

---

# Differentiable Community Detection with Graph Neural Networks and Stochastic Block Models

---

**William Arliss**  
Leidos  
william.f.arliss@leidos.com

**Graham Mueller**  
Leidos  
william.g.mueller@leidos.com

## Abstract

We propose a set of loss functions adapted from Stochastic Block Model (SBM) likelihood functions to train Graph Neural Networks (GNNs) for the task of unsupervised community detection. Identifying latent community structures is a prominent challenge for many graph applications. SBMs are classical models that describe the generating process of random graphs and are commonly used to infer community structure. The likelihood functions associated with SBMs are well-defined, differentiable, and measure the quality of inferred community partitions; this makes them particularly useful for unsupervised learning with GNNs. Our proposed loss functions are independent of any specific GNN architecture and demonstrate competitive or improved community detection performance against several alternatives. Evaluation is carried out on multiple architectures and datasets, offering a thorough empirical analysis of the state of community detection with GNNs.

## 1 Introduction

Graphs provide a rich source of relational information on which to perform a variety of machine learning tasks. Unsupervised community detection (also known as node clustering) refers to the problem of partitioning graph nodes into groups based on attributes and structural information. Methods for analyzing community structure are essential to applications in cybersecurity, social sciences, and e-commerce. For example, many Recommender Systems predict which product to recommend to a customer based on the inferred category (community) the product or customer belongs to.

A common tool for identifying and analyzing communities is the Stochastic Block Model (SBM) [1]. SBMs are statistical models of community structure in networks, parameterized by a node partition and a community structure matrix. Variations of the SBM have been proposed to address alternative assumptions about the generating processes of graphs [2–7].

The progress of Graph Neural Networks (GNNs) in representation learning on graphs has motivated their use for community detection as well. Several GNN-based frameworks have been proposed [8–12] and demonstrate impressive performance for both semi-supervised and unsupervised community detection. Most unsupervised methods involve estimating the partition matrix through modularity maximization, link prediction, or solving the minimum-cut problem; typically a combination of custom GNN architectures and training routines are used.

Significant work has also been done to integrate the strong theoretical foundations of SBMs with the expressive power of GNNs [13–17]. Usually, these approaches incorporate GNNs as a component in a mixture model or as a Bayesian prior.

A natural synthesis of the two approaches is to incorporate an SBM likelihood function as a loss function for training a GNN [18]. Adapting the likelihood functions for different types of SBMs offers a set of configurable, fully differentiable objectives which can be used for unsupervised training of an arbitrary neural network.

This paper has two main contributions: (i) A set of loss functions derived from SBM likelihood functions and (ii) an extensive comparison of unsupervised loss functions for the task of community detection. The loss functions are motivated by maximum likelihood estimation and the Graph Matching problem [19]. We compare the performance of GNN models trained with the SBM loss functions to several state-of-the-art alternatives on synthetic and real-world graphs.

The proposed loss functions are fully differentiable and do not require custom architectures or training routines. So for a fair empirical analysis, the same GNN architecture and training loop are used for each loss function. Consequently, comparison approaches that do not use stand-alone loss functions are not considered.

## 2 Related Work

### 2.1 Stochastic Block Models

Stochastic Block Models, introduced in [1], are a family of generative models which assume that the existence of any edge in a graph is dependent only on the partition (community membership) of the two component nodes. SBMs are identified by a node partition and a structure matrix, which defines the expectation of an edge between each partition.

Several variants of the SBM have been proposed to address different challenges: the Degree-Corrected SBM [2] mitigates the problem of skewed degree distributions by directly modeling degree heterogeneity; the Mixed Membership SBM [3] and the closely related Overlapping SBM [4] allow nodes to be members of more than one community, leading to more flexible interpretations of community structure; the Contextual SBM [6] incorporates node attributes, which are assumed to be generated conditionally on node communities; the Microcanonical SBM [5] enforces strict structural constraints in the model directly, rather than in expectation only. An in-depth review of these (and other) variants is given in [7].

The task of inference with SBMs typically involves identifying the process that generates a given graph and estimating the relevant parameters [7]. Inferring the partition of a graph from an SBM is sufficient for the task of community detection. A full survey of statistical community detection techniques related to SBM inference is given in [20]. While there are several ways to define community detection, this paper considers it to be the task of clustering nodes such that within group connectivity is high and between group connectivity is low [7, 9, 21, 22].

The success of Graph Neural Networks in representation learning on graphs (see [23–26]) has motivated several deep learning frameworks for SBMs. GNNs are used to parameterize Contextual SBMs in [13, 15, 17], where node features are assumed to be a function of community membership. Conversely, [14] uses a single-layer perceptron to model community membership as a function of node features. In [16] the authors design a variational auto-encoder GNN to parameterize an Overlapping SBM.

The SBM likelihood function is used as part of a composite loss function in [18]. In this work, the authors combine an approximate SBM log-likelihood, a custom link prediction loss, an entropy term, and task-specific losses in order to optimize a neural network with a custom training routine. The framework is evaluated on the tasks of community detection, graph alignment, and anomalous correlation detection. This approach differs from ours in the formulation of the adapted loss function and training routine as well as in the extent of evaluation.

### 2.2 Deep Community Detection

Deep community detection refers to unsupervised or semi-supervised community detection performed with deep neural networks. Much work has been done (orthogonally to the SBM class) on GNN-based community detection.

In [11], a framework consisting of multiple GNN layers is proposed, where one module generates node embeddings and the other pools node features according to (predicted) community structure. A link prediction objective is used to guide the pooling function. In [8], the authors apply a GNN to the minimum-cut problem, which seeks to partition the set of nodes into disjoint (i.e., minimally connected) subsets by maximizing the average ratio of edges within communities to edges between

communities. They do this by directly minimizing the negative of the minimum-cut metric plus a custom orthogonality constraint.

Both of the above approaches depend on custom architectures for task-specific problems. The focus of this paper, though, is on stand-alone objective functions that do not require custom architectures. One such general approach is taken in [9], where the authors propose using GNN embeddings to parameterize a Bernoulli-Poisson model [27, 28] of a graph. A likelihood-based loss function is derived from this model and edge sampling is used to address imbalance.

In [10] it is proposed to directly optimize modularity, a metric that measures the quality of community partitions. To train a GNN, the authors minimize the negative of estimated modularity plus a novel regularization term. The authors suggest that the orthogonality constraint from [8] tends to trap the optimization routine in local minima and instead devise their own ‘‘collapse regularization’’ meant to penalize trivial partitions. As a generalization of modularity optimization, [12] propose using the negative of the trace of the Markov Stability matrix [29–31] to train a GNN. Markov Stability is a dynamic quality metric that measures the probability that a random walk starting in one community will end in another after a certain number of time steps.

### 3 Methods

#### 3.1 Preliminary

Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  be a directed graph with  $n = |\mathcal{V}|$  nodes and  $m = |\mathcal{E}|$  unique edges. Furthermore, let  $\mathcal{D} = \mathcal{V} \times \mathcal{V}$  be the set of all possible node dyads so that  $\mathcal{E} \subseteq \mathcal{D}$ . The adjacency matrix  $\mathbf{A} \in \{0, 1\}^{n \times n}$  represents the edge structure of  $\mathcal{G}$  and the  $n$ -vector  $\mathbf{d}$  measures the node degrees such that  $\mathbf{d}_u = \sum_{v=1}^n \mathbf{A}_{uv}$ . Community memberships are represented in the  $k$ -partition matrix  $\mathbf{Z} \in \{0, 1\}^{n \times k}$ , where  $\mathbf{Z}_{ui} = 1$  if node  $u$  is a member of community  $i$  and 0 otherwise. It is also assumed that  $\sum_{i=1}^k \mathbf{Z}_{ui} = 1$  for all  $u \in \mathcal{V}$ .

#### 3.2 Likelihood Functions

We now consider the likelihood functions of several Stochastic Block Model variants induced by different assumptions of the underlying generating process of  $\mathcal{G}$ . All models have parameters  $\mathbf{Z}$  and  $\Theta$ , where  $\Theta$  is a  $k \times k$  structure matrix (also known as the block matrix). The block matrix is defined such that  $\Theta_{ij}$  is the expected value of any edge from a node in community  $i$  to a node in community  $j$ . That is,  $E(\mathbf{A}_{uv} | \mathbf{Z}) = \mathbf{Z}'_u \Theta \mathbf{Z}_v$  for all  $(u, v) \in \mathcal{D}$ .

##### 3.2.1 Poisson

One common form of the SBM assumes that elements of  $\mathbf{A}$  are Poisson distributed, conditional upon the community membership of the incident nodes [2, 7]. This gives the formal assumption  $\mathbf{A}_{uv} | \mathbf{Z} \sim \text{Pois}(\mathbf{Z}'_u \Theta \mathbf{Z}_v)$ . Here,  $\Theta_{ij}$  is interpreted as the average number of edges between nodes in communities  $i$  and  $j$ . The likelihood of this model is

$$\mathcal{L}_P(\mathbf{Z}, \Theta; \mathbf{A}) = \prod_{u,v} \frac{(\mathbf{Z}'_u \Theta \mathbf{Z}_v)^{\mathbf{A}_{uv}}}{\mathbf{A}_{uv}!} \exp(-\mathbf{Z}'_u \Theta \mathbf{Z}_v) \quad (1)$$

which is similar to the formulation in [2].

Recall that because  $\mathcal{G}$  is unweighted and has no parallel edges,  $\mathbf{A}$  only takes the values zero or one. This implies that equation 1 is in fact a *partial* likelihood, as it is defined over the  $\{0, 1\}$  sub-region of the standard Poisson support. The partial log-likelihood of this model is

$$\ell_P(\mathbf{Z}, \Theta; \mathbf{A}) = \sum_{u,v} [\mathbf{A}_{uv} \ln(\mathbf{Z}'_u \Theta \mathbf{Z}_v) - \mathbf{Z}'_u \Theta \mathbf{Z}_v]. \quad (2)$$

Note that the factorial term has been dropped because it is always equal to zero.

The maximum likelihood estimate (MLE) for  $\Theta$ , derived in [2], is the ratio of the number of edges between two communities to the product of their community sizes. That is,

$$\hat{\Theta}_{ij} = \left[ \hat{\Theta}(\mathcal{G}, \mathbf{Z}) \right]_{ij} = \frac{\mathbf{M}_{ij}}{\mathbf{n}_i \mathbf{n}_j} \quad (3)$$

where  $\mathbf{M} = \sum_{u,v \in \mathcal{E}} \mathbf{Z}_u \mathbf{Z}_v'$  is a  $k \times k$  matrix such that  $\mathbf{M}_{ij}$  is the number of edges from nodes in community  $i$  to nodes in community  $j$  (or twice that if  $i = j$ ) and  $\mathbf{n} = \sum_{u=1}^n \mathbf{Z}_u$  is a  $k$ -vector representing the number of nodes in each community. Note that  $\sum_{i=1}^k \sum_{j=1}^k \mathbf{M}_{ij} = m$  and  $\sum_{i=1}^k \mathbf{n}_i = n$ .

### 3.2.2 Bernoulli

In focusing on graphs without weighted or parallel edges, it is useful to consider a model that fully aligns with the restriction on the adjacency matrix. An intuitive choice is to assume that elements in  $\mathbf{A}$  are conditionally Bernoulli distributed [7]. Thus, the distribution assumption is modified to  $\mathbf{A}_{uv} | \mathbf{Z} \sim \text{Bern}(\mathbf{Z}'_u \boldsymbol{\Theta} \mathbf{Z}_v)$ . This new model interprets  $\boldsymbol{\Theta}_{ij}$  as the probability of an edge between nodes in communities  $i$  and  $j$ . The log-likelihood of this model is

$$\ell_B(\mathbf{Z}, \boldsymbol{\Theta}; \mathbf{A}) = \sum_{u,v} [\mathbf{A}_{uv} \ln(\mathbf{Z}'_u \boldsymbol{\Theta} \mathbf{Z}_v) + (1 - \mathbf{A}_{uv}) \ln(1 - \mathbf{Z}'_u \boldsymbol{\Theta} \mathbf{Z}_v)]. \quad (4)$$

The advantage of this model over the partial Poisson model is that it is defined over the full support of the assumed distribution, not just a sub-region. The MLE of  $\boldsymbol{\Theta}$  is the same as that of the Poisson model (see Appendix F.1).

### 3.2.3 Degree Correction

Another variety of the SBM, introduced in [2], seeks to incorporate degree heterogeneity into the model. The standard SBM assumes that nodes within the same block have the same expected degree. This assumption can lead to a sub-optimal solutions on real-world networks, where degree distributions are observed to be highly skewed.

To address this, the  $n$ -vector  $\phi$  is introduced, allowing heterogeneous degree expectations. The expected value of the dyad  $(u, v)$  is now  $\phi_u \phi_v \mathbf{Z}'_u \boldsymbol{\Theta} \mathbf{Z}_v$ . The partial log-likelihood for the degree-corrected Poisson model is

$$\ell_{\text{P-DC}}(\mathbf{Z}, \boldsymbol{\Theta}, \phi; \mathbf{A}) = \sum_{u,v} [\mathbf{A}_{uv} \ln(\phi_u \phi_v \mathbf{Z}'_u \boldsymbol{\Theta} \mathbf{Z}_v) - \phi_u \phi_v \mathbf{Z}'_u \boldsymbol{\Theta} \mathbf{Z}_v] \quad (5)$$

and the log-likelihood for the degree-corrected Bernoulli model, referred to as  $\ell_{\text{B-DC}}$ , is obtained by the same substitution.

The MLE for  $\phi$ , given in [2], is the ratio of a node's degree to the sum of degrees in that node's community. So the scaled degree correction of node  $u$  is

$$\hat{\phi}_u = \left[ \hat{\phi}(\mathcal{G}, \mathbf{Z}) \right]_u = (\mathbf{Z}'_u \mathbf{n}) \frac{\mathbf{d}_u}{\mathbf{Z}'_u \boldsymbol{\delta}} \quad (6)$$

where  $\boldsymbol{\delta} = \sum_{u=1}^n \mathbf{Z}_u \mathbf{d}_u$  is a  $k$ -vector representing the sum of degrees in each community. Also note that  $\mathbf{Z}'_u \boldsymbol{\delta}$  is the sum of degrees in the community that node  $u$  belongs to and  $\mathbf{Z}'_u \mathbf{n}$  is the size of that community.

In the Bernoulli model, the constraint that  $0 < \phi_u \phi_v \mathbf{Z}'_u \boldsymbol{\Theta} \mathbf{Z}_v < 1$  for all  $(u, v)$  must be observed. If  $\phi_u \phi_v$  scales the quantity to a value greater than 1, then the log-likelihood will be undefined. Therefore, it is necessary in  $\ell_{\text{B-DC}}$  to impose the boundary  $\phi_u \leq 1$  for all  $u$ . We achieve this in practice by simply clamping the values to 1.

### 3.2.4 Overlap

In settings where nodes may belong to more than one community, alternative assumptions are required [3, 4]. Consider interaction-specific community memberships [3], where the community of a node varies depending on the edge it is observed in. To understand this model, we define an expanded membership matrix  $\mathbf{Z}^* \in \{0, 1\}^{n \times n \times k}$  such that  $\mathbf{Z}^*_{uvi} = 1$  if node  $u$  is a member of community  $i$  when it transmits an edge to node  $v$ . The expected value of the dyad  $(u, v)$  is therefore  $\mathbf{Z}^*_{uv} \boldsymbol{\Theta} \mathbf{Z}^*_{vu}$ .

To reduce the dimensionality of the expanded model, let  $\mathbf{P} \in [0, 1]^{n \times k}$  be the collapsed membership matrix. This matrix is a summary of overlapping community memberships, defined as the degree-normalized sum over the second axis of the expanded membership matrix:  $\mathbf{P}_u = \mathbf{d}_u^{-1} \sum_{v=1}^n \mathbf{Z}^*_{uv}$ . In the overlapping setting,  $\mathbf{P}_u$  is interpreted as the frequency of node  $u$ 's membership in each of the  $k$  communities [4]. Here, the expected value of the dyad  $(u, v)$  is  $\mathbf{P}'_u \boldsymbol{\Theta} \mathbf{P}_v$ . In the non-overlapping setting, the collapsed membership matrix  $\mathbf{P}$  is equivalent to the partition matrix  $\mathbf{Z}$ . In both cases,  $\sum_{i=1}^k \mathbf{P}_{ui} = 1$ . The collapsed membership matrix will be relevant to neural network optimization.

### 3.3 Graph Neural Networks

With estimates of  $\Theta$  and  $\phi$  in place, we now turn to estimating the partition matrix. To begin, note that none of the log-likelihood functions described above are differentiable with respect to  $\mathbf{Z}$ , as it is a collection of discrete vectors. Because of this, gradient-based optimization methods are unavailable and Monte Carlo methods are commonly used to find the likelihood-maximizing partition [7, 20].

#### 3.3.1 Parameter Specification

Neural networks are generally optimized with some variation of the gradient descent algorithm. For gradient descent to be applicable in the SBM setting, an estimate of  $\mathbf{Z}$  that is differentiable with respect to the neural network parameters is required.

We consider the collapsed membership matrix  $\mathbf{P}$ , which is useful for both the standard and overlapping settings because it is a generalization of  $\mathbf{Z}$ . The choice of  $\mathbf{P}$  is convenient, as an estimate can be obtained by differentiable functions such as Softmax applied to GNN embeddings. Therefore, the partition estimator takes the following form:

$$\hat{\mathbf{P}} = \text{Softmax}(\text{GNN}(\mathcal{G}, \mathbf{X}; \mathbf{W})) \quad (7)$$

where GNN is any standard graph neural network with parameters  $\mathbf{W}$ . Here,  $\hat{\mathbf{P}}$  is considered a relaxation of the partition matrix  $\mathbf{Z}$  to a soft partition.

It should be noted that the output dimension of the GNN is the assumed number of communities  $k$  in the graph of interest. In practice, the exact number of communities need not be known a priori. Our experiments suggest that setting the output dimension to a reasonable overestimate typically allows the model to learn the optimal number of communities (see Appendix D). Thus,  $k$  can be inferred as the number of unique elements in  $\mathcal{K} = \{\arg \max \hat{\mathbf{P}}_u : u \in \mathcal{V}\}$ . This is convenient for real-world graphs where the true number of communities may be unknown.

#### 3.3.2 Loss Functions

The loss function associated with the Poisson model is the negative of the log-likelihood function in equation 2, with  $\mathbf{Z}$  replaced by  $\hat{\mathbf{P}}$  (equation 7) and  $\Theta$  replaced  $\hat{\Theta}(\mathcal{G}, \hat{\mathbf{P}})$ . That is,

$$J_P(\mathbf{W}) = - \sum_{u,v} [\mathbf{A}_{uv} \ln(\hat{\mathbf{P}}'_u \hat{\Theta} \hat{\mathbf{P}}_v) - \hat{\mathbf{P}}'_u \hat{\Theta} \hat{\mathbf{P}}_v]. \quad (8)$$

For brevity, the scalar  $\hat{\pi}_{uv} = \hat{\mathbf{P}}'_u \hat{\Theta} \hat{\mathbf{P}}_v$  is used for the remainder of this section. The loss function associated with the Bernoulli model is derived in the same way from equation 4, resulting in

$$J_B(\mathbf{W}) = - \sum_{u,v} [\mathbf{A}_{uv} \ln(\hat{\pi}_{uv}) + (1 - \mathbf{A}_{uv}) \ln(1 - \hat{\pi}_{uv})]. \quad (9)$$

Both losses are expressed as functions of the GNN parameters,  $\mathbf{W}$ .

Recall that  $\mathbf{A}_{uv}$  is equal to 1 if  $(u, v) \in \mathcal{E}$  and 0 if  $(u, v) \notin \mathcal{E}$ . Therefore, the loss function can be broken out into a summation over the positive edge set  $\mathcal{E}$  and the negative edge set  $\mathcal{N} = \mathcal{D} \setminus \mathcal{E}$ . Doing so highlights the difference in cardinality of both sets. Often, real-world graphs are highly sparse, meaning that  $|\mathcal{N}| \gg |\mathcal{E}|$ . Such imbalance can be problematic for optimizing a GNN. A common approach to address this is to under-sample [32] the majority class (usually the negative edge set) [9, 23, 33]. With this sampling approach, the loss functions are rewritten

$$J_P(\mathbf{W}) = - \sum_{u,v \in \mathcal{E}} [\ln(\hat{\pi}_{uv}) - \hat{\pi}_{uv}] + \sum_{u,v \notin \mathcal{E}} \hat{\pi}_{uv} \quad (10)$$

$$\approx - \sum_{u,v \sim P_{\mathcal{E}}} [\ln(\hat{\pi}_{uv}) - \hat{\pi}_{uv}] + \eta^{-1} \sum_{u,v \sim P_{\mathcal{N}}} \hat{\pi}_{uv} \quad (11)$$

$$J_B(\mathbf{W}) = - \sum_{u,v \in \mathcal{E}} \ln(\hat{\pi}_{uv}) - \sum_{u,v \notin \mathcal{E}} \ln(1 - \hat{\pi}_{uv}) \quad (12)$$

$$\approx - \sum_{u,v \sim P_{\mathcal{E}}} \ln(\hat{\pi}_{uv}) - \eta^{-1} \sum_{u,v \sim P_{\mathcal{N}}} \ln(1 - \hat{\pi}_{uv}). \quad (13)$$

The summations over  $P_{\mathcal{E}}$  and  $P_{\mathcal{N}}$  represent uniform samples from  $\mathcal{E}$  and  $\mathcal{N}$  with  $\eta$  as a scaling constant. For our experiments, all  $m$  positive edges are drawn deterministically and  $\eta m$  samples from the negative set are drawn randomly at each training step.

The degree-corrected versions of both models are achieved by including  $\hat{\phi}(\mathcal{G}, \hat{\mathbf{P}})$  in each loss function. Thus, the degree-corrected (DC) loss functions  $J_{P\text{-}DC}$  and  $J_{B\text{-}DC}$  are derived by substituting  $\hat{\pi}_{uv} = \hat{\phi}_u \hat{\phi}_v \hat{\mathbf{P}}'_u \hat{\Theta} \hat{\mathbf{P}}_v$  for  $\hat{\pi}_{uv}$  in the above equations.

### 3.3.3 Graph Matching

Another objective function is motivated by the Graph Matching problem [19, 34]. Graph Matching refers to the (approximate or exact) alignment of nodes across two graphs of possibly different sizes. Node alignment is usually defined by some real-world mechanism. Rather than matching two arbitrary graphs, our approach involves matching a graph to its community representation.

Let the block graph  $\mathcal{G}_{\Theta} = (\mathcal{V}_{\Theta}, \mathcal{E}_{\Theta})$  be defined with its (weighted) adjacency structure given by  $\Theta$ . Each of its  $k$  nodes is a subset of nodes from  $\mathcal{V}$ ; that is, node  $i \in \mathcal{V}_{\Theta}$  corresponds to the set  $\{u \in \mathcal{V} : \mathbf{Z}_{ui} = 1\}$ . Thus,  $\mathbf{Z}$  is considered a mapping matrix which transforms the block matrix to the expectation of  $\mathbf{A}$ ; that is,  $\mathbf{Z}\Theta\mathbf{Z}' \mapsto \mathbb{E}(\mathbf{A})$ . Finding the optimal mapping matrix  $\tilde{\mathbf{Z}}$  is the optimization problem

$$\underset{\tilde{\mathbf{Z}} \in \mathcal{Z}}{\operatorname{argmin}} \|\mathbf{A} - \tilde{\mathbf{Z}}\Theta\tilde{\mathbf{Z}}'\|_F = \underset{\tilde{\mathbf{Z}} \in \mathcal{Z}}{\operatorname{argmin}} -\operatorname{tr}(\mathbf{A}'\tilde{\mathbf{Z}}\Theta\tilde{\mathbf{Z}}') \quad (14)$$

where  $\mathcal{Z} = \{\tilde{\mathbf{Z}} \in \{0, 1\}^{n \times k} : \sum_{i=1}^k \tilde{\mathbf{Z}}_{ui} = 1 \text{ for all } u \in \mathcal{V}\}$  and  $\operatorname{tr}(\cdot)$  is the matrix trace.

Because  $\Theta$  is an unknown parameter in the SBM formulation, we use its MLE:  $\hat{\Theta}(\mathcal{G}, \tilde{\mathbf{Z}})$ . Thus, we are seeking to find the (inverse) mapping of  $\mathcal{G}$  to its estimated stochastic block representation  $\mathcal{G}_{\hat{\Theta}}$  by minimizing the quantity  $-\operatorname{tr}(\mathbf{A}'\tilde{\mathbf{Z}}\hat{\Theta}\tilde{\mathbf{Z}}')$ . A proof of this statement is provided in Appendix F.2.

The optimization problem above is addressed with GNNs by substituting  $\tilde{\mathbf{Z}}$  with the predicted membership matrix  $\hat{\mathbf{P}}$  from equation 7 and using the block matrix estimate  $\hat{\Theta}(\mathcal{G}, \hat{\mathbf{P}})$ . The loss function in this setting is

$$J_{\text{Match}}(\mathbf{W}) = -\operatorname{tr}(\mathbf{A}'\hat{\mathbf{P}}\hat{\Theta}\hat{\mathbf{P}}') \quad (15)$$

and is again expressed as a function of the GNN parameters,  $\mathbf{W}$ . Note that  $\mathbf{A}$  can be represented as a compressed sparse matrix, making the product  $\mathbf{A}'\hat{\mathbf{P}}$  relatively efficient to compute. Another gain in computational efficiency comes from the relation  $\operatorname{tr}(\mathbf{A}'\hat{\mathbf{P}}\hat{\Theta}\hat{\mathbf{P}}') = \sum_{u=1}^n \sum_{i=1}^k (\mathbf{A}'\hat{\mathbf{P}})_{ui} (\hat{\Theta}\hat{\mathbf{P}}')_{iu}$ . The use of sparse matrix multiplication results in a significant speedup compared to the edge sampling loss functions, as will be shown empirically in Section 4.

### 3.3.4 Regularization

The final component of the SBM loss framework is a regularization term. This regularization is meant to encourage the model to distribute the predicted partition across a sufficient number of communities and to ensure that the predicted partition is assortative (i.e., edges occur mostly between nodes of the same community [7]).

We propose a term that helps minimize the distance of the structure matrix diagonal from unity. The regularized form of an arbitrary loss function  $J$  is

$$J^*(\mathbf{W}) = m^{-1} J(\mathbf{W}) + \alpha \|\mathbf{1}_k - \hat{\theta}_d\|_F \quad (16)$$

where  $\hat{\theta}_d$  is the diagonal of  $\hat{\Theta}$ ,  $\mathbf{1}_k$  is a  $k$ -vector of ones, and  $\|\cdot\|_F$  is the Frobenius norm. The hyper-parameter  $\alpha$  controls the strength of the regularization.

This is analogous to the ‘‘collapse regularization’’ term introduced in [10] (see Appendix B.3). Both functions help to avoid trivial solutions to the optimization problem which arise when all nodes are assigned to one partition. Also, both are conveniently bounded in the interval  $[0, \sqrt{k}]$ .

## 4 Experiments

The proposed loss functions are evaluated based on the performance of the neural networks they are optimized with respect to. Community detection performance is measured by Normalized Mutual



Information (NMI) and Pairwise-F1 (PF1) scores; where relevant, the overlapping variants [35] are used. Evaluation is carried out on a variety of synthetic and real-world datasets.

For baseline comparison, the loss functions used in several alternative approaches are also evaluated: Neural Overlapping Community Detection  $J_{\text{NOCD}}$  [9], Deep Modularity Network  $J_{\text{DMoN}}$  [10], Minimum-Cut Pooling  $J_{\text{MCP}}$  [8], Markov Stability  $J_{\text{CDMG}}$  [12], and link prediction  $J_{\text{LP}}$  [23]. These loss functions are described in greater detail in Appendix B. To ensure a fair comparison, a standard Graph Neural Network architecture is used for all experiments regardless of what was used in the original papers.

The architecture of choice is a two-layer Graph Convolutional Network (GCN) [24]. Implementation details are provided in Appendix A. The model configuration and hyperparameters were determined based on preliminary experimentation and are kept the same for all loss functions considered. An examination of (SBM-specific) hyperparameter sensitivity is provided in Appendix D. We also compare against the architecture described in [10]; these results, as well as the aggregated results of several other architectures, are provided in Appendix E.

The models are tested on a number of synthetic datasets generated to demonstrate both the overlapping and non-overlapping settings. We also test on several real-world benchmarks for overlapping and non-overlapping community detection. Further detail for all datasets is provided in Appendix C.

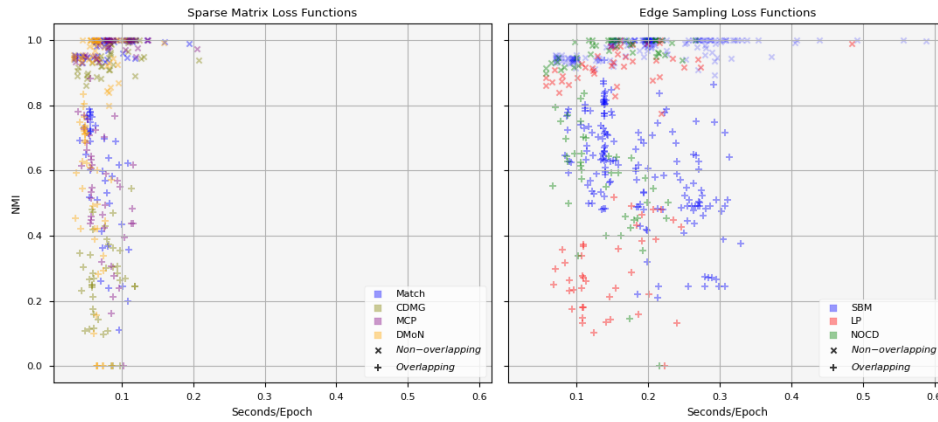
#### 4.1 Synthetic Results

	SBM-4		SBM-8		SBM-16		SBM-32		Avg.	
	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1
$J_{\text{MCP}}$	<u>99.2</u>	<u>99.2</u>	99.3	99.0	98.6	95.8	<u>94.5</u>	<b>78.8</b>	<u>97.9</u>	<u>93.2</u>
$J_{\text{CDMG}}$	98.0	97.7	97.4	94.5	93.9	81.0	89.3	63.2	94.7	84.1
$J_{\text{DMoN}}$	89.3	87.6	98.7	98.1	<b>99.7</b>	<b>99.3</b>	94.5	<u>78.1</u>	95.5	90.8
$J_{\text{LP}}$	93.0	93.1	95.6	91.6	91.4	76.0	87.9	61.6	92.0	80.5
$J_{\text{NOCD}}$	97.5	97.9	99.2	98.7	98.9	96.4	92.3	71.5	97.0	91.1
$J_{\text{B}}$	<b>99.3</b>	<b>99.5</b>	99.3	98.5	98.7	95.2	94.1	75.1	97.9	92.1
$J_{\text{P}}$	98.3	98.3	<b>99.8</b>	<b>99.7</b>	99.2	97.0	93.7	74.8	97.7	92.4
$J_{\text{B-DC}}$	97.2	97.0	98.8	97.9	98.4	95.5	93.7	74.6	97.0	91.3
$J_{\text{P-DC}}$	98.5	98.6	<u>99.4</u>	<u>99.2</u>	<u>99.7</u>	<u>99.1</u>	94.3	75.9	<b>98.0</b>	<b>93.2</b>
$J_{\text{Match}}$	98.4	98.5	98.1	96.8	99.4	98.8	<b>94.7</b>	76.9	97.7	92.7

**Table 1:** Community detection performance on synthetic data with *non-overlapping* communities using GCN. Results are averaged over ten trials. The best scores (NMI and PF1) for each dataset are in bold and second best are underlined.

	OSBM-4		OSBM-8		OSBM-16		OSBM-32		Avg.	
	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1
$J_{\text{MCP}}$	38.7	86.3	45.8	83.4	57.8	78.8	<u>69.7</u>	72.2	53.0	80.2
$J_{\text{CDMG}}$	27.2	69.5	26.2	57.1	28.3	57.1	23.5	55.3	26.3	59.7
$J_{\text{DMoN}}$	19.0	53.5	37.1	75.8	<b>74.6</b>	<b>90.6</b>	61.9	70.8	48.1	72.7
$J_{\text{LP}}$	34.7	57.6	29.3	68.2	23.7	62.2	26.2	52.7	28.5	60.2
$J_{\text{NOCD}}$	38.3	57.6	55.0	77.3	68.8	76.3	63.2	66.3	56.3	69.4
$J_{\text{B}}$	<b>56.8</b>	84.8	56.5	<u>85.2</u>	73.3	82.8	<b>70.5</b>	<b>73.7</b>	<b>64.3</b>	<u>81.6</u>
$J_{\text{P}}$	46.8	90.7	51.8	<b>89.7</b>	71.9	81.2	61.6	65.8	58.0	<b>81.8</b>
$J_{\text{B-DC}}$	<u>50.8</u>	<u>93.4</u>	<b>58.9</b>	84.8	<u>73.8</u>	77.5	67.2	69.4	<u>62.7</u>	81.3
$J_{\text{P-DC}}$	37.9	<b>97.2</b>	<u>56.7</u>	82.5	68.7	76.5	65.4	70.0	57.1	81.6
$J_{\text{Match}}$	36.0	85.0	54.2	75.4	72.2	<u>86.4</u>	69.3	<u>73.4</u>	57.9	80.0

**Table 2:** Community detection performance on synthetic data with *overlapping* communities using GCN. Results are averaged over ten trials. The best scores (overlapping NMI and PF1) for each dataset are in bold and second best are underlined.



**Figure 1:** Accuracy (NMI) vs. training time (seconds-per-epoch) for each loss function on synthetic data. The left panel shows the loss functions that exploit matrix sparsity; the right panel shows loss functions that use edge sampling. Each SBM loss variant is labeled “SBM” in the right panel. Results are marked by  $\times$  for non-overlapping datasets and by  $+$  for overlapping datasets.

We first evaluate the community detection performance of each loss function on synthetic graphs without community overlap. Community predictions are the row-wise argmax of the model output and scoring is done with standard NMI and PF1. Results are shown in table 1. The top half of the table are comparison loss functions and the bottom half are our proposed SBM loss functions. The Bernoulli loss function performs best in terms of both NMI and PF1 on average. The MCP loss is also competitive.

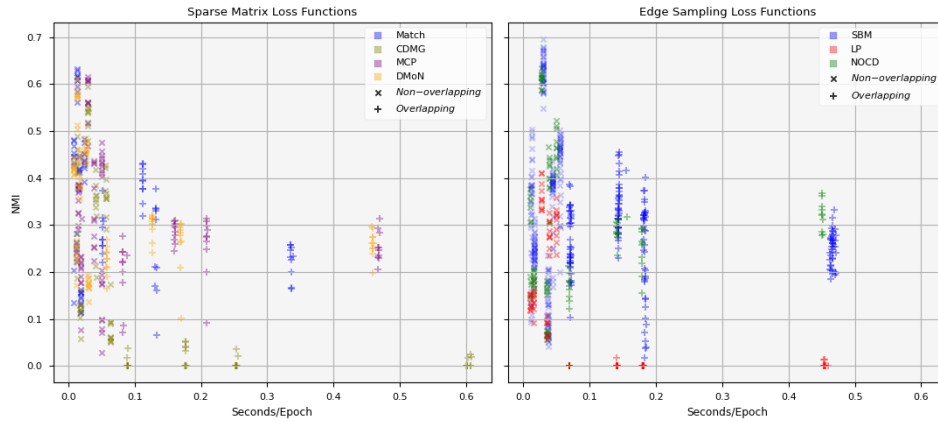
For overlapping community detection, a threshold is applied to model outputs and the arguments exceeding that threshold are the predicted overlapping communities for each node. The threshold for each model is chosen to maximize accuracy on the validation set. For evaluation, overlapping NMI and PF1 are used. Results are shown in table 2. The SBM loss functions generally outperform most comparison losses, with the DMoN loss function being competitive. On average, the Bernoulli variants perform best in overlapping NMI and PF1.

Figure 1 shows accuracy (measured as NMI) plotted against training time (measured in seconds-per-epoch) for each loss function. The left-hand panel shows the loss functions that exploit matrix sparsity:  $J_{\text{Match}}$ ,  $J_{\text{CDMG}}$ ,  $J_{\text{MCP}}$ , and  $J_{\text{DMoN}}$ . The right-hand panel shows the loss functions that employ (negative) edge sampling:  $J_{\text{B}}$ ,  $J_{\text{P}}$ ,  $J_{\text{LP}}$ , and  $J_{\text{NOCD}}$ . The Bernoulli and Poisson SBM loss functions (and their degree-corrected variants) have all been labeled “SBM” for simplicity in the right-hand panel. Note how the loss functions that exploit matrix sparsity are significantly faster than those that use edge sampling. The DMoN loss function is the fastest and our Graph Matching loss function is competitive in terms of both speed and accuracy.

## 4.2 Real-World Results

	Cora		Citeseer		Pubmed		Wiki		ACB-Comp		ACB-Photo		Avg.	
	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1
$J_{\text{MCP}}$	34.2	22.6	19.4	16.4	13.0	19.5	33.1	26.7	43.4	45.8	57.1	52.6	33.4	30.6
$J_{\text{CDMG}}$	37.0	29.3	11.8	14.9	7.9	25.9	32.9	34.1	36.7	43.6	53.8	51.2	30.0	33.2
$J_{\text{DMoN}}$	<b>42.3</b>	25.9	23.9	19.1	<b>17.4</b>	20.4	<b>40.8</b>	36.0	<u>45.9</u>	34.6	57.0	47.1	<b>37.9</b>	30.5
$J_{\text{LP}}$	13.7	16.8	14.4	16.0	7.7	12.0	28.4	23.0	30.4	26.4	36.4	28.8	21.8	20.5
$J_{\text{NOCD}}$	27.3	19.8	18.6	19.8	12.7	19.1	40.6	34.6	<b>46.6</b>	41.7	61.1	56.9	34.5	32.0
$J_{\text{B}}$	29.9	27.9	24.8	<b>31.6</b>	12.4	<u>33.1</u>	36.5	35.8	42.8	<b>52.2</b>	<u>64.2</u>	<u>60.8</u>	35.1	<u>40.2</u>
$J_{\text{P}}$	28.2	25.8	22.8	<u>30.0</u>	13.5	<b>36.6</b>	36.2	36.5	40.2	<u>51.3</u>	<b>64.7</b>	<b>63.6</b>	34.3	<b>40.6</b>
$J_{\text{B-DC}}$	29.5	26.7	22.9	28.5	13.3	29.0	38.9	<b>38.4</b>	42.8	47.9	62.3	60.2	34.9	38.5
$J_{\text{P-DC}}$	<u>41.3</u>	<b>37.4</b>	<b>27.9</b>	28.2	10.8	28.3	33.6	31.1	45.8	32.6	60.1	53.4	36.6	35.2
$J_{\text{Match}}$	39.7	<u>30.3</u>	<u>25.0</u>	20.2	<u>14.2</u>	20.5	<u>40.8</u>	36.3	44.1	37.8	60.7	53.5	<u>37.4</u>	33.1

**Table 3:** Community detection performance on real-world data with *non-overlapping* communities using GCN. Results are averaged over ten trials. The best scores (NMI and PF1) for each dataset are in bold and second best are underlined.



**Figure 2:** Accuracy (NMI) vs. training time (seconds-per-epoch) for each loss function on real-world data. The left panel shows the loss functions that exploit matrix sparsity; the right panel shows loss functions that use edge sampling. Each SBM loss variant is labeled “SBM” in the right panel. Results are marked by  $\times$  for non-overlapping datasets and by  $+$  for overlapping datasets.



	MAG-Chem		MAG-CS		MAG-Eng		MAG-Med		Avg.	
	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1
$J_{MCP}$	25.6	28.9	28.5	27.6	19.7	29.9	25.5	23.3	24.8	27.4
$J_{CDMG}$	0.6	<u>34.6</u>	2.2	28.7	0.6	<b>48.5</b>	0.6	23.4	1.0	33.8
$J_{DMoN}$	25.9	28.4	29.8	29.5	22.4	26.8	26.0	23.6	26.0	27.1
$J_{LP}$	0.0	9.7	0.2	10.6	0.0	10.0	0.3	17.1	0.1	11.8
$J_{NOCD}$	24.5	27.6	28.9	29.1	15.1	29.7	<b>32.0</b>	24.0	25.1	27.6
$J_B$	19.7	30.9	29.7	34.8	<b>29.4</b>	47.1	24.9	<u>28.9</u>	25.9	35.4
$J_P$	18.5	28.2	30.3	36.5	<u>28.3</u>	47.9	24.0	26.9	25.3	34.9
$J_{B-DC}$	<u>26.4</u>	32.8	35.7	<u>44.2</u>	25.7	<u>48.3</u>	<u>26.9</u>	<b>30.1</b>	<u>28.7</u>	<b>38.9</b>
$J_{P-DC}$	<b>30.8</b>	<b>39.4</b>	<b>40.5</b>	<b>44.9</b>	22.1	36.1	24.2	27.8	<b>29.4</b>	<u>37.0</u>
$J_{Match}$	25.0	32.9	<u>38.9</u>	42.1	28.0	37.4	22.3	25.3	28.6	34.4

**Table 4:** Community detection performance on real-world data with *overlapping* communities using GCN. Results are averaged over ten trials. The best scores (overlapping NMI and PF1) for each dataset are in bold and second best are underlined.

NMI and PF1. The Graph Matching loss demonstrates better performance (relatively) on real-world graphs than on synthetic graphs. This could suggest that it is better suited for sparse graphs.

Figure 2 shows accuracy plotted against training time for each loss function. Notice how the gain in speed of the sparse matrix losses over the edge sampling losses is not as pronounced as in the synthetic datasets. This is likely due to higher feature dimensions in the real-world graphs.

Results that show significant disagreement between the NMI and PF1 scores could be attributed to the sensitivity of NMI to community size imbalance. Many of the real-world graphs in this study are made up of communities that vary considerably in size. Such imbalance can cause misleading results when comparing the NMI of multiple predicted partitions [36, 37]. On the other hand, PF1 is an aggregation of pairwise precision and pairwise recall, making it preferable when evaluating partitions on imbalanced data.

### 4.3 Discussion

Figures 1 and 2 suggest that, in general, the loss functions that exploit matrix sparsity are faster and just as accurate as those that use edge sampling. It should be noted that most losses are computed more quickly on non-overlapping datasets. This is simply because the (real-world) overlapping datasets happen to contain more nodes and edges than the non-overlapping datasets (see table E6). In fact, almost every instance where “Seconds/Epoch” is greater than 0.25 corresponds to the MAG-Med dataset, which has the greatest number of nodes and edges.

Another general trend is that the overlapping datasets usually produce lower accuracy scores than the non-overlapping ones. This limitation is observed not only for the SBM-based loss functions, but also for the comparison losses. Further effort to adapt Overlapping and Mixed Membership SBMs to GNN loss functions is a direction for future work.

It is interesting to note that while all the graphs used in evaluation exhibit skewed degree distributions, there is no clear indication that the degree-corrected SBM variants outperform their degree-free counterparts. It is possible that GNNs are powerful enough to effectively fit these networks without the help of degree-correction parameters. It is also possible that the degree-correction parameters require a more effective estimation approach, rather than simply the plug-in MLE.

A key limitation of the proposed framework is that it is only designed for graphs with binary adjacency matrices. For multi-graphs or weighted graphs, the Bernoulli and partial Poisson loss functions will not be applicable. Our Poisson variant can easily be extended to support multi-graphs (see Appendix B.6), which allows the structure matrix to take values greater than 1. However, such a structure matrix does not conform with the regularization term in equation 16, which penalizes matrices with diagonals that are far from unity.

Another limitation is the stability of the structure matrix MLE. Note that the gradient  $\frac{\partial \hat{\Theta}_{ij}}{\partial \mathbf{n}_i} = -\frac{\mathbf{M}_{ij}}{\mathbf{n}_i^2 \mathbf{n}_j}$  has a limit of 0 as  $\mathbf{n}_i \rightarrow \infty$  and approaches  $-\infty$  as  $\mathbf{n}_i \rightarrow 0$ . This means that  $\hat{\Theta}$  is less stable for communities that are small in size (or unused). This also suggests that SBM losses are more suitable

The community detection performance of each loss function is evaluated on real-world data in the same way as done on synthetic data. We first consider graphs without community overlap; results are shown in table 3. The SBM loss functions outperform many of the comparison losses in terms of PF1 and NMI. The DMoN loss does best in average NMI and the Poisson loss does best in average PF1. The NOCD loss is also competitive.

Overlapping community detection performance is evaluated in the same way as the synthetic datasets. Results are shown in table 4. The SBM loss functions again outperform most comparison losses, with the degree-corrected variants doing best on average in terms of overlapping

for networks with more nodes per partition. A direction for future work is identifying a more stable estimator of the structure matrix.

## 5 Conclusion

In this paper, a collection of loss functions are derived from Stochastic Block Model likelihood functions. These loss functions are configurable, fully differentiable, and theoretically grounded. They show strong performance in unsupervised training of Graph Neural Networks for community detection. An additional loss function is adapted from the Graph Matching problem and shows significant speed improvements. The proposed framework is subjected to extensive evaluation and shows competitive or improved performance against state-of-the-art loss functions.

## References

- [1] Krzysztof Nowicki and Tom A. B. Snijders. Estimation and prediction for stochastic block-structures. *Journal of the American Statistical Association*, 96:1077 – 1087, 2001. URL <https://api.semanticscholar.org/CorpusID:9478789>. 1, 2
- [2] Brian Karrer and M. E. J. Newman. Stochastic blockmodels and community structure in networks. *Phys. Rev. E*, 83:016107, 1 2011. doi: 10.1103/PhysRevE.83.016107. URL <https://link.aps.org/doi/10.1103/PhysRevE.83.016107>. 1, 2, 3, 4
- [3] Edoardo M. Airolidi, David M. Blei, Stephen E. Fienberg, and Eric P. Xing. Mixed membership stochastic blockmodels. *Journal of Machine Learning Research*, 9(65):1981–2014, 2008. URL <http://jmlr.org/papers/v9/airolidi08a.html>. 2, 4
- [4] Pierre Latouche, Etienne Birmelé, and Christophe Ambroise. Overlapping stochastic block models with application to the french political blogosphere. *The Annals of Applied Statistics*, 5(1), March 2011. ISSN 1932-6157. doi: 10.1214/10-aos382. URL <http://dx.doi.org/10.1214/10-AOS382>. 2, 4
- [5] Tiago P. Peixoto. Nonparametric bayesian inference of the microcanonical stochastic block model. *Physical Review E*, 95(1), January 2017. ISSN 2470-0053. doi: 10.1103/physreve.95.012317. URL <http://dx.doi.org/10.1103/PhysRevE.95.012317>. 2
- [6] Yash Deshpande, Andrea Montanari, Elchanan Mossel, and Subhabrata Sen. Contextual stochastic block models, 2018. URL <https://arxiv.org/abs/1807.09596>. 2
- [7] Clement Lee and Darren J. Wilkinson. A review of stochastic block models and extensions for graph clustering. *Applied Network Science*, 4(1), December 2019. ISSN 2364-8228. doi: 10.1007/s41109-019-0232-2. URL <http://dx.doi.org/10.1007/s41109-019-0232-2>. 1, 2, 3, 4, 5, 6
- [8] Filippo Maria Bianchi, Daniele Grattarola, and Cesare Alippi. Spectral clustering with graph neural networks for graph pooling. In *Proceedings of the 37th International Conference on Machine Learning*, ICML’20. JMLR.org, 2020. 1, 2, 3, 7, 14, 15, 17
- [9] Oleksandr Shchur and Stephan Günnemann. Overlapping community detection with graph neural networks. *Deep Learning on Graphs Workshop, KDD*, abs/1909.12201, 2019. URL <https://api.semanticscholar.org/CorpusID:88492570>. 2, 3, 5, 7, 13, 16
- [10] Anton Tsitsulin, John Palowitch, Bryan Perozzi, and Emmanuel Müller. Graph clustering with graph neural networks. *J. Mach. Learn. Res.*, 24(1), 3 2024. ISSN 1532-4435. 3, 6, 7, 14, 15, 17
- [11] Rex Ying, Jiaxuan You, Christopher Morris, Xiang Ren, William L. Hamilton, and Jure Leskovec. Hierarchical graph representation learning with differentiable pooling. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, NIPS’18, page 4805–4815, Red Hook, NY, USA, 2018. Curran Associates Inc. 2, 17
- [12] Shunjie Yuan, Chao Wang, Qi Jiang, and Jianfeng Ma. Community detection with graph neural network using markov stability. In *2022 International Conference on Artificial Intelligence in Information and Communication (ICAIIIC)*, pages 437–442, 2022. doi: 10.1109/ICAIIIC54071.2022.9722614. 1, 3, 7, 14
- [13] Eli Chien, Jianhao Peng, Pan Li, and Olgica Milenkovic. Adaptive universal generalized pagerank graph neural network, 2021. URL <https://arxiv.org/abs/2006.07988>. 1, 2

- [14] O Duranthon and L Zdeborová. Neural-prior stochastic block model. *Machine Learning: Science and Technology*, 4(3):035017, August 2023. ISSN 2632-2153. doi: 10.1088/2632-2153/ace60f. URL <http://dx.doi.org/10.1088/2632-2153/ace60f>. 2
- [15] Kimon Fountoulakis, Dake He, Silvio Lattanzi, Bryan Perozzi, Anton Tsitsulin, and Shenghao Yang. On classification thresholds for graph attention with edge features, 2022. URL <https://arxiv.org/abs/2210.10014>. 2
- [16] Nikhil Mehta, Lawrence Carin Duke, and Piyush Rai. Stochastic blockmodels meet graph neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 4466–4474. PMLR, 6 2019. URL <https://proceedings.mlr.press/v97/mehta19a.html>. 2
- [17] Cheng Shi, Liming Pan, Hong Hu, and Ivan Dokmanić. Homophily modulates double descent generalization in graph convolution networks, 2024. URL <https://arxiv.org/abs/2212.13069>. 1, 2
- [18] Zheng Chen, Xinli Yu, Yuan Ling, and Xiaohua Hu. Neural stochastic block model & scalable community-based graph learning, 2020. URL <https://arxiv.org/abs/2005.07855>. 1, 2
- [19] S. Umeyama. An eigendecomposition approach to weighted graph matching problems. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 10(5):695–703, 1988. doi: 10.1109/34.6778. 2, 6
- [20] Emmanuel Abbe. Community detection and stochastic block models: Recent developments. *Journal of Machine Learning Research*, 18(177):1–86, 2018. URL <http://jmlr.org/papers/v18/16-480.html>. 2, 5
- [21] M. E. J. Newman and M. Girvan. Finding and evaluating community structure in networks. *Physical Review E*, 69(2), February 2004. ISSN 1550-2376. doi: 10.1103/physreve.69.026113. URL <http://dx.doi.org/10.1103/PhysRevE.69.026113>. 2, 14
- [22] Santo Fortunato. Community detection in graphs. *Physics Reports*, 486(3–5):75–174, February 2010. ISSN 0370-1573. doi: 10.1016/j.physrep.2009.11.002. URL <http://dx.doi.org/10.1016/j.physrep.2009.11.002>. 2
- [23] William L. Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS’17, page 1025–1035, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964. 2, 5, 7, 14
- [24] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks, 2017. 7, 13, 17
- [25] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks, 2018. 17
- [26] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks?, 2019. 2, 17
- [27] Jaewon Yang and Jure Leskovec. Overlapping community detection at scale: a nonnegative matrix factorization approach. In *Proceedings of the Sixth ACM International Conference on Web Search and Data Mining*, WSDM ’13, page 587–596, New York, NY, USA, 2013. Association for Computing Machinery. ISBN 9781450318693. doi: 10.1145/2433396.2433471. URL <https://doi.org/10.1145/2433396.2433471>. 3, 13
- [28] Mingyuan Zhou. Infinite Edge Partition Models for Overlapping Community Detection and Link Prediction. In Guy Lebanon and S. V. N. Vishwanathan, editors, *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, volume 38 of *Proceedings of Machine Learning Research*, pages 1135–1143, San Diego, California, USA, 5 2015. PMLR. URL <https://proceedings.mlr.press/v38/zhou15a.html>. 3
- [29] J. C. Delvenne, S. N. Yaliraki, and M. Barahona. Stability of graph communities across time scales, 2009. URL <https://arxiv.org/abs/0812.1811>. 3, 14
- [30] Jean-Charles Delvenne, Michael T. Schaub, Sophia N. Yaliraki, and Mauricio Barahona. *The Stability of a Graph Partition: A Dynamics-Based Framework for Community Detection*, page 221–242. Springer New York, 2013. ISBN 9781461467298. doi: 10.1007/978-1-4614-6729-8\_11. URL [http://dx.doi.org/10.1007/978-1-4614-6729-8\\_11](http://dx.doi.org/10.1007/978-1-4614-6729-8_11).

- [31] Renaud Lambiotte, Jean-Charles Delvenne, and Mauricio Barahona. Random walks, markov processes and the multiscale modular organization of complex networks. *IEEE Transactions on Network Science and Engineering*, 1(2):76–90, 2014. doi: 10.1109/TNSE.2015.2391998. 3, 14
- [32] Haibo He and Eduardo A. Garcia. Learning from imbalanced data. *IEEE Transactions on Knowledge and Data Engineering*, 21(9):1263–1284, 2009. doi: 10.1109/TKDE.2008.239. 5
- [33] Zhen Yang, Ming Ding, Chang Zhou, Hongxia Yang, Jingren Zhou, and Jie Tang. Understanding negative sampling in graph representation learning. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD '20*, page 1666–1676, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450379984. doi: 10.1145/3394486.3403218. URL <https://doi.org/10.1145/3394486.3403218>. 5, 14
- [34] Joshua T Vogelstein, John M Conroy, Vince Lyzinski, Louis J Podrazik, Steven G Kratzer, Eric T Harley, Donniell E Fishkind, R Jacob Vogelstein, and Carey E Priebe. Fast approximate quadratic programming for graph matching. *PLOS one*, 10(4):e0121002, 2015. 6
- [35] Aaron F. McDaid, Derek Greene, and Neil Hurley. Normalized mutual information to evaluate overlapping community finding algorithms, 2013. URL <https://arxiv.org/abs/1110.2515>. 7
- [36] Hanneke van der Hoef and Matthijs J Warrens. Understanding information theoretic measures for comparing clusterings. *Behaviormetrika*, 46(2):353–370, 2019. 9
- [37] Marcilio C.P. de Souto, André L.V. Coelho, Katti Faceli, Tiemi C. Sakata, Viviane Bonadia, and Ivan G. Costa. A comparison of external clustering evaluation indices in the context of imbalanced data sets. In *2012 Brazilian Symposium on Neural Networks*, pages 49–54, 2012. doi: 10.1109/SBRN.2012.25. 9
- [38] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL <https://arxiv.org/abs/1412.6980>. 13
- [39] Hidetaka Kamigaito and Katsuhiko Hayashi. Unified interpretation of softmax cross-entropy and negative sampling: With case study for knowledge graph embedding. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Association for Computational Linguistics, 2021. doi: 10.18653/v1/2021.acl-long.429. URL <http://dx.doi.org/10.18653/v1/2021.acl-long.429>. 14
- [40] Karol A. Bacik, Michael T. Schaub, Mariano Beguerisse-Díaz, Yazan N. Billeh, and Mauricio Barahona. Flow-based network analysis of the *caenorhabditis elegans* connectome. *PLOS Computational Biology*, 12(8):e1005055, 3 2016. ISSN 1553-7358. doi: 10.1371/journal.pcbi.1005055. URL <http://dx.doi.org/10.1371/journal.pcbi.1005055>. 14
- [41] Yu and Shi. Multiclass spectral clustering. In *Proceedings Ninth IEEE International Conference on Computer Vision*, pages 313–319 vol.1, 2003. doi: 10.1109/ICCV.2003.1238361. 14
- [42] Andrew McCallum, Kamal Nigam, Jason D. M. Rennie, and Kristie Seymore. Automating the construction of internet portals with machine learning. *Information Retrieval*, 3:127–163, 2000. URL <https://api.semanticscholar.org/CorpusID:349242>. 15
- [43] C. Lee Giles, Kurt D. Bollacker, and Steve Lawrence. Citeseer: an automatic citation indexing system. In *Proceedings of the Third ACM Conference on Digital Libraries, DL '98*, page 89–98, New York, NY, USA, 1998. Association for Computing Machinery. ISBN 0897919653. doi: 10.1145/276675.276685. URL <https://doi.org/10.1145/276675.276685>. 15
- [44] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Gallagher, and Tina Eliassi-Rad. Collective classification in network data. In *The AI Magazine*, 2008. URL <https://api.semanticscholar.org/CorpusID:62016134>. 15
- [45] Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. Pitfalls of graph neural network evaluation. *Relational Representation Learning Workshop, NeurIPS 2018*, 2018. 15

## Acknowledgements

This research was conducted under the Graph Artificial Intelligence initiative of the Leidos Office of Technology. Export approval number: 25-LEIDOS-0428-29484. *This document does not contain*

*export-controlled information as defined under the U.S. International Traffic in Arms Regulations or the U.S. Export Administration Regulations.*

## A Implementation

The neural network used for all experiments consists of two GNN layers with feature dropout, batch normalization, and ReLU activation in-between. The default GNN layer is GCN [24] unless otherwise specified. The dimension of the output layer is 25 for all experiments and the activation is determined by the loss function. The hidden dimension is 250 for synthetic graphs and 500 for real-world graphs. The dropout rate is 0.5.

Weight decay is applied with a strength of 0.0001. Adam optimization [38] is used with a learning rate of 0.0001. Training is carried out for a maximum of 500 epochs; early stopping is evaluated every 5 epochs and is engaged after 10 evaluations with no improvement. The SBM loss functions all take the form of equation 16 with regularization strength  $\alpha = 1.0$ . The collapse regularization strength for the DMoN loss is the same. For loss functions that support negative sampling, 3 negative edges are drawn for every positive edges, and balanced weighting is applied.

The common model/hyperparameter configuration was determined based on the performance of all models (considered simultaneously) during preliminary experimentation. Hyperparameters are kept the same for all loss functions considered rather than chosen for each loss separately; this is intended to allow us to study the effect of changing only the loss function.

In our experiments, nodes are split into train (60%), validation (20%), and test (20%) sets. To train the models, we use a “training subgraph” induced by nodes in the train and validation sets. Node embeddings are computed on just the training subgraph, then loss, early stopping, and thresholds are computed on the training/validation nodes separately. To evaluate the models, node embeddings are computed on the full graph (containing all train, validation, and test nodes), then metrics are computed on just the nodes in the test set. Training and evaluation is repeated over ten trials and the average test metrics are reported in this paper. In this way, the models’ ability to inductively generalize to nodes not seen during training is evaluated.

Experiments are done in Python 3.8 with DGL<sup>1</sup> and PyTorch<sup>2</sup>. Synthetic experiments are conducted with a 2.4 GHz Intel Core i9 processor and 32 GB of memory. Real-world experiments are conducted with a 2.5 GHz Intel(R) Xeon(R) Platinum 8259CL processor and 15 GB of memory and a Tesla T4 GPU and Cuda version 12.4.

## B Additional Objective Functions

### B.1 Bernoulli-Poisson

The Bernoulli-Poisson model, explored in [9, 27], assumes  $\mathbf{A}_{uv}|\mathbf{Z} \sim \text{Bern}(1 - e^{-\mathbf{Z}'_u \mathbf{Z}_v})$ . In [9], the “Neural Overlapping Community Detection” (NOCD) model is introduced, which derives a loss function from the Bernoulli-Poisson likelihood function. The authors propose

$$\begin{aligned} \hat{\mathbf{Z}} &= \text{ReLU}(\text{GNN}(\mathcal{G}, \mathbf{X}; \mathbf{W})) \\ J_{\text{NOCD}}(\mathbf{W}) &= - \sum_{u,v \sim P_{\mathcal{E}}} \ln(1 - \exp(-\hat{\mathbf{Z}}'_u \hat{\mathbf{Z}}_v)) + \eta^{-1} \sum_{u,v \sim P_{\mathcal{N}}} \hat{\mathbf{Z}}'_u \hat{\mathbf{Z}}_v. \end{aligned} \quad (17)$$

This is similar to the Poisson SBM (equation 11), with one difference being that there is no block matrix. Also, a complete derivation of the Poisson log-likelihood is used instead of a partial Poisson log-likelihood, as in equation 2.

### B.2 Link Prediction

A common objective function used for graph representation learning is cross-entropy-based link prediction loss. This model seeks to estimate the adjacency structure of a graph and is motivated by

<sup>1</sup><https://www.dgl.ai/>

<sup>2</sup><https://pytorch.org/>



the assumption  $\mathbf{A}_{uv}|\mathbf{Z} \sim \text{Bern}(\sigma(\mathbf{z}'_u \mathbf{z}_v))$ . One popular formulation [23, 33, 39] is

$$\begin{aligned} \hat{\mathbf{Z}} &= \text{GNN}(\mathcal{G}, \mathbf{X}; \mathbf{W}) \\ J_{\text{LP}}(\mathbf{W}) &= - \sum_{u,v \sim P_{\mathcal{E}}} \ln(\sigma(\hat{\mathbf{Z}}'_u \hat{\mathbf{Z}}_v)) - \eta^{-1} \sum_{u,v \sim P_{\mathcal{N}}} \ln(1 - \sigma(\hat{\mathbf{Z}}'_u \hat{\mathbf{Z}}_v)) \end{aligned} \quad (18)$$

where  $\sigma$  is the sigmoid function. This can be viewed as an analog to the Bernoulli SBM (equation 13) with the key difference being the absence of a block matrix.

### B.3 Modularity

In [10], it is proposed to directly maximize a graph partition quality metric. The approach, referred to as “Deep Modularity Network” (DMoN), is focused on modularity [21], defined as

$$Q = \frac{1}{2m} \text{tr} \left( \mathbf{Z}' \left( \mathbf{A} - \frac{\mathbf{d}\mathbf{d}'}{2m} \right) \mathbf{Z} \right). \quad (19)$$

The objective is to maximize modularity, or equivalently minimize the negative of modularity. Thus, the loss function is

$$\begin{aligned} \hat{\mathbf{Z}} &= \text{Softmax}(\text{GNN}(\mathcal{G}, \mathbf{X}; \mathbf{W})) \\ J_{\text{DMoN}}(\mathbf{W}) &= - \frac{1}{2m} \text{tr} \left( \hat{\mathbf{Z}}' \left( \mathbf{A} - \frac{\mathbf{d}\mathbf{d}'}{2m} \right) \hat{\mathbf{Z}} \right) + \alpha \left( \frac{\sqrt{k}}{n} \left\| \sum_u^n \hat{\mathbf{Z}}'_u \right\|_F - 1 \right) \end{aligned} \quad (20)$$

where the first term minimizes negative modularity and the second term is a “collapse regularization” meant to prevent all nodes from being assigned to the same block (with regularization strength given by hyperparameter  $\alpha$ ) [10].

### B.4 Markov Stability

Another perspective of modularity is taken in [12]. The approach, referred to as “Community Detection based on Markov Stability and Graph Neural Network” (CDMG), seeks to maximize a dynamic graph partition quality metric. Markov Stability [29–31] is defined

$$R_t = \text{tr}(\mathbf{Z}'(\mathbf{\Pi}\mathbf{M}^t - \mathbf{\pi}'\mathbf{\pi})\mathbf{Z}) \quad (21)$$

where  $\mathbf{\pi} = \mathbf{d}'(2m)^{-1}$ ,  $\mathbf{\Pi} = \text{diag}(\mathbf{\pi})$ ,  $\mathbf{M} = \mathbf{D}^{-1}\mathbf{A}$ , and  $\mathbf{D} = \text{diag}(\mathbf{d})$ . The matrix in this expression represents the probability that a random walk starting in one community will end in another after a certain number of time steps. This is also equivalent to the modularity in equation 19 when  $t = 1$  [31].

The CDMG approach seeks to maximize Markov Stability by minimizing its negative:

$$\begin{aligned} \hat{\mathbf{Z}} &= \text{ReLU}(\text{GNN}(\mathcal{G}, \mathbf{X}; \mathbf{W})) \\ J_{\text{CDMG}}(\mathbf{W}) &= -\text{tr}(\hat{\mathbf{Z}}'(\mathbf{\Pi}\mathbf{M}^t - \mathbf{\pi}'\mathbf{\pi})\hat{\mathbf{Z}}) \end{aligned} \quad (22)$$

where  $t$  is a hyperparameter that can be tuned according to the resolution of the graph communities. Larger values of  $t$  detect coarser communities [12, 29–31, 40]. We use  $t = 1$  for all experiments.

### B.5 Minimum-Cut

An objective function motivated by the minimum-cut problem is proposed in [8]. The minimum-cut problem is a task that seeks to find the partition that minimizes the number of edges between groups. This is done by maximizing the ratio of the number of edges within a group to the number of edges between groups in the rest of the graph. Formally,

$$\max \frac{1}{k} \sum_{i=1}^k \frac{\mathbf{Z}'_i \mathbf{A} \mathbf{Z}_{.i}}{\mathbf{Z}'_i \mathbf{D} \mathbf{Z}_{.i}} \quad (23)$$

where  $\mathbf{Z}_{.i} \in \{0, 1\}^n$  is the  $i^{\text{th}}$  column of  $\mathbf{Z}$  [41].



The proposed approach, known as MinCutPool [8], uses a pooling-based architecture with the softmax function applied to outputs to estimate community assignments. The objective function used for training is

$$\begin{aligned} \hat{\mathbf{Z}} &= \text{Softmax}(\text{GNN}(\mathcal{G}, \mathbf{X}; \mathbf{W})) \\ J_{\text{MCP}}(\mathbf{W}) &= -\frac{\text{tr}(\hat{\mathbf{Z}}' \tilde{\mathbf{A}} \hat{\mathbf{Z}})}{\text{tr}(\hat{\mathbf{Z}}' \tilde{\mathbf{D}} \hat{\mathbf{Z}})} + \left\| \frac{\hat{\mathbf{Z}}' \hat{\mathbf{Z}}}{\|\hat{\mathbf{Z}}' \hat{\mathbf{Z}}\|_F} - \frac{\mathbf{I}_k}{\sqrt{k}} \right\|_F \end{aligned} \quad (24)$$

where  $\mathbf{I}_k$  is the  $k \times k$  identity matrix,  $\tilde{\mathbf{A}} = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{1/2}$  is the normalized adjacency matrix, and  $\tilde{\mathbf{D}} = \text{diag}(\tilde{\mathbf{A}} \mathbf{1}_n)$  is the normalized degree matrix. The second term is an orthogonality penalty, meant to encourage orthogonality between communities and uniformity in community sizes.

## B.6 Poisson

The partial Poisson loss function in Section 3.3.2 can be extended to support multi-graphs, where  $\mathbf{A}_{uv} \in \mathbb{N}_0$ . To do so, the same log-likelihood function (equation 2) is considered. Note that the factorial term from equation 1 is dropped because its derivative is zero with respect to the parameters. Consequently, equation 8 is a valid loss function for both graphs and multi-graphs. In order to incorporate negative sampling, equation 11 is modified to

$$J_{p^*}(\mathbf{W}) \approx - \sum_{u,v \sim P_{\mathcal{E}}} [\mathbf{A}_{uv} \ln(\hat{\pi}_{uv}) - \hat{\pi}_{uv}] + \eta^{-1} \sum_{u,v \sim P_{\mathcal{N}}} \hat{\pi}_{uv} \quad (25)$$

where  $\eta$  is scaling constant meant to balance the contribution of the positive and negative edge sets to the total loss. Note that  $J_{p^*}$  is identical to  $J_p$  on binary adjacency matrices.

## C Datasets

Name	Nodes	Edges	Part.	Dim.	Overlap
SBM-4	2,000	202,102	4	100	No
SBM-8	2,000	150,648	8	100	No
SBM-16	2,000	96,436	16	100	No
SBM-32	2,000	56,526	32	100	No
OSBM-4	2,000	240,242	4	100	Yes
OSBM-8	2,000	196,510	8	100	Yes
OSBM-16	2,000	132,976	16	100	Yes
OSBM-32	2,000	72,422	32	100	Yes

Table 5: Synthetic graph summaries.

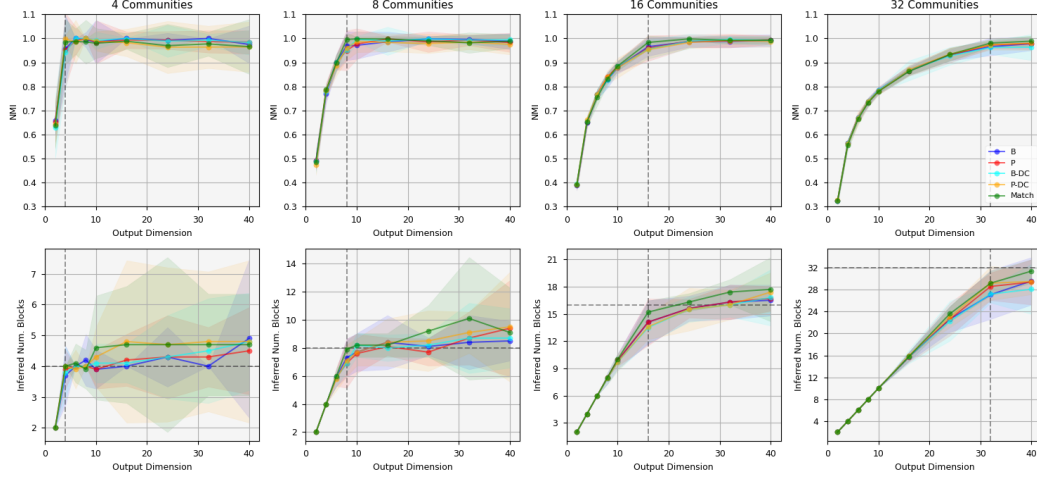
Name	Nodes	Edges	Part.	Dim.	Overlap
Cora	2,708	10,556	7	1,433	No
Citeseer	3,327	9,228	6	3,703	No
Pubmed	19,717	88,651	3	500	No
Wiki	11,367	431,726	10	300	No
ACB-Comp	13,752	491,722	10	767	No
ACB-Photo	7,650	238,163	8	745	No
MAG-Chem	35,409	314,716	14	4,877	Yes
MAG-CS	21,957	193,500	18	7,793	Yes
MAG-Eng	14,927	98,610	16	4,839	Yes
MAG-Med	63,282	1,620,628	17	5,538	Yes

Table 6: Real-world graph summaries.

For experiments on synthetic data, graphs are generated according to the degree-corrected Poisson SBM (equation 1). The generated graphs are all assortative — that is,  $\Theta_{ii} > \Theta_{ij}$  for all  $i \neq j$  — with density between 0.01 and 0.06. Furthermore, degree distributions are sampled from a scale free distribution, where  $\Pr(\mathbf{d}_u) \propto \mathbf{d}_u^{-\gamma}$  and  $\gamma = 2.5$ . Each graph has 2,000 nodes and 4, 8, 16, or 32 communities. The number of nodes in each community is normally distributed. For the overlapping setting, the partition and block matrices are randomly augmented.

To generate node features, length-100 vectors of means and variances are drawn from a multivariate normal distribution (squared for variances) for each community. Each node is given a feature vector sampled from a multivariate normal parameterized by the mean and variance of its assigned community. For overlapping communities, a node’s feature vector is the average of the vectors sampled for each of its communities. This method of attribute generation is similar to that of [10].

For experiments on real-world data, ten common benchmark graphs are considered. Cora [42], Citeseer [43], and Pubmed [44] are all citation networks, where nodes are publications and edges indicate citations. Node attributes are vector representations of text associated with each publication and community partitions are the category of the publication. AmazonCoBuy (as presented in [45]) is a dataset of co-purchase graphs, where nodes are products and edges indicate products that are purchased together. Nodes are attributed by vector representations of user reviews and community



**Figure 3:** Top row: Accuracy (NMI) vs. GNN output dimension for different numbers of communities. Bottom row: Inferred number of communities vs. GNN output dimension. Columns: Actual number of communities in the graph.

partitions are the category of the products. The dataset is split into two graphs: computer products (ACB-Comp) and photo products (ACB-Photo). These graphs are accessed through DGL<sup>3</sup>.

The Microsoft Academic Graph (as presented in [9]) is a dataset of co-authorship graphs, where nodes are authors and edges indicate co-authored publications. Node attributes are vectors of keywords associated with each author. Research areas form the overlapping community partition, as authors can research in multiple areas. The dataset is split into four graphs: chemistry (MAG-Chem), computer science (MAG-CS), engineering (MAG-Eng), and medicine (MAG-Med). These graphs are accessed through the GitHub repository<sup>4</sup> associated with [9].

## D Hyperparameter Sensitivity

This section provides an examination of SBM-specific hyperparameters. For these experiments, consideration is restricted to the non-overlapping setting and synthetically generated graphs are used.

To begin, we consider the relation of the output dimension  $d$  of the GNN embeddings to the number of effective clusters found in a given graph. The top row of figure 3 displays accuracy (measured as NMI) plotted against output dimension. The bottom row shows the inferred number of clusters (computed as  $|\mathcal{K}|$ ) against output dimension. Results are averaged over ten trials for graphs with  $k = 4, 8, 16, 32$  actual clusters.

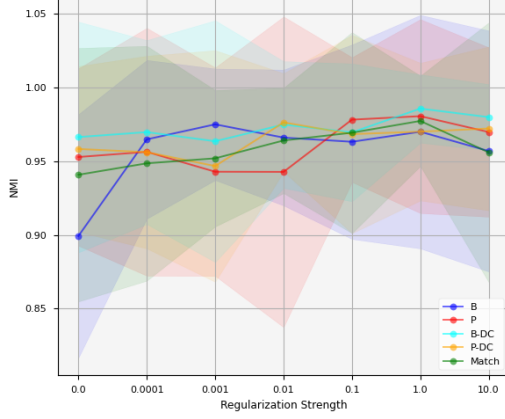
When the output dimension is smaller than the true number of clusters, there is generally poor accuracy and  $|\mathcal{K}| = d$ . When the output dimension is greater than the true number of clusters, accuracy is much higher and  $|\mathcal{K}| \approx k$ . The variance of community detection performance (in terms of accuracy and ability to recover the true number of clusters) tends to be greater when the output dimension is slightly larger than  $k$ . It is also worth noting that  $J_{\text{Match}}$  tends to overestimate the number of clusters.

Next, we look the impact of the regularization strength parameter  $\alpha$  from equation 16. Figure 4 shows accuracy measured against different values of  $\alpha$  averaged over ten trials. Note that the horizontal axis is not to scale. There is a moderate upward trend, suggesting that the regularization term does contribute to community detection performance. Indeed, all five loss functions exhibit better performance when  $\alpha = 1$  versus when the regularization term is left out ( $\alpha = 0$ ).

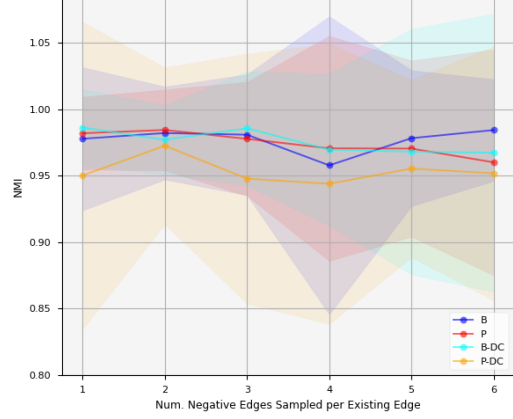
Finally, the effect of negative sampling is studied. Figure 5 plots accuracy measured against different numbers of negative edges sampled per each existing edge (referred to as  $\eta$ ). Accuracy appears fairly

<sup>3</sup>[https://www.dgl.ai/dgl\\_docs/api/python/dgl.data.html](https://www.dgl.ai/dgl_docs/api/python/dgl.data.html)

<sup>4</sup><https://github.com/shchur/overlapping-community-detection>



**Figure 4:** Accuracy (NMI) vs. regularization strength. Horizontal axis is not to scale.



**Figure 5:** Accuracy (NMI) vs. number of negative edges sampled per each existing edge.

consistent across different levels of negative sampling, with a slight drop-off after  $\eta = 3$ . For this experiment, the loss function is weighted such that negative and positive edges have equal importance.

## E Additional Experiments

In [10], a custom GNN architecture is implemented to go along with  $J_{\text{DMoN}}$ . The layer, which we refer to as SkipGCN, introduces an additional weight matrix for nodes' own features instead of a self-loop augmentation on the adjacency matrix. Additionally, the SELU activation function is used instead of ReLU. To provide a better comparison with this approach, we substitute SkipGCN layers for the standard GCN layers in our implementation. While the other implementation details (e.g., hidden dimension, weight decay, etc.) are not identical, the component GNN layer is the same as that described in [10]. It should be noted that in [10], the authors compare their approach to several other approaches that are not GNN based, require custom training routines, or depend on very specific architectures (e.g., [8, 11]). They show that their DMoN approach is superior on a number of real-world graphs, including Cora, Citeseer, and Pubmed.

	Cora		Citeseer		Pubmed		Wiki		ACB-Comp		ACB-Photo		Avg.	
	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1
$J_{\text{MCP}}$	32.0	18.2	18.2	12.8	<b>19.1</b>	26.1	<b>38.6</b>	<b>37.8</b>	44.2	<u>49.1</u>	<b>58.9</b>	<b>56.9</b>	<b>35.1</b>	33.5
$J_{\text{CDMG}}$	38.7	29.5	12.7	13.2	6.1	15.6	26.2	29.4	27.8	36.2	44.5	46.5	26.0	28.4
$J_{\text{DMoN}}$	37.1	21.3	21.3	15.7	16.9	15.4	33.7	24.9	31.8	27.2	49.9	39.1	31.8	23.9
$J_{\text{LP}}$	13.2	11.2	6.9	8.4	5.0	9.4	21.8	16.2	23.1	16.2	25.4	19.8	15.9	13.5
$J_{\text{NOCD}}$	14.6	11.4	7.2	9.6	6.0	12.9	<b>38.6</b>	<b>34.9</b>	<u>47.4</u>	45.2	<b>63.2</b>	<b>62.9</b>	29.5	29.5
$J_{\text{B}}$	39.0	33.6	24.2	<u>27.7</u>	17.8	<b>30.7</b>	27.8	28.8	38.0	43.1	56.0	53.0	33.8	<b>36.1</b>
$J_{\text{P}}$	34.5	27.3	<u>26.3</u>	<b>30.2</b>	14.8	25.8	27.2	26.3	38.8	43.1	54.3	52.2	32.6	34.1
$J_{\text{B-DC}}$	36.6	31.2	24.6	25.4	13.6	25.5	25.3	27.1	33.8	38.9	53.6	51.2	31.2	33.2
$J_{\text{P-DC}}$	<b>43.3</b>	<b>36.0</b>	<b>26.5</b>	23.6	<u>18.1</u>	<u>30.5</u>	21.0	24.6	<b>47.8</b>	<b>53.5</b>	52.5	48.6	<u>34.9</u>	<u>36.1</u>
$J_{\text{Match}}$	<u>41.7</u>	29.0	25.8	23.7	14.2	18.0	29.8	23.2	39.5	33.7	54.9	49.7	34.3	29.5

**Table 7:** Community detection performance on real-world data with *non-overlapping* communities using SkipGCN, as in [10]. Results are averaged over ten trials. The best scores (NMI and PF1) for each dataset are in bold.

Table 7 shows the results of this comparison with the SkipGCN layer on non-overlapping graphs. The SBM loss functions generally outperform or are competitive with the DMoN loss function. Interestingly, the MCP and NOCD loss functions also perform well in this setting.

To emphasize that this analysis is independent of specific GNN architectures, aggregated results of additional experiments are provided in tables 8 and 9 with the following architectures: GraphSAGE [24], Graph Attention Network (GAT) [25], and Graph Isomorphism Network (GIN) [26]. Table 8 shows the results of each loss function and each architecture averaged over all ten (synthetic and real-world) non-overlapping datasets. Table 9 shows the results over all eight (synthetic and real-world) overlapping datasets.

	GCN		GraphSAGE		GAT		GIN		SkipGCN		Avg.	
	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1
$J_{MCP}$	59.2	55.6	54.1	46.5	39.7	37.2	49.3	48.9	<b>60.5</b>	<b>58.0</b>	52.6	49.3
$J_{CDMG}$	55.9	53.5	45.2	44.9	38.5	37.8	51.6	50.8	52.5	49.9	48.7	47.4
$J_{DMoN}$	60.9	54.6	<u>54.2</u>	48.3	41.2	37.9	<u>57.9</u>	51.4	56.5	49.3	<u>54.1</u>	48.3
$J_{LP}$	49.9	44.5	48.2	44.4	40.3	37.1	52.0	48.4	42.4	35.6	46.6	42.0
$J_{NOCD}$	59.5	55.6	38.5	36.5	36.7	35.0	57.0	53.2	56.0	53.3	49.5	46.7
$J_B$	60.2	<u>61.0</u>	<b>54.5</b>	<u>53.8</u>	<b>42.7</b>	<b>41.7</b>	53.6	<b>54.8</b>	58.6	<u>57.8</u>	53.9	<b>53.8</b>
$J_P$	59.7	<b>61.4</b>	53.9	<b>54.0</b>	42.4	<u>41.6</u>	53.8	<u>54.7</u>	58.1	57.1	53.6	<u>53.8</u>
$J_{B-DC}$	59.8	59.6	53.0	52.9	40.8	39.0	51.8	52.7	57.1	56.2	52.5	52.1
$J_{P-DC}$	<u>61.2</u>	58.4	53.8	53.4	41.9	40.4	51.3	52.1	<u>59.0</u>	57.7	53.4	52.4
$J_{Match}$	<b>61.5</b>	57.0	51.9	50.5	<u>42.5</u>	40.6	<b>59.5</b>	54.5	58.4	53.3	<b>54.8</b>	51.2

**Table 8:** Community detection performance for different GNN architectures. Results are averaged over ten trials for all ten datasets with *non-overlapping* communities. The best scores (NMI and PF1) for each model are in bold.

	GCN		GraphSAGE		GAT		GIN		SkipGCN		Avg.	
	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1
$J_{MCP}$	38.9	53.8	42.6	52.6	32.9	42.8	36.4	53.4	46.9	54.0	39.6	51.3
$J_{CDMG}$	13.6	46.8	17.0	54.9	18.5	<u>49.7</u>	18.0	57.4	10.4	40.7	15.5	49.9
$J_{DMoN}$	37.1	49.9	44.6	53.4	<b>36.3</b>	45.4	38.5	54.7	43.3	53.5	40.0	51.4
$J_{LP}$	14.3	36.0	20.0	49.5	16.9	<b>50.9</b>	20.0	48.3	13.9	31.0	17.0	43.1
$J_{NOCD}$	40.7	48.5	17.2	42.6	16.3	42.8	<u>40.3</u>	57.5	42.3	50.6	31.4	48.4
$J_B$	<u>45.1</u>	58.5	<b>50.7</b>	<b>58.9</b>	<u>34.9</u>	47.3	37.4	58.8	<b>51.7</b>	60.2	<b>44.0</b>	<u>56.8</u>
$J_P$	41.6	58.4	<u>50.3</u>	58.3	32.8	45.6	38.7	58.6	48.3	61.4	42.3	56.4
$J_{B-DC}$	<b>45.7</b>	<b>60.1</b>	49.5	<u>58.8</u>	34.2	46.8	37.3	58.1	<u>48.4</u>	<b>63.1</b>	<u>43.0</u>	<b>57.4</b>
$J_{P-DC}$	43.3	<u>59.3</u>	46.3	56.1	30.1	43.1	34.3	<b>59.7</b>	48.4	<u>61.7</u>	40.5	56.0
$J_{Match}$	43.3	57.2	44.1	54.4	31.3	43.2	<b>40.7</b>	<u>59.4</u>	46.9	61.3	41.2	55.1

**Table 9:** Community detection performance for different GNN architectures. Results are averaged over ten trials for all eight datasets with *overlapping* communities. The best scores (overlapping NMI and PF1) for each model are in bold.

## F Proofs

### F.1 Maximum Likelihood Estimates

The Maximum Likelihood Estimate of  $\Theta$  is the same for the Poisson  $\ell_P$  and Bernoulli  $\ell_B$  models:

$$\begin{aligned}
 \ell_P(\mathbf{Z}, \Theta; \mathbf{A}) &= \sum_{u,v} [\mathbf{A}_{uv} \ln(\mathbf{Z}'_u \Theta \mathbf{Z}_v) - \mathbf{Z}'_u \Theta \mathbf{Z}_v] \\
 &= \sum_{i,j} [\mathbf{M}_{ij} \ln(\Theta_{ij}) - \mathbf{n}_i \mathbf{n}_j \Theta_{ij}] \\
 \frac{\partial \ell_P}{\partial \Theta_{ij}} &= \frac{\mathbf{M}_{ij}}{\Theta_{ij}} - \mathbf{n}_i \mathbf{n}_j \stackrel{\text{set}}{=} 0 \\
 \implies \hat{\Theta}_{ij}^P &= \frac{\mathbf{M}_{ij}}{\mathbf{n}_i \mathbf{n}_j} \\
 \ell_B(\mathbf{Z}, \Theta; \mathbf{A}) &= \sum_{u,v} [\mathbf{A}_{uv} \ln(\mathbf{Z}'_u \Theta \mathbf{Z}_v) + (1 - \mathbf{A}_{uv}) \ln(1 - \mathbf{Z}'_u \Theta \mathbf{Z}_v)] \\
 &= \sum_{i,j} [\mathbf{M}_{ij} \ln(\Theta_{ij}) + (\mathbf{n}_i \mathbf{n}_j - \mathbf{M}_{ij}) \ln(1 - \Theta_{ij})] \\
 \frac{\partial \ell_B}{\partial \Theta_{ij}} &= \frac{\mathbf{M}_{ij}}{\Theta_{ij}} - \frac{\mathbf{n}_i \mathbf{n}_j - \mathbf{M}_{ij}}{1 - \Theta_{ij}} \stackrel{\text{set}}{=} 0 \\
 \implies \hat{\Theta}_{ij}^B &= \frac{\mathbf{M}_{ij}}{\mathbf{n}_i \mathbf{n}_j}
 \end{aligned}$$

so  $\hat{\Theta}_{ij}^P = \hat{\Theta}_{ij}^B$ .

□

## F.2 Graph Matching

Let  $\mathcal{Z} = \{\tilde{\mathbf{Z}} \in \{0, 1\}^{n \times k} : \sum_{i=1}^k \tilde{\mathbf{Z}}_{ui} = 1 \text{ for all } u \in \mathcal{V}\}$ . Consider the distance quantity

$$q = \left\| \mathbf{A} - \tilde{\mathbf{Z}} \tilde{\boldsymbol{\Theta}} \tilde{\mathbf{Z}}' \right\|_F^2.$$

Minimizing  $q$  with respect to  $\boldsymbol{\Theta}$  gives

$$\begin{aligned} \frac{\partial q}{\partial \boldsymbol{\Theta}} &= 2\tilde{\mathbf{Z}} \tilde{\boldsymbol{\Theta}}' \tilde{\mathbf{Z}}' \tilde{\mathbf{Z}} \tilde{\mathbf{Z}}' - 2\mathbf{A}' \tilde{\mathbf{Z}} \tilde{\mathbf{Z}}' \stackrel{\text{set}}{=} 0 \\ \implies \tilde{\boldsymbol{\Theta}} &= (\tilde{\mathbf{Z}}' \tilde{\mathbf{Z}})^{-1} \tilde{\mathbf{Z}}' \mathbf{A} \tilde{\mathbf{Z}} (\tilde{\mathbf{Z}}' \tilde{\mathbf{Z}})^{-1}. \end{aligned}$$

Now recall the MLE  $\hat{\boldsymbol{\Theta}}$  from equation 3:

$$\hat{\boldsymbol{\Theta}}(\mathcal{G}, \tilde{\mathbf{Z}}) = \mathbf{M} \oslash \mathbf{n} \mathbf{n}' \quad \mathbf{M} = \sum_{u,v \in \mathcal{E}} \tilde{\mathbf{Z}}_u \tilde{\mathbf{Z}}_v' = \tilde{\mathbf{Z}}' \mathbf{A} \tilde{\mathbf{Z}} \quad \mathbf{n} = \sum_{u=1}^n \tilde{\mathbf{Z}}_u$$

where  $\oslash$  represents element-wise division. When  $\tilde{\mathbf{Z}} \in \mathcal{Z}$ , note that  $\tilde{\mathbf{Z}}' \tilde{\mathbf{Z}} = \text{diag}(\mathbf{n})$  and  $(\tilde{\mathbf{Z}}' \tilde{\mathbf{Z}})^{-1} = \text{diag}(\mathbf{n}^{-1})$  where  $\mathbf{n}^{-1}$  is the element-wise inverse of  $\mathbf{n}$ . Furthermore,  $[(\tilde{\mathbf{Z}}' \tilde{\mathbf{Z}})^{-1} (\tilde{\mathbf{Z}}' \tilde{\mathbf{Z}})^{-1}]_{ii} = 1/(\mathbf{n} \mathbf{n}')_{ii}$  for all  $i = 1, \dots, k$ . It can also be verified that

$$(\tilde{\mathbf{Z}}' \tilde{\mathbf{Z}})^{-1} \mathbf{M} (\tilde{\mathbf{Z}}' \tilde{\mathbf{Z}})^{-1} = \mathbf{M} \oslash \mathbf{n} \mathbf{n}'.$$

Therefore, the distance minimizer  $\tilde{\boldsymbol{\Theta}}$  is equivalent to the maximum likelihood estimate  $\hat{\boldsymbol{\Theta}}$  in this case.

We can now substitute the MLE into  $q$  as done in Section 3.3.3:

$$\begin{aligned} q &= \left\| \mathbf{A} \right\|_F^2 + \left\| \tilde{\mathbf{Z}} \hat{\boldsymbol{\Theta}} \tilde{\mathbf{Z}}' \right\|_F^2 - 2 \text{tr}(\mathbf{A}' \tilde{\mathbf{Z}} \hat{\boldsymbol{\Theta}} \tilde{\mathbf{Z}}') \\ &= \left\| \mathbf{A} \right\|_F^2 + \left\| \tilde{\mathbf{Z}} (\tilde{\mathbf{Z}}' \tilde{\mathbf{Z}})^{-1} \tilde{\mathbf{Z}}' \mathbf{A} \tilde{\mathbf{Z}} (\tilde{\mathbf{Z}}' \tilde{\mathbf{Z}})^{-1} \tilde{\mathbf{Z}}' \right\|_F^2 - 2 \text{tr}(\mathbf{A}' \tilde{\mathbf{Z}} \hat{\boldsymbol{\Theta}} \tilde{\mathbf{Z}}') \\ &= \left\| \mathbf{A} \right\|_F^2 + \text{tr}(\tilde{\mathbf{Z}} (\tilde{\mathbf{Z}}' \tilde{\mathbf{Z}})^{-1} \tilde{\mathbf{Z}}' \mathbf{A}' \tilde{\mathbf{Z}} \hat{\boldsymbol{\Theta}} \tilde{\mathbf{Z}}') - 2 \text{tr}(\mathbf{A}' \tilde{\mathbf{Z}} \hat{\boldsymbol{\Theta}} \tilde{\mathbf{Z}}') \\ &= \left\| \mathbf{A} \right\|_F^2 - \text{tr}(\mathbf{A}' \tilde{\mathbf{Z}} \hat{\boldsymbol{\Theta}} \tilde{\mathbf{Z}}'). \end{aligned}$$

The final equality comes from the cyclic property of the matrix trace:

$$\text{tr}(\tilde{\mathbf{Z}} (\tilde{\mathbf{Z}}' \tilde{\mathbf{Z}})^{-1} \tilde{\mathbf{Z}}' \mathbf{A}' \tilde{\mathbf{Z}} \hat{\boldsymbol{\Theta}} \tilde{\mathbf{Z}}') = \text{tr}(\tilde{\mathbf{Z}}' \tilde{\mathbf{Z}} (\tilde{\mathbf{Z}}' \tilde{\mathbf{Z}})^{-1} \tilde{\mathbf{Z}}' \mathbf{A}' \tilde{\mathbf{Z}} \hat{\boldsymbol{\Theta}}) = \text{tr}(\mathbf{A}' \tilde{\mathbf{Z}} \hat{\boldsymbol{\Theta}} \tilde{\mathbf{Z}}')$$

Dropping the constant term, we have the result

$$\underset{\tilde{\mathbf{Z}} \in \mathcal{Z}}{\text{argmin}} \left\| \mathbf{A} - \tilde{\mathbf{Z}} \hat{\boldsymbol{\Theta}} \tilde{\mathbf{Z}}' \right\|_F^2 = \underset{\tilde{\mathbf{Z}} \in \mathcal{Z}}{\text{argmin}} -\text{tr}(\mathbf{A}' \tilde{\mathbf{Z}} \hat{\boldsymbol{\Theta}} \tilde{\mathbf{Z}}').$$

which (noting that  $\|\cdot\|_F^2$  is a monotonic transformation of  $\|\cdot\|_F$  in equation 14) proves the statement in Section 3.3.3.  $\square$

Substituting predicted membership  $\hat{\mathbf{P}}$  for the mapping matrix  $\tilde{\mathbf{Z}}$  is a useful approximation of this result when training a GNN.

## G Standard Deviations

	SBM-4		SBM-8		SBM-16		SBM-32		Avg.	
	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1
$J_{MCP}$	1.7	1.7	1.1	1.5	2.2	6.2	0.8	2.6	1.4	3.0
$J_{CDMG}$	3.6	4.6	3.7	8.5	1.3	3.0	1.8	4.4	2.6	5.1
$J_{DMoN}$	6.1	8.5	1.8	2.5	0.6	1.8	0.8	2.4	2.3	3.8
$J_{LP}$	6.6	6.4	4.2	8.2	4.2	9.8	2.3	5.4	4.3	7.4
$J_{NOCD}$	2.5	2.1	1.3	2.1	1.6	4.7	1.6	4.6	1.8	3.4
$J_B$	1.7	1.3	1.4	3.6	1.3	4.6	0.5	2.5	1.2	3.0
$J_P$	2.4	2.6	0.6	1.0	1.1	4.1	1.1	3.6	1.3	2.8
$J_{B-DC}$	4.1	4.9	1.6	3.7	2.4	6.3	1.4	4.5	2.3	4.8
$J_{P-DC}$	2.4	2.4	0.7	1.1	0.5	1.9	0.6	2.5	1.1	2.0
$J_{Match}$	2.4	2.6	2.5	4.9	0.6	1.7	0.7	3.6	1.6	3.2

**Table 10:** Standard deviations associated with table 1. Community detection on synthetic data with *non-overlapping* communities using GCN.

	OSBM-4		OSBM-8		OSBM-16		OSBM-32		Avg.	
	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1
$J_{MCP}$	21.8	12.3	16.0	13.4	13.3	9.8	7.3	6.8	14.6	10.6
$J_{CDMG}$	17.4	27.9	14.4	11.9	13.1	9.9	8.8	16.8	13.4	16.6
$J_{DMoN}$	25.9	39.9	15.6	6.9	5.1	13.1	11.7	10.3	14.6	17.5
$J_{LP}$	18.0	23.8	12.7	13.2	10.1	11.1	5.1	9.3	11.5	14.4
$J_{NOCD}$	17.4	25.1	9.6	6.3	9.2	5.9	13.9	9.9	12.5	11.8
$J_B$	14.6	14.7	8.9	11.3	11.0	4.8	5.5	4.8	10.0	8.9
$J_P$	18.1	11.5	13.1	8.7	9.3	5.9	11.4	6.7	13.0	8.2
$J_{B-DC}$	19.1	7.0	14.7	9.2	14.0	7.2	8.4	5.4	14.0	7.2
$J_{P-DC}$	12.3	4.0	13.3	13.1	11.9	13.9	10.1	10.1	11.9	10.3
$J_{Match}$	16.3	11.9	14.7	7.6	5.1	8.6	10.4	7.0	11.6	8.8

**Table 11:** Standard deviations associated with table 2. Community detection on synthetic data with *overlapping* communities using GCN.

	Cora		Citeseer		Pubmed		Wiki		ACB-Comp		ACB-Photo		Avg.	
	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1
$J_{MCP}$	3.6	2.5	4.6	2.5	6.2	3.6	8.2	3.3	2.2	4.1	4.8	3.8	5.0	3.3
$J_{CDMG}$	10.4	9.2	2.6	1.8	3.8	5.0	7.3	5.9	4.2	6.6	3.0	5.3	5.2	5.7
$J_{DMoN}$	1.2	2.0	2.7	1.6	1.6	1.2	3.5	3.7	1.7	2.4	3.6	3.4	2.4	2.4
$J_{LP}$	1.4	2.9	2.5	3.4	1.7	1.2	2.8	3.2	3.6	4.8	3.0	4.7	2.5	3.4
$J_{NOCD}$	9.0	5.3	1.5	2.0	5.1	3.5	4.6	4.2	3.8	4.6	1.5	3.7	4.2	3.9
$J_B$	11.0	5.6	2.1	2.4	5.9	7.9	7.2	4.3	6.7	4.7	2.6	4.2	5.9	4.9
$J_P$	11.0	6.5	2.5	2.5	3.8	6.6	6.8	4.9	9.6	7.1	3.4	4.0	6.2	5.3
$J_{B-DC}$	11.0	8.7	2.5	3.6	5.4	9.5	3.8	5.2	6.1	5.0	2.4	3.1	5.2	5.9
$J_{P-DC}$	11.9	9.2	2.2	4.1	6.6	6.4	7.9	5.2	2.1	2.3	2.8	2.1	5.6	4.9
$J_{Match}$	12.3	7.0	1.7	2.2	1.7	2.9	2.7	6.0	2.6	3.5	2.1	3.5	3.9	4.2

**Table 12:** Standard deviations associated with table 3. Community detection on real-world data with *non-overlapping* communities using GCN.

	MAG-Chem		MAG-CS		MAG-Eng		MAG-Med		Avg.	
	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1	NMI	PF1
$J_{MCP}$	6.7	4.0	2.2	0.7	6.8	7.9	3.3	1.1	4.7	3.4
$J_{CDMG}$	1.2	10.0	2.3	5.6	1.3	15.5	1.0	3.7	1.5	8.7
$J_{DMoN}$	6.2	3.2	2.7	0.9	3.1	1.8	2.8	1.2	3.7	1.8
$J_{LP}$	0.0	1.3	0.5	1.0	0.0	1.0	0.6	22.5	0.3	6.4
$J_{NOCD}$	4.6	2.3	2.4	1.6	6.0	7.0	3.2	1.1	4.0	3.0
$J_B$	11.6	5.4	3.5	1.3	5.8	3.5	3.9	2.4	6.2	3.1
$J_P$	8.8	6.8	3.1	3.2	6.6	5.1	3.8	1.3	5.6	4.1
$J_{B-DC}$	10.8	5.2	2.9	2.0	5.8	5.4	3.0	2.0	5.6	3.6
$J_{P-DC}$	8.1	6.5	4.3	3.5	6.2	11.8	3.4	5.3	5.5	6.8
$J_{Match}$	10.1	4.4	3.6	4.1	4.4	4.7	3.5	2.5	5.4	3.9

**Table 13:** Standard deviations associated with table 4. Community detection on real-world data with *overlapping* communities using GCN.



	Cora		Citeseer		Pubmed		Wiki		ACB-Comp		ACB-Photo		Avg.	
	<i>NMI</i>	<i>PFI</i>	<i>NMI</i>	<i>PFI</i>	<i>NMI</i>	<i>PFI</i>	<i>NMI</i>	<i>PFI</i>	<i>NMI</i>	<i>PFI</i>	<i>NMI</i>	<i>PFI</i>	<i>NMI</i>	<i>PFI</i>
$J_{MCP}$	1.4	1.7	2.2	1.3	1.4	4.7	8.2	6.2	7.7	6.4	9.2	8.3	5.0	4.8
$J_{CDMG}$	2.6	3.5	3.9	2.6	2.7	2.2	4.1	2.9	7.9	9.9	3.6	4.3	4.1	4.2
$J_{DMoN}$	9.7	4.5	2.7	1.3	2.6	0.8	2.5	1.6	2.8	2.1	1.4	3.2	3.6	2.3
$J_{LP}$	1.3	2.3	0.9	0.8	0.6	0.7	1.5	1.8	1.9	1.7	3.4	3.0	1.6	1.7
$J_{NOCD}$	5.8	2.5	2.1	1.4	4.5	3.6	5.9	6.2	4.0	3.2	2.6	3.8	4.1	3.4
$J_B$	4.8	4.4	2.4	3.1	2.0	4.7	5.7	3.1	3.7	3.5	3.7	4.7	3.7	3.9
$J_P$	9.1	6.4	3.0	3.2	4.3	6.1	7.2	2.8	2.5	4.9	4.5	5.0	5.1	4.8
$J_{B-DC}$	5.2	6.1	2.5	3.5	3.7	5.0	2.3	1.8	5.9	4.6	3.7	5.0	3.9	4.3
$J_{P-DC}$	3.4	3.9	2.5	3.5	6.9	5.7	5.0	6.2	3.1	3.5	4.3	4.5	4.2	4.5
$J_{Match}$	1.9	1.6	3.7	5.4	1.4	3.6	2.3	1.7	4.0	7.4	3.6	5.9	2.8	4.3

**Table 14:** Standard deviations associated with table 7. Community detection on real-world data with *non-overlapping* communities using SkipGCN.

	GCN		GraphSAGE		GAT		GIN		SkipGCN		Avg.	
	<i>NMI</i>	<i>PFI</i>	<i>NMI</i>	<i>PFI</i>	<i>NMI</i>	<i>PFI</i>	<i>NMI</i>	<i>PFI</i>	<i>NMI</i>	<i>PFI</i>	<i>NMI</i>	<i>PFI</i>
$J_{MCP}$	34.0	33.2	34.0	34.7	37.5	31.1	37.1	29.7	33.5	33.4	35.2	32.4
$J_{CDMG}$	34.5	28.6	38.2	30.0	37.1	28.9	34.5	27.1	34.6	29.4	35.8	28.8
$J_{DMoN}$	30.5	31.3	34.3	34.0	35.1	28.6	32.7	32.1	32.1	33.0	32.9	31.8
$J_{LP}$	35.7	31.5	35.0	29.3	35.3	28.1	34.2	29.0	33.9	28.3	34.8	29.2
$J_{NOCD}$	33.6	32.2	39.7	32.2	36.3	28.2	32.3	28.4	36.9	34.2	35.8	31.0
$J_B$	33.7	28.2	36.5	31.8	38.4	32.4	37.0	29.5	32.2	28.4	35.6	30.1
$J_P$	34.1	28.5	37.1	32.0	38.9	32.9	35.8	28.0	33.0	30.1	35.8	30.3
$J_{B-DC}$	33.3	28.9	37.4	32.9	38.4	31.9	37.3	29.8	33.5	30.0	36.0	30.7
$J_{P-DC}$	32.9	30.4	36.7	32.2	38.8	32.9	36.8	29.4	31.8	28.8	35.4	30.7
$J_{Match}$	32.1	31.5	36.5	33.6	38.2	33.0	32.0	29.0	31.6	30.9	34.1	31.6

**Table 15:** Standard deviations associated with table 8. Community detection for different GNN architectures on data with *non-overlapping* communities.

	GCN		GraphSAGE		GAT		GIN		SkipGCN		Avg.	
	<i>NMI</i>	<i>PFI</i>	<i>NMI</i>	<i>PFI</i>	<i>NMI</i>	<i>PFI</i>	<i>NMI</i>	<i>PFI</i>	<i>NMI</i>	<i>PFI</i>	<i>NMI</i>	<i>PFI</i>
$J_{MCP}$	20.0	28.0	25.6	29.1	21.8	25.0	16.7	24.1	24.7	28.8	21.8	27.0
$J_{CDMG}$	15.9	20.6	18.6	31.7	20.0	27.7	20.1	30.7	13.1	22.1	17.5	26.6
$J_{DMoN}$	22.2	28.9	25.1	28.6	18.4	23.2	16.6	29.5	22.4	27.6	20.9	27.6
$J_{LP}$	16.8	27.9	20.9	26.6	18.9	25.5	20.8	31.8	15.3	21.6	18.5	26.7
$J_{NOCD}$	20.5	24.0	16.5	33.6	17.5	28.9	15.0	26.7	24.4	26.7	18.8	28.0
$J_B$	22.0	25.1	27.2	26.1	25.8	26.9	17.6	25.0	22.1	25.7	22.9	25.8
$J_P$	20.6	26.2	25.6	25.2	26.2	27.0	19.1	25.1	21.7	26.9	22.6	26.1
$J_{B-DC}$	21.3	23.6	26.9	25.6	27.5	28.7	17.5	25.6	21.4	25.7	22.9	25.8
$J_{P-DC}$	19.3	25.5	29.1	28.9	27.5	29.1	18.5	28.1	20.8	27.9	23.0	27.9
$J_{Match}$	20.8	24.6	29.5	30.1	27.0	28.9	15.3	26.7	21.7	28.2	22.8	27.7

**Table 16:** Standard deviations associated with table 9. Community detection for different GNN architectures on data with *overlapping* communities.