} \quad (2)$$

is the population counterpart of  $\mathbf{A}$ .  $\mathbf{P}$  is referred to as the edge probability matrix. Note that  $\boldsymbol{\Theta}$  and  $\mathbf{B}$  are only identifiable up to scaling. As a remedy, we use the following normalization rule

$$\sum_{i, g_i=k} \boldsymbol{\theta}_i = N_k, \quad k = 1, \dots, K.$$

With this normalization rule, the SBM is nested by the DCSBM by letting  $\boldsymbol{\theta}_i = 1$  for  $i = 1, \dots, N$ . To estimate  $\mathbf{P}$ , we should estimate  $\mathbf{Z}$ ,  $\mathbf{B}$  and  $\boldsymbol{\Theta}$ , respectively. Before that, we first recall and introduce some notation. Let  $\hat{d}_i = \sum_{j=1}^N \mathbf{A}_{ij}$  be the degree of node  $i$  and  $\mathbf{D} = \text{diag}\{\hat{d}_1, \dots, \hat{d}_N\}$ . The graph Laplacian is defined as  $\mathbf{L} = \mathbf{I} + \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$ . Note that this graph Laplacian matrix has the same eigenspaces as that of  $\mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$ .

Under SBMs, we estimate  $\mathbf{Z}$  using the standard spectral clustering on the graph Laplacian matrix  $\mathbf{L}$ . That is, conducting  $k$ -means on the top- $K$  eigenvectors of  $\mathbf{L}$ . The estimator is denoted by  $\hat{\mathbf{Z}}$ . WLOG, we assume that

$\hat{\mathbf{Z}}$  has been orthogonally transformed to align with  $\mathbf{Z}$ . We estimate  $\mathbf{B}$  by the following  $\hat{\mathbf{B}} = (\hat{\mathbf{B}}_{ql})_{1 \leq q, l \leq K}$ ,

$$\hat{\mathbf{B}}_{ql, q \neq l} := \frac{\sum_{\hat{g}_i=q, \hat{g}_j=l} \mathbf{A}_{ij}}{\hat{N}_q \hat{N}_l} \quad \text{and} \quad \hat{\mathbf{B}}_{qq} := \frac{\sum_{\hat{g}_i=q, \hat{g}_j=q} \mathbf{A}_{ij}}{\hat{N}_q (\hat{N}_q - 1)}.$$

Thereby, we obtain  $\hat{\mathbf{P}} := \hat{\mathbf{Z}} \hat{\mathbf{B}} \hat{\mathbf{Z}}^\top$ . Under the DC SBMs, we estimate  $\mathbf{Z}$  using the spherical spectral clustering on the graph Laplacian matrix  $\mathbf{L}$ . That is, conducting  $k$ -means on the  $L_2$ -row-normalized top- $K$  eigenvectors of  $\mathbf{L}$ . With a light abuse of notation, the estimator is also denoted by  $\hat{\mathbf{Z}}$ . We estimate  $\mathbf{B}$  by  $\hat{\mathbf{B}} = (\hat{\mathbf{B}}_{ql})_{1 \leq q, l \leq K}$ ,

$$\hat{\mathbf{B}}_{ql} := \frac{\sum_{1 \leq i \neq j \leq N} \mathbf{A}_{ij} \hat{\mathbf{Z}}_{iq} \hat{\mathbf{Z}}_{jl}}{\sum_{1 \leq i \neq j \leq N} \hat{\mathbf{Z}}_{iq} \hat{\mathbf{Z}}_{jl}} = \frac{\sum_{\hat{g}_i=q, \hat{g}_j=l} \mathbf{A}_{ij}}{\hat{N}_q \hat{N}_l}.$$

We estimate  $\theta_i$  by  $\hat{\theta}_i$  defined as

$$\hat{\theta}_i = \frac{\hat{N}_{\hat{g}_i} \sum_j \mathbf{A}_{ij}}{\sum_{\hat{g}_l=\hat{g}_i} \sum_{j=1}^N \mathbf{A}_{lj}},$$

where  $\hat{N}_{\hat{g}_i}$  is the number of nodes in the estimated community  $\hat{g}_i$ . We also denote  $\hat{\Theta} = \text{diag}(\hat{\theta})$ . Finally, we obtain  $\hat{\mathbf{P}} := \hat{\Theta} \hat{\mathbf{Z}} \hat{\mathbf{B}} \hat{\mathbf{Z}}^\top \hat{\Theta}$ . The detailed derivations of the estimators can be found in the Appendix.

### 3.3 Multi-View Topological Graph Neural Networks for Undirected Graph

**Multi-View Topological Convolutional Layer.** To capture the underlying topological features of the subgraph  $\mathcal{G}_u$  of each node  $u$ , we employ  $\mathcal{K}$  filtration functions:  $f_i : \mathcal{V} \mapsto \mathbb{R}$  for  $i = \{1, \dots, \mathcal{K}\}$ . Each filtration function  $f_i$  gradually reveals one specific topological structure at different levels of connectivity, e.g., degree centrality score, betweenness centrality score, closeness centrality score, and other node centrality measurements. With each filtration function  $f_i$ , we construct a set of  $Q$  persistence images of resolution  $P \times P$  using tools in persistent homology analysis. Combining  $Q$  persistence images of resolution  $P \times P$  from  $\mathcal{K}$  different filtration functions, we construct a *multi-view* topological representation, i.e., the set of persistence images (PIs)  $[\text{PI}_1, \text{PI}_2, \dots, \text{PI}_{\mathcal{K}}]$  with the dimension  $\mathcal{K} \times Q \times P \times P$ . We design a multi-view topological convolutional layer  $f_{\text{MV-GCL}}$  to (i) jointly extract and learn the latent topological features and (ii) leverage and preserve the multi-modal structure. Firstly, hidden representations of the set of PIs are achieved through a combination of a CNN-based model and global pooling, which can be defined as

$$\mathcal{Z}_{u, \mathcal{T}} = \xi_{\text{POOL}}(f_{\text{CNN}}([\text{PI}_1, \text{PI}_2, \dots, \text{PI}_{\mathcal{K}}])), \quad (3)$$

where  $f_{\text{CNN}}$  is a CNN-based neural network,  $\xi_{\text{POOL}}$  is a pooling layer that preserves the information of the input in a fixed-size representation (in general, we consider either global average pooling or global max pooling). In Eq. equation 3, we first apply a CNN-based model to learn the latent feature of PIs, and then employ a global pooling layer over the latent feature and obtain an image-level feature.

**Graph Convolutional Layer.** Our third representation learning module is the Graph Convolutional Layer (GCL). It utilizes the graph structure of  $\mathcal{G}_{\mathcal{M}}$  with its node feature matrix  $\mathbf{X}$  through the graph convolution operation and a multi-layer perceptron (MLP). The representation learned at the  $\ell$ -th layer is given by

$$\mathcal{Z}_{\mathcal{G}}^{(\ell+1)} = f_{\text{MLP}}([\sigma(\mathbf{A}_{\mathcal{M}} \mathbf{H}_{\mathcal{M}}^{(\ell)} \mathbf{W}_{\mathcal{M}}^{(\ell)}), \sigma(\hat{\mathbf{P}} \mathbf{H}_{\text{SBM}}^{(\ell)} \mathbf{W}_{\text{SBM}}^{(\ell)})]),$$

where  $\mathbf{H}_{\mathcal{M}}^{(0)} = \mathbf{H}_{\text{SBM}}^{(0)} = \mathbf{X}$ ,  $f_{\text{MLP}}$  is an MLP which has 2 layers with batch normalization,  $\sigma(\cdot)$  is the non-linear activation function,  $\mathbf{W}_{\mathcal{M}}^{(\ell)}$  and  $\mathbf{W}_{\text{SBM}}^{(\ell)}$  are trainable weight matrices of  $\ell$ -th layer. Then, we obtain the final embedding  $\mathcal{Z}$  by combining embeddings from the above modules, i.e.,  $\mathcal{Z} = [\mathcal{Z}_{\mathcal{T}}, \mathcal{Z}_{\mathcal{G}}]$ , where  $[\cdot, \cdot]$  denotes the concatenation operation and  $\mathcal{Z}_{\mathcal{G}}$  represents the final output of the graph convolutional layer.

## 4 Stochastic Block Model-Aware Topological Neural Networks: Knowledge Graph

A knowledge graph (KG) is defined as a directed graph that stores structured information about real-world entities and relations. Let  $\mathcal{G} = \{\mathcal{V}, \mathcal{R}, \mathcal{L}\}$  be an instance of a KG, where  $\mathcal{V}$ ,  $\mathcal{R}$ , and  $\mathcal{L}$  denote the entity (i.e.,

node), relation, and edge sets respectively. Each edge  $e \in \mathcal{L}$  is presented as a triple  $(h, r, t) \in \mathcal{V} \times \mathcal{R} \times \mathcal{V}$ , describing that there is a relationship  $r \in \mathcal{R}$  from head entity  $h$  to tail entity  $t$ . In order to apply the stochastic co-block model (which is introduced in Section 4.1) to classify nodes into  $K$  clusters, we first transform the KG into a directed graph and then generate community information based on the graph's adjacency matrix instead of the knowledge graph by using SBM method. Note that, for the KG, we incorporate its information into the model and we don't have an additional assumption for KG.

#### 4.1 Stochastic Co-Block Models for Directed Graph

Similar to the undirected networks, we consider two classes of probabilistic models for generating directed networks with co-clusters, namely, row clusters (communities) and column clusters (communities). The first is the stochastic co-block model (ScBM), and the second is the degree-corrected stochastic co-block model (DCScBM) Rohe et al. (2016).

Different from SBMs and DCScBMs, the models for directed networks incorporate two kinds of clusters. Suppose the  $N$  nodes are assigned to  $K^y$  non-overlapping row clusters and  $K^z$  non-overlapping column clusters. WLOG, suppose  $K^y \leq K^z$ . The  $k^y$ th (resp.  $k^z$ th) row (resp. column) cluster has  $N_k^y$  (resp.  $N_k^z$ ) nodes with  $\sum_{k=1}^{K^y} N_k^y = N$  (resp.  $\sum_{k=1}^{K^z} N_k^z = N$ ), and denote  $\pi_k^y := N_k^y/N$  (resp.  $\pi_k^z := N_k^z/N$ ). Let  $g_i^y \in \{1, \dots, K^y\}$  (resp.  $g_i^z \in \{1, \dots, K^z\}$ ) be the row (resp. column) community assignment of node  $i$ . Following the same logic as undirected networks, the community assignments can be represented by a row membership matrix  $\mathbf{Y} \in \{0, 1\}^{N \times K^y}$  and a column membership matrix  $\mathbf{Z} \in \{0, 1\}^{N \times K^z}$ . Let  $\mathbf{B} \in \mathbb{R}^{K^y \times K^z}$  be the block probability matrix.

Given  $\mathbf{B}$ ,  $\mathbf{Y}$  and  $\mathbf{Z}$ , the ScBM assume that each entry  $\mathbf{A}_{ij}$  ( $i \neq j$ ) of  $\mathbf{A}$  is generated independently by  $\mathbf{A}_{ij} \sim \text{Bernoulli}(\mathbf{B}_{g_i^y g_j^z})$ . In ScBMs, the nodes in a common row (resp. column) cluster are stochastically equivalent senders (resp. receivers) in the sense that they send (resp. receive) out an edge to a third node with equal probabilities. To incorporate the node heterogeneity in sending and receiving edges, the more general model DCScBM is considered as follows. Let  $\boldsymbol{\theta}^y = (\theta_1^y, \dots, \theta_N^y)^\top \in \mathbb{R}^N$  be the node propensity parameters in sending edges and denote  $\boldsymbol{\Theta}^y = \text{diag}(\boldsymbol{\theta}^y)$ . Let  $\boldsymbol{\theta}^z = (\theta_1^z, \dots, \theta_N^z)^\top \in \mathbb{R}^N$  be the node propensity parameters in receiving edges and denote  $\boldsymbol{\Theta}^z = \text{diag}(\boldsymbol{\theta}^z)$ . Given  $\mathbf{B}$ ,  $\mathbf{Y}$ ,  $\mathbf{Z}$ ,  $\boldsymbol{\Theta}^y$  and  $\boldsymbol{\Theta}^z$ , the DCScBM assumes each entry  $\mathbf{A}_{ij}$  ( $i \neq j$ ) of  $\mathbf{A}$  is generated independently by  $\mathbf{A}_{ij} \sim \text{Bernoulli}(\theta_i^y \theta_j^z \mathbf{B}_{g_i^y g_j^z})$ . Then, under the DCScBM, it is easy to see

$$\mathbf{P} := \boldsymbol{\Theta}^y \mathbf{Y} \mathbf{B} \mathbf{Z}^\top \boldsymbol{\Theta}^z \in \mathbb{R}^{N \times N} \quad (4)$$

is the population adjacency matrix of  $\mathbf{A}$ , referred to as the edge probability matrix later. To ensure identifiability, we use the following normalization rule

$$\begin{aligned} \sum_{i, g_i^y = k^y} \theta_i^y &= N_k^y, \quad k^y = 1, \dots, K^y, \\ \sum_{i, g_i^z = k^z} \theta_i^z &= N_k^z, \quad k^z = 1, \dots, K^z. \end{aligned}$$

The ScBM is nested by the DCScBM by letting  $\theta_i^y = 1$  and  $\theta_i^z = 1$  for  $i = 1, \dots, N$ . To estimate  $\mathbf{P}$ , we introduce some notation now. Let  $\hat{d}_i^y = \sum_{j=1}^N \mathbf{A}_{ij}$  (resp.  $\hat{d}_i^z = \sum_{j=1}^N \mathbf{A}_{ji}$ ) be the out-degree (resp. in-degree) of node  $i$  and  $\mathbf{D}^y = \text{diag}\{\hat{d}_1^y, \dots, \hat{d}_N^y\}$  (resp.  $\mathbf{D}^z = \text{diag}\{\hat{d}_1^z, \dots, \hat{d}_N^z\}$ ). Define the graph Laplacian by  $\mathbf{L} = \mathbf{I} + (\mathbf{D}^y)^{-1/2} \mathbf{A} (\mathbf{D}^z)^{-1/2}$ . Under ScBMs, we compute the SVD of  $\mathbf{L}$  and then conduct the  $k$ -means on the top- $K^y$  left (resp. top- $K^z$  right) singular vectors of  $\mathbf{L}$ , to obtain  $K^y$  row (resp.  $K^z$  column) clusters, denoted by  $\hat{\mathbf{Y}}$  (resp.  $\hat{\mathbf{Z}}$ ). WLOG, we assume that  $\hat{\mathbf{Y}}$  ( $\hat{\mathbf{Z}}$ ) has been orthogonally transformed to align with  $\mathbf{Y}$  (resp.  $\mathbf{Z}$ ). Similar to the undirected set-up, we estimate  $\mathbf{B}$  by the following  $\hat{\mathbf{B}} = (\hat{\mathbf{B}}_{ql})_{1 \leq q \leq K^y, 1 \leq l \leq K^z}$ ,

$$\hat{\mathbf{B}}_{ql} := \frac{\sum_{1 \leq i \neq j \leq N} \mathbf{A}_{ij} \hat{\mathbf{Y}}_{iq} \hat{\mathbf{Z}}_{jl}}{\sum_{1 \leq i \neq j \leq N} \hat{\mathbf{Y}}_{iq} \hat{\mathbf{Z}}_{jl}} = \frac{\sum_{\hat{g}_i^y = q, \hat{g}_j^z = l} \mathbf{A}_{ij}}{\hat{N}_q^y \hat{N}_l^z},$$

where  $\hat{N}_q^y$  (resp.  $\hat{N}_l^z$ ) denotes the number of nodes in the estimated row cluster  $q$  (resp. column cluster  $l$ ). Thereby, we obtain  $\hat{\mathbf{P}} := \hat{\mathbf{Y}} \hat{\mathbf{B}} \hat{\mathbf{Z}}^\top$ . Under the DCScBMs, we estimate  $\mathbf{Z}$  by applying the spherical spectral

clustering to the graph Laplacian matrix  $\mathbf{L}$ . That is, we obtain the top- $K^y$  row and column singular vectors by computing the SVD of  $\mathbf{L}$ , and then conduct  $k$ -means on the  $L_2$ -row-normalized left and right singular vectors to obtain the  $K^y$  row clusters and  $K^z$  column clusters, respectively. With a light abuse of notation, the estimators are denoted by  $\hat{\mathbf{Y}}$  and  $\hat{\mathbf{Z}}$ . Following the same logic as in DCSBM, we obtain the estimators  $\hat{\mathbf{B}} = (\hat{\mathbf{B}}_{ql})$ ,  $\hat{\Theta}^y = \text{diag}(\hat{\theta}^y)$  and  $\hat{\Theta}^z = \text{diag}(\hat{\theta}^z)$  as follows.

$$\hat{\mathbf{B}}_{ql} := \frac{\sum_{\hat{g}_i^y=q, \hat{g}_j^z=l} \mathbf{A}_{ij}}{\hat{N}_q^y \hat{N}_l^z}, \quad \hat{\theta}_i^y = \frac{\hat{N}_{\hat{g}_i^y} \sum_j \mathbf{A}_{ij}}{\sum_{\hat{g}_i^y=\hat{g}_i^y} \sum_{j=1}^N \mathbf{A}_{lj}},$$

$$\hat{\theta}_j^z = \frac{\hat{N}_{\hat{g}_j^z} \sum_i \mathbf{A}_{ij}}{\sum_{\hat{g}_i^z=\hat{g}_j^z} \sum_{i=1}^N \mathbf{A}_{il}},$$

where  $\hat{N}_{\hat{g}_i^y}$  is the number of nodes in the estimated row cluster  $\hat{g}_i^y$  and  $\hat{N}_{\hat{g}_j^z}$  is the number of nodes in the estimated column cluster  $\hat{g}_j^z$ . With these estimators at hand, we obtain  $\hat{\mathbf{P}} := \hat{\Theta}^y \hat{\mathbf{Y}} \hat{\mathbf{B}} \hat{\mathbf{Z}}^\top \hat{\Theta}^z$ .

## 4.2 Attention-Based Topological Neural Networks for Knowledge Graph

In this section, we propose a novel attention-based topological neural network (A-TNN) to combine the community information and topological signatures learned from ScBM for entity representation learning.

$$\mathbf{z}_h^{(\ell+1)} = \sigma\left(\sum_{(r,t) \in \mathcal{N}_h} \alpha_{h,r,t}^{(\ell)} \tilde{\mathbf{X}}_{h,r,t}^{(\ell)}\right),$$

$$\tilde{\mathbf{X}}_{h,r,t}^{(\ell)} = \mathbf{W}_{\text{KG}_1} [\mathbf{X}_h^{(\ell)}, \mathbf{X}_t^{(\ell)}, \mathbf{X}_{C_k}^{(\ell)}, \mathbf{X}_h^{\text{topo}}, \mathbf{X}_r],$$

$$\alpha_{h,r,t}^{(\ell)} = \frac{\exp(\sigma(\mathbf{W}_{\text{KG}_2}^{(\ell)} \tilde{\mathbf{X}}_{h,r,t}^{(\ell)}))}{\sum_{(r,t') \in \mathcal{N}_h} \exp(\sigma(\mathbf{W}_{\text{KG}_2}^{(\ell)} \tilde{\mathbf{X}}_{h,r,t'}^{(\ell)}))},$$
(5)

where  $\mathbf{X}_h^{(\ell)}$ ,  $\mathbf{X}_t^{(\ell)}$ , and  $\mathbf{X}_r^{(\ell)}$  denote embeddings of head entity  $h$ , tail entity  $t$ , and relation  $r$  respectively,  $\mathbf{W}_{\text{KG}_1}$  and  $\mathbf{W}_{\text{KG}_2}$  denote the linear transformation matrices,  $\mathbf{X}_{C_k}^{(\ell)}$ ,  $\mathcal{N}_h$  denotes the set of neighboring tuples  $(r, t)$  for entity  $h$ ,  $\mathbf{X}_k^{(\ell)} = \sum_{h \in C_k} \mathbf{X}_h^{(\ell)}$  which aggregates node embeddings from the community  $C_k$  (i.e., results from Section 4.1), and  $\mathbf{X}_h^{\text{topo}} = f_{\text{MV-GCL}}(\mathcal{G}_h)$  (where  $\mathcal{G}_h$  denotes the subgraph of the node  $h$ ).

## 5 Theoretical Guarantees

To establish the statistical property of  $\hat{\mathbf{P}}$ , we need the following assumptions. In order to provide the results uniformly and reduce redundancy, we use the notations in the more general ScBMs and DCScBMs. The SBMs and DCScBMs can be regarded as the special case of ScBMs and DCScBMs.

**Assumption 5.1.** *Suppose  $K^y$  and  $K^z$  are both fixed and both the row and column clusters are balanced such that for any  $k^y \in \{1, \dots, K^y\}$ ,  $cN/K^y \leq N_k^y \leq CN/K^y$  for some constants  $0 \leq c \leq C$ ; for any  $k^z \in \{1, \dots, K^z\}$ ,  $c'N/K^z \leq N_k^z \leq C'N/K^z$  for some constants  $0 \leq c' \leq C'$ .*

**Assumption 5.2.** *Suppose the edge probability matrix  $\mathbf{B} \in [0, 1]^{K^y \times K^z}$  ( $K^y \leq K^z$ ) is of rank  $K^z$ . The entries of  $\mathbf{B}$  are of the same magnitude  $\rho_N$  with  $N\rho_N \geq c \log N$ .*

**Assumption 5.3.** *Define  $\sigma_{K^y}$  be the  $K^y$ th singular value of  $\mathcal{L} = (\mathcal{D}^y)^{-1/2} \mathbf{P} (\mathcal{D}^z)^{-1/2}$ , where  $\mathcal{D}^y = \text{diag}\{d_1^y, \dots, d_N^y\}$  with  $d_i^y = \sum_{j=1}^N \mathbf{P}_{ij}$ , and  $\mathcal{D}^z = \text{diag}\{d_1^z, \dots, d_N^z\}$  with  $d_i^z = \sum_{j=1}^N \mathbf{P}_{ji}$ . Suppose  $\liminf_N |\sigma_{K^y}| > 0$ .*

**Assumption 5.4.** *Define  $\bar{\theta}^y = \max_i \theta_i^y$  and  $\underline{\theta}^y = \min_i \theta_i^y$ . Suppose  $c \leq \liminf_N \underline{\theta}^y \leq \limsup_N \bar{\theta}^y \leq C$  for some constants  $0 \leq c \leq C$ . Similarly, define  $\bar{\theta}^z = \max_i \theta_i^z$  and  $\underline{\theta}^z = \min_i \theta_i^z$ . Suppose  $c' \leq \liminf_N \underline{\theta}^z \leq \limsup_N \bar{\theta}^z \leq C'$  for some constants  $0 \leq c' \leq C'$ .*

**Assumption 5.5.** *Define  $\mathbf{H} = (\mathbf{Y}^\top \Theta^y \mathbf{Y})^{1/2} \mathbf{B}_L (\mathbf{Z}^\top \Theta^z \mathbf{Z})^{1/2} \in \mathbb{R}^{K^y \times K^z}$  with  $\mathbf{B}_L := \mathbf{O}_B^{-1/2} \mathbf{B} \mathbf{P}_B^{-1/2}$ , where  $\mathbf{O}_B$  is a  $K^y \times K^y$  diagonal matrix with  $[\mathbf{O}_B]_{ss} = \sum_t \mathbf{B}_{st} n_t^z$  and  $\mathbf{P}_B$  is a  $K^z \times K^z$  diagonal matrix with*

Table 1: AUC-ROC score on different link prediction benchmarks.

Model	Cora-ML	Citeseer	PubMed	Photo	Computers
GCN Kipf & Welling (2017)	90.5±0.2	82.6±1.9	89.6±3.7	91.8±0.0	87.8±0.0
GAT Veličković et al. (2018b)	72.8±0.2	74.8±1.5	80.3±0.0	92.9±0.3	86.4±0.0
HGCN Chami et al. (2019)	93.8±0.1	96.6±0.1	96.3±0.0	95.4±0.0	93.6±0.0
P-GNN You et al. (2019)	74.1±2.4	73.9±2.6	79.6±0.5	90.9±0.7	88.3±1.0
SEAL Zhang & Chen (2018)	91.3±5.7	89.8±2.3	92.4±1.2	97.8±1.3	96.8±1.5
BScNets Chen et al. (2022)	94.9±0.7	95.5±0.5	97.6±0.1	96.6±0.3	97.0±0.3
TLC-GNN Yan et al. (2021)	94.9±0.4	95.1±0.7	97.0±0.1	98.2±0.1	97.9±0.1
<b>SBM-TNN (ours)</b>	<b>96.2±0.2</b>	<b>97.1±0.3</b>	<b>98.2±0.1</b>	<b>98.8±0.2</b>	<b>99.0±0.2</b>

$[\mathbf{P}_B]_{tt} = \sum_s \mathbf{B}_{st} n_s^y$ . Suppose there exists gap between any two columns of  $\mathbf{H}$  such that  $\min_{i \neq j} \|\mathbf{H}_{\cdot i} - \mathbf{H}_{\cdot j}\|_2 \geq \xi$  for some constant  $\xi > 0$ .

Assumptions 5.1-5.4 are generally required for undirected network models SBMs and DCSBMs. For directed network models ScBMs and DCScBMs, we need additional Assumption 5.5. Assumption 5.1 implies that the communities are balanced and it is a common assumption in the SBM literature. Assumption 5.2 requires that the network is not too sparse. This condition is the minimal requirement for strong consistency of SBM Abbe et al. (2015); Su et al. (2019); Ma et al. (2021). Assumption 5.3 implies that the singular value of the population Laplacian matrix is lower bounded by a constant. For an SBM with cross-block probability being  $r$ , and within-block probability being  $r + p$ , then  $\sigma_K = p/(Kr + p)$  is a constant provided that  $r$  and  $p$  are of the same order. Assumption 5.4 requires that the node propensity parameters are upper and lower bounded. More generally, the lower bound can be relaxed to  $n^{-\alpha}$  for some positive constant  $\alpha$  with a sacrifice of the simplicity of other conditions Su et al. (2019). This assumption is only needed when networks follow DCScBMs and DCSBMs. Assumption 5.5 is required to ensure that there exists a gap between two rows of the population right singular vectors when the two nodes are in different column clusters. This is a remedy assumption for the invalidity of the full column rank under directed network models.

Table 2: KGC results (%) with different scoring functions.

	FB15k-237				WN18RR				NELL-995			
	MRR	Hits@1	Hits@3	Hits@10	MRR	Hits@1	Hits@3	Hits@10	MRR	Hits@1	Hits@3	Hits@10
CompGCN Vashishth et al. (2019)	35.5±0.1	26.4±0.1	39.0±0.2	53.6±0.3	47.2±0.2	43.7±0.3	48.5±0.3	54.0±0.0	38.1±0.4	30.4±0.5	42.2±0.3	52.9±0.1
RGCN Schlichtkrull et al. (2018)	29.6±0.3	19.1±0.5	34.0±0.2	50.1±0.2	43.0±0.2	38.6±0.3	45.0±0.1	50.8±0.3	27.8±0.2	19.9±0.2	31.4±0.0	43.0±0.3
KBGAT Nathani et al. (2019)	35.0±0.3	26.0±0.3	38.5±0.3	53.1±0.3	46.4±0.2	42.6±0.2	47.9±0.3	53.9±0.2	37.4±0.6	29.7±0.7	41.4±0.8	52.0±0.4
AcrE Ren et al. (2020)	35.8±0.3	26.6±0.2	39.3±0.3	54.5±0.2	45.9±0.2	42.2±0.3	47.3±0.2	53.2±0.1	-	-	-	-
RelInceptionE Xie et al. (2020)	34.9±0.2	-	-	52.8±0.2	48.3±0.3	-	-	58.2±0.3	-	-	-	-
SE-GNN Li et al. (2022)	36.1±0.3	23.4±0.2	37.0±0.3	51.5±0.2	48.4±0.4	43.6±0.2	47.9±0.2	57.2±0.3	39.3±0.5	30.2±0.3	43.0±0.2	52.8±0.3
GNNs + NBFNet Zhu et al. (2021)	41.5±0.1	32.1±0.1	45.6±0.4	59.9±0.3	55.1±0.1	49.7±0.1	57.2±0.3	66.6±0.3	40.5±0.3	32.7±0.2	44.9±0.3	55.0±0.4
<b>SBM-TNN (ours)</b>	36.9±0.2	27.3±0.3	40.8±0.2	55.6±0.1	49.1±0.2	44.1±0.2	48.9±0.1	59.3±0.2	41.7±0.2	34.7±0.2	45.2±0.1	55.2±0.1
<b>SBM-TNN + NBFNet (ours)</b>	<b>41.9±0.1</b>	<b>33.0±0.1</b>	<b>47.1±0.2</b>	<b>60.1±0.2</b>	<b>55.5±0.1</b>	<b>49.9±0.1</b>	<b>59.2±0.1</b>	<b>67.0±0.0</b>	<b>42.0±0.1</b>	<b>34.9±0.2</b>	<b>45.5±0.3</b>	<b>55.9±0.1</b>

**Lemma 5.6** (Strong consistency). Suppose Assumptions 5.1-5.5 hold. Then for large enough  $N$ , it holds for both ScBMs and DCScBMs that

$$\sup_{1 \leq i \leq N} \mathbf{1}\{\hat{g}_i^y \neq g_i^y\} = 0 \text{ and } \sup_{1 \leq i \leq N} \mathbf{1}\{\hat{g}_i^z \neq g_i^z\} = 0 \text{ a.s.} \quad (6)$$

and

$$\sup_{1 \leq q \leq K^y, 1 \leq l \leq K^z} |\hat{\mathbf{B}}_{ql} - \mathbf{B}_{ql}| = O_{a.s.}\left(\frac{\sqrt{\rho \log N}}{N}\right). \quad (7)$$

For DCScBMs, it also holds that

$$\begin{aligned} \sup_{1 \leq i \leq N} |\hat{\boldsymbol{\theta}}_i^y - \boldsymbol{\theta}_i^y| &= O_{a.s.}\left(\frac{\log N}{N\rho}\right) \\ \sup_{1 \leq i \leq N} |\hat{\boldsymbol{\theta}}_i^z - \boldsymbol{\theta}_i^z| &= O_{a.s.}\left(\frac{\log N}{N\rho}\right). \end{aligned} \quad (8)$$



**Remark 1.** Lemma 5.6 also holds for SBMs and DCSBMs with the notation slightly changed and without the requirement of Assumption 5.5. The strong consistency results are critical for deriving the error bound and asymptotic normality of  $\hat{\mathbf{P}}$ . equation 6 and equation 8 can be implied by Corollary III.1 and Theorem III.6 in Su et al. Su et al. (2019).

**Theorem 5.7** (Error bound). Suppose Assumptions 5.1-5.5 hold. Then for large enough  $N$ , it holds for ScSBMs that

$$\sup_{1 \leq i, j \leq N} |\hat{\mathbf{P}}_{ij} - \mathbf{P}_{ij}| = O_{a.s.} \left( \frac{\sqrt{\rho \log N}}{N} \right), \quad (9)$$

and for DCScSBMs that

$$\sup_{1 \leq i, j \leq N} |\hat{\mathbf{P}}_{ij} - \mathbf{P}_{ij}| = O_{a.s.} \left( \frac{\log N}{N} \right). \quad (10)$$

**Remark 2.** For SBMs and DCSBMs, equation 9 and equation 21 hold without Assumption 5.5. Most of the current literature on SBM and its variants study the high probability spectral norm bound of  $\mathbf{A}$  from  $\mathbf{P}$  Gao et al. (2015); Lei & Rinaldo (2015). Our results are novel in that we provide the uniform entrywise bound of  $\hat{\mathbf{P}}$  from  $\mathbf{P}$  which holds almost surely. The less tightness of DCSBMs and DCScSBMs comes from the estimation of the node propensity parameters.

**Theorem 5.8** (Asymptotic normality). Suppose Assumptions 5.1-5.5 hold. Then for ScSBMs, (11) holds with  $c_0 = 1, \theta_i^y = \theta_j^z = 1$  and  $\mathbf{E} = \mathbf{0}$ . For DCScSBMs, (11) holds with  $c_0 = 1$  and  $\mathbf{E}_{ij} = O_{a.s.}(\log N/N)$  (where  $1 \leq i, j \leq N$ ).

$$N(\hat{\mathbf{P}} - \mathbf{P} + \mathbf{E})_{ij} \rightarrow N(0, c_0 \cdot (\theta_i^y)^2 \cdot \frac{\mathbf{B}_{g_i^y g_j^z} (1 - \mathbf{B}_{g_i^y g_j^z})}{\pi_{g_i^y} \pi_{g_j^z}} \cdot (\theta_j^z)^2), \quad (11)$$

**Remark 3.** (11) also holds for SBMs and DCSBMs without Assumption 5.5 and with the same order of  $\mathbf{E}$  and notation slightly modified except that  $c_0 = 2$  when  $g_i = g_j$  and  $c_0 = 1$  when  $g_i \neq g_j$ . For SBMs, the estimator  $\hat{\mathbf{P}}$  is asymptotic efficient (Bickel et al., 2013). For DCSBMs and DCScSBMs, the bias  $\mathbf{E}_{ij}$  and  $\mathbf{F}_{ij}$  come from the estimation of  $\theta$ . Note that for denser networks, namely,  $\rho = \Omega(\log N/N)$ , the bias is dominated by the signal  $\mathbf{P}_{ij}$ . To the best of our knowledge, this is the first result to show the asymptotic normality of an estimator against the edge probability matrix  $\mathbf{P}$ .

## 6 Experiments

**Datasets and Baselines.** We experiment on 2 types of networks for link prediction (i) citation networks: Cora-ML, Citeseer, and PubMed (Sen et al., 2008) and (ii) graphs related to Amazon shopping records: Photo and Computers (Shchur et al., 2018). For link prediction tasks, we compare against 10 state-of-the-art (SOA) baselines, including (i) Graph convolution network (GCN) (Kipf & Welling, 2017); (ii) Graph Attention Networks (GAT) (Veličković et al., 2018b); (iii) Hyperbolic Graph Convolutional Neural Networks (HGNC) (Chami et al., 2019); (iv) Position-aware Graph Neural Networks (P-GNN) (You et al., 2019); (v) SEAL (Zhang & Chen, 2018); (vi) Block Simplicial Complex Neural Networks (BScNets) (Chen et al., 2022); (vii) Topological Loop-Counting Graph Neural Network (TLC-GNN) (Yan et al., 2021). For knowledge graph completion tasks, we conduct experiments on 3 well-known KG datasets including (i) FB15k-237 (Toutanova et al., 2015; Toutanova & Chen, 2015), WN18RR Dettmers et al. (2018), and NELL-995 (Xiong et al., 2017), and use the following popular models as baselines: (i) CompGCN (Vashishth et al., 2019); (ii) Relational Graph Convolutional Network (RGCN) (Schlichtkrull et al., 2018); (iii) KBGAT (Nathani et al., 2019); (iv) Atrous Convolution and Residual Embedding (AcrE) (Ren et al., 2020); (v) ReInceptionE (Xie et al., 2020); (vi) Semantic Evidence aware Graph Neural Network (SE-GNN) (Li et al., 2022); and (vii) Neural Bellman-Ford Network (NBFNet) (Zhu et al., 2021).

**Experiment Settings.** We implement our proposed SBM-TNN with Pytorch framework on two NVIDIA RTX A5000 GPUs with 24 GiB RAM. Following Chami et al. (2019), for graph link prediction tasks, we randomly split edges into 85%/5%/10% for training, validation, and testing, and we evaluate link prediction using the ROC-AUC score on the test set. For KG completion tasks, we follow the settings in previous works (Vashishth et al., 2019; Schlichtkrull et al., 2018), i.e., triplets in these datasets are randomly split into

Table 3: Link prediction (AUC-ROC) of SBM-TNN with different  $\tau$ .

Dataset	$\tau = 1$	$\tau = 2$	$\tau = 3$	$\tau = 4$
Cora-ML	93.7±0.8	<b>96.2±0.2</b>	95.4±0.3	93.2±0.9
Citeseer	94.9±0.8	97.0±0.3	<b>97.1±0.3</b>	95.1±0.6
PubMed	97.9±0.3	<b>98.2±0.2</b>	98.7±0.6	95.5±0.7

Table 4: Performance comparison for link prediction (AUC-ROC) between  $\mathbf{P}$  and  $\hat{\mathbf{P}}$ .

Architecture	Cora-ML	Citeseer	PubMed
SBM-TNN ( $\mathbf{P}$ )	95.3±0.3	95.0±0.5	96.3±0.3
SBM-TNN ( $\hat{\mathbf{P}}$ )	<b>96.2±0.2</b>	<b>97.1±0.3</b>	<b>98.2±0.1</b>

training, validation, and test sets respectively, and we evaluate the KG completion performance by using Mean Reciprocal Rank (MRR) and Hits@ $N$  (here we consider  $N \in \{1, 3, 10\}$ ). Our datasets and codes are available on [Yuzhouwilladd](https://github.com/yuzhouwilladd). For further details on the experiment settings, please refer to Appendix.

**Experiment Results.** The link prediction and KG completion results are summarized in Tables 1 and 2. From Table 1, the results indicate that our SBM-TNN consistently achieves the best performance on all datasets. More specifically, we find that (i) Compared to the spectral-based GNN (i.e., GCN), our SBM-TNN yields up to 12.8% relative improvements for all 5 datasets; (ii) Compared to the spatial-based GNNs (i.e., GAT, P-GNN, and SEAL), SBM-TNN improves upon the runner-up by a margin of 5.4%, 8.1%, 6.3%, 1.0%, and 2.3% on datasets Cora-ML, Citeseer, PubMed, Photo and Computers; (iii) SBM-TNN outperforms the hyperbolic-based NNs, i.e., HGCN with a statistically significant margin; (iv) SBM-TNN further improves topology-based GNN (i.e., BScNets and TLC-GNN) with a significant margin on all 5 datasets. We also compare with two additional state-of-the-art baselines, i.e., NCNC (Wang et al., 2024) and LPFormer (Shomer et al., 2024) on Cora-ML and Citeseer datasets with MRR, and we observe that our SBM-TNN achieves average improvements of 23.1% and 3.3% over NCNC and LPFormer respectively (further details can be found in the Appendix). Additionally, Table 2 shows the performance of SBM-TNN and baseline methods on 3 KG datasets. From Table 2, we observe that SBM-TNN surpasses the baselines in terms of the MRR, Hits@1, Hits@3, and Hits@10 on all datasets. Furthermore, we have conducted an additional comparison with the GNNs + NBFNet on ogbl-wikikg2 data (Hu et al., 2020). The test MRR of SBM-TNN and GNNs + NBFNet are  $0.7121 \pm 0.0009$  and  $0.7086 \pm 0.0028$ , i.e., our SBM-TNN is significantly better than this state-of-the-art method. Overall, the results show that SBM-TNN can accurately capture and learn the key structural and local topological information, and achieve highly promising performance in both link prediction and KG completion tasks.

**Sensitivity Analysis** To evaluate the link prediction performance of SBM-TNN with different  $\tau$ , we conduct experiments on the Cora-ML, Citeseer, and PubMed datasets. As shown in Table 3, we observe that our approach achieves optimal performance under a specific power of the normalized adjacency matrix. In addition, we have run additional experiments on Cora-ML and Citeseer. As shown in Table 4, SBM-TNN with estimated  $\hat{\mathbf{P}}$  outperforms SBM-TNN with actual  $\mathbf{P}$ . In particular, the average relative gain of SBM-TNN with estimated  $\hat{\mathbf{P}}$  over SBM-TNN with actual  $\mathbf{P}$  is 1.58%.

**Ablation Study** We have also conducted ablation studies to explore the importance of different components, and considered two ablated variants, i.e., (i) SBM-GNN represents that replacing TNN by a graph neural network (GNN), and (ii) TNN represents that SBM-TNN without adding SBM method, (i.e., we do not incorporate community-level information into the model architecture). From Table 5, we observe that our SBM-TNN always outperforms both the GNN model equipped with SBM (i.e., SBM-GNN) and TNN model on Cora-ML and Citeseer data. That is, when ablating the components (i.e., SBM and TNN), the ROC AUC score of SBM-TNN drops significantly. Our results indicate that community and topological information consistently boost the performance of link prediction.

Table 5: Ablation study.

Architecture	Cora-ML	Citeseer	PubMed
SBM-TNN	<b>96.2±0.2</b>	<b>97.1±0.3</b>	<b>98.2±0.1</b>
SBM-GNN	90.7±0.3	83.7±0.4	93.7±0.2
TNN	93.2±0.5	96.0±0.5	95.3±0.6

**Computational Complexity** The topological complexity of the standard persistent homology (PH) matrix reduction algorithm runs in time at most  $\mathcal{O}(Q^3)$ , where  $Q$  is the number of simplices in a filtration. For 0-dimensional PH, it can be computed efficiently using disjoint sets with complexity  $\mathcal{O}(Q\alpha^{-1}(Q))$ , where  $\alpha^{-1}(\cdot)$  is the inverse Ackermann function. In our study, for graph representation learning, we only consider dimension 0 (connected components) and dimension 1 (cycles) due to the fact that we cannot observe enough higher-order (sub)structures in target datasets. If we consider high-dimensional topological features, the time complexity will grow large with the worst-case complexity  $\mathcal{O}(m^d)$  for  $d$ -dimensional topological features (where  $m$  denotes the number of edges).

## 7 Conclusion

We propose a new Stochastic Block Model-Aware Topological Neural Networks (SBM-TNN) method for both link prediction and knowledge graph completion tasks. By leveraging the topological information and estimated probability matrix with communities from different network topologies, SBM-TNN achieves state-of-the-art results on all datasets and the experimental evaluation confirms that SBM-TNN is accurate, flexible, and scalable. We also provide a theoretical guarantee for statistical inference based on the estimated edge probabilities. Interesting future directions include extending how the SBM-TNN can be used for non-attributed, overlapping, and dynamic network community detection.

## References

- Emmanuel Abbe. Community detection and stochastic block models: recent developments. *Journal of Machine Learning Research*, 18(177):1–86, 2018.
- Emmanuel Abbe, Afonso S Bandeira, and Georgina Hall. Exact recovery in the stochastic block model. *IEEE Transactions on Information Theory*, 62(1):471–487, 2015.
- Peter Bickel, David Choi, Xiangyu Chang, and Hai Zhang. Asymptotic normality of maximum likelihood and its variational approximation for stochastic blockmodels. *The Annals of Statistics*, 41(4):1922 – 1943, 2013. doi: 10.1214/13-AOS1124. URL <https://doi.org/10.1214/13-AOS1124>.
- Aleksandar Bojchevski and Stephan Günnemann. Deep gaussian embedding of graphs: Unsupervised inductive learning via ranking. In *ICLR*, 2018.
- Michael M Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Vandergheynst. Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine*, 34(4):18–42, 2017.
- Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. Spectral networks and locally connected networks on graphs. In *ICLR*, 2014.
- Joshua Cape, Minh Tang, and Carey E Priebe. Signal-plus-noise matrix models: eigenvector deviations and fluctuations. *Biometrika*, 106(1):243–250, 2019.
- Ines Chami, Zhitao Ying, Christopher Ré, and Jure Leskovec. Hyperbolic graph convolutional neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- Yuzhou Chen, Baris Coskunuzer, and Yulia Gel. Topological relational learning on graphs. *Advances in Neural Information Processing Systems*, 34:27029–27042, 2021.

- Yuzhou Chen, Yulia R Gel, and H Vincent Poor. BScNets: Block simplicial complex neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 6333–6341, 2022.
- Yuzhou Chen, Roshni Anna Jacob, Yulia R Gel, Jie Zhang, and H Vincent Poor. Learning power grid outages with higher-order topological neural networks. *IEEE Transactions on Power Systems*, 2023.
- Zhengdao Chen, Lisha Li, and Joan Bruna. Supervised community detection with line graph neural networks. In *ICLR*, 2018.
- Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. *NeurIPS*, 29, 2016.
- Tim Dettmers, Pasquale Minervini, Pontus Stenetorp, and Sebastian Riedel. Convolutional 2d knowledge graph embeddings. In *Proceedings of the AAAI conference on Artificial Intelligence*, volume 32, 2018.
- Herbert Edelsbrunner. Persistent homology: theory and practice. 2013.
- Chao Gao, Yu Lu, and Harrison H. Zhou. Rate-optimal graphon estimation. *The Annals of Statistics*, 43(6): 2624 – 2652, 2015. doi: 10.1214/15-AOS1354. URL <https://doi.org/10.1214/15-AOS1354>.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. *NeurIPS*, 30, 2017.
- Felix Hensel, Michael Moor, and Bastian Rieck. A survey of topological machine learning methods. *Frontiers in Artificial Intelligence*, 4:681108, 2021.
- Paul W Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. Stochastic blockmodels: First steps. *Social Networks*, 5(2):109–137, 1983.
- Max Horn, Edward De Brouwer, Michael Moor, Yves Moreau, Bastian Rieck, and Karsten Borgwardt. Topological graph neural networks. In *International Conference on Learning Representations*, 2021.
- Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. *Advances in Neural Information Processing Systems*, 33:22118–22133, 2020.
- Antony Joseph and Bin Yu. Impact of regularization on spectral clustering. *The Annals of Statistics*, 44(4): 1765–1791, 2016.
- Brian Karrer and Mark EJ Newman. Stochastic blockmodels and community structure in networks. *Physical Review E*, 83(1):016107, 2011.
- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *Proceedings of the International Conference on Learning Representations*, 2017.
- Jing Lei and Alessandro Rinaldo. Consistency of spectral clustering in stochastic block models. *The Annals of Statistics*, 43(1):215 – 237, 2015. doi: 10.1214/14-AOS1274. URL <https://doi.org/10.1214/14-AOS1274>.
- Ren Li, Yanan Cao, Qiannan Zhu, Guanqun Bi, Fang Fang, Yi Liu, and Qian Li. How does knowledge graph embedding extrapolate to unseen data: a semantic evidence view. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pp. 5781–5791, 2022.
- Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. Diffusion convolutional recurrent neural network: Data-driven traffic forecasting. In *ICLR*, 2018.
- Shujie Ma, Liangjun Su, and Yichong Zhang. Determining the number of communities in degree-corrected stochastic block models. *The Journal of Machine Learning Research*, 22(1):3217–3279, 2021.
- Nikhil Mehta, Lawrence Carin Duke, and Piyush Rai. Stochastic blockmodels meet graph neural networks. In *Proceedings of the International Conference on Machine Learning*, pp. 4466–4474. PMLR, 2019.

- Deepak Nathani, Jatin Chauhan, Charu Sharma, and Manohar Kaul. Learning attention-based embeddings for relation prediction in knowledge graphs. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4710–4723, 2019.
- Chi Seng Pun, Si Xian Lee, and Kelin Xia. Persistent-homology-based machine learning: a survey and a comparative study. *Artificial Intelligence Review*, 55(7):5169–5213, 2022.
- Feiliang Ren, Juchen Li, Huihui Zhang, Shilei Liu, Bochao Li, Ruicheng Ming, and Yujia Bai. Knowledge graph embedding with atrous convolution and residual learning. *arXiv preprint arXiv:2010.12121*, 2020.
- Karl Rohe, Sourav Chatterjee, and Bin Yu. Spectral clustering and the high-dimensional stochastic blockmodel. *The Annals of Statistics*, 39(4):1878 – 1915, 2011. doi: 10.1214/11-AOS887. URL <https://doi.org/10.1214/11-AOS887>.
- Karl Rohe, Tai Qin, and Bin Yu. Co-clustering directed graphs to discover asymmetries and directional communities. *Proceedings of the National Academy of Sciences*, 113(45):12679–12684, 2016.
- Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne Van Den Berg, Ivan Titov, and Max Welling. Modeling relational data with graph convolutional networks. In *The semantic web: 15th international conference, ESWC 2018, Heraklion, Crete, Greece, June 3–7, 2018, proceedings 15*, pp. 593–607. Springer, 2018.
- Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. *AI Magazine*, 29(3):93–93, 2008.
- Oleksandr Shchur and Stephan Günnemann. Overlapping community detection with graph neural networks. *arXiv preprint arXiv:1909.12201*, 2019.
- Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. Pitfalls of graph neural network evaluation. *arXiv preprint arXiv:1811.05868*, 2018.
- Harry Shomer, Yao Ma, Haitao Mao, Juanhui Li, Bo Wu, and Jiliang Tang. LPFormer: An adaptive graph transformer for link prediction. In *Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 2686–2698, 2024.
- Liangjun Su, Wuyi Wang, and Yichong Zhang. Strong consistency of spectral clustering for stochastic block models. *IEEE Transactions on Information Theory*, 66(1):324–338, 2019.
- Minh Tang and Carey E Priebe. Limit theorems for eigenvectors of the normalized laplacian for random graphs. *The Annals of Statistics*, 46(5):2360–2415, 2018.
- Minh Tang, Joshua Cape, and Carey E Priebe. Asymptotically efficient estimators for stochastic blockmodels: The naive mle, the rank-constrained mle, and the spectral estimator. *Bernoulli*, 28(2):1049–1073, 2022.
- Kristina Toutanova and Danqi Chen. Observed versus latent features for knowledge base and text inference. In *Proceedings of the Workshop on Continuous Vector Space Models and Their Compositionality*, pp. 57–66, 2015.
- Kristina Toutanova, Danqi Chen, Patrick Pantel, Hoifung Poon, Pallavi Choudhury, and Michael Gamon. Representing text for joint embedding of text and knowledge bases. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pp. 1499–1509, 2015.
- Shikhar Vashishth, Soumya Sanyal, Vikram Nitin, and Partha Talukdar. Composition-based multi-relational graph convolutional networks. *arXiv preprint arXiv:1911.03082*, 2019.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *Proceedings of the International Conference on Learning Representations*, 2018a.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *ICLR*, 2018b.

- U. von Luxburg. A tutorial on spectral clustering. *Statistics and Computing*, 17(4):395–416, 2007.
- Xiyuan Wang, Haotong Yang, and Muhan Zhang. Neural common neighbor with completion for link prediction. In *Proceedings of the International Conference on Learning Representations*, 2024.
- Larry Wasserman. Topological data analysis. *Annual Review of Statistics and Its Application*, 5:501–532, 2018.
- Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. A comprehensive survey on graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 32(1):4–24, 2020.
- Feng Xia, Ke Sun, Shuo Yu, Abdul Aziz, Liangtian Wan, Shirui Pan, and Huan Liu. Graph learning: A survey. *IEEE Trans AI*, 2(2):109–127, 2021.
- Fangzheng Xie. Entrywise limit theorems for eigenvectors of signal-plus-noise matrix models with weak signals. *Bernoulli*, 30(1):388–418, 2024.
- Zhiwen Xie, Guangyou Zhou, Jin Liu, and Xiangji Huang. ReInceptionE: relation-aware inception network with joint local-global structural information for knowledge graph embedding. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 5929–5939, 2020.
- Wenhan Xiong, Thien Hoang, and William Yang Wang. DeepPath: A reinforcement learning method for knowledge graph reasoning. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2017.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *ICLR*, 2018.
- Zuoyu Yan, Tengfei Ma, Liangcai Gao, Zhi Tang, and Chao Chen. Link prediction with persistent homology: An interactive view. In *International Conference on Machine Learning*, pp. 11659–11669. PMLR, 2021.
- Zhitao Ying, Jiaxuan You, Christopher Morris, Xiang Ren, Will Hamilton, and Jure Leskovec. Hierarchical graph representation learning with differentiable pooling. *NeurIPS*, 31, 2018.
- Jiaxuan You, Rex Ying, and Jure Leskovec. Position-aware graph neural networks. In *International Conference on Machine Learning*, pp. 7134–7143. PMLR, 2019.
- Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks. *Advances in Neural Information Processing Systems*, 31, 2018.
- Ziwei Zhang, Peng Cui, and Wenwu Zhu. Deep learning on graphs: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 34(1):249–270, 2020.
- Qi Zhao and Yusu Wang. Learning metrics for persistence-based summaries and applications for graph classification. *Advances in Neural Information Processing Systems*, 32, 2019.
- Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. Graph neural networks: A review of methods and applications. *AI open*, 1:57–81, 2020.
- Zhaocheng Zhu, Zuobai Zhang, Louis-Pascal Xhonneux, and Jian Tang. Neural bellman-ford networks: A general graph neural network framework for link prediction. *Advances in Neural Information Processing Systems*, 34:29476–29490, 2021.
- Afra Zomorodian and Gunnar Carlsson. Computing persistent homology. In *Proceedings of the twentieth annual symposium on Computational geometry*, pp. 347–356, 2004.

## A Derivation of Estimators under Undirected Network Models

Under the SBMs, to estimate  $\mathbf{B}$ , we note that for  $1 \leq q < l \leq K$ ,

$$\mathbf{B}_{ql} := \frac{\sum_{1 \leq i \neq j \leq N} \mathbf{P}_{ij} \mathbf{Z}_{iq} \mathbf{Z}_{jl}}{\sum_{1 \leq i \neq j \leq N} \mathbf{Z}_{iq} \mathbf{Z}_{jl}} = \frac{\sum_{g_i=q, g_j=l} \mathbf{P}_{ij}}{N_q N_l} \quad \text{and} \quad \mathbf{B}_{qq} := \frac{\sum_{g_i=q, g_j=q} \mathbf{P}_{ij}}{N_q(N_q - 1)}.$$

Thus, it is reasonable to estimate  $\mathbf{B}$  by the following  $\hat{\mathbf{B}} = (\hat{\mathbf{B}}_{ql})_{1 \leq q \leq l \leq K}$ ,

$$\hat{\mathbf{B}}_{ql, q \neq l} := \frac{\sum_{\hat{g}_i=q, \hat{g}_j=l} \mathbf{A}_{ij}}{\hat{N}_q \hat{N}_l} \quad \text{and} \quad \hat{\mathbf{B}}_{qq} := \frac{\sum_{\hat{g}_i=q, \hat{g}_j=q} \mathbf{A}_{ij}}{\hat{N}_q(\hat{N}_q - 1)}.$$

Under the DCSBMs, to estimate  $\mathbf{B}$ , we note that

$$\sum_{g_i=q, g_j=l} \mathbf{P}_{ij} = \sum_{g_i=q} \boldsymbol{\theta}_i \sum_{g_j=l} \boldsymbol{\theta}_j \mathbf{B}_{ql} = N_q N_l \mathbf{B}_{ql},$$

Hence, we can estimate  $\mathbf{B}$  by  $\hat{\mathbf{B}} = (\hat{\mathbf{B}}_{ql})_{1 \leq q, l \leq K}$ ,

$$\hat{\mathbf{B}}_{ql} := \frac{\sum_{1 \leq i \neq j \leq N} \mathbf{A}_{ij} \hat{\mathbf{Z}}_{iq} \hat{\mathbf{Z}}_{jl}}{\sum_{1 \leq i \neq j \leq N} \hat{\mathbf{Z}}_{iq} \hat{\mathbf{Z}}_{jl}} = \frac{\sum_{\hat{g}_i=q, \hat{g}_j=l} \mathbf{A}_{ij}}{\hat{N}_q \hat{N}_l}.$$

To estimate  $\boldsymbol{\theta}$ , we note that

$$\sum_j \mathbf{P}_{ij} = \boldsymbol{\theta}_i \sum_k \sum_{g_j=k} \boldsymbol{\theta}_j \mathbf{B}_{g_i g_j} = \boldsymbol{\theta}_i \sum_k N_k \mathbf{B}_{g_i k} = \boldsymbol{\theta}_i \sum_{g_l=g_i} \sum_{j=1}^N \mathbf{P}_{lj} / N_{g_i},$$

where we used  $\mathbf{B}_{g_i k} = \sum_{g_l=g_i, g_j=k} \mathbf{P}_{lj} / (N_{g_i} N_k)$ . Hence, we estimate  $\boldsymbol{\theta}_i$  by  $\hat{\boldsymbol{\theta}}_i$  defined as

$$\hat{\boldsymbol{\theta}}_i = \frac{\hat{N}_{\hat{g}_i} \sum_j \mathbf{A}_{ij}}{\sum_{\hat{g}_l=\hat{g}_i} \sum_{j=1}^N \mathbf{A}_{lj}},$$

where  $\hat{N}_{\hat{g}_i}$  is the number of nodes in the estimated community  $\hat{g}_i$ . We also denote  $\hat{\boldsymbol{\Theta}} = \text{diag}(\hat{\boldsymbol{\theta}})$ .

## B Proofs and Lemmas

For simplicity, we provide the proofs for the undirected network models SBMs and DCSBMs. In most cases, the proofs for the directed network models follow similarly. We highlight the differences if any.

### B.1 Proof of Lemma 5.6

*Proof.* First, we consider the strong consistency of community detection. Under the undirected network model SBM, the strong consistency of  $\hat{\mathbf{Z}}$  follows from the Corollary II.1 in Su et al. (2019) and Assumptions 5.1 - 5.3. Under the undirected network model DCSBM, the strong consistency of  $\hat{\mathbf{Z}}$  follows from Corollary III.1 in Su et al. (2019), where they considered the regularized DCSBMs with regularization parameter  $\tau$ . In our set-up,  $\tau = 0$ . Assumptions 5.1 - 5.4 imply that Assumptions 11-13 in Su et al. (2019) are satisfied. Hence the result follows. Under the directed network model ScBM, the strong consistency of  $\hat{\mathbf{Y}}$  is similarly derived as that under the SBM. The strong consistency of  $\hat{\mathbf{Z}}$  mainly depends on the success of Theorem II.1 in Su et al. (2019). In particular, Assumption 5.5 (with the notation simplified to the ScBMs) implies that the rows of the population singular vectors  $\mathbf{V} = (v_i) \in \mathbb{R}^{n \times K^z}$  satisfies that

$$n^{1/2} \|v_i - v_j\|_2 \geq C > 0$$

for  $g_i^z \neq g_j^z$  and some constant  $C > 0$  (see details in Theorem 2.1 in the first version of Su et al. (2019)). Hence, Theorem II.1 in Su et al. (2019) holds and the strong consistency of  $\hat{\mathbf{Z}}$  follows from Corollary II.1 in Su et al. (2019). Under the directed network model DCScBM, the strong consistency of  $\hat{\mathbf{Y}}$  is similarly derived as that under the DCScBM. The strong consistency of  $\hat{\mathbf{Z}}$  mainly depends on the success of Theorem III.4 in Su et al. (2019). In particular, Assumption 5.5 implies that  $L_2$  normalized rows of the population singular vectors  $\mathbf{V} = (v_i) \in \mathbb{R}^{n \times K^z}$  satisfies that

$$\left\| \frac{v_i}{\|v_i\|_2} - \frac{v_j}{\|v_j\|_2} \right\|_2 \geq C' > 0$$

for  $g_i^z \neq g_j^z$  and some constant  $C' > 0$  Rohe et al. (2016). Hence, Theorem III.4 in Su et al. (2019) holds and the strong consistency of  $\hat{\mathbf{Z}}$  follows from Corollary III.1 in Su et al. (2019).

The strong consistency of  $\hat{\boldsymbol{\theta}}$  (resp.  $\hat{\boldsymbol{\theta}}^y$  and  $\hat{\boldsymbol{\theta}}^z$ ) under DCScBMs (resp. DCScBMs) follows from Theorem III.6 in Su et al. (2019). Actually, the strong consistency of  $\hat{\mathbf{Z}}$  (resp.  $\hat{\mathbf{Y}}$  and  $\hat{\mathbf{Z}}$ ) and Assumption 5.4 implies Assumption 15 in Su et al. (2019). And the result follows by noting that the minimal average degree in our set-up is  $O(N\rho)$ .

Now we proceed to show the strong consistency of  $\hat{\mathbf{B}}$ . The proof holds for both undirected and directed network models. For simplicity, we use the notation under undirected network models. For an  $\epsilon_N$  to be selected, we have

$$\begin{aligned} & \mathbb{P} \left( \sup_{1 \leq q < l \leq K} |\hat{\mathbf{B}}_{ql} - \mathbf{B}_{ql}| \geq \epsilon_N \quad i.o. \right) \\ & \leq \mathbb{P} \left( \sup_{1 \leq q < l \leq K} |\hat{\mathbf{B}}_{ql} - \mathbf{B}_{ql}| \geq \epsilon_N \quad i.o., \sup_{1 \leq i \leq N} \mathbf{1}\{\hat{g}_i \neq g_i\} = 0 \right) + \mathbb{P} \left( \sup_{1 \leq i \leq N} \mathbf{1}\{\hat{g}_i \neq g_i\} > 0 \quad i.o. \right) \\ & \leq \mathbb{P} \left( \sup_{1 \leq q < l \leq K} \left| \sum_{g_i=q, g_j=l} \left( \frac{\mathbf{A}_{ij}}{N_q N_l} - \frac{\mathbf{P}_{ij}}{N_q N_l} \right) \right| \geq \epsilon_N \quad i.o. \right), \end{aligned} \quad (12)$$

where the last inequality follows from the strong consistency of  $\hat{g}_i$ 's. To make the RHS of equation 12 zero, it suffices to show that

$$\sum_{K=1}^{\infty} \sum_{1 \leq q < l \leq K} \mathbb{P} \left( \left| \sum_{g_i=q, g_j=l} \left( \frac{\mathbf{A}_{ij}}{N_q N_l} - \frac{\mathbf{P}_{ij}}{N_q N_l} \right) \right| \geq \epsilon_N \right) < \infty \quad (13)$$

for some  $\epsilon_N$ . To this end, we use the Bernstein inequality. Define  $X^{(ij)} = \frac{\mathbf{A}_{ij} - \mathbf{P}_{ij}}{N_q N_l}$ , we have  $E(X^{(ij)}) = 0$ ,  $|X^{(ij)}| \leq \frac{1}{N_q N_l}$  and

$$\sum_{g_i=q, g_j=l} \mathbb{E}[(X^{(ij)})^2] = \sum_{g_i=q, g_j=l} \mathbb{E}\left[\frac{(\mathbf{A}_{ij} - \mathbf{P}_{ij})^2}{N_q^2 N_l^2}\right] = \sum_{g_i=q, g_j=l} \frac{\mathbf{P}_{ij}(1 - \mathbf{P}_{ij})}{N_q^2 N_l^2} = \frac{\mathbf{B}_{ql}(1 - \mathbf{B}_{ql})}{N_q N_l}.$$

Then by the Bernstein equality, we have

$$\mathbb{P} \left( \left| \sum_{g_i=q, g_j=l} X^{(ij)} \right| \geq \epsilon_N \right) \leq \exp \left( - \frac{\frac{1}{2} \epsilon_N^2}{\frac{\mathbf{B}_{ql}(1 - \mathbf{B}_{ql})}{N_q N_l} + \frac{\epsilon_N}{3N_q N_l}} \right).$$

Choosing  $\epsilon_N = C \max_{ql} \sqrt{\frac{\mathbf{B}_{ql}(1 - \mathbf{B}_{ql})}{N_q N_l}} \cdot \sqrt{\log N} = O\left(\frac{\sqrt{\rho \log N}}{N}\right)$ , it is easy to see that

$$\mathbf{B}_{ql}(1 - \mathbf{B}_{ql}) \asymp \rho \gtrsim \frac{\sqrt{\rho \log N}}{N} \asymp \epsilon_N,$$

where the inequality follows from Assumption 5.2. We thus have

$$\mathbb{P} \left( \left| \sum_{g_i=q, g_j=l} X^{(ij)} \right| \geq \epsilon_N \right) \leq N^{-\alpha}$$



for some constant  $\alpha > 0$ . As a result, equation 13 is met because of fixed  $K$ . Finally, we obtain the strong consistency of  $\hat{\mathbf{B}}$  that

$$\sup_{1 \leq q < l \leq K} |\hat{\mathbf{B}}_{ql} - \mathbf{B}_{ql}| = O_{a.s.} \left( \frac{\sqrt{\rho \log N}}{N} \right).$$

The proof for  $q = l$  goes similarly by noting that we can represent  $\hat{\mathbf{B}}_{qq} - \mathbf{B}_{qq}$  as the following summation of independent terms,

$$\hat{\mathbf{B}}_{qq} - \mathbf{B}_{qq} = \sum_{g_i=q, g_j=q, i < j} \frac{\mathbf{A}_{ij} - \mathbf{P}_{ij}}{N_q(N_q - 1)/2}.$$

Hence we omit it.  $\square$

## B.2 Proof of Theorem 5.8

*Proof.* We first provide the proof for the more general model DCSBMs, and then present the results for SBMs as a special case.

It is easy to see that

$$\begin{aligned} \hat{\mathbf{P}} - \mathbf{P} &= \hat{\Theta} \hat{\mathbf{Z}} \hat{\mathbf{B}} \hat{\mathbf{Z}}^\top \hat{\Theta} - \mathbf{P} = (\hat{\Theta} - \Theta) \hat{\mathbf{Z}} \hat{\mathbf{B}} \hat{\mathbf{Z}}^\top \hat{\Theta} + \Theta \hat{\mathbf{Z}} \hat{\mathbf{B}} \hat{\mathbf{Z}}^\top (\hat{\Theta} - \Theta) + \Theta \hat{\mathbf{Z}} \hat{\mathbf{B}} \hat{\mathbf{Z}}^\top (\Theta) - \mathbf{P} \\ &:= I + II + III - \mathbf{P}. \end{aligned} \quad (14)$$

Now we proceed to bound  $I$ . We have by Lemma 5.6 that

$$\sup_i |\hat{\theta}_i - \theta_i| = O_{a.s.}(\log N / (N\rho)),$$

and thus

$$\sup_i \hat{\theta}_i \leq \sup_i |\hat{\theta}_i - \theta_i| + \sup_i \theta_i = O_{a.s.}(1)$$

by noting Assumption 5.4. By Eqs. 8 and 9 in Lemma 5.6, We can easily have

$$\begin{aligned} \sup_{\hat{g}_i, \hat{g}_j} \hat{\mathbf{B}}_{\hat{g}_i \hat{g}_j} &=_{a.s.} \sup_{g_i, g_j} \hat{\mathbf{B}}_{g_i g_j} \leq \sup_{g_i, g_j} |\hat{\mathbf{B}}_{g_i g_j} - \mathbf{B}_{g_i g_j}| + \sup_{g_i, g_j} \mathbf{B}_{g_i g_j} \\ &= O_{a.s.} \left( \frac{\sqrt{\rho \log N}}{N} \right) + O_{a.s.}(\rho) = O_{a.s.}(\rho), \end{aligned}$$

where the last equality is implied by Assumption 5.2. As a result,

$$\begin{aligned} \sup_{1 \leq i, j \leq N} |(I)_{ij}| &= \sup_{1 \leq i, j \leq N} |(\text{diag}(\hat{\theta} - \theta) \hat{\mathbf{Z}} \hat{\mathbf{B}} \hat{\mathbf{Z}}^\top \text{diag}(\hat{\theta}))_{ij}| \\ &= \sup_{1 \leq i, j \leq N} |(\hat{\theta}_i - \theta_i) \hat{\mathbf{B}}_{\hat{g}_i \hat{g}_j} \hat{\theta}_j| = O_{a.s.}(\log N / N). \end{aligned} \quad (15)$$

Similarly, we can show that

$$\sup_{1 \leq i, j \leq N} |(II)_{ij}| = O_{a.s.}(\log N / N). \quad (16)$$

It remains to bound  $III - \mathbf{P}$ . Noting

$$\Theta \hat{\mathbf{Z}} \hat{\mathbf{B}} \hat{\mathbf{Z}}^\top \Theta - \mathbf{P} = \Theta (\hat{\mathbf{Z}} \hat{\mathbf{B}} \hat{\mathbf{Z}}^\top - \mathbf{Z} \mathbf{B} \mathbf{Z}^\top) \Theta \quad (17)$$

and the boundness of  $\theta$  by Assumption 5.4, we only need to bound  $\hat{\mathbf{Z}} \hat{\mathbf{B}} \hat{\mathbf{Z}}^\top - \mathbf{Z} \mathbf{B} \mathbf{Z}^\top$ . It is easy to observe that

$$\begin{aligned} \hat{\mathbf{Z}} \hat{\mathbf{B}} \hat{\mathbf{Z}}^\top - \mathbf{Z} \mathbf{B} \mathbf{Z}^\top &= (\hat{\mathbf{Z}} - \mathbf{Z}) \hat{\mathbf{B}} \hat{\mathbf{Z}}^\top + \mathbf{Z} \mathbf{B} (\hat{\mathbf{Z}} - \mathbf{Z})^\top + \hat{\mathbf{Z}} (\hat{\mathbf{B}} - \mathbf{B}) \hat{\mathbf{Z}}^\top \\ &:= E_1 + E_2 + E_3. \end{aligned}$$

For  $E_3$ , we further have

$$\begin{aligned} E_3 &= \hat{\mathbf{Z}}(\hat{\mathbf{B}} - \mathbf{B})\hat{\mathbf{Z}}^\top \\ &= \mathbf{Z}(\hat{\mathbf{B}} - \mathbf{B})\mathbf{Z}^\top + (\hat{\mathbf{Z}} - \mathbf{Z})(\hat{\mathbf{B}} - \mathbf{B})\hat{\mathbf{Z}}^\top + \mathbf{Z}(\hat{\mathbf{B}} - \mathbf{B})(\hat{\mathbf{Z}} - \mathbf{Z})^\top. \end{aligned}$$

Combining the above two facts, we obtain

$$\hat{\mathbf{Z}}\hat{\mathbf{B}}\hat{\mathbf{Z}}^\top - \mathbf{Z}\mathbf{B}\mathbf{Z}^\top = \mathbf{Z}(\hat{\mathbf{B}} - \mathbf{B})\mathbf{Z}^\top + R, \quad (18)$$

where

$$R := (\hat{\mathbf{Z}} - \mathbf{Z})\mathbf{B}\hat{\mathbf{Z}}^\top + \mathbf{Z}\mathbf{B}(\hat{\mathbf{Z}} - \mathbf{Z})^\top + (\hat{\mathbf{Z}} - \mathbf{Z})(\hat{\mathbf{B}} - \mathbf{B})\hat{\mathbf{Z}}^\top + \mathbf{Z}(\hat{\mathbf{B}} - \mathbf{B})(\hat{\mathbf{Z}} - \mathbf{Z})^\top.$$

By the strong consistency of  $\hat{\mathbf{Z}}$ ,

$$\hat{\mathbf{Z}} - \mathbf{Z} = \mathbf{0}, \quad a.s., \quad (19)$$

and thus  $R =_{a.s.} \mathbf{0}$ . It remains to bound  $\mathbf{Z}(\hat{\mathbf{B}} - \mathbf{B})\mathbf{Z}^\top$ . Noting

$$(\mathbf{Z}(\hat{\mathbf{B}} - \mathbf{B})\mathbf{Z}^\top)_{ij} = (\hat{\mathbf{B}} - \mathbf{B})_{g_i g_j}$$

and equation 7 in Lemma 5.6, we have

$$\sup_{i,j} |\mathbf{Z}(\hat{\mathbf{B}} - \mathbf{B})\mathbf{Z}^\top|_{ij} = O_{a.s.} \left( \frac{\sqrt{\rho \log N}}{N} \right).$$

As a result, we have

$$\sup_{1 \leq i,j \leq N} |(III - \mathbf{P})_{ij}| = O_{a.s.} \left( \frac{\sqrt{\rho \log N}}{N} \right). \quad (20)$$

Combining equation 15, equation 16 with equation 20, we have for DCSBMs that

$$\sup_{1 \leq i,j \leq N} |\hat{\mathbf{P}}_{ij} - \mathbf{P}_{ij}| = O_{a.s.} \left( \frac{\log N}{N} \right). \quad (21)$$

While for SBMs, the error terms  $I$  and  $II$  are exactly zero. Hence,

$$\sup_{1 \leq i,j \leq N} |\hat{\mathbf{P}}_{ij} - \mathbf{P}_{ij}| = O_{a.s.} \left( \frac{\sqrt{\rho \log N}}{N} \right).$$

The proofs for the directed network models ScBMs and DCScBMs go similarly provided that Assumption 5.5 is additionally required for Lemma 5.6.  $\square$

### B.3 Proof of Theorem 5.8

*Proof.* Recall equation 14, we can decompose  $\hat{\mathbf{P}} - \mathbf{P}$  by  $I + II + III - \mathbf{P}$ . For DCSBMs, we have shown in the proof of Theorem 5.7 that  $\sup |(I)_{ij}| = O_{a.s.}(\log N/N)$  and  $\sup |(II)_{ij}| = O_{a.s.}(\log N/N)$ . For SBMs,  $I = II = \mathbf{0}$ .

We now show the asymptotic normality of  $III - \mathbf{P}$ , for which by equation 17, equation 18 and equation 19, we only need to show the asymptotic normality of  $\mathbf{Z}(\hat{\mathbf{B}} - \mathbf{B})\mathbf{Z}^\top$ . Note that

$$(\mathbf{Z}(\hat{\mathbf{B}} - \mathbf{B})\mathbf{Z}^\top)_{ij} = (\hat{\mathbf{B}} - \mathbf{B})_{g_i g_j} = \hat{\mathbf{B}}_{ql} - \mathbf{B}_{ql}$$

provided that  $g_i = q$  and  $g_j = l$ . By strong consistency of  $\hat{\mathbf{Z}}$ , we have  $\hat{g}_i =_{a.s.} g_i$  and  $\hat{N}_q =_{a.s.} N_q$  for all  $i \in [N]$  and  $q \in [K]$ . So for  $q \neq l$ ,

$$\hat{\mathbf{B}}_{ql} - \mathbf{B}_{ql} = \sum_{g_i=q, g_j=l} \frac{\mathbf{A}_{ij} - \mathbf{P}_{ij}}{N_q N_l}.$$

We use the Lindeberg-Feller Central Limit Theorem to derive the limit distribution of  $\hat{\mathbf{B}}_{ql} - \mathbf{B}_{ql}$ . First, note that

$$\begin{aligned} s_N^2 &:= \text{Var}\left(\sum_{g_i=q, g_j=l} \frac{\mathbf{A}_{ij} - \mathbf{P}_{ij}}{N_q N_l}\right) = \sum_{g_i=q, g_j=l} \text{Var}\left(\frac{\mathbf{A}_{ij}}{N_q N_l}\right) \\ &= \sum_{g_i=q, g_j=l} \mathbf{P}_{ij}(1 - \mathbf{P}_{ij})/(N_q^2 N_l^2) = \mathbf{B}_{ql}(1 - \mathbf{B}_{ql})/(N_q N_l) \end{aligned}$$

We only need to show

$$\frac{1}{s_N^2} \sum_{g_i=q, g_j=l} \mathbb{E}\left\{\left(\frac{\mathbf{A}_{ij} - \mathbf{P}_{ij}}{N_q N_l}\right)^2 \mathbb{I}\left(\left|\frac{\mathbf{A}_{ij} - \mathbf{P}_{ij}}{N_q N_l}\right| \geq \epsilon s_N\right)\right\} \rightarrow 0$$

for every  $\epsilon > 0$ , which holds sufficiently if  $\left|\frac{\mathbf{A}_{ij} - \mathbf{P}_{ij}}{N_q N_l}\right| \leq \epsilon s_N$ . Indeed,

$$\left|\frac{\mathbf{A}_{ij} - \mathbf{P}_{ij}}{N_q N_l}\right| \lesssim \frac{1}{N^2} \lesssim \frac{\rho^{1/2}}{N} \asymp \epsilon \sqrt{\mathbf{B}_{ql}(1 - \mathbf{B}_{ql})/N_q N_l}.$$

Therefore, by central limit theorem, we can obtain in SBMs that

$$N(\hat{\mathbf{B}}_{ql} - \mathbf{B}_{ql}) \rightarrow N\left(0, \frac{\mathbf{B}_{ql}(1 - \mathbf{B}_{ql})}{\pi_q \pi_l}\right)$$

for  $q \neq l$ . When  $q = l$ , considering the dependency of  $\mathbf{A}_{ij}$  and  $\mathbf{A}_{ji}$ , we consider the following halved term

$$\hat{\mathbf{B}}_{qq} - \mathbf{B}_{qq} = \sum_{g_i=q, g_j=q, i < j} \frac{\mathbf{A}_{ij} - \mathbf{P}_{ij}}{N_q(N_q - 1)/2}. \quad (22)$$

Similarly, we can derive

$$N(\hat{\mathbf{B}}_{qq} - \mathbf{B}_{qq}) \rightarrow N\left(0, \frac{2\mathbf{B}_{qq}(1 - \mathbf{B}_{qq})}{\pi_q^2}\right).$$

For DCSBMs, we thus have that

$$N(\hat{\mathbf{P}} - \mathbf{P} + \mathbf{E})_{ij} \rightarrow N\left(0, \boldsymbol{\theta}_i^2 \cdot \frac{\mathbf{B}_{g_i g_j}(1 - \mathbf{B}_{g_i g_j})}{\pi_{g_i} \pi_{g_j}} \cdot \boldsymbol{\theta}_j^2\right)$$

with  $\mathbf{E}_{ij} = O_{a.s.}(\log N/N)$  and  $g_i \neq g_j$ , and

$$N(\hat{\mathbf{P}} - \mathbf{P} + \mathbf{F})_{ij} \rightarrow N\left(0, 2\boldsymbol{\theta}_i^2 \cdot \frac{\mathbf{B}_{g_i g_i}(1 - \mathbf{B}_{g_i g_i})}{\pi_{g_i}^2} \cdot \boldsymbol{\theta}_j^2\right)$$

with  $\mathbf{F}_{ij} = O_{a.s.}(\log N/N)$  and  $g_i = g_j$ .

The proofs for the directed network models ScBMs and DCScBMs are similar except that when  $g_i^y = g_j^z$ , the directed networks do not involve dependency of pairs  $\mathbf{A}_{ij}$  and  $\mathbf{A}_{ji}$  like in (22), hence the constant 2 in the asymptotic variances are replaced by 1.

□

## B.4 Auxiliary lemmas

**Lemma B.1** (Bernstein inequality). *Let  $X_1, \dots, X_N$  be independent zero-mean random variables. Suppose that  $|X_i| \leq M$  almost surely, for all  $i$ . Then, for all positive  $t$ ,*

$$\mathbb{P}\left(\sum_{i=1}^N X_i \geq t\right) \leq \exp\left(-\frac{\frac{1}{2}t^2}{\sum_{i=1}^N \mathbb{E}[X_i^2] + \frac{1}{3}Mt}\right).$$

Table 6: A summary of the notes and notations.

Notation	Definition
<b>Notations for undirected network models:</b>	
$N$	Number of communities
$g_i \in \{1, \dots, K\}$	Community assignment of node $i$
$N_k$	Number of nodes within the $k$ -th community
$\pi_k := N_k/N$	Proportion of nodes in the $k$ -th community
$\mathbf{Z} \in \{0, 1\}^{N \times K}$	Membership matrix
$\mathbf{B} \in \mathbb{R}^{K \times K}$	Block probability matrix
$\boldsymbol{\theta} \in \mathbb{R}^N$	Node propensity vector
$\boldsymbol{\Theta} = \text{diag}\{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N\}$	Node propensity matrix
$\mathbf{P} := \boldsymbol{\Theta} \mathbf{Z} \mathbf{B} \mathbf{Z}^\top \boldsymbol{\Theta} \in \mathbb{R}^{N \times N}$	Edge probability matrix
$d_i = \sum_{j=1}^N \mathbf{P}_{ij}$	Population degree of node $i$
$\mathcal{D} = \text{diag}\{d_1, \dots, d_N\}$	Population degree matrix
$\rho$	Network sparsity
$\mathbf{A} \in \{0, 1\}^{N \times N}$	Symmetric Adjacency matrix
$\hat{d}_i = \sum_{j=1}^N \mathbf{A}_{ij}$	Degree of node $i$
$\mathbf{D} = \text{diag}\{\hat{d}_1, \dots, \hat{d}_N\}$	Degree matrix
$\mathbf{L} = \mathbf{I} + \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$	Graph Laplacian
<b>Notations for directed network models:</b>	
$N$	Number of nodes
$K^y (K^z)$	Number of row (column) clusters
$g_i^y \in \{1, \dots, K^y\} (g_i^z \in \{1, \dots, K^z\})$	Row (column) cluster assignment of node $i$
$N_k^y (N_k^z)$	Number of nodes in the $k^y$ th ( $k^z$ th) row (column) cluster
$\pi_k^y := N_k^y/N (\pi_k^z := N_k^z/N)$	Proportion of nodes in the $k^y$ th ( $k^z$ th) row (column) cluster
$\mathbf{Y} \in \{0, 1\}^{N \times K^y} (\mathbf{Z} \in \{0, 1\}^{N \times K^z})$	Row (Column) membership matrix
$\mathbf{B} \in \mathbb{R}^{K^y \times K^z} (K^y \leq K^z)$	Block probability matrix
$\boldsymbol{\theta}^y \in \mathbb{R}^N (\boldsymbol{\theta}^z \in \mathbb{R}^N)$	Node propensity vector for sending (receiving) edges
$\boldsymbol{\Theta}^y = \text{diag}\{\boldsymbol{\theta}^y\} (\boldsymbol{\Theta}^z = \text{diag}\{\boldsymbol{\theta}^z\})$	Node propensity matrix for sending (receiving) edges
$\mathbf{P} := \boldsymbol{\Theta}^y \mathbf{Y} \mathbf{B} \mathbf{Z}^\top \boldsymbol{\Theta}^z \in \mathbb{R}^{N \times N}$	Edge probability matrix
$d_i^y = \sum_{j=1}^N \mathbf{P}_{ij} (d_i^z = \sum_{j=1}^N \mathbf{P}_{ji})$	Population out-degree (in-degree) of node $i$
$\mathcal{D}^y = \text{diag}\{d_1^y, \dots, d_N^y\} (\mathcal{D}^z = \text{diag}\{d_1^z, \dots, d_N^z\})$	Population out-degree (in-degree) matrix
$\rho$	Network sparsity
$\mathbf{A} \in \{0, 1\}^{N \times N}$	Asymmetric adjacency matrix
$\hat{d}_i^y = \sum_{j=1}^N \mathbf{A}_{ij} (\hat{d}_i^z = \sum_{j=1}^N \mathbf{A}_{ji})$	Out-degree (in-degree) of node $i$
$\mathbf{D}^y = \text{diag}\{\hat{d}_1^y, \dots, \hat{d}_N^y\} (\mathbf{D}^z = \text{diag}\{\hat{d}_1^z, \dots, \hat{d}_N^z\})$	Out-degree (In-degree) matrix
$\mathbf{L} = \mathbf{I} + (\mathbf{D}^y)^{-1/2} \mathbf{A} (\mathbf{D}^z)^{-1/2}$	Graph Laplacian

## C Notes and notations

The notations are summarized in Table 6.

## D Additional Experiment Settings

For all 5 datasets (i.e., Cora-ML, Citeseer, PubMed, Photo, and Computers), SBM-TNN is trained by the Adam optimizer with the Cross Entropy Loss function. Additionally, all baseline methods are initialized with the parameters suggested in their respective works, we carefully tune the parameters during training to ensure that the baseline model achieves optimal performance. Here we treat the resulting topological features in dimension 0 (connected components) and 1 (cycles) (i.e.,  $Q = 2$ ). For link prediction, we perform an extensive grid search for learning rate among  $\{0.001, 0.005, 0.008, 0.01, 0.1\}$ , the dropout rate among  $\{0.1, 0.2, \dots, 0.9\}$ , the number of hidden units among  $\in \{8, 16, 32, 64, 128\}$ , and the model is trained for 5,000 epochs with early

stopping applied when the metric (i.e., validation loss) starts to drop. For KG completion, we set the batch size to be 512 and the model is trained for 500 epochs, and we perform an extensive grid search for learning rate among  $\{0.00001, 0.001, 0.01, 0.1\}$ .