## Appendix - Community Detection Guarantees using Embeddings Learned by Node2Vec

The Appendix consists of the proofs of the results stated within the paper, along with some extra discussions which would detract from the flow of the main paper. We also provide some additional simulation results relating to node classification, and further simulated and real data experiments examining community detection.

## A  Additional Experimental Results

Here we provide additional details describing the experimental results presented in the main paper. We also describe additional experiments. All experiments were run on a computing cluster utilising 4 cores of an Intel E5-2683 v4 Broadwell 2.1GHz CPU or similar with 2 GB of memory per core. Each individual experimental run required at most 2 hours of computing time. All experiments, including initial preliminary experiments, required approximately 25k CPU hours. All code required to reproduce all results is included in the code repository in the supplemental files.

**Additional Simulation, Node Classification**  We provide a simple experiment to support the theoretical results on node classification demonstrated in Section D of the appendix. We simulate data from a $\text{SBM}(n/\kappa, \kappa, \tilde{p}, \tilde{q}, \rho_n)$ as before with $\tilde{q} = \tilde{p}\beta$ as in the main text. We learn an embedding of each node using node2vec with embedding dimension of 64 and all other parameters set at their default values. We then use the true community labels of 10% of these nodes to train a (multinomial) logistic regression classifier, and predict the class label for the remaining 90% of nodes in the network. We examine the performance of this classification tool using the node2vec embeddings in terms of classification accuracy. We show these results in Figure S1 for $\rho_n = \log(n)/n$, with 10 simulations for each setting, with the mean across these simulations and error bars indicating one standard error. This classifier has excellent accuracy at predicting the labels of other nodes.
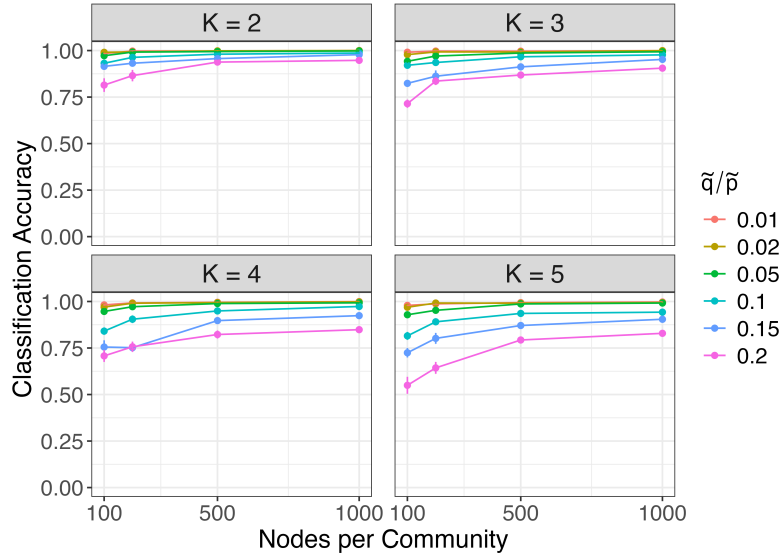


Figure S1: Classification accuracy using 10% of the node embeddings to learn a multinomial logistic regression classifier. Mean and one standard error shown.

**Additional Results, Community Detection**  Here we include additional simulation results which were omitted from the main text. In particular, for the simulations considered in the main manuscript we now examine the community recovery performance in terms of the normalized mutual information [9]. We show the average NMI score across these simulations, along with error bars corresponding to one standard error. In each case, the NMI metric is similar to the proportion of nodes correctly recovered. As we increase the number of nodes this performance improves.
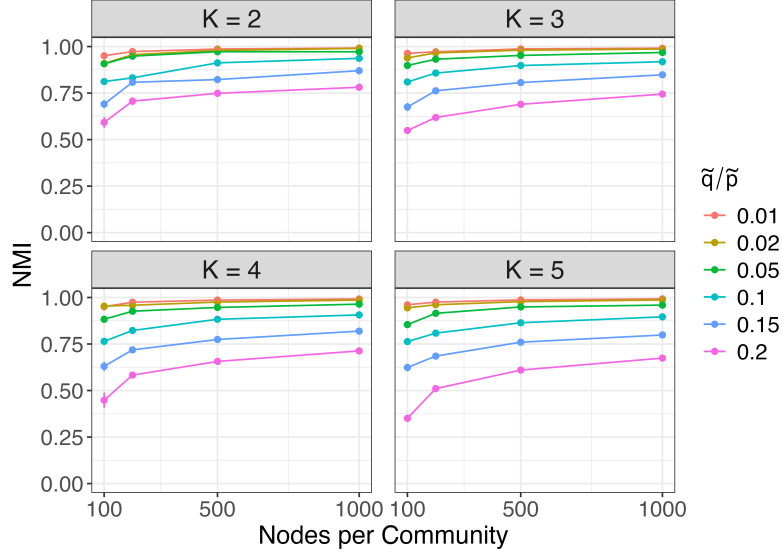
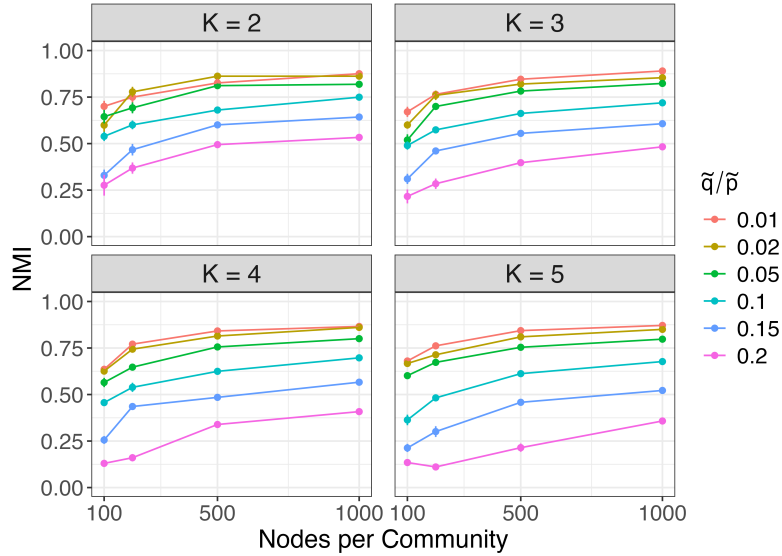Figure S2: NMI for relatively sparse SBM. Mean and one standard error shown.



Figure S3: NMI for relatively sparse DC-SBM. Mean and one standard error shown.

**Rates of Convergence**    We can also investigate the empirical convergence of these methods. Here, we consider the same simulated SBM data as above, and examine the convergence in the proportion of nodes correctly recovered, as we increase the number of nodes in the network, for $\kappa = 2, 3, 4, 5$. We empirically investigate this convergence using a log-log plot, which is shown in Figure S5 for a relatively sparse SBM. Our node2vec procedures demonstrates empirical convergence which is super-linear for dense networks while being sub-linear for relatively sparse networks.

**Varying the node2vec walk parameters**    We also wish to examine the performance of our proposed clustering procedure when the parameters of the random walk are varied. While $p$ and $q$ are both commonly chosen to be $1$, resulting in a simple random walk, other values are possible. We consider data simulated from the relatively sparse DC-SBM considered previously with $\kappa = 2$ communities and consider the within between community probability ratio $\beta = .01$ and $\beta = 0.2$, corresponding to an easier and harder setting to recover the communities respectively. We then consider $p, q \in \{0.5, 1, 2\}$, the common possible values and vary the number of nodes in each community as before. For each
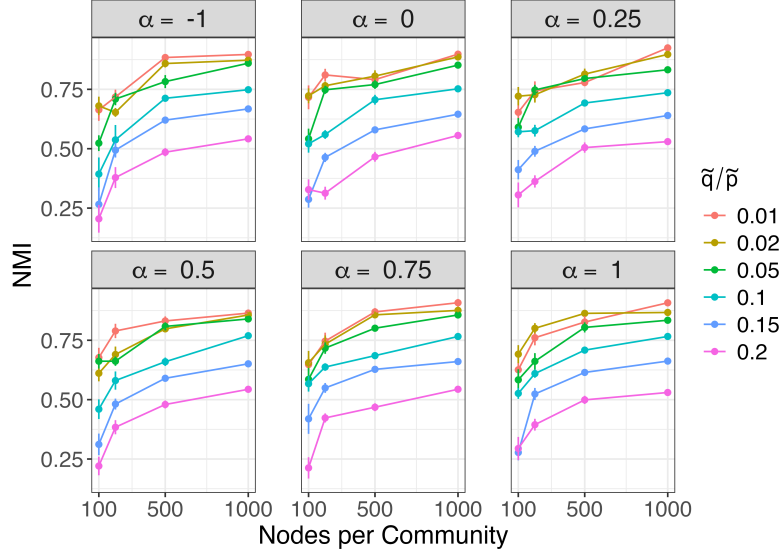
14

Figure S4: NMI varying $\alpha$ for relatively sparse DC-SBM. Mean and one standard error shown.
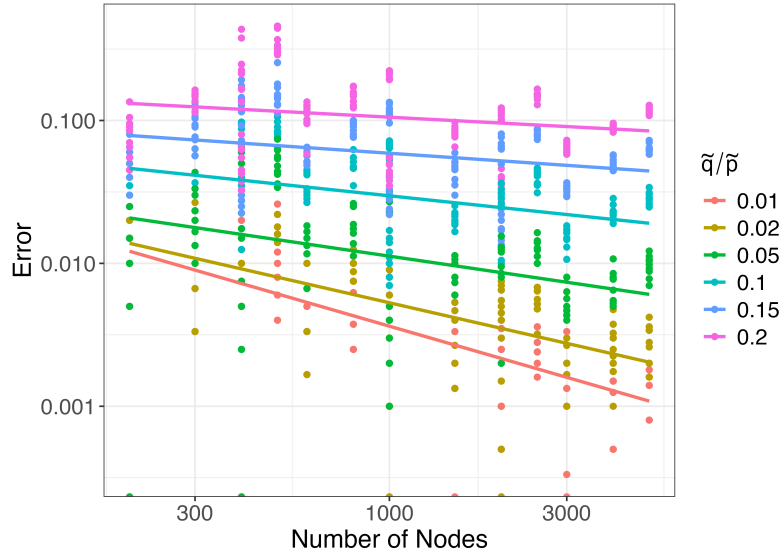


Figure S5: Log-Log plot showing the rate of convergence as we increase the number of nodes in the network. We show a fitted regression for each of the values of $\beta$, showing better convergence when the difference between the within and between community edge probabilities is higher.

of these settings we perform community detection using node2vec and spectral clustering. When $\beta = 0.01$ weobtain excellent community recovery for all values of $p$ and $q$, as shown in Figure S6(a). When $\beta = 0.2$ community recovery is more challenging for small networks for all values of $p$ and $q$. As the number of nodes increases, Figure S6(b) shows that all choices of $p$ and $q$ result in good performance.

### A.1   Performance on Real Networks

We wish to further examine the performance of this community detection procedure for real networks, with known community structure. We also wish to compare this procedure to spectral clustering, which is widely used in practice for community detection. We use two publicly available networks containing known community structure. We first consider a network of emails between 1005 members
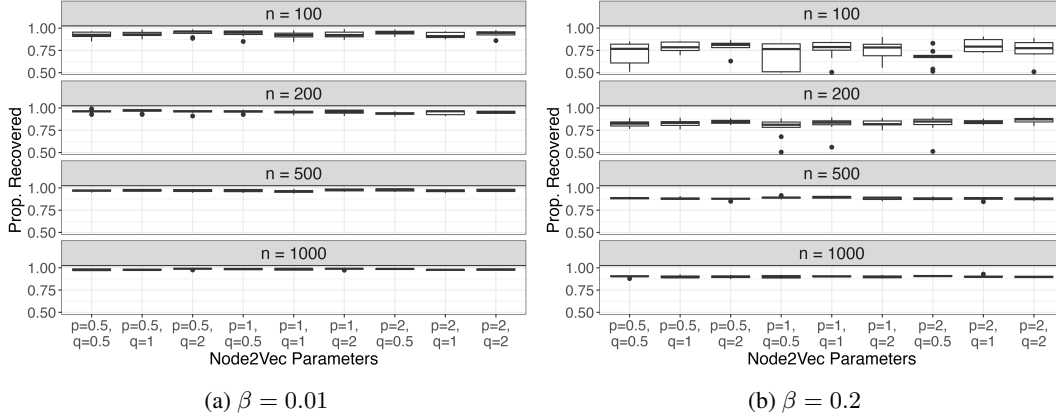
15

Figure S6: Varying the node2vec sampling parameters for DC-SBMs with $\beta = 0.01$ (left) and $\beta = 0.2$ (right). Community recovery is harder when $\beta$ is larger and this is seen for all values of $p$ and $q$ for small networks. As the number of nodes increases we get good community recovery for all choices of $p$ and $q$.

of a large research institution, available as part of the Stanford Network Analysis Project [30]. There are 25571 directed edges between the nodes in this network, with known ground truth communities consisting of 42 departments present in this institution. We also consider a widely used dataset of directed edges between 1490 U.S political blogs, collected before the 2004 elections [2]. Here the directed edges correspond to hyperlinks, with ground truth communities corresponding to whether the blogs has been identified as liberal or conservative.

For each of these datasets we compare the community recovery of Node2Vec and traditional spectral clustering, using the normalized graph Laplacian. As is common in the literature, we remove the direction from these edges and take the largest connected component, forming symmetric adjacency matrices with 986 and 1222 nodes respectively. We then use the previously described procedure to perform community detection using Node2Vec. We consider a range of embedding dimensions ($d = 16, 32, 64, 128, 256$) and unigram sampling parameter ($\alpha = -1, 0.0, 0.25, 0.5, 0.75, 1.0$), while keeping all other parameters fixed at the defaults considered before. With the true number of communities known, we then compare the estimated communities from 10 simulations for each of these parameter settings, along with performing 10 simulations of spectral clustering for each of these settings.

In Figure S7 we compare the performance of Node2Vec and spectral clustering for the Email network and in Figure S8 we use the Political Blogs network. We measure community recovery in terms of the normalized mutual information (NMI) between the estimated and true communities. Other metrics such as the adjusted rand index (ARI) showing similar results. In each case the communities estimated by Node2Vec are substantially closer to the true communities than those estimated by spectral clustering. As highlighted by Karrer and Newman [22] for the political blog data, models which do not account for degree heterogeneity can struggle to recover the underlying community structure. As shown in Figure S8, spectral clustering is unable to recover the communities due to this heterogeneity, while clustering using the Node2Vec embedding shows strong performance at community recovery.

We also further expand on the role of the embedding parameters in the performance of Node2Vec on these real networks. In Figure S9 we examine community recovery for the Email data as we vary the embedding dimension $d$ and the unigram sampling parameter $\alpha$. As we vary each of these parameters we see good community recovery in all settings. For this dataset all choices of embedding dimension and unigram parameter give good NMI scores.

# B  Additional Notation

We give a brief recap of some of the notation introduced in the main paper, along with some more notation which is used purely within the Supplemntary Material. Throughout, we will suppose that
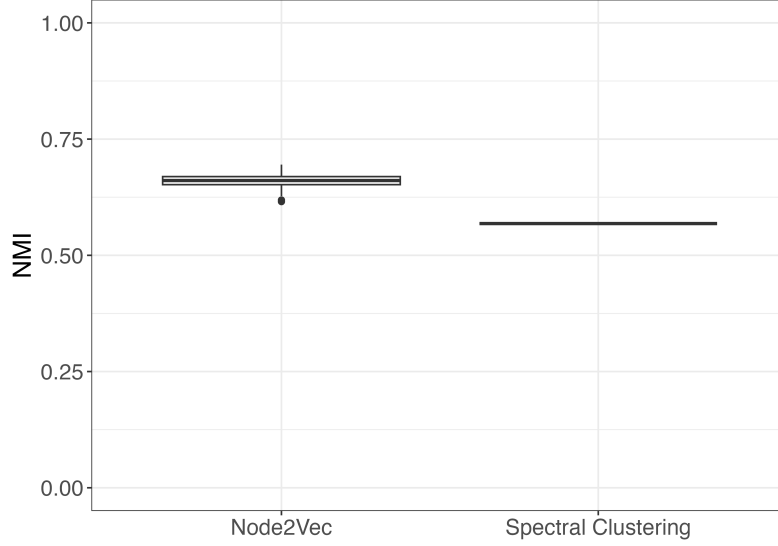
16

Figure S7: Community recovery for the Email data, using both Node2Vec and Spectral Clustering. Node2Vec can better recover the true communities.
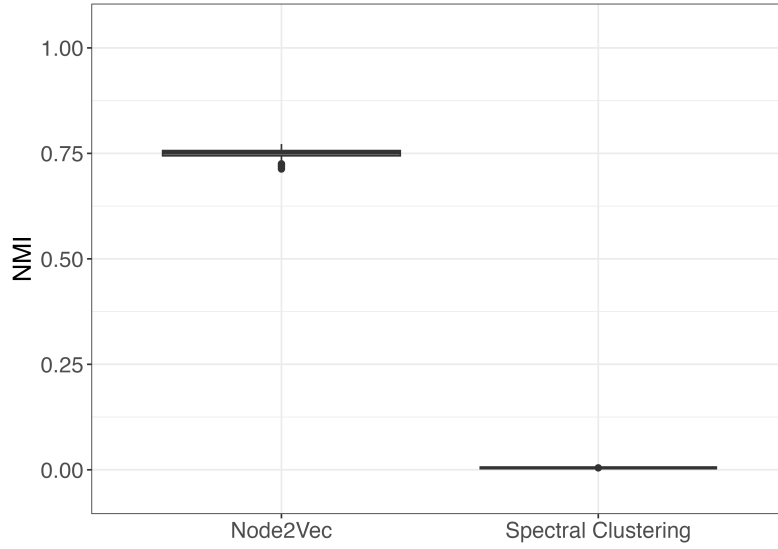


Figure S8: Community recovery for the Political Blog data, using both Node2Vec and Spectral Clustering. Node2Vec can better recover the true communities.

the graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is drawn according to the following generative model: each vertex $u \in \mathcal{V}$ have latent variables $\lambda_u = (c(u), \theta_u)$ where $c(u) \in [\kappa]$ is a community assignment, and $\theta_u$ is a degree-heterogenity correction factor. We then suppose that the edges $a_{uv} \in \{0, 1\}$ in the graph $\mathcal{G}_n$ on $n$ vertices arise independently with probability

$$\mathbb{P}(a_{uv} = 1 \mid \lambda_u, \lambda_v) = \rho_n \theta_u \theta_v P_{c(u), c(v)} \tag{S1}$$

for $u < v$, with $a_{uv} = a_{vu}$ by symmetry for $u > v$[2]. The factor $\rho_n$ accounts for sparsity in the network. The above model corresponds to a degree corrected stochastic block model [22]; we

---

[2]To prevent notation overloading when $A$ is used to indicate constants, we use $a_{uv}$ to describe the presence or absence of an edge between nodes $u$ and $v$ in the supplement, rather than $A_{uv}$ which was used in the main text.

(a) Varying the embedding dimension used.      (b) Varying the Unigram Parameter $\alpha$.
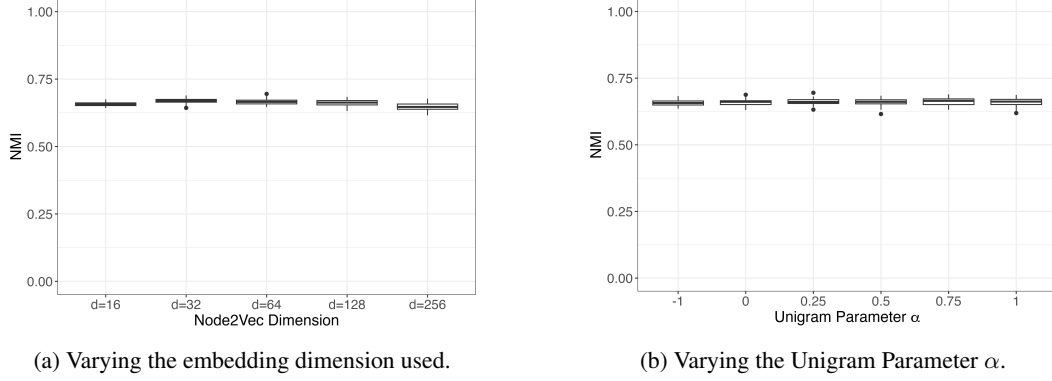
Figure S9: The effect of different Node2Vec parameters on community recovery, measure in terms of Normalized Mutual Information (NMI), for the Email Data.

highlight that the case where $\theta_u$ is constant across all $u \in \mathcal{V}$ corresponds to the original stochastic block model [19]. For convenience, we will write

$$W(\lambda_u, \lambda_v) = \theta_u \theta_v P_{c(u),c(v)} \qquad \text{so} \qquad \mathbb{P}(A_{uv} = 1 \mid \lambda_u, \lambda_v) = \rho_n W(\lambda_u, \lambda_v). \qquad \text{(S2)}$$

We then introduce the notation

$$W(\lambda_i, \cdot) := \mathbb{E}[W(\lambda_i, \lambda_j) \mid \lambda_i], \qquad \mathcal{E}_W(\alpha) := \mathbb{E}[W(\lambda_i, \cdot)^\alpha] \text{ for } \alpha > 0. \qquad \text{(S3)}$$

Note that under the assumptions that the community assignments are drawn i.i.d from a Categorical$(\pi)$ random variable, and the degree correction factors are drawn i.i.d from a distribution $\vartheta$ independently of the community assignments, we have

$$W(\lambda_i, \cdot) = \theta_i \cdot \mathbb{E}[\theta] \cdot \mathbb{E}_{j \sim \text{Cat}(\pi)}[P_{c(i),j} \mid c(i)] = \theta_i \cdot \mathbb{E}[\theta] \cdot \sum_{j=1}^{\kappa} \pi_j P_{c(i),j}, \qquad \text{(S4)}$$

$$\mathcal{E}_W(\alpha) = \mathbb{E}[\theta^\alpha] \cdot \mathbb{E}[\theta]^\alpha \cdot \sum_{i=1}^{\kappa} \pi_i \Big( \sum_{j=1}^{\kappa} \pi_j P_{i,j} \Big)^\alpha \qquad \text{(S5)}$$

For convenience, we will write $\widetilde{P}_{c(i)} = \sum_{j=1}^{\kappa} \pi_j P_{c(i),j}$.

Recall that node2vec attempts to minimize the objective

$$\mathcal{L}_n(U, V) := \sum_{i \neq j} \Big\{ -\mathbb{P}\big((i,j) \in \mathcal{P}(\mathcal{G}_n) \mid \mathcal{G}_n\big) \log(\sigma(\langle u_i, v_j \rangle))$$

$$- \mathbb{P}\big((i,j) \in \mathcal{N}(\mathcal{G}_n) \mid \mathcal{G}_n\big) \log(1 - \sigma(\langle u_i, v_j \rangle)) \Big\}$$

where $U, V \in \mathbb{R}^{n \times d}$, with $u_i, v_j \in \mathbb{R}^d$ denoting the $i$-th and $j$-th rows of $U$ and $V$ respectively, and $\sigma(x) := (1 + e^{-x})^{-1}$ denoting the sigmoid function. Here $\mathcal{P}$ and $\mathcal{N}$ correspond to the positive and negative sampling schemes induced by the random walk and unigram mechanisms respectively.

# C    Proof of Theorems 2 and 3

## C.1    Proof overview

To give an overview of the proof approach, we work by forming successive approximations to the function $\mathcal{L}_n(U, V)$ where we have uniform convergence of the approximation error as $n \to \infty$ over either level sets of the function considered, or the overall domain of optimization of the embedding matrices $U$ and $V$. We break these approximations up into multiple steps:

1. Theorems S1, S2, S3 and Proposition S4 - We begin by working with an approximation $\widehat{\mathcal{L}}_n(U, V)$ of $\mathcal{L}_n(U, V)$, where the sampling weights $\mathbb{P}\big((i,j) \in \mathcal{P}(\mathcal{G}_n) \mid \mathcal{G}_n\big)$ and $\mathbb{P}\big((i,j) \in \mathcal{N}(\mathcal{G}_n) \mid \mathcal{G}_n\big)$ are replaced by functions of the latent variables $(\lambda_i, \lambda_j)$ of the vertices $i$ and $j$, along with $a_{ij}$ in the case of $f_{\mathcal{P}}(\lambda_i, \lambda_j)$.

18

2. The resulting approximation $\widehat{\mathcal{L}}_n(U, V)$ has a dependence on the adjacency matrix of the network. We argue that this loss function converges uniformly to its average over the adjacency matrix when the vertex latent variables remain fixed; this is the contents of Theorem S5.

3. So far, the loss function only looks between interactions of $u_i$ and $v_j$ for $i \neq j$. For theoretical purposes, it is more convenient to work with a loss function where the term with $i = j$ is included. This is handled within Lemma S6.

4. Now that we have an averaged version of the loss function to work with, we are able to examine the minima of this loss function, and find that there is a unique minima (in the sense that for any pair of optima matrices $U^*$ and $V^*$, the matrix $U^*(V^*)^T$ is unique). Moreover, in certain circumstances we can give closed forms for these minima. This is the contents of Section C.6.

5. This is then all combined together in order to give Theorems S13 and S14, which correspond to Theorems 1 and 2 of the main text.

We recap that we consider three scenarios - referred to as Scenario (i), (ii) and (iii) throughout - when proving the following result:

(i) We use DeepWalk ($p = q = 1$ in node2vec), and the graph is drawn according to a SBM with $\rho_n \gg \log(n)/n$;

(ii) We use node2vec, and the graph is drawn according to a SBM with $\rho_n = n^{-\alpha}$ for some $\alpha < \alpha'$, where $\alpha'$ depends on node2vec's hyperparameters;

(iii) We use DeepWalk and a unigram parameter of $\alpha = 1$, and the graph is drawn according to a DCSBM with $\rho_n \gg \log(n)/n$ where the degree heterogeneity parameters $\theta_u \in [C^{-1}, C]$ for some $C > \infty$.

Generally speaking, the approach is the exact same for all three scenarios. As we have a closed formula in the case where we examine DeepWalk, we will consistently provide the details for the DeepWalk case first, and then discuss afterwards how the results and proofs change (if at all) when considering node2vec in generality. Throughout, we also contextualize the proof by examining what it says for a $\mathrm{SBM}(n, \kappa, \tilde{p}, \tilde{q}, \rho_n)$ model. This corresponds to a balanced network with $\pi = (\kappa^{-1}, \dots, \kappa^{-1})$.

## C.2   Replacing the sampling weights

Before giving an approximation to $\mathcal{L}_n(U, V)$, we need to first come up with approximate forms of $\mathbb{P}\big((i, j) \in \mathcal{P}(\mathcal{G}_n) \,|\, \mathcal{G}_n\big)$ and $\mathbb{P}\big((i, j) \in \mathcal{N}(\mathcal{G}_n) \,|\, \mathcal{G}_n\big)$. The next three results give examples of this. In this section we prove three main results. The first two give us guarantees for the sampling probabilities of vertex pairs $(u, v)$ for node2vec for any choice of the hyperparameters $(p, q)$. In particular they will allow us to argue that when the underlying graph arises from a SBM, the sampling probabilities asymptotically depend only on the underlying communities. The last specializes this to the case of DeepWalk (where $p = q = 1$), which has enough structure to allow us to get some additional information, such as closed formula for these sampling probabilities, which can be used in the case where the graph arises through a DCSBM.

**Theorem S1.** *There exists $\alpha$ sufficiently small, depending on the walk length $k$, such that if $\rho_n = n^{-\alpha}$ then there exists a symmetric measurable (with respect to the sigma field generated by $W$) function $f_{\mathcal{P}}(\lambda, \lambda')$ which is bounded below away from zero, and bounded above by $C\rho_n^{-1}$ for some constant $C < \infty$, such that*

$$\max_{i \neq j} \left| \frac{n^2 \mathbb{P}\big((i, j) \in \mathcal{P}(\mathcal{G}_n) \,|\, \mathcal{G}_n\big)}{a_{ij} f_{\mathcal{P}}(\lambda_i, \lambda_j)} - 1 \right| = o_p(1). \tag{S6}$$

**Theorem S2.** *There exists $\alpha$ sufficiently small, depending on the walk length $k$, such that if $\rho_n = n^{-\alpha}$ then there exists a symmetric measurable (with respect to the sigma field generated by $W$) function $f_{\mathcal{P}}(\lambda, \lambda')$ which is bounded below away from zero, and bounded above by some constant $C < \infty$, such that*

$$\max_{i \neq j} \left| \frac{n^2 \mathbb{P}\big((i, j) \in \mathcal{P}(\mathcal{G}_n) \,|\, \mathcal{G}_n\big)}{f_{\mathcal{N}}(\lambda_i, \lambda_j)} - 1 \right| = o_p(1). \tag{S7}$$

19

The proof of these two results are given in Appendix E.1.1 and E.1.2 respectively. We note that while in principle we could give a closed formula for $f_{\mathcal{P}}$ and $f_{\mathcal{N}}$ in this scenario, they are sufficiently intractable to inspection that doing so would not provide any benefit.

In the case of DeepWalk where $p = q = 1$, the calculations involved are tractable enough such that we can improve the sparsity constraints, give closed forms for the measurable functions discussed above, and also provide rates of convergence.

**Theorem S3.** *Denote*

$$f_{\mathcal{P}}(\lambda_i, \lambda_j) := \frac{2k}{\rho_n \mathcal{E}_W(1)}, \tag{S8}$$

$$f_{\mathcal{N}}(\lambda_i, \lambda_j) := \frac{l(k+1)}{\mathcal{E}_W(1)\mathcal{E}_W(\alpha)} \big(W(\lambda_i, \cdot)W(\lambda_j, \cdot)^\alpha + W(\lambda_i, \cdot)^\alpha W(\lambda_j, \cdot)\big). \tag{S9}$$

*Then we have that*

$$\max_{i \neq j} \left| \frac{n^2 \mathbb{P}\big((i, j) \in \mathcal{P}(\mathcal{G}_n) \,|\, \mathcal{G}_n\big)}{a_{ij} f_{\mathcal{P}}(\lambda_i, \lambda_j)} - 1 \right| = O_p\Big(\Big(\frac{\log n}{n\rho_n}\Big)^{1/2}\Big), \tag{S10}$$

$$\max_{i \neq j} \left| \frac{n^2 \mathbb{P}\big((i, j) \in \mathcal{N}(\mathcal{G}_n) \,|\, \mathcal{G}_n\big)}{f_{\mathcal{N}}(\lambda_i, \lambda_j)} - 1 \right| = O_p\Big(\Big(\frac{\log n}{n\rho_n}\Big)^{1/2}\Big). \tag{S11}$$

*Proof.* This is a consequence of [11, Proposition 26]. We highlight the referenced result supposes that for the negative sampling scheme, vertices for which $a_{ij} = 0$ are rejected, whereas this does not happen here. Other than for the factor of $(1 - a_{ij})$ in the quoted result, the proof is otherwise unchanged, which gives the statement above for $\mathbb{P}\big((i, j) \in \mathcal{N}(\mathcal{G}_n) \,|\, \mathcal{G}_n\big)$. $\qquad\square$

With this, we then get the following result:

**Proposition S4.** *Denote*

$$\widehat{\mathcal{L}}_n(U, V) := \frac{1}{n^2} \sum_{i \neq j} \Big\{ -f_{\mathcal{P}}(\lambda_i, \lambda_j) a_{ij} \log(\sigma(\langle u_i, v_j \rangle)) - f_{\mathcal{N}}(\lambda_i, \lambda_j) \log(1 - \sigma(\langle u_i, v_j \rangle)) \Big\} \tag{S12}$$

*and define the set*

$$\Psi_A := \Big\{ U, V \in \mathbb{R}^{n \times d} \,|\, \mathcal{L}_n(U, V) \leq A\mathcal{L}_n(0_{n \times d}, 0_{n \times d}) \Big\} \subseteq \mathbb{R}^{n \times d} \times \mathbb{R}^{n \times d} \tag{S13}$$

*for any constant $A > 1$, where $0_{n \times d}$ denotes the zero matrix in $\mathbb{R}^{n \times d}$. Then for any set $X \subseteq \mathbb{R}^{n \times d} \times \mathbb{R}^{n \times d}$ containing the pair of zero matrices $O_{n \times d}$, we have under Scenario i) and iii) that*

$$\sup_{(U,V) \in \Psi_A \cap X} \big| \mathcal{L}_n(U, V) - \widehat{\mathcal{L}}_n(U, V) \big| = O_p\Big(A \cdot \Big(\frac{\log n}{n\rho_n}\Big)^{1/2}\Big), \tag{S14}$$

$$\mathbb{P}\Big( \underset{(U,V) \in X}{\arg\min}\, \mathcal{L}_n(U, V) \cup \underset{(U,V) \in X}{\arg\min}\, \widehat{\mathcal{L}}_n(U, V) \subseteq \Psi_A \cap X \Big) = 1 - o(1). \tag{S15}$$

*In Scenario (ii), the $O_p(\cdot)$ bound is replaced by an $o_p(1)$ bound.*

*Proof.* The proof is essentially equivalent to Lemma 32 of Davison and Austern [11] up to changes in notation, and so we do not repeat the details. $\qquad\square$

Note that in practice we can choose $A$ to be any constant greater than 1 but fixed with $n$ - e.g $A = 10$, and have the result hold. We will do so going forward.

## C.3 Averaging over the adjacency matrix of the graph

Following the proof outline, the next step is to argue that $\mathcal{L}_n(U, V)$ is close to its expectation when we average over the adjacency matrix of the graph $\mathcal{G}_n$. We begin with showing what occurs in the

DeepWalk case (Scenarios (i) and (iii)), and at the end of the section we discuss how the proof changes for the more general node2vec case. Note that we have

$$\mathbb{E}[\widehat{\mathcal{L}}_n(U,V)\,|\,\lambda] = \frac{1}{n^2}\sum_{i\neq j}\left\{-f_{\mathcal{P}}(\lambda_i,\lambda_j)\rho_n W(\lambda_u,\lambda_v)\log(\sigma(\langle u_i,v_j\rangle))-f_{\mathcal{N}}(\lambda_i,\lambda_j)\log(1-\sigma(\langle u_i,v_j\rangle))\right\} \tag{S16}$$

and so

$$E_n(U,V) := \frac{\mathcal{E}_W(1)}{2k}\left(\widehat{\mathcal{L}}_n(U,V) - \mathbb{E}[\widehat{\mathcal{L}}_n(U,V)\,|\,\lambda]\right) \tag{S17}$$

$$= \frac{1}{n^2}\sum_{i\neq j}\left(\rho_n^{-1}a_{ij} - W(\lambda_i,\lambda_j)\right)\cdot(-\log\sigma(\langle u_i,v_j\rangle)). \tag{S18}$$

Note that $\mathbb{E}[E_n(U,V)\,|\,\lambda] = 0$, and so it therefore suffices to control $E_n(U,V) - \mathbb{E}[E_n(U,V)\,|\,\lambda]$ uniformly over embedding matrices $U,V \in \mathbb{R}^{n\times d}$. This is the contents of the next theorem.

**Theorem S5.** *Begin by defining the set*

$$B_{2,\infty}(A_{2,\infty}) := \left\{U \in \mathbb{R}^{n\times d} \,:\, \|U\|_{2,\infty} \leq A_{2,\infty}\right\}. \tag{S19}$$

*Then we have the bound*

$$\sup_{U,V\in B_{2,\infty}(A_{2,\infty})}\left|E_n(U,V)\right| = O_p\left(A_{2,\infty}^2\left(\frac{d}{n\rho_n}\right)^{1/2}\right). \tag{S20}$$

*In particular, we also have that*

$$\sup_{U,V\in B_{2,\infty}(A_{2,\infty})}\left|\widehat{\mathcal{L}}_n(U,V) - \mathbb{E}[\widehat{\mathcal{L}}_n(U,V)\,|\,\lambda]\right| = O_p\left(\frac{A_{2,\infty}^2 k}{\mathcal{E}_W(1)}\left(\frac{d}{n\rho_n}\right)^{1/2}\right). \tag{S21}$$

*Proof.* Begin by noting that for any set $C \subseteq \mathbb{R}^{n\times d}\times\mathbb{R}^{n\times d}$ for which $0_{n\times d}\times 0_{n\times d}\in C$, we have that

$$\sup_{(U,V)\in C}|E_n(U,V)| \leq \sup_{(U,V)\in C}\left|E_n(U,V) - E_n(0_{n\times d},0_{n\times d})\right| + |E_n(0_{n\times d},0_{n\times d})| \tag{S22}$$

$$\leq \sup_{(U,V),(\tilde{U},\tilde{V})\in C}\left|E_n(U,V) - E_n(\tilde{U},\tilde{V})\right| + |E_n(0_{n\times d},0_{n\times d})|. \tag{S23}$$

We therefore need to control these two terms. We begin with the second; note that as

$$E_n(0_{n\times d},0_{n\times d}) = \frac{1}{n^2}\sum_{i\neq j}\left(\rho_n^{-1}a_{ij} - W(\lambda_i,\lambda_j)\right)\cdot\frac{1}{n^2} \tag{S24}$$

it follows by Lemma S30 that this term is $O_p((n^2\rho_n)^{-1/2})$. For the first term, we make use of a chaining bound. Note that if we write $T_{ij} = -\log\sigma(\langle u_i,v_j\rangle)$ and $S_{ij} = -\log\sigma(\langle\tilde{u}_i,\tilde{v}_j\rangle)$ for $i,j\in[n]$, then we have that

$$E_n(U,V) - E_n(\tilde{U},\tilde{V}) = \frac{1}{n^2}\sum_{i\neq j}\left(\rho_n^{-1}a_{ij} - W(\lambda_i,\lambda_j)\right)\cdot(T_{ij} - S_{ij}). \tag{S25}$$

Because the function $x\mapsto -\log\sigma(x)$ is 1-Lipschitz, it follows that

$$\|T-S\|_F^2 \leq \|UV^T - \tilde{U}\tilde{V}^T\|_F^2, \qquad \|T-S\|_\infty \leq \|UV^T - \tilde{U}\tilde{V}^T\|_\infty \tag{S26}$$

and consequently we have that

$$\mathbb{P}\left(|E_n(U,V) - E_n(\tilde{U},\tilde{V})| \geq u\right) \tag{S27}$$

$$\leq 2\exp\left(-\min\left\{\frac{u^2}{128\rho_n^{-1}n^{-4}\|UV^T - \tilde{U}\tilde{V}^T\|_F^2},\frac{u}{16\rho_n^{-1}n^{-2}\|UV^T - \tilde{U}\tilde{V}^T\|_\infty}\right\}\right) \tag{S28}$$

21

as a result of Lemma S30. Now, as $U, V \in B_F(A_F) \cap B_{2,\infty}(A_{2,\infty})$, by Lemma S19 if we define the metrics

$$d_F((U_1, V_1), (U_2, V_2)) := \|U_1 - U_2\|_F + \|V_1 - V_2\|_F, \tag{S29}$$

$$d_{2,\infty}((U_1, V_1), (U_2, V_2)) := \|U_1 - U_2\|_{2,\infty} + \|V_1 - V_2\|_{2,\infty}, \tag{S30}$$

then we have that

$$\mathbb{P}\big(|E_n(U, V) - E_n(\tilde{U}, \tilde{V})| \geq u\big) \tag{S31}$$

$$\leq 2 \exp \Big( - \min \Big\{ \frac{u^2}{128 \rho_n^{-1} n^{-4} A_F^2 d_F((U, V), (\tilde{U}, \tilde{V}))^2}, \frac{u}{16 \rho_n^{-1} n^{-2} A_{2,\infty} d_{2,\infty}((U, V), (\tilde{U}, \tilde{V}))} \Big\}\Big). \tag{S32}$$

As a result of Corollary S22, it therefore follows that

$$\sup_{(U,V),(\tilde{U},\tilde{V}) \in T \times T} \big|E_n(U, V) - E_n(\tilde{U}, \tilde{V})\big| = O_p\Big(A_{2,\infty}^2 \Big(\frac{d}{n\rho_n}\Big)^{1/2} + A_{2,\infty}^2 \frac{d}{n\rho_n}\Big) \tag{S33}$$

The desired conclusion follows by combining the bounds (S24) and (S33). $\qquad \square$

For the more abstract node2vec case under Scenario (ii), we highlight that we can take

$$E_n(U, V) = \frac{1}{n^2} \sum_{i \neq j} \rho_n f_{\mathcal{P}}(\lambda_i, \lambda_j)\Big(\rho_n^{-1} a_{ij} - W(\lambda_i, \lambda_j)\Big) \cdot (- \log \sigma(\langle u_i, v_j \rangle)). \tag{S34}$$

Now, as $f_{\mathcal{P}}(\lambda_u, \lambda_v)$ is a function of the community assignments only within the SBM case, we can replace this by a matrix of constants $f_{\mathcal{P},c,c'}$ for $c, c' \in [\kappa]$, and therefore the error term can be decomposed into a sum

$$\sum_{c_1, c_2} (\rho_n f_{\mathcal{P}, c_1, c_2}) \sum_{\substack{i \neq j \\ i: c(u) = c_1 \\ j: c(u) = c_2}} \Big( \rho_n^{-1} a_{ij} - W(\lambda_i, \lambda_j)\Big) \cdot (- \log \sigma(\langle u_i, v_j \rangle)), \tag{S35}$$

where we recall that $\max_{c_1, c_2}(\rho_n f_{\mathcal{P}, c_1, c_2}) < \infty$ as guaranteed by Theorem S1. Each of these terms (of which there are finitely many) can be controlled using the exact same argument as in Theorem S5, and so the conclusion of the Theorem also holds with the same overall rate of convergence in Scenario (ii).

## C.4 Adding in a diagonal term

Currently the sum in $\mathbb{E}[\widehat{\mathcal{L}}_n(U, V) \mid \lambda]$ is defined only terms $i, j$ with $i \neq j$ - it is more convenient to work with the version where the diagonal term is added in:

$$\mathcal{R}_n(U, V) := \frac{1}{n^2} \sum_{i, j \in [n]} \Big\{ - f_{\mathcal{P}}(\lambda_i, \lambda_j) \rho_n W(\lambda_u, \lambda_v) \log(\sigma(\langle u_i, v_j \rangle)) \tag{S36}$$

$$- f_{\mathcal{N}}(\lambda_i, \lambda_j) \log(1 - \sigma(\langle u_i, v_j \rangle)) \Big\}. \tag{S37}$$

We show that this does not significantly change the size of the loss function.

**Lemma S6.** *With the same notation as in Theorem S5, we have that*

$$\sup_{U, V \in B_{2,\infty}(A_{2,\infty})} \big|\mathcal{R}_n(U, V) - \mathbb{E}[\widehat{\mathcal{L}}_n(U, V) \mid \lambda]\big|$$

$$= O_p\Big(\frac{1}{n} A_{2,\infty}^2 \Big(\|\rho_n f_{\mathcal{P}}(\lambda, \lambda') W(\lambda, \lambda')\|_\infty + \|f_{\mathcal{N}}(\lambda, \lambda')\|_\infty\Big)\Big).$$

*In particular, in the case of DeepWalk we have that*

$$\sup_{U, V \in B_{2,\infty}(A_{2,\infty})} \big|\mathcal{R}_n(U, V) - \mathbb{E}[\widehat{\mathcal{L}}_n(U, V) \mid \lambda]\big| = O_p\Big(\frac{1}{n} A_{2,\infty}^2 \Big(\frac{2k\|W\|_\infty}{\mathcal{E}_W(1)} + \frac{2l(k+1)\|W\|_\infty^2}{\mathcal{E}_W(1)\mathcal{E}_W(\alpha)}\Big)\Big).$$

22

*Proof.* Begin by noting that

$$0 \leq \mathcal{R}_n(U, V) - \mathbb{E}[\widehat{\mathcal{L}}_n(U, V) \mid \lambda]$$

$$= \frac{1}{n^2} \sum_{i=1}^{n} \left\{ -f_{\mathcal{P}}(\lambda_i, \lambda_j) \rho_n W(\lambda_u, \lambda_v) \log(\sigma(\langle u_i, v_i \rangle)) - f_{\mathcal{N}}(\lambda_i, \lambda_j) \log(1 - \sigma(\langle u_i, v_i \rangle)) \right\}.$$

(S38)

Note that we can bound

$$-\log(\sigma(\langle u_i, v_j \rangle)) \leq |\langle u_i, v_i \rangle \leq \|u_i\|_2 \|v_i\|_2$$

(S39)

and similarly $-\log(1 - \sigma(\langle u_i, v_i \rangle)) \leq |\langle u_i, v_i \rangle| \leq \|u_i\|_2 \|v_i\|_2$. Moreover, we have the bounds

$$f_{\mathcal{P}}(\lambda_i, \lambda_j) \rho_n W(\lambda_i, \lambda_j) \leq \|\rho_n f_{\mathcal{P}}(\lambda, \lambda') W(\lambda, \lambda')\|_\infty < \infty, \quad f_{\mathcal{N}}(\lambda_i, \lambda_j) \leq \|f_{\mathcal{N}}(\lambda, \lambda')\|_\infty < \infty$$

(S40)

under our assumptions. As a result, because $U, V \in \mathcal{B}_{2,\infty}(A_{2,\infty})$, we end up with the final bound

$$\left| \mathcal{R}_n(U, V) - \mathbb{E}[\widehat{\mathcal{L}}_n(U, V) \mid \lambda] \right| \leq \frac{1}{n} A_{2,\infty}^2 \left( \|\rho_n f_{\mathcal{P}}(\lambda, \lambda') W(\lambda, \lambda')\|_\infty + \|f_{\mathcal{N}}(\lambda, \lambda')\|_\infty \right)$$ (S41)

which gives the stated result as the RHS is free of $U$ and $V$. $\qquad \square$

## C.5 Chaining up the loss function approximations

By chaining up the prior results, we end up with the following result:

**Proposition S7.** *There exists a non-empty set $\Psi_n$ for each $n$ such that, for any set $X \subseteq \mathbb{R}^{n \times d} \times \mathbb{R}^{n \times d}$ containing $0_{n \times d} \times 0_{n \times d}$, we have for DeepWalk that*

$$\sup_{(U,V) \in \Psi_n \cap B_{2,\infty}(A_{2,\infty})} \left| \mathcal{L}_n(U, V) - \mathcal{R}_n(U, V) \right| = O_p\left( \left( \frac{\log n}{n \rho_n} \right)^{1/2} + A_{2,\infty}^2 \left( \frac{d}{n \rho_n} \right)^{1/2} \right) \quad \text{(S42)}$$

*and*

$$\mathbb{P}\left( \mathop{\arg\min}_{(U,V) \in B_{2,\infty}(A_{2,\infty}) \cap X} \mathcal{L}_n(U, V) \, cup \mathop{\arg\min}_{(U,V) \in B_{2,\infty}(A_{2,\infty}) \cap X} \mathcal{R}_n(U, V) \subseteq \Psi_A \cap B_{2,\infty}(A_{2,\infty}) \cap X \right)$$
$$= 1 - o(1).$$

(S43)

*For node2vec, the same result holds when we replace the $(\log n / n \rho_n)^{1/2}$ term with an $o_p(1)$ term and add the constraint that $d \ll n \rho_n$. The same result also holds when we constrain $U = V$, but otherwise keep everything else unchanged.*

## C.6 Minimizers of $\mathcal{R}_n(U, V)$

Recall that we have earlier defined

$$\mathcal{R}_n(U, V) := \frac{1}{n^2} \sum_{i,j \in [n]} \left\{ -f_{\mathcal{P}}(\lambda_i, \lambda_j) \rho_n W(\lambda_u, \lambda_v) \log(\sigma(\langle u_i, v_j \rangle)) \right.$$

(S44)

$$\left. -f_{\mathcal{N}}(\lambda_i, \lambda_j) \log(1 - \sigma(\langle u_i, v_j \rangle)) \right\}.$$

We now want to reason about the minima of these functions. To do so, note that the optimization domain is non-convex - firstly due to the rank constraints on the matrix $UV^T$, and secondly due to the fact that the loss function is invariant to any mapping $(U, V) \to (UM, VM^{-1})$ for any invertible $d \times d$ matrix $M$. To handle the second part, we consider the global minima of this function when parameterized only in term of the matrix $UV^T$. We will then see that the minima matrix is already low rank.

We first begin by giving some basic facts about the function $\mathcal{R}_n(U, V)$ when parameterized as a function of $UV^T$.

**Lemma S8.** *Define the modified function*

$$\mathcal{R}_n(M) := \frac{1}{n^2} \sum_{i,j \in [n]} \left\{ -f_{\mathcal{P}}(\lambda_i, \lambda_j) \rho_n W(\lambda_u, \lambda_v) \log(\sigma(M_{ij})) - f_{\mathcal{N}}(\lambda_i, \lambda_j) \log(1 - \sigma(M_{ij})) \right\}.$$

(S45)

*over all matrices $M \in \mathbb{R}^{n \times n}$. Then we have the following:*

23

a) *The function $\mathcal{R}_n(M)$ is strictly convex in $M$.*

b) *The global minimizer of $\mathcal{R}_n(M)$ is given by*

$$M_{ij}^* = \log\left(\frac{f_{\mathcal{P}}(\lambda_i, \lambda_j)\rho_n W(\lambda_i, \lambda_j)}{f_{\mathcal{N}}(\lambda_i, \lambda_j)}\right) \tag{S46}$$

*and satisfies $\nabla_M \mathcal{R}_n(M) = 0$.*

c) *When restricted to a cone of semi-positive definite matrices $M \in \mathcal{M}_n^{\succeq 0}$, there exists a unique minimizer to $\mathcal{R}_n(M)$ over this set, which we call $M^{\succeq 0}$. Moreover, $M^{\succeq 0}$ has the property that $\langle \nabla_M \mathcal{R}_n(M^{\succeq 0}), M^{\succeq 0} - M \rangle \leq 0$ for all $M \in \mathcal{M}_n^{\succeq 0}$.*

*Proof.* For part a), this follows by the fact that the functions $-\log(\sigma(x))$ and $-\log(1 - \sigma(x))$ are positive and strictly convex functions of $x \in \mathbb{R}$, the fact that $f_{\mathcal{P}}(\lambda_i, \lambda_j)\rho_n W(\lambda_i, \lambda_j)$ and $f_{\mathcal{N}}(\lambda_i, \lambda_j)$ are positive quantities which are bounded above (see e.g Lemma S6), and the fact that the sum of strictly convex functions is strictly convex. For part b), this follows by noting that each of the $M_{ij}^*$ are pointwise minima of the functions

$$r_{ij}(x) = -f_{\mathcal{P}}(\lambda_i, \lambda_j)\rho_n W(\lambda_u, \lambda_v)\log(\sigma(x))) - f_{\mathcal{N}}(\lambda_i, \lambda_j)\log(1 - \sigma(x)) \tag{S47}$$

defined over $x \in \mathbb{R}$. Indeed, note that

$$\frac{dr_{ij}}{dx} = (-1 + \sigma(x))f_{\mathcal{P}}(\lambda_i, \lambda_j)\rho_n W(\lambda_u, \lambda_v) + \sigma(x)f_{\mathcal{N}}(\lambda_i, \lambda_j), \tag{S48}$$

so setting this equal to zero, rearranging and making use of the equality $\sigma^{-1}(a/(a+b)) = \log(a/b)$ gives the stated result. Part c) is a consequence of strong convexity, the optimization domain being convex and self dual, and the KKT conditions. $\qquad\square$

To understand the form of the the global minimizer of $\mathcal{R}_n(M)$ in the DeepWalk case, by substituting in the values for $f_{\mathcal{P}}(\lambda_i, \lambda_j)$ and $f_{\mathcal{N}}(\lambda_i, \lambda_j)$ we end up with

$$M_{ij}^* = \log\left(\frac{2P_{c(i),c(j)}\mathcal{E}_W(\alpha)}{(1+k^{-1})\mathbb{E}[\theta]\mathbb{E}[\theta]^\alpha\left(\theta_j^{\alpha-1}\widetilde{P}_{c(i)}\widetilde{P}_{c(j)}^\alpha + \theta_i^{\alpha-1}\widetilde{P}_{c(i)}^\alpha\widetilde{P}_{c(j)}\right)}\right) \tag{S49}$$

$$= \log\left(\frac{2\mathcal{E}_W(\alpha)}{(1+k^{-1})\mathbb{E}[\theta]\mathbb{E}[\theta]^\alpha} \cdot \frac{P_{c(i),c(j)}}{\widetilde{P}_{c(i)}\widetilde{P}_{c(j)} \cdot \left(\theta_i^{\alpha-1}\widetilde{P}_{c(i)}^{\alpha-1} + \theta_j^{\alpha-1}\widetilde{P}_{c(j)}^{\alpha-1}\right)}\right) \tag{S50}$$

In particular, from the above formula we get the following lemma as a consequence:

**Lemma S9.** *Suppose that Scenarios (i) or (iii) holds, so that either a) $\theta_i$ is constant for all $i$, or b) $\alpha = 1$. Then if we write $\Pi_C \in \mathbb{R}^{n \times \kappa}$ for the matrix where $(\Pi_C)_{il} = 1[c(i) = l]$, and define the matrix*

$$(\widetilde{M}_\alpha^*)_{lm} = \log\left(\frac{2\mathcal{E}_W(\alpha)}{(1+k^{-1})\mathbb{E}[\theta]\mathbb{E}[\theta]^\alpha} \cdot \frac{P_{lm}}{\widetilde{P}_m\widetilde{P}_l^\alpha + \widetilde{P}_m^\alpha\widetilde{P}_l}\right) \text{ for } l, m \in [\kappa], \tag{S51}$$

*then we have that $M^* = \Pi_C \widetilde{M}_\alpha^* \Pi_C^T$. In particular, as soon as the matrix $\Pi_C$ is of full rank (which occurs with asymptotic probability 1), then the rank of $M^*$ equals the rank of $\widetilde{M}_\alpha^*$. Moreover, as soon as $d$ is greater than or equal to the rank of $\widetilde{M}_\alpha^*$, $(U, V)$ is a minimizer of $\mathcal{R}_n(U, V)$ if and only if $UV^T = M^*$.*

*Under Scenario (ii), the same result applies noting that $f_{\mathcal{P}}$ and $f_{\mathcal{N}}$ are functions only of the underling communities, and so if we abuse notation and write e.g $f_{\mathcal{P}}(l, m)$ to indicate the value of $f_{\mathcal{P}}(\lambda_i, \lambda_j)$ when $c(i) = l$ and $c(j) = m$, one can take*

$$(\widetilde{M}^*)_{lm} = \log\left(\frac{f_{\mathcal{P}}(l, m)\rho_n P_{l,m}}{f_{\mathcal{N}}(l, m)}\right) \tag{S52}$$

*and have the above result hold.*

We discuss in Appendix F what happens when we apply DeepWalk in the DCSBM regime when $\alpha \neq 1$. To give an example of what $M^*$ looks like, we write it down in the case of a $\mathrm{SBM}(n, \kappa, \tilde{p}, \tilde{q}, \rho_n)$ model, which is frequently used to illustrate the behavior of various community detection algorithms. Such a model assumes that the community assignments $\pi_l = 1/\kappa$ for all $l \in [\kappa]$, and that

$$P_{kl} = \begin{cases} \tilde{p} & \text{if } k = l, \\ \tilde{q} & \text{if } k \neq l. \end{cases} \tag{S53}$$

In this case, we have that

$$\widetilde{P}_l = \frac{\tilde{p} + \kappa(\tilde{q} - 1)}{\kappa} \text{ for } l \in [\kappa], \qquad \mathcal{E}_W(\alpha) = \mathbb{E}[\theta]^\alpha \mathbb{E}[\theta^\alpha] \cdot \left(\frac{\tilde{p} + (\kappa - 1)\tilde{q}}{\kappa}\right)^\alpha. \tag{S54}$$

Substituting these values into the matrix $\widetilde{M}^*_\alpha$ gives

$$(\widetilde{M}^*_\alpha)_{lm} = \log\left(\frac{\mathbb{E}[\theta^\alpha]}{\mathbb{E}[\theta](1 + k^{-1})} \cdot \frac{\kappa\tilde{p}}{\tilde{p} + (\kappa - 1)\tilde{q}}\right)\delta_{lm} + \log\left(\frac{\mathbb{E}[\theta^\alpha]}{\mathbb{E}[\theta](1 + k^{-1})} \cdot \frac{\kappa\tilde{q}}{\tilde{p} + (\kappa - 1)\tilde{q}}\right)(1 - \delta_{lm}). \tag{S55}$$

We highlight this is a matrix of the form $\alpha\delta_{lm} + \beta(1 - \delta_{lm})$, and so it is straightforward to describe the spectral behavior of the matrix (see Lemma S31).

### C.6.1  Minimizers in the constrained regime $U = V$

In the case where we have constrained $U = V$, it is not possible in general to write down the closed form of the minimizer of $\mathcal{R}_n(M)$ over $\mathcal{M}_n^{\succeq 0}$. However, it is still possible to draw enough conclusions about the form of the minimizer in order to give guarantees for community detection. We begin with the proposition below. We state the next two results for DeepWalk only, but note that the first generalizes to the node2vec case immediately.

**Proposition S10.** *Suppose that $\theta_i$ is constant across all $i$. Supposing that $\widetilde{M} \in \mathbb{R}^{\kappa \times \kappa}$ is of the form $\widetilde{M} = \widetilde{U}\widetilde{U}^T$ for matrices $\widetilde{U} \in \mathbb{R}^{\kappa \times d}$, define the function*

$$\widetilde{\mathcal{R}}_n(\widetilde{M}) = \sum_{l,m \in [\kappa]} \hat{p}_n(l)\hat{p}_n(m)\left\{-2kP_{lm}\log\sigma(\langle u_l, u_m\rangle) - \{\widetilde{P}_l\widetilde{P}_m^\alpha + \widetilde{P}_m\widetilde{P}_l^\alpha\}\log(1 - \sigma(\langle u_l, u_m\rangle))\right\} \tag{S56}$$

*where we define $\hat{p}_n(l) := n^{-1}|\{i : c(i) = l\}|$ for $l \in [\kappa]$. Then $\widetilde{\mathcal{R}}_n(\widetilde{M})$ is strongly convex, and moreover has a unique minimizer as soon as $d \geq \kappa$.*

*Moreover, any minimizer of $\mathcal{R}_n(M)$ over matrices $M$ of the form $M = UU^T$ where $U \in \mathbb{R}^{n \times d}$ must take the form $M = \Pi_C M^* \Pi_C^T$ where $(\Pi_C)_{il} = 1[c(i) = l]$ where $M^*$ is a minimizer of $\widetilde{\mathcal{R}}_n(\widetilde{M})$. In particular, once $d \geq \kappa$, there is a unique minimizer to $\mathcal{R}_n(M)$.*

*Proof.* The properties of $\widetilde{\mathcal{R}}_n(\widetilde{M})$ are immediate by similar arguments to Lemma S8 and standard facts in convex analysis. We begin by noting that if we substitute in the values

$$\rho_n W(\lambda_i, \lambda_j) f_\mathcal{P}(\lambda_i, \lambda_j) = \frac{2k P_{c(i),c(j)}}{\mathcal{E}_W(1)}, \tag{S57}$$

$$f_\mathcal{N}(\lambda_i, \lambda_j) = \frac{l(k + 1)}{\mathcal{E}_W(1)\mathcal{E}_W(\alpha)}\left(\widetilde{P}_{c(i)}\widetilde{P}_{c(j)}^\alpha + \widetilde{P}_{c(j)}\widetilde{P}_{c(i)}^\alpha\right), \tag{S58}$$

for $f_\mathcal{P}(\lambda_i, \lambda_j)$ and $f_\mathcal{N}(\lambda_i, \lambda_j)$, then we can write that (recalling that $M_{ij} = \langle u_i, u_j\rangle$)

$$\mathcal{R}_n(M) := \frac{1}{n^2}\sum_{i,j \in [n]}\left\{-2kP_{c(i),c(j)}\log\sigma(\langle u_i, u_j\rangle)\right. \tag{S59}$$

$$\left. - \frac{l(k + 1)}{\mathcal{E}_W(1)\mathcal{E}_W(\alpha)}\left(\widetilde{P}_{c(i)}\widetilde{P}_{c(j)}^\alpha + \widetilde{P}_{c(j)}\widetilde{P}_{c(i)}^\alpha\right)\log(1 - \sigma(\langle u_i, u_j\rangle))\right\} \tag{S60}$$

$$:= \sum_{l,m \in [\kappa]}\hat{p}_n(l)\hat{p}_n(m)\left\{-2kP_{lm}\frac{1}{|\mathcal{C}_l||\mathcal{C}_m|}\sum_{i \in \mathcal{C}_l, j \in \mathcal{C}_m}\log\sigma(\langle u_i, u_j\rangle)\right. \tag{S61}$$

$$\left. - \{\widetilde{P}_{c(i)}\widetilde{P}_{c(j)}^\alpha + \widetilde{P}_{c(j)}\widetilde{P}_{c(i)}^\alpha\}\frac{1}{|\mathcal{C}_l||\mathcal{C}_m|}\sum_{i \in \mathcal{C}_l, j \in \mathcal{C}_m}\log(1 - \sigma(\langle u_i, u_j\rangle))\right\} \tag{S62}$$

25

where for $l \in [\kappa]$ we define $\hat{p}_n(l) := n^{-1}|\{i : c(i) = l\}|$, along with the sets $\mathcal{C}_l = \{i : c(i) = l\}$. Now, note that as the functions $-\log(\sigma(x))$ and $-\log(1 - \sigma(x))$ are strictly convex, by Jensen's inequality we have that e.g

$$\frac{1}{|\mathcal{C}_l||\mathcal{C}_m|} \sum_{i \in \mathcal{C}_l, j \in \mathcal{C}_m} -\log \sigma(\langle u_i, u_j \rangle) \geq -\log \sigma\Big(\Big\langle \frac{1}{|\mathcal{C}_l|} \sum_{i \in \mathcal{C}_l} u_i, \frac{1}{|\mathcal{C}_m|} \sum_{j \in \mathcal{C}_m} u_j \Big\rangle\Big) \tag{S63}$$

(where we also used bilinearity of the inner product) where equality holds above if and only if the $u_i$ are constant are across all indices $i$. In particular, any minimizer of $\mathcal{R}_n(M)$ must have the $u_i$ constant across $i \in \mathcal{C}_l$ for each $l \in [\kappa]$, which defines the function $\tilde{\mathcal{R}}_n(\widetilde{M})$. This gives the claimed statement. $\qquad\square$

In certain cases, we are able to give a closed form to the minimizer. We illustrate this for the case of the $SBM(n, \kappa, \tilde{p}, \tilde{q}, \rho_n)$ model.

**Proposition S11.** *Let $\widetilde{M}^*$ be the unique minimizer of $\tilde{\mathcal{R}}_n(\widetilde{M})$ as introduced in Proposition S10. In the case of a $SBM(n, \kappa, \tilde{p}, \tilde{q}, \rho_n)$ model, we have that $\kappa^{-2}\|\widetilde{M}^* - M^*\|_1 = O_p((\kappa \log \kappa/n)^{1/4})$, where $M^*$ is of the form*

$$(M^*)_{ij} = \alpha^* \delta_{ij} - \frac{\alpha^*}{\kappa - 1}(1 - \delta_{ij}) \tag{S64}$$

*for some $\alpha^* = \alpha^*(\tilde{p}, \tilde{q}) \geq 0$. Moreover, $\alpha^* > 0$ iff $\tilde{p} > \tilde{q}$.*

*Proof.* We begin by arguing that the objective function $\widetilde{\mathcal{R}}_n(\widetilde{M})$ converges uniformly to the objective

$$\bar{\mathcal{R}}_n(\widetilde{M}) := \frac{1}{\kappa^2} \sum_{l,m \in [\kappa]} \Big\{ -2kP_{lm} \log \sigma(\langle u_l, u_m \rangle) - \{\widetilde{P}_m \widetilde{P}_l^\alpha + \widetilde{P}_l \widetilde{P}_m^\alpha\} \log(1 - \sigma(\langle u_l, u_m \rangle)) \Big\} \tag{S65}$$

over a set containing the minimizers of both functions. Note that this function is also strictly convex, and has a unique minimizer as soon as $d \geq \kappa$. To do so, we highlight that as we have that

$$\max_{k \neq l} \left| \frac{\hat{p}_n(l)\hat{p}_n(k) - \kappa^{-2}}{\kappa^{-2}} \right| = O_p\Big(\Big(\frac{\kappa \log \kappa}{n}\Big)^{1/2}\Big) \tag{S66}$$

by standard concentration results for Binomial random variables (e.g Proposition 47 of [11]), it follows that

$$\left| \bar{\mathcal{R}}_n(\widetilde{M}) - \widetilde{\mathcal{R}}_n(\widetilde{M}) \right| \leq \bar{\mathcal{R}}_n(\widetilde{M}) \cdot O_p\Big(\Big(\frac{\kappa \log \kappa}{n}\Big)^{1/2}\Big). \tag{S67}$$

Consequently, $\widetilde{\mathcal{R}}_n(\widetilde{M})$ converges to $\bar{\mathcal{R}}_n(\widetilde{M})$ uniformly over any level set of $\bar{\mathcal{R}}_n(\widetilde{M})$, which necessarily contains the minima of $\bar{\mathcal{R}}_n(\widetilde{M})$. If one does so over the set (for example)

$$A = \{\widetilde{M} : \bar{\mathcal{R}}_n(\widetilde{M}) \leq 10\bar{\mathcal{R}}_n(0)\} \tag{S68}$$

(for example), then as $\bar{\mathcal{R}}_n(0)$ is constant across $n$, we have uniform convergence of (S67) over the set $A$ at a rate of $O_p((\log \kappa/np)^{1/2})$. This argument can be reversed, which therefore ensures uniform convergence (over the same set) which contains the minimizers (with the minimizer of $\tilde{\mathcal{R}}_n(M)$ being contained within this set with asymptotic probability 1) at a rate of $O_p((\kappa \log \kappa/n)^{1/2})$.

With this, we note that an application of Lemma S33 gives that for any matrices $\widetilde{M}_1$ and $\widetilde{M}_2$ we have that

$$\bar{\mathcal{R}}_n(\widetilde{M}_1) \geq \bar{\mathcal{R}}_n(\widetilde{M}_2) + \langle \Delta \bar{\mathcal{R}}_n(\widetilde{M}_2), \widetilde{M}_1 - \widetilde{M}_2 \rangle \tag{S69}$$

$$+ \frac{C}{\kappa^2} \sum_{i,j \in [\kappa]} \min\{|(\widetilde{M}_2)_{ij} - (\widetilde{M}_1)_{ij}|^2, 2|(\widetilde{M}_2)_{ij} - (\widetilde{M}_1)_{ij}|\}. \tag{S70}$$

where to save on notation, we define

$$C := \frac{1}{4} e^{-\|\widetilde{M}_2\|_\infty} \min_{l,m}\{2kP_{lm}, \widetilde{P}_m \widetilde{P}_l^\alpha\}. \tag{S71}$$

26

In particular, if $\widetilde{M}_2 = \bar{M}^*$ is an optimum of $\bar{\mathcal{R}}_n(\widetilde{M})$, then by the KKT conditions (similarly as in Lemma S8) we have that

$$\bar{\mathcal{R}}_n(\widetilde{M}_1) - \bar{\mathcal{R}}_n(\bar{M}^*) \geq \frac{C}{\kappa^2} \sum_{i,j \in [\kappa]} \min\{|(\bar{M}^*)_{ij} - (\widetilde{M}_1)_{ij}|^2, 2|(\bar{M}^*)_{ij} - (\widetilde{M}_1)_{ij}|\}. \quad \text{(S72)}$$

In particular, if we then let $\widetilde{M}^*$ be any minimizer of $\widetilde{\mathcal{R}}_n(\widetilde{M})$, then we have that

$$\frac{C}{\kappa^2} \sum_{i,j \in [\kappa]} \min\{|(\bar{M}^*)_{ij} - (\widetilde{M}_1)_{ij}|^2, 2|(\bar{M}^*)_{ij} - (\widetilde{M}_1)_{ij}|\} \quad \text{(S73)}$$

$$\leq \bar{\mathcal{R}}_n(\widetilde{M}_1) - \bar{\mathcal{R}}_n(\bar{M}^*) \leq \bar{\mathcal{R}}_n(\widetilde{M}_1) - \tilde{\mathcal{R}}_n(\bar{M}^*) + \tilde{\mathcal{R}}_n(\widetilde{M}^*) - \bar{\mathcal{R}}_n(\bar{M}^*) \quad \text{(S74)}$$

$$\leq 2 \sup_{M \in A} \left| \tilde{\mathcal{R}}_n(M) - \bar{\mathcal{R}}_n(M) \right| \quad \text{(S75)}$$

on an event of asymptotic probability 1. Consequently, it follows by Lemma S34 that

$$\frac{1}{\kappa^2} \|\bar{M}^* - \widetilde{M}^*\|_1 = O_p\big((\kappa \log \kappa/n)^{1/4}\big). \quad \text{(S76)}$$

We now need to find the minimizing positive semi-definite matrix which optimizes $\bar{\mathcal{R}}_n(\widetilde{M})$. To do so, we will argue that one can find $\alpha$ for which

$$\widehat{M}_{ij} = \alpha \delta_{ij} - \frac{\alpha}{\kappa - 1}(1 - \delta_{ij}), \quad \nabla \bar{\mathcal{R}}_n(\widehat{M}) = C 1_\kappa 1_\kappa^T, \quad 1_\kappa = (1, \cdots, 1)^T$$

for some positive constant $C$, as then the KKT conditions for the constrained optimization problem will hold. Indeed, for any positive definite matrix $M$, as by definition of $\widehat{M}$ we have that $\langle \nabla \bar{\mathcal{R}}_n(\widehat{M}), \widehat{M} \rangle = 0$ as all of the eigenvectors of $\widehat{M}$ are orthogonal to the unit vector $1_\kappa$ (Lemma S31). It consequently follows that as $\nabla \bar{\mathcal{R}}_n(\widehat{M})$ is itself positive definite, we get that $\langle -\nabla \bar{\mathcal{R}}_n(\widehat{M}), \widehat{M} - M \rangle = \langle \nabla \bar{\mathcal{R}}_n(\widehat{M}), M \rangle \geq 0$. We now need to verify the existence of a constant $\alpha$ for which this condition holds. We note that as $\widehat{M}_{ij}$ is constant across $i = j$, and also constant across $i \neq j$, to verify the condition that $\nabla \bar{\mathcal{R}}_n(\widehat{M})$ is proportional to $1_\kappa 1_\kappa^T$, it suffices to check whether the on and off diagonal terms of $\nabla \bar{\mathcal{R}}_n(\widehat{M})$ are equal to each other. This gives the equation

$$\sigma(\alpha) \cdot \left( k\tilde{p} + l(k+1)\frac{\tilde{p} + (\kappa - 1)\tilde{q}}{\kappa} \right)$$

$$= k(\tilde{p} - \tilde{q}) + \sigma(-\alpha/(\kappa - 1))\left( k\tilde{q} + l(k+1)\frac{\tilde{p} + (\kappa - 1)\tilde{q}}{\kappa} \right)$$

By applying Lemma S32, this has a singular positive solution in $\alpha$ if and only if $k(\tilde{p} - \tilde{q}) \geq k(\tilde{p} - \tilde{q})/2$, which holds iff $\tilde{p} \geq \tilde{q}$. In the case where $\tilde{p} < \tilde{q}$, it follows that the solution has $\alpha = 0$. $\qquad \square$

## C.7 Strong convexity properties of the minima matrix

**Proposition S12.** *Define the modified function*

$$\mathcal{R}_n(M) := \frac{1}{n^2} \sum_{i,j \in [n]} \Big\{ -f_{\mathcal{P}}(\lambda_i, \lambda_j)\rho_n W(\lambda_u, \lambda_v) \log(\sigma(M_{ij})) - f_{\mathcal{N}}(\lambda_i, \lambda_j) \log(1 - \sigma(M_{ij})) \Big\}.$$

$$\text{(S77)}$$

*over all matrices $M \in \mathbb{R}^{n \times n}$. Then we have for any matrices $M_1, M_2 \in \mathbb{R}^{n \times n}$ with $\|M_1\|_\infty, \|M_2\|_\infty \leq A_\infty$ that*

$$\mathcal{R}_n(M_1) \geq \mathcal{R}_n(M_2) + \langle \nabla \mathcal{R}_n(M_2), M_1 - M_2 \rangle + \frac{\widetilde{C}e^{-A_\infty}}{2} \cdot \frac{1}{n^2}\|M_1 - M_2\|_F^2 \quad \text{(S78)}$$

*where $\widetilde{C} = \min_{l,m}\{2kP_{l,m}, \tilde{P}_l^\alpha \tilde{P}_m\}$ for Scenarios (i) and (iii), and $\widetilde{C} = \min\{\|\rho_n f_{\mathcal{P}}(\lambda, \lambda')\|_{-\infty}, \|f_{\mathcal{N}}(\lambda, \lambda')\|_{-\infty}\} > 0$ for Scenario (ii). Moreover,*

    *i) If $\mathcal{R}_n(M)$ is constrained over a set $\mathcal{X} = \{M = UV^T : U, V \in \mathbb{R}^{n \times d}, \|M\|_\infty \leq A_\infty\}$, and there exists $M^*$ in $\mathcal{X}$ such that $\nabla \mathcal{R}_n(M^*) = 0$, then we have that*

$$\frac{1}{n^2}\|M^* - M\|_F^2 \leq 2\widetilde{C}^{-1}e^{A_\infty} \cdot \big(\mathcal{R}_n(M) - \mathcal{R}_n(M^*)\big) \text{ for all } M \in \mathcal{X}. \quad \text{(S79)}$$

*ii)* *If $\mathcal{R}_n(M)$ is constrained over a set $\mathcal{X}^{\geq 0} = \{M = UU^T \ : \ U \in \mathbb{R}^{n \times d}, \|M\|_\infty \leq A_\infty \}$, and there exists $M^*$ in $\mathcal{X}^{\geq 0}$ such that $\langle \nabla \mathcal{R}_n(M^*), M - M^* \rangle \geq 0$ for all $M \in \mathcal{X}^{\geq 0}$, then we get the same inequality as in part i) above.*

*Proof.* The first inequality follows by an application of Lemma S33, with the second and third parts following by applying the conditions stated and rearranging. $\qquad\square$

## C.8 Convergence of the gram matrices of the embeddings

By combining together Proposition S12 and Proposition S7 we end up with the following result:

**Theorem S13.** *Suppose that the conditions of Lemma S9 hold. (In particular, recall that $d \geq \kappa$.) Then there exist constants $A_\infty$ and $A_{2,\infty}$ (depending on the parameters of the model and sampling scheme) and a matrix $M^* \in \mathbb{R}^{\kappa \times \kappa}$ (also depending on the parameters of the model and the sampling scheme) such that for any minimizer $(U^*, V^*)$ of $\mathcal{L}(U,V)$ over the set*

$$X = \{(U,V) \ : \ \|U\|_\infty, \|V\|_\infty \leq A_\infty, \|U\|_{2,\infty}, \|V\|_{2,\infty} \leq A_{2,\infty}\}, \tag{S80}$$

*we have that*

$$\frac{1}{n^2} \sum_{i,j\in[n]} \left( \langle u_i^*, v_j^* \rangle - M_{c(i),c(j)}^* \right)^2 = C \cdot \begin{cases} O_p(( \frac{\max\{\log n, d\}}{n\rho_n})^{1/2}) & \text{under Scenarios (i) and (iii);} \\ o_p(1) & \text{under Scenario (ii);} \end{cases} \tag{S81}$$

*for some constant $C$ depending on the model, the node2vec hyperparameters, $A_\infty$ and $A_{2,\infty}$. In the case where we constrain $U = V$, the same result holds provided the conditions of Proposition S10 hold.*

*Proof.* We note that by Lemma S9, there exists a minimizer $\widetilde{M}^*$ for $\mathcal{R}_n(M)$ of the form $\widetilde{M}^* = \Pi M^* \Pi^T$ for a matrix $M^* \in \mathbb{R}^{\kappa \times \kappa}$. We can then take $A_\infty$ and $A_{2,\infty}$ as $2\|M^*\|_\infty$ and $2\|M^*\|_{2,\infty}$. We highlight that we can do this even when $d > \kappa$, as we can embed $M^*$ into the block diagonal matrix $\mathrm{diag}(M^*, O_{d-\kappa,d-\kappa})$, which preserves both the norms above. Lemma S8 and Proposition S12 then guarantee that

$$\frac{1}{n^2} \|U^*(V^*)^T - \widetilde{M}^*\|_F^2 \leq \tilde{C} \cdot \left( \mathcal{R}_n(UV^T) - \mathcal{R}_n(\widetilde{M}^*) \right) \tag{S82}$$

for some constant $\tilde{C}$ depending only on the quantities mentioned in the theorem statement. As $\mathcal{X}$ is a subset of $\mathcal{B}_{2,\infty}(A_{2,\infty})$, and $(U^*, V^*)$ is a minimizer of $\mathcal{L}(U,V)$, we end up getting that

$$\left( \mathcal{R}_n(UV^T) - \mathcal{R}_n(\widetilde{M}^*) \right) \tag{S83}$$

$$\leq \mathcal{R}_n(UV^T) - \mathcal{L}_n(U^*, V^*) + \mathcal{L}_n(M^*) - \mathcal{R}_n(\widetilde{M}^*) \tag{S84}$$

$$\leq 2 \sup_{(U,V)\in X} \left| \mathcal{R}_n(U,V) - \mathcal{L}_n(U,V) \right| \tag{S85}$$

from which we can apply Proposition S7 to then give the claimed result. $\qquad\square$

We give some brief intuition as to the size of the constants involved here, to understand any potential hidden dependencies involved in them. Of greatest concern are the constants $A_\infty$ and $A_{2,\infty}$ (as the remaining constants are explicit throughout the proof, and depend only on the hyperparameters of the sampling schema and the model in a polynomial fashion). Note that in the case where $k$ is large and we have a SBM$(n, \kappa, \tilde{p}, \tilde{q}, \rho_n)$ model and we apply the DeepWalk scheme, from the discussion after Lemma S9, the minimizing matrix $M^*$ takes the form

$$(M^*)_{lm} \approx \log \left( \frac{\kappa \tilde{p}}{\tilde{p} + (\kappa - 1)\tilde{q}} \right)\delta_{lm} + \log \left( \frac{\kappa \tilde{q}}{\tilde{p} + (\kappa - 1)\tilde{q}} \right)(1 - \delta_{lm}). \tag{S86}$$

Supposing for simplicity that $\tilde{p} > \tilde{q}$, it follows that we can take can take $A_\infty$ to be of the order $O(\log(\tilde{p}/\tilde{q}))$ when $\kappa$ is large. In the rate from Proposition S12, this gives a rate of $O(\tilde{p}/\tilde{q})$ from the $e^{A_\infty}$ factor; note that the dependence on the parameters of the models here are not unreasonable. As for $A_{2,\infty}$, we first highlight the fact that

$$(\kappa - 1) \log \left( \frac{\kappa \tilde{q}}{\tilde{p} + (\kappa - 1)\tilde{q}} \right) \to \frac{\tilde{p} - \tilde{q}}{\tilde{q}} \text{ as } \kappa \to \infty. \tag{S87}$$

By Lemma S31 we can therefore take $A_{2,\infty}$ to be a scalar multiple of $|\log(\tilde{p}/\tilde{q})|^{1/2}$, avoiding any implicit dependence on $\kappa$ or the embedding dimension $d$.

28

## C.9 Convergence of the embedding vectors

We can then get results guaranteeing the convergence of the individual embedding vectors (rather than their gram matrix) up to rotations, as stated by the following theorem.

**Theorem S14.** *Suppose that the conclusion of Theorem S13 holds, and further suppose that $d$ equals the rank of the matrix $M^*$. Then there exists a matrix $\tilde{U}^* \in \mathbb{R}^{\kappa \times d}$ such that*

$$\min_{Q \in O(d)} \frac{1}{n} \sum_{i=1}^{n} \|u_i^* - \tilde{u}_{c(i)}^* Q\|_2^2 = C \cdot \begin{cases} O_p\left(\left(\frac{\max\{\log n, d\}}{n\rho_n}\right)^{1/2}\right) & \text{under Scenarios (i) and (iii);} \\ o_p(1) & \text{under Scenario (ii);} \end{cases} \tag{S88}$$

*Proof.* We handle the cases where $U \neq V$ and $U = V$ separately. For the case where $U \neq V$, we note that without loss of generality we can suppose that $UU^T = VV^T$, in which case we can apply Lemma S23 and Theorem S13 to give the stated result. To do so, we note that by Lemma S25 we have that $n^{-1}\sigma_d(\Pi M^* \Pi^T) \geq c\sigma_d(M^*)$ for some constant $c$ with asymptotic probability 1, as a result of the fact that $n_k(\Pi) \geq 1/2n\pi_k$ with asymptotic probability 1 uniformly across all communities $k \in [\kappa]$. As moreover we have that $n^{-1}\|UV^T - \Pi M^* \Pi^T\|_{\text{op}} \leq n^{-1}\|UV^T - \Pi M^* \Pi^T\|_F = o_p(1)$, the condition that $\|UV^T - \Pi M^* \Pi^T\|_{\text{op}} \leq 1/2\sigma_d(\Pi M^* \Pi^T)$ holds with asymptotic probability 1, we have verified the conditions in Lemma S23, giving the desired result. In the case where we constrain $U = V$, the same argument holds, except we no longer need to verify the condition that $\|UU^* - M^*\|_{\text{op}}$ is sufficiently small, and so we have concluded in this case also. $\square$

In the case of a SBM$(n, \kappa, \tilde{p}, \tilde{q}, \rho_n)$ model it is actually able to give closed form expressions for the embedding vectors which are converged to by factorizing the minima matrix $M^*$ in the way described by the above proof. These details are given in Lemma S31.

# D  Proof of Theorem 4 and Corollary 5

## D.1  Guarantees for community detection

We begin with a discussion of how we can get guarantees for community detection via approximate k-means clustering method, using the convergence criteria for embeddings we have derived already. To do so, suppose we have a matrix $U \in \mathbb{R}^{n \times d}$ corresponding of $n$ columns of $d$-dimensional vectors. Defining the set

$$M_{n,K} := \{\Pi \in \{0, 1\}^{n \times K} : \text{each row of } \Pi \text{ has exactly } K - 1 \text{ zero entries}\}, \tag{S89}$$

we seek to find a factorization $U \approx \Pi X$ for matrices $\Pi \in M_{n,K}$ and $X \in \mathbb{R}^{K \times d}$. To do so, we minimize the objective

$$\mathcal{L}_k(\Pi, X) = \frac{1}{n}\|U - \Pi X\|_F^2 \tag{S90}$$

In practice, this minimization problem is NP-hard [5], but we can find $(1 + \epsilon)$-approximate solutions in polynomial time [24]. As a result, we consider any minimizers $\hat{\Pi}$ and $\hat{X}$ such that

$$\mathcal{L}_k(\hat{\Pi}, \hat{X}) \leq (1 + \epsilon) \min_{\Pi, X} \mathcal{L}_k(\Pi, X). \tag{S91}$$

We want to examine the behavior of k-means clustering on the matrix $U$, when it is close to a matrix $U^*$ which has an exact factorization $U^* = \Pi^* X^*$ for some matrices $\Pi^* \in M_{n,K}$ and $X^* \in \mathbb{R}^{K \times d}$. We introduce the notation

$$G_k(\Pi) := \{i \in [n] : \Pi_{ik} = 1\}, \qquad n_k(\Pi) := |G_k(\pi)| \tag{S92}$$

for the columns of $U$ which are assigned as closest to the $k$-th column of $X$ as according to the matrix $\Pi$.

We make use of the following theorem from Lei and Rinaldo [29], which we restate for ease of use.

**Proposition S15** (Lemma 5.3 of Lei and Rinaldo [29]). *Let $(\hat{\Pi}, \hat{X})$ be any $(1 + \epsilon)$-approximate minimizer to the k-means problem given a matrix $U \in \mathbb{R}^{n \times d}$. Suppose that $U^* = \Pi^* X^*$ for some matrices $\Pi^* \in M_{n,\kappa}$ and $X^* \in \mathbb{R}^{\kappa \times d}$. Fix any $\delta_k \leq \min_{l \neq k}\|X_l^* - X_k^*\|_2$, and suppose that the condition*

$$(16 + 8\epsilon)\|U - U^*\|_F^2/\delta_k^2 < n_k(\Pi^*) \text{ for all } k \in [\kappa] \tag{S93}$$

*holds. Then there exist subsets $S_k \subseteq G_k(\Pi^*)$ and a permutation matrix $\sigma \in \mathbb{R}^{\kappa \times \kappa}$ such that the following holds:*

    *i) For $G = \bigcup_k (G_k(\Pi^*) \setminus S_k)$, we have that $(\Pi^*)_{G\cdot} = \sigma \Pi_{G\cdot}$. In words, outside of the sets $S_k$ we recover the assignments given by $\Pi^*$ up to a re-labelling of the clusters.*

    *ii) The inequality $\sum_{k=1}^{\kappa} |S_k| \delta_k^2 \leq (16 + 8\epsilon) \|U - U^*\|_F^2$ holds.*

In particular, we can then apply this to our consistency results with the embeddings learned by node2vec. Recall that we are interested in the following metrics measuring recovery of communities by any given procedure:

$$L(c, \hat{c}) := \min_{\sigma \in \mathrm{Sym}(\kappa)} \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}[\hat{c}(i) \neq \sigma(c(i))], \tag{S94}$$

$$\widetilde{L}(c, \hat{c}) := \max_{k \in [\kappa]} \min_{\sigma \in \mathrm{Sym}(\kappa)} \frac{1}{|\mathcal{C}_k|} \sum_{i \in \mathcal{C}_k} \mathbb{1}[\hat{c}(i) \neq \sigma(k)]. \tag{S95}$$

These measure the overall misclassification rate and worst-case class misclassification rate respectively.

**Corollary S16.** *Suppose that we have embedding vectors $\omega_i \in \mathbb{R}^d$ for $i \in [n]$ such that*

$$\min_{Q \in O(d)} \frac{1}{n} \sum_{i=1}^{n} \|\omega_i - \eta_{C(i)} Q\|_2^2 = O_p(r_n) \tag{S96}$$

*for some rate function $r_n \to 0$ as $n \to \infty$ and vectors $\eta_l \in \mathbb{R}^d$ for $l \in [\kappa]$. Moreover suppose that $\delta := \min_{l \neq k} \|\eta_l - \eta_k\|_2 > 0$. Then if $\hat{c}(i)$ are the community assignments produced by applying a $(1 + \epsilon)$-approximate k-means clustering to the matrix whose columns are the $\omega_i$, we have that $L(c, \hat{c}) = O_p(\delta^{-2} r_n)$ and $\widetilde{L}(c, \hat{c}) = O_p(\delta^{-2} r_n)$. If the RHS of (S96) is instead $o_p(1)$, then we replace $O_p(r_n)$ by $o_p(1)$ in the statements for $L(c, \hat{c})$ and $\widetilde{L}(c, \hat{c})$.*

*Proof.* We apply Proposition S15 with $\Pi^*$ corresponding to the matrix of community assignments according to $c(\cdot)$, and $X^*$ the matrix whose columns are the $Q\eta_l$ for $l \in [\kappa]$ where $Q \in O(d)$ attains the minimizer in (S96). Letting $U$ be the matrix whose columns are the $\omega_i$ and taking $\delta_k = \delta$, the condition (S93) to verify becomes

$$\frac{16 + 8\epsilon}{\delta^2} \frac{1}{n} \sum_{i=1}^{n} \|\omega_i - Q\eta_{c(i)}\|_2^2 < \frac{|\mathcal{C}_k|}{n} \text{ for all } k \in [\kappa]. \tag{S97}$$

As $r_n \to 0$ and $|\mathcal{C}_l|/n > c > 0$ for some constant $c$ uniformly across vertices $l \in [\kappa]$ with asymptotic probability 1 (as a result of the community generation mechanism, the communities are balanced), the above event will be satisfied with asymptotic probability 1. The desired conclusion follows by making use of the inequalities

$$L(c, \hat{c}) \leq \frac{1}{n} \sum_{k \in [\kappa]} |S_k|, \qquad \widetilde{L}(c, \hat{c}) \leq \max_{k \in [\kappa]} \frac{1}{|\mathcal{C}_k|} |S_k| \leq \left( \max_{k \in [\kappa]} \frac{n}{|\mathcal{C}_k|} \right) \cdot \frac{1}{n} \sum_{l \in [\kappa]} |S_l| \tag{S98}$$

which hold by the first consequence in Proposition S15, and then applying the bound

$$\frac{1}{n} \sum_{k \in [\kappa]} |S_k| \leq \frac{16 + 8\epsilon}{\delta^2} \cdot \frac{1}{n} \sum_{i=1}^{n} \|\omega_i - Q\eta_{c(i)}\|_2^2. \tag{S99}$$

$\square$

We note that in order to apply this theorem, we require the further separation criterion of $\delta > 0$. As a result of Lemma S31, we can guarantee this for the $\mathrm{SBM}(n, \kappa, \tilde{p}, \tilde{q}, \rho_n)$ model when either a) DeepWalk is trained in the unconstrained setting, or b) we are in the constrained setting with $\tilde{p} > \tilde{q}$. As we know that the embedding vectors converge to the zero vector on average when we are in the constrained setting with $\tilde{p} \leq \tilde{q}$, as a result we know that community detection is possible in the constrained setting iff $\tilde{p} > \tilde{q}$, which gives Corollary 5 of the main paper.

## D.2 Guarantees for node classification and link prediction

We now discuss what guarantees we can make when using the embedding vectors for classification. In this section, we suppose that we have a guarantee

$$\frac{1}{n} \min_{Q \in O(d)} \sum_{i=1}^{n} \|u_i - \eta_{C(i)} Q\|_2^2 \leq C(\tau) r_n \qquad \text{holds with probability } \geq 1 - \tau \qquad \text{(S100)}$$

for some constant $C(\tau)$ and rate function $r_n \to 0$ as $n \to \infty$. This is the same as saying that the LHS is $O_p(r_n)$ - it will happen to be more convenient to use this formulation. We also suppose that there exists a positive constant $\delta > 0$ for which

$$\delta \leq \min_{k \neq l} \|\eta_k - \eta_l\|_2. \qquad \text{(S101)}$$

We begin with a lemma which discusses the underlying geometry when we take a small sample of the embedding vectors.

**Lemma S17.** *Suppose we sample $K$ embeddings from the set $(u_i)_{i \in [n]}$, which we denote as $u_{i_1}, \ldots, u_{i_K}$. Define the sets*

$$S_l = \{i \in \mathcal{C}_l \; : \; \|u_i - \eta_{C(i)}\|_2 < \delta/4\}. \qquad \text{(S102)}$$

*Then there exists $n_0(K, \delta, \tau')$ such that if $n \geq n_0$, with probability $1 - \tau'$ we have that $u_{i_j} \in S_{c(i_j)}$ for all $j \in [K]$.*

*Proof.* Without loss of generality, we will suppose that $Q = I$. For each $l \in [\kappa]$, define the sets $S_l = \{i \in \mathcal{C}_l \; : \; \|u_i - \eta_l\|_2 \leq \delta/4\}$. Then by the condition (S100), by Markov's inequality we know that with probability $1 - \tau$ we have that

$$\frac{1}{n} \sum_{l \in [\kappa]} |\mathcal{C}_l \setminus S_l| \leq 4\delta^{-2} C(\tau/2) r_n. \qquad \text{(S103)}$$

We now suppose that we sample $K$ embeddings uniformly at random; for convenience, we suppose that they are done so with replacement. Then the probability that all of the embeddings are outside the set $\bigcup_l (\mathcal{C}_l \setminus S_l)$ is given by $(1 - \frac{1}{n} \sum_l |\mathcal{C}_l \setminus S_l|)^K \geq 1 - \frac{K}{n} \sum_l |\mathcal{C}_l \setminus S_l|$. In particular, this means with probability no less than $1 - \tau - 4K\delta^{-1}C(\tau)r_n$, if we sample $K$ embeddings with indices $i_1, \ldots, i_K$ at random from the set of $n$ embeddings, they lie within the sets $S_{C(i_1)}, \ldots, S_{C(i_K)}$ respectively. The desired result then follows by noting that we take $\tau = \tau'/2$, and choose $n$ such that $4\delta^{-2}C(\tau/2)r_n < \tau'/2$. $\qquad \square$

To understand how this lemma can give insights into the downstream use of embeddings, suppose that we have access to an oracle which provides the community assignments of a vertex when requested, but otherwise the community assignments are unseen.

We note that in practice, only a small number of labels are needed to be provided to embedding vectors in order to achieve good classification results (see e.g the experiments in Hamilton et al. [16], Veličković et al. [46]). As a result, we can imagine keeping $K$ fixed in the regime where $n$ is large. Moreover, the constant $\delta$ simply reflects the underlying geometry of the learned embeddings, and $\tau'$ is a tolerance we can choose such that the stated result is very likely to hold (by e.g choosing $\tau' = 10^{-2}$ or $10^{-3}$). As a consequence, the above lemma tells us with high probability, we can

    i) learn a classifier which is able to distinguish between the sets $S_l$ given use of the sampled embeddings $u_{i_1}, \ldots, u_{i_K}$ and the labels $c(i_1), \ldots, c(i_K)$, provided the classifier is flexible enough to separate $\kappa$ disjoint convex sets; and

    ii) as a consequence of (S103), this classifier will correctly classify a large proportion of vertices within the correct sets $S_l$.

The same argument applies if instead we have classes assigned to embedding vectors which form a coarser partitioning of the underlying community assignments. The importance of the above result is that in order to understand the behavior of embedding methods for classification, it suffices to understand which geometries particular classifiers are able to separate - for example, when the number

of classes equals 2, this reduces down to the classic concept of linear separability, in which case a logistic classifier would suffice.

We end with a discussion as to the task of link prediction, which asks to predict whether two vertices are connected or not given a partial observation of the network. To do so, we suppose that from the observed network, we delete half of the edges in the network, and then train node2vec on the resulting network. Note that the node2vec mechanism only makes explicit use of known edges within the network. This corresponds to training the node2vec model on the data with sparsity factor $\rho_n \to \rho_n/2$; in particular, this leaves the underlying asymptotic representations unchanged and slows the rate of convergence by a factor of 2. With this, a link prediction classifier is formed by the following process:

1. Take a set of edges $J \subseteq \{(i,j) : a_{ij} = 1\}$ for which the node2vec algorithm was not trained on, and a set of non-edges $\tilde{J} \subseteq \{(i,j) : a_{ij} = 0\}$. As in practice networks are sparse, these sets are not sampled randomly from the network, but are assumed to be sampled in a balanced fashion so that the sets $J$ and $\tilde{J}$ are roughly balanced in size. One way of doing so is to pick a number of edges in advance, say $E$, and then sample $E$ elements from the set of edges and non-edges in order to form $J$ and $\tilde{J}$ respectively.

2. Form edge embeddings $e_{ij} = f(u_i, u_j)$ given some symmetric function $f(x,y)$ and node embeddings $u_i$. Two popular choices of functions are the average function $f(x,y) = (x+y)/2$ and the Hadamard product $f(x,y) = (x_i y_i)_{i \in [d]}$.

3. Using the features $e_{ij}$ and the labels provided by the sets $J$ and $\tilde{J}$, build a classifier using your favorite ML algorithm.

By our convergence guarantees, we know that the asymptotic distribution of the edge embeddings $e_{ij}$ will approach some vectors $\eta_{c(i),c(j)} \in \mathbb{R}^d$, giving at most $\kappa^2$ distinct vectors overall. Note that these embedding vectors in of themselves contain little information about whether the edges are connected; that said, even given perfect information of the communities and the connectivity matrix $P$, one can only form probabilistic guesses as to whether two vertices are connected. That said, by clustering together the link embeddings we can identify together edges as having vertices belonging to a particular pair of communities. With knowledge of the sampling mechanism, it is then possible to backout estimates for $p$ and $q$ by counting the overlap of the sets $J$ and $\tilde{J}$ in the neighbourhoods of the clustered node embeddings.

We note that in practice, ML classification algorithms such as logistic regression are used instead. This instead depends on the typical geometry of the sets $J$ and $\tilde{J}$. Suppose we have a SBM$(n, 2, \tilde{p}, \tilde{q}, \rho_n)$ model. In this case, the set $J$ will approximately consist of $\tilde{p}/2(\tilde{p} + \tilde{q}) \times E$ vectors from $\eta_{11}$, $\tilde{p}/2(\tilde{p} + \tilde{q}) \times E$ vectors from $\eta_{22}$, $\tilde{q}/2(\tilde{p} + \tilde{q}) \times E$ vectors from $\eta_{12}$ and $\tilde{q}/2(\tilde{p} + \tilde{q}) \times E$ vectors from $\eta_{21}$. In contrast, the set $\tilde{J}$ will approximately have $E/4$ of each of $\eta_{11}, \eta_{12}, \eta_{21}$ and $\eta_{22}$. As a result, in the case where $\tilde{p} \gg \tilde{q}$, a linear classifier (for example) will be biased towards classifying more frequently vectors with $c(i) = c(j)$, which is at least directionally correct.

So far, we have not talked about the particular mechanism used to form link embeddings from the node embeddings. The Hadamard product is popular, but particularly difficult to analyze given our results, as it does not remain invariant to an orthogonal rotation of the embedding vectors. In contrast, the average link function retains this information. In the SBM$(n, 2, \tilde{p}, \tilde{q}, \rho_n)$, it ends up giving embeddings which will asymptotically depend on only whether $c(i) = c(j)$ or not (i.e, whether the vertices belong to the same community or not).

# E    Intermediate results

## E.1    Sampling probabilities for node2vec

In this section, we derive asymptotic results for the sampling probabilities of edges within node2vec. We begin by recapping the second-order random walk defined for node2vec. To do so, we define a

32

1027 random process $(X_n)_{n \geq 1}$ via the second-order Markov property

$$\mathbb{P}\big(X_n = u \mid X_{n-1} = s, X_{n-2} = v\big) \propto \begin{cases} 0 & \text{if } (u, s) \notin \mathcal{E}, \\ 1/p & \text{if } d_{u,v} = 0 \text{ and } (u, s) \in \mathcal{E}, \\ 1 & \text{if } d_{u,v} = 1 \text{ and } (u, s) \in \mathcal{E}, \\ 1/q & \text{if } d_{u,v} = 2 \text{ and } (u, s) \in \mathcal{E}. \end{cases} \tag{S104}$$

1028 where $d_{u,s}$ denotes the length of the shortest path between $u$ and $s$. Given the extra information
1029 that $(u, s)$ is an edge, $d_{u,v} = 0$ occurs iff $u = v$, $d_{u,v} = 1$ occurs iff $(u, v)$ is an edge, and $d_{u,v} = 2$
1030 occurs iff $(u, v)$ is not an edge (as given that $(v, s)$ is an edge, the shortest path must be $v \to s \to u$).
1031 With this, we select positive samples by selecting $k$ concurrent edges within the walk (via taking a
1032 walk of length $k + 1$).

1033 To initialize the random walk, we note that for the second order walk we need to specify a distribution
1034 on the first two vertices; for DeepWalk where this collapses down to a first order walk, we only need
1035 to specify a distribution on ther first vertex. To do so generally, we consider an initial distribution of
1036 selecting the first vertex via $\pi(u) = \deg(u)/\sum_v \deg(v) = \deg(u)/2E_n$ with $E_n$ being the number
1037 of edges in the graph (single counting $(u, v) \in \mathcal{E}$ and $(v, u) \in \mathcal{E}$), and select the second vertex
1038 uniformly at random from those connected to the first. (Note that this is the transition kernel used
1039 for DeepWalk, and so we handle both cases via this argument.) One can show this is equivalent to
1040 selecting an edge uniformly at random.

1041 For the negative sampling mechanism, we consider the vertices which arose as part of the positive
1042 sampling process - which we denote $V(\mathcal{P})$ - and then sample $l$ vertices independently according to
1043 the unigram distribution

$$\text{Ug}_\alpha(v \mid u, \mathcal{G}_n) = \frac{\deg(v)^\alpha}{\sum_{v' \neq u} \deg(v)^\alpha} \tag{S105}$$

1044 where $u \in V(\mathcal{P})$. We note that the case where $\alpha \to 0$ corresponds to the uniform distribution on
1045 vertices not equal to $u$.

### E.1.1 Proof of Theorem S1

1047 In this section and the next, it will be convenient to use the notation $\sim_p$ to indicate that two
1048 positive random variables $X_n$ and $Y_n$ are asymptotic in the sense that $|X_n/Y_n - 1| = o_p(1)$ when
1049 $n \to \infty$. If we say such a bound happens uniformly over some free variables - say $X_{n,k} \sim_p Y_{n,k}$
1050 uniformly over $k$ - then this means $\max_k |X_{n,k}/Y_{n,k} - 1| = o_p(1)$. We also make extensive
1051 use of the result that if $X_n^{(i)} \sim_p r_n Y_n^{(i)}$ for $i \in \{0, 1\}$ and $Y_n^{(i)} \in [C^{-1}, C]$ for $C > 1$, then
1052 $X_n^{(0)} + X_n^{(1)} \sim_p r_n(Y_n^{(0)} + Y_n^{(1)})$. Indeed, if we write $X_n^{(i)} = Y_n^{(i)} r_n(1 + \epsilon_n^{(i)})$ where $\epsilon_n^{(1)} = o_p(1)$,
1053 then

$$X_n^{(0)} + X_n^{(1)} = r_n(Y_n^{(0)} + Y_n^{(1)}) \cdot \left(1 + \frac{Y_n^{(0)}}{Y_n^{(0)} + Y_n^{(1)}}\epsilon_n^{(0)} + \frac{Y_n^{(1)}}{Y_n^{(0)} + Y_n^{(1)}}\epsilon_n^{(1)}\right) \tag{S106}$$

1054 from which the claimed result follows as the terms weighting the $\epsilon_n^{(1)}$ can be bounded below away
1055 from zero, and are bounded above by 1. We also note that $X_n^{(0)} - X_n^{(1)} = O_p(r_n)$, meaning that the
1056 order of magnitude of terms cannot increase (only decrease) by subtracting them.

1057 As we are interested in the sampling probability of edges within node2vec, it will be convenient
1058 to instead study the first order Markov process $Y_n = (X_n, X_{n-1})$, as then we instead study the
1059 sampling probability of individual states in a regular Markov chain. We note that normally we use
1060 the notation $(u, v)$ to refer an unordered pair belonging to an edge in a graph, but for the Markov
1061 process $(Y_n)_{n \geq 1}$ the order matters, we will write $Y_n = e_{v \to u}$ whenever $X_n = u$ and $X_{n-1} = v$. In
1062 such a scenario, the random walk is therefore defined on the state space

$$S = \bigcup_{(u,v) \in \mathcal{E}} \{e_{u \to v}, e_{v \to u}\}.$$

1063 with the law of $Y$ given by

$$\mathbb{P}\big(Y_n = e_{t \to u} \mid Y_{n-1} = e_{v \to s}\big) = 0 \text{ if } t \neq s, \tag{S107}$$

$$\mathbb{P}\big(Y_n = e_{s \to u} \mid Y_{n-1} = e_{v \to s}\big) \propto \begin{cases} 0 & \text{if } (s, u) \notin \mathcal{E} \\ \frac{1[u=v]}{p} + 1[u \neq v](a_{uv} + \frac{1 - a_{uv}}{q}) & \text{otherwise.} \end{cases} \tag{S108}$$

33

One can calculate the normalizing factor for the probability distribution as being

$$\left(\frac{1}{p} - \frac{1}{q}\right) + \frac{1}{q}\deg(s) + \left(1 - \frac{1}{q}\right)\sum_{u \in \mathcal{V}\setminus\{v\}} a_{su}a_{uv}, \tag{S109}$$

from which we observe that when $p = q = 1$ we recover the simple random walk defined by DeepWalk, as then the probability an edge is selected with source node $u$ is uniform over edges $(u, v)$ where $v$ is a neighbour of $u$.

With this in mind, we define the transition matrix

$$P_{v \to s, s \to u} = \frac{a_{su} \cdot \{1[u = v] \cdot 1/p + 1[u \neq v](a_{uv} + 1/q \cdot (1 - a_{uv}))\}}{\left(\frac{1}{p} - \frac{1}{q}\right) + \frac{1}{q}\deg(s) + \left(1 - \frac{1}{q}\right)\sum_{u \in \mathcal{V}\setminus\{v\}} a_{su}a_{uv}} \tag{S110}$$

governing the transition probabilities on the above chain. We note that by [11, Proposition 72] and Theorem S26 respectively that

$$\deg(s) \sim_p n\rho_n W(\lambda_s, \cdot), \tag{S111}$$

$$\sum_{u \in \mathcal{V}\setminus\{v\}} a_{su}a_{uv} \sim_p n\rho_n^2 T(\lambda_s, \lambda_v) \text{ where } T(\lambda_s, \lambda_v) := \mathbb{E}_{\lambda \sim \mathrm{Unif}[0,1]}[W(\lambda_u, \lambda)W(\lambda, \lambda_v) \,|\, \lambda_u, \lambda_v] \tag{S112}$$

uniformly over all $s, u, v$. As a result, we define

$$\widetilde{P}_{v \to s, s \to u} = \frac{a_{su} \cdot \{q^{-1} + (1 - q^{-1})a_{vu} + \delta_{uv}(p^{-1} - q^{-1})\}}{\left(\frac{1}{p} - \frac{1}{q}\right) + \frac{1}{q}n\rho_n W(\lambda_s, \cdot) + \left(1 - \frac{1}{q}\right)n\rho_n^2 T(\lambda_s, \lambda_v)}. \tag{S113}$$

where $\delta_{uv} := 1[u = v]$ and the numerator is the same as in $P_{v \to s, s \to u}$ (only written in a more convenient to use fashion), and the denominator makes use of the asymptotic statements (S111) and (S112). As a result, we have that $P_{v \to s, s \to u} \sim_p \widetilde{P}_{v \to s, s \to u}$ uniformly over $v, s, u$. In particular, we have that $\widetilde{P}_{v \to s, s \to u} = \Theta_p(a_{su}(n\rho_n)^{-1})$ uniformly over all triples of indices $(v, s, u)$.

Let $A_j(u \to v) = \{Y_j = e_{u \to v}\}$. We then note that the sampling probability of $(u, v)$ being sampled within the first $k + 1$ steps of the second order random walk is given by

$$\mathbb{P}\Big(\bigcup_{j \leq k} A_j(u \to v) \cup A_j(v \to u) \,|\, \mathcal{G}_n\Big). \tag{S114}$$

To ease on the notation going forward, we write $\mathbb{P}_n(\cdot) := \mathbb{P}(\cdot \,|\, \mathcal{G}_n)$. By the inclusion-exclusion principle, we can write this probability as equalling

$$\sum_{\substack{l,m \geq 1 \\ l+m \leq k}} (-1)^{k+m+1} \sum_{\substack{1 \leq i_1 < i_2 < \cdots < i_l \leq k \\ 1 \leq j_1 < j_2 < \cdots < j_m \leq m}} \mathbb{P}_n\Big(\bigcap_{k \leq l} A_{i_k}(u \to v) \cap \bigcap_{k \leq m} A_{j_k}(v \to u)\Big). \tag{S115}$$

We note that the number of terms in this sum is bounded above by $(2k)!$ (some terms will be zero, as we cannot select $e_{u \to v}$ two times in a row), and so for asymptotic purposes we can focus on the individual terms.

We now address the individual probabilities making up this sum. Intuitively, we want to show the following: that the terms for which $(l, m) \neq (1, 0)$ or $(0, 1)$ are asymptotically negligible, and that asymptotically these terms are functions only of $(\lambda_u, \lambda_v)$. We fix a particular instance of the $i_1, \ldots, i_l$ and $j_1, \ldots, j_m$, and denote $\beta_1 < \beta_2 < \cdots < \beta_{l+m}$ for the ordering of these indices. As we use indices $i_k$ to denote the direction $u \to v$ and $j_k$ for the direction $v \to u$, we write

$$A_i(u \to v) =: A_\beta(u, v, 0), \qquad A_j(v \to u) =: A_\beta(u, v, 1) \tag{S116}$$

where the third argument (which we refer to as the orientation herein) indicates which of the first two arguments are used as the source node for the edge. For each $\beta_k$ for $k \leq l + m$, we write $o_k$ to denote this orientation. As a result, it suffices for us to analyze

$$\mathbb{P}_n\Big(\bigcap_{k \leq l+m} A_{\beta_k}(u, v, o_k)\Big) \tag{S117}$$

34

over all sequences $1 \leq \beta_1 < \beta_2 < \cdots < \beta_{l+m} \leq k$ and orientations $(o_k)_{k=1}^{l+m}$. For this, we then note that by the Markov property of the random walk, we are able to write this probability as

$$\left[ \prod_{k \leq l+m-1} \mathbb{P}_n \Big( A_{\beta_{k+1}}(u, v, o_{k+1}) \mid A_{\beta_k}(u, v, o_k) \Big) \right] \cdot \mathbb{P}_n \big( A_{\beta_1}(u, v, o_1) \big) \tag{S118}$$

$$= \left[ \prod_{k \leq l+m-1} \mathbb{P}_n \Big( A_{\beta_{k+1}-\beta_k+1}(u, v, o_{k+1}) \mid A_1(u, v, o_k) \Big) \right] \cdot \mathbb{P}_n \big( A_{\beta_1}(u, v, o_1) \big) \tag{S119}$$

Focusing now on the terms in the product, if $\beta_{k+1} - \beta_k = 1$, then this term equals zero if $o_k = o_{k=1}$, or otherwise equals e.g $P_{u \to v, v \to u}$ which is $O_p((n\rho_n)^{-1})$ as discussed above. If the walk is longer, then by the same argument as in [11, Proposition 73], by conditioning on the second step in the walk one can show this probability is asymptotically of the same order of a walk of length $\beta_{k+1} - \beta_k - 1$ initialized from the uniform distribution on the edges of $\mathcal{G}_n$. As a result, we therefore only need to analyze events of the form

$$\mathbb{P}_n \big( A_\beta(u, v, o) \big) \tag{S120}$$

which will allow us to then show that the events of the form $(l, m) = (1, 0)$ or $(0, 1)$ are the only ones we need to consider in the asymptotic expansion. Going forward, we assume that $o = 0$, as the sum (S115) is symmetric in the orientation $o$ and the arguments are unchanged.

To do so, we begin by writing writing $\pi' = (a_{uv}/|\mathcal{E}|)_{u,v}$ for the initial distribution provided to $Y_1$. To analyze $p_n(u, v, \beta) := \mathbb{P}_n \big( A_\beta(u, v, 0) \big)$, note that when $\beta = 1$ we trivially have that this probability equals $a_{uv}/|\mathcal{E}|$ and we know that $|\mathcal{E}| \sim_p n^2 \rho_n \mathcal{E}_W(1)$. In the case where $\beta \geq 2$, we consider the set of sequences $\alpha = (\alpha_0, \ldots, \alpha_{\beta-2}) \in \mathcal{V}^{\beta-1}$, where we then have that

$$p_n(u, v, 2) = \frac{1}{|\mathcal{E}|} \sum_{\alpha_0} a_{\alpha_0, u} P_{\alpha_0 \to u, u \to v} \tag{S121}$$

$$p_n(u, v, \beta) = \frac{1}{|\mathcal{E}|} \sum_\alpha a_{\alpha_0, \alpha_1} \cdot \prod_{j=1}^\beta P_{\alpha_{j-1} \to \alpha_j, \alpha_j \to \alpha_{j+1}} \cdot P_{\alpha_{\beta-2} \to \alpha_{\beta-1}, \alpha_{\beta-1} \to u} P_{\alpha_{\beta-1} \to u, u \to v} \tag{S122}$$

for $\beta \geq 3$.

To study these sums, we begin by noting that they are asymptotic to their versions where we replace $P \to \widetilde{P}$. Indeed, we note that if we have positive sequences $(a_i)$ and $(b_i)$, then

$$\left| \frac{\sum_j a_j}{\sum_j b_j} - 1 \right| = \frac{|\sum_j b_j(a_j/b_j - 1)|}{\sum_j b_j} \leq \max_j \left| \frac{a_j}{b_j} - 1 \right|, \tag{S123}$$

and so the fact that we know $P \sim_p \widetilde{P}$ uniformly, means that we can apply this to obtain asymptotic formulae for their sums also. With this, if we write $N(\lambda_s, \lambda_t)$ for the denominator of $\widetilde{P}_{t \to s, s \to u}$, $p_n(u, v, \beta)$ can be asymptotically be decomposed into a linear combination of terms (bounded in number by a function of $k$ independent of $n$) of the form

$$\frac{c(p, q) a_{uv}}{|\mathcal{E}|} \sum_{\alpha \in \mathcal{V}^{\beta-1}} \left\{ \Big( \prod_{2 \leq i \leq \beta} N(\lambda_{\tilde{\alpha}_{i-1}}, \lambda_{\tilde{\alpha}_i}) \Big)^{-1} \cdot \prod_{i \leq \beta-1} a_{\tilde{\alpha}_{i-1}, \tilde{\alpha}_i} \cdot \prod_{j \in J} a_{\tilde{\alpha}_{j-1}, \tilde{\alpha}_{j+1}} \cdot \prod_{k \in K} \delta_{\tilde{\alpha}_{k-1}, \tilde{\alpha}_{k+1}} \right\} \tag{S124}$$

where:

- we write $\tilde{\alpha}$ for the concatenation $(\alpha, u, v)$, meaning $\tilde{\alpha}$ is of length $\beta + 1$, with $\tilde{\alpha}_k = \alpha_k$ for $k \leq \beta - 1$, $\tilde{\alpha}_\beta = u$ and $\tilde{\alpha}_{\beta+1} = v$;

- $c(p, q) = (q^{-1})^{\beta - |J| - |K|}(1 - q^{-1})^{|J|}(p^{-1} - q^{-1})^{|K|}$ is a polynomial in $p^{-1}$ and $q^{-1}$;

- $J$ and $K$ are possibly empty subsets of $\{1, \ldots, \beta\}$ which are disjoint.

35

The more tedious part to handle is when the set $K$ is non-empty; as each delta function acts to contract the sum along one variable, doing so allows us to rewrite (S124) as

$$\frac{a_{uv}}{|\mathcal{E}|}c(p,q)\sum_{\alpha\in\mathcal{V}^{\beta-1-|K|}}\left\{\left(\prod_{2\leq i\leq\beta-|K|}N(\lambda_{\tilde{\alpha}_{i-1}},\lambda_{\tilde{\alpha}_i})^{n_i}\right)^{-1}\cdot\prod_{i\leq\beta-1-|K|}a_{\tilde{\alpha}_{i-1},\tilde{\alpha}_i}\cdot\prod_{j\in\tilde{J}}a_{\tilde{\alpha}_{j-1},\tilde{\alpha}_{j+1}}\right\} \tag{S125}$$

after a) performing some relabeling of the indices and modification to the set $J$, to give a new set $\tilde{J}$ which is a subset of $\{1,\ldots,\beta-|K|\}$ and b) introducing some multiplicities $n_i$ which sum to $\beta-1$. By Theorem S26 we uniformly have that this quantity is asymptotic, uniformly over all the free variables in the expression, to

$$\frac{\rho_n^{|\tilde{J}|}}{(n\rho_n)^{|K|}}\cdot\frac{a_{uv}c(p,q)\rho_n^{-1}}{n^2\mathcal{E}_W(1)}\cdot\mathbb{E}\left[\frac{\prod_{i\leq\beta-1-|K|}W(\lambda'_{i-1},\lambda'_i)\prod_{j\in\tilde{J}}W(\lambda'_{j-1},\lambda'_{j+1})}{\prod_{2\leq i\leq\beta-|K|}N'(\lambda'_{i-1},\lambda'_i)^{n_i}}\,\Big|\,\lambda_u,\lambda_v\right] \tag{S126}$$

where we write $\lambda'=(\widetilde{\lambda}_0,\ldots,\widetilde{\lambda}_{\beta-2-|K|},\lambda_u,\lambda_v)$ and $\widetilde{\lambda}$ is an independent copy of $\lambda$, and $N'(\lambda_u,\lambda_v):=(n\rho_n)^{-1}N(\lambda_u,\lambda_v)$. As $n\rho_n\to\infty$ under the prescribed conditions, we only need to consider leading terms of the order $\rho_n^{-1}\,n^2$, which shows that the sampling probability is asymptotic (uniformly over all vertices) to $\rho_n^{-1}\,n^2$ for some function $g_{\mathcal{P}}(\lambda_u,\lambda_v)$. To argue that this function is bounded above away from zero, we note that the terms where $|J|+|K|>0$ will be asymptotically negligible, and the remainder of the terms give a positive weighted sum.

### E.1.2  Proof of Theorem S2

To understand the selection probability for the vertex pair $(u,v)$ to be selected via negative sampling, define the events

$$A_i(u)=\{X_i=u\},\qquad B_i(v|u)=\{v\text{ selected via negative sampling from u}\} \tag{S127}$$

so then

$$\mathbb{P}((u,v)\in\mathcal{N}(\mathcal{G}_n)\,|\,\mathcal{G}_n)=\mathbb{P}\Big(\bigcup_{i=0}^{k}(A_i(u)\cap B_i(v|u))\cup(A_i(v)\cap B_i(u|v))\,|\,\mathcal{G}_n\Big). \tag{S128}$$

We note that

$$\mathbb{P}(A_i(u)\cap B_i(v|u)\,|\,\mathcal{G}_n)=\mathbb{P}(A_i(u)\,|\,\mathcal{G}_n)\cdot\mathbb{P}(\text{Binomial}(l,\text{Ug}_\alpha(v|u))\geq 1\,|\,\mathcal{G}_n). \tag{S129}$$

As a result, we need to begin by understanding the asymptotic probabilities of $\mathbb{P}(A_i(v)\,|\,\mathcal{G}_n)$ and the unigram sampling probability. We begin with understanding the first probability. If $i\in\{0,1\}$, then we have that $\mathbb{P}(A_i(v)\,|\,\mathcal{G}_n)=\deg(v)/2E_n\sim_p W(\lambda_v,\cdot)/n\mathcal{E}_W(1)$ uniformly in $v$ [11, Proposition 72]. For $i\geq 2$, we have that

$$\mathbb{P}(A_i(v)\,|\,\mathcal{G}_n)=\sum_u\mathbb{P}(A_i(u\to v)\,|\,\mathcal{G}_n) \tag{S130}$$

using the same notation as in Appendix E.1.1. Consequently, via the same arguments as in Appendix E.1.1, it will be asymptotic to a positive linear combination of statistics of the form

$$\frac{c(p,q)}{|\mathcal{E}|}\sum_{\alpha\in\mathcal{V}^\beta}\left\{\left(\prod_{2\leq i\leq\beta-|K|}N(\lambda_{\tilde{\alpha}_{i-1}},\lambda_{\tilde{\alpha}_i})^{n_i}\right)^{-1}\cdot\prod_{i\leq\beta-|K|}a_{\tilde{\alpha}_{i-1},\tilde{\alpha}_i}\cdot\prod_{j\in\tilde{J}}a_{\tilde{\alpha}_{j-1},\tilde{\alpha}_{j+1}}\right\} \tag{S131}$$

where we write $\tilde{\alpha}=(\alpha,v)$ for $\alpha\in\mathcal{V}^\beta$. Using the same relabeling and arguments as given in Appendix E.1.1 will be asymptotic to

$$\frac{\rho_n^{|\tilde{J}|}}{(n\rho_n)^{|K|}}\cdot\frac{c(p,q)}{n\mathcal{E}_W(1)}\cdot\mathbb{E}\left[\frac{\prod_{i\leq\beta-|K|}W(\lambda'_{i-1},\lambda'_i)\prod_{j\in\tilde{J}}W(\lambda'_{j-1},\lambda'_{j+1})}{\prod_{2\leq i\leq\beta-|K|}N'(\lambda'_{i-1},\lambda'_i)^{n_i}}\,\Big|\,\lambda_v\right] \tag{S132}$$

uniformly in all the free variables involved, where $\lambda'=(\widetilde{\lambda}_0,\ldots,\widetilde{\lambda}_{\beta-1-|K|},\lambda_v)$ and $\widetilde{\lambda}$ is an independent copy of $\lambda$. (We note that while Theorem S26 is expressed in terms of concentration of quantities around functions which depend on both $\lambda_u$ and $\lambda_v$, the exact same reasoning will apply for statistics which only end up depending on $\lambda_v$.) In particular by taking the highest order terms of this expansion,

we have that there exists some measurable function $g_i(\cdot)$ which is bounded below and above, for each $i$, such that $\mathbb{P}(A_i(u) \mid \mathcal{G}_n) \sim_p n^{-1} g_i(\lambda_u)$ uniformly in $u$.

As for the unigram sampling term, we note that by [11, Proposition 77] we have that

$$\mathbb{P}(\text{Binomial}(l, \text{Ug}_\alpha(v|u)) \sim_p \frac{l W(\lambda_u, \cdot)^\alpha}{n \mathcal{E}_W(\alpha)} \tag{S133}$$

uniformly in the vertices $v, u$. With this, we note that the same arguments via self-intersection allow us to argue that

$$\mathbb{P}((u,v) \in \mathcal{N}(\mathcal{G}_n) \mid \mathcal{G}_n) \sim_p \frac{l}{n^2} \sum_{i=0}^{k} \frac{l}{\mathcal{E}_W(\alpha)} (g_i(\lambda_u) W(\lambda_v, \cdot)^\alpha + g_i(\lambda_v) W(\lambda_u, \cdot)^\alpha) \tag{S134}$$

which gives the claimed result.

### E.2 Chaining and bounds on Talagrand functionals

In this section, let $L > 0$ denote a universal constant (which may differ across occurrences) and $K(\alpha)$ a universal constant which depends on a variable $\alpha$ (but for fixed $\alpha$ also differs across occurrences). For a metric space $(T, d)$, we define the *diameter* of $T$ as

$$\Delta(T) := \sup_{t_1, t_2 \in T} d(t_1, t_2). \tag{S135}$$

We also define the entropy and covering numbers respectively by

$$N(T, d, \epsilon) := \min \left\{ n \in \mathbb{N} \mid F \subseteq T, |F| \leq n, d(t, F) \leq \epsilon \text{ for all } t \in T \right\}, \tag{S136}$$

$$e_n(T) := \inf \left\{ \sup_{t \in T} d(t, T_n) \mid T_n \subseteq T, |T_n| \leq 2^{2^n} \right\} = \inf \left\{ \epsilon > 0 \mid N(t, d, \epsilon) \leq 2^{2^n} \right\}. \tag{S137}$$

We then define the *Talagrand $\gamma_\alpha$ functional* [42] of the metric space $(T, d)$ by

$$\gamma_\alpha(T, d) = \inf \sup_{t \in T} \sum_{n \geq 0} 2^{n/\alpha} \Delta(A_n(t)) \tag{S138}$$

where the infimum is taking over all *admissable sequences*; these are increasing sequences $(\mathcal{A}_n)_{n \geq 0}$ of $T$ such that $|\mathcal{A}_0| = 1$ and $|\mathcal{A}_n| \leq 2^{2^n}$ for all $n$, with $A_n(t)$ being the unique element of $\mathcal{A}_n$ which contains $t$. We will shortly see that this quantity helps to control the supremum of empirical processes on the metric space $(T, d)$. We first give some generic properties for the above functional.

**Lemma S18.** *a) Suppose that $d$ is a metric on $T$, and $M > 0$ is a constant. Then $\gamma_\alpha(T, Md) = M\gamma_\alpha(T, d)$. If $U \subseteq T$, then $\gamma_\alpha(U, d) \leq \gamma_\alpha(T, d)$.*

*b) Suppose that $(T_1, d_1)$ and $(T_2, d_2)$ are metric spaces, so $d = d_1 + d_2$ is a metric on the product space $T = T_1 \times T_2$. Then $\gamma_\alpha(T, d) \leq K(\alpha)(\gamma_\alpha(T_1, d_1) + \gamma_\alpha(T_2, d_2))$.*

*c) We have the upper bounds*

$$\gamma_\alpha(T, d) \leq K(\alpha) \sum_{n \geq 0} 2^{n/\alpha} e_n(T) \leq K(\alpha) \int_0^\infty \left( \log N(T, d, \epsilon) \right)^{1/\alpha} d\epsilon. \tag{S139}$$

*d) Suppose that $\| \cdot \|$ is a norm on $\mathbb{R}^m$, $d$ is the metric induced by $\| \cdot \|$, and $B_A = \{x : \|x\| \leq A\}$. Then one has the bound $N(B_A, d, \epsilon) \leq \max\{(3A/\epsilon)^m, 1\}$, and consequently $\gamma_\alpha(B_A, d) \leq K(\alpha) A m^{1/\alpha}$.*

*Proof.* The first statement in a) is immediate, and the second part is Theorem 2.7.5 a) of Talagrand [42].

For part b), suppose that $\mathcal{A}_n^i$ are admissable sequences for $(T_i, d_i)$ such that

$$\sup_{t_i \in T_i} \sum_{n \geq 0} 2^{n/\alpha} \Delta(A_n^i(t)) \leq 2\gamma_\alpha(T_i, d_i) \text{ for } i = 1, 2. \tag{S140}$$

If we then form the sequence of sets $\mathcal{B}_n := \{A_1 \times A_2 \,:\, A_i \in \mathcal{A}_{n-1}^i\}$ for $n \geq 1$ and $\mathcal{B}_0 = T_1 \times T_2$, we have that $\mathcal{B}_n$ is a partition of $T$ for each $n$, $|\mathcal{B}_0| = 1$ and $|\mathcal{B}_n| = |\mathcal{A}_{n-1}^1| \cdot |\mathcal{A}_{n-1}^2| \leq 2^{2^n}$ for each $n$, meaning that $\mathcal{B}_n$ is an admissable sequence for the metric space $(T, d)$. Moreover, note that we have

$$\Delta((A_1 \times A_2)(t_1, t_2)) = \Delta(A_1(t_1)) + \Delta(A_2(t_2)) \tag{S141}$$

for all sets $A_1 \subseteq T_1$, $A_2 \subseteq T_2$ and $t_1 \in T_1$, $t_2 \in T_2$. As a result, if write $B_n(t_1, t_2) = A_{n-1}^1(t_1) \times A_{n-1}^2(t_2)$ for the unique set in $\mathcal{B}_n$ for which the point $(t_1, t_2)$ lies within it, then we have that

$$\sum_{n \geq 0} 2^{n/\alpha} \Delta(B_n(t_1, t_2)) \leq 2^\alpha \Big( \sum_{n \geq 1} 2^{(n-1)/\alpha} \Delta(A_{n-1}^i(t_1)) + \sum_{n \geq 1} 2^{(n-1)/\alpha} \Delta(A_{n-1}^i(t_2)) \Big). \tag{S142}$$

In particular, taking supremum over all $t \in T$ then gives the result, as the resuling LHS is lower bounded by $\gamma_\alpha(T, d)$, and the resulting RHS is upper bounded by $2(\gamma_\alpha(T_1, d_1) + \gamma_\alpha(T_2, d_2))$.

For part c), the first inequality is Corollary 2.3.2 in Talagrand [42]. As for the second inequality, note that if $\epsilon \leq e_n(T)$, then $N(T, d, \epsilon) > 2^{2^n}$ and consequently $N(T, d, \epsilon) \geq 2^{2^n} + 1$ (recall that both quantities are integers). Writing $N_n = 2^{2^n}$, this implies that

$$\big( \log(1 + N_n) \big)^{1/\alpha} (e_n(T) - e_{n+1}(T)) \leq \int_{e_{n+1}(T)}^{e_n(T)} \big( \log N(T, d, \epsilon) \big)^\alpha \, d\epsilon. \tag{S143}$$

As $\log(1 + N_n) \leq 2^n \log(2)$ for all $n \geq 0$, summation over all $n \geq 0$ implies that

$$(\log 2)^{1/\alpha} \sum_{n \geq 0} 2^{n/\alpha} (e_n(T) - e_{n+1}(T)) \leq \int_0^{e_0(T)} \big( \log N(T, d, \epsilon) \big)^\alpha \, d\epsilon. \tag{S144}$$

As we have that

$$\sum_{n \geq 0} 2^{n/\alpha} \big( e_n(T) - e_{n+1}(T) \big) \geq (1 - 2^{1/\alpha}) \sum_{n \geq 0} 2^{n/\alpha} e_n(T), \tag{S145}$$

combining this and the prior inequality gives the stated result.

For part d), we can calculate that

$$\int_0^\infty \big( \log N(B_A, d, \epsilon) \big)^{1/\alpha} \, d\epsilon \leq \int_0^{3A} m^{1/\alpha} \big( \log(3A/\epsilon) \big)^{1\alpha} \, d\epsilon \leq 3Am^{1/\alpha} \int_0^1 (\log(1/y))^{1/\alpha} \, dy. \tag{S146}$$

For the remaining integral, note that if we make the substitution $y = \exp(-t^\alpha)$, then the integral equals

$$\int_0^1 (\log(1/y))^{1/\alpha} \, dy = \alpha \int_0^\infty t^\alpha e^{-t^\alpha} \, dt, \tag{S147}$$

which we recognize as the mean of an $\text{Exp}(1)$ random variabe in the case where $\alpha = 1$, and the variance of an unnormalized $N(0, 2)$ density in the case where $\alpha = 2$, and so in both cases the integral is finite. The desired conclusion follows. $\qquad\square$

Before stating a corollary of this result involving bounds on the $\gamma$-functional of some of the sets introduced in Theorem S5, we discuss some of the properties of these sets.

**Lemma S19.** *Define the sets*

$$\mathcal{B}_F(A) := \big\{ U \in \mathbb{R}^{n \times d} \,|\, \|U\|_F \leq A \big\}, \tag{S148}$$

$$\mathcal{B}_{2,\infty}(A) := \big\{ U \in \mathbb{R}^{n \times d} \,|\, \|U\|_{2,\infty} \leq A \big\}. \tag{S149}$$

*Moreover, define the metrics*

$$d_F((U_1, V_1), (U_2, V_2)) := \|U_1 - U_2\|_F + \|V_1 - V_2\|_F \tag{S150}$$

$$d_{2,\infty}((U_1, V_1), (U_2, V_2)) := \|U_1 - U_2\|_{2,\infty} + \|V_1 - V_2\|_{2,\infty} \tag{S151}$$

*defined on the space $\mathbb{R}^{n \times d} \times \mathbb{R}^{n \times d}$ of pairs of $n \times d$ matrices. Then we have that for $U_1, U_2, V_1, V_2 \in \mathcal{B}_F(A_F) \cap \mathcal{B}_{2,\infty}(A_{2,\infty})$ that*

$$\|U_1 V_1^T - U_2 V_2^T\|_F \leq A_F d_F((U_1, V_1), (U_2, V_2)), \qquad \|U_1 V_1^T - U_2 V_2^T\|_\infty \leq A_{2,\infty} d_{2,\infty}((U_1, V_1), (U_2, V_2)). \tag{S152}$$

*Moreover, if $U \in \mathcal{B}_{2,\infty}(A)$, then $U \in \mathcal{B}_F(\sqrt{n}A)$ also, and consequently if $U \in \mathcal{B}_{2,\infty}(A_{2,\infty})$ then we have that $U \in \mathcal{B}_{2,\infty}(A_{2,\infty}) \cap \mathcal{B}_F(\sqrt{n}A_{2,\infty})$.*

*Proof.* Begin by noting that, if $U_1, V_1, U_2, V_2 \in \mathbb{R}^{n \times d}$ are matrices, then we have that

$$\|U_1 V_1^T - U_2 V_2^T\|_F = \|U_1(V_1 - V_2)^T + (U_1 - U_2)V_2^T\|_F \leq \|U_1\|_F \|V_1 - V_2\|_F + \|U_1 - U_2\|_F \|V_2\|_F$$

and similarly

$$\|U_1 V_1^T - U_2 V_2^T\|_\infty = \|U_1(V_1 - V_2)^T + (U_1 - U_2)V_2^T\|_\infty \leq \|U_1\|_{2,\infty} \|V_1 - V_2\|_{2,\infty} + \|U_1 - U_2\|_{2,\infty} \|V_2\|_{2,\infty}.$$

As a result, we therefore have that in the case where $U_1, V_1, U_2, V_2$ all have $\| \cdot \|_F \leq A_F$, then

$$\|U_1 V_1^T - U_2 V_2^T\|_F \leq A_F \big(\|U_1 - U_2\|_F + \|V_1 - V_2\|_F\big) \tag{S153}$$

and similarly if each of $U_1, V_1, U_2, V_2$ have $\| \cdot \|_{2,\infty} \leq A_{2,\infty}$ then

$$\|U_1 V_1^T - U_2 V_2^T\| \leq A_{2,\infty} \big(\|U_1 - U_2\|_{2,\infty} + \|V_1 - V_2\|_{2,\infty}\big), \tag{S154}$$

giving the first result of the lemma. The second part follows by noting that

$$\sum_{i=1}^n \sum_{j=1}^d |u_{ij}|^2 \leq n \max_{i \in [n]} \sum_{j=1}^d |u_{ij}|^2 \tag{S155}$$

and taking square roots. $\qquad\square$

**Corollary S20.** *With the same notation as in Lemma S19, and writing $T = \mathcal{B}_F(A_F) \cap \mathcal{B}_{2,\infty}(A_{2,\infty})$, we have that for any constant $C > 0$ that*

$$\gamma_\alpha(T \times T, C d_F) \leq \gamma_\alpha(B_F(A_F), C d_F) \leq K(\alpha) \cdot C A_F (nd)^{1/\alpha} \leq K(\alpha) \cdot C A_{2,\infty} n^{1/2 + 1/\alpha} d^{1/\alpha}, \tag{S156}$$

$$\gamma_\alpha(T \times T, C d_{2,\infty}) \leq \gamma_\alpha(B_{2,\infty}(A_{2,\infty}), C d_F) \leq K(\alpha) \cdot C A_{2,\infty} (nd)^{1/\alpha}. \tag{S157}$$

*Proof.* This is a combination of Lemma S18 and Lemma S19 $\qquad\square$

We now state a result which illustrates the usefulness of the above quantity when trying to control the supremum of empirical processes on a metric space $(T, d)$.

**Theorem S21.** *Suppose $(X_t) t \in T$ is a mean-zero stochastic process, where $d_1$ and $d_2$ are two metrics on $T$. Suppose for all $s, t \in T$ we have the inequality*

$$\mathbb{P}\big(|X_s - X_t| \geq u\big) \leq 2 \exp\Big( - \min\Big\{ \frac{u^2}{d_2(s,t)^2}, \frac{u}{d_1(s,t)} \Big\} \Big). \tag{S158}$$

*Then we have that*

$$\mathbb{P}\big( \sup_{s,t \in T} |X_s - X_t| \geq L u \big(\gamma_2(T, d_2) + \gamma_1(T, d_1)\big)\big) \leq L \exp(-u). \tag{S159}$$

*Proof.* This can be found within the proof of Theorem 2.2.23 in Talagrand [42]. $\qquad\square$

**Corollary S22.** *With the notation of Theorem S5, Lemma S19 and Corollary S20, if we have the bound*

$$\mathbb{P}\big(|E_n(U, V) - E_n(\tilde{U}, \tilde{V})| \geq u\big) \tag{S160}$$

$$\leq 2 \exp\Big( - \min\Big\{ \frac{u^2}{128 \rho_n^{-1} n^{-4} A_F^2 d_F((U,V), (\tilde{U}, \tilde{V}))^2}, \frac{u}{16 \rho_n^{-1} n^{-2} A_{2,\infty} d_{2,\infty}((U,V), (\tilde{U}, \tilde{V}))} \Big\} \Big) \tag{S161}$$

*then as a consequence we can deduce that*

$$\sup_{(U,V),(\tilde{U},\tilde{V}) \in T \times T} \big|E_n(U, V) - E_n(\tilde{U}, \tilde{V})\big| = O_p\Big(A_{2,\infty}^2 \Big(\frac{d}{n \rho_n}\Big)^{1/2} + A_{2,\infty}^2 \frac{d}{n \rho_n}\Big) \tag{S162}$$

*Proof.* This is a consequence of Corollary S20 and Theorem S21. $\qquad\square$

39

**E.3   Matrix Algebra**

**Proposition S23.** *Suppose that we have matrices $U, X \in \mathbb{R}^{n \times d}$ with $n \geq d$, and suppose that $X$ is a full rank matrix so $\sigma_d(XX^T) > 0$. Then we have that*

$$\min_{Q \in O(d)} \frac{1}{n} \|U - XQ\|_F^2 \leq \frac{n^{-2}\|UU^T - XX^T\|_F^2}{\sqrt{2}(\sqrt{2}-1)n^{-1}\sigma_d(XX^T)}. \tag{S163}$$

*Now instead suppose we have matrices $U, V \in \mathbb{R}^{n \times d}$ and a matrix $M \in \mathbb{R}^{n \times d}$ of rank $d$. Let $M = U_M \Sigma V_M^T$ be a SVD of $M$. Moreover suppose that $U^T U = V^T V$, and $\|UV^T - M\|_{op} \leq \sigma_d(M)/2$. Then we have that*

$$\min_{Q \in O(d)} \frac{1}{n} \|U - U_M \Sigma^{1/2} Q\|_F^2 \leq \frac{2n^{-2}\|UV^T - M\|_F^2}{(\sqrt{2}-1)n^{-1}\sigma_d(M)}. \tag{S164}$$

*Proof.* The first part of the theorem statement is Lemma 5.4 of Tu et al. [43]. For the second part, we note that by Proposition S24, we can let $U = U_M \Sigma^{1/2} Q$ and $V = V_M \Sigma^{1/2} Q$ for some orthonormal matrix $Q$, where $\tilde{U}\tilde{\Sigma}\tilde{V}^T$ is the SVD of $UV^T$. As a result, we can therefore apply without loss of generality Lemma 5.14 of Tu et al. [43], which then gives the desired statement. □

**Proposition S24.** *Suppose that $U, V \in \mathbb{R}^{n \times d}$ are matrices such that $UV^T = M$ for some rank $d$ matrix $M \in \mathbb{R}^{n \times n}$. Moreover suppose that $U^T U = V^T V$. Let $M = U_M \Sigma V_M^T$ be the SVD of $M$. Then there exists an orthonormal matrix $Q \in O(d)$ such that $V = V_M \Sigma^{1/2} Q$. In particular, the symmetry group of the mapping $(U, V) \to UV^T$ under the constraint $U^T U = V^T V$ is exactly the orthogonal group $O(d)$.*

*Proof.* Begin by noting that the condition $U^T U = V^T V$ forces there to exist an orthonormal matrix $R \in O(n)$ such that $RU = V$ (e.g by Theorem 7.3.11 of Horn and Johnson [20]). As a consequence, we therefore have that $M = R^{-1}VV^T$. This is a polar decomposition of $M$, and therefore as the semi-positive definite factor is unique, we have that $VV^T = (V_M \Sigma^{1/2})(V_M \Sigma^{1/2})^T$, where $M = U_M \Sigma V_M^T$ is the SVD of $M$, and we highlight that the polar decomposition of $M$ is usually represented by $M = (U_M V_M^{-1}) \cdot (V_M \Sigma V_M^T)$. As $VV^T = (V_M \Sigma^{1/2})(V_M \Sigma^{1/2})^T$, again by e.g Theorem 7.3.11 of Horn and Johnson [20] we have that there exists an orthonormal matrix $Q \in O(d)$ such that $V = V_M \Sigma^{1/2} Q$, giving the desired result. □

**Lemma S25.** *Suppose $X \in \mathbb{R}^{n \times n}$ is a symmetric matrix such that $X = \Pi A \Pi^T$ where $A \in \mathbb{R}^{d \times d}$ is of full rank, and $\Pi \in \mathbb{R}^{n \times d}$ is the assignment matrix for a partition of $[n]$; that is, there exists a partition of $[n]$ into $d$ sets $B(1), \ldots, B(d)$ such that $\Pi_{il} = 1[i \in B(l)]$. Suppose further that $\Pi$ is of full rank. Then we have that $\sigma_d(X) \geq \sigma_d(A) \times \min_l |B(l)|$.*

*Proof.* Let $\Delta = \text{diag}(|B(1)|^{1/2}, \ldots, |B(d)|^{1/2})$. Then note that we can write

$$X = (\Pi \Delta^{-1}) \cdot \Delta A \Delta \cdot (\Pi \Delta)^{-1} \tag{S165}$$

where $(\Pi \Delta^{-1})$ is an orthonormal matrix. As a result, we can simply concentrate on the spectrum of the matrix $\Delta A \Delta$. As the smallest singular value of a matrix product is less than the product of the smallest singular values, the stated result follows. □

**E.4   Concentration inequalities**

**Theorem S26.** *Suppose that $H$ is a graph on a vertex set $\{r_1, \ldots, r_l, v_1, \ldots, v_m\}$ where the vertices $r_i$ are referred to as root vertices, and the remaining vertices as free vertices. We refer to such a graph as a rooted graph. Suppose that all the edges in $H$ have at least one free vertex as an endpoint. Write $\mathbf{x} = (x_1, \ldots, x_m)$ for the collection of $m$ variables $x_i$, and let $Y$ be a statistic of the form*

$$Y = \sum_{x_1, \ldots, x_m \in [n]} g_{\mathbf{x}} \prod_{i \sim_H j} t_{x_i, x_j} \tag{S166}$$

*where the random variables $t_{x_i, x_j}$ are independent and $\{0, 1\}$ valued with $c_p \leq \mathbb{P}(t_{x_i, x_j} = 1) \leq 1 - c_p$ for all $x_i, x_j$; the coefficients $c_g \leq g_x \leq \|g\|_\infty < \infty$ for some $c_g > 0$; and $i \sim_H j$ iff*

$(i, j)$ *is an edge within the graph $H$. Suppose that $\rho_n = n^{-\alpha}$ for some $\alpha < 1/m'(H)$ where* $m'(H) = \max_{2 \leq j \leq k}(j-1)/(v(j)-2)$, $v(j) = \min_{|A| \geq j} v(A)$ *and $v(A)$ for a set of edges $A$ indicates the number of vertices in $A$. Then there exist constants $c, \delta, \Delta$ which depend only on $c_g$, $c_p$, $\|g\|_\infty$, $H$ and $\alpha$ such that*

$$\mathbb{P}\big(\big|Y - \mathbb{E}[Y]\big| \geq \mathbb{E}[Y]\sqrt{\lambda(n^2 \rho_n)^{-1}}\big) \leq \exp(-c\lambda) \tag{S167}$$

*for all $\Delta \leq \lambda \leq n^\delta$.*

*Proof.* Without loss of generality suppose that $\|g\|_\infty = 1$. The proof is essentially the same as Vu [47, Corollary 6.4], where we extend the result derived for the asymptotics of subgraph counts to that of a weighted count of rooted subgraph counts. To do so, we introduce some notation introduced within [47]. If $H$ has $k$ edges, and $A$ is a set of pairs $\{x_i, x_j\}$, we write $\partial_A T$ for the polynomial $\prod_{x \in A} \partial_x T$ when interpreting $T$ as a formal sum in the variables $a_{x_i, x_j}$ (which we recall are $\{0, 1\}$ valued. We then define for $1 \leq j \leq k$ the quantities

$$\mathbb{E}_j[Y] = \max_{|A| \geq j} \mathbb{E}[\partial_A Y], M_j(Y) = \max_{t, |A| \geq j} \partial_A Y(t). \tag{S168}$$

Let $v(A)$ denote the number of vertices specified within the set $A$, and let $v(j) - \min_{|A| \geq j} v(A)$. With this, we note that $\mathbb{E}[Y] = \Theta(n^m \rho_n^k)$ and $\mathbb{E}[\partial_A Y] = \Theta(n^{m-v(A)}\rho_n^{k-|A|})$. Consequently, we have that

$$\mathbb{E}_j[Y] = \max_{h \geq j} \Theta(n^{m-v(h)}\rho_n^{k-h}), \mathbb{E}[Y]/\mathbb{E}_j[Y] = \Theta(\min_{h \geq j} n^{v(h)}\rho_n^h) \tag{S169}$$

where the implied constants depend only on $k$, $c_g$ and $c_p$. The same arguments as given in Claim 6.2 and Corollary 6.4 in [47] can then be applied verbatim to give the claimed result. □

**Lemma S27.** *Let $T$ be a statistic of the form*

$$T' = \sum_{x_1 \neq x_2 \neq \cdots \neq x_m} g(\lambda_{x_1}, \ldots, \lambda_{x_m}) \tag{S170}$$

*where $c_g \leq g(\cdot) \leq \|g\|_\infty < \infty$. Then we have that*

$$\mathbb{P}\big(|T' - \mathbb{E}[T']| \geq \epsilon \mathbb{E}[T']\big) \leq 2\exp\Big(\frac{-\epsilon^2 c_g^2 \lfloor n/m \rfloor}{2\|g\|_\infty^2}\Big). \tag{S171}$$

*Consequently, if we define*

$$T_{l,k} = \sum_{x_1, x_2, \ldots, x_m} g(\lambda_{x_1}, \ldots, \lambda_{x_m}, \lambda_l, \lambda_k), \quad T'_{l,k} = \sum_{x_1 \neq x_2 \neq \cdots \neq x_m} g(\lambda_{x_1}, \ldots, \lambda_{x_m}, \lambda_l, \lambda_k) \tag{S172}$$

*where $c_g \leq g(\cdot) \leq \|g\|_\infty < \infty$ as above, then we have that*

$$\max_{l,k}\Big|\frac{T_{l,k}}{\mathbb{E}[T'_{l,k} \mid \lambda_l, \lambda_k]} - 1\Big| = O_p\Big(\Big(\frac{\log n}{n}\Big)^{1/2}\Big) \tag{S173}$$

*where the implied constant depends only on $m$ and $c_g$.*

*Proof.* The first part is an immediate consequence of Hoeffding's inequality for U-statistics [38], which states that for $U = ((n-m)!/n!) \cdot T$ that

$$\mathbb{P}\Big(|U - \mathbb{E}[U]| \geq t\Big) \leq 2\exp\Big(\frac{-t^2 \lfloor n/m \rfloor}{2\|g\|_\infty^2}\Big), \tag{S174}$$

by substituting in $t \mapsto t\mathbb{E}[U]$ and making use of the bound $\mathbb{E}[U] \geq c_g$.

For the second part, we work conditionally on $\lambda_l, \lambda_k$ and note we can decompose $T_{l,m}$ for each $l, m$ into a sum of statistics of the form $T'$, one of order $\Theta_p(n^m)$ and $\binom{m}{k}$ of order $\Theta_p(n^{m-k})$ (corresponding to when some of the indices $x_i$ are equal) for $1 \leq k \leq m$. By applying the first concentration inequality to these $m! \cdot n^2$ random variables, conditional on the $(\lambda_l, \lambda_k)$, we note the RHS is independent of these quantities, and so the probability bounds hold unconditionally. Consequently, we know that asymptotically $T_{l,k}$ is asymptotic to $T'_{l,k}$, from which we can then apply the resulting concentration bound for this term. □

41

**Theorem S28.** *Suppose we have a statistic of the form*

$$T_{n,\beta,J}(\lambda_u, \lambda_v) = \rho_n^{-\beta-|J|} \sum_{\alpha \in \mathcal{V}^{\beta-1}} g(\lambda_{\tilde\alpha_0}, \dots, \lambda_{\tilde\alpha_{\beta-1}}, \lambda_u, \lambda_v) \prod_{i \le \beta} a_{\tilde\alpha_{i-1}, \tilde\alpha_i} \cdot \prod_{j \in J} a_{\tilde\alpha_{j-1}, \tilde\alpha_{j+1}} \quad \text{(S175)}$$

*where $\tilde\alpha = (\alpha, u, v)$ is a concatenation of $\alpha$, $u$ and $v$ in order, $g : \mathbb{R}^{\beta+1} \to \mathbb{R}$ is a positive function which satisfies $c_g \le g \le \|g\|_\infty < \infty$ for some constant $c_g$, and $J$ is a possibly empty set of indices. Define $\lambda' = (\widetilde\lambda_0, \dots, \widetilde\lambda_{\beta-1}, \lambda_u, \lambda_v)$ where $\widetilde\lambda$ is an independent copy of $\lambda$. Further define the statistic*

$$T'_{n,\beta,J}(\lambda_u, \lambda_v) := \frac{(n-\beta)!}{n!} \cdot \mathbb{E}\Big[ g(\lambda') \prod_{i \le \beta} W(\lambda'_{i-1}, \lambda'_i) \prod_{j \in J} W(\lambda'_{j-1}, \lambda'_{j+1}) \,\big|\, \lambda_u, \lambda_v \Big]. \quad \text{(S176)}$$

*Then for any $\rho_n = n^{-\alpha}$ for $\alpha$ sufficiently small, we have that*

$$\max_{\beta, J, u, v} \Big| \frac{T_{n,\beta,J}(\lambda_u, \lambda_v)}{T'_{n,\beta,J}(\lambda_u, \lambda_v)} - 1 \Big| = O_p\Big( \Big( \frac{(\log n)^k}{n \cdot (n\rho_n)} \Big)^{1/2} \Big). \quad \text{(S177)}$$

*Proof.* For this, we apply the above results. We begin by working conditionally on all of the $\lambda$, whose collection we denote $\lambda$, and note that by Theorem S26 by taking $\lambda = (\log n)^k$ for some $k > 1$ and a union bound, we have that

$$T_{n,\beta,J}(\lambda_u, \lambda_v) = \mathbb{E}[T_{n,\beta,J}(\lambda_u, \lambda_v) \,|\, \lambda] \cdot (1 + E_n^{(1)}) \text{ where } E_n^{(1)} = O\Big( \Big( \frac{(\log n)^k}{n \cdot (n\rho_n)} \Big)^{1/2} \Big) \quad \text{(S178)}$$

uniformly over all $O(m^2 m! \cdot n^2)$ random variables with probability $1 - \exp(O((\log n)^k))$. As we have that

$$\mathbb{E}[T_{n,\beta,J}(\lambda_u, \lambda_v) \,|\, \lambda] = \sum_{\alpha \in \mathcal{V}^{\beta-1}} g(\lambda_{\tilde\alpha_0}, \dots, \lambda_{\tilde\alpha_{\beta-1}}, \lambda_u, \lambda_v) \prod_{i \le \beta} W(\lambda_{\tilde\alpha_{i-1}}, \lambda_{\tilde\alpha_i}) \cdot \prod_{j \in J} W(\lambda_{\tilde\alpha_{j-1}}, \lambda_{\tilde\alpha_{j+1}})$$

$$\text{(S179)}$$

where the function is bounded below by $c_g \cdot c_p^{\beta+|J|}$ and is bounded above by $\|g\|_\infty$, we can make use of Lemma S27 to show that

$$\max_{\beta, J, u, v} \Big| \frac{\mathbb{E}[T_{n,\beta,J}(\lambda_u, \lambda_v) \,|\, \lambda]}{T'_{n,\beta,J}(\lambda_u, \lambda_v)} - 1 \Big| = O_p\Big( \Big( \frac{\log n}{n} \Big)^{1/2} \Big) \quad \text{(S180)}$$

from which the claimed result follows. $\qquad\square$

**Remark 1.** *One natural question to ask about the necessity of the range of values of $\rho_n$ specified above. Generally speaking, one can show for Erdos-Renyi graphs $G(n, p)$ that the number of subgraphs $Y_H$ of $H$ in $\mathcal{G}_n$ satisfy a zero-one law, where*

$$\mathbb{P}(Y_H = 0) = \begin{cases} 1 - o(1) \text{ if } p \ll n^{-c(H)}, \\ o(1) \text{ if } p \gg n^{-c(H)} \end{cases} \quad \text{(S181)}$$

*for some constant $c(H)$ which relates to the geometry of the graph $G$ [6]. In the latter regime, one can then show that $Y_H \sim E[Y_H]$ asymptotically again, and in the former this shows that the term is asymptotically negligible. As the purpose of this result is to derive an asymptotic expansion for the sum of various statistics of the form of $T$ to the highest order, provided $\rho_n$ is of an order which avoids any of the "phase transition" stages of the form above we could eventually generalize our results further. As this involves even more additional book-keeping, we do not do so here.*

**Lemma S29.** *Let $I$ be a finite index set of size $|I| = m$. Suppose that there exist constants $\tau > 0$, a bounded non-negative sequence $(p_i)_{i \in I}$ such that $p_i \le \tau^{-1}$ for all $i$, and a real sequence $(t_i)_{i \in I}$. Define the random variable*

$$X = \frac{1}{m} \sum_{i \in I} \big( \tau^{-1} a_i - p_i \big) t_i \qquad \text{where} \qquad a_i \stackrel{indep}{\sim} \text{Bernoulli}(\tau p_i) \text{ for } i \in I. \quad \text{(S182)}$$

*Then for all $u > 0$, we have that*

$$\mathbb{P}\big( |X| \ge u \big) \le 2 \exp\Big( -\min\Big\{ \frac{u^2}{4\tau^{-1}m^{-2}\|t\|_2^2}, \frac{u}{2\tau^{-1}m^{-1}\|t\|_\infty} \Big\} \Big). \quad \text{(S183)}$$

42

*Proof.* This follows by an application of Bernstein's inequality, by noting that $X$ is a sum of independent mean zero random variables $X_i = m^{-1}(\tau^{-1}a_i - p_i)t_i$ which satisfy

$$|X_i| \leq \tau^{-1}m^{-1}|t_i| \leq \tau^{-1}m^{-1}\|t\|_\infty \text{ for all } i, \qquad \mathbb{E}[X_i^2] \leq m^{-2}\tau^{-1}t_i^2. \qquad \square$$

**Lemma S30.** *Define the random variable*

$$Y = \frac{1}{n(n-1)} \sum_{i \neq j} \left(\rho_n^{-1}a_{ij} - W(\lambda_i, \lambda_j)\right)T_{ij} \tag{S184}$$

*for some constants $(T_{ij})$. Write $\|T\|_2^2 = \sum_{i \neq j} T_{ij}^2$ and $\|T\|_\infty = \max_{i \neq j}|T_{ij}|$. Then we have that*

$$\mathbb{P}\Big(|Y| \geq u\Big) \leq 2\exp\Big(-\min\Big\{\frac{u^2}{128\rho_n^{-1}n^{-4}\|T\|_2^2}, \frac{u}{16\rho_n^{-1}n^{-2}\|T\|_\infty}\Big\}\Big) \tag{S185}$$

*In particular, when $T_{ij} = 1$ for all $i \neq j$, we have that $Y = O_p((n^2\rho_n)^{-1/2})$.*

*Proof.* Note that under the assumptions on the model (where we have that $a_{ij} = a_{ji}$ and $W(\lambda_i, \lambda_j) = W(\lambda_j, \lambda_i)$ for all $i \neq j$), we can write

$$Y = \frac{2}{n(n-1)/2} \sum_{i < j} \left(\rho_n^{-1}a_{ij} - W(\lambda_i, \lambda_j)\right)(T_{ij} + T_{ji}). \tag{S186}$$

Note that

$$\sum_{i < j}(T_{ij} + T_{ji})^2 \leq 2\sum_{i < j}\left(T_{ij}^2 + T_{ji}^2\right) \leq 2\|T\|_2^2, \tag{S187}$$

$$\max_{i < j}|T_{ij} + T_{ji}| \leq \max_{i < j}|T_{ij}| + \max_{i < j}|T_{ji}| \leq 2\|T\|_\infty, \tag{S188}$$

where we have used the inequality $(a+b)^2 \leq 2(a^2+b^2)$ which holds for all $a, b \in \mathbb{R}$. Consequently, as a result of Lemma S29, we have conditional on $\lambda$ that

$$\mathbb{P}\Big(|Y| \geq u \,|\, \lambda\Big) \leq 2\exp\Big(-\min\Big\{\frac{u^2}{128\rho_n^{-1}n^{-4}\|T\|_2^2}, \frac{u}{16\rho_n^{-1}n^{-2}\|T\|_\infty}\Big\}\Big) \tag{S189}$$

As the right hand side has no dependence on $\lambda$, taking expectations gives the first part of the lemma statement. For the second part, note that if $T_{ij} = 1$ for all $i \neq j$, then we have that $\|T\|_2^2 \leq n^2$ and $\|T\|_\infty = 1$, and consequently

$$\mathbb{P}\big(|Y| \geq u\big) \leq 2\exp\Big(-\min\Big\{\frac{u^2}{128\rho_n^{-1}n^{-2}}, \frac{u}{16\rho_n^{-1}n^{-2}}\Big\}\Big) \tag{S190}$$

In particular, this implies that $Y = O_p((n^2\rho_n)^{-1/2})$. $\qquad \square$

## E.5 Miscellaneous results

**Lemma S31.** *Suppose that $A \in \mathbb{R}^{m \times m}$ is a matrix whose diagonal entries are $\alpha$, and off-diagonal entries are $\beta$, so $A_{ij} = \alpha\delta_{ij} + \beta(1 - \delta_{ij})$, where $\delta_{ij}$ is the Kronecker delta. Then $A$ has an eigenvalue $\alpha + (m-1)\beta$ of multiplicity one with eigenvector $1_m$, and an eigenvalue $\alpha - \beta$ of multiplicity $m - 1$, whose eigenvectors form an orthonormal basis of the subspace $\{v : \langle v, 1_m \rangle = 0\}$. For the subspace $\{v : \langle v, 1_m \rangle = 0\}$, we can take the eigenvectors to be*

$$v_i = \frac{1}{\sqrt{2}}(e_{m,1} - e_{m,i+1}) \text{ for } i \in [m-1]$$

*where $e_{m,i}$ are the unit column vectors in $\mathbb{R}^m$, The singular values of $A$ are $|\alpha - \beta|$ and $|\alpha + (\kappa - 1)\beta|$. Consequently, we can write $A = UV^T$ for matrices $U, V \in \mathbb{R}^{m \times m}$ with $UU^T = VV^T$, where the rows of $U$ satisfy*

$$U_{1\cdot} = \frac{|\alpha + \beta(m-1)|^{1/2}}{\sqrt{m}}e_{m,1} + \frac{|\alpha - \beta|^{1/2}}{\sqrt{2}}e_{m,2} \tag{S191}$$

$$U_{i\cdot} = \frac{|\alpha + \beta(m-1)|^{1/2}}{\sqrt{m}}e_{m,1} - \frac{|\alpha - \beta|^{1/2}}{\sqrt{2}}e_{m,i} \text{ for } i \in \{2, \ldots, m\}. \tag{S192}$$

43

*Consequently, we then have that $\|U_{i\cdot}\|_2 \leq \left(2|\alpha + \beta(m-1)|/m + |\alpha - \beta|/2\right)^{1/2}$ for all $i$, and $\min_{i \neq j} \|U_{i\cdot} - U_{j\cdot}\|_2 = (|\alpha - \beta|)^{1/2}$.*

*Further suppose that $\beta = -\alpha/(m-1)$. Then provided $\alpha > 0$, $A$ is positive semi-definite, is of rank $m-1$, with a singular non-zero eigenvalue $\alpha m/(m-1)$ of multiplicity $m-1$. Consequently one can write $A = UU^T$ where $U \in \mathbb{R}^{m \times (m-1)}$ and whose columns equal the $\sqrt{\alpha m/(m-1)}v_i$. In particular, the rows of $U$ equal*

$$U_{1\cdot} = \left(\frac{\alpha m}{2(m-1)}\right)^{1/2} e_{m-1,1}^T, \quad U_{i\cdot} = -\left(\frac{\alpha m}{2(m-1)}\right)^{1/2} e_{m-1,i-1}^T \text{ for } i \in [2, m].$$

*Consequently, one has that $\|U_{i\cdot}\|_2 = \sqrt{\alpha m/(m-1)}$ for all $i$, and moreover we have the separability condition $\min_{1 \leq i < j \leq m} \|U_{i\cdot} - U_{j\cdot}\|_2 = (\alpha m/(m-1))^{1/2}$.*

*Proof.* It is straightforward to verify that $A$ has an eigenvalue of $\alpha + (n-1)\beta$ with the claimed eigenvector. For the second part, we note that the characteristic polynomial of $A$ is

$$\det(A - tI) = (\alpha - \beta - t)^{n-1} \cdot (\alpha + (n-1)\beta - t)$$

and so $A$ has $m-1$ eigenvalues equal to $\alpha - \beta$; as $A$ is symmetric, we know that we can always take eigenvectors to be orthogonal to each other, and consequently the eigenspace associated with such an eigenvalue must be a subspace of $\{v : \langle v, 1_m \rangle = 0\}$. As both of these subspaces are of dimension $m-1$, it consequently follows that they are equal. We then highlight that if $A$ is a symmetric matrix with eigendecomposition $A = Q\Lambda Q^T$ for an orthogonal matrix $Q$, then the SVD is given by $Q|\Lambda|\mathrm{sgn}(\Lambda)Q^T$, and we can write $A = UV^T$ with $U = Q|\Lambda|^{1/2}$ and $V = Q\mathrm{sgn}(\Lambda)|\Lambda|^{1/2}$ such that $UU^T = VV^T$. This allows us to derive the remaining statements about the matrix $A$ which hold in generality. The remaining parts discussing what occurs when $\beta = -\alpha/(m-1)$ follow by routine calculation. $\qquad\square$

**Lemma S32.** *Let $\sigma(x) = (1 + \exp(-x))^{-1}$ be the sigmoid function. Then there exists a unique $y \in \mathbb{R}$ which solves the equation*

$$\alpha\sigma(y) = \beta + \gamma\sigma(-y/s) \tag{S193}$$

*for $\alpha, \gamma, s > 0$ and $\beta \in \mathbb{R}$ if and only if $\beta < \alpha$ and $\beta + \gamma > 0$. Moreover, $y > 0$ if and only if $\beta + \gamma/2 > \alpha/2$.*

*Proof.* Note that $\alpha\sigma(x)$ is a function whose range is $(0, \alpha)$ on $x \in (-\infty, \infty)$, and is strictly monotone increasing on the domain. Similarly, $\beta + \gamma\sigma(-y/s)$ is strictly monotone decreasing with range $(\beta, \beta + \gamma)$, and so simple geometric considerations of the graphs of the two functions gives the existence result. For the second part, note that the ranges of the functions on the LHS and the RHS on the range $y > 0$ are $[\alpha/2, \alpha)$ and $(\beta, \beta + \gamma/2]$ respectively, and so the same considerations as above give the second claim. $\qquad\square$

**Lemma S33.** *Let $\sigma(x) = (e^x)/(1 + e^x)$ be the sigmoid function. Then for any $x, y \in \mathbb{R}$, we have that*

$$-\log(1 - \sigma(x)) \geq -\log(1 - \sigma(y)) + \sigma(y)(x - y) + E(x - y) \tag{S194}$$

*where*

$$E(z) = \begin{cases} \frac{1}{2}e^{-A}z^2 & \text{if } |x|, |y| \leq A, \\ \frac{1}{4}e^{-A}\min\{z^2, 2|z|\} & \text{if either } |x| \leq A \text{ or } |y| \leq A. \end{cases} \tag{S195}$$

*Proof.* Note that by the integral version of Taylor's theorem, for a twice differentiable function $f$ one has for all $x, y \in \mathbb{R}$ that

$$f(x) = f(y) + f'(y)(x - y) + \int_0^1 (1 - t)f''(tx + (1 - t)y)(x - y)^2 \, dt. \tag{S196}$$

Applying this to $f(x) = -\log\sigma(x)$ gives

$$-\log\sigma(x) = -\log\sigma(y) + (-1 + \sigma(y))(x - y) + \int_0^1 (1 - t)(x - y)^2\sigma'(tx + (1 - t)y) \, dt \tag{S197}$$

44

where $\sigma'(x) = e^x/(1 + e^x)^2$. Applying this to $f(x) = \log(1 - \sigma(x))$ gives

$$-\log(1 - \sigma(x)) = -\log(1 - \sigma(y)) + \sigma(y)(x - y) + \int_0^1 (1 - t)(x - y)^2 \sigma'(tx + (1 - t)y)\, dt \quad \text{(S198)}$$

As the integral terms are the same, we concentrate on lower bounding this quantity. To do so, we make use of the lower bound $\sigma'(x) \geq e^{-|x|}/4$ (Lemma 68 of Davison and Austern [11]) which holds for all $x \in \mathbb{R}$. We then note that if $|x|, |y| \leq A$, then we have that

$$-\log(1 - \sigma(x)) = -\log(1 - \sigma(y)) + \sigma(y)(x - y) + \int_0^1 (1 - t)(x - y)^2 \sigma'(tx + (1 - t)y)\, dt$$
$$\text{(S199)}$$

$$\geq -\log(1 - \sigma(y)) + \sigma(y)(x - y) + \frac{e^{-|A|}}{2}(x - y)^2. \quad \text{(S200)}$$

Alternatively, if we only make use of the fact that $|x| \leq A$ (without loss of generality - the argument is essentially equivalent if we only assume that $|y| \leq A$), then we have that

$$\int_0^1 (1 - t)\sigma'(tx + (1 - t)y)(x - y)^2\, dt \geq \int_0^1 (1 - t)e^{-|tx+(1-t)y|}(x - y)^2\, dt \quad \text{(S201)}$$

$$\geq \int_0^1 (1 - t)e^{-|x|}e^{-(1-t)|x-y|}(x - y)^2\, dt \quad \text{(S202)}$$

$$= e^{-|x|}\{|x - y| + e^{-|x-y|} - 1\} \quad \text{(S203)}$$

$$\geq \frac{1}{4}e^{-A}\min\{(x - y)^2, 2|x - y|\}, \quad \text{(S204)}$$

and consequently we get that

$$-\log(1 - \sigma(x)) \geq -\log(1 - \sigma(y)) + \sigma(y)(x - y) + \frac{1}{4}e^{-A}\min\{|x - y|^2, 2|x - y|\} \quad \text{(S205)}$$

as claimed. □

**Lemma S34.** *Suppose that we have a function*

$$f(X) = \frac{1}{m^2} \sum_{i,j=1}^m \min\{X_{ij}^2, 2|X_{ij}|\}. \quad \text{(S206)}$$

*Then if $f(X) \leq r$, we have that $m^{-2} \sum_{i,j=1}^m |X_{ij}| \leq r + r^{1/2}$.*

*Proof.* To proceed, note that if we have that

$$\mathbb{E}[\min\{X^2, 2X\}] \leq r \quad \text{(S207)}$$

for a non-negative random variable $X$, then by Jensen's inequality we get that

$$\left(\mathbb{E}[X1[X < 2]]\right)^2 + \mathbb{E}[X1[X \geq 2]] \leq \mathbb{E}[\min\{X^2, 2X\}] \leq r \quad \text{(S208)}$$

and consequently $\mathbb{E}[X] \leq r + r^{1/2}$ by decomposing $\mathbb{E}[X]$ into the parts where $X \geq 2$ and $X < 2$. Applying this result to the empirical measure on the $|X_{ij}|$ across indices $i, j \in [m]$ gives the desired result. □

# F Minimizers for degree corrected SBMs when $\alpha \neq 1$

In this section, we give an informal discussion of how to study the minimizers of $\mathcal{R}_n(M)$ for degree corrected SBMs when the unigram parameter $\alpha \neq 1$. We begin by highlighting that $\mathcal{R}_n(M)$ does not concentrate around its expectation when averaging over only the degree heterogenity parameters $\theta_i$, which rules out using a similar proof approach as to what was carried out earlier in Appendix 1.

Recall that we were able to derive that the global minima of $\mathcal{R}_n(M)$ was the matrix

$$M_{ij}^* = \log\left(\frac{2\mathcal{E}_W(\alpha)}{(1+k^{-1})\mathbb{E}[\theta]\mathbb{E}[\theta]^\alpha} \cdot \frac{P_{c(i),c(j)}}{\widetilde{P}_{c(i)}\widetilde{P}_{c(j)} \cdot \left(\theta_i^{\alpha-1}\widetilde{P}_{c(i)}^{\alpha-1} + \theta_j^{\alpha-1}\widetilde{P}_{c(j)}^{\alpha-1}\right)}\right). \tag{S209}$$

When $\alpha = 1$ or the $\theta_i$ are constant, this allows us to write $M^* = \Pi M \Pi^T$ where $\Pi$ is the matrix of community assignments for the network and $M$ is some matrix, which allows us to simplify the problem. If we supposed that the $\theta$ actually had some dependence on the $c(i)$ and were discrete - in that $\theta_i | c(i) = l \sim Q_l$ for some discrete distributions $Q_l$ for $l \in [\kappa]$, then we could in fact employ the same type of argument as done throughout the paper. The major change is that then the embedding vectors would each concentrate around a vector decided by both a) their community assignment, and b) the particular degree correction parameter they were assigned. This would then potentially effect our ability to perform community detection depending on the underlying geometry of these vectors. One possible idea would be to explore $\mathcal{R}_n(M)$ partially averaged over the $\theta_i$ - we divide the $\theta_i$ into $B$ bins where $B = n^\beta$ for some $\beta \in (0, 1)$, and average over only over the refinement of the $\theta_i$ as belonging to the different bins. This would be similar to the argument employed in Davison and Austern [11].

An alternative perspective to give some type of guarantee on the concentration of the embedding vectors is to study the rank of the matrix $M^*$. If we are able to prove that is of finite rank $r$ even as $n$ grows large, then we are able to give a convergence result for the embeddings as soon as the embedding dimension $d$ is greater than or equal to $r$. To study this, it suffices to look at the matrix

$$(M_E^*)_{ij} = \log\left(\theta_i^{\alpha-1}\widetilde{P}_{c(i)}^{\alpha-1} + \theta_j^{\alpha-1}\widetilde{P}_{c(j)}^{\alpha-1}\right) \tag{S210}$$

and argue that this is low rank (due to the logarithm, we can write $M^*$ as the difference between this matrix and a matrix of rank $\kappa$, which is therefore also low rank). The entry-wise logarithm is a complicating factor here, as otherwise it is straightforward to argue that the entry-wise exponential of this matrix is of rank 2. One can reduce studying the rank of the matrix $M_E^*$ to studying the rank of the kernel

$$K_M\big((x, c_x), (y, c_y)\big) = \log\left(x^{\alpha-1}\tilde{P}_{c_x}^{\alpha-1} + y^{\alpha-1}\tilde{P}_{c_y}^{\alpha-1}\right) \tag{S211}$$

of an operator $L^2(P) \to L^2(P)$, where $P$ is the product measure induced by $\theta$ and the community assignment mechanism $c$. As $K_M$ is of finite rank $r$ if and only if it can be written as

$$K_M\big((x, c_x), (y, c_y)\big) = \sum_{i=1}^r \phi_i(x, c_x)\psi_i(y, c_y) \tag{S212}$$

for some functions $\phi_i, \psi_i$, it follows that the matrix $(M_E^*)_{ij}$ will be of finite rank $r$ also. Indeed, this representation forces that $M_E^* = \Phi\Psi^T$ for some matrices $\Phi, \Psi \in \mathbb{R}^{n \times r}$, meaning that $M_E^*$ is of rank $\leq r$; Corollary 5.5 of Koltchinskii and Giné [23] then guarantees convergence of the eigenvalues of the matrix $M_E^*$ to the operator $K_M$ so that $M_E^*$ is actually of full rank.

# NeurIPS Paper Checklist

1. **Claims**

   Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

   Answer: [Yes]

   Justification: The claims in the abstract and introduction match the theoretical contributions of the paper. These are supported by experimental verification. Similarly, where theoretical results are not obtained we investigate these scenarios using simulation.

   Guidelines:

   - The answer NA means that the abstract and introduction do not include the claims made in the paper.
   - The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
   - The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
   - It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. **Limitations**

   Question: Does the paper discuss the limitations of the work performed by the authors?

   Answer: [Yes]

   Justification: We highlight and discuss the assumptions required for the theoretical results presented within the paper and in detail in the appendix. We demonstrate empirically the performance of our procedure when these assumptions are relaxed, if theoretical results were not obtained.

   Guidelines:

   - The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
   - The authors are encouraged to create a separate "Limitations" section in their paper.
   - The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
   - The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
   - The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
   - The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
   - If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
   - While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. **Theory Assumptions and Proofs**

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: Due to space constraints all proofs appear in the supplemental material. We provide intuition for these proofs in the paper where space allows. Complete proofs are included in the supplemental material, along with all required Lemmas and exact assumptions.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. **Experimental Result Reproducibility**

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We detail the experimental setup used in this work in the supplemental procedure, along with providing all code required to run and replicate these experiments also.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general. releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).

(d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. **Open access to data and code**

   Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

   Answer: [Yes]

   Justification: We have included an anonymized version of the code repository used to create all experimental results in this paper.

   Guidelines:

   - The answer NA means that paper does not include experiments requiring code.
   - Please see the NeurIPS code and data submission guidelines (`https://nips.cc/public/guides/CodeSubmissionPolicy`) for more details.
   - While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
   - The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (`https://nips.cc/public/guides/CodeSubmissionPolicy`) for more details.
   - The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
   - The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
   - At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
   - Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. **Experimental Setting/Details**

   Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

   Answer: [Yes]

   Justification: We provide sufficient detail in the main text to understand the experimental results presented. In the appendix, we completely detail all experimental details, along with providing the exact code used as supplemental material.

   Guidelines:

   - The answer NA means that the paper does not include experiments.
   - The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
   - The full details can be provided either with the code, in appendix, or as supplemental material.

7. **Experiment Statistical Significance**

   Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

   Answer: [Yes]

   Justification: For all experimental results we either show error bars corresponding to one standard error or all simulation results (in the case of box plots).

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. **Experiments Compute Resources**

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The computation required for individual experiments was relatively small (in terms of both memory and time) and is detailed in the appendix. These were run on a computing cluster.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. **Code Of Ethics**

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: We have ensured the research conforms with the code of ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. **Broader Impacts**

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: This paper provides theoretical guarantees for community detection in a specific class of statistical network models. Any potential societal impacts, positive or negative, will be ancillary from the theoretical focus of this paper.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. **Safeguards**

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: As this work theoretical guarantees for community detection in a specific class of statistical network models, such safeguards are not applicable.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. **Licenses for existing assets**

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We credit the original owners of code and data used.

Guidelines:

- The answer NA means that the paper does not use existing assets.

- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, `paperswithcode.com/datasets` has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. **New Assets**

   Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

   Answer: [Yes]

   Justification: We provide an anonymized zip file which details the code used to generate all results.

   Guidelines:
   - The answer NA means that the paper does not release new assets.
   - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
   - The paper should discuss whether and how consent was obtained from people whose asset is used.
   - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and Research with Human Subjects**

   Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

   Answer: [NA]

   Justification: Crowdsourcing or human subjects were not used in this research.

   Guidelines:
   - The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
   - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
   - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

   Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

   Answer: [NA]

   Justification: Crowdsourcing or human subjects were not used in this research.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.