

# Ergodic limits, relaxations, and geometric properties of random walk node embeddings

---

C. Lin, D. Sussman, P. Ishwar. "Ergodic Limits, Relaxations, and Geometric Properties of Random Walk Node Embeddings." IEEE Transactions on Network Science and Engineering, <https://doi.org/10.48550/arXiv:1803.09882>, <https://hdl.handle.net/2144/44494>

*"Downloaded from OpenBU. Boston University's institutional repository."*

# Ergodic Limits, Relaxations, and Geometric Properties of Random Walk Node Embeddings

Christy Lin, Daniel Sussman, and Prakash Ishwar

**Abstract**—Random walk based node embedding algorithms learn vector representations of nodes by optimizing an objective function of node embedding vectors and skip-bigram statistics computed from random walks on the network. They have been applied to many supervised learning problems such as link prediction and node classification and have demonstrated state-of-the-art performance. Yet, their properties remain poorly understood. This paper studies properties of random walk based node embeddings in the unsupervised setting of discovering hidden block structure in the network, i.e., learning node representations whose cluster structure in Euclidean space reflects their adjacency structure within the network. We characterize the ergodic limits of the embedding objective, its generalization, and related convex relaxations to derive corresponding non-randomized versions of the node embedding objectives. We also characterize the optimal node embedding Grammians of the non-randomized objectives for the expected graph of a two-community Stochastic Block Model (SBM). We prove that the solution Grammian has rank 1 for a suitable nuclear norm relaxation of the non-randomized objective. Comprehensive experimental results on SBM random networks reveal that our non-randomized ergodic objectives yield node embeddings whose distribution is Gaussian-like, centered at the node embeddings of the expected network within each community, and concentrate in the linear degree-scaling regime as the number of nodes increases.



## 1 INTRODUCTION

MOST statistical and computational tools originally developed for vector-valued data do not leverage the unique structured form of network data. Tools that exploit the graph-structure of network data could be custom-made for each network problem. A powerful alternative, however, is to develop a Euclidean-space embedding of a network that enables methods and tools developed for Euclidean-space data to effectively reason about various network properties.

**Node embedding** algorithms [1] aim to map nodes of a given graph into points in Euclidean space (i.e., vectors in  $\mathbb{R}^d$ ) such that their relative positions capture their propensities for adjacency within the network. These embeddings make it possible to apply to network data, tools and algorithms from multivariate statistics and machine learning that were developed for Euclidean-space data. For example, with suitable embeddings, node classification, community detection, and vertex nomination problems reduce, respectively, to standard classification, clustering, and ranking problems. Therefore, developing new node embedding algorithms, establishing the theoretical properties of these embeddings, and demonstrating how connectivity properties are reflected in the embedding space is fundamental to developing principled network inference procedures.

**Random walk embeddings** [2], [3], [4], [5] are a class of recently developed node embedding techniques which use

random walks on graphs to capture notions of proximity between nodes. They may be viewed as network counterparts of techniques used for learning **word embeddings** [6], [7] in the field of natural language processing. In fact, by viewing samples of random walks in the network as sentences, with nodes playing the role of words, word embeddings can be directly applied to networks to yield node embeddings. Nodes which appear nearby within a sample of a random walk are analogous to words that appear nearby within a sentence. Word embeddings have been found to accurately capture the relationships between words and have been highly successful in several natural language processing tasks such as topic modeling, translation, and word analogy [8]. Random walk node embeddings too have been applied to a number of supervised and unsupervised learning problems such as link prediction, node classification and community detection and have demonstrated state-of-the-art performance [2], [3], [4], [5].

Unfortunately, despite excellent empirical performance in a number of supervised learning problems, random walk embeddings remain poorly understood. This is in stark contrast to the well-known spectral embeddings whose properties for the *unsupervised* learning problem of community detection have been extensively studied and characterized under a variety of statistical network models, specifically the Stochastic Block Model (SBM) and its generalizations [9], [10], [11], [12], [13], [14]. Attempts of theoretical analysis so far have focused on building connections between random walk embeddings algorithms and matrix factorization [15]. The properties of the resulting embedding vectors, however, still remain unexplored.

**Contributions:** This paper proposes a framework for random walk based node-embedding algorithms for graphs. This is based on learning node embeddings by optimizing objective functions involving skip-bigram statistics computed from random walks on a graph. This framework

- Christy Lin is with the Division of Systems Engineering, College of Engineering, Boston University, 15 St Marys St, Boston, MA, 02215.

cy93lin@bu.edu

- Daniel Sussman is with the Department of Mathematics & Statistics, College of Arts and Sciences, Boston University, 111 Cummington Mall, Boston, MA, 02215. sussman@bu.edu
- Prakash Ishwar is with the Division of Systems Engineering and Department of Electrical & Computer Engineering, College of Engineering, Boston University, 8 St Marys St, Boston, MA, 02215. pi@bu.edu

subsumes several existing algorithms as special cases and introduces extensions and techniques that simplify theoretical analysis. We establish ergodic limits of the proposed node-embeddings. We analyze Grammian re-parameterized convex relaxations and characterize the solution for the expected graph of a two-community SBM and the unconstrained solution for any graph. We prove that the solution of the expected graph of a two-community SBM has rank at most 2. We develop algorithms for computing solutions to our proposed embedding objectives for general graphs and conduct numerical experiments to understand the geometric structure of embedding vectors (community clustering and separation properties) for SBM random graphs. We also empirically study the concentration properties of node embeddings for SBM random graphs in the linear and logarithmic scaling regimes. We find empirically that the distribution of embeddings are Gaussian-like, centered at the node embeddings of the expected graph within each community, and that they concentrate in the linear degree scaling regime as the number of nodes increases.

**Paper organization:** Section 2 overviews recent work on random walk embeddings, sets up basic notation, and provides background on SBMs. Section 3 describes our proposed theoretical framework, results on ergodic limits (Section 3.1), various relaxations (Section 3.3), and the characterization of the solution for the expected graph of a two-community SBM (Section 3.4). Section 4 describes the setting of all our experiments in full detail. The geometric and concentration properties of the distribution of embedding vectors of our proposed algorithms under 2-community SBM are presented and discussed in Section 5. Concluding remarks appear in Section 6.

**Notation:** In this work we consider graphs that are undirected and simple with a node set  $\mathcal{V} = [n] := \{1, 2, \dots, n\}$  and an edge set  $\mathcal{E} \subset \{\{i, j\} : i, j \in \mathcal{V}, i \neq j\}$ . The edges may be possibly weighted. We denote such a graph by  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  and its adjacency matrix by  $A \in \{0, 1\}^{n \times n}$ , where  $A_{ij} = 0$  if, and only if,  $\{i, j\} \in \mathcal{E}$ .

We denote the set of all real numbers by  $\mathbb{R}$ , the set of all natural numbers by  $\mathbb{N}$ , the set of all  $n \times n$  real symmetric matrices by  $\mathbb{S}^n$ , the set of all real, symmetric, and positive semidefinite matrices by  $\mathbb{S}_+^n$ , and the natural logistic-loss function by  $\sigma(t) := \ln(1 + e^{-t}), t \in \mathbb{R}$ . Matrix transpose is denoted by  $\top$ .

## 2 BACKGROUND AND RELATED WORK

In this section we overview recent work on random walk node embeddings with a focus on the unsupervised algorithm VEC. We also summarize key aspects of the Stochastic Block Model (SBM) used in our experiments.

### 2.1 Random walk node embedding algorithms

A random walk node embedding algorithm typically consists of three steps: 1) Generating multiple random walks over the graph via Markov chains with the set of nodes as the state space, specified probability transition matrices at each step, and specified initial distributions. 2) Computing various statistics from the sample paths of the random walks. 3) Generating embeddings by optimizing a function

that only involves the computed statistics and node embedding variables of the input graph.

Among the random walk node embedding algorithms, [2], [4], [16] make use of node embeddings within the context of *supervised* learning problems such as node attribute prediction and link prediction and accordingly design probability transition matrices that depend on the supervised labels. In contrast, the VEC algorithm [5] focuses on the *unsupervised* community detection problem [17]. The unsupervised setting of [5] is ideal for studying random walk node embeddings that capture pure network connectivity properties unsullied by node labels. We therefore select VEC as our prototypical algorithm for analysis and introduce it in detail in the next subsection.

While our focus is on unsupervised setting, the general Markov-Chain based framework we develop can be used to analyze the supervised setting as well through transition matrices that are label dependent.

In addition to the node embedding algorithms discussed above, the use of a random walks and their steady-state-distributions for graph clustering has been studied in [18] and [19]. Subsequent work [20] further proposed to exploit multi-step transition probabilities between nodes for clustering.

In terms of theoretical results, [21] have analyzed the stationary distribution of second-order random walks in [4] for specific types of networks. We provide a complete characterization of the ergodic limits for general random walk node embedding objectives in Section 3. For the task of community detection, [22] have provided large-sample error bounds for consistent community recovery from the perspective of matrix factorization. Their setting is a special, unconstrained case of our general problem stated in Definition 6 of Section 3.3.

### 2.2 VEC: unsupervised random walk node embedding

VEC learns a low-dimensional vector representation for each node of a graph such that the local neighborhood structures of the graph are encoded within the Euclidean geometry of node vectors. Specifically, the inner product between the embedding vectors of node pairs encode their propensity to appear nearby in random walks on the graph.

VEC generates  $r$  random walks on  $\mathcal{G}$  of fixed length  $\ell$  starting from each node. We let  $\{X_s^{(m,p)}\}_{s=1}^{\ell}, p = 1, \dots, r$ , denote the  $p$ -th random walk starting from node  $m$ . All random walks follow the “natural” transition matrix  $W$  where the next node is chosen from the immediate neighbors of the current node with probability proportional to the edge weight between them.

VEC learns node embedding using the negative-sampling framework of noise-contrastive estimation [6]. The statistics used for learning node embeddings are based on two multisets of node pairs that are computed from the sample paths of the random walks as follows. The positive multiset  $\mathcal{D}_+$  consists of all node pairs  $(X_s^{(m,p)}, X_{s'}^{(m,p)})$ , including repetitions, that occur within  $w$  steps of each other, i.e.,  $|s - s'| \leq w$ , in all the generated sample paths. Such node pairs are called  $w$ -skip bigrams in Natural Language Processing with words viewed as nodes and sentences as sample paths of random walks. The algorithm parameter

$w$  controls the size of the local neighborhood of a node in the given graph. The negative multiset  $\mathcal{D}_-$  is constructed as follows. For each node pair  $(i, j)$  in  $\mathcal{D}_+$ , we append  $k$  node pairs  $(i, j_1), \dots, (i, j_k)$  to  $\mathcal{D}_-$ , where the  $k$  nodes  $j_1, \dots, j_k$  are drawn in an IID manner from *all* the nodes according to the empirical unigram node distribution computed from all the sample paths. Let  $n_{ij}^+$  and  $n_{ij}^-$  denote the number of  $(i, j)$  pairs, counting repetitions, in  $\mathcal{D}_+$  and  $\mathcal{D}_-$  respectively.

VEC finds the embedding vector  $\mathbf{u}_i \in \mathbb{R}^d$  for each node  $i$  by solving the following minimization problem:

**Definition 1** (VEC optimization problem).

$$\arg \min_{\{\mathbf{u}_i \in \mathbb{R}^d, i \in \mathcal{V}\}} \sum_{(i,j) \in \mathcal{V}^2} \left[ n_{ij}^+ \sigma(\mathbf{u}_i^\top \mathbf{u}_j) + n_{ij}^- \sigma(-\mathbf{u}_i^\top \mathbf{u}_j) \right] \quad (1)$$

One approach to solve Eq. (1) is via stochastic gradient descent (SGD) [23], [24]. This approach is followed in [6] and implemented in Python `gensim` package. Besides its conceptual simplicity, SGD can be parallelized and nicely scaled to large datasets [25]. The per-iteration computational complexity of the SGD algorithm used to solve Eq. (1) is  $O(d)$ , i.e., linear in the embedding dimension. The number of iterations is  $O(r\ell wk)$ .

### 2.3 Stochastic Block Model

The Stochastic Block Model (SBM) [26], [27], [28] is a canonical generative probabilistic model for random graphs that reflects block (community) structures among the nodes wherein nodes within the same block have the same tendencies for connecting to all other nodes. Free of node or edge labels, it serves as a clean platform for generating graphs to empirically study and compare the properties of various node embedding algorithms and conduct a theoretical analysis. For example, SBM has helped in understanding the behavior of spectral embeddings [12].

For any given  $K \in \mathbb{N}$ , a  $K$ -block SBM is parameterized by the latent block membership labels  $y_1, \dots, y_n \in [K]$ , and the edge probability matrix, a symmetric matrix  $B \in [0, 1]^{K \times K}$ . The latent labels  $\{y_i\}$  partition the nodes into communities indexed by each  $k \in [K]$ . We note that there are versions of SBM in which the  $y_i$ 's are treated as random. This, however, poses minor additional difficulties. To ease the subsequent discussion, unless noted otherwise, the  $y_i$ 's will always be viewed as fixed deterministic unknowns throughout this work. For a node in block  $k_1$  and a *different* node in block  $k_2$  (where  $k_2$  may equal  $k_1$ ), the probability that an edge is present between the two nodes is  $B_{k_1 k_2}$ , and all edges appear independently. We use this model for generating graphs in all our experiments.

The goal of any community detection algorithm is to learn the latent communities of nodes purely from the graph structure. Thus community detection is an unsupervised learning problem which aims to uncover the underlying block structure. A series of work [29], [30], [31], [32], [33], [34] characterizes the information-theoretic limits of community detection in SBMs in different degree-scaling regimes. Some of our experiments are designed to operate with respect to these information-theoretic limits.

## 3 ANALYTICAL FRAMEWORK AND RESULTS

There are three distinct challenges which complicate the analysis of VEC embedding vectors and their relationship to the latent graph community structure. First, the objective function Eq. (1) is nonlinear due to the logistic loss function. Second, even though the function  $\sigma(t)$  is strictly convex, the overall objective is not convex with respect to the node embedding vectors. Finally, the objective function is itself random, partly due to intrinsic randomness in network connectivity, but also due to algorithmic randomness from the random walks and the Stochastic Gradient Descent algorithm.

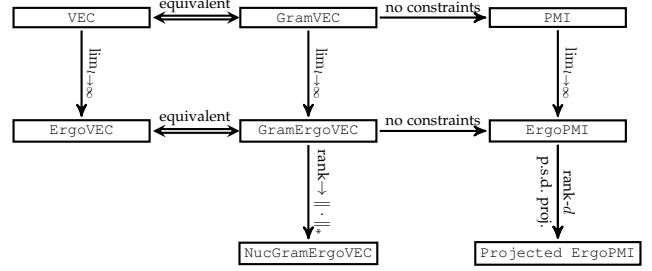


Fig. 1: Relationships between analysis strategies.

To tackle these challenges, in this section we introduce and develop techniques, generalized formulations, and their extensions which are more amenable to theoretical analysis. We leverage three distinct strategies whose inter-relationships are succinctly depicted in Fig. 1. These are:

- (1) *ErgoVEC: Ergodic limits of random walks* ( $\lim_{\ell \rightarrow \infty}$ ). We begin by noting that the sampled coefficients,  $n_{ij}^+$ 's and  $n_{ij}^-$ 's in Eq. (1), inherit the randomness of the random walks and depend on a number of algorithm parameters that are described in Sec. 2.2. Previous empirical results [5] demonstrate that the parameters such as the number of random walks  $r$  and their length  $\ell$  do not substantially impact performance. Motivated by this observation, as a first step, in Sec. 3.1, we eliminate algorithmic randomness by taking the ergodic limits ( $\ell \rightarrow \infty$ ) of the coefficients. This gives rise to a more principled formulation, which we call *ErgoVEC*, that removes dependence on sampled random walks and parameters  $r$  and  $\ell$ .
- (2) *GramErgoVEC and PMI: reparameterize, unconstrain and project*. Like VEC, ErgoVEC is a nonconvex optimization problem since the objective is a nonconvex function of the embedding vectors  $\mathbf{u}_i$ . We leverage a reparameterization trick which is similar in spirit to that used in [35] to arrive at an equivalent problem, named GramErgoVEC, that has a convex objective function with respect to new matrix variables and additional constraints. GramErgoVEC has a convex objective, but is still a nonconvex optimization problem due to the rank constraint. In order to gain insight into the structure of the solution, we characterize the solution to GramErgoVEC without any constraints and then project the unconstrained solution onto the constraint set. It turns out that the solution to the unconstrained GramErgoVEC objective is directly related to the so-

called Pointwise Mutual Information (PMI) matrix [36]. We study GramErgoVEC and PMI in Sec. 3.3

- (3) *NucGramErgoVEC: reparameterize and convexify.* Another strategy to convexify GramErgoVEC is to replace the non-convex rank constraint by a convex nuclear norm constraint. We term the resulting optimization problem *NucGramErgoVEC* and study its properties in the later part of Sec. 3.3

In the rest of this section, we will formally study and establish important theoretical properties of these alternative formulations and their inter-relationships.

### 3.1 Ergodic limits

As described in Sec. 2.2,  $n_{ij}^+$  and  $n_{ij}^-$  are the number of the  $(i, j)$  node pairs in the positive and negative multisets,  $\mathcal{D}_+$  and  $\mathcal{D}_-$ , respectively. These depend on 5 algorithm parameters:  $r$  (number of random walks per node),  $\ell$  (length of each walk),  $w$  (context window size) and  $k$  (number of negative  $w$ -skip bigrams per positive  $w$ -skip bigram). Specifically,  $n_{ij}^+$ , as a  $w$ -skip bigram count over  $r$  IID sets of  $n$  random walks, increases proportionally with  $r$  and  $n$  and roughly proportionately with  $\ell$ , for large  $\ell$ , since the number of segments of  $w$  consecutive steps in a length- $\ell$  walk equals  $(\ell - w + 1)$ . They also increase as  $w$  increases however their distribution can change substantially with  $w$ . As for the negative multiset, note that  $|\mathcal{D}_-| = k|\mathcal{D}_+|$ , so  $n_{ij}^-$  increases proportional to  $k, r$ , and  $\ell$ . Among these parameters, the results in [5] show that  $r$  and  $\ell$  have little effect on the final performance of VEC, while  $w$  plays a more important role.

Besides their dependence on the algorithm parameters,  $n_{ij}^+$ 's and  $n_{ij}^-$ 's inherit the randomness intrinsic to the random walks. Additionally, the number of negative  $(i, j)$  pairs  $n_{ij}^-$  also inherit randomness from the categorical sampling of appended nodes  $j_1, \dots, j_k$  in the negative pairs. In order to gain an algorithmic-randomness-free understanding of network properties captured by  $n_{ij}^+$  and  $n_{ij}^-$ , we study their *ergodic limits*.

**Definition 2** (Ergodic limits of  $n_{ij}^+$  and  $n_{ij}^-$ ). *Let  $n_{ij}^+$  and  $n_{ij}^-$  be defined as above. The (normalized) ergodic limits of  $n_{ij}^+$  and  $n_{ij}^-$  are defined as*

$$\bar{n}_{ij}^+ := \frac{1}{rn} \lim_{\ell \rightarrow \infty} \frac{n_{ij}^+}{\ell}, \quad (2)$$

$$\bar{n}_{ij}^- := \frac{1}{rn} \lim_{\ell \rightarrow \infty} \frac{n_{ij}^-}{\ell}, \quad (3)$$

whenever these limits exist in the almost sure sense.

The Ergodic limits in Definition 2 provide, for a given graph, a deterministic version of  $n_{ij}^+$  and  $n_{ij}^-$ , normalized by the cumulative length of all random walks. We note that letting  $\ell$  go to infinity may seem like incorporating global information about the entire graph instead of the more useful local connectivity patterns, but this is not the case. Regardless of the value of  $\ell$ ,  $\mathcal{D}_+$  only contains pairs of nodes which appear within  $w$  steps from each other. Therefore, the positive pairs sampled still reflect local information.

In VEC we launch  $r$  random walks starting deterministically from each node which yields a total of  $rn$  random walks. Dividing the  $w$ -skip bigram counts by  $rn$  averages

them across all random walks. The averaged counts can be loosely viewed as arising from a single random walk with a uniform initial distribution over nodes, i.e., with probability  $1/n$  for each node. If the Markov Chain underlying the random walk is ergodic, as  $\ell$  tends to infinity, the  $n_{ij}^+$ 's and  $n_{ij}^-$ 's, suitably normalized, will converge to their respective expected values under the sampling distribution of random walks. This intuition is formalized in Theorem 1 below. The theorem encompasses disconnected graphs that consist of several connected components that are often encountered in practice. In such cases, the random walks can be launched within and confined to each connected component. The theorem also covers the case where edges in the graph have real-valued (non-binary) nonnegative weights. The theorem provides explicit closed-form expressions for the ergodic limits  $\bar{n}_{ij}^+$  and  $\bar{n}_{ij}^-$ .

**Theorem 1** (Ergodic limits of  $n_{ij}^+$  and  $n_{ij}^-$ )

*Let  $\mathcal{G}$  be a weighted graph with connected components  $\{\mathcal{G}_t\}_{t=1}^m$ , where for each  $t$ ,  $\mathcal{G}_t$  has  $n_t$  nodes and a nonnegative weighted adjacency matrix  $A_t$ . Let the VEC algorithm be executed on  $\mathcal{G}$  with parameters  $w$  and  $k$  and transition matrix  $W_t := D_t^{-1}A_t$  in component  $t$ , where  $D_t$  is a diagonal matrix with  $i$ th diagonal element  $D_{t,ii} = \sum_j A_{t,ij} =: d_i$ , i.e., the degree of node  $i$ . Then the ergodic limits  $\bar{n}_{ij}^+$ 's and  $\bar{n}_{ij}^-$ 's in Definition 2 exist and are given by*

$$\bar{n}_{ij}^+ = \begin{cases} \pi_i \sum_{v=1}^w (W_t^v)_{ij}, & \text{if } i, j \in \mathcal{G}_t; \\ 0 & \text{otherwise.} \end{cases}, \quad (4)$$

$$\bar{n}_{ij}^- = kw\pi_i\pi_j, \quad (5)$$

where  $\pi$  is a stationary distribution of the random walk with  $\pi_i = \frac{n_t}{\sum_t n_t} \frac{D_{t,ii}}{\sum_i D_{t,ii}}$  for each  $t$  and all  $i \in \mathcal{G}_t$ .

The proof of Theorem 1 is based on convergence results for irreducible Markov chains and is presented in Appendix A.1. The key ideas are as follows. For the positive pairs we expand the state-space of the Markov chain and show that it is irreducible. This implies that the long term average of distributions converges to the stationary distribution. For the negative pairs the major obstacle is to deal with the second-layer of randomness conditioned on the positive samples. We overcome this difficulty by applying McDiarmid's inequality conditionally to establish almost complete convergence.

Theorem 1 states that the ergodic limits  $\bar{n}_{ij}^+$  and  $\bar{n}_{ij}^-$  can be evaluated directly without having to actually launch any random walks. The additional randomness from the random walks and dependence on the algorithm parameters  $r$  and  $\ell$  are removed. As a result, the coefficients, in the form of ergodic limits, are deterministic functions of the graph adjacency matrix and two algorithm parameters  $w$  and  $k$ . Replacing the coefficients in Eq. (1) with their limiting values (scaled down by the factor  $1/(rn\ell)$ ) yields the following optimization problem that we name ErgoVEC:

**Definition 3** (ErgoVEC optimization problem).

$$\arg \min_{\{\mathbf{u}_i \in \mathbb{R}^d, i \in \mathcal{V}\}} \sum_{(i,j) \in \mathcal{V}^2} \left[ \bar{n}_{ij}^+ \sigma(+\mathbf{u}_i^\top \mathbf{u}_j) + \bar{n}_{ij}^- \sigma(-\mathbf{u}_i^\top \mathbf{u}_j) \right] \quad (6)$$

A practical approach to compute the embedding vectors of ErgoVEC can be described as follows. Given a graph

and algorithm parameters  $w$  and  $k$ , first calculate  $\bar{n}_{ij}^+$ 's and  $\bar{n}_{ij}^-$ 's using Theorem 1. Then use them to solve the ErgoVEC optimization problem in Definition 3 via stochastic gradient descent to find embedding vectors  $\mathbf{u}_i$ 's. A neural-network implementation is described in Section 4.2 and Appendix B.

ErgoVEC calculates the coefficients of the optimization objective in a more principled way compared to VEC and completely bypasses the random walk sampling process. The  $r$  and  $\ell$  algorithm parameters of VEC are not needed at all in ErgoVEC. However, when the graph is dense or  $w$  is large, evaluating  $\bar{n}_{ij}^+$  from Eq. (4) can be computationally very expensive. In these cases,  $n_{ij}^+$  computed from random walks could serve as an approximation. Thus VEC can be viewed as a practical approximation to the more principled ErgoVEC.

**Relationship to modularity maximization.** When the graph is connected and we set  $w = 1$ , Eq. (6) reduces to

$$\arg \min_{\{\mathbf{u}_i, i \in \mathcal{V}\}} \sum_{(i,j) \in \mathcal{V}^2} \left[ A_{ij} \sigma(\mathbf{u}_i^\top \mathbf{u}_j) + k \frac{d_i d_j}{\sum_k d_k} \sigma(-\mathbf{u}_i^\top \mathbf{u}_j) \right], \quad (7)$$

where  $d_i$  denotes the degree of node  $i$ . If instead we set  $\sigma(t) := t$ , and constrain the embedding vectors so that for all  $i, j$ ,  $\mathbf{u}_i^\top \mathbf{u}_j \in \{0, 1\}$ , then the minimization becomes equivalent to the modularity maximization problem [37] for two communities defined by

$$\arg \max_{\{y_i \in \{0, 1\}, i \in \mathcal{V}\}} \sum_{(i,j) \in \mathcal{V}^2} \left( A_{ij} - \frac{d_i d_j}{2|\mathcal{E}|} \right) 1(y_i = y_j)$$

where  $y_i$  denotes the community assignment for node  $i$ ,  $1(\cdot)$  is the indicator function, and  $|\mathcal{E}|$  is the number of edges. This is often relaxed to (7) and solved via spectral approaches followed by clustering [38].

### 3.2 Walk-distance weighting and large $r$ asymptotics

**Walk-distance-weighted count statistics:** In VEC,  $n_{ij}^+$  is the count of all instances where node  $i$  appears within  $w$  steps of node  $j$  in all the random walks. Instances where nodes  $i$  and  $j$  appear exactly 1 step from each other and instances where they appear in exactly  $w$  steps from each other, both contribute a count of 1 to the value of  $n_{ij}^+$ . A nuanced alternative must account for the number of steps between appearances of nodes.

As a general approach to construct such a statistic, we propose associating a walk-distance weight  $\alpha_v$  to the counts of instances of node pairs that occur *exactly*  $v$  steps from each other. With this modification, the walk-distance-weighted positive-pair counts will become  $n_{ij}^+ := \sum_{v=1}^{\infty} \alpha_v n_{ij}^+(v)$ , where  $n_{ij}^+(v)$  is the count of instances where node  $i$  appears *exactly*  $v$  steps from node  $j$  in all the random walks. The original count statistic for positive pairs can be recovered as a special case of our proposed general framework by choosing  $\alpha_v = 1$  for all  $v \leq w$  and  $\alpha_v = 0$  for all  $v > w$ . Choosing a nonnegative decreasing sequence of walk-distance weights  $\alpha_v$  can be viewed as providing a ‘‘soft cutoff’’ for the bigram counts when compared to the ‘‘hard cutoff’’ of the original counts.

To compute walk-distance weighted counts for negative-pairs, we propose the following modification to the original negative sampling process. For each positive pair of nodes

that occur exactly  $v$  steps apart, we append  $k$  node pairs drawn in an IID manner exactly as in the original sampling process. However, these  $k$  negative node pairs will now contribute the value  $\alpha_v$  to the walk-distance weighted negative-pair counts as opposed to the value of 1 previously.

**Large  $r$  asymptotics:** The effect of increasing  $\ell$  is similar to that of increasing  $r$ . In a random walk on a graph, the choice of the next node depends only on the current node. From this point of view, we may loosely visualize a long random walk as being formed by joining many shorter segments which are nearly independent random walks. In this sense, an infinitely long random walk is similar to an infinite sequence of short random walks with each starting node chosen from the stationary distribution of the Markov chain. Thus, in addition to the large  $\ell$  asymptotics characterized in Theorem 1, we can also study other types of asymptotics such as  $r \rightarrow \infty$  or, more generally,  $\ell$  and  $r$  both going to infinity together in some manner.

The counterpart of Theorem 1 for the proposed walk-distance-weighted counts is the following general result which is proved in Appendix A.2.

#### Theorem 2 (Limits of walk-distance weighted counts)

Let  $\mathcal{G}$  be a weighted connected graph with  $n$  nodes and  $W$  be the probability transition matrix of the natural random walk on  $\mathcal{G}$  with stationary distribution  $\pi$ . Let the VEC algorithm be executed with  $\mathcal{G}$  as input, walk-distance weights  $\{\alpha_v\}_{v=1}^{\infty}$ , and negative sampling rate  $k$ . If  $\{\alpha_v\}_{v=1}^{\infty}$  is absolutely convergent, i.e.,  $\sum_{v=1}^{\infty} |\alpha_v| < \infty$ , the following limits of  $\bar{n}_{ij}^+$ 's and  $\bar{n}_{ij}^-$ 's exist in the almost sure sense:

1) When  $r$  is fixed and  $\ell \rightarrow \infty$  (ergodic limits):

$$\frac{1}{rn} \lim_{\ell \rightarrow \infty} \frac{n_{ij}^+}{\ell} = \pi_i \sum_{v=1}^{\infty} \alpha_v (W^v)_{ij}, \quad (8)$$

$$\frac{1}{rn} \lim_{\ell \rightarrow \infty} \frac{n_{ij}^-}{\ell} = k\pi_i \pi_j \sum_{v=1}^{\infty} \alpha_v. \quad (9)$$

2) When  $\ell$  is fixed and  $r \rightarrow \infty$ :

$$\frac{1}{\ell n} \lim_{r \rightarrow \infty} \frac{n_{ij}^+}{r} = \frac{1}{\ell n} \sum_{m=1}^n \sum_{v=1}^{\infty} \alpha_v (W^v)_{ij} \sum_{s=1}^{\ell-v} (W^s)_{mi}^{s-1}, \quad (10)$$

$$\frac{1}{\ell n} \lim_{r \rightarrow \infty} \frac{n_{ij}^-}{r} = \frac{k\pi_j^{(\ell)}}{\ell n} \sum_{m=1}^n \sum_{v=1}^{\infty} \alpha_v \sum_{s=1}^{\ell-v} (W^s)_{mi}^{s-1} \quad (11)$$

where  $\pi_j^{(\ell)} = \frac{1}{\ell} \sum_{u=1}^{\ell} \frac{1}{n} \mathbf{1}_n^\top W^{u-1} \mathbf{e}_j$ .

3) Double limits:

$$\frac{1}{n} \lim_{r \rightarrow \infty} \lim_{\ell \rightarrow \infty} \frac{n_{ij}^+}{r\ell} = \frac{1}{n} \lim_{\ell \rightarrow \infty} \lim_{r \rightarrow \infty} \frac{n_{ij}^+}{r\ell} = \pi_i \sum_{v=1}^{\infty} \alpha_v (W^v)_{ij}, \quad (12)$$

$$\frac{1}{n} \lim_{r \rightarrow \infty} \lim_{\ell \rightarrow \infty} \frac{n_{ij}^-}{r\ell} = \frac{1}{n} \lim_{\ell \rightarrow \infty} \lim_{r \rightarrow \infty} \frac{n_{ij}^-}{r\ell} = k\pi_i \pi_j \sum_{v=1}^{\infty} \alpha_v. \quad (13)$$

Theorem 2 is stated for a connected graph, but it holds for each connected component of a disconnected graph. The main changes are that the expressions  $\pi_i \sum_{v=1}^w (W^v)_{ij}$  and  $k\pi_i \pi_j w$  in Lemma 1 change to  $\pi_i \sum_{v=1}^{\infty} \alpha_v (W^v)_{ij}$  and

1. We follow the convention that when the upper limit of a summation is smaller than its lower limit, the sum is 0.

$k\pi_i\pi_j(\sum_{v=1}^{\infty}\alpha_v)$  respectively and the large- $r$  asymptotic limits of the normalized count statistics are also characterized. To establish these results, we need to assume that the walk-distance weight series is absolutely convergent. Walk-distance weighting makes it possible to realize different nonlinear functions of the transition matrix as the ergodic limit of count statistics, not just polynomial functions. For instance, if we choose  $\alpha_v = 1/v!$  for all  $v$ , then for all  $i, j$ ,  $\bar{n}_{ij}^+ = \pi_i(\exp\{W\})_{ij}$ , where  $\exp\{W\}$  denotes matrix exponential, and  $\bar{n}_{ij}^- = k\pi_i\pi_j e$ . We note that the large- $\ell$  asymptotic limits are independent of  $r$  but the large- $r$  asymptotic limits depend on  $\ell$ . Yet, the double limits where  $r$  is sent to infinity first before  $\ell$  equal the corresponding large- $\ell$  limits.

### 3.3 Reparameterized relaxations and their properties

In this subsection, we study matrix re-parameterizations and convex relaxations of the VEC and ErgoVEC optimization problems. We follow [35] and begin by defining the  $n \times n$  matrix  $X$  to be the Gram matrix of the node embedding vectors, i.e., for all  $i, j$ ,  $X_{ij} := \mathbf{u}_i^\top \mathbf{u}_j$ . Let  $N^+$  and  $N^-$  denote the  $n \times n$  matrices of the positive-pair and negative-pair counts respectively, i.e., for all  $i, j$ ,  $[N^+]_{ij} = n_{ij}^+$  and  $[N^-]_{ij} := n_{ij}^-$ . If

$$f(N^+, N^-, X) := \sum_{(i,j) \in \mathcal{V}^2} n_{ij}^+ \sigma(+X_{ij}) + \sum_{(i,j) \in \mathcal{V}^2} n_{ij}^- \sigma(-X_{ij}),$$

then the VEC optimization problem (Eq. (1)) reduces to the following equivalent optimization problem in the matrix variable  $X$  named GramVEC:

**Definition 4** (GramVEC optimization objective).

$$\operatorname{argmin}_{X \in \mathbb{S}_+^n, \operatorname{rank}(X) \leq d} f(N^+, N^-, X). \quad (14)$$

In Eq. (14), the constraint  $X \in \mathbb{S}_+^n$  arises from the fact that the Gram matrix of embedding vectors is real, symmetric, and positive semi-definite. The rank constraint comes from the fact that  $\mathbf{u}_i \in \mathbb{R}^d$ . The equivalence of the VEC (Eq. (1)) and GramVEC (Eq. (14)) optimization problems can be seen as follows. For any set of feasible  $\mathbf{u}_i$ 's in (1), setting  $X_{ij} = \mathbf{u}_i^\top \mathbf{u}_j$  for all  $i, j$ , makes  $X$  a rank- $d$  matrix in  $\mathbb{S}_+^n$  and yields the same cost as in (14). In the other direction, for any feasible choice of  $X$  in (14), let  $X = V_d^\top \Sigma_d V_d$  denote its rank- $d$  reduced SVD with diagonal  $\Sigma_d \in \mathbb{S}_+^d$  and define  $U = [\mathbf{u}_1, \dots, \mathbf{u}_n] := \sqrt{\Sigma_d} V_d$ . Then,  $X = U^\top U$  and for all  $i, j$ ,  $X_{ij} = \mathbf{u}_i^\top \mathbf{u}_j$  and  $\mathbf{u}_i \in \mathbb{R}^d$ , and we obtain the same cost in (1). The choices for the  $\mathbf{u}_i$ 's are not unique since  $X = U^\top F^\top F U$  for any real orthonormal matrix  $F$ . The  $\mathbf{u}_i$ 's are unique only up to a real orthonormal transformation, just as in (1).

The same re-parameterization can also be applied to ErgoVEC (Eq. (6)). Let  $\bar{N}^+$  and  $\bar{N}^-$  be  $n \times n$  matrices such that for all  $i, j$ ,  $[\bar{N}^+]_{ij} := \bar{n}_{ij}^+$  and  $[\bar{N}^-]_{ij} := \bar{n}_{ij}^-$ . Then the GramErgoVEC optimization problem is defined as follows.

**Definition 5** (GramErgoVEC optimization problem).

$$\operatorname{argmin}_{X \in \mathbb{S}_+^n, \operatorname{rank}(X) \leq d} f(\bar{N}^+, \bar{N}^-, X). \quad (15)$$

Although Eq. (14) and Eq. (15) are equivalent to Eq. (1) and Eq. (6), respectively, they are more convenient to work with and analyze. The matrix re-parameterization transfers the non-convexity from the objective function to the constraint set which makes it possible to relax or convexify the problem as we do next.

Relaxing all constraints on  $X$  in GramErgoVEC leads to the following optimization problem named GramErgoPMI (relaxing constraints in GramVEC similarly will yield a corresponding optimization problem GramPMI):

**Definition 6** (GramErgoPMI optimization problem).

$$\operatorname{argmin}_{X \in \mathbb{R}^{n \times n}} f(\bar{N}^+, \bar{N}^-, X). \quad (16)$$

In general, GramErgoPMI is not equivalent to GramErgoVEC. The relaxation enlarges the feasible set and the optimal solution may not satisfy the constraints in Eq.(15). Nonetheless, Eq. (16) admits a closed-form solution:

**Proposition 1**

Let  $X^*$  be the solution to Eq. (16). Then,  $X^*$  is unique, symmetric, and given by

$$X_{ij}^* = X_{ji}^* = \begin{cases} \ln \left( \frac{\bar{n}_{ij}^+}{\bar{n}_{ij}^-} \right) & \text{if } \bar{n}_{ij}^+ \neq 0; \\ -\infty & \text{if } \bar{n}_{ij}^+ = 0. \end{cases} \quad (17)$$

Let  $p_\ell(i, j)$  denote the probability that a randomly sampled pair from the positive set  $\mathcal{D}_+$  equals  $(i, j)$  and let  $p_{\ell 1}(i)$  and  $p_{\ell 2}(j)$  denote, respectively, the first- and second-component marginal probabilities of  $i$  and  $j$  that are consistent with the joint pmf  $p_\ell(i, j)$ .<sup>2</sup> Let  $\text{PMI}_\ell(i, j) := \ln \left( \frac{p_\ell(i, j)}{p_{\ell 1}(i)p_{\ell 2}(j)} \right)$  denote the **Pointwise Mutual Information (PMI)** of  $(i, j)$  [36]. Then for all  $i, j$ ,

$$X_{ij}^* = \lim_{\ell \rightarrow \infty} \text{PMI}_\ell(i, j) - \ln k. \quad (18)$$

The proof of Proposition 1 is presented in Appendix A.3. Although  $X^*$  is symmetric, there is no guarantee that  $X^*$  will satisfy the positive semi-definiteness constraint of GramErgoVEC, let alone the rank constraint. Without positive semi-definiteness, the square root of  $X^*$ 's nonzero singular values will be imaginary and there will not exist any real-valued embedding vectors, even in  $\mathbb{R}^n$ , whose Gram matrix equals  $X^*$ . A practical solution then is to compute the  $\ell_2$  projection of  $X^*$  into the rank- $d$  real positive semi-definite cone of  $n \times n$  matrices and factorize the projected matrix  $\text{Proj}(X^*, d)$  to get embeddings. Still, there is no guarantee that the projected matrix  $\text{Proj}(X^*, d)$  or the embeddings  $\mathbf{u}_i$ 's obtained from it will be a solution to GramErgoVEC. We compare the GramErgoVEC and GramErgoPMI embedding vectors experimentally in Section 5.

An alternative approach to deal with the non-convexity of the GramErgoVEC is to replace the non-convex rank constraint with a nuclear norm constraint which is convex. The nuclear norm  $\|X\|_*$  of a matrix  $X$  is defined as the sum of its singular values. Its relationship with the rank of a matrix has been extensively studied in the literature. For example, nuclear norm level sets have been shown to be the convex-envelope of rank level sets [39]. A bounded

<sup>2</sup> That is,  $p_{\ell 1}(i) := \sum_{j \in \mathcal{V}} p_\ell(i, j)$  and  $p_{\ell 2}(j) := \sum_{i \in \mathcal{V}} p_\ell(i, j)$ . Note that  $p_\ell(i, j)$  may not be symmetric when  $\ell$  is finite.

nuclear norm constraint has been used as a proxy for a bounded rank constraint in a number of problems such as low rank matrix completion [40], tensor robust PCA [41], and compressed sensing [42]. For some problems there exists an exact equivalence between these constraints but the conditions under which this occurs varies from problem to problem. Relaxing the rank constraint of GramErgoVEC via a bound on the nuclear norm leads to the following optimization problem that we term NucGramErgoVEC (we can similarly define NucGramVEC):

**Definition 7** (NucGramErgoVEC optimization problem).

$$\operatorname{argmin}_{X \in \mathbb{S}_+^n, \|X\|_* \leq \nu_n} f(\bar{N}^+, \bar{N}^-, X) \quad (19)$$

A larger  $\nu_n$  implies a looser nuclear norm constraint. When  $\nu_n$  goes to  $\infty$ , the solution to NucGramErgoVEC will approach the  $\ell_2$  projection of the GramErgoPMI solution onto the real positive semi-definite cone of  $n \times n$  matrices. However, when  $\nu_n$  goes to 0, the solution to NucGramErgoVEC will reduce to the all zeros matrix which has rank 0. In general, we should expect smaller values of  $\nu_n$  to yield solutions with approximately lower rank. This is corroborated by our experiments in Section 5.1. Thus the rank of the solution matrix, and consequently the dimension of the embedding vectors, can be controlled by the value of  $\nu_n$ . Note that we allow the nuclear norm threshold  $\nu_n$  to depend on  $n$ . In Section 3.4 we show that the nuclear norms of the solutions to GramErgoPMI and GramErgoVEC for idealized graphs that have community structure scale with  $n$  as  $\nu_n = \Theta(n)$ .

There are numerous algorithms available to numerically compute a solution to the NucGramErgoVEC optimization problem. We modified and implemented Hazan's algorithm [43] to generate all our experimental results.

### 3.4 Embeddings of expected SBM graphs

As an important step toward analyzing concentration properties of embeddings, in this section we study the embedding solutions of different optimization problems focusing on the expected graph of a two-community SBM. Such graphs have an ideal community structure: all edges between nodes belonging to any specified pair of communities have the same edge weight. Our main result is summarized in the following theorem:

**Theorem 3** (Embeddings of an expected SBM graph)

Let  $\mathcal{G}$  be an SBM graph with  $n = 2m$ ,  $m \geq 2$ , nodes and two balanced communities. For all  $i \in \mathcal{V}$ , let  $y_i \in \{0, 1\}$  denote the community label of node  $i$ . Let  $a$  and  $b$  be the edge forming probabilities for within- and cross-community edges, respectively, with  $a > \frac{m}{m-1}b$ . Let  $\mathbf{E}[\mathcal{G}]$  denote the expected graph and  $\bar{n}_{ij}^+$ 's and  $\bar{n}_{ij}^-$ 's the ergodic limits for  $\mathbf{E}[\mathcal{G}]$  as in Definition 2, with  $k \geq 1$  and  $w \geq 1$ . Let

$$X^*(\mathcal{H}) := \operatorname{argmin}_{X \in \mathcal{H}} f(\bar{N}^+, \bar{N}^-, X), \quad (20)$$

where  $\mathcal{H} \subset \mathbb{R}^{n \times n}$ . Let  $\mathcal{E}_0 := \{(i, j) \in \mathcal{E} : y_i = y_j\}$  and  $\mathcal{E}_1 := \mathcal{E} \setminus \mathcal{E}_0$ . Then:

1) *Structure of ergodic limits: The values of  $\bar{n}_{ij}^+$  and  $\bar{n}_{ij}^-$  depend only on the community membership of  $(i, j)$ , i.e.,*

$$\bar{n}_{ij}^+ = \begin{cases} \alpha_1, & \text{if } (i, j) \in \mathcal{E}_0 \\ \alpha_2, & \text{if } (i, j) \in \mathcal{E}_1 \\ \alpha_3, & \text{if } i = j \end{cases}$$

$$\bar{n}_{ij}^- = \beta, \quad \forall (i, j),$$

where  $\beta = \frac{kw}{n^2}$  and  $\alpha_i = C_i/n^2 + o(1/n^2)$  for  $i = 1, 2, 3$ , where  $C_i$ 's are functions of only  $a, b$  and  $w$ .

2) *GramErgoPMI solution: Let  $\mathcal{H} = \mathbb{R}^{n \times n}$ . Then  $X^*(\mathcal{H})$  has full rank with the same structure as  $\bar{n}_{ij}^+$  with*

$$X^*(\mathbb{R}^{n \times n}) = \begin{cases} \ln\left(\frac{\alpha_1}{\beta}\right), & \text{if } (i, j) \in \mathcal{E}_0 \\ \ln\left(\frac{\alpha_2}{\beta}\right), & \text{if } (i, j) \in \mathcal{E}_1 \\ \ln\left(\frac{\alpha_3}{\beta}\right), & \text{if } i = j \end{cases}$$

3) *NucGramErgoVEC solution for  $\nu_n = \infty$ : Let  $\mathcal{H} = \mathbb{S}_+^n$  and  $\nu_1 := \ln\left(\frac{\bar{\alpha}_{13} + \beta}{\alpha_2 + \beta}\right)$  where  $\bar{\alpha}_{13} := \frac{m-1}{m}\alpha_1 + \frac{1}{m}\alpha_3$ . Then  $X^*(\mathcal{H})$  has rank 1. Moreover, if  $\nu_n = \nu_1 n$ , then*

$$X^*(\mathbb{S}_+^n) = \begin{cases} \nu_1, & \text{if } (i, j) \in \mathcal{E}_0 \text{ or } i = j \\ -\nu_1, & \text{if } (i, j) \in \mathcal{E}_1. \end{cases}$$

4) *The nuclear norms of  $X^*(\mathbb{R}^{n \times n})$  and  $X^*(\mathbb{S}_+^n)$  scale with  $n$  as  $\Theta(n)$ .*

The full proof of Theorem 3 is presented in Appendix A.4, but the key ideas are as follows. Part 1) holds as a result of Theorem 1 which characterizes the ergodic limits of normalized bigram counts in terms of the random walk transition matrix. Since the transition matrix of the expected graph has block-wise constant entries, this property is carried forward to the normalized ergodic counts. Part 2) follows directly from part 1) and Proposition 1. Part 3) is the major piece of this theorem and is proved via an intricate analysis of the structure of the solution.

The  $\Theta(n)$  scaling of nuclear norms of  $X^*(\mathbb{R}^{n \times n})$  and  $X^*(\mathbb{S}_+^n)$  arises from the block structure and the fact that all entries are of constant order. Theorem 3 also shows that  $\bar{N}^+$  and  $X^*(\mathbb{R}^{n \times n})$ , the solution to GramErgoPMI, have the structure of a rank 2 matrix minus a scalar multiple of the identity matrix making them full rank. In contrast,  $X^*(\mathbb{S}_+^n)$ , the solution to NucGramErgoVEC, has rank 1 due to the positive semi-definite constraint.

Part 3) of Theorem 3 may seem surprising on first glance because we are getting a rank 1 solution without any rank or nuclear norm constraints. The surprise dissipates when we note that the solution is for an expected graph which has an ideal community structure. Such a result would not hold true for a random graph realization. Nonetheless, we can directly obtain the GramErgoVEC and NucGramErgoVEC solutions for the expected graph from part 3) of Theorem 3:

**Corollary 1.** *Under the assumptions of Theorem 3,*

1) *GramErgoVEC solution for any positive rank: Let  $\mathcal{H} = \{X \in \mathbb{S}_+^n : \operatorname{rank}(X) \leq d\}$  with  $d \geq 1$ . Then  $X^*(\mathcal{H})$  has rank 1 and equals the GramErgoVEC solution for  $\nu_n = \infty$  characterized in part 3) of Theorem 3.*

2) *NucGramErgoVEC* solution for all  $\nu_n \geq \nu_1 n$ : Let  $\mathcal{H} = \{X \in \mathbb{S}_+^n : \|X\|_* \leq \nu_n\}$ . If  $\nu_n \geq \nu_1 n$ , then  $X^*(\mathcal{H})$  has rank 1 and equals the *GramErgoVEC* solution for  $\nu_n = \infty$  characterized in part 3) of Theorem 3.

When given the expected SBM graph as input, *GramErgoVEC* and *NucGramErgoVEC* will return a Gram matrix of rank 1 or 2 which when factorized will provide two distinct embedding vectors, each representing one community in the original graph. In short, the algorithms will give embeddings that are perfectly separated across communities and perfectly concentrated within communities.

In part 2) of Corollary 1, with  $\nu_n \geq \nu_1 n$ , the nuclear norm constraint becomes inactive. Suppose that  $\nu_n = \nu_0 n$ . If  $\nu_0 \leq \nu_1$ , we conjecture that the solution will scale proportionally with  $\nu_0$ :

### Conjecture 1

Under the assumptions of Theorem 3, let  $\mathcal{H} = \{X \in \mathbb{S}_+^n : \|X\|_* \leq \nu_0 n\}$ . If  $\nu_0 < \nu_1$ , then

$$X^*(\mathcal{H}) = \begin{cases} \nu_0, & \text{if } (i, j) \in \mathcal{E}_0 \text{ or } i = j \\ -\nu_0, & \text{if } (i, j) \in \mathcal{E}_1. \end{cases}$$

Conjecture 1 makes an assertion about the solution to the *NucGramErgoVEC* optimization problem when the nuclear norm constraint is active. For a suitable nuclear norm constraint, we conjecture that the solution will be a scaling of the solution in part 3) of Theorem 3. If this conjecture holds, then, we can conclude that the solution to *NucGramErgoVEC* is always rank 1 with perfect separation of the communities regardless of the sparsity level of the graph. This would provide a solid starting point for analyzing the concentration properties of solutions as  $n$  increases to  $\infty$ .

## 4 EXPERIMENTAL SETUP

In the following sections, we compare the node embeddings from different algorithms qualitatively and quantitatively through extensive experiments. This section details our experimental setup, including the generation of random graphs, details of implementation and parameter choices for algorithms, and evaluation metrics for embedding vectors. Section 5.1 explores the geometric properties of embedding vectors for a fixed graph size. Specifically, we study how the nuclear norm linear scaling factor  $\nu_0$  of the nuclear norm limit  $\nu_n = \nu_0 n$  impacts the embedding geometry in *NucGramErgoVEC*. In Section 5.2 we investigate how embedding vectors change as the graph size  $n$  increases and whether they tend to concentrate.

### 4.1 SBM graph generation

For simplicity, we focus on assortative, equal-sized, planted-partition SBM graphs with 2 communities. We generate random graphs with  $n = 100, 200, 500, 1000$  nodes. We consider two scaling regimes for the within-community edge forming probability  $p$  and the cross-community edge forming probability  $q$ : 1) *Linear regime*: Here,  $p$  and  $q$  are held constant, with values  $p = 0.6$  and  $q = 0.06$ , for all graph sizes. As a result, the expected node degree scales linearly with  $n$ . 2) *Logarithmic regime*: Here,  $p$  and  $q$  diminish with

increasing graph size  $n$  as a multiple of  $(\ln n)/n$ , specifically as  $p = 9(\ln n)/n$  and  $q = 2(\ln n)/n$ . The expected node degree then increases proportionally with  $\ln n$ . Our choices of scaling factors  $\bar{p} = 9$  and  $\bar{q} = 2$  in the logarithmic regime ensure that the information-theoretic threshold for exact community recovery for two communities, given by  $\sqrt{\bar{p}} - \sqrt{\bar{q}} > \sqrt{2}$  [30], is slightly surpassed.

To ensure graph connectivity and improve community detection performance, we follow the prescription in [44] and apply  $\varepsilon$ -smoothing to all generated graphs. For a given  $\varepsilon \geq 0$ , the  $\varepsilon$ -smoothing of  $\mathcal{G}$ , denoted by  $\mathcal{G}^\varepsilon$ , is the weighted complete graph with adjacency matrix  $A^\varepsilon$ , where for all  $i, j$ ,  $A_{ij}^\varepsilon := A_{ij} + \varepsilon$ . In addition to analytical convenience, graph-smoothing also improves the performance of spectral clustering in practice [44]. The ergodic limit of the coefficients under  $\varepsilon$ -smoothed graphs can be computed by Theorem 1 with a modified probability transition matrix  $W_\varepsilon$  that corresponds to the new graph.

The optimal choice of  $\varepsilon$  for various performance measures such as signal-to-noise ratio, community detection accuracy, etc., varies across different random graph realizations and embedding algorithms. Since our focus is on embedding algorithms, we apply  $\varepsilon$ -smoothing to all random graphs that we generate with the fixed choice  $\varepsilon = 1/10n$  instead of optimizing  $\varepsilon$  for each algorithm and each performance measure. This is a relatively small value of  $\varepsilon$  as it changes the degrees from the original graph by at most  $1/10$  whereas the expected degrees in the linear and logarithmic regimes scale as  $n$  and  $\ln(n)$ , respectively.

### 4.2 Algorithm parameter choices and implementation

We implement and compare *ErgoVEC*, *GramErgoPMI*, and *NucGramErgoVEC*, that were proposed in Section 3, with *VEC* and *Spectral Clustering (SC)*. *SC* serves as a classical benchmark due to its widespread usage.

**Parameter choices.** For ease of visualization, we compute and plot 2-dimensional embedding vectors for all algorithms, i.e.  $d = 2$ .

For all algorithms other than *SC*, we set the window size to  $w = 8$ , and the negative sampling rate to  $k = 5$ . Settings that are specific to *SC*, *VEC*, and *NucGramErgoVEC* are as follows:

- 1) **SC**: We use the first two eigenvectors of the symmetrically normalized Laplacian, i.e.,  $D^{-1/2}AD^{-1/2}$ , [45]. Since the unit-norm eigenvectors are in  $\mathbb{R}^n$ , their components, and therefore also the node embedding vectors, scale as  $O(1/\sqrt{n})$ . In order to simultaneously visualize and compare embedding vectors across different values of  $n$ , we scale them by  $\sqrt{n}$ .
- 2) **VEC**: We launch  $r = 10$  walks starting from each node, each of length  $\ell = 100$ .
- 3) **NucGramErgoVEC**: In light of the linear scaling of the nuclear norm of the gram matrix of node embedding vectors for an expected SBM graph, (cf. Theorem 3 and Corollary 1), we set  $\nu_n = \nu_0 n$  and sweep  $\nu_0$  over the range 0.018 through 0.216, in steps of 0.018, in order to illustrate changes in the geometric structure of embedding vectors (cf. Fig. 4 and Fig. 5).

**Implementation of VEC and ErgoVEC.** Both *VEC* and *ErgoVEC* have non-convex objectives for which there is

no optimization procedure available which guarantees convergence to a global minimizer. A practical way forward is to use Stochastic Gradient Descent (SGD). We can consider two distinct approaches for implementing SGD in VEC or ErgoVEC 1) Map them to an equivalent Word2Vec problem by identifying nodes as words and random walks as sentences and then obtain word embeddings using a Word2Vec package such as Gensim [46]. 2) Reformulate each optimization problem as the training of a neural network with an appropriate architecture and cost and then train the neural network using a neural network package such as Keras [47]. Since the Gensim package cannot handle non-integer coefficients that arise in ErgoVEC, we use Keras to implement ErgoVEC and VEC. Details of our neural network implementation are presented in Appendix B. In section 5.2 we also compare the Gensim and Keras implementations of VEC.

Since both algorithms involve optimizing non-convex objectives, convergence is not guaranteed. We assess the convergence of the objective function value and the solution by evaluating the change in the objective function value and the embeddings after each epoch. To quantify the change in embeddings, we first perform a Procrustes alignment of the embedding solutions from successive epochs and then compute the Frobenius norm of the difference between the aligned sets of embeddings. We found that the change in the objective function value diminishes as the number of epochs increases, but the change in the corresponding embeddings retains a small fluctuation after diminishing initially (*cf.* Appendix B). This suggests that although the objective function value converges, the embeddings may be oscillating around a local optimizer. Note that the Keras implementation, which implements SGD, is not designed to minimize changes in the solution (arguments), but rather in the objective function.

**Implementation of NucGramErgoVEC.** We use Hazan’s algorithm [43] (suitably modified to handle inequality constraints) to solve the NucGramErgoVEC optimization problem. The algorithm is iterative and is guaranteed to converge to the global minimum. We also empirically confirmed the convergence of both the objective function value and the solution matrix using the approach used for VEC and ErgoVEC described above.

We note, however, that even though Hazan’s algorithm is guaranteed to converge to a global minimum, its convergence rate is slow. In order to improve convergence speed, we initialize with the GramErgoPMI solution suitably scaled to fit the nuclear norm limit. We also terminate the algorithm after 1000 iterations which is adequate for all our experiments.

### 4.3 Visualization and performance evaluation

**Representation and alignment of embeddings:** The absolute positions and orientations of embedding vectors may vary across algorithms, graph realizations, and graph sizes. Even for a given graph and algorithm the embedding vectors are not unique due to invariance of the objective function under orthogonal transformations. However clustering and separation properties of embeddings only depend on the *relative* positions and orientations of embedding vectors.

Thus, in order to visualize and simultaneously compare different embeddings qualitatively and quantitatively, we first represent the embedding vectors using their SVD coordinates and then align them with Procrustes analysis. Specifically, let  $U$  be an  $n \times d$  matrix whose  $i$ -th row  $u_i^\top$  is the embedding of node  $i$ . Let  $U = \tilde{U}\Sigma\tilde{V}^\top$  be the SVD decomposition of  $U$ . Then, the SVD coordinates of the embedding vectors are given by  $U\tilde{V}$ . To align two sets of embedding points  $U_1$  and  $U_2$ , we do Procrustes analysis, which finds the orthogonal matrix  $P$  that minimizes  $\|U_1 - U_2P\|_F^2$ . The aligned points are given by  $U_1$  and  $U_2P$ .

**Quantifying community separation.** We quantify the separation of nodes belonging to different communities using a signal-to-noise ratio (SNR) measured along the line joining the embedding centroids of the two communities. Specifically, for embeddings of nodes in community  $i$  ( $i = 1, 2$ ), let  $\hat{\mu}_i$  denote their empirical mean and  $\hat{K}_i$  their empirical covariance matrix. Then we define SNR-1D as follows:

$$\text{SNR-1D} := \frac{\|\hat{\mu}_1 - \hat{\mu}_2\|_2^2}{\frac{1}{2}(\hat{\eta}_1^2 + \hat{\eta}_2^2)}$$

where  $\hat{\eta}_i^2 := (\hat{\mu}_1 - \hat{\mu}_2)^\top \hat{K}_i (\hat{\mu}_1 - \hat{\mu}_2)$  is the empirical variance of the embeddings of nodes in community  $i$  when projected onto the line joining the embedding centroids of the two communities.

## 5 NODE EMBEDDING GEOMETRY OF SBM GRAPHS

In this section, we present and compare embeddings for SBM graphs produced by different algorithms and how they are positioned relative to the embeddings of the expected graph (indicated by black crosses in all our plots). Section 5.1 focuses on the comparing the geometry of embeddings across different algorithms and parameter choices whereas Section 5.2 focuses on the large graph asymptotic behavior of embeddings.

### 5.1 Geometry of embeddings

The geometry of node embedding vectors from SC, VEC, GramErgoPMI and ErgoVEC are shown in Fig. 2. We plot the 2D 95% confidence ellipse for each embedding cluster (red-colored elliptical curves) based on a maximum likelihood Gaussian fit to the data. For SC, ErgoVEC, and GramErgoPMI, the embeddings of the expected SBM graph are two distinct points (characterized in Theorem 3 and Corollary 1) which are marked as black crosses in Fig. 2. Since the VEC objective is based on empirical skip bigram counts from random walks, for a finite random walk length  $\ell$ , the embedding vectors of the expected graph will not collapse to two just distinct points, but will be distributed around the embedding vectors of ErgoVEC which are indicated as black crosses in the VEC subplot of Fig. 2.

In Fig. 2, we observe that in all four algorithms, the node embeddings in each cluster have an elliptical distribution around the cluster centroid and they can be perfectly separated linearly by the bisector of the line joining the two cluster centroids. However, the major axes of the SC embedding ellipses are nearly parallel to their inter-centroid line whereas the major axes of embedding ellipses in the other three algorithms are nearly perpendicular to their respective inter-centroid lines.

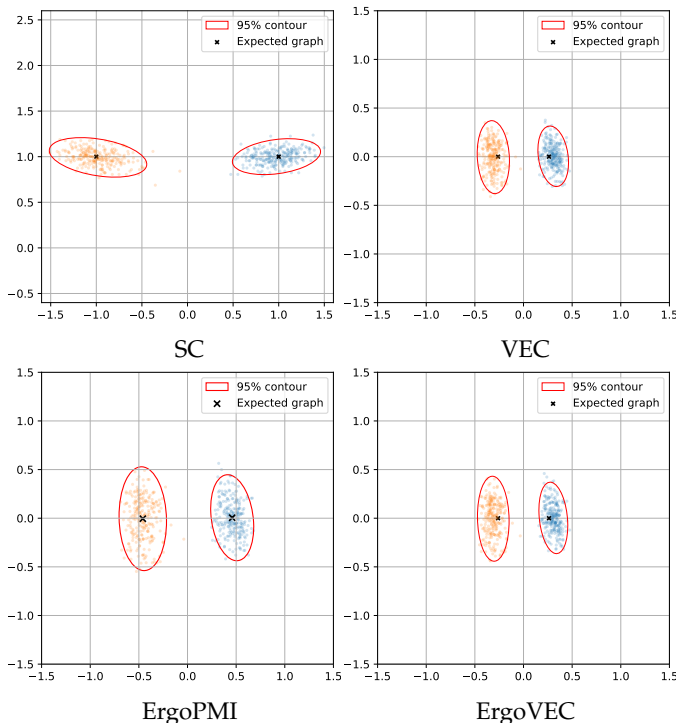


Fig. 2: 2D-visualization of embeddings for SC, VEC, ErgoPMI and ErgoVEC. All four algorithms receive the same graph input with  $n = 500$  nodes generated using within community edge probabilities  $p = 9 \ln n/n$  and across community edge probabilities  $q = 2 \ln n/n$ .

We also notice that the embedding ellipses of VEC and ErgoVEC in Fig. 2 look very similar. This is to be expected as the ErgoVEC objective is exactly the large  $\ell$  ergodic limit of the VEC objective introduced in Section 3.1. To empirically confirm that the ErgoVEC embeddings converge to the VEC embeddings in the large  $\ell$  limit, in Fig. 3 we plot the distance between VEC and ErgoVEC embeddings for increasing values of  $\ell$  and different graph sizes.

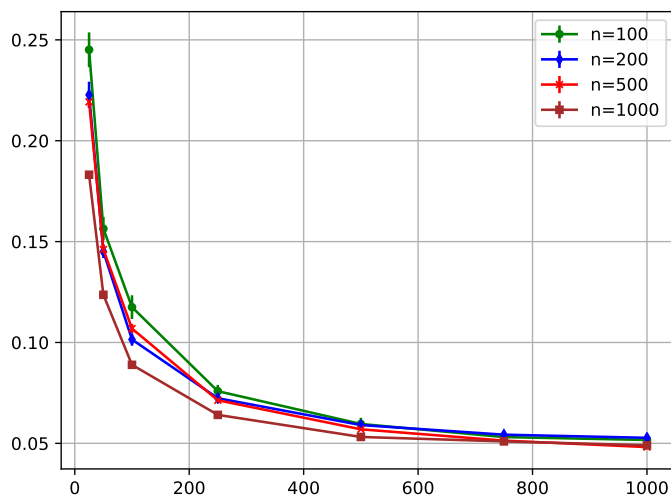


Fig. 3: Convergence of VEC embeddings: Frobenius norm distance between the Gram matrices of VEC and ErgoVEC versus random walk length  $\ell$ .

In order to measure the distance between embeddings up to any orthogonal transformation, we use the normalized

Frobenius norm distance between the Gram matrices of the embeddings. For each  $n$ , graphs are generated using within-community edge probabilities  $p = 9 \ln n/n$  and across-community edge probabilities  $q = 2 \ln n/n$ . The plot depicts the mean distance and associated error bar averaged over 5 independent graph realizations. Observe that for all graph sizes  $n = 100, 200, 500, 1000$ , as the length of the random walk increases, the distance between VEC and ErgoVEC Gram matrices shrinks. However, due to the non-convexity of the VEC and lack of global convergence guarantees for SGD methods used to optimize VEC and ErgoVEC objectives (cf. Section 4.2), the distance seems to be strictly bounded away from zero even at  $\ell = 1000$ . However, the positive and negative  $w$ -skip bigram counts and the objective function of VEC do converge to their respective ErgoVEC counterparts as  $\ell$  increases to infinity.

We now discuss the embedding geometry of NucGramErgoVEC. We separated this discussion from the previous four algorithms because although the embeddings of NucGramErgoVEC are also elliptically distributed and separate well into two clusters, the specific shape depends on the nuclear norm linear scaling factor  $\nu_0$  as shown in Fig. 4.

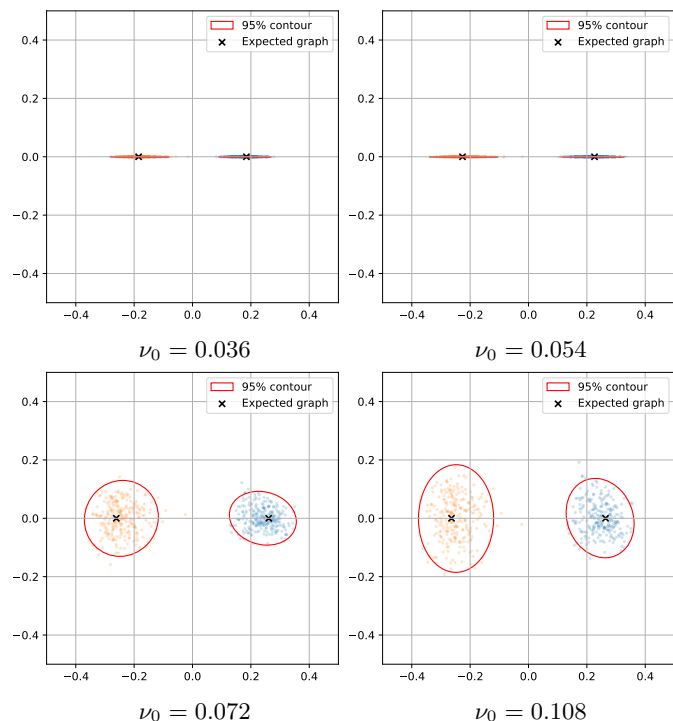


Fig. 4: 2D-visualization of NucGramErgoVEC embeddings for different nuclear norm linear scaling factors. The input graph is the same as the one used in Fig. 2.

When  $\nu_0$  is very small, the embeddings are one dimensional (cf. Fig. 4(a)). As  $\nu_0$  increases slightly, the embeddings remains one dimensional but spread out within each cluster and the cluster centroids move apart (cf. Fig. 4(b)). This continues until  $\nu_0$  reaches a threshold. When  $\nu_0$  increases beyond the threshold, the embeddings stop extending in the first dimension and start to spread in the second dimension (cf. Fig. 4(c)(d)).

In order to obtain a more quantitative understanding of

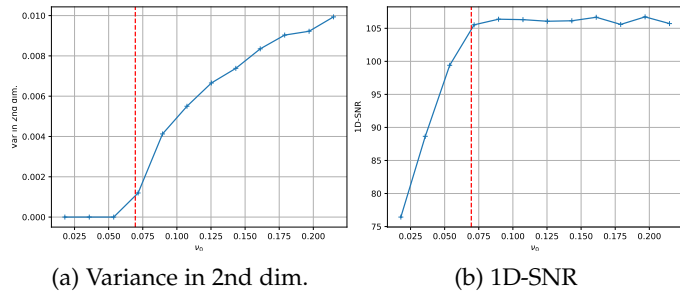


Fig. 5: The change in (a) variance in second dimension and (b) 1D-SNR of NucGramErgoVEC embeddings as the nuclear norm linear scaling factor  $\nu_0$  increases. The input graph has  $n = 200$  nodes generated with within community edge probabilities  $p = 9 \ln n/n$  and across community edge probabilities  $q = 2 \ln n/n$ .

how  $\nu_0$  influences the embedding geometry, we plot the 1D-SNR of embeddings and their variance in the second dimension for a range of values of  $\nu_0$  in Fig. 5.

Fig 5(a) shows how the variance of embeddings in the second dimension changes as the nuclear norm linear scaling factor  $\nu_0$  increases. When  $\nu_0$  is very small, the variance in the second dimension is 0, suggesting that embeddings are 1 dimensional. As  $\nu_0$  increases, the variance in the second dimension remains zero until  $\nu_0$  crosses a threshold that lies somewhere between  $\nu_0 = 0.054$  and  $\nu_0 = 0.072$  and then thereafter the variance increases monotonically. The exact value of  $\nu_0$  where the second dimension variance emerges depends on the input graph in general and specifically on the edge forming probability.

Fig 5(b) shows how  $\nu_0$  affects 1D-SNR of the embeddings. Here, we see a clear increase of 1D-SNR as  $\nu_0$  increase from 0.018 to 0.072. A relative maximum level is reached when the  $\nu_0$  is around the transition point where the second dimension variance emerges. Beyond the transition point, the 1D-SNR holds steady around the maximum level. These properties are consistent with our observations for Fig. 4.

## 5.2 Concentration of embeddings

After understanding how the geometry of embeddings of a single graph differs across embedding algorithms and changes with  $\nu_0$ , in this section, we explore how the embeddings change as  $n$ , the number of nodes, increases. To focus on the effect of increasing the number of nodes, throughout this section, we fix the scaling factors of edge forming probabilities within and across communities in each set of experiments. In addition, we omit the results of VEC because of their similarity to ErgoVEC (*cf.* Fig. 3). As we will see, the asymptotic behavior of embeddings largely depends on the edge forming probability.

To gain a qualitative perspective, we first plot the embeddings and their 95% Gaussian contours for graph sizes  $n = 100, 500, 1000$  for each algorithm. Fig. 6 shows the embedding contours in the linear degree scaling regime. We can see that all the contours shrink as  $n$  increases. This suggests that empirically, the embeddings of all the four algorithms concentrate to their centroids. In the logarithmic degree scaling regime, as shown in Fig. 7, the Gaussian contours for different graph sizes mostly overlap on top of each other, suggesting a convergence in distribution as

opposed to a concentration that we observed in the linear regime.

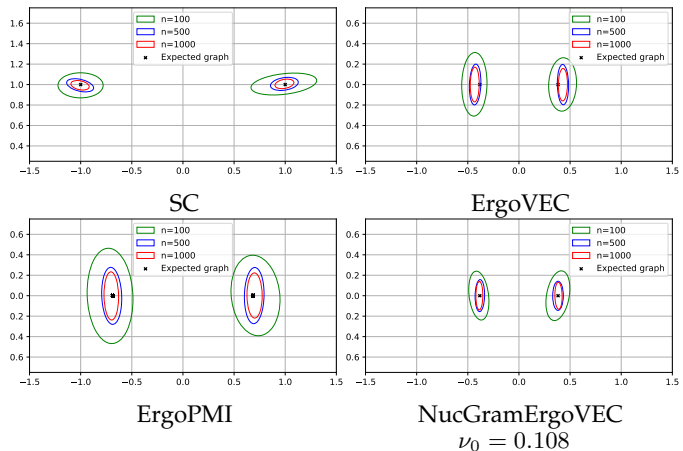


Fig. 6: 95% Gaussian contours of 2D-embeddings from four algorithms in the linear regime. All algorithms receive the same sets of graphs with  $n = 100, 500$  and  $1000$  nodes generated using within-community edge probabilities  $p = 0.6$  and across community edge probabilities  $q = 0.06$ .

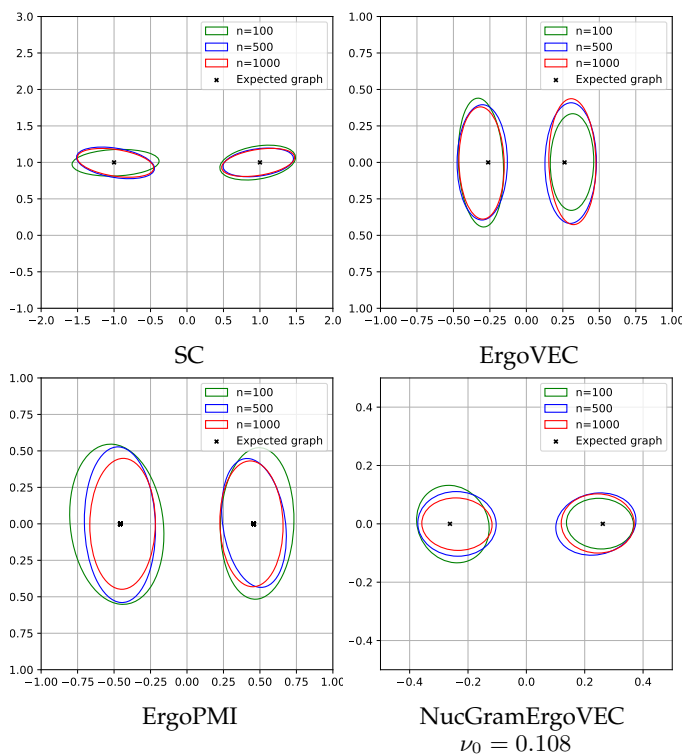


Fig. 7: 95% Gaussian contours of 2D-embeddings from four algorithms in the logarithmic regime. All algorithms receive the same sets of graphs with  $n = 100, 500$  and  $1000$  nodes generated using within-community edge probabilities  $p = 9 \ln n/n$  and across community edge probabilities  $q = 2 \ln n/n$ .

We turn to quantitative metrics to gain a more nuanced understanding. In Fig. 8, we plot the 1D-SNR of embeddings for increasing values of  $n$  for each algorithm. Note that a higher 1D-SNR indicates either a smaller within group variance along the line that joins the cluster centroids or a greater distance between cluster centroids. Fig. 8(a) shows

results for the linear degree scaling regime, where we see that 1D-SNR increases as  $n$  increases, with NucGramErgoVEC leading, followed by SC, ErgoVEC and ErgoPMI. The VEC embeddings obtained from both implementations (Keras and Gensim) reside at the bottom. In the logarithmic degree scaling regime, as shown in Fig. 8(b), the 1D-SNR is relatively steady across different  $n$  as opposed to a clear a growth trend observed in the linear degree scaling regime. This is consistent with our observations for Fig. 6 and Fig. 7 that the embeddings concentrate in linear regime but converges to a fixed distribution in the logarithmic regime. While VEC embeddings still under perform, ErgoPMI and ErgoVEC surpasses SC and catches NucGramErgoVEC’s lead.

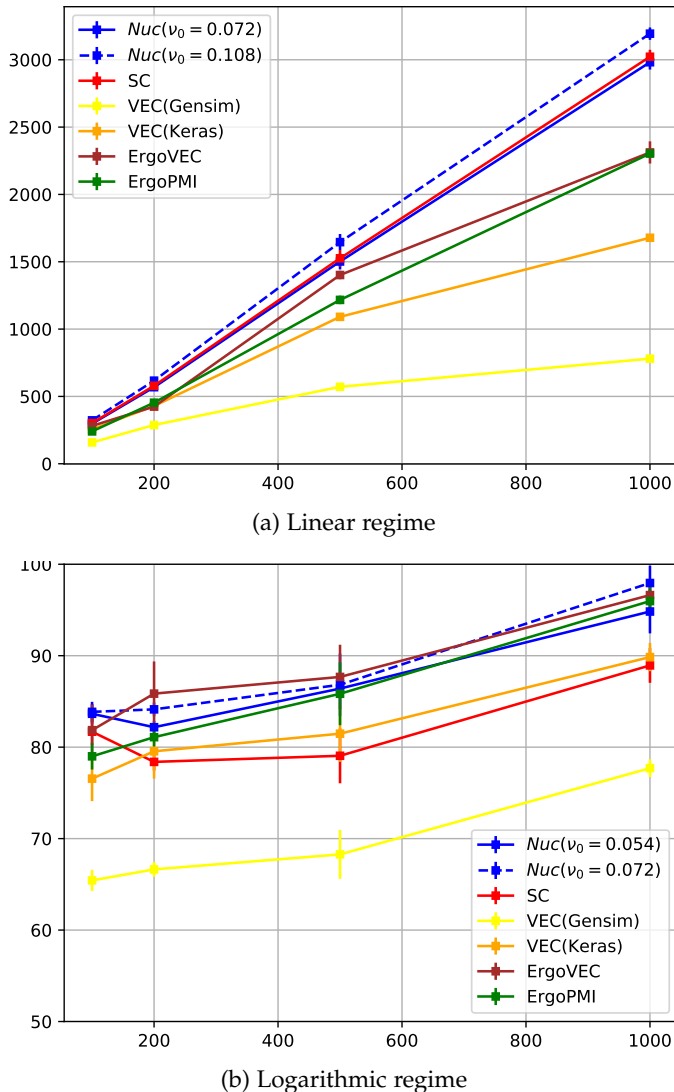


Fig. 8: 1D-SNR versus graph size in the linear and logarithmic scaling regimes.

## 6 CONCLUDING REMARKS

In this paper, we proposed a novel framework consisting of ergodic limits of random walks and a Grammian re-parameterization of the embedding objective to analyze a

large class of random walk based node-embedding algorithms. In particular, we derived a closed-form expression for the ergodic limit of the random walk node embedding objective and proved that under the positive semi-definite constraint, the Gram matrix of optimum embedding vectors for two-community expected SBM graphs has either rank 1 or rank 2. In addition, through an empirical study we demonstrated that the embeddings based on ergodic limits, while forming better clusters, in terms of 1D-SNR, compared to the original random walk embeddings, concentrate to the embeddings of the expected graph in the linear degree scaling regime and seem to converge to a fixed distribution in the logarithmic regime.

Computational costs can vary substantially across different algorithms. For example, the Gram matrix of the optimal embedding vectors in ErgoPMI has a simple closed-form solution, whereas the better performing NucGramErgoVEC requires a computationally expensive iterative optimization procedure to compute the optimal Grammian. This suggests a possible trade-off between computational cost and accuracy of algorithms. Although not the focus of this paper, understanding these trade-offs would benefit the end users of these methods.

The results of this paper can be further improved and extended on both theoretical and practical fronts. For simplicity we focused on SBM graphs with two balanced communities. Our theoretical and experimental results can be potentially extended to more complex graph models that have a community structure. On the theoretical side, although we have shown perfect separation of the embeddings of the expected graph, there is no theoretical guarantee that the embeddings of SBM random graphs will concentrate to those of the expected graph. Further analysis of random walk embedding algorithms, especially the concentration properties of their solutions in various degree scaling regimes, would bring us more insight and understanding. On the practical side, the convergence of our Keras implementations for VEC and ErgoVEC depend highly on tuning parameters and may not converge very well, and the Hazan’s algorithm for NucGramErgoVEC suffers from slow convergence. Developing more scalable implementations of algorithms with faster and more stable convergence can bring these generalized formulations into large-scale real-world problems and also guide the theoretical analysis endeavor.

## APPENDIX A

### A.1 Proof of Theorem 1

Natural random walks will remain within the connected components in which they start. Since only pairs of nodes within the same connected component will occur in any random walk, we can analyze each connected component separately. Within each connected component  $\mathcal{G}_t$ , the random walk has transition matrix  $W_t = D_t^{-1}A_t$ . The proof of the theorem will follow immediately from the following lemma which focuses on connected graphs.

**Lemma 1.** *Let  $W$  be the probability transition matrix of an irreducible Markov chain on the (finite) node space of  $\mathcal{G}$ . Let the VEC algorithm be executed on  $\mathcal{G}$  with random walk transition*

matrix  $W$  and parameters  $w$  and  $k$ . Then for all  $i, j$ , the ergodic limits  $\bar{n}_{ij}^+$  and  $\bar{n}_{ij}^-$  in Definition 2 exist and are given by

$$\begin{aligned}\bar{n}_{ij}^+ &= \pi_i \sum_{v=1}^w (W^v)_{ij}, \\ \bar{n}_{ij}^- &= kw\pi_i\pi_j,\end{aligned}$$

where  $\pi$  is the unique stationary distribution of the random walk. Moreover, the Ergodic limiting coefficients are symmetric, i.e.,

$$\begin{aligned}\bar{n}_{ij}^+ &= \bar{n}_{ji}^+, \\ \bar{n}_{ij}^- &= \bar{n}_{ji}^-.\end{aligned}$$

The proof of Lemma 1 is based on convergence results for irreducible Markov chains.

First, we prove the result for positive pairs. Let  $\{X_s^{(m,p)}\}_{s=1}^\infty$  be the  $p$ -th random walk starting from node  $m$  following the transition matrix  $W$ . We examine the first  $\ell$  steps in each random walk. Since the  $n_{ij}^+$ 's consist of positive pairs extracted from all  $rn$  random walks ( $r$  walks from each of the  $n$  nodes), we have

$$\frac{n_{ij}^+}{nr\ell} = \frac{1}{nr\ell} \sum_{m=1}^n \sum_{p=1}^r \sum_{v=1}^w \sum_{s=1}^{\ell-v} \mathbf{1}_{\{X_s^{(m,p)}=i, X_{s+v}^{(m,p)}=j\}}. \quad (21)$$

Letting  $\ell$  go to infinity on both sides, we have

$$\begin{aligned}\frac{1}{nr} \lim_{\ell \rightarrow \infty} \frac{n_{ij}^+}{\ell} \\ = \frac{1}{nr} \sum_{m=1}^n \sum_{p=1}^r \sum_{v=1}^w \lim_{\ell \rightarrow \infty} \frac{1}{\ell} \sum_{s=1}^{\ell-v} \mathbf{1}_{\{X_s^{(m,p)}=i, X_{s+v}^{(m,p)}=j\}}.\end{aligned}$$

The key step is to compute  $\lim_{\ell \rightarrow \infty} \frac{1}{\ell} \sum_{s=1}^{\ell-v} \mathbf{1}_{\{X_s^{(m,p)}=i, X_{s+v}^{(m,p)}=j\}}$ . To begin, we define a new Markov Chain  $\{Y_s^{(m,p)}\}_{s=1}^\infty$ , where  $Y_s^{(m,p)} = (X_s^{(m,p)}, X_{s+1}^{(m,p)}, \dots, X_{s+v}^{(m,p)})$ . The state space of  $\{Y_s^{(m,p)}\}_{s=1}^\infty$  is the set of all length- $(v+1)$  walks under  $W$ , i.e.,  $S_v = \{(i_1, i_2, \dots, i_{v+1}) \mid i_{k+1} \text{ is accessible from } i_k, k = 1, \dots, v\} \subset [n]^{v+1}$ .

We claim that  $\{Y_s^{(m,p)}\}_{s=1}^\infty$  is a positive recurrent Markov Chain. To see this, we first note that the state space is finite as  $|S_v| \leq |[n]^{v+1}| = n^{v+1} < \infty$ . Then, we notice that  $\forall \mathbf{a}, \mathbf{b} \in S_v$ , since  $W$  is irreducible,  $b_1$  is reachable from  $a_{v+1}$  in  $X_s^{(m,p)}$ . Therefore,  $\mathbf{b}$  is also reachable from  $\mathbf{a}$ , which shows that  $Y_s^{(m,p)}$  is irreducible. An irreducible Markov chain on a finite state space must be positive recurrent.

Applying standard results from renewal theory, specifically [48, Proposition 3.3.1, p.102] to Markov chain  $\{X_s^{(m,p)}\}_{s=1}^\infty$  and [48, Theorem 3.3.4, p.107] to Markov chain  $\{Y_s^{(m,p)}\}_{s=1}^\infty$ , we get

$$\begin{aligned}\lim_{\ell \rightarrow \infty} \frac{\sum_{s=1}^{\ell} \mathbf{1}_{\{X_s^{(m,p)}=i\}}}{\ell} &= \pi_i, \quad \text{a.s.}, \\ \lim_{\ell \rightarrow \infty} \frac{\mathbb{E} \sum_{s=1}^{\ell} \mathbf{1}_{\{X_s^{(m,p)}=i\}}}{\ell} &= \pi_i, \quad (22)\end{aligned}$$

$$\begin{aligned}\lim_{\ell \rightarrow \infty} \frac{\sum_{s=1}^{\ell} \mathbf{1}_{\{Y_s^{(m,p)}=\mathbf{a}\}}}{\ell} &= \eta_{\mathbf{a}}, \quad \text{a.s.}, \quad (23) \\ \lim_{\ell \rightarrow \infty} \frac{\mathbb{E} \sum_{s=1}^{\ell} \mathbf{1}_{\{Y_s^{(m,p)}=\mathbf{a}\}}}{\ell} &= \eta_{\mathbf{a}},\end{aligned}$$

where  $\pi$  and  $\eta$  are the stationary distributions of  $X_s^{(m,p)}$  and  $Y_s^{(m,p)}$ , respectively, and do not depend on  $m, p$  because of the positive recurrence and irreducibility of the Markov chains. Note that the relationship between state counts of  $X_s^{(m,p)}$  and  $Y_s^{(m,p)}$  is given by

$$\mathbf{1}_{\{X_s^{(m,p)}=i, X_{s+v}^{(m,p)}=j\}} = \sum_{\mathbf{a}: a_1=i, a_{v+1}=j} \mathbf{1}_{\{Y_s^{(m,p)}=\mathbf{a}\}}.$$

And thus,

$$\begin{aligned}\lim_{\ell \rightarrow \infty} \frac{1}{\ell} \sum_{s=1}^{\ell} \mathbf{1}_{\{X_s^{(m,p)}=i, X_{s+v}^{(m,p)}=j\}} \\ = \lim_{\ell \rightarrow \infty} \frac{1}{\ell} \sum_{\mathbf{a}: a_1=i, a_{v+1}=j} \mathbf{1}_{\{Y_s^{(m,p)}=\mathbf{a}\}} \\ \stackrel{(23)}{=} \sum_{\mathbf{a}: a_1=i, a_{v+1}=j} \eta_{\mathbf{a}} \\ = \sum_{\mathbf{a}: a_1=i, a_{v+1}=j} \lim_{s \rightarrow \infty} \mathbb{P}\{Y_s^{(m,p)} = \mathbf{a}\} \\ = \lim_{s \rightarrow \infty} \mathbb{P}\{X_s^{(m,p)} = i, X_{s+v}^{(m,p)} = j\} \\ = \lim_{\ell \rightarrow \infty} \frac{\sum_{s=1}^{\ell} \mathbb{P}\{X_s^{(m,p)} = i, X_{s+v}^{(m,p)} = j\}}{\ell} \\ = (W)_{ij}^v \lim_{\ell \rightarrow \infty} \frac{\sum_{s=1}^{\ell} \mathbb{P}\{X_s^{(m,p)} = i\}}{\ell} \\ = (W)_{ij}^v \lim_{\ell \rightarrow \infty} \frac{\mathbb{E} \sum_{s=1}^{\ell} \mathbf{1}_{\{X_s^{(m,p)}=i\}}}{\ell} \\ \stackrel{(22)}{=} (W^v)_{ij} \pi_i.\end{aligned}$$

Therefore, we have

$$\lim_{\ell \rightarrow \infty} \frac{1}{\ell} \sum_{s=1}^{\ell} \mathbf{1}_{\{X_s^{(m,p)}=i, X_{s+v}^{(m,p)}=j\}} = (W^v)_{ij} \pi_i,$$

and

$$\frac{1}{nr} \lim_{\ell \rightarrow \infty} \frac{n_{ij}^+}{\ell} = \frac{1}{nr} \sum_{m=1}^n \sum_{p=1}^r \sum_{v=1}^w (W^v)_{ij} \pi_i = \pi_i \sum_{v=1}^w (W^v)_{ij}.$$

We now analyze the ergodic limits of negative pairs. We first count the number of  $(i, j)$  pairs in the negative multi-set.

Let  $Z_c^{(m,p,v,s)}$  be the second node in  $c$ -th negative pair generated from the positive pair  $(X_s^{(m,p)}, X_{s+v}^{(m,p)})$  (for each positive pair we generate  $k$  negative pairs). Since all negative pairs are generated in an i.i.d. manner, for all  $m, p, v, s, c, Z_c^{(m,p,v,s)}$ ,  $c = 1, \dots, k$  are i.i.d. random variables with a distribution specified by the unigram node frequencies computed from the collection of random walks  $X = \bigcup_{m=1}^n \bigcup_{p=1}^r \{X^{(m,p)}\}$ . As in Equation (21), the counts of negative pairs is given by

$$\frac{n_{ij}^-}{nr\ell} = \frac{1}{nr\ell} \sum_{m=1}^n \sum_{p=1}^r \sum_{v=1}^w \sum_{s=1}^{\ell-v} \mathbf{1}_{\{X_s^{(m,p)}=i\}} \sum_{c=1}^k \mathbf{1}_{\{Z_c^{(m,p,v,s)}=j\}}. \quad (24)$$

Letting  $\ell$  go to infinity on both sides, we get

$$\frac{1}{nr} \lim_{\ell \rightarrow \infty} \frac{n_{ij}^-}{\ell}$$

$$= \sum_{v=1}^w \sum_{c=1}^k \lim_{\ell \rightarrow \infty} \frac{1}{nr\ell} \sum_{m=1}^n \sum_{p=1}^r \sum_{s=1}^{\ell-v} \mathbf{1}_{\{X_s^{(m,p)}=i\}} \mathbf{1}_{\{Z_c^{(m,p,s,v)}=j\}}. \quad (25)$$

The remainder of the proof will focus on calculating the right hand side. For this purpose, we introduce the following proposition:

**Notation:**  $[n] := \{1, \dots, n\}$  and  $X_{[n]} := X_1, \dots, X_n$ .

**Proposition A.1.1**

Let  $\{X_\ell, \ell \in \mathbb{N}\}$  be a sequence of random variables with  $X_\ell \in [n]$  for every  $\ell$ . Let  $\varphi : [n] \rightarrow [0, 1]$ . For every  $L \in \mathbb{N}$ , let  $\hat{q}_L : [n]^L \rightarrow [0, 1]$  and  $V_{[L]}^{(L)} \in \{0, 1\}$  be a sequence of random variables such that

$$V_{[L]}^{(L)} \mid X_{[L]} \stackrel{i.i.d.}{\sim} \text{Ber}(\hat{q}_L(X_{[L]})),$$

If for some  $p, q \in [0, 1]$ ,

$$\frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) \xrightarrow[L \rightarrow \infty]{a.s.} p, \quad (26)$$

and

$$\hat{q}_L(X_{[L]}) \xrightarrow[L \rightarrow \infty]{a.s.} q \quad (27)$$

then

$$\frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) V_\ell^{(L)} \xrightarrow[L \rightarrow \infty]{a.s.} pq.$$

*Proof.*

$$\begin{aligned} \frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) V_\ell^{(L)} &= \left( \frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) V_\ell^{(L)} - \frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) \hat{q}_L \right) \\ &\quad + \hat{q}_L \frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell). \end{aligned}$$

Due to Equations (26) and (27) we immediately have

$$\hat{q}_L \left( \frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) \right) \xrightarrow[L \rightarrow \infty]{a.s.} pq.$$

We will prove that

$$\frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) V_\ell^{(L)} - \frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) \hat{q}_L \xrightarrow[L \rightarrow \infty]{a.s.} 0.$$

For any fixed  $x_{[L]} \in [n]$ , define  $g : [0, 1]^L \rightarrow [0, 1]$  as

$$g(v_{[L]}) := \frac{1}{L} \sum_{\ell=1}^L \varphi(x_\ell) v_\ell.$$

We can show that  $g(\cdot)$  satisfies the so-called coordinate-wise bounded difference property. In fact, for any  $i \in [L]$  and any  $v_{[L]}, v'_i \in [0, 1]$ , since  $\varphi(x) \in [0, 1]$ , we have

$$\begin{aligned} &|g(v_{[L]}) - g(v_{[L] \setminus \{i\}}, v'_i)| \\ &= \frac{1}{L} |\varphi(x_i)| |v_i - v'_i| \\ &\leq \frac{1}{L}. \end{aligned}$$

Since  $V_{[L]}^{(L)}$  are i.i.d. conditioned on  $X_{[L]}$  and  $g(\cdot)$  is coordinate-wise bounded, we can apply McDiarmid's inequality [49] to  $V_{[L]}^{(L)}$  and  $g(\cdot)$  under the conditional proba-

bility measure:  $\forall \varepsilon > 0$ ,

$$\mathbb{P} \left[ \left| g(V_{[L]}^{(L)}) - \mathbb{E} \left[ g(V_{[L]}^{(L)}) \mid X_{[L]} \right] \right| \geq \varepsilon \mid X_{[L]} \right] \leq 2e^{-2L\varepsilon^2}.$$

Since the right hand side is constant and independent of  $X_{[L]}$ , the above bound also holds for unconditional probability:

$$\mathbb{P} \left[ \left| g(V_{[L]}^{(L)}) - \mathbb{E} \left[ g(V_{[L]}^{(L)}) \mid X_{[L]} \right] \right| \geq \varepsilon \right] \leq 2e^{-2L\varepsilon^2}.$$

Since

$$g(V_{[L]}^{(L)}) = \frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) V_\ell^{(L)},$$

we have

$$\begin{aligned} \mathbb{E} \left[ g(V_{[L]}^{(L)}) \mid X_{[L]} \right] &= \mathbb{E} \left[ \frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) V_\ell^{(L)} \mid X_{[L]} \right] \\ &= \frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) \mathbb{E} \left[ V_\ell^{(L)} \mid X_{[L]} \right] \\ &= \frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) \hat{q}_L. \end{aligned}$$

In other words, we have shown

$$\mathbb{P} \left[ \left| \frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) V_\ell^{(L)} - \frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) \hat{q}_L \right| \geq \varepsilon \right] \leq 2e^{-2L\varepsilon^2}.$$

Therefore,

$$\begin{aligned} &\sum_{L=1}^{\infty} \mathbb{P} \left[ \left| \frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) V_\ell^{(L)} - \frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) \hat{q}_L \right| \geq \varepsilon \right] \\ &\leq \sum_{L=1}^{\infty} 2e^{-2L\varepsilon^2} \\ &= \frac{2e^{-2\varepsilon^2}}{1 - e^{-2\varepsilon^2}} \\ &< \infty \end{aligned}$$

which proves that  $\frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) V_\ell^{(L)} - \frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) \hat{q}_L$  is converges to zero completely. Since complete convergence implies almost sure convergence [50, Theorem 4 (c), p.310], it follows that

$$\frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) V_\ell^{(L)} - \frac{1}{L} \sum_{\ell=1}^L \varphi(X_\ell) \hat{q}_L \xrightarrow[L \rightarrow \infty]{a.s.} 0,$$

which completes the proof of this proposition.  $\square$

To compute the right hand side of Equation (25), for each fixed  $c, v$ , we apply Proposition A.1.1 in the following way:

**Identifying variables:** Let  $\Gamma := \{1, \dots, n\} \times \{1, \dots, r\} \times \{1, \dots, \ell\}$ . For  $\gamma = (m, p, s)$ , we define

$$L := |\Gamma| = nr\ell,$$

$$X_\gamma := X_s^{(m,p)}$$

$$V_\gamma^{(L)} := \mathbf{1}_{\{Z_c^{(m,p,s,v)}=j\}}$$

$$\varphi(X_\gamma) := \mathbf{1}_{\{X_\gamma=i\}}$$

$$\widehat{q}_L(X_{[L]}) := \frac{1}{L} \sum_{\gamma \in \Gamma} \mathbf{1}_{\{X_\gamma=j\}}. \quad = d_j(W^{s+1})_{ji}.$$

**Verification of assumptions:**

$$\begin{aligned} V_{[L]}^{(L)} \mid X_{[L]} &\stackrel{\text{i.i.d.}}{\sim} \text{Ber}(\widehat{q}_L(X_L)), \\ \frac{1}{L} \sum_{\gamma \in \Gamma} \varphi(X_\gamma) &\xrightarrow[L \rightarrow \infty]{\text{a.s.}} \pi_i, \\ \widehat{q}_L(X_{[L]}) &= \xrightarrow[L \rightarrow \infty]{\text{a.s.}} \pi_j \end{aligned}$$

Therefore, by Equation (25), we have

$$\frac{1}{L} \sum_{\gamma \in \Gamma} \varphi(X_\gamma) V_{[L]}^{(L)} \xrightarrow[\ell \rightarrow \infty]{\text{a.s.}} \pi_i \pi_j.$$

Or equivalently, for each  $c, v$ ,

$$\frac{1}{nr\ell} \sum_{m=1}^n \sum_{p=1}^r \sum_{s=1}^{\ell} \mathbf{1}_{\{X_s^{(m,p)}=i\}} \mathbf{1}_{\{Z_c^{(m,p,s,v)}=j\}} \xrightarrow[\ell \rightarrow \infty]{\text{a.s.}} \pi_i \pi_j.$$

Dropping a finite number of terms in the summation will not affect the limit as  $\ell \rightarrow \infty$ . Thus,

$$\frac{1}{nr\ell} \sum_{m=1}^n \sum_{p=1}^r \sum_{s=1}^{\ell-v} \mathbf{1}_{\{X_s^{(m,p)}=i\}} \mathbf{1}_{\{Z_c^{(m,p,s,v)}=j\}} \xrightarrow[\ell \rightarrow \infty]{\text{a.s.}} \pi_i \pi_j.$$

Together with Equation (25), we have

$$\begin{aligned} &\frac{1}{nr} \lim_{\ell \rightarrow \infty} \frac{n_{ij}^-}{\ell} \\ &= \sum_{v=1}^w \sum_{c=1}^k \lim_{\ell \rightarrow \infty} \frac{1}{nr\ell} \sum_{m=1}^n \sum_{p=1}^r \sum_{s=1}^{\ell-v} \mathbf{1}_{\{X_s^{(m,p)}=i\}} \mathbf{1}_{\{Z_c^{(m,p,s,v)}=j\}} \\ &= \sum_{v=1}^w \sum_{c=1}^k \pi_i \pi_j \quad \text{a.s.} \\ &= kw\pi_i \pi_j. \quad \text{a.s.} \end{aligned}$$

This concludes the proof of expressions for  $\bar{n}_{ij}^+$  and  $\bar{n}_{ij}^-$  in Lemma 1 and also shows that  $\bar{n}_{ij}^- = \bar{n}_{ji}^-$ . In order to show that  $\bar{n}_{ij}^+$  is symmetric, it suffices to show that for any  $v$ ,

$$\pi_i(W^v)_{ij} = \pi_j(W^v)_{ji}.$$

Since  $\pi_i$  is proportional to the node degree, this is equivalent to showing that

$$d_i(W^v)_{ij} = d_j(W^v)_{ji}$$

We will prove this via induction. For  $v = 1$ , by definition,  $d_i W_{ij} = A_{ij} = d_j W_{ji}$  (initial case). If  $d_i(W^s)_{ij} = d_j(W^s)_{ji}$ , for  $v = s + 1$  (induction hypothesis), then

$$\begin{aligned} d_i(W^{s+1})_{ij} &= \sum_{k=1}^n d_i(W^s)_{ik} W_{kj} \\ &= \sum_{k=1}^n d_k(W^s)_{ki} W_{kj} \\ &= \sum_{k=1}^n d_k W_{kj} (W^s)_{ki} \\ &= \sum_{k=1}^n d_j W_{jk} (W^s)_{ki} \end{aligned}$$

which proves the inductive step and concludes the proof of symmetry of  $\bar{n}_{ij}^+$ .  $\square$

## A.2 Proof of Theorem 2

We will follow the same ideas as in the proof of Lemma 1. With walk-distance weights  $\{\alpha_v\}_{v=1}^\infty$ , the positive pair count Equation (21) becomes:

$$\frac{n_{ij}^+}{nr\ell} = \frac{1}{nr\ell} \sum_{m=1}^n \sum_{p=1}^r \sum_{v=1}^\infty \alpha_v \sum_{s=1}^{\ell-v} \mathbf{1}_{\{X_s^{(m,p)}=i, X_{s+v}^{(m,p)}=j\}}. \quad (28)$$

And the negative pair count Equation (24) becomes:

$$\frac{n_{ij}^-}{nr\ell} = \frac{1}{nr\ell} \sum_{m=1}^n \sum_{p=1}^r \sum_{v=1}^\infty \alpha_v \sum_{s=1}^{\ell-v} \mathbf{1}_{\{X_s^{(m,p)}=i\}} \sum_{c=1}^k \mathbf{1}_{\{Z_c^{(m,p,s,v)}=j\}}. \quad (29)$$

This provides the starting point for our proof.

1) The proof closely parallels the proof of Lemma 1 with minor modifications to account for the walk-distance weighting. We note that the exchange of the limit and the infinite sum is ensured by the dominated convergence theorem.

2) From (28), we have

$$\begin{aligned} &\frac{1}{\ell n} \lim_{r \rightarrow \infty} \frac{n_{ij}^+}{\ell} \\ &= \frac{1}{rn} \sum_{m=1}^n \sum_{v=1}^\infty \alpha_v \sum_{s=1}^{\ell-v} \lim_{r \rightarrow \infty} \frac{1}{r} \sum_{p=1}^r \mathbf{1}_{\{X_s^{(m,p)}=i, X_{s+v}^{(m,p)}=j\}} \\ &= \frac{1}{rn} \sum_{m=1}^n \sum_{v=1}^\infty \alpha_v \\ &\quad \left( \sum_{s=1}^{\ell-v} \lim_{r \rightarrow \infty} \frac{1}{r} \sum_{p=1}^r \mathbf{1}_{\{X_s^{(m,p)}=i\}} \mathbf{1}_{\{X_{s+v}^{(m,p)}=j \mid X_s^{(m,p)}=i\}} \right). \end{aligned}$$

Note that

$$\begin{aligned} &\lim_{r \rightarrow \infty} \frac{1}{r} \sum_{p=1}^r \mathbf{1}_{\{X_s^{(m,p)}=i\}} \mathbf{1}_{\{X_{s+v}^{(m,p)}=j \mid X_s^{(m,p)}=i\}} \\ &= \lim_{r \rightarrow \infty} \frac{1}{r} \sum_{p=1}^r \mathbf{1}_{\{X_s^{(m,p)}=i \mid X_1^{(m,p)}=m\}} \mathbf{1}_{\{X_{s+v}^{(m,p)}=j \mid X_s^{(m,p)}=i\}} \\ &= \mathbb{E} \mathbf{1}_{\{X_s^{(m,p)}=i \mid X_1^{(m,p)}=m\}} \mathbf{1}_{\{X_{s+v}^{(m,p)}=j \mid X_s^{(m,p)}=i\}} \\ &= \mathbb{E} \mathbf{1}_{\{X_s^{(m,p)}=i \mid X_1^{(m,p)}=m\}} \mathbb{E} \mathbf{1}_{\{X_{s+v}^{(m,p)}=j \mid X_s^{(m,p)}=i\}} \\ &= (W^{s-1})_{mi} (W^v)_{ij} \end{aligned}$$

Therefore, we have

$$\frac{1}{\ell n} \lim_{r \rightarrow \infty} \frac{n_{ij}^+}{r} = \frac{1}{\ell n} \sum_{k=1}^n \sum_{v=1}^w \alpha_v (W^v)_{ij} \sum_{s=1}^{\ell-v} (W)_{mi}^{s-1},$$

where the exchange of the limit and infinite sum is allowed by the dominated convergence theorem.

For the negative terms, from (29), we have

$$\begin{aligned} &\frac{1}{\ell n} \lim_{r \rightarrow \infty} \frac{1}{r} n_{ij}^- = \frac{1}{\ell n} \sum_{m=1}^n \sum_{v=1}^\infty \alpha_v \\ &\quad \left( \sum_{s=1}^{\ell-v} \sum_{c=1}^k \lim_{r \rightarrow \infty} \frac{1}{r} \sum_{p=1}^r \mathbf{1}_{\{X_s^{(m,p)}=i\}} \mathbf{1}_{\{Z_c^{(m,p,s,v)}=j\}} \right) \end{aligned}$$

Proceeding as we did in the proof of Lemma 1,, we apply Proposition A.1.1 to obtain

$$\begin{aligned} & \lim_{r \rightarrow \infty} \frac{1}{r} \sum_{p=1}^r \mathbf{1}_{\{X_s^{(m,p)}=i\}} \mathbf{1}_{\{Z_c^{(m,p,s,v)}=j\}} \\ &= (W)_{mi}^{s-1} \left( \frac{1}{\ell} \sum_{u=1}^{\ell} \frac{1}{n} \mathbf{1}_n^\top W^{u-1} e_j \right) \\ &:= (W)_{mi}^{s-1} \pi_j^{(\ell)} \end{aligned}$$

Therefore,

$$\frac{1}{\ell n} \lim_{r \rightarrow \infty} \frac{n_{ij}^-}{r} = \frac{k \pi_j^{(\ell)}}{\ell n} \sum_{m=1}^n \sum_{v=1}^{\infty} \alpha_v \sum_{s=1}^{\ell-v} (W)_{mi}^{s-1}$$

3) From 1), we know that  $\frac{1}{rn} \lim_{\ell \rightarrow \infty} \frac{n_{ij}^+}{\ell}$  and  $\frac{1}{rn} \lim_{\ell \rightarrow \infty} \frac{n_{ij}^-}{\ell}$  does not depend on  $r$ , and therefore the first part of the equality holds.

An irreducible Markov chain on a finite state space with a time-homogeneous transition matrix  $W$  has a unique stationary distribution  $\pi$ . Moreover, for any initial distribution on states  $\pi_0$ , the Cesaro-average:  $\bar{\pi}_\ell := \pi_0^\top \frac{1}{\ell} \sum_{s=1}^{\ell} W^s$ ,  $\ell = 1, 2, \dots$ , converges to the unique stationary distribution  $\pi$  (even if the Markov chain is not aperiodic). While this is a somewhat well-known result, we were unable to find a reliable reference that explicitly states or proves it. So for completeness we briefly sketch its proof. We argue that  $\bar{\pi}_\ell$  must converge to  $\pi$ . If not, there is an  $\epsilon > 0$  and a subsequence that lies strictly outside an  $\epsilon$ -ball around  $\pi$ . But, the probability simplex in finite-dimensional Euclidean space is compact and has the Bolzano-Weierstrass property: there is a subsequence of the subsequence (a sub-subsequence) which converges. Below we will show that the limit of this sub-subsequence must be  $\pi$  which would result in a contradiction (since the subsequence is outside an  $\epsilon$  ball around  $\pi$ ). Therefore,  $\bar{\pi}_\ell$  must converge to the unique stationary distribution  $\pi$ . We will now show that any convergent subsequence of  $\bar{\pi}_\ell$  (a convergent sub-subsequence is also convergent subsequence) must converge to  $\pi$ . Let  $\bar{\pi}_{\ell_t}$  denote a convergent subsequence and  $\bar{\pi}_\infty$  its limit. Then,

$$\begin{aligned} \bar{\pi}_\infty W &= \lim_{t \rightarrow \infty} (\bar{\pi}_{\ell_t} W) = \lim_{t \rightarrow \infty} \left( \pi_0^\top \frac{1}{\ell_t} \sum_{s=1}^{\ell_t} W^s \cdot W \right) \\ &= \lim_{t \rightarrow \infty} \left( \frac{\ell_t + 1}{\ell_t} \bar{\pi}_{\ell_t+1} - \frac{1}{\ell_t} \pi_0^\top W \right) \\ &= \bar{\pi}_\infty \end{aligned}$$

where in the first equality we made use of the fact that linear maps between finite-dimensional Euclidean spaces are continuous. Thus,  $\bar{\pi}_\infty$  is a stationary distribution of  $W$  since the above analysis shows that  $\bar{\pi}_\infty W = \bar{\pi}_\infty$ . Since  $W$  has a unique stationary distribution  $\pi$ , we have  $\bar{\pi}_\infty = \pi$ . We reiterate that aperiodicity is not needed. This is important since it is not guaranteed that the connected component subgraphs of a given graph will be aperiodic.

The following results follow immediately:

$$\lim_{\ell \rightarrow \infty} \frac{1}{\ell} \sum_{s=1}^{\ell-v} (W)_{ki}^s = \pi_i.$$

and

$$\lim_{\ell \rightarrow \infty} \pi_j^{(\ell)} = \lim_{\ell \rightarrow \infty} \frac{1}{\ell} \sum_{u=1}^{\ell} \frac{1}{n} \mathbf{1}_n^\top W^{u-1} e_j = \pi_j.$$

And therefore,

$$\begin{aligned} \frac{1}{n} \lim_{\ell \rightarrow \infty} \lim_{r \rightarrow \infty} \frac{n_{ij}^+}{r \ell} &= \pi_i \sum_{v=1}^{\infty} \alpha_v (W^v)_{ij}, \\ \frac{1}{n} \lim_{\ell \rightarrow \infty} \lim_{r \rightarrow \infty} \frac{n_{ij}^-}{r \ell} &= k \pi_i \pi_j \sum_{v=1}^{\infty} \alpha_v. \end{aligned}$$

□

**Remark:** The characterization of ergodic limits of walk-distance weighted counts stated in Theorem 2 are for connected graphs. For disconnected graphs the characterization is similar, but confined to each connected component, as in Theorem 1, and can be proved similarly.

### A.3 Proof of Proposition 1

*Proof.* Equation (16) is separable with respect to the  $X_{ij}$  variables, and for each  $X_{ij}$ , the problem reduces to the following univariate optimization problem:

$$\operatorname{argmin}_{x \in \mathbb{R}} f_{ij}(x) := \bar{n}_{ij}^+ \ln(1 + e^{-x}) + \bar{n}_{ij}^- \ln(1 + e^{+x}). \quad (30)$$

Since

$$\frac{d^2 f_{ij}}{dx^2} = \bar{n}_{ij}^+ \frac{e^{-x}}{(1 + e^{-x})^2} + \bar{n}_{ij}^- \frac{e^x}{(1 + e^x)^2} > 0,$$

it follows that  $f_{ij}$  is a twice-differentiable convex function and therefore attains a global minimum at values of  $x$  where the derivative vanishes, i.e.,

$$\frac{df_{ij}}{dx} = -\bar{n}_{ij}^+ \frac{e^{-x}}{1 + e^{-x}} + \bar{n}_{ij}^- \frac{e^x}{1 + e^x} = 0,$$

or equivalently

$$\bar{n}_{ij}^- e^{2x} + (\bar{n}_{ij}^+ - \bar{n}_{ij}^-) e^x - \bar{n}_{ij}^+ = 0.$$

Note that from Equation (5) of Theorem 1 we know  $\bar{n}_{ij}^- > 0$ . Therefore, when  $\bar{n}_{ij}^+ \neq 0$ , we have a unique solution  $e^x = \frac{\bar{n}_{ij}^+}{\bar{n}_{ij}^-}$ , i.e.,  $x = \ln\left(\frac{\bar{n}_{ij}^+}{\bar{n}_{ij}^-}\right)$ . When  $\bar{n}_{ij}^+ = 0$ ,  $f_{ij}(x)$  is monotonically increasing over the entire real line, and we take  $x = -\infty$  as the solution. Thus,

$$X_{ij}^* = \begin{cases} \ln\left(\frac{\bar{n}_{ij}^+}{\bar{n}_{ij}^-}\right) & \text{if } \bar{n}_{ij}^+ \neq 0; \\ -\infty & \text{if } \bar{n}_{ij}^+ = 0. \end{cases}$$

From Lemma 1, we have  $\bar{n}_{ij}^+ = \bar{n}_{ji}^+$  and  $\bar{n}_{ij}^- = \bar{n}_{ji}^-$ . Therefore we have  $X_{ij}^* = X_{ji}^*$ .

For  $rn$  random walks each of length  $\ell$ , the total number of node pairs within  $w$  steps of each other is  $|\mathcal{D}_{\ell,+}| = rn \left( \ell w - \frac{w(w+1)}{2} \right)$ . First note that, when  $\bar{n}_{ij}^+ = 0$ ,  $(i, j)$  are in two different connected components, and thus  $n_{ij}^+ = 0$ , or equivalently,  $p_\ell(i, j) = 0$ . Therefore,  $\text{PMI}_\ell(i, j) = -\infty = X_{ij}^* - \ln k$  holds. For the rest of the proof, we only consider the case when  $\bar{n}_{ij}^+ \neq 0$  or, equivalently, when  $(i, j)$  are in the same connected component. For the joint distribution, we have

$$p_\ell(i, j) = \frac{n_{ij}^+}{|\mathcal{D}_{\ell,+}|} \xrightarrow[\text{Eq. (2)}]{\ell \rightarrow \infty \text{ a.s.}} \frac{\bar{n}_{ij}^+}{w}.$$

For the marginal distributions,

$$p_{\ell 1}(i) = \sum_{j \in \mathcal{V}} p_{\ell}(i, j) = \frac{\sum_{j \in \mathcal{V}} n_{ij}^+}{|\mathcal{D}_{\ell,+}|} \xrightarrow[\text{Eq. (2)}]{\ell \rightarrow \infty \text{ a.s.}} \frac{\sum_{j \in \mathcal{V}} \bar{n}_{ij}^+}{w} \stackrel{\text{Eq. (4)}}{=} \pi_i \mathcal{E}_0 := \{(i, j) : i \neq j, i, j \leq m \text{ or } i, j \geq m + 1\}$$

and

$$p_{\ell 2}(j) = \sum_{i \in \mathcal{V}} p_{\ell}(i, j) = \frac{\sum_{i \in \mathcal{V}} n_{ij}^+}{|\mathcal{D}_{\ell,+}|} \xrightarrow[\text{Eq. (2)}]{\ell \rightarrow \infty \text{ a.s.}} \frac{\sum_{i \in \mathcal{V}} \bar{n}_{ij}^+}{w}.$$

Note that since  $\pi$  is the stationary distribution of the random walk, for any  $v$ ,  $\sum_{i \in \mathcal{V}} \pi_i (W^v)_{ij} = \pi_j$ . Combining this with Eq. (4), we get

$$\sum_{i \in \mathcal{V}} \bar{n}_{ij}^+ = \sum_{v=1}^w \sum_{i \in \mathcal{V}} \pi_i (W^v)_{ij} = \sum_{v=1}^w \pi_j = w \pi_j.$$

Therefore,  $p_{\ell 2}(j) \xrightarrow[\text{a.s.}]{\ell \rightarrow \infty} \pi_j$  and

$$\begin{aligned} \text{PMI}_{\ell}(i, j) &= \ln \left( \frac{p_{\ell}(i, j)}{p_{\ell 1}(i) p_{\ell 2}(j)} \right) \\ &\xrightarrow[\text{a.s.}]{\ell \rightarrow \infty} \ln \left( \frac{\bar{n}_{ij}^+}{w \pi_i \pi_j} \right) \\ &\stackrel{\text{Eq. (5)}}{=} \ln \left( \frac{\bar{n}_{ij}^+}{\bar{n}_{ij}} \right) + \ln k \\ &= X_{ij}^* + \ln k \end{aligned}$$

□

#### A.4 Proof of Theorem 3

The main structure of the proof is as follows:

- 1) Part 1 will be shown using Lemma 1 and the eigenvalue decomposition of the *diagonal-blockwise-constant* (DBC) matrices (defined below).
- 2) Part 2 is a direct consequence of combining Part 1) and Proposition 1.
- 3) Part 3 is intricate and will be proved in 3-steps:
  - 1) First, we will show that the solution must be a DBC matrix, and thus can be re-parameterized by the three scalars that define a DBC matrix. This will be established by showing that for any feasible solution, a DBC matrix can be constructed that is both feasible and yields a lower objective cost.
  - 2) Second, we will prove that among the 3 scalar variables in the re-parameterized problem, the optimal value of two of them must equal. This implies that the matrix solution must have a block structure. We will then eliminate one variable and re-parameterize the optimization problem in terms of the remaining two variables.
  - 3) Lastly, we will show that the optimal values of the two variables will be opposite numbers of each other which will imply that the solution matrix has rank 1.
- 4) Part 4 is a direct consequence of parts (1)–(3).

Before getting into the derivations, we set up some notation and define *diagonal-blockwise-constant* (DBC) matrices.

Without loss of generality, we assume that nodes in the two balanced communities are  $\{1, \dots, m\}$  and  $\{m +$

$1, \dots, 2m\}$ . Under this labeling, we define the following two subsets of node pairs (edges)

$$\mathcal{E}_0 := \{(i, j) : i \neq j, i, j \leq m \text{ or } i, j \geq m + 1\}$$

$$\mathcal{E}_1 := \{(i, j) : i \leq m, j \geq m + 1\}$$

$$\bigcup \{(i, j) : i \geq m + 1, j \leq m\}.$$

Then,  $|\mathcal{E}_0| = 2m(m - 1)$  and  $|\mathcal{E}_1| = 2m^2$ .

Next, we define *diagonal-blockwise-constant* (DBC) matrices.

**Definition A.4.1** (DBC matrix)

Let  $\mathbf{1}_m := (1, 1, \dots, 1)^\top \in \mathbb{R}^m$ . Let  $\mathbf{y}_1 = (\mathbf{1}_m^\top, \mathbf{1}_m^\top)^\top$  and  $\mathbf{y}_2 = (\mathbf{1}_m^\top, -\mathbf{1}_m^\top)^\top$ . For  $m \geq 2$ , a  $2m \times 2m$  matrix is called *diagonal-blockwise-constant* (DBC) if it has the form

$$Z_{2m}(c_1, c_2, c_3) := \frac{c_1 + c_2}{2} \mathbf{y}_1 \mathbf{y}_1^\top + \frac{c_1 - c_2}{2} \mathbf{y}_2 \mathbf{y}_2^\top + (c_3 - c_1) I_{2m}. \quad (31)$$

Certain key properties of DBC matrices that we use in our proof are described in the following proposition.

**Proposition A.4.1** (Properties of DBC matrices)

Let  $\mathcal{E}_0$  and  $\mathcal{E}_1$  be as stated above and let  $X$  be a  $2m \times 2m$  matrix for  $m \geq 2$ . Then,  $X = Z_{2m}(c_1, c_2, c_3)$  if, and only if, any one of the following holds:

- 1)  $X$  has the following block structure:

$$X_{ij} = \begin{cases} c_1 & \text{if } (i, j) \in \mathcal{E}_0; \\ c_2 & \text{if } (i, j) \in \mathcal{E}_1; \\ c_3 & \text{if } i = j \end{cases}$$

- 2) The eigenvalues and eigenvectors of  $X$  satisfy:

$$\begin{aligned} a) & \lambda_3 = \lambda_4 = \dots = \lambda_{2m}; \\ b) & \mathbf{u}_1 = \frac{1}{\sqrt{2m}} \mathbf{y}_1, \mathbf{u}_2 = \frac{1}{\sqrt{2m}} \mathbf{y}_2. \end{aligned}$$

In addition, the set of all DBC matrices is closed under matrix addition and multiplication operations.

*Proof.*

**Proof of equivalence.**

1) Both if and only if parts can be obtained directly from Equation (31) in Definition A.4.1:

$$X_{ij} = \begin{cases} \frac{c_1 + c_2}{2} + \frac{c_1 - c_2}{2} = c_1 & \text{if } (i, j) \in \mathcal{E}_0; \\ \frac{c_1 + c_2}{2} - \frac{c_1 - c_2}{2} = c_2 & \text{if } (i, j) \in \mathcal{E}_1; \\ \frac{c_1 + c_2}{2} + \frac{c_1 - c_2}{2} + (c_3 - c_1) = c_3 & \text{if } i = j \end{cases}$$

2) If  $X = Z_{2m}(c_1, c_2, c_3)$ , directly from equation (31), we can compute the spectral decomposition of  $X$ . Let  $\mathbf{u}_1 = \frac{1}{\sqrt{2m}} \mathbf{y}_1$ ,  $\mathbf{u}_2 = \frac{1}{\sqrt{2m}} \mathbf{y}_2$  and  $\mathbf{u}_3, \dots, \mathbf{u}_{2m}$  be any set of orthonormal vectors that together with  $\mathbf{u}_1$  and  $\mathbf{u}_2$  form an orthonormal basis for  $\mathbb{R}^{2m}$ . Then,

$$\begin{aligned} X &= m(c_1 + c_2) \mathbf{u}_1 \mathbf{u}_1^\top + m(c_1 - c_2) \mathbf{u}_2 \mathbf{u}_2^\top + (c_3 - c_1) \sum_{i=1}^{2m} \mathbf{u}_i \mathbf{u}_i^\top \\ &= (m(c_1 + c_2) + (c_3 - c_1)) \mathbf{u}_1 \mathbf{u}_1^\top + \\ &\quad (m(c_1 - c_2) + (c_3 - c_1)) \mathbf{u}_2 \mathbf{u}_2^\top + \sum_{i=3}^{2m} (c_3 - c_1) \mathbf{u}_i \mathbf{u}_i^\top. \end{aligned}$$

Therefore,  $\mathbf{u}_1, \dots, \mathbf{u}_{2m}$  are the eigenvectors of  $X$  and the eigenvalues satisfy  $\lambda_3 = \lambda_4 = \dots = \lambda_{2m} = c_3 - c_1$ .

Reversely, if the eigenvalues and eigenvectors of  $X$  have the given property, letting  $U = [\mathbf{u}_1, \dots, \mathbf{u}_{2m}]$ , we have

$$\begin{aligned} X &= U \text{Diag}\{\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_3\}U^\top \\ &= U \text{Diag}\{\lambda_1 - \lambda_3, \lambda_2 - \lambda_3, 0, \dots, 0\}U^\top + U\lambda_3 I_{2m}U^\top \\ &= \frac{\lambda_1 - \lambda_3}{2m} \mathbf{y}_1 \mathbf{y}_1^\top + \frac{\lambda_2 - \lambda_3}{2m} \mathbf{y}_2 \mathbf{y}_2^\top + \lambda_3 I_{2m} \\ &= Z_{2m} \left( \frac{\lambda_1 + \lambda_2 - 2\lambda_3}{2m}, \frac{\lambda_1 - \lambda_2}{2m}, \frac{\lambda_1 + \lambda_2 + (2m-2)\lambda_3}{2m} \right) \end{aligned}$$

### Proof of set closure

Let  $X_1, X_2$  be two DBC matrices. From part 2), defining  $U = [\mathbf{u}_1, \dots, \mathbf{u}_{2m}]$  where  $\mathbf{u}_1 = \frac{1}{\sqrt{2m}} \mathbf{y}_1$ ,  $\mathbf{u}_2 = \frac{1}{\sqrt{2m}} \mathbf{y}_2$  and  $\{\mathbf{u}_3, \dots, \mathbf{u}_{2m}\}$  is any set of orthonormal vectors that together with  $\mathbf{u}_1$  and  $\mathbf{u}_2$  form an orthonormal basis for  $\mathbb{R}^{2m}$ , we have

$$\begin{aligned} X_1 &= U \text{Diag}\{\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_3\}U^\top \\ X_2 &= U \text{Diag}\{\mu_1, \mu_2, \mu_3, \dots, \mu_3\}U^\top. \end{aligned}$$

Therefore,

$$\begin{aligned} X_1 + X_2 &= U \text{Diag}\{\lambda_1 + \mu_1, \lambda_2 + \mu_2, \lambda_3 + \mu_3, \dots, \lambda_3 + \mu_3\}U^\top \\ X_1 X_2 &= U \text{Diag}\{\lambda_1 \mu_1, \lambda_2 \mu_2, \lambda_3 \mu_3, \dots, \lambda_3 \mu_3\}U^\top \end{aligned}$$

satisfy the conditions a) and b) in part 2), and they are both DBC matrices.  $\square$

**Notation.** For ease of reference, for a DBC matrix  $X$ , we denote  $\lambda_i(X)$  as its eigenvalues and  $c_i(X)$  ( $i = 1, 2, 3$ ) as its entry values in  $\mathcal{E}_0$ ,  $\mathcal{E}_1$  and diagonal, respectively. I.e.,  $X = Z_{2m}(c_1(X), c_2(X), c_3(X))$ . The derivation in the proof above gives the transformation formula between them. Specifically, given  $X = Z_{2m}(c_1, c_2, c_3)$ , we have

$$\lambda_1(X) = (m-1)c_1 + c_3 + mc_2, \quad (32)$$

$$\lambda_2(X) = (m-1)c_1 + c_3 - mc_2, \quad (33)$$

$$\lambda_3(X) = c_3 - c_1. \quad (34)$$

And given the eigenvalues  $\lambda_1, \lambda_2, \lambda_3 = \dots = \lambda_{2m}$  of  $X$ , we have

$$c_1(X) = \frac{\lambda_1 + \lambda_2 - 2\lambda_3}{2m} \quad (35)$$

$$c_2(X) = \frac{\lambda_1 - \lambda_2}{2m} \quad (36)$$

$$c_3(X) = \frac{\lambda_1 + \lambda_2 + (2m-2)\lambda_3}{2m}. \quad (37)$$

### Proposition A.4.2 (P.S.D. condition of DBC matrices)

Let  $X = Z_{2m}(c_1, c_2, c_3)$  be a DBC matrix with  $c_3 \geq c_1$ . Denote  $\bar{c}_{13} := \frac{m-1}{m}c_1 + \frac{1}{m}c_3$  and let  $Y = Z_{2m}(\bar{c}_{13}, c_2, \bar{c}_{13})$ . Then, if  $Y \succeq 0$ , we have  $X \succeq 0$ .

*Proof.*

By Equations (32) and (33), we have

$$\begin{aligned} \lambda_1(Y) &= (m-1)\bar{c}_{13} + \bar{c}_{13} + mc_2 \\ &= (m-1)c_1 + c_3 + mc_2 \\ &= \lambda_1(X) \end{aligned}$$

$$\begin{aligned} \lambda_2(Y) &= (m-1)\bar{c}_{13} + \bar{c}_{13} - mc_2 \\ &= (m-1)c_1 + c_3 + mc_2 \\ &= \lambda_2(X) \end{aligned}$$

Since  $Y \succeq 0$ , we have  $\lambda_1(X) \geq 0$  and  $\lambda_2(X) \geq 0$ . Since  $c_3 \geq c_1$ , we have  $\lambda_3(X) = c_3 - c_1 \geq 0$ . And therefore,  $X \succeq 0$ .  $\square$

Now, we are ready to prove Theorem 3.

### Part 1)

Note that for expected graph, the adjacency matrix  $A$  and random walk transition matrix  $W$  are both DBC matrices, and the stationary distribution  $\pi$  is uniform distribution. Specifically, we have

$$\begin{aligned} A &= Z_{2m}(a, b, 0) \\ W &= Z_{2m} \left( \frac{a}{(m-1)a + mb}, \frac{a}{(m-1)a + mb}, 0 \right) \\ \pi &= \frac{1}{2m} \mathbf{1}_{2m} \end{aligned}$$

By Lemma 1, we can compute the positive and negative coefficient matrices  $\bar{N}^+$  and  $\bar{N}^-$  as

$$\bar{N}^+ = \frac{1}{2m} \sum_{v=1}^w W^v, \quad (38)$$

$$\bar{N}^- = kw\pi\pi^\top = \frac{kw}{4m^2} \mathbf{1}_{2m} \mathbf{1}_{2m}^\top. \quad (39)$$

Equation (39) gives us  $\bar{n}_{ij}^- = \frac{kw}{4m^2} = \frac{kw}{n^2} =: \beta$ . It remains to show  $\bar{n}_{ij}^+$ .

Since  $\bar{N}^+$  is a sum of products of DBC matrices, by closure of DBC set (Proposition A.4.1),  $\bar{N}^+$  is a DBC matrix. In order to compute  $c_1(\bar{N}^+)$ ,  $c_2(\bar{N}^+)$ , and  $c_3(\bar{N}^+)$ , we begin from its eigenvalues. Since  $W$  is a DBC matrix, by Equations (32)–(34), we have

$$\begin{aligned} \lambda_1(W) &= 1, \\ \lambda_2(W) &= \frac{(m-1)a - mb}{(m-1)a + mb}, \\ \lambda_3(W) &= \dots = \lambda_{2m}(W) = -\frac{a}{(m-1)a + mb}. \end{aligned}$$

Note that since we assumed  $a > \frac{m}{m-1}b$ , we have  $\lambda_1(W) > 0$ ,  $\lambda_2(W) > 0$ ,  $\lambda_3(W) < 0$ .

From Equation (38), we obtain the eigenvalues of  $\bar{N}^+$

$$\begin{aligned} \lambda_1(\bar{N}^+) &= \frac{w}{2m}, \\ \lambda_2(\bar{N}^+) &= \frac{1}{2m} \sum_{v=1}^w \lambda_2(W)^v, \\ \lambda_3(\bar{N}^+) &= \dots = \lambda_{2m}(\bar{N}^+) = \frac{1}{2m} \sum_{v=1}^w \lambda_3(W)^v. \end{aligned}$$

Given the sign of  $\lambda_i(W)$ , we have  $\lambda_3(\bar{N}^+) < 0 < \lambda_2(\bar{N}^+) < \lambda_1(\bar{N}^+)$ . With Equations (35)–(37), the entry values are given as

$$\begin{aligned} c_1(\bar{N}^+) &= \frac{1}{4m^2} \left[ w + \sum_{v=1}^w \lambda_2(W)^v - 2 \sum_{v=1}^w \lambda_3(W)^v \right] \\ &:= \alpha_1, \end{aligned} \quad (40)$$

$$\begin{aligned} c_2(\bar{N}^+) &= \frac{1}{4m^2} \left[ w - \sum_{v=1}^w \lambda_2(W)^v \right] \\ &:= \alpha_2, \end{aligned} \quad (41)$$

$$c_3(\bar{N}^+) = \frac{1}{4m^2} \left[ w + \sum_{v=1}^w \lambda_2(W)^v + (2m-2) \sum_{v=1}^w \lambda_3(W)^v \right] \quad (42)$$

$$:= \alpha_3.$$

This completes the proof of Part 1). Note that  $\lambda_1(W) = 1$ ,  $\lambda_2(W) = 1 - O(1/n)$ ,  $\lambda_3(W) = O(1/n)$ . Therefore, from Equations (40)–(42), we have  $\alpha_i = C_i/n^2 + o(1/n^2)$  for  $i = 1, 2, 3$ , where  $C_i$ 's are functions of only  $a, b$  and  $w$ . Given the sign of  $\lambda_i(W)$ , we have

$$\begin{aligned} \alpha_1 &> \alpha_3 > 0, \\ \alpha_1 &> \alpha_2 > 0. \end{aligned} \quad (43)$$

### Part 2)

Applying Proposition 1, since  $\bar{n}_{ij}^+ > 0$  and  $\bar{n}_{ij}^- > 0$  hold for all  $i, j$ , we have

$$X_{ij}^* = \ln \left( \frac{\bar{n}_{ij}^+}{\bar{n}_{ij}^-} \right) = \begin{cases} \ln \left( \frac{\alpha_1}{\beta} \right), & \text{if } (i, j) \in \mathcal{E}_0 \\ \ln \left( \frac{\alpha_2}{\beta} \right), & \text{if } (i, j) \in \mathcal{E}_1 \\ \ln \left( \frac{\alpha_3}{\beta} \right), & \text{if } i = j \end{cases}$$

### Part 3)

When  $\mathcal{H} = \{X \mid X \succeq 0\}$ , we first establish structures that  $X^*$  must have, and then solve it explicitly. We take three major steps:

Step 1 We show that  $X^*$  must be a DBC matrix, and thus we can re-parameterize the optimization problem into three scalar variables:  $c_1, c_2$ , and  $c_3$ .

Step 2 We prove that among the optimal solution of this re-parameterized problem must satisfy  $c_1^* = c_3^*$ . Then, we substitute  $c_3$  by  $c_1$  and only keep  $c_1$  and  $c_2$  as optimizing variables.

Step 3 We show that  $c_1^* = -c_2^*$  must hold. After eliminating  $c_2$ , we solve the optimization explicitly.

#### Step 1.

For any matrix  $X \in S_+$ , let  $c_1, c_2$ , and  $c_3$  be the average of its entries in region  $\mathcal{E}_0, \mathcal{E}_1$  and on diagonal, respectively. I.e.,

$$x_1 := \frac{1}{2m^2 - 2m} \sum_{(i,j) \in \mathcal{E}_0} X_{ij}, \quad (44)$$

$$x_2 := \frac{1}{2m^2} \sum_{(i,j) \in \mathcal{E}_1} X_{ij}, \quad (45)$$

$$x_3 := \frac{1}{2m} \sum_{i=1}^{2m} X_{ii}. \quad (46)$$

Then, we construct a DBC matrix  $\tilde{X}$  as

$$\tilde{X} = Z_{2m}(x_1, x_2, x_3).$$

Denoting our objective function in Equation (20) as  $f$ , i.e.,

$$f(X) := \sum_{(i,j)} \left[ \bar{n}_{ij}^+ \ln(1 + e^{-X_{ij}}) + \bar{n}_{ij}^- \ln(1 + e^{X_{ij}}) \right],$$

we claim that

- $\tilde{X} \succeq 0$ . I.e.,  $\tilde{X} \in \mathcal{H}$  is feasible.
- $f(\tilde{X}) \leq f(X)$ . I.e.,  $\tilde{X}$  will be no worse than  $X$ .

Combining a) and b) will show that the optimal solution  $X^*$  must be a DBC matrix. Below, we will prove these claims.

a) We begin by computing the eigenvalues of the DBC matrix  $\tilde{X}$  and substituting the Equations (44)–(46):

$$\begin{aligned} \lambda_1(\tilde{X}) &= (m-1)x_1 + x_3 + mx_2 \\ &= \frac{1}{2m} \sum_{(i,j) \in \mathcal{E}_0} X_{ij} + \frac{1}{2m} \sum_{i=1}^{2m} X_{ii} + \frac{1}{2m} \sum_{(i,j) \in \mathcal{E}_1} X_{ij} \\ &= \frac{1}{2m} \sum_{i,j} X_{ij} \\ &= \frac{1}{2m} \mathbf{1}_{2m}^\top X \mathbf{1}_{2m} \\ &\geq 0 \end{aligned}$$

$$\begin{aligned} \lambda_2(\tilde{X}) &= (m-1)x_1 + x_3 - mx_2 \\ &= \frac{1}{2m} \sum_{(i,j) \in \mathcal{E}_0} X_{ij} + \frac{1}{2m} \sum_{i=1}^{2m} X_{ii} - \frac{1}{2m} \sum_{(i,j) \in \mathcal{E}_1} X_{ij} \\ &= \frac{1}{2m} \begin{bmatrix} \mathbf{1}_m^\top & -\mathbf{1}_m^\top \end{bmatrix} X \begin{bmatrix} \mathbf{1}_m \\ -\mathbf{1}_m \end{bmatrix} \\ &\geq 0 \end{aligned}$$

$$\begin{aligned} \lambda_3(\tilde{X}) &= x_3 - x_1 \\ &= \frac{1}{2m^2 - 2m} \left[ (m-1) \sum_{i=1}^{2m} X_{ii} - \sum_{(i,j) \in \mathcal{E}_0} X_{ij} \right]. \end{aligned}$$

To show that  $\lambda_3(\tilde{X}) \geq 0$ , we first prove the below proposition:

#### Proposition A.4.3

If an  $m \times m$  matrix  $X \succeq 0$ , then

$$\text{Tr}(X) \geq \frac{1}{m} \mathbf{1}_m^\top X \mathbf{1}_m$$

*Proof.* Let the eigen-decomposition of  $X$  be given as follows

$$X = U \Lambda U^\top,$$

where  $U = [\mathbf{u}_1, \dots, \mathbf{u}_m]$  and  $\Lambda = \text{Diag}\{\lambda_1, \dots, \lambda_m\}$ . Then, we have

$$\text{Tr}(X) = \text{Tr}(U \Lambda U^\top) = \text{Tr}(\Lambda U^\top U) = \sum_{i=1}^m \lambda_i.$$

And

$$\begin{aligned} \frac{1}{n} \mathbf{1}_m^\top X \mathbf{1}_m &= \left( \frac{1}{\sqrt{n}} \mathbf{1}_m^\top U \right) \Lambda \left( U^\top \frac{1}{\sqrt{n}} \mathbf{1}_m \right) \\ &= \sum_{i=1}^m \lambda_i \left( \frac{1}{\sqrt{n}} \mathbf{1}_m^\top \mathbf{u}_i \right)^2 \\ &\leq \sum_{i=1}^m \lambda_i \left\| \frac{1}{\sqrt{n}} \mathbf{1}_m \right\| \|\mathbf{u}_i\| \\ &= \sum_{i=1}^m \lambda_i. \end{aligned}$$

Therefore,

$$\text{Tr}(X) \geq \frac{1}{n} \mathbf{1}_m^\top X \mathbf{1}_m.$$

□

We divide  $X$  into 4  $m \times m$  block matrices as

$$X = \begin{pmatrix} X_{11} & X_{12} \\ X_{21} & X_{22} \end{pmatrix}.$$

Note that  $X_{11} \succeq 0$  and  $X_{22} \succeq 0$ . To see this, for any  $\mathbf{a} \in \mathbb{R}^m$ , we have  $\mathbf{a}^\top X_{11} \mathbf{a} = [\mathbf{a}^\top, \mathbf{0}] X \begin{bmatrix} \mathbf{a} \\ \mathbf{0} \end{bmatrix} \geq 0$  and  $\mathbf{a}^\top X_{22} \mathbf{a} = [\mathbf{0}, \mathbf{a}^\top] X \begin{bmatrix} \mathbf{0} \\ \mathbf{a} \end{bmatrix} \geq 0$ . Therefore, by Proposition A.4.3, we have

$$\text{Tr}(X_{11}) \geq \frac{1}{m} \mathbf{1}_m^\top X \mathbf{1}_m.$$

Or equivalently,

$$(m-1) \sum_{i=1}^m X_{ii} \geq \sum_{i \neq j, i, j \leq m} X_{ij}.$$

Similarly with  $X_{22}$ , we have

$$(m-1) \sum_{i=m+1}^{2m} X_{ii} \geq \sum_{i \neq j, i, j \geq m+1} X_{ij}.$$

Note that  $\mathcal{E}_0 = \{i \neq j \mid i, j \leq m \text{ or } i, j \geq m+1\}$ . Adding the above two equations yields

$$(m-1) \sum_{i=1}^{2m} X_{ii} \geq \sum_{(i,j) \in \mathcal{E}_0} X_{ij},$$

which shows that  $\lambda_3(\tilde{X}) \geq 0$ . This concludes our proof of  $\tilde{X} \succeq 0$ .

b) To show that  $\tilde{X}$  has a better cost, we will use convexity. Specifically, we define

$$\Psi(x; \alpha, \beta) := \alpha \ln(1 + e^{-x}) + \beta \ln(1 + e^x).$$

And we can rewrite  $f(X)$  as

$$\begin{aligned} f(X) &= \sum_{(i,j)} \left[ \bar{n}_{ij}^+ \ln(1 + e^{-X_{ij}}) + \bar{n}_{ij}^- \ln(1 + e^{X_{ij}}) \right] \\ &= \sum_{(i,j) \in \mathcal{E}_0} \Psi(X_{ij}; \alpha_1, \beta) + \sum_{(i,j) \in \mathcal{E}_1} \Psi(X_{ij}; \alpha_2, \beta) + \\ &\quad \sum_{i=1}^{2m} \Psi(X_{ii}; \alpha_3, \beta). \end{aligned}$$

Since for any  $\alpha, \beta > 0$ ,

$$\Psi''(x; \alpha, \beta) = \frac{e^x(\alpha + \beta)}{(1 + e^x)^2} > 0.$$

We know that  $\Psi(x; \alpha, \beta)$  is strictly convex with respect to  $x$  for any positive  $\alpha$  and  $\beta$ . With Equations (35)–(37) in mind, we have

$$\begin{aligned} \frac{1}{2m^2 - 2m} \sum_{(i,j) \in \mathcal{E}_0} \Psi(X_{ij}; \alpha_1, \beta) &\geq \Psi(x_1; \alpha_1, \beta), \\ \frac{1}{2m^2} \sum_{(i,j) \in \mathcal{E}_1} \Psi(X_{ij}; \alpha_2, \beta) &\geq \Psi(x_2; \alpha_2, \beta), \end{aligned}$$

$$\frac{1}{2m} \sum_{i=1}^{2m} \Psi(X_{ii}; \alpha_3, \beta) \geq \Psi(x_3; \alpha_3, \beta).$$

Therefore

$$\begin{aligned} f(\tilde{X}) &= \sum_{(i,j) \in \mathcal{E}_0} \Psi(x_1; \alpha_1, \beta) + \sum_{(i,j) \in \mathcal{E}_1} \Psi(x_2; \alpha_2, \beta) + \\ &\quad \sum_{i=1}^{2m} \Psi(x_3; \alpha_3, \beta) \\ &= (2m^2 - 2m) \Psi(x_1; \alpha_1, \beta) + 2m^2 \Psi(x_2; \alpha_2, \beta) + \\ &\quad 2m \Psi(x_3; \alpha_3, \beta) \\ &\leq \sum_{(i,j) \in \mathcal{E}_0} \Psi(X_{ij}; \alpha_1, \beta) + \sum_{(i,j) \in \mathcal{E}_1} \Psi(X_{ij}; \alpha_2, \beta) + \\ &\quad \sum_{i=1}^{2m} \Psi(X_{ii}; \alpha_3, \beta) \\ &= f(X). \end{aligned}$$

By far, we have shown that the DBC matrix  $\tilde{X}$  we constructed is in the feasible set and has a lower cost. Therefore, we conclude that the optimal solution matrix  $X^*$  must be a DBC matrix. Without the loss of generality, we can assume  $X = Z_{2m}(c_1, c_2, c_3)$ , and  $f(X)$  reduces to (up to a constant scaling)

$$\begin{aligned} f_3(c_1, c_2, c_3) &:= \\ &(m-1) \Psi(c_1, \alpha_1, \beta) + m \Psi(c_2, \alpha_2, \beta) + \Psi(c_3, \alpha_3, \beta). \end{aligned}$$

The optimization problem Equation (20) is equivalently transformed into

$$\begin{aligned} (c_1^*, c_2^*, c_3^*) &= \operatorname{argmin} f_3(c_1, c_2, c_3) \quad (47) \\ \text{s.t.} \quad &c_1 \leq c_3 \\ &|c_2| \leq \frac{m-1}{m} c_1 + \frac{1}{m} c_3. \end{aligned}$$

**Step 2.**

In this step, we will prove that the optimal solution to (47) must satisfy  $c_1^* = c_3^*$ . Specifically, we have the below proposition:

**Proposition A.4.4**

Let  $\bar{c}_{13} := \frac{m-1}{m} c_1 + \frac{1}{m} c_3$ . If  $(c_1, c_2, c_3)$  is a feasible solution to optimization problem (47), then  $(\bar{c}_{13}, c_2, \bar{c}_{13})$  is also feasible, and its cost is no worse than  $(c_1, c_2, c_3)$ . I.e.,

$$f_3(\bar{c}_{13}, c_2, \bar{c}_{13}) \leq f_3(c_1, c_2, c_3).$$

*Proof.*

We first show that  $(\bar{c}_{13}, c_2, \bar{c}_{13})$  is feasible. The first constraint of (47) holds as we have the same value in the first and third argument. It remains to verify the second constraint

$$\frac{m-1}{m} \bar{c}_{13} + \frac{1}{m} \bar{c}_{13} = \bar{c}_{13} = \frac{m-1}{m} c_1 + \frac{1}{m} c_3 \geq |c_2|,$$

where the last inequality is exactly the second constraint for  $(c_1, c_2, c_3)$  and holds because of its feasibility.

Next, we show that  $f_3(\bar{c}_{13}, c_2, \bar{c}_{13}) \leq f_3(c_1, c_2, c_3)$ . Expanding both sides, our goal is equivalent to

$$\begin{aligned} (m-1) \Psi(\bar{c}_{13}, \alpha_1, \beta) + \Psi(\bar{c}_{13}, \alpha_3, \beta) \\ \leq (m-1) \Psi(c_1, \alpha_1, \beta) + \Psi(c_3, \alpha_3, \beta). \end{aligned}$$

Collecting terms, it is equivalent to show that

$$(m-1)(\Psi(\bar{c}_{13}, \alpha_1, \beta) - \Psi(c_1, \alpha_1, \beta)) \leq \Psi(c_3, \alpha_3, \beta) - \Psi(\bar{c}_{13}, \alpha_3, \beta). \quad (48)$$

Let  $\delta := \bar{c}_{13} - c_1$ , and expanding  $\bar{c}_{13}$  we can verify that  $c_3 - \bar{c}_{13} = (m-1)\delta$ . The right hand side of Equation (48) can be rewritten as

$$\begin{aligned} & \Psi(c_3, \alpha_3, \beta) - \Psi(\bar{c}_{13}, \alpha_3, \beta) \\ &= \Psi(\bar{c}_{13} + (m-1)\delta, \alpha_3, \beta) - \Psi(\bar{c}_{13}, \alpha_3, \beta) \\ &= \sum_{i=1}^{m-1} \Psi(\bar{c}_{13} + i\delta, \alpha_3, \beta) - \Psi(\bar{c}_{13} + (i-1)\delta, \alpha_3, \beta). \end{aligned}$$

In order to show that it is greater or equal than the left hand side of Equation (48), it suffices to show that  $\forall i \in \{1, \dots, m-1\}$ ,

$$\Psi(\bar{c}_{13} + i\delta, \alpha_3, \beta) - \Psi(\bar{c}_{13} + (i-1)\delta, \alpha_3, \beta) \geq \Psi(\bar{c}_{13}, \alpha_1, \beta) - \Psi(\bar{c}_{13} - \delta, \alpha_1, \beta). \quad (49)$$

Both sides of Equation (49) are in the form of the difference between the  $\Psi()$  function value of two points. Since  $\Psi()$  is smooth with respect to  $x$ , the difference can be written as an integral of the derivative  $\Psi'()$  between the two points. Specifically, the left hand side of (49)

$$\begin{aligned} & \Psi(\bar{c}_{13} + i\delta, \alpha_3, \beta) - \Psi(\bar{c}_{13} + (i-1)\delta, \alpha_3, \beta) \\ &= \int_{\bar{c}_{13} + (i-1)\delta}^{\bar{c}_{13} + i\delta} \Psi'(t, \alpha_3, \beta) dt \\ &= \int_{\bar{c}_{13} - \delta}^{\bar{c}_{13}} \Psi'(t + i\delta, \alpha_3, \beta) dt. \end{aligned}$$

And the right hand side

$$\begin{aligned} & \Psi(\bar{c}_{13}, \alpha_1, \beta) - \Psi(\bar{c}_{13} - \delta, \alpha_1, \beta) \\ &= \int_{\bar{c}_{13} - \delta}^{\bar{c}_{13}} \Psi'(t, \alpha_1, \beta) dt. \end{aligned}$$

Thus, Equation (49) is equivalent to

$$\int_{\bar{c}_{13} - \delta}^{\bar{c}_{13}} \Psi'(t + i\delta, \alpha_3, \beta) dt \geq \int_{\bar{c}_{13} - \delta}^{\bar{c}_{13}} \Psi'(t, \alpha_1, \beta) dt. \quad (50)$$

To prove (50), it suffices to show that  $\forall t \in [\bar{c}_{13} - \delta, \bar{c}_{13}]$ ,

$$\Psi'(t + i\delta, \alpha_3, \beta) \geq \Psi'(t, \alpha_1, \beta). \quad (51)$$

We can compute  $\Psi'(x, \alpha, \beta)$  explicitly as

$$\Psi'(x, \alpha, \beta) = \beta - \frac{\alpha + \beta}{1 + e^x}.$$

Thus, (51) is equivalent to

$$\beta - \frac{\alpha_3 + \beta}{1 + e^{t+i\delta}} \geq \beta - \frac{\alpha_1 + \beta}{1 + e^t}$$

or

$$\frac{\alpha_3 + \beta}{1 + e^{t+i\delta}} \leq \frac{\alpha_1 + \beta}{1 + e^t}.$$

Given  $\alpha_1 \geq \alpha_3$  (cf. Equation (43)) and  $t \in [\bar{c}_{13} - \delta, \bar{c}_{13}]$ , this inequality holds, which concludes the proof.  $\square$

Proposition A.4.4 shows that the optimal solution to (47) must satisfy  $c_1^* = c_3^*$ . Therefore, we can substitute  $c_1 = c_3$

and remove  $c_3$  in (47). This reduces  $f_3(c_1, c_2, c_3)$  to (up to a constant scaling)

$$f_2(c_1, c_2) := (m-1)\Psi(c_1, \alpha_1, \beta) + m\Psi(c_2, \alpha_2, \beta) + \Psi(c_1, \alpha_3, \beta).$$

And the optimization problem Equation (47) is equivalently transformed into

$$(c_1^*, c_2^*) = \operatorname{argmin} f_2(c_1, c_2) \quad (52)$$

s.t:  $|c_2| \leq c_1$ .

**Step 3.**

We first consider the unconstrained optimal solution  $(\tilde{c}_1, \tilde{c}_2)$  of optimization problem (52). Since  $\alpha_1, \alpha_2, \alpha_3, \beta > 0$ , all the  $\Psi()$  functions are strictly convex and thus  $f_2(c_1, c_2)$  is strictly convex. The unconstrained optimal solution is unique and can be computed by the  $\nabla f_2(c_1, c_2) = 0$ . Denote  $\bar{\alpha}_{13} := \frac{m-1}{m}\alpha_1 + \frac{1}{m}\alpha_3$ , we have

$$\begin{aligned} \frac{\partial f_2(c_1, c_2)}{\partial c_1} &= m\beta - \frac{m\bar{\alpha}_{13} + m\beta}{1 + e^{c_1}} = 0 \\ \frac{\partial f_2(c_1, c_2)}{\partial c_2} &= m\beta - \frac{m\alpha_2 + m\beta}{1 + e^{c_2}} = 0, \end{aligned}$$

which gives us

$$\begin{aligned} \tilde{c}_1 &= \ln\left(\frac{\bar{\alpha}_{13}}{\beta}\right) \\ \tilde{c}_2 &= \ln\left(\frac{\alpha_2}{\beta}\right). \end{aligned}$$

We claim that  $(\tilde{c}_1, \tilde{c}_2)$  is infeasible. I.e.,  $|\tilde{c}_2| > \tilde{c}_1$ . Given  $k \geq 1$ , with Equation (41), we have

$$\alpha_2 = \frac{1}{4m^2} \left[ w - \sum_{v=1}^w \lambda_2(W)^v \right] < \frac{w}{4m^2} \leq \frac{kw}{4m^2} = \beta.$$

Thus,  $\tilde{c}_2 < 0$  and  $|\tilde{c}_2| = -\tilde{c}_2$ . Therefore, to show that  $(\tilde{c}_1, \tilde{c}_2)$  is infeasible, we only need to prove

$$\ln\left(\frac{\beta}{\alpha_2}\right) > \ln\left(\frac{\bar{\alpha}_{13}}{\beta}\right).$$

Or equivalently,

$$\beta^2 \geq \alpha_2 \bar{\alpha}_{13}. \quad (53)$$

Recall the definition of  $\alpha_1, \alpha_2$ , and  $\alpha_3$  in Equations (40)–(42), we have

$$\bar{\alpha}_{13} := \frac{m-1}{m}\alpha_1 + \frac{1}{m}\alpha_3 = \frac{1}{4m^2} \left[ w + \sum_{v=1}^w \lambda_2(W)^v \right] \quad (54)$$

and

$$\bar{\alpha}_2 = \frac{1}{4m^2} \left[ w - \sum_{v=1}^w \lambda_2(W)^v \right].$$

Therefore,

$$\begin{aligned} \bar{\alpha}_{13}\alpha_2 &= \frac{1}{(4m^2)^2} \left[ w^2 - \left( \frac{\lambda_2(W) - \lambda_2(W)^{w+1}}{1 - \lambda_2(W)} \right)^2 \right] \\ &< \frac{w^2}{(4m^2)^2} \\ &\leq \frac{k^2 w^2}{(4m^2)^2} \end{aligned}$$

$$= \beta^2.$$

Thus, we have shown Equation (53) and thus,  $(\tilde{c}_1, \tilde{c}_2)$  is infeasible.

Since the unconstrained optimal solution  $(\tilde{c}_1, \tilde{c}_2)$  is infeasible, the constrained optimal solution  $(c_1^*, c_2^*)$  must activate the constraint. Next, we will show that the activated constraint must be  $c_1 = -c_2$ .

Denote  $\mathcal{L}$  the line segment joining  $(c_1^*, c_2^*)$  and  $(\tilde{c}_1, \tilde{c}_2)$ , and  $\mathcal{G}$  the feasible set of (52). We first claim that  $\mathcal{L} \cap \mathcal{G} = \{(c_1^*, c_2^*)\}$  must hold. If not, assume there exists  $(c_1^0, c_2^0) \neq (c_1^*, c_2^*)$  and  $(c_1^0, c_2^0) \in \mathcal{L} \cap \mathcal{G}$ . Since  $(c_1^0, c_2^0) \in \mathcal{L}$ , there exist  $\gamma \in (0, 1)$  such that

$$(c_1^0, c_2^0) = \gamma(c_1^*, c_2^*) + (1 - \gamma)(\tilde{c}_1, \tilde{c}_2).$$

Then, by convexity of  $f_2$  and global optimality of  $(\tilde{c}_1, \tilde{c}_2)$ ,

$$f_2(c_1^0, c_2^0) \leq \gamma f_2(c_1^*, c_2^*) + (1 - \gamma)f_2(\tilde{c}_1, \tilde{c}_2) < f_2(c_1^*, c_2^*).$$

It gives us a feasible  $(c_1^0, c_2^0)$  that has a lower cost, which contradicts with the constrained optimality of  $(c_1^*, c_2^*)$ . And thus, by contradiction, we have shown that  $\mathcal{L} \cap \mathcal{G} = \{(c_1^*, c_2^*)\}$ .

Note that, for the global optimizer  $(\tilde{c}_1, \tilde{c}_2)$ , we have  $\tilde{c}_1 > \tilde{c}_2$ . To see this, note it is equivalent to  $\bar{\alpha}_{13} > \alpha_2$ , which is shown from Equations (41) and (54). For any points on the  $\{(c_1, c_2) \mid c_1 = c_2 > 0\}$ , the line segment joining  $(c_1, c_2)$  and  $(\tilde{c}_1, \tilde{c}_2)$  will intersect the feasible set  $\mathcal{G}$  on infinite points, which contradicts with the claim we just proved above. Therefore, the constrained optimizer  $(c_1^*, c_2^*)$  must satisfy  $c_1^* = -c_2^*$ .

Therefore, we can substitute  $c_1 = -c_2$  and remove  $c_2$  in (52). This reduces  $f_2(c_1, c_2)$  to (up to a constant scaling)

$$f_1(c_1) := (m - 1)\Psi(c_1, \alpha_1, \beta) + m\Psi(c_1, \alpha_2, \beta) + \Psi(c_1, \alpha_3, \beta).$$

And the optimization problem reduces to

$$c_1^* = \operatorname{argmin} f_1(c_1). \quad (55)$$

Optimization problem (55) has a unique optimal solution given by  $f_1'(c_1) = 0$ :

$$c_1^* = \ln \left( \frac{\bar{\alpha}_{13} + \beta}{\alpha_2 + \beta} \right).$$

This gives the optimal solution  $X^*$  to (20) when  $\mathcal{H} = \{X \mid X \succeq 0\}$ :

$$X^* = \begin{cases} \ln \left( \frac{\bar{\alpha}_{13} + \beta}{\alpha_2 + \beta} \right), & \text{if } (i, j) \in \mathcal{E}_0 \text{ or } i = j \\ -\ln \left( \frac{\bar{\alpha}_{13} + \beta}{\alpha_2 + \beta} \right), & \text{if } (i, j) \in \mathcal{E}_1. \end{cases}$$

#### Part 4)

Since both  $X^*(\mathbb{R}^{n \times n})$  and  $X^*(\mathbb{S}_+^n)$  are DBC matrices, we can compute their nuclear norms from the proof of Proposition A.4.1. Specifically, for a DBC matrix  $X = Z_{2m}(c_1, c_2, c_3)$ , we have

$$\|X\|_* = |m(c_1 + c_2) + (c_3 - c_1)| + |m(c_1 - c_2) + (c_3 - c_1)| + (2m - 2)|c_3 - c_1|.$$

Since from part 2) we have

$$X^*(\mathbb{R}^{n \times n}) = Z_{2m} \left( \ln \left( \frac{\alpha_1}{\beta} \right), \ln \left( \frac{\alpha_2}{\beta} \right), \ln \left( \frac{\alpha_3}{\beta} \right) \right)$$

with

$$\begin{aligned} \alpha_1 &> \alpha_3 > 0 \\ \alpha_1 &> \alpha_2 > 0, \end{aligned}$$

we have

$$\begin{aligned} \|X^*(\mathbb{R}^{n \times n})\|_* &= m \ln \left( \frac{\alpha_1 \alpha_2}{\beta^2} \right) + \ln \left( \frac{\alpha_3}{\alpha_1} \right) + \\ &\quad \left| m \ln \left( \frac{\alpha_1}{\alpha_2} \right) + \ln \left( \frac{\alpha_3}{\alpha_1} \right) \right| + \\ &\quad (2m - 2) \ln \left( \frac{\alpha_1}{\alpha_3} \right). \end{aligned}$$

Note that  $\alpha_i = \Theta(1/n^2)$  and  $\beta = \Theta(1/n^2)$ . Therefore,  $\|X^*(\mathbb{R}^{n \times n})\|_* = \Theta(n)$ . Similarly, from part 3) we have

$$X^*(\mathbb{S}_+^n) = Z_{2m}(\nu_1, -\nu_1, \nu_1)$$

where  $\nu_1 := \ln \left( \frac{\bar{\alpha}_{13} + \beta}{\alpha_2 + \beta} \right)$  with  $\bar{\alpha}_{13} := \frac{m-1}{m}\alpha_1 + \frac{1}{m}\alpha_3$ . We can get  $\bar{\alpha}_{13} = \Theta(1/n^2)$  which leads to  $\nu_1 = \Theta(1)$ . And

$$\|X^*(\mathbb{S}_+^n)\|_* = 2m\nu_1 = \Theta(n).$$

□

## APPENDIX B NEURAL NETWORK IMPLEMENTATION

In this appendix we detail our neural network implementation of VEC (or ErgoVEC). We first list the structure of the neural network. After that, we will detail our construction of training set (samples and labels). Next, we show that the neural network optimization objective is exactly the objective of VEC (or ErgoVEC). Lastly, we provide the learning and optimization settings we used in training.

**Structure of the neural network.** Figure 9 illustrates the structure of our neural network. There are four layers in the neural network.

- 1) Input layer. This layer receives a one-hot vector encoders for each nodes in a pair  $(i, j) \in \mathcal{V}^2$  as the input of this neural network
- 2) Embedding layer. The embedding layer is a  $n \times d$  matrix where the rows represent the  $d$  dimensional embedding vectors of nodes. These vectors are updated in the optimization iteration after each epoch. After the optimization process, they will be used as final output of the VEC (or ErgoVEC) algorithm. Please note that this is the only layer that will be updated in the entire optimization process. In the neural network, this layer takes the two one-hot vectors from the input layer and returns the two corresponding row vectors to the next layer.
- 3) Dot Product layer. This layer takes two embedding vectors and returns the dot product between them.
- 4) Output layer. This layer takes a scalar (the dot product from previous layer), and returns the sigmoid function value  $S(x) := \frac{1}{1+e^{-x}}$  of it as the output of this neural network.

To sum up, this neural network takes a pair of nodes  $(i, j)$  as input and returns the sigmoid function of their dot product  $\hat{y}_{(i,j)}$  as output.

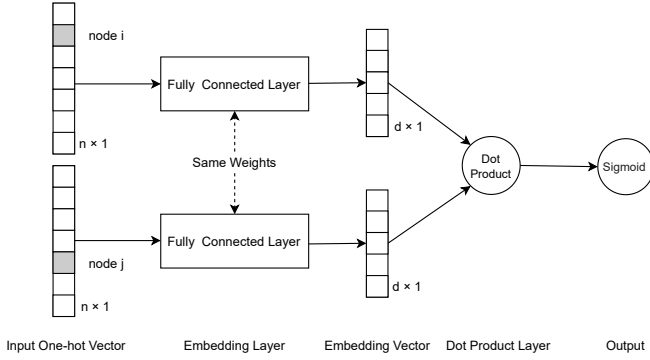


Fig. 9: Structure diagram of the neural network implemented with Keras package.

**Training set and loss function.** We use a weighted training set  $D = \{(i, j), w_{(i,j)}, y_{(i,j)}\}$  obtained from the union of two parts: a weighted positive set and a weighted negative set. Both sets contain all node pairs  $(i, j) \in \mathcal{V}^2$ , but label and weigh them differently. All node pairs in the positive set are labeled 1 with weights equal  $n_{ij}^+$  ( $\bar{n}_{ij}^+$  for ErgoVEC), whereas node pairs in the negative set are labeled as 0 with weights equal  $n_{ij}^-$  ( $\bar{n}_{ij}^-$  for ErgoVEC). The training set is randomly shuffled and fed in the neural network during each epoch. For the loss function, we choose binary cross entropy

$$H(D) := \frac{1}{N} \sum_{(i,j) \in D} H(i, j),$$

where

$$H(i, j) = -w_{(i,j)} \left[ y_{(i,j)} \ln(\hat{y}_{(i,j)}) + (1 - y_{(i,j)}) \ln(1 - \hat{y}_{(i,j)}) \right]. \quad (56)$$

**Equivalence proof.** Here we show that the neural network equipped with this training set and loss function has the exact same objective as VEC. (For ErgoVEC, the same holds after replacing  $n_{ij}^+$  with  $\bar{n}_{ij}^+$  in the following equations.) First note that, for  $(i, j)$  in positive set,  $y_{(i,j)} = 1$ ,

$$\begin{aligned} H(i, j) &= -n_{ij}^+ \ln(\hat{y}_{(i,j)}) = -n_{ij}^+ \ln(S(\mathbf{u}_i^\top \mathbf{u}_j)) \\ &= n_{ij}^+ \sigma(+\mathbf{u}_i^\top \mathbf{u}_j) \end{aligned}$$

and for  $(i, j)$  in negative set,  $y_{(i,j)} = 0$ ,

$$\begin{aligned} H(i, j) &= -n_{ij}^- \ln(1 - \hat{y}_{(i,j)}) = -n_{ij}^- \ln(1 - S(\mathbf{u}_i^\top \mathbf{u}_j)) \\ &= n_{ij}^- \sigma(-\mathbf{u}_i^\top \mathbf{u}_j). \end{aligned}$$

Therefore,

$$\begin{aligned} H(D) &:= \frac{1}{N} \sum_{(i,j) \in D} H(i, j) \\ &= \sum_{(i,j) \in \mathcal{V}^2} \left[ n_{ij}^+ \sigma(\mathbf{u}_i^\top \mathbf{u}_j) + n_{ij}^- \sigma(-\mathbf{u}_i^\top \mathbf{u}_j) \right], \end{aligned}$$

which is the same as (1).

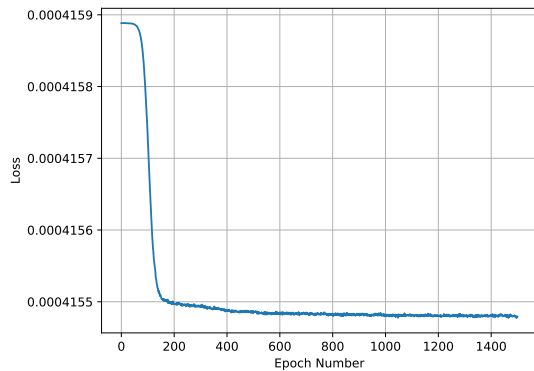
**Optimization paramters.** We used the Adam optimizer with default parameter choice except for learning rate. We set learning rate as described in Table 1, although we want to make a note that the optimal learning rates do depend on specific graph realizations.

**Remarks on convergence.** In our experiments, we note

	Algorithm	$n$	l.r.	# epochs
Linear Degree Regime	VEC	100	0.001	400
	VEC	200	0.001	200
	VEC	500	0.001	80
	VEC	1000	0.001	40
	ErgoVEC	100	0.02	400
	ErgoVEC	200	0.02	200
	ErgoVEC	500	0.02	80
	ErgoVEC	1000	0.02	40
Logarithmic Degree Regime	VEC	100	0.001	400
	VEC	200	0.00021	1500
	VEC	500	0.001	200
	VEC	1000	0.001	200
	ErgoVEC	100	0.0025	400
	ErgoVEC	200	0.00021	1500
	ErgoVEC	500	0.0025	200
	ErgoVEC	1000	0.0025	200

TABLE 1: Learning rates and number of epochs used in each experiment.

that the objective functions seem to converge after a number of epochs, but the embedding vectors do not. The convergence behavior over epochs is shown in Fig. 10 with the plot of the loss function as a function of number of epochs displayed in Fig. 10 (a). Changes in the embedding vectors measured by the ratio of the Procrustes distance between embedding vectors in consecutive epochs and the Frobenius norm of the embedding vectors in the previous epoch are displayed in Fig. 10 (b). We observe that the loss function drops quickly after the first few epochs and remains essentially flat after 1000 epochs, but the change in the embedding vectors is bounded away from 0 even after 1500 epochs. A possible explanation for this behavior is that many neural network implementations and optimization procedures, including the Keras package that we used, focus on the convergence of the objective loss rather than the convergence of layer weights. Although this is very useful in various applications, it may be inadequate for finding the optimal numerical solution (the minimizing weights). Future work could attempt improving our implementation to overcome such limitations.



(a) Loss function value as a function of number of epochs.

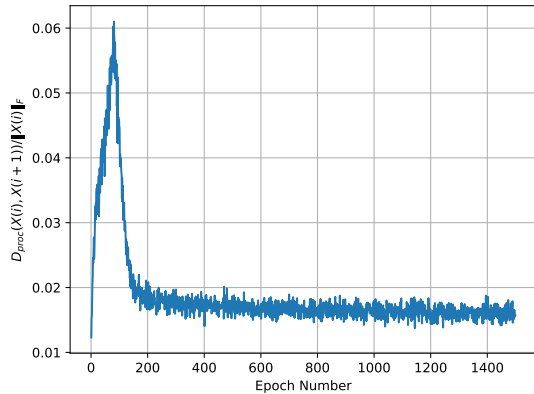
(b) Change in embedding vectors versus epochs. The change is computed as the ratio of the Procrustes distance between embedding vectors in epoch  $i$  and  $i + 1$  and the Frobenius norm of the embedding vectors in epoch  $i$ .

Fig. 10: Illustrating potential convergence issues associated with neural-network-based optimization of node embedding objectives.

## ACKNOWLEDGMENTS

This work was supported in part by the U.S. National Science Foundation under grant 1527618, the Department of Electrical and Computer Engineering, and the Division of Systems Engineering at Boston University. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the supporting institutions.

## REFERENCES

- [1] H. Cai, V. W. Zheng, and K. C.-C. Chang, "A comprehensive survey of graph embedding: Problems, techniques, and applications," *IEEE Transactions on Knowledge and Data Engineering*, vol. 30, no. 9, pp. 1616–1637, 2018.
- [2] B. Perozzi, R. Al-Rfou, and S. Skiena, "Deepwalk: Online learning of social representations," in *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2014, pp. 701–710.
- [3] Z. Yang, W. Cohen, and R. Salakhudinov, "Revisiting semi-supervised learning with graph embeddings," in *International conference on machine learning*. PMLR, 2016, pp. 40–48.
- [4] A. Grover and J. Leskovec, "node2vec: Scalable feature learning for networks," in *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2016, pp. 855–864.
- [5] W. Ding, C. Lin, and P. Ishwar, "Node embedding via word embedding for network community discovery," *IEEE Transactions on Signal and Information Processing over Networks*, vol. 3, no. 3, pp. 539–552, 2017.
- [6] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Advances in neural information processing systems*, 2013, pp. 3111–3119.
- [7] J. Pennington, R. Socher, and C. Manning, "Glove: Global vectors for word representation," in *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, 2014, pp. 1532–1543.
- [8] A. Bakarov, "A survey of word embeddings evaluation methods," *arXiv preprint arXiv:1801.09536*, 2018.
- [9] K. Rohe, S. Chatterjee, B. Yu *et al.*, "Spectral clustering and the high-dimensional stochastic blockmodel," *Annals of Statistics*, vol. 39, no. 4, pp. 1878–1915, 2011.
- [10] D. L. Sussman, M. Tang, D. E. Fishkind, and C. E. Priebe, "A consistent adjacency spectral embedding for stochastic blockmodel graphs," *Journal of the American Statistical Association*, vol. 107, no. 499, pp. 1119–1128, 2012.
- [11] T. Qin and K. Rohe, "Regularized spectral clustering under the degree-corrected stochastic blockmodel," in *Advances in Neural Information Processing Systems*, 2013, pp. 3120–3128.
- [12] A. Athreya, D. E. Fishkind, M. Tang, C. E. Priebe, Y. Park, J. T. Vogelstein, K. Levin, V. Lyzinski, and Y. Qin, "Statistical inference on random dot product graphs: a survey," *Journal of machine learning research: JMLR*, vol. 18, no. 1, pp. 8393–8484, Jan. 2017.
- [13] K. Chaudhuri, F. Chung, and A. Tsiatas, "Spectral partitioning of graphs with general degrees and the extended planted partition model," in *Proceedings of the 25th conference on learning theory*, vol. 2906, 2012.
- [14] J. Cape, M. Tang, and C. E. Priebe, "On spectral embedding performance and elucidating network structure in stochastic blockmodel graphs," *Network Science*, vol. 7, no. 3, pp. 269–291, 2019.
- [15] J. Qiu, Y. Dong, H. Ma, J. Li, K. Wang, and J. Tang, "Network embedding as matrix factorization: Unifying deepwalk, line, pte, and node2vec," in *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*. ACM, 2018, pp. 459–467.
- [16] J. Tang, M. Qu, M. Wang, M. Zhang, J. Yan, and Q. Mei, "Line: Large-scale information network embedding," in *Proceedings of the 24th international conference on world wide web*. International World Wide Web Conferences Steering Committee, 2015, pp. 1067–1077.
- [17] M. Girvan and M. Newman, "Girvan, m. & newman, m. e. j. community structure in social and biological networks. *proc. natl acad. sci. usa* 99, 7821-7826," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 99, pp. 7821–6, 07 2002.
- [18] D. A. Spielman and S.-H. Teng, "Nearly-linear time algorithms for graph partitioning, graph sparsification, and solving linear systems," in *Proceedings of the thirty-sixth annual ACM symposium on Theory of computing*. ACM, 2004, pp. 81–90.
- [19] R. Andersen, F. Chung, and K. Lang, "Local graph partitioning using pagerank vectors," in *Foundations of Computer Science, 2006. FOCS'06. 47th Annual IEEE Symposium on*. IEEE, 2006, pp. 475–486.
- [20] R. Lambiotte, J.-C. Delvenne, and M. Barahona, "Random walks, markov processes and the multiscale modular organization of complex networks," *IEEE Transactions on Network Science and Engineering*, vol. 1, no. 2, pp. 76–90, 2014.
- [21] L. Meng and N. Masuda, "Analysis of node2vec random walks on networks," *Proceedings of the Royal Society A*, vol. 476, no. 2243, p. 20200447, 2020.
- [22] Y. Zhang and M. Tang, "Consistency of random-walk based network embedding algorithms," *arXiv preprint arXiv:2101.07354*, 2021.
- [23] L. Bottou, "Large-scale machine learning with stochastic gradient descent," in *Proceedings of COMPSTAT 2010*. Springer, 2010, pp. 177–186.
- [24] —, "Online algorithms and stochastic approximations," in *Online Learning and Neural Networks*, D. Saad, Ed. Cambridge, UK: Cambridge University Press, 1998, revised, oct 2012. [Online]. Available: <http://leon.bottou.org/papers/bottou-98x>
- [25] B. Recht, C. Re, S. Wright, and F. Niu, "Hogwild: A lock-free approach to parallelizing stochastic gradient descent," in *Advances in Neural Information Processing Systems*, 2011, pp. 693–701.

- [26] H. White, S. Boorman, and R. Breiger, "Social structure from multiple networks, blockmodels of roles and positions," *American Journal of Sociology*, pp. 730–780, 1976.
- [27] P. Holland, K. Laskey, and S. Leinhardt, "Stochastic blockmodels: First steps," *Social Networks*, vol. 5, no. 2, pp. 109–137, 1983.
- [28] R. Boppana, "Eigenvalues and graph bisection: An average-case analysis," in *Proc. of the 28th Annual Symposium on Foundations of Computer Science (FOCS)*, 1987, pp. 280–285.
- [29] E. Abbe, A. S. Bandeira, and G. Hall, "Exact recovery in the stochastic block model," *IEEE transactions on information theory / Professional Technical Group on Information Theory*, vol. 62, no. 1, pp. 471–487, Jan. 2016.
- [30] E. Abbe and C. Sandon, "Community detection in general stochastic block models: Fundamental limits and efficient algorithms for recovery," in *Proc. of the 56th Annual Symposium on Foundations of Computer Science (FOCS)*, Sep. 2015, pp. 670–688.
- [31] —, "Detection in the stochastic block model with multiple clusters: proof of the achievability conjectures, acyclic bp, and the information-computation gap," in *Advances in Neural Information Processing Systems (NIPS)*, Dec. 2016.
- [32] A. Decelle, F. Krzakala, C. Moore, and L. Zdeborová, "Asymptotic analysis of the stochastic block model for modular networks and its algorithmic applications," *Physical Review E*, vol. 84, no. 6, p. 066106, 2011.
- [33] E. Mossel, J. Neeman, and A. Sly, "Belief propagation, robust reconstruction and optimal recovery of block models," in *Proc. of the 27th Conference on Learning Theory (COLT)*, 2014, pp. 356–370.
- [34] —, "Reconstruction and estimation in the planted partition model," *Probability Theory and Related Fields*, vol. 162, no. 3-4, pp. 431–461, Aug. 2015.
- [35] O. Levy and Y. Goldberg, "Neural word embedding as implicit matrix factorization," in *Advances in neural information processing systems*, 2014, pp. 2177–2185.
- [36] K. Church and P. Hanks, "Word association norms, mutual information, and lexicography," *Computational linguistics*, vol. 16, no. 1, pp. 22–29, 1990.
- [37] M. E. Newman and M. Girvan, "Finding and evaluating community structure in networks," *Physical review E*, vol. 69, no. 2, p. 026113, 2004.
- [38] M. E. J. Newman, "Spectral methods for community detection and graph partitioning," *Physical review. E, Statistical, nonlinear, and soft matter physics*, vol. 88, no. 4, p. 042822, Oct. 2013.
- [39] M. Fazel, H. Hindi, and S. Boyd, "Rank minimization and applications in system theory," in *American Control Conference, 2004. Proceedings of the 2004*, vol. 4. IEEE, 2004, pp. 3273–3278.
- [40] J. Dong, Z. Xue, J. Guan, Z.-F. Han, and W. Wang, "Low rank matrix completion using truncated nuclear norm and sparse regularizer," *Signal Processing: Image Communication*, vol. 68, pp. 76–87, 2018.
- [41] C. Lu, J. Feng, Y. Chen, W. Liu, Z. Lin, and S. Yan, "Tensor robust principal component analysis with a new tensor nuclear norm," *IEEE transactions on pattern analysis and machine intelligence*, vol. 42, no. 4, pp. 925–938, 2019.
- [42] B. Recht, M. Fazel, and P. A. Parrilo, "Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization," *SIAM review*, vol. 52, no. 3, pp. 471–501, 2010.
- [43] E. Hazan, "Sparse approximate solutions to semidefinite programs," in *Latin American symposium on theoretical informatics*. Springer, 2008, pp. 306–316.
- [44] A. Joseph, B. Yu *et al.*, "Impact of regularization on spectral clustering," *The Annals of Statistics*, vol. 44, no. 4, pp. 1765–1791, 2016.
- [45] U. Von Luxburg, "A tutorial on spectral clustering," *Statistics and computing*, vol. 17, no. 4, pp. 395–416, 2007.
- [46] R. Rehůřek and P. Sojka, "Software Framework for Topic Modelling with Large Corpora," in *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*. Valletta, Malta: ELRA, May 2010, pp. 45–50.
- [47] F. Chollet *et al.*, "Keras," <https://keras.io>, 2015.
- [48] S. M. Ross, J. J. Kelly, R. J. Sullivan, W. J. Perry, D. Mercer, R. M. Davis, T. D. Washburn, E. V. Sager, J. B. Boyce, and V. L. Bristow, *Stochastic processes*. Wiley New York, 1996, vol. 2.
- [49] C. McDiarmid, "On the method of bounded differences," *Surveys in combinatorics*, vol. 141, no. 1, pp. 148–188, 1989.
- [50] G. Grimmett and D. Stirzaker, *Probability and random processes, 3rd Edition*. Oxford university press, 2001.