

The Multilayer Random Dot Product Graph

Andrew Jones and Patrick Rubin-Delanchy

University of Bristol

Abstract

We present a comprehensive extension of the latent position network model known as the random dot product graph to accommodate multiple graphs—both undirected and directed—which share a common subset of nodes, and propose a method for jointly embedding the associated adjacency matrices, or submatrices thereof, into a suitable latent space. Theoretical results concerning the asymptotic behaviour of the node representations thus obtained are established, showing that after the application of a linear transformation these converge uniformly in the Euclidean norm to the latent positions with Gaussian error. Within this framework, we present a generalisation of the stochastic block model to a number of different multiple graph settings, and demonstrate the effectiveness of our joint embedding method through several statistical inference tasks in which we achieve comparable or better results than rival spectral methods. Empirical improvements in link prediction over single graph embeddings are exhibited in a cyber-security example.

Contents

Contents	1
1 Introduction	2
2 The multilayer random dot product graph	5
2.1 Asymptotics and sparsity	6
2.2 Theoretical results	8
3 The generalised multilayer stochastic block model	11
3.1 Undirected graphs	11
3.2 Directed graphs	13
3.3 Bipartite multilayer SBMs	14
3.4 Rank considerations	15
4 Multiple graph inference: comparison with existing methods	16
4.1 Recovery of latent positions	16
4.2 Estimation of invariant subspaces	17
4.3 Model estimation	18
4.4 Two-graph hypothesis testing	19
5 Real data: Link prediction on a computer network	21
5.1 Dynamic link prediction	21
5.2 Port-specific link prediction	22
5.3 Link prediction using mixed data sources	22
6 Chernoff information and the GMSBM	23
7 Conclusion	27

Bibliography	27
Appendix	29
Proof of Theorem 2	41
Proof of Theorem 3	43

1 Introduction

Networks permeate the world in which we live, and so developing an accurate understanding of them is a matter of great interest to many branches of academia and industry, with applications as diverse as identifying patterns in brain scans [12] and the detection of fraudulent behaviour in the financial services sector [1]. The mathematical study of random graphs has its roots in the work of E. N. Gilbert [13] and, contemporaneously, Erdős and Rényi [10], who considered graphs in which edges between nodes occur independently according to Bernoulli random variables with a fixed probability p , in what can be considered the simplest probabilistic model of a naturally occurring network (this type of graph now being universally referred to as an *Erdős-Rényi graph*).

Among the more modern statistical treatments of networks is the concept of a *latent position model* [15], in which the i th node of a graph is mapped to a vector \mathbf{X}_i in some underlying latent space $\mathcal{X} \subseteq \mathbb{R}^d$, and (conditional on this choice of latent positions) the i th and j th nodes connect independently with probability $\kappa(\mathbf{X}_i, \mathbf{X}_j)$, where the *kernel function* $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow [0, 1]$. The *random dot product graph* (RDPG, [25], [40], [5]), which uses the kernel function $\kappa(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{y}$, and its generalisation (GRDPG, [29]), using the kernel function $\kappa(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{I}_{p,q} \mathbf{y}$ (where $\mathbf{I}_{p,q}$ is the diagonal matrix whose entries are p ones followed by q minus ones, where $p + q = d$) are of particular interest here due to their computational tractability and associated statistical estimation theory: spectral embedding [38] the observed graph based on largest eigenvalues (respectively, largest-in-magnitude) produces uniformly consistent estimates of the latent positions \mathbf{X}_i of the RDPG (respectively, GRDPG) model up to orthogonal (respectively, indefinite orthogonal) transformation, with asymptotically Gaussian error [32, 4, 23, 8, 6, 33, 29, 5]. The GRDPG can be used to effectively model networks which exhibit disassortative connectivity behaviour [18] in which dissimilar nodes are the more likely to connect to each other, and is therefore the preferred model for studying biological or technological networks, which typically exhibit such behaviour [24].

Recently, attention has turned to the joint study of *multiple* graphs. Often, the graphs of interest share a common set of nodes but have different edges, and such a collection of graphs is known as a *multilayer* network [19] (more precisely, it is an example of a *multiplex* network). This framework is of interest in the study of *dynamic networks*, in which each of the graphs may represent a “snapshot” of a network at a given point in time, and has been used, for instance, for link prediction in cyber-security applications [28]. Alternatively, one may be interested in detecting differences in node behaviour between graphs, and this approach was used to identify regions of the brain associated with schizophrenia by comparing brain scans of both healthy and schizophrenic patients [21].

Latent position models readily extend to the study of multiple graphs by allowing the kernel functions κ_r to vary, while retaining a common set of latent positions \mathbf{X}_i across the graphs, and in particular there exist several RDPG-based methods for working with multiple graphs. If each graph is drawn from the same distribution (that is, $\kappa_r(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{y}$ for each r) then one can consider the *mean embedding* [35] by spectrally embedding the average of the adjacency matrices $\bar{\mathbf{A}} = \frac{1}{k} \sum_{r=1}^k \mathbf{A}^{(r)}$, or the *omnibus embedding* [21], in which each graph is assigned a different embedding in a common latent space. The mean embedding is known to perform well asymptotically at the task of estimating the latent positions [35], while the omnibus embedding is particularly suited to the task of testing whether the graphs are drawn from the same underlying distribution.

Other RDPG-based methods are more general, allowing different kernel functions κ_r across the graphs. In the *multiple random eigengraph* (MREG, [39]) model, the kernel function $\kappa_r(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{\Lambda}_r \mathbf{y}$ is used, where the matrices $\mathbf{\Lambda}_r \in \mathbb{R}^{d \times d}$ are diagonal with non-negative entries. The *multiple random dot product graph* (multi-RDPG, [26]) loosens these restrictions by allowing the matrices $\mathbf{\Lambda}_r$ to be non-diagonal (but symmetric and positive definite), while requiring that the matrix

of latent positions \mathbf{X} have orthonormal columns. Expanding on this is the *common subspace independent edge* (COSIE, [3]) graph model which allows the matrices $\mathbf{\Lambda}_r$ to have positive and negative eigenvalues, while still requiring them to be symmetric. Each of these models is proposed along with a spectral embedding technique for latent position estimation. Under the COSIE model, the adjacency matrix of each component graph is embedded separately, into a dimension d_r , say. A second, joint, spectral decomposition is then applied to the point clouds, to find a common embedding of dimension d . The approach requires estimation of each d_r as well as d , for which a generic method based on the ‘elbow’ of the scree-plot is suggested [42].

Statistical discourse on network embedding is often inspired by the stochastic block model [16], in which an unknown partition of the nodes exists so that the nodes of any group (or community) are statistically indistinguishable. Under this model, a network embedding procedure can reasonably be expected to ascribe identical positions to the nodes of one group, up to statistical error, and different embedding techniques can therefore be compared through the theoretical performance of an appropriate clustering algorithm at recovering the communities [30, 33, 7].

Of the approaches referred to above, only the COSIE model allows estimation of a generic multilayer stochastic block model [16]. For example, if one graph has assortative community structure (“birds of a feather flock together”) and the other disassortative (“opposites attract”), then the mean embedding can evidently eradicate all community structure visible in either individual graph.

As a model for multiple undirected graph embedding, the Multilayer Random Dot Product Graph model (MRDPG), presented here, is equivalent to the COSIE model in terms of its likelihood given latent positions, but the latent positions are themselves defined differently. The spectral embedding method to which this leads is materially different and simpler, while estimation performance is apparently superior (by numerical experiments). However, the MRDPG in fact allows for far greater generalization than the models we have discussed; it not only extends naturally to accommodate both symmetric and *asymmetric* adjacency matrices (and thus allow us to extract information from *directed* graphs) but also *non-square* binary matrices, such as the non-zero off-diagonal blocks appearing in the adjacency matrix of a bipartite graph, by allowing the columns of such matrices to correspond to a second set of latent positions \mathbf{Y}_i . Moreover, provided that the rows of each matrix correspond to a common set of latent positions \mathbf{X}_i , we may in fact allow the columns of each individual matrix to correspond to *different* sets of latent positions $\mathbf{Y}_i^{(r)}$ (allowing us, for example, to study the behaviour of a particular collection of nodes in a network through observing their interactions—potentially in a number of different settings—with other collections of nodes).

We retain the use of the kernel functions $\kappa_r(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{\Lambda}_r \mathbf{y}$ for each graph, but now note that $\kappa_r : \mathcal{X} \times \mathcal{Y}_r \rightarrow [0, 1]$ for subsets $\mathcal{Y}_r \subseteq \mathbb{R}^{d_r}$, where we allow *arbitrary* matrices $\mathbf{\Lambda}_r \in \mathbb{R}^{d \times d_r}$ and impose no restriction of orthogonality on the matrices of latent positions \mathbf{X} and $\mathbf{Y}^{(r)}$. These latent position matrices are either independent of each other, or else satisfy certain dependence criteria; such as the $\mathbf{Y}^{(r)}$ being equal—potentially up to some linear transformation—with probability one (a full discussion of this is presented in Section 2.1). Given a collection of binary matrices $\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(k)}$ with each $\mathbf{A}^{(r)} \in \{0, 1\}^{n \times n_r}$, those latent positions are then estimated by the following procedure: we form a matrix \mathbf{A} —which we refer to as the *unfolding* of the matrices $\mathbf{A}^{(r)}$ —by adjoining the matrices $\mathbf{A}^{(r)}$, and obtain left and right spectral embeddings of \mathbf{A} by scaling its left and right singular vectors by the square roots of the corresponding singular values. We refer to these embeddings as the Unfolded Adjacency Spectral Embeddings (UASEs) of the $\mathbf{A}^{(r)}$. The left-sided embedding $\mathbf{X}_\mathbf{A}$ is our proposed estimate of \mathbf{X} , that is, a *single* embedding of the nodes that is common to all graphs; the right-sided embedding $\mathbf{Y}_\mathbf{A}$ can be split into k distinct embeddings $\mathbf{Y}_\mathbf{A}^{(r)}$, which can be shown (under certain criteria) to provide estimates of the latent position matrices $\mathbf{Y}^{(r)}$.

We allow the matrices $\mathbf{\Lambda}_r$ to be of non-maximal rank, requiring instead that the matrix $\mathbf{\Lambda} = [\mathbf{\Lambda}_1 | \dots | \mathbf{\Lambda}_k]$ be of maximal rank. This allows for the situation in which information about the latent positions can be *obscured* in individual graphs. As a simple example, consider a three-party political system in which members always vote along party lines, with graphs representing the outcome of votes on particular motions (in which members can either support or oppose the motion, and two members are linked if they vote in the same way). Suppose further that there are no coalitions (that is, every pair of parties has at least one motion on which they vote differently). Then any individual vote will only highlight two groups (those who support and those who oppose

that particular motion) and it is only with knowledge of *multiple* votes that one can correctly identify the individual parties. In our method, no intermediate estimate of the rank d_r of $\mathbf{\Lambda}_r$ is required, in contrast to the COSIE-based approach.

We investigate the asymptotic behaviour of the left- and right-sided embeddings of the unfolding \mathbf{A} under an MRDPG model. It is shown that, up to linear transformation, the rows of each embedding converge uniformly in the Euclidean norm to the latent positions \mathbf{X}_i (Theorem 2) and, through the derivation of a central limit theorem (Theorem 3), that these rows are distributed around their corresponding latent positions according to a Gaussian mixture model, and thereby significantly extending the existing results of [5] and [29] to our more general model. These distributional results show that, in particular, if the graphs are identically distributed then the transformed rows of the left-sided embedding have the same limiting distribution as those of the mean embedding (Corollaries 4 and 5). Consequently, if multiple graphs are identically drawn according to a stochastic block model [16] then joint embedding will *always* be more effective at the task of cluster separation than any individual graph embedding (Proposition 6), where we evaluate this effectiveness via the *Chernoff information* [14] of the limiting Gaussian distributions of the embeddings. The Chernoff information belongs to the class of f -divergences [2, 9] and is therefore invariant under invertible linear transformations [22], an important requirement here since distributional results hold only up to such transformation.

Using the framework of the MRDPG, we propose an extension to the multilayer stochastic block model [16] which we refer to as the *generalised* multilayer stochastic block model (GMSBM). Given a collection of graphs—which may be either directed or undirected—each of which follows a stochastic blockmodel, the GMSBM allows us to study one or more communities common to all graphs by observing their interactions with *all* communities. The point cloud obtained by jointly embedding the unfolding of the resulting adjacency submatrices then exhibits the asymptotic behaviour detailed in our main results; allowing us, for example, to determine whether we can further divide these communities given the information extracted from the multiple embeddings. We provide empirical evidence of the effectiveness of our embedding method at the task of community detection under the GMSBM, demonstrating that fitting a Gaussian mixture model to the UASE achieves better results than rival spectral methods in both the directed and undirected cases.

We assess the effectiveness of unfolded adjacency spectral embedding for the general MRDPG at the inference tasks of recovery of latent positions, estimation of the common invariant subspaces, estimation of the underlying probabilistic model and two-graph hypothesis testing in simulated data. We demonstrate that performance at the estimation tasks is often better than that of the multiple adjacency spectral embedding (the method proposed in [3] to embed multiple graphs distributed according to a COSIE model, which demonstrably yields state of the art performance at such tasks), while its performance at the latter task is comparable with that of the omnibus embedding for reasonably-sized graphs (those with at least 500 nodes). We also apply the UASE to the task of link prediction, using connectivity data from the Los Alamos National Laboratory computer network [37] to predict connections between computers across the entire network as an example of a dynamic link prediction inference task, before restricting our attention to connections occurring through specific ports, demonstrating that the majority of the time the UASE yields greater accuracy than individual adjacency spectral embeddings. As a final example, we show how incorporating connection data between computers and ports into our model can increase the accuracy of our link prediction method.

The remainder of this article is structured as follows. In Section 2 we present our model, and corresponding asymptotic results. In Section 3 these results are explored within the context of a stochastic block model, appropriately extended to different multiple graph settings, all special cases of the MRDPG. Section 4 presents a series of experiments comparing the performance of unfolded adjacency spectral embedding with rival methods at different inference tasks. Section 5 presents an example from a real computer network, from which multiple graphs are extracted by considering different time windows or different port numbers (indicating different types of network service), and are used together to improve link prediction. In Section 6, we investigate the statistical gain of joint versus individual spectral embedding in terms of Chernoff information, proving there is definite improvement in one special (but interesting) case, leaving open a wider conjecture. Finally, Section 7 concludes.

2 The multilayer random dot product graph

Definition 1. (The multilayer random dot product graph model).

For a positive integer k , fix matrices $\mathbf{\Lambda}_r \in \mathbb{R}^{d \times d_r}$ for each $r \in \{1, \dots, k\}$ such that the matrix $\mathbf{\Lambda} = [\mathbf{\Lambda}_1 | \dots | \mathbf{\Lambda}_k]$ is of rank d , and fix bounded subsets $\mathcal{X}, \mathcal{Y}_1, \dots, \mathcal{Y}_k$ of $\mathbb{R}^d, \mathbb{R}^{d_1}, \dots, \mathbb{R}^{d_k}$ respectively such that for each r , $\mathbf{x}^\top \mathbf{\Lambda}_r \mathbf{y}_r \in [0, 1]$ for all $\mathbf{x} \in \mathcal{X}$ and $\mathbf{y}_r \in \mathcal{Y}_r$.

Fix a joint distribution \mathcal{F} supported on $\mathcal{X}^n \times \mathcal{Y}_1^{n_1} \times \dots \times \mathcal{Y}_k^{n_k}$ for positive integers n, n_1, \dots, n_k and let $(\mathbf{X}_1, \dots, \mathbf{X}_n, \mathbf{Y}_1^{(1)}, \dots, \mathbf{Y}_{n_k}^{(k)}) \sim \mathcal{F}$. Define $\mathbf{X} = [\mathbf{X}_1 | \dots | \mathbf{X}_n]^\top \in \mathbb{R}^{n \times d}$ and $\mathbf{Y} = \mathbf{Y}^{(1)} \oplus \dots \oplus \mathbf{Y}^{(k)}$, where each matrix $\mathbf{Y}^{(r)} = [\mathbf{Y}_1^{(r)} | \dots | \mathbf{Y}_{n_r}^{(r)}]^\top \in \mathbb{R}^{n_r \times d_r}$. Finally, define matrices $\mathbf{P}^{(r)} = \mathbf{X} \mathbf{\Lambda}_r \mathbf{Y}^{(r)\top} \in [0, 1]^{n \times n_r}$ for each r , and define the *unfolding* $\mathbf{P} = [\mathbf{P}^{(1)} | \dots | \mathbf{P}^{(k)}]$ (which we note satisfies $\mathbf{P} = \mathbf{X} \mathbf{\Lambda} \mathbf{Y}^\top$).

Given a set of matrices $\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(k)}$ with each $\mathbf{A}^{(r)} \in \{0, 1\}^{n \times n_r}$, we similarly define the unfolding $\mathbf{A} = [\mathbf{A}^{(1)} | \dots | \mathbf{A}^{(k)}]$. We say that $(\mathbf{A}, \mathbf{X}, \mathbf{Y}) \sim \text{MRDPG}(\mathcal{F}, \mathbf{\Lambda})$ if each matrix $\mathbf{A}^{(r)}$ satisfies one of the following:

- If, with probability one, there exists a matrix $\mathbf{G}_r \in \mathbb{R}^{d \times d_r}$ such that $\mathbf{Y}^{(r)} = \mathbf{X} \mathbf{G}_r$, then conditional on \mathbf{X} and $\mathbf{Y}^{(r)}$, either:
 - $\mathbf{A}^{(r)}$ is hollow and symmetric, satisfying $\mathbf{A}_{ij}^{(r)} \sim \text{Bern}(\mathbf{P}_{ij}^{(r)})$ for all $i > j$, or
 - $\mathbf{A}^{(r)}$ is hollow and asymmetric, satisfying $\mathbf{A}_{ij}^{(r)} \sim \text{Bern}(\mathbf{P}_{ij}^{(r)})$ for all $i \neq j$.
- If, for all matrices $\mathbf{G}_r \in \mathbb{R}^{d \times d_r}$, the probability that $\mathbf{Y}^{(r)} = \mathbf{X} \mathbf{G}_r$ is strictly less than one then, conditional on \mathbf{X} and $\mathbf{Y}^{(r)}$, $\mathbf{A}_{ij}^{(r)} \sim \text{Bern}(\mathbf{P}_{ij}^{(r)})$ for all i and j ;

We say that the graph corresponding to the matrix $\mathbf{A}^{(r)}$ is “*undirected*”, “*directed*” or “*bipartite*” to distinguish between these cases, and allow a mixture of these cases among the matrices $\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(k)}$ (we note that for the sake of brevity—particularly within the proofs located in the Appendix—we will occasionally abuse our terminology and refer to the *matrices* $\mathbf{A}^{(r)}$ using these terms).

The distributional results that we will obtain later can be refined if we restrict our attention to the case in which the matrices $\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(k)}$ are identically distributed (corresponding, say, to a multiple graph embedding in which we expect identical behaviour across the graphs). To this end, for a fixed matrix $\mathbf{\Lambda} \in \mathbb{R}^{d \times d'}$ of rank d and a distribution \mathcal{F} , we write $(\mathbf{A}, \mathbf{X}, \mathbf{Y}) \stackrel{\text{id}}{\sim} \text{MRDPG}(\mathcal{F}, \mathbf{\Lambda})$ if, under the distribution \mathcal{F} , all of the matrices $\mathbf{Y}^{(r)}$ are equal with probability one to $\mathbf{Y} \in \mathbb{R}^{n' \times d'}$, and the matrices $\mathbf{\Lambda}_r$ are all equal to $\mathbf{\Lambda}$.

We note that there is a degree of ambiguity in the choice of latent positions for the MRDPG, as the following result shows:

Proposition 1. *Let $(\mathbf{A}, \mathbf{X}, \mathbf{Y}) \sim \text{MRDPG}(\mathcal{F}, \mathbf{\Lambda})$. Then:*

- (i) $(\mathbf{A}, \mathbf{X} \mathbf{G}^\top, \mathbf{Y} \mathbf{H}^\top) \sim \text{MRDPG}(\mathcal{F}_{\mathbf{G}, \mathbf{H}}, \mathbf{G}^{-\top} \mathbf{\Lambda} \mathbf{H}^{-1})$ for any matrices $\mathbf{G} \in \text{GL}(d)$ and $\mathbf{H} = \mathbf{H}_1 \oplus \dots \oplus \mathbf{H}_k \in \text{GL}(d_1) \times \dots \times \text{GL}(d_k)$, where the distribution $\mathcal{F}_{\mathbf{G}, \mathbf{H}}$ is derived from \mathcal{F} by multiplying elements of \mathcal{X} by \mathbf{G} and elements of \mathcal{Y}_r by \mathbf{H}_r for each r .
- (ii) There is a joint distribution $\tilde{\mathcal{F}}$ and matrices of latent positions $\tilde{\mathbf{X}}$ and $\tilde{\mathbf{Y}}$ where each vector $\tilde{\mathbf{X}}_i, \tilde{\mathbf{Y}}_j^{(r)} \in \mathbb{R}^d$ such that $(\mathbf{A}, \tilde{\mathbf{X}}, \tilde{\mathbf{Y}}) \sim \text{MRDPG}(\tilde{\mathcal{F}}, \mathbf{I}_{d,k})$, where $\mathbf{I}_{d,k} = [\mathbf{I}_d | \dots | \mathbf{I}_d]$.

Proof.

- (i) This follows from the fact that $(\mathbf{X} \mathbf{G}^\top)(\mathbf{G}^{-\top} \mathbf{\Lambda} \mathbf{H}^{-1})(\mathbf{Y} \mathbf{H}^\top)^\top = \mathbf{\Lambda}$.
- (ii) Let $\mathbf{\Lambda}$ admit the singular value decomposition $\mathbf{\Lambda} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$ with matrices $\mathbf{U} \in \text{O}(d)$, $\mathbf{\Sigma} \in \mathbb{R}^{d \times d}$ and $\mathbf{V} \in \text{O}((d_1 + \dots + d_k) \times d)$. Then setting $\tilde{\mathbf{X}} = \mathbf{X} \mathbf{U} \mathbf{\Sigma}^{1/2}$ and $\tilde{\mathbf{Y}} = \mathbf{Y} \mathbf{V} \mathbf{\Sigma}^{1/2}$ gives the result. □

Note that under the second transformation, while $\widetilde{\mathbf{X}}$ is always of maximal rank, the rank of $\widetilde{\mathbf{Y}}^{(r)}$ is equal to that of $\mathbf{\Lambda}_r$ and so it is possible to “lose” information about the latent positions $\mathbf{Y}^{(r)}$ in some sense by applying such a transformation.

Key to our study of the MRDPG will be the spectral embeddings of the unfolding \mathbf{A} . Unlike the GRDPG, in which one considers a single symmetric adjacency matrix, whose left and right singular vectors therefore coincide, we obtain distinct embeddings by considering each side of \mathbf{A} :

Definition 2. (Unfolded adjacency spectral embeddings).

Let $(\mathbf{A}, \mathbf{X}, \mathbf{Y}) \sim \text{MRDPG}(\mathcal{F}, \mathbf{\Lambda})$, and let \mathbf{A} and \mathbf{P} admit singular value decompositions

$$\mathbf{A} = \mathbf{U}_\mathbf{A} \mathbf{\Sigma}_\mathbf{A} \mathbf{V}_\mathbf{A}^\top + \mathbf{U}_{\mathbf{A},\perp} \mathbf{\Sigma}_{\mathbf{A},\perp} \mathbf{V}_{\mathbf{A},\perp}^\top, \quad \mathbf{P} = \mathbf{U}_\mathbf{P} \mathbf{\Sigma}_\mathbf{P} \mathbf{V}_\mathbf{P}^\top, \quad (1)$$

where $\mathbf{U}_\mathbf{A}, \mathbf{U}_\mathbf{P} \in \mathbb{O}(n \times d)$, $\mathbf{V}_\mathbf{A}, \mathbf{V}_\mathbf{P} \in \mathbb{O}((n_1 + \dots + n_k) \times d)$, and $\mathbf{\Sigma}_\mathbf{A}, \mathbf{\Sigma}_\mathbf{P} \in \mathbb{R}^{d \times d}$ are diagonal containing the largest singular values of \mathbf{A} and \mathbf{P} respectively. We then define the following:

- The left UASE is the matrix $\mathbf{X}_\mathbf{A} \in \mathbb{R}^{n \times d}$ given by $\mathbf{X}_\mathbf{A} = \mathbf{U}_\mathbf{A} \mathbf{\Sigma}_\mathbf{A}^{1/2}$.
- For $r \in \{1, \dots, k\}$, the r th right UASE is the matrix $\mathbf{Y}_\mathbf{A}^{(r)} \in \mathbb{R}^{n_r \times d}$ obtained by dividing $\mathbf{Y}_\mathbf{A} = \mathbf{V}_\mathbf{A} \mathbf{\Sigma}_\mathbf{A}^{1/2}$ into k blocks of sizes $n_1 \times d, \dots, n_k \times d$ (equivalently, $\mathbf{Y}_\mathbf{A}^{(r)} = \mathbf{V}_\mathbf{A}^{(r)} \mathbf{\Sigma}_\mathbf{A}^{1/2}$, where we divide $\mathbf{V}_\mathbf{A}$ into k blocks $\mathbf{V}_\mathbf{A}^{(1)}, \dots, \mathbf{V}_\mathbf{A}^{(k)}$).

We define the matrices $\mathbf{X}_\mathbf{P}$, $\mathbf{Y}_\mathbf{P}$ and $\mathbf{Y}_\mathbf{P}^{(r)}$ analogously.

The reader with a passing knowledge of tensor theory may wish to draw parallels between our notion of an unfolding and that of the matrix unfoldings of a tensor. For the uninitiated, any 3-tensor $\mathcal{M} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ can be represented as a matrix in one of 3 standard ways, each of which is known as an unfolding of \mathcal{M} (a precise description of these can be found in [20]). The unfoldings of a tensor provide a more accessible means for studying its properties; in particular, there is a notion of a singular value decomposition for tensors (see [20], Theorem 2) in which the “higher order” analogues of singular values and vectors correspond to the standard matrix singular values and vectors of the unfoldings.

In the special case in which the matrices of latent positions $\mathbf{Y}^{(r)}$ are equal with probability one (and so we may consider the $\mathbf{A}^{(r)}$ to be adjacency submatrices for some fixed subset of nodes across a series of graphs) then the matrices $\mathbf{A}^{(r)}$ give rise to a 3-tensor \mathcal{A} by setting $\mathcal{A}_{ijr} = \mathbf{A}_{ij}^{(r)}$. In this case, our notion of the unfolding of the matrices $\mathbf{A}^{(r)}$ is a column permutations of the first standard unfolding of the tensor \mathcal{A} , which therefore shares the same left singular values and vectors, and is the natural choice to consider to allow us extend to the general case in which the matrices $\mathbf{Y}^{(r)}$ differ (the second standard unfolding corresponds to the matrix $[\mathbf{A}^{(1)\top} | \dots | \mathbf{A}^{(k)\top}]$, while the third is the matrix whose (i, j) th entry is the Frobenius inner product $\langle \mathbf{A}^{(i)}, \mathbf{A}^{(j)} \rangle_F$).

2.1 Asymptotics and sparsity

In our asymptotic analysis of the behaviour of the UASE, we will assume that the number of graphs k is either fixed, or at most grows at a rate much slower than n , in the sense that $\lim_{k, n \rightarrow \infty} \frac{k}{n} = 0$. We will also assume that the number of latent positions $\mathbf{Y}_i^{(r)}$ grows at a comparable rate to the number of latent positions \mathbf{X}_i , in the sense that for each $r \in \{1, \dots, k\}$, there exists a positive constant c_r such that $\lim_{n_r, n \rightarrow \infty} \frac{n_r}{n} = c_r$.

Before proceeding further, we establish some (standard) notation relating to the asymptotic growth of various functions. In the following, f and g are real-valued functions of n, n_1, \dots, n_k , and X and Y are real-valued random variables. We say that:

- $f = \Omega(g)$ if there is a constant $c > 0$ and integers N, N_1, \dots, N_k such that for all $n \geq N$ and $n_r \geq N_r$, $f(n, n_1, \dots, n_k) \geq cg(n, n_1, \dots, n_k)$;
- $f = \mathcal{O}(g)$ if there is a constant $c > 0$ and integers N, N_1, \dots, N_k such that for all $n \geq N$ and $n_r \geq N_r$, $f(n, n_1, \dots, n_k) \leq cg(n, n_1, \dots, n_k)$;
- $f = \Theta(g)$ if both $f = \Omega(g)$ and $f = \mathcal{O}(g)$;

- $f = \omega(g)$ if there is a constant $c > 0$ and integers N, N_1, \dots, N_k such that for all $n \geq N$ and $n_r \geq N_r$, $f(n, n_1, \dots, n_k) \geq cg(n, n_1, \dots, n_k)$ and $\lim_{n, n_1, \dots, n_k \rightarrow \infty} \left| \frac{f(n, n_1, \dots, n_k)}{g(n, n_1, \dots, n_k)} \right| = \infty$;
- $|X| = O(f)$ *almost surely* if, for any $\alpha > 0$, there is a constant $c > 0$ and integers N, N_1, \dots, N_k such that for all $n \geq N$ and $n_r \geq N_r$, $|X| \leq cf(n, n_1, \dots, n_k)$ with probability at least $1 - n^{-\alpha}$;
- $|X| = O(f)$ and $|Y| = O(g)$ *mutually almost surely* if, for any $\alpha > 0$, there is a constant $c > 0$ and integers N, N_1, \dots, N_k such that for all $n \geq N$ and $n_r \geq N_r$, the probability that both $|X| \leq cf(n, n_1, \dots, n_k)$ and $|Y| \leq cg(n, n_1, \dots, n_k)$ is at least $1 - n^{-\alpha}$.

The inter- and intra-dependence of \mathbf{X} and the $\mathbf{Y}^{(r)}$ under the joint distribution \mathcal{F} in our definition of the MRDPG has so far been left open. In order to establish asymptotic distributional results for the UASE, however, we must impose certain restrictions on \mathcal{F} . The first of these is that:

- Each of the collections $(\mathbf{X}_1, \dots, \mathbf{X}_n)$ and $(\mathbf{Y}_1^{(r)}, \dots, \mathbf{Y}_{n_r}^{(r)})$ is marginally i.i.d., with marginal distributions \mathcal{F}_X on \mathcal{X} and $\mathcal{F}_{Y,r}$ on \mathcal{Y}_r for each r .

Given such a joint distribution \mathcal{F} , let $\xi \sim \mathcal{F}_X$ and $v_r \sim \mathcal{F}_{Y,r}$ for each r . Our next requirement is that these marginal distributions be non-degenerate, in the sense that:

- The second moment matrices $\Delta_X = \mathbb{E}[\xi\xi^\top] \in \mathbb{R}^{d \times d}$ and $\Delta_{Y,r} = \mathbb{E}[v_r v_r^\top] \in \mathbb{R}^{d_r \times d_r}$ for each r are all invertible.

Our third and final requirement is a technical condition which ensures that the growth of the singular values of the unfolding \mathbf{P} is regulated. We note first that a standard application of Hoeffding's inequality shows that the spectral norm $\|\mathbf{X}^\top \mathbf{X} - n\Delta_X\|$ is of order $O(n^{1/2} \log(n))$ almost surely, and one can argue similarly that the spectral norms $\|\mathbf{Y}^{(r)\top} \mathbf{Y}^{(r)} - n_r \Delta_{Y,r}\|$ are of order $O(n_r^{1/2} \log(n_r))$ almost surely. Our final requirement is that these bounds are attained simultaneously with high probability, or in other words that:

- The matrices $\mathbf{X}, \mathbf{Y}^{(1)}, \dots, \mathbf{Y}^{(k)}$ satisfy $\|\mathbf{X}^\top \mathbf{X} - n\Delta_X\| = O(n^{1/2} \log(n))$ and, for each r , $\|\mathbf{Y}^{(r)\top} \mathbf{Y}^{(r)} - n_r \Delta_{Y,r}\| = O(n_r^{1/2} \log(n_r))$ mutually almost surely.

This condition is satisfied, for example, when \mathbf{X} and the $\mathbf{Y}^{(r)}$ are independent, or when some or all of the $\mathbf{Y}^{(r)}$ are equal (possibly to \mathbf{X}) with probability one. Due to submultiplicativity of the spectral norm, it also holds if we extend the latter example to allow equality up to linear transformation.

When studying the asymptotic behaviour of graph embeddings, it is typical to impose some control over how sparse or dense the graphs we are considering are allowed to be. This is often done through the introduction of a *sparsity factor*, which is a real-valued function (depending on the size of the graph) which is either equal to 1, or which tends to zero as the size of the graph increases (corresponding to dense and sparse regimes respectively). This factor can then be incorporated into the model *either* by scaling the latent positions themselves, *or* by scaling the kernel function. In the case of a single graph which follows a GRDPG, these two options are equivalent (and for example in [29] the sparsity factor is used to scale the latent positions \mathbf{X}_i). For multiple graphs, however, scaling the latent positions produces the same sparsity behaviour for *all* graphs, and does not allow individual control over the density of each graph.

Our approach is to incorporate both a *global* sparsity function ρ and *local* sparsity functions ϵ_r into our model, where ρ applies a uniform scaling to the latent positions \mathbf{X} and $\mathbf{Y}^{(r)}$ and ϵ_r an individual scaling to the matrix $\mathbf{\Lambda}_r$. Given a joint distribution \mathcal{F} satisfying the above conditions, let $(\xi_1, \dots, \xi_n, v_1^{(1)}, \dots, v_{n_k}^{(k)}) \sim \mathcal{F}$, and let $\rho, \epsilon_1, \dots, \epsilon_k : \mathbb{Z}_+ \rightarrow [0, 1]$, where each such function is either constant (and equal to 1) or tends to zero as n (and consequently the n_r) tend to infinity. Given matrices $\mathbf{\Lambda}_1, \dots, \mathbf{\Lambda}_k$, we define $\mathbf{\Lambda}_{\epsilon,r} = \epsilon_r \mathbf{\Lambda}_r$ for each r and $\mathbf{\Lambda}_\epsilon = [\mathbf{\Lambda}_{\epsilon,1} \dots \mathbf{\Lambda}_{\epsilon,k}]$, and set $\epsilon = (\epsilon_1^2 + \dots + \epsilon_k^2)^{1/2}$. Note that by submultiplicativity of the spectral norm, we have $\|\mathbf{\Lambda}_\epsilon\| \leq \epsilon \|\mathbf{\Lambda}\|$, and that we have the upper bound $\epsilon \leq k^{1/2}$.

Uniform scaling is then overlaid onto the model by defining latent positions $\mathbf{X}_i = \rho^{1/2}\xi_i$ and $\mathbf{Y}_j^{(r)} = \rho^{1/2}\nu_j^{(r)}$. We denote by \mathcal{F}_ρ this scaled distribution.

Throughout this paper, we will impose certain restrictions on the sparsity factors which are necessary for our results to hold. If *all* the graphs are too sparsely connected, then the UASE is unable to give us any useful information about the latent positions \mathbf{X}_i and $\mathbf{Y}_j^{(r)}$, and so we impose the global condition that:

- The sparsity factor ρ satisfies $\rho = \omega\left(\frac{\log^c(n)}{n^{1/2}}\right)$ for some constant c .

As mentioned previously, we require a sufficient number of graphs to be dense enough for us to be able to recover the latent positions \mathbf{X}_i ; if the only sufficiently dense graphs are degenerate, then it is entirely possible that we will be unable to do so. To avoid this, we impose the following local conditions (where we assume that we have reordered the matrices $\mathbf{A}^{(r)}$ appropriately):

- There exists an integer $1 \leq \kappa \leq k$ such that $\epsilon_r = 1$ for all $r \leq \kappa$ and $\epsilon_r \rightarrow 0$ for all $r > \kappa$, and that the matrix $\mathbf{\Lambda}_* = [\mathbf{\Lambda}_1 | \dots | \mathbf{\Lambda}_\kappa]$ is of rank d .

This condition guarantees the existence of a sufficiently dense subset of graphs that will allow us to recover the latent positions \mathbf{X}_i (although not necessarily the positions $\mathbf{Y}_j^{(r)}$). If each $\mathbf{\Lambda}_r$ is of rank d , then having $\kappa = 1$ is sufficient for our purposes.

2.2 Theoretical results

The main aim of this paper is to accurately describe the asymptotic behaviour of both sides of the UASE, which we do by establishing two key results. The first of these shows that the rows of the left embedding approximate some invertible linear transformation (of bounded spectral norm) of the latent positions \mathbf{X}_i , and that by inverting this transformation we obtain a good approximation to the latent positions themselves, in the sense that the maximum error vanishes. Similarly, we show that the rows of the r th right embedding approximate some (not necessarily invertible) linear transformation of the latent positions $\mathbf{Y}_i^{(r)}$, and that under appropriate conditions we can invert this transformation to obtain a good approximation to the latent positions themselves. This result is stated using the two-to-infinity norm [6] of the associated error matrix, which is the maximum Euclidean norm of any of its rows.

Theorem 2 (Two-to-infinity norm bound for the UASE).

Let $(\mathbf{A}, \mathbf{X}, \mathbf{Y}) \sim \text{MRDPG}(\mathcal{F}_\rho, \mathbf{\Lambda}_\epsilon)$ for a distribution \mathcal{F} and sparsity factors ρ and ϵ_r satisfying the criteria stated in Section 2.1. Then there exist sequences of matrices $\mathbf{L} = \mathbf{L}(n, n_1, \dots, n_k) \in \text{GL}(d)$ and $\mathbf{R}_r = \mathbf{R}_r(n, n_1, \dots, n_k) \in \mathbb{R}^{d_r \times d}$ for each r satisfying $\mathbf{L}\mathbf{R}_r^\top = \mathbf{\Lambda}_{\epsilon_r}$ such that

$$\|\mathbf{X}_\mathbf{A} - \mathbf{X}\mathbf{L}\|_{2 \rightarrow \infty} = O\left(\frac{\log^{1/2}(n)}{\rho^{1/2}n^{1/2}}\right), \quad \|\mathbf{Y}_\mathbf{A}^{(r)} - \mathbf{Y}^{(r)}\mathbf{R}_r\|_{2 \rightarrow \infty} = O\left(\frac{\log^{1/2}(n)}{\rho^{1/2}n^{1/2}}\right) \quad (2)$$

almost surely. Moreover, we may invert the matrices \mathbf{L} and (if $\mathbf{\Lambda}_r$ is of maximal rank) \mathbf{R}_r , and consequently find that

$$\|\mathbf{X}_\mathbf{A}\mathbf{L}^{-1} - \mathbf{X}\|_{2 \rightarrow \infty} = O\left(\frac{\epsilon^{1/2} \log^{1/2}(n)}{\rho^{1/2}n^{1/2}}\right), \quad \|\mathbf{Y}_\mathbf{A}^{(r)}\mathbf{R}_r^+ - \mathbf{Y}^{(r)}\|_{2 \rightarrow \infty} = O\left(\frac{\log^{1/2}(n)}{\epsilon_r^{1/2} \rho^{1/2}n^{1/2}}\right) \quad (3)$$

almost surely, where $\mathbf{R}_r^+ = \mathbf{R}_r^\top(\mathbf{R}_r\mathbf{R}_r^\top)^{-1}$ is the Moore-Penrose inverse of \mathbf{R}_r .

The matrices \mathbf{L} (and consequently \mathbf{R}_r) can be described explicitly, and are constructed by a two-step process; first using a modified Procrustes-style argument to simultaneously align $\mathbf{X}_\mathbf{A}$ with $\mathbf{X}_\mathbf{P}$ and $\mathbf{Y}_\mathbf{A}$ with $\mathbf{Y}_\mathbf{P}$ via an *orthogonal* transformation, and then applying a second *linear* transformation which maps $\mathbf{X}_\mathbf{P}$ directly to \mathbf{X} . For the first step, let $\mathbf{U}_\mathbf{P}^\top \mathbf{U}_\mathbf{A} + \mathbf{V}_\mathbf{P}^\top \mathbf{V}_\mathbf{A}$ admit the singular value decomposition

$$\mathbf{U}_\mathbf{P}^\top \mathbf{U}_\mathbf{A} + \mathbf{V}_\mathbf{P}^\top \mathbf{V}_\mathbf{A} = \mathbf{W}_1 \mathbf{\Sigma} \mathbf{W}_2^\top, \quad (4)$$

and let $\mathbf{W} = \mathbf{W}_1 \mathbf{W}_2^\top$. The matrix \mathbf{W} solves the one mode *orthogonal* Procrustes problem

$$\mathbf{W} = \arg \min_{\mathbf{Q} \in \mathbb{O}(d)} \|\mathbf{U}_\mathbf{A} - \mathbf{U}_\mathbf{P} \mathbf{Q}\|_F^2 + \|\mathbf{V}_\mathbf{A} - \mathbf{V}_\mathbf{P} \mathbf{Q}\|_F^2, \quad (5)$$

and we use \mathbf{W}^\top to align $\mathbf{X}_\mathbf{A}$ with $\mathbf{X}_\mathbf{P}$.

For the second step, we construct a matrix $\tilde{\mathbf{L}}$ which satisfies $\mathbf{X}_\mathbf{P} = \mathbf{X} \tilde{\mathbf{L}}$ and whose inverse can be shown to have spectral norm of order $O(\epsilon^{1/2})$ (see the Appendix for full details). The matrix \mathbf{L} is then given by $\mathbf{L} = \tilde{\mathbf{L}} \mathbf{W}$. A similar process is repeated for the right-sided embedding.

The second of our main results centres on the error distribution of the estimates for the latent positions \mathbf{X}_i and (if $\mathbf{\Lambda}_r$ is invertible) $\mathbf{Y}_j^{(r)}$ established in Theorem 2, and shows that conditional on the true position it is asymptotically Gaussian:

Theorem 3 (Central limit theorem for the UASE).

Let $(\mathbf{A}, \mathbf{X}, \mathbf{Y}) \sim \text{MRDPG}(\mathcal{F}_\rho, \mathbf{\Lambda}_\epsilon)$ for a distribution \mathcal{F} and sparsity factors ρ and ϵ_r satisfying the criteria stated in Section 2.1, and let \mathbf{L} and \mathbf{R}_r be the transformation matrices specified in Theorem 2.

Let $\xi \sim \mathcal{F}_X$ and $v_r \sim \mathcal{F}_{Y,r}$ for each $r \in \{1, \dots, \kappa\}$, where \mathcal{F}_X and $\mathcal{F}_{Y,r}$ are the marginal distributions of \mathcal{F} , and define $\Delta_Y = c_1 \Delta_{Y,1} \oplus \dots \oplus c_\kappa \Delta_{Y,\kappa}$, where the constants c_r are as stated in Section 2.1. Given $\mathbf{x} \in \mathcal{X}$, define

$$\Sigma_Y^{(r)}(\mathbf{x}) = \begin{cases} \mathbb{E}[\mathbf{x}^\top \mathbf{\Lambda}_r v_r (1 - \mathbf{x}^\top \mathbf{\Lambda}_r v_r) \cdot v_r v_r^\top] & \text{if } \rho = 1 \\ \mathbb{E}[\mathbf{x}^\top \mathbf{\Lambda}_r v_r \cdot v_r v_r^\top] & \text{if } \rho \rightarrow 0 \end{cases} \quad (6)$$

for each r , and let $\Sigma_Y = c_1 \Sigma_Y^{(1)} \oplus \dots \oplus c_\kappa \Sigma_Y^{(\kappa)}$. Then, for all $\mathbf{z} \in \mathbb{R}^d$ and for any fixed i ,

$$\mathbb{P}\left(n^{1/2}(\mathbf{X}_\mathbf{A} \mathbf{L}^{-1} - \mathbf{X})_i^\top \leq \mathbf{z} \mid \xi_i = \mathbf{x}\right) \rightarrow \Phi(\mathbf{z}, \Delta_{\mathbf{\Lambda}, Y}^{-1} \mathbf{\Lambda}_* \Sigma_Y(\mathbf{x}) \mathbf{\Lambda}_*^\top \Delta_{\mathbf{\Lambda}, Y}^{-1}) \quad (7)$$

almost surely, where $\Delta_{\mathbf{\Lambda}, Y} = \mathbf{\Lambda}_* \Delta_Y \mathbf{\Lambda}_*^\top$.

Moreover, for each $r \in \{1, \dots, \kappa\}$, if the matrix $\mathbf{\Lambda}_r$ is invertible, then given $\mathbf{y} \in \mathcal{Y}_r$ define

$$\Sigma_X^{(r)}(\mathbf{y}) = \begin{cases} \mathbb{E}[\xi^\top \mathbf{\Lambda}_r \mathbf{y} (1 - \xi^\top \mathbf{\Lambda}_r \mathbf{y}) \cdot \xi \xi^\top] & \text{if } \rho = 1 \\ \mathbb{E}[\xi^\top \mathbf{\Lambda}_r \mathbf{y} \cdot \xi \xi^\top] & \text{if } \rho \rightarrow 0 \end{cases} \quad (8)$$

Then, for all $\mathbf{z} \in \mathbb{R}^{d_r}$ and for any fixed i ,

$$\mathbb{P}\left(n^{1/2}(\mathbf{Y}_\mathbf{A}^{(r)} \mathbf{R}_r^{-1} - \mathbf{Y}^{(r)})_i^\top \leq \mathbf{z} \mid v_i^{(r)} = \mathbf{y}\right) \rightarrow \Phi(\mathbf{z}, \mathbf{\Lambda}_r^{-1} \Delta_X^{-1} \Sigma_X^{(r)}(\mathbf{y}) \Delta_X^{-1} \mathbf{\Lambda}_r^{-\top}) \quad (9)$$

almost surely.

The theorems, of which single-graph analogs were derived in [32, 4, 23, 8, 6, 33, 29, 5], also have analogous methodological implications. Under a multilayer stochastic block model, discussed in Section 3, the left UASE asymptotically follows a Gaussian mixture model with non-circular components. Fitting this model is preferable to using K -means, which is implicitly fitting circular components. Apart from shape considerations, note that (as with the GRDPG) latent positions under the MRDPG are only identifiable up to a distance-distorting transformation (here invertible linear, there indefinite orthogonal) to which partitions obtained using a Gaussian mixture model are invariant but those obtained using K -means are not. Consistency in the two-to-infinity norm should imply the consistency of many subsequent statistical analyses: usually, if a method is consistent, it remains so under a perturbation of the data of vanishing maximal error; one need only then worry about the effect of an unidentifiable linear transformation on conclusions; if the estimand is invariant, this effect will often vanish on account of the transformation having bounded spectral norm.

In the special case in which the graphs $\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(k)}$ are identically distributed, we can *always* recover the latent positions:

Corollary 4. Let $(\mathbf{A}, \mathbf{X}, \mathbf{Y}) \stackrel{\text{id}}{\sim} \text{MRDPG}(\mathcal{F}_\rho, \mathbf{\Lambda})$ for a distribution \mathcal{F} and sparsity factor ρ satisfying the criteria stated in Section 2.1. Then there exist sequences of matrices $\mathbf{L} = \mathbf{L}(n, n') \in \text{GL}(d)$ and $\mathbf{R} = \mathbf{R}(n, n') \in \mathbb{R}^{d' \times d}$ satisfying $\mathbf{L}\mathbf{R}^\top = \mathbf{\Lambda}$ such that

$$\|\mathbf{X}_\mathbf{A} - \mathbf{X}\mathbf{L}\|_{2 \rightarrow \infty} = O\left(\frac{\log^{1/2}(n)}{\rho^{1/2}n^{1/2}}\right), \quad \|\mathbf{Y}_\mathbf{A}^{(r)} - \mathbf{Y}\mathbf{R}\|_{2 \rightarrow \infty} = O\left(\frac{\log^{1/2}(n)}{\rho^{1/2}n^{1/2}}\right) \quad (10)$$

almost surely. Moreover, we may invert the matrices \mathbf{L} and (if $\mathbf{\Lambda}$ is of maximal rank) \mathbf{R} , and consequently find that

$$\|\mathbf{X}_\mathbf{A}\mathbf{L}^{-1} - \mathbf{X}\|_{2 \rightarrow \infty} = O\left(\frac{\log^{1/2}(n)}{\rho^{1/2}n^{1/2}}\right), \quad \|\mathbf{Y}_\mathbf{A}^{(r)}\mathbf{R}^+ - \mathbf{Y}\|_{2 \rightarrow \infty} = O\left(\frac{\log^{1/2}(n)}{\rho^{1/2}n^{1/2}}\right) \quad (11)$$

almost surely, where $\mathbf{R}^+ = \mathbf{R}^\top(\mathbf{R}\mathbf{R}^\top)^{-1}$ is the Moore-Penrose inverse of \mathbf{R} .

We can similarly derive a more precise version of the Central Limit Theorem of Theorem 3:

Corollary 5. Let $(\mathbf{A}, \mathbf{X}, \mathbf{Y}) \stackrel{\text{id}}{\sim} \text{MRDPG}(\mathcal{F}_\rho, \mathbf{\Lambda})$ for a distribution \mathcal{F} and sparsity factor ρ satisfying the criteria stated in Section 2.1, and let \mathbf{L} be the transformation matrix specified in Corollary 4.

Let $\xi \sim \mathcal{F}_X$ and $v \sim \mathcal{F}_Y$, where \mathcal{F}_X and \mathcal{F}_Y are the marginal distributions of \mathcal{F} . Given $\mathbf{x} \in \mathcal{X}$, define

$$\Sigma_Y(\mathbf{x}) = \begin{cases} \mathbb{E}[\mathbf{x}^\top \mathbf{\Lambda} v (1 - \mathbf{x}^\top \mathbf{\Lambda} v) \cdot v v^\top] & \text{if } \rho = 1 \\ \mathbb{E}[\mathbf{x}^\top \mathbf{\Lambda} v \cdot v v^\top] & \text{if } \rho \rightarrow 0 \end{cases} \quad (12)$$

Then, for all $\mathbf{z} \in \mathbb{R}^d$ and for any fixed i ,

$$\mathbb{P}\left(n^{1/2}(\mathbf{X}_\mathbf{A}\mathbf{L}^{-1} - \mathbf{X})_i^\top \leq \mathbf{z} \mid \xi_i = \mathbf{x}\right) \rightarrow \Phi\left(\mathbf{z}, \frac{c}{k} \mathbf{\Lambda}^{-\top} \Delta_Y^{-1} \Sigma_Y(\mathbf{x}) \Delta_Y^{-1} \mathbf{\Lambda}^{-1}\right) \quad (13)$$

almost surely, where $\lim_{n, n' \rightarrow \infty} \frac{n'}{n} = c$.

Moreover, if the matrix $\mathbf{\Lambda}_r$ is invertible, then given $\mathbf{y} \in \mathcal{Y}$ define

$$\Sigma_X(\mathbf{y}) = \begin{cases} \mathbb{E}[\xi^\top \mathbf{\Lambda} \mathbf{y} (1 - \xi^\top \mathbf{\Lambda} \mathbf{y}) \cdot \xi \xi^\top] & \text{if } \rho = 1 \\ \mathbb{E}[\xi^\top \mathbf{\Lambda} \mathbf{y} \cdot \xi \xi^\top] & \text{if } \rho \rightarrow 0 \end{cases} \quad (14)$$

Then, for all $\mathbf{z} \in \mathbb{R}^{d'}$ and for any fixed i ,

$$\mathbb{P}\left(n^{1/2}(\mathbf{Y}_\mathbf{A}^{(r)}\mathbf{R}^{-1} - \mathbf{Y})_i^\top \leq \mathbf{z} \mid v_i = \mathbf{y}\right) \rightarrow \Phi\left(\mathbf{z}, \mathbf{\Lambda}^{-1} \Delta_X^{-1} \Sigma_X(\mathbf{y}) \Delta_X^{-1} \mathbf{\Lambda}^{-\top}\right) \quad (15)$$

almost surely.

If $\mathbf{Y} = \mathbf{X}$ with probability one and $\mathbf{\Lambda} = \mathbf{I}_{p,q}$ (the diagonal matrix whose first p entries are equal to 1 and remaining q entries are equal to -1) then the limiting distribution for the rows UASE are the same as that for the ASE stated in [29], scaled by a factor of $\frac{1}{k}$, and in particular coincides with that obtained by spectrally embedding the average $\bar{\mathbf{A}} = \frac{1}{k} \sum_{r=1}^k \mathbf{A}^{(r)}$ of the adjacency matrices (see for example [35]).

In Corollaries 4 and 5 the matrices \mathbf{R} are common across graphs, allowing a direct comparison of their right-sided embeddings. By contrast, independently embedded graphs first need to be aligned before a meaningful comparison is possible. In [34] this is achieved using Procrustes, the appropriate method of alignment under an RDPG model where latent positions are identifiable only up to orthogonal transformation. Under the GRDPG, finding a best indefinite orthogonal alignment is less straightforward. Computational issues aside, alignment comes at a statistical cost [21] when testing whether two point clouds differ statistically, and the empirical performance of using right-sided unfolded adjacency spectral embeddings for two-graph testing are investigated in Section 4.4.

3 The generalised multilayer stochastic block model

Definition 3. (The generalised multilayer stochastic block model).

We say that the matrices $\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(k)}$ are distributed as a \mathbf{K} -community generalised multilayer stochastic block model for a tuple $\mathbf{K} = (K, K_1, \dots, K_k)$ of positive integers if $(\mathbf{A}, \mathbf{X}, \mathbf{Y}) \sim \text{MRDPG}(\mathcal{F}, \mathbf{B})$ for a set of matrices $\mathbf{B}^{(r)} \in [0, 1]^{K \times K_r}$ and a distribution \mathcal{F} whose marginal distributions \mathcal{F}_X and \mathcal{F}_{Y_r} are supported on the sets $\{\mathbf{e}_1, \dots, \mathbf{e}_K\}$ and $\{\mathbf{e}_1, \dots, \mathbf{e}_{K_r}\}$ of standard basis vectors of \mathbb{R}^K and \mathbb{R}^{K_r} respectively. If this is the case we write $(\mathbf{A}, \mathbf{X}, \mathbf{Y}) \sim \text{GMSBM}(\mathcal{F}, \mathbf{B})$.

If each of the $\mathbf{Y}^{(r)}$ is equal to \mathbf{X} with probability one (and consequently $K_r = K$ for each r), the distribution \mathcal{F} assigns each latent position to the i th basis vector of \mathbb{R}^K with probability π_i for some tuple $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K)$ whose entries sum to 1, and each of the matrices $\mathbf{A}^{(r)}$ and $\mathbf{B}^{(r)}$ is symmetric, then the GMSBM reduces to the standard K -community multilayer SBM [16]. In this case, the matrices $\mathbf{A}^{(r)}$ are the adjacency matrices of a set of graphs generated independently from a common set of vertices, in which, independently conditional on a partition of these vertices into K disjoint communities, an edge is generated between the i th and j th vertices in the r th graph with probability $\mathbf{B}_{z_i, z_j}^{(r)}$, where $z_i \in \{1, \dots, K\}$ denotes the community membership of the i th vertex.

Note that, as per Proposition 1, we may consider alternative choices of latent positions (and thus matrices \mathbf{A}_r) for the GMSBM. We often use the second choice posited in Proposition 1, in which we consider the singular value decomposition $\mathbf{B} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^\top$ with $\mathbf{U} \in O(K)$ and $\mathbf{V} \in O((K_1 + \dots + K_k) \times K)$, and let the positions \mathbf{X}_i be chosen from the rows of $\mathbf{U}\boldsymbol{\Sigma}^{1/2}$ and $\mathbf{Y}_j^{(r)}$ be chosen from the rows of $\mathbf{V}_r\boldsymbol{\Sigma}$ (where we split \mathbf{V} into k distinct blocks $\mathbf{V}_r \in \mathbb{R}^{K_r \times K}$). In this case, each matrix \mathbf{A}_r is equal to the identity matrix. If $k = 1$, this choice closely resembles the model presented in [29], with the signs of the eigenvalues of the matrix \mathbf{B} being absorbed into the latent positions \mathbf{Y}_j (if we impose the additional condition that \mathbf{Y} is equal to \mathbf{X} with probability one, then $\mathbf{A} = \mathbf{I}_{p,q}$, exactly as in [29]).

3.1 Undirected graphs

As a first demonstration of the UASE for the GMSBM, we consider the case in which the latent positions $\mathbf{Y}^{(r)} = \mathbf{X}$, and the matrices $\mathbf{B}^{(r)}$ and $\mathbf{A}^{(r)}$ are symmetric (that is, the standard multilayer SBM). We begin with two examples. For the first, we assume that the adjacency matrices $\mathbf{A}^{(1)}$ and $\mathbf{A}^{(2)}$ are *identically* distributed according to a multilayer SBM with parameters

$$\mathbf{B} = \begin{pmatrix} 0.42 & 0.42 \\ 0.42 & 0.5 \end{pmatrix}, \quad \boldsymbol{\pi} = (0.6, 0.4) \quad (16)$$

(the Laplacian spectral embedding of a single graph generated with these parameters was studied in [33]). Figure 1 plots the estimated latent positions for the ASE of the matrix $\mathbf{A}^{(1)}$ (first row) and the UASE of the matrix $\mathbf{A} = [\mathbf{A}^{(1)} | \mathbf{A}^{(2)}]$ (second row) for $n = 1000, 2000$ and 4000 . Also displayed are the 95% level curves of the empirical distributions (dashed curves) and the theoretical distributions specified by Theorem 3 (solid curves).

For the second example, we take a pair of graphs with adjacency matrices $\mathbf{A}^{(1)}$ and $\mathbf{A}^{(2)}$ generated according to a multilayer SBM with parameters

$$\mathbf{B}^{(1)} = \begin{pmatrix} 0.58 & 0.58 \\ 0.58 & 0.5 \end{pmatrix}, \quad \mathbf{B}^{(2)} = \begin{pmatrix} 0.42 & 0.42 \\ 0.42 & 0.5 \end{pmatrix}, \quad \boldsymbol{\pi} = (0.6, 0.4). \quad (17)$$

Since the matrices $\mathbf{B}^{(1)}$ and $\mathbf{B}^{(2)}$ have signatures $(2, 0)$ and $(1, 1)$ respectively, they exhibit different assortativity behaviours. Figure 2 plots the estimated latent positions for the ASE of the matrix $\mathbf{A}^{(1)}$ (first row) and the UASE (second row) for $n = 1000, 2000$ and 4000 , with the 95% level curves displayed as in Figure 1.

In each example, the UASE demonstrates greater cluster separation over the ASE. To test empirically whether this behaviour holds in general, we performed the following experiment: for each value of $n \in \{50, 100, 250, 500, 750, 1000, 1500, 2000\}$ we performed 500 trials in which we generate two matrices $\mathbf{B}^{(r)}$ with entries $\mathbf{B}_{ij}^{(r)} \sim \text{Uniform}[0, 1]$, and probability vectors $\boldsymbol{\pi} \sim \text{Dirichlet}(1, 1)$

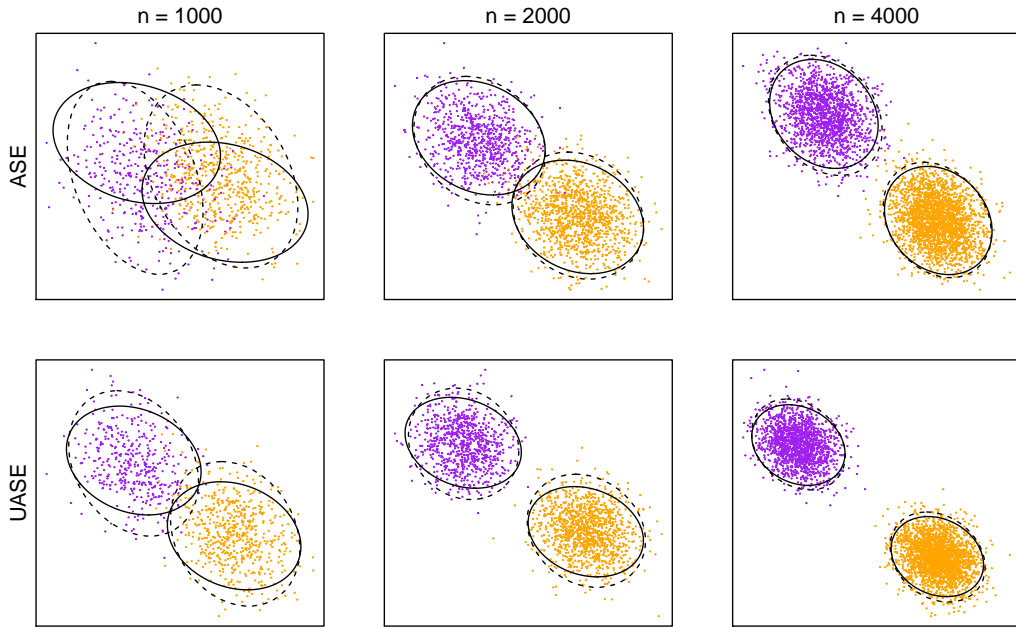


Figure 1: Plots of the latent position estimates by the ASE and UASE of a pair of identically distributed graphs drawn from a 2-community SBM on the same set of n nodes. Points are coloured according to the community membership of the corresponding vertices. Ellipses give the 95% level curves of the *empirical* (dashed curves) and *theoretical* (solid curves) distributions specified by Theorem 3.

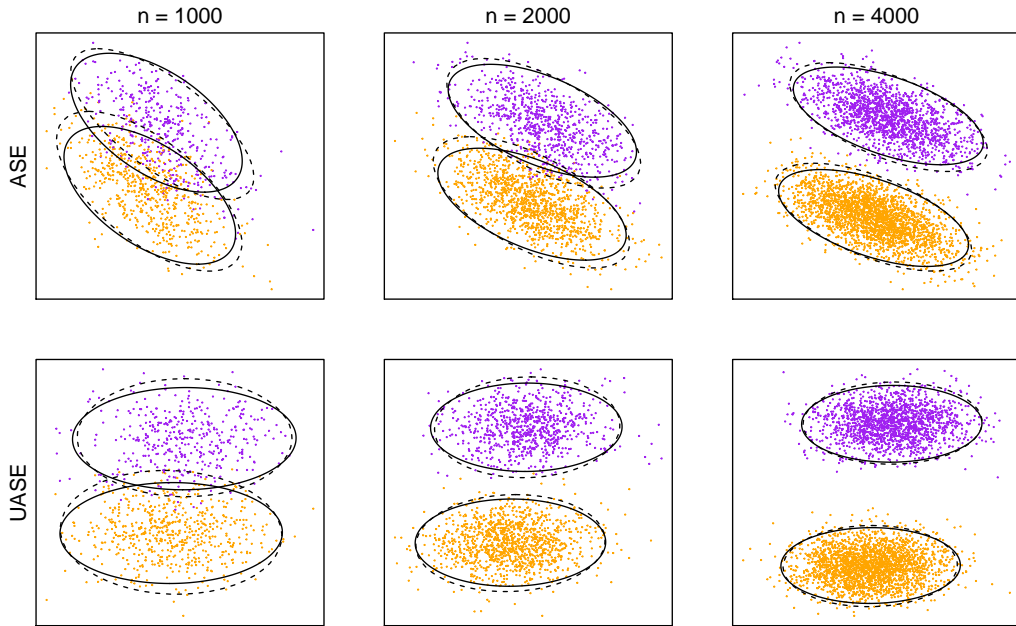


Figure 2: Plots of the latent position estimates by the ASE and UASE of a pair of graphs drawn from a 2-community SBM on the same set of n nodes with differing distributions. Points are coloured according to the community membership of the corresponding vertices. Ellipses give the 95% level curves of the *empirical* (dashed curves) and *theoretical* (solid curves) distributions specified by Theorem 3.

(to ensure that both clusters were of a reasonable size, we discard vectors π for which either of the π_i was less than 0.2). Matrices $\mathbf{A}^{(r)}$ were then generated according to the resulting multilayer SBM, the UASE calculated, and nodes were then assigned to the most likely cluster predicted by the Gaussian mixture model obtained via the MCLUST algorithm (see [31]). These were then compared against the known cluster assignments given by the latent positions \mathbf{X}_i , and the average

classification error rate calculated across all samples of a given size n (for the ASE, the average error rate of the two embeddings was used). For added differentiation, we performed this test for 3 separate cases: one in which the matrices $\mathbf{B}^{(r)}$ were identical; one in which they had a common signature; and one in which they had different signatures. Figure 3 displays these error rates in the case of the identical and mixed parameter examples, as well as the average error rates for the *mean embedding*—that is, the spectral embedding of the matrix $\bar{\mathbf{A}} = \frac{1}{2}(\mathbf{A}^{(1)} + \mathbf{A}^{(2)})$, which we denote by ASE(mean)—and, in the identically distributed case, the *omnibus embedding*—denoted by OMNI (see [21]).

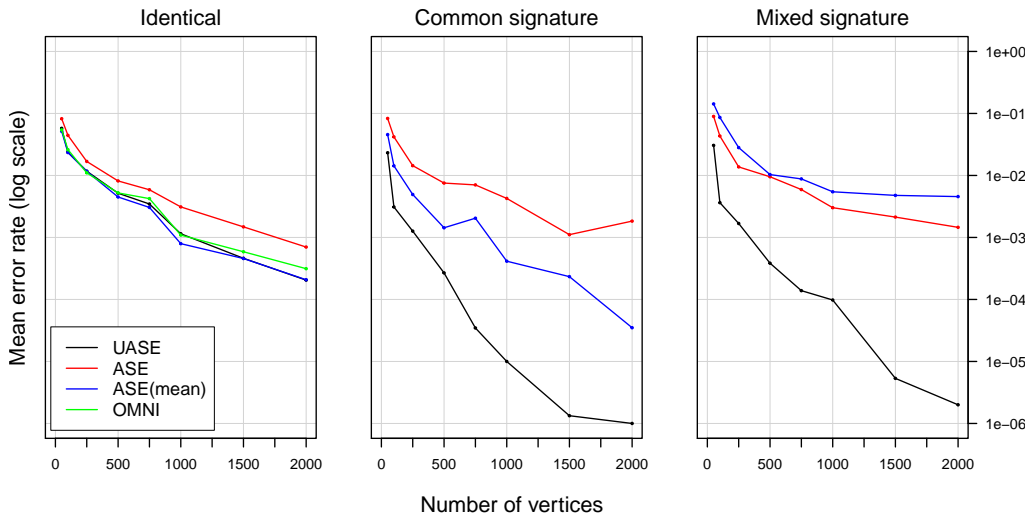


Figure 3: Classification error rates for assignment of nodes to the most likely cluster predicted by the Gaussian mixture model obtained via the MCLUST algorithm. Error rates are plotted on a logarithmic scale. See main text for details.

As expected, the UASE (black line) clearly outperforms the ASE (red line) in all cases. When the two graphs are identically distributed, the UASE is comparable with the mean embedding (blue line) and slightly outperforms the omnibus embedding (green line). If the matrices $\mathbf{B}^{(r)}$ differ then the UASE significantly outperforms the mean embedding; in particular, if the matrices $\mathbf{B}^{(r)}$ have different signatures (resulting in different assortativity behaviours in the graphs $\mathbf{A}^{(r)}$) then the mean embedding performs worse than even the ASE. This is not particularly surprising; if the adjacency matrices of different graphs have different signatures, then it is entirely possible for the matrix $\bar{\mathbf{P}}$ to have non-maximal rank, causing some of the information in the system to be lost when we spectrally embed the matrix $\bar{\mathbf{A}}$. Conversely, one finds that the embedding $\mathbf{X}_{\bar{\mathbf{A}}}$ is the same as that of the positive-definite square root of the matrix $\sum_{r=1}^k (\mathbf{A}^{(r)})^2$, which will *always* be of maximal rank.

3.2 Directed graphs

Unlike the standard multilayer SBM, the GMSBM does not require that the matrices $\mathbf{B}^{(r)}$ or $\mathbf{A}^{(r)}$ be symmetric, and so it allows us to consider *directed* graphs. In particular, given a *single* directed graph whose adjacency matrix $\mathbf{A} \in \{0, 1\}^{n \times n}$ follows a GMSBM, Theorems 2 and 3 show us that taking the standard ASE of \mathbf{A} will provide us with consistent estimates of the latent positions \mathbf{X} .

However, the ability of the UASE to evaluate multiple adjacency matrices presents an alternative method for working with directed graphs. Let $(\mathbf{A}, \mathbf{X}, \mathbf{X}) \sim \text{GMSBM}(\mathcal{F}, \mathbf{B})$ for some \mathcal{F} and \mathbf{B} , and let $\text{Sym}(\mathbf{A})$ be the hollow symmetric matrix whose upper-triangular part is equal to that of \mathbf{A} . Then we observe that $(\text{Sym}(\mathbf{A}), \mathbf{X}, \mathbf{X}) \sim \text{GMSBM}(\mathcal{F}, \mathbf{B})$, where we view $\text{Sym}(\mathbf{A})$ as the adjacency matrix of an *undirected* graph drawn on the same set of nodes as our original graph. Similarly, we see that $(\text{Sym}(\mathbf{A}^\top), \mathbf{X}, \mathbf{X}) \sim \text{GMSBM}(\mathcal{F}, \mathbf{B}^\top)$, and thus $(\tilde{\mathbf{A}}, \mathbf{X}, \mathbf{X} \oplus \mathbf{X}) \sim \text{GMSBM}(\tilde{\mathcal{F}}, \tilde{\mathbf{B}})$, where $\tilde{\mathbf{A}} = [\text{Sym}(\mathbf{A}) | \text{Sym}(\mathbf{A}^\top)]$, $\tilde{\mathbf{B}} = [\mathbf{B} | \mathbf{B}^\top]$ and $\tilde{\mathcal{F}}$ is the natural extension of the distribution \mathcal{F} .

We demonstrate the effectiveness of this proposed method through the following experiment: for each value of $n \in \{50, 100, 250, 500, 750, 1000, 1500, 2000\}$, we performed 1000 trials in which a directed graph was generated according to a 2-community stochastic block model, where the

(not necessarily symmetric) probability matrix \mathbf{B} and community probabilities $\boldsymbol{\pi}$ were randomly generated as before. The adjacency matrix \mathbf{A} was then constructed, and the ASE of \mathbf{A} , the UASE of $\tilde{\mathbf{A}}$, and the mean embedding we calculated. Using each of these, we assign nodes to their most likely cluster using the MCLUST algorithm as before, and calculate the mean error rate across all the trials of a given sample size. The results (in the form of the average classification error rate) are plotted in Figure 4, and indicate that on average using the UASE offers significant improvement over the ASE, and a minor improvement over the mean embedding (we note however, that this does not guarantee that the UASE will *always* offer the best performance).

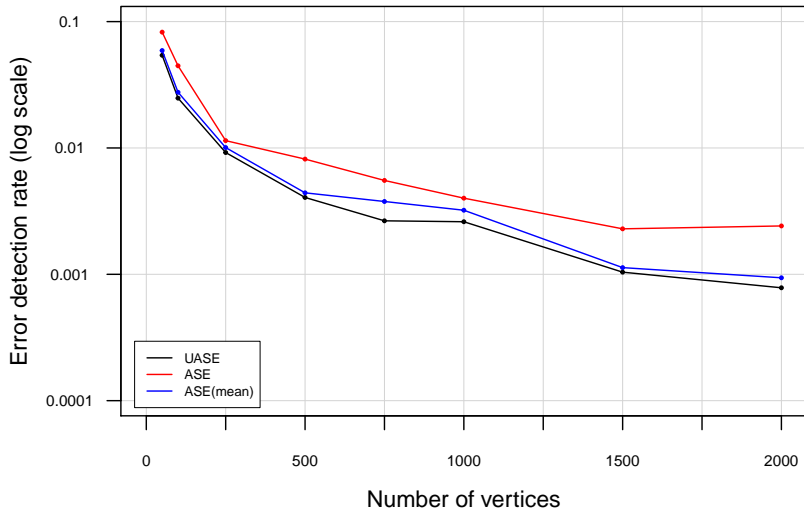


Figure 4: Classification error rates for assignment of nodes from a directed graph to the most likely cluster predicted by the Gaussian mixture model obtained via the MCLUST algorithm. Error rates are plotted on a logarithmic scale. See main text for details.

w

3.3 Bipartite multilayer SBMs

The flexibility of the GMSBM means that we do not have to restrict our attention to a standard multiple graph embedding; we may restrict our attention to submatrices of the adjacency matrices $\mathbf{A}^{(r)}$ for each graph in order to focus only on the interactions involving a given set of nodes. As an illustrative example, consider the “bipartite” situation in which we have a GMSBM with two sets of latent positions $\mathbf{X}_i \in \mathbb{R}^2$ and $\mathbf{Y}_j \in \mathbb{R}^3$, with probability matrices

$$\mathbf{B}^{(1)} = \begin{pmatrix} 0.45 & 0.45 & 0.52 \\ 0.51 & 0.51 & 0.54 \end{pmatrix}, \quad \mathbf{B}^{(2)} = \begin{pmatrix} 0.46 & 0.51 & 0.43 \\ 0.42 & 0.47 & 0.52 \end{pmatrix} \quad (18)$$

and community membership probabilities $\boldsymbol{\pi}_X = (1/2, 1/2)$ and $\boldsymbol{\pi}_Y = (1/3, 1/3, 1/3)$.

Figure 5 plots the left embedding \mathbf{X}_A and the right embeddings $\mathbf{Y}_A^{(r)}$ for the pairs $(n, n') = (1000, 1500)$, $(2000, 3000)$ and $(4000, 6000)$, where we use the latent positions $\tilde{\mathbf{X}}, \tilde{\mathbf{Y}}^{(r)} \in \mathbb{R}^2$ posited in Proposition 1. This demonstrates an important point: while Theorems 2 and 3 guarantee that the embeddings $\mathbf{Y}_A^{(r)}$ provide consistent estimates of the latent positions $\tilde{\mathbf{Y}}^{(r)}$, we cannot guarantee that they will distinguish between different communities, as the rank of $\tilde{\mathbf{Y}}^{(r)}$ is equal to the rank of $\mathbf{B}^{(r)}$, which may be less than the rank of $\text{rank}(\mathbf{Y})$. In our example, the embedding $\mathbf{Y}_A^{(1)}$ fails to distinguish between all three communities, while $\mathbf{Y}_A^{(2)}$ does distinguish between them. In general, if the columns of the matrix $\mathbf{B}^{(r)}$ are distinct, then the latent positions corresponding to different communities will be distinct.

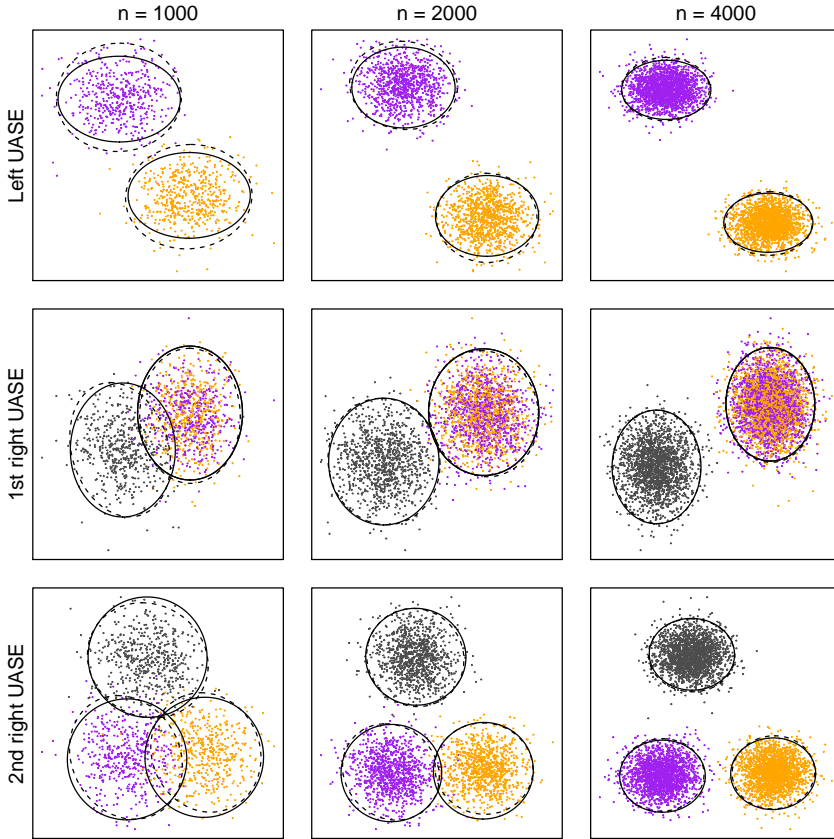


Figure 5: Plots of the latent position estimates by the UASE of a pair of graphs drawn from a $(2, 3)$ -community GMSBM. Points are coloured according to the community membership of the corresponding vertices. Ellipses give the 95% level curves of the *empirical* (dashed curves) and *theoretical* (solid curves) distributions specified by Theorem 3.

3.4 Rank considerations

One advantage of studying the MRDPG is that we do *not* require the matrices $\mathbf{\Lambda}_r$ to have maximal rank; this can lead to situations in which information about latent positions is obscured in individual graphs, but becomes apparent when considering the joint embedding. As an example, consider a dynamic network which can be modeled as a two graph 3-community multilayer SBM, in which the matrices $\mathbf{B}^{(r)}$ of probabilities take the form

$$\mathbf{B}^{(1)} = \begin{pmatrix} p_1 & p_1 & q_1 \\ p_1 & p_1 & q_1 \\ q_1 & q_1 & r_1 \end{pmatrix}, \quad \mathbf{B}^{(2)} = \begin{pmatrix} p_2 & q_2 & q_2 \\ q_2 & r_2 & r_2 \\ q_2 & r_2 & r_2 \end{pmatrix} \quad (19)$$

for values $p_i, q_i, r_i \in [0, 1]$. We could view this as a simple model for time-dependent snapshots of the communication preferences between two departments in a company, in which a third team moves from the first department to the second in between snapshots, and inherits the communication preferences of the department to which they are assigned at the time. The matrices $\mathbf{B}^{(r)}$ are both of non-maximal rank, but $\mathbf{B} = [\mathbf{B}_1 | \mathbf{B}_2]$ has maximal rank, provided that the p_i, q_i and r_i are distinct.

We demonstrate this with an example. Let $n = 4000$ and suppose that we have three communities, containing 1750, 500 and 1750 nodes respectively. Let the probability matrices be given by

$$\mathbf{B}^{(1)} = \begin{pmatrix} 0.47 & 0.47 & 0.39 \\ 0.47 & 0.47 & 0.39 \\ 0.39 & 0.39 & 0.56 \end{pmatrix}, \quad \mathbf{B}^{(2)} = \begin{pmatrix} 0.53 & 0.61 & 0.61 \\ 0.61 & 0.44 & 0.44 \\ 0.61 & 0.44 & 0.44 \end{pmatrix}. \quad (20)$$

Figure 6 shows the embedded point clouds generated by the individual ASEs and the UASE in

this situation. As one would expect, the individual ASEs display only the two communities that one observes at that given snapshot in time, while the UASE clearly displays all three communities..

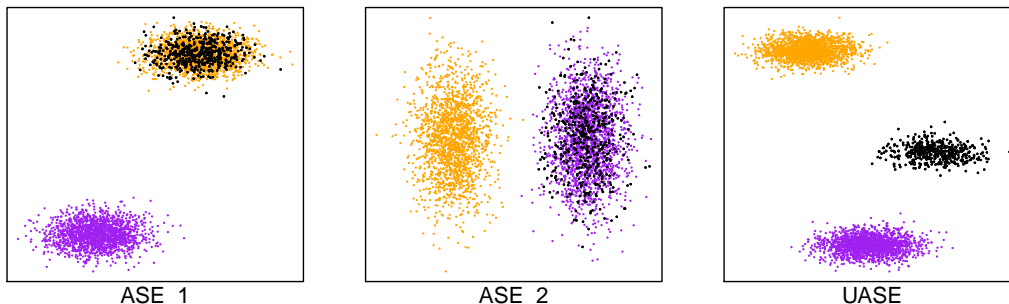


Figure 6: Plots of the ASEs of the adjacency matrices \mathbf{A}_1 and \mathbf{A}_2 and the UASE, with nodes coloured according to community membership. Points are coloured according to the community to which their corresponding nodes belong, with points in black representing those nodes that switch between communities.

4 Multiple graph inference: comparison with existing methods

The performance of unfolded adjacency spectral embedding is now compared with alternative spectral approaches on the inference tasks of latent position recovery, subspace estimation, model estimation, and two-graph hypothesis testing. We restrict our attention to the case in which the matrices $\mathbf{A}^{(r)}$ are true adjacency matrices of graphs, and thus work under the assumption that the latent position matrices $\mathbf{Y}^{(r)}$ are equal to \mathbf{X} with probability one for all r .

4.1 Recovery of latent positions

An important estimation problem for the data of a random dot product graph is that of estimating the latent positions \mathbf{X}_i , and so we shall investigate the performance of the MRDPG in the context of such an estimation problem. For comparison, we will consider the *multiple adjacency spectral embedding* (MASE) [3], which is an alternative method of jointly embedding adjacency matrices which follow a model that is essentially identical to the MRDPG, known as the *common subspace independent edge graph model*. In [3], the authors demonstrate that the MASE yields state-of-the-art performance on subsequent inference tasks, ahead of other competing models for studying multiple graph embeddings such as the multi-RDPG [26] and MREG [39] models, making it an ideal method to compare the UASE against.

Definition 4. (Common Subspace Independent Edge graphs).

Let $\mathbf{U} = [\mathbf{U}_1 | \dots | \mathbf{U}_n]^\top \in \mathbb{R}^{n \times d}$ have orthonormal columns, and let $\mathbf{R}^{(1)}, \dots, \mathbf{R}^{(k)} \in \mathbb{R}^{d \times d}$ be symmetric matrices such that, for each $r \in \{1, \dots, k\}$, $\mathbf{U}_i^\top \mathbf{R}^{(r)} \mathbf{U}_j \in [0, 1]$ for all $i, j \in \{1, \dots, n\}$. The random adjacency matrices $\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(k)}$ are said to be jointly distributed according to the common subspace independent-edge graph model with bounded rank d and parameters \mathbf{U} and $\mathbf{R}^{(1)}, \dots, \mathbf{R}^{(k)}$ if for each $r \in \{1, \dots, k\}$, conditional upon \mathbf{U} and $\mathbf{R}^{(r)}$ we have $\mathbf{A}_{ij}^{(r)} \sim \text{Bern}(\mathbf{P}_{ij}^{(r)})$, where $\mathbf{P}^{(r)} = \mathbf{U} \mathbf{R}^{(r)} \mathbf{U}^\top$, in which case we write $(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(k)}) \sim \text{COSIE}(\mathbf{U}; \mathbf{R}^{(1)}, \dots, \mathbf{R}^{(k)})$.

For all intents and purposes, the COSIE and MRDPG models are equivalent. Any COSIE model gives rise to a MRDPG by simply setting the latent positions \mathbf{X}_i to be equal to the rows \mathbf{U}_i , and the matrices $\mathbf{A}_r = \mathbf{R}^{(r)}$ for each r . Conversely, given a MRDPG such that the matrix \mathbf{X} of latent positions is of rank d , we can define $\mathbf{U} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1/2}$, where we have taken the positive-definite matrix square root of the matrix $\mathbf{X}^\top \mathbf{X}$. It is clear that the columns of \mathbf{U} are orthonormal, and we obtain a COSIE model by setting $\mathbf{R}^{(r)} = (\mathbf{X}^\top \mathbf{X})^{1/2} \mathbf{A}_r (\mathbf{X}^\top \mathbf{X})^{1/2}$. In both cases the two definitions of the matrices $\mathbf{P}^{(r)}$ coincide.

Definition 5. (Multiple adjacency spectral embedding).

Let $(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(m)}) \sim \text{COSIE}(\mathbf{U}; \mathbf{R}^{(1)}, \dots, \mathbf{R}^{(k)})$. For each $r \in \{1, \dots, k\}$ let d_r denote the rank

of $\mathbf{R}^{(r)}$, let $\mathbf{X}_{\mathbf{A}^{(r)}} \in \mathbb{R}^{n \times d_r}$ be the adjacency spectral embedding of $\mathbf{A}^{(r)}$ and define the matrix of concatenated spectral embeddings $\mathbf{M}_{\mathbf{A}} = [\mathbf{X}_{\mathbf{A}^{(1)}} | \dots | \mathbf{X}_{\mathbf{A}^{(k)}}]$. The multiple adjacency spectral embedding of $\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(k)}$ is the matrix $\widehat{\mathbf{U}}_{\mathbf{A}} \in \mathbb{R}^{n \times d}$ containing the d leading left singular vectors of $\mathbf{M}_{\mathbf{A}}$.

We note that it is possible to use the MASE to produce estimates for the latent positions in a MRDPG in an analogous way to the method used for the UASE. Let $\widehat{\mathbf{X}}_{\mathbf{A}} = \widehat{\mathbf{U}}_{\mathbf{A}} \boldsymbol{\Sigma}_{\mathbf{A}}$, where $\boldsymbol{\Sigma}_{\mathbf{A}} \in \mathbb{R}^{d \times d}$ is the diagonal matrix of the leading d singular values of $\mathbf{M}_{\mathbf{A}}$, and define $\mathbf{M}_{\mathbf{P}}, \widehat{\mathbf{U}}_{\mathbf{P}}$ and $\widehat{\mathbf{X}}_{\mathbf{P}}$ analogously for the matrices $\mathbf{P}^{(r)}$. Firstly, we note that adding columns of zeros to any of the $\mathbf{X}_{\mathbf{P}^{(r)}}$ will not alter $\widehat{\mathbf{X}}_{\mathbf{P}}$, and so we may assume without loss of generality that the matrices $\mathbf{X}_{\mathbf{P}^{(r)}} \in \mathbb{R}^{n \times d}$, and thus that there exist matrices $\mathbf{L}^{(r)} \in \mathbb{R}^{d \times d}$ of rank d_r such that $\mathbf{X}_{\mathbf{P}^{(r)}} = \mathbf{X} \mathbf{L}^{(r)}$. One can prove the existence of a matrix $\mathbf{L} \in \text{GL}(d)$ such that $\widehat{\mathbf{X}}_{\mathbf{P}} = \mathbf{X} \mathbf{L}$, and so performing a Procrustes-style alignment between $\widehat{\mathbf{X}}_{\mathbf{A}}$ and $\widehat{\mathbf{X}}_{\mathbf{P}}$ and multiplying by \mathbf{L}^{-1} produces a set of points that in practice are a good approximation to the latent positions \mathbf{X}_i .

We first tested the performance of the two embeddings on graphs of different sizes by performing, for each value of $n \in \{10, 25, 50, 75, 100, 250, 500, 750, 1000\}$, 1000 independent trials in which the latent positions \mathbf{X}_i are drawn i.i.d. from a Dirichlet distribution with parameter $(1, 1, 1)^\top \in \mathbb{R}^3$. In each trial, we generate two graphs $\mathbf{A}^{(1)}, \mathbf{A}^{(2)} \in \mathbb{R}^{n \times n}$, where $\mathbf{A}_{ij}^{(r)} \sim \text{Bern}(\mathbf{X}_i^\top \boldsymbol{\Lambda}_r \mathbf{X}_j)$ for $i < j$, where each $\boldsymbol{\Lambda}_r$ is a randomly chosen matrix. We then calculate estimates $\widehat{\mathbf{X}}$ using the UASE and MASE as described previously, and compare the average mean squared error $\frac{1}{n} \sum_{i=1}^n |\widehat{\mathbf{X}}_i - \mathbf{X}_i|^2$ of the two embeddings across the 1000 trials.

We then investigated the effect of changing the number of graphs to be embedded on the accuracy of each embedding type. Fixing $n = 750$, we again performed 1000 independent trials as above for $m = 2, \dots, 10$ embeddings, using the same procedure for generating the latent positions and adjacency matrices, and again compared the average mean squared error between the estimated and actual latent positions.

Figure 7 plots the results of the two experiments. While the MASE outperforms the UASE for values of $n < 75$ (with the joint embedding performing significantly worse for $n < 50$) the UASE clearly demonstrates superior accuracy as the size of the graph grows, a trend which continues as we increase the number of graphs to be embedded.

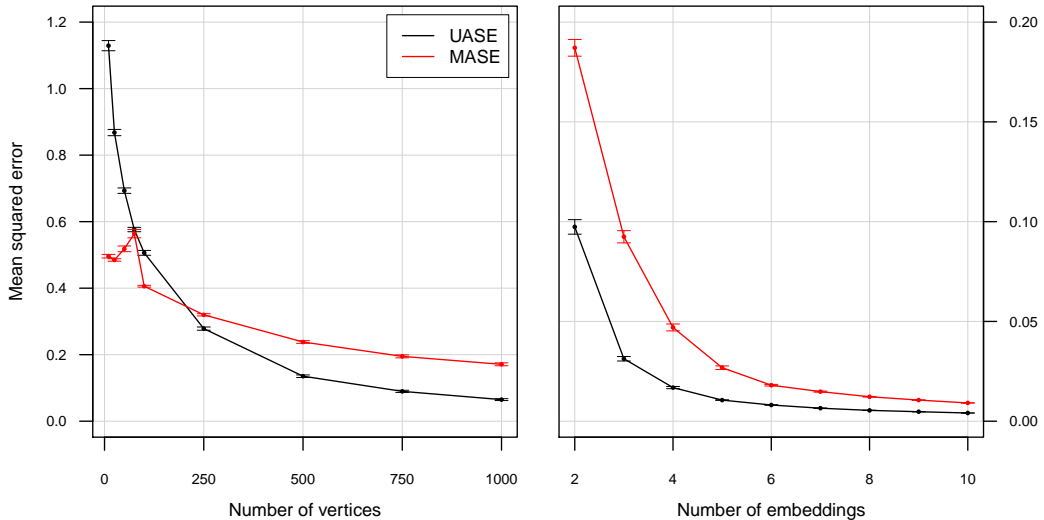


Figure 7: Mean squared error in recovery of latent positions in a 2-graph MRDPG model as a function of the number of vertices (left-hand graph) and the number of embeddings (right-hand graph).

4.2 Estimation of invariant subspaces

We next investigate the performance of UASE at estimating the invariant subspace \mathbf{U} in the COSIE model. We do this by setting the matrix \mathbf{X} of latent positions in the MRDPG to be equal to \mathbf{U} ,

and considering the *unscaled* UASE, \mathbf{U}_A . Unlike the scaled embedding, which approximates the latent positions only up to linear transformation, the unscaled UASE approximates the invariant subspace \mathbf{U} up to *orthogonal* transformation. Indeed, from our results for the scaled embedding the matrix $\mathbf{U}_P = \mathbf{U}\mathbf{Q}_X$ for some $\mathbf{Q}_X \in \text{GL}(d)$, whence the requirement that both \mathbf{U} and \mathbf{U}_P have orthonormal columns forces \mathbf{Q}_X to in fact belong to $\mathcal{O}(d)$, while the transformation applied in the Procrustes alignment between \mathbf{U}_A and \mathbf{U}_P is by definition orthogonal.

We can measure the distance between the estimate \mathbf{U}_A and the true invariant subspace \mathbf{U} using the spectral norm of the difference between the projections $\|\mathbf{U}_A\mathbf{U}_A^\top - \mathbf{U}\mathbf{U}^\top\|$ (and similarly for the estimate $\hat{\mathbf{U}}$ produced by the MASE). This distance is zero only when there exists an orthogonal matrix $\mathbf{W} \in \mathcal{O}(d)$ such that $\mathbf{U}_A = \mathbf{U}\mathbf{W}$ (respectively $\hat{\mathbf{U}} = \mathbf{U}\mathbf{W}$).

As in the previous example, we investigated the effect of changing both the size of the graphs and the number of graphs to be embedded on the performance of the UASE and MASE. Again, we began by performing 1000 independent trials for each value of $n \in \{10, 25, 50, 75, 100, 250, 500, 750, 1000\}$, but this time the adjacency matrices $\mathbf{A}^{(1)}$ and $\mathbf{A}^{(2)}$ were distributed according to a 3-community multilayer stochastic block model, where the matrices $\mathbf{B}^{(r)}$ were randomly chosen, and vertices assigned to a community uniformly at random, discarding any trials for which the matrix \mathbf{X} of community assignments was not of full rank. We then calculated and compared the average of the subspace distances $\|\mathbf{U}_A\mathbf{U}_A^\top - \mathbf{U}\mathbf{U}^\top\|$ and $\|\hat{\mathbf{U}}\hat{\mathbf{U}}^\top - \mathbf{U}\mathbf{U}^\top\|$ across each of the 1000 trials. For the second experiment, we again fixed $n = 750$, performed 1000 independent trials as above for $k = 2, \dots, 10$ embeddings, and compared the subspace distance between the estimated and actual invariant subspaces.

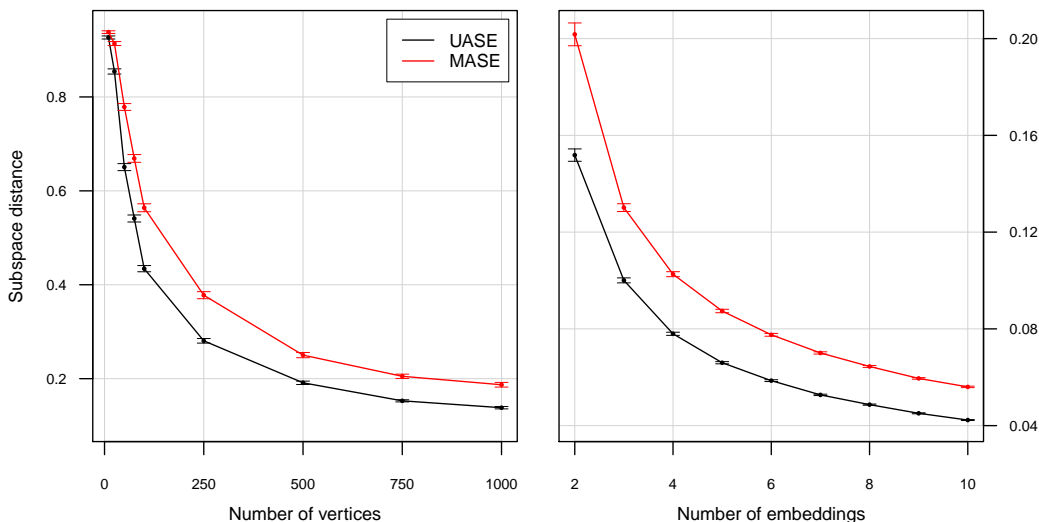


Figure 8: Average distance between the estimated and actual invariant subspaces in a 3-community multilayer stochastic block model as a function of the number of vertices (left-hand graph) and the number of embeddings (right-hand graph).

Figure 8 plots the results of the two experiments. For this task, although the performance of the two embedding types is almost indistinguishable for very small graphs, as the number of vertices grows the UASE consistently outstrips the MASE. As in the previous example, increasing the number of embedded graphs results in greater accuracy for both methods, where again the UASE offers the best performance of the two.

4.3 Model estimation

As a final comparison of the UASE and MASE methods, we investigate the efficiency of both at the task of estimating the underlying matrices $\mathbf{P}^{(r)}$ in the MRDPG and COSIE models, which is of particular practical interest for link prediction tasks. To establish an appropriate estimate, we first consider the case of the standard GRDPG (that is, when $k = 1$). In this case, an estimate $\hat{\mathbf{P}}$ for the matrix \mathbf{P} can be obtained by setting $\hat{\mathbf{P}} = \mathbf{X}_A \mathbf{I}_{p,q} \mathbf{X}_A^\top$. Note that due to orthogonality

of the singular vectors, the matrix $\mathbf{X}_A \in \mathbb{R}^{n \times d}$ of the leading d left singular vectors of \mathbf{A} is the projection of the *full* matrix of left singular vectors onto the d -dimensional subspace spanned by \mathbf{U}_A . Since this projection corresponds to left multiplication by the matrix $\mathbf{U}_A \mathbf{U}_A^\top$, we have the alternative description $\widehat{\mathbf{P}} = \mathbf{U}_A \mathbf{U}_A^\top \mathbf{A} \mathbf{U}_A \mathbf{U}_A^\top$.

Returning to the general case, we obtain an estimate $\widehat{\mathbf{P}}^{(r)} = \mathbf{U}_A \mathbf{U}_A^\top \mathbf{A}^{(r)} \mathbf{U}_A \mathbf{U}_A^\top$ for the matrix $\mathbf{P}^{(r)}$ for each $r \in \{1, \dots, k\}$ using the unscaled UASE. For the MASE, we use the matrix $\widehat{\mathbf{U}} \widehat{\mathbf{U}}^\top \mathbf{A}^{(r)} \widehat{\mathbf{U}} \widehat{\mathbf{U}}^\top$ as our estimate. For each of the trials in the previous example, we calculated these estimates, and measured the model estimation error in each case using the normalised mean squared error

$$\frac{\|\widehat{\mathbf{P}}^{(r)} - \mathbf{P}^{(r)}\|_F}{\|\mathbf{P}^{(r)}\|_F}. \quad (21)$$

Figure 9 plots the results of the two experiments, in which we see that once again the UASE consistently demonstrates greater accuracy than the MASE for all but the smallest of graphs, and for all numbers of embedded graphs.

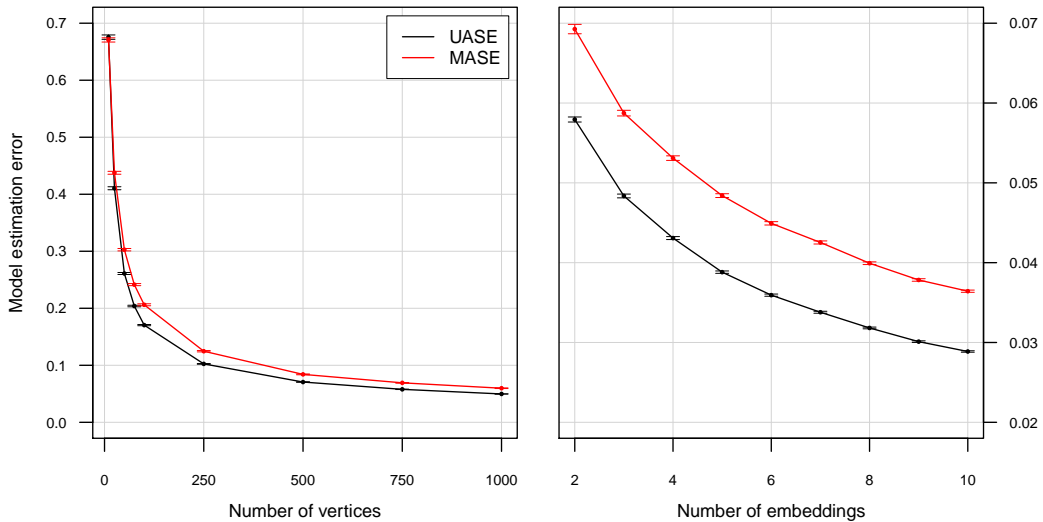


Figure 9: Model estimation error for the UASE and MASE in a 3-community multilayer stochastic block model as a function of the number of vertices (left-hand graph) and the number of embeddings (right-hand graph).

4.4 Two-graph hypothesis testing

When $\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(k)}$ are identically distributed, the right embeddings $\mathbf{Y}_A^{(r)}$ are identically distributed too and each subject to the *same* unidentifiable linear transformation (Corollaries 4 and 5). It is therefore natural to consider the effectiveness of the UASE at testing the semiparametric hypothesis that two observed graphs are drawn from the same underlying latent positions. This problem was considered for the omnibus embedding in [21], and we shall use the framework established there to test the UASE. Suppose, then, that we have points $\mathbf{X}_1, \dots, \mathbf{X}_n, \mathbf{Y}_1, \dots, \mathbf{Y}_n \in \mathbb{R}^d$, and that we have two graphs G_1 and G_2 whose adjacency matrices $\mathbf{A}^{(1)}$ and $\mathbf{A}^{(2)}$ satisfy $\mathbf{A}_{ij}^{(1)} \sim \text{Bern}(\mathbf{X}_i^\top \mathbf{I}_{p,q} \mathbf{X}_j)$ and $\mathbf{A}_{ij}^{(2)} \sim \text{Bern}(\mathbf{Y}_i^\top \mathbf{I}_{p,q} \mathbf{Y}_j)$. The UASE allows us to test the hypothesis:

$$H_0 : \mathbf{X}_i = \mathbf{Y}_i \quad \forall i \in \{1, \dots, n\} \quad (22)$$

by comparing the right embeddings $\mathbf{Y}_A^{(1)}$ and $\mathbf{Y}_A^{(2)}$. If H_0 holds, then the rows of the $\mathbf{Y}_A^{(r)}$ are identically (although not independently) distributed, whereas if H_0 fails to hold then for some k the k th row of $\mathbf{Y}_A^{(1)}$ and $\mathbf{Y}_A^{(2)}$ should be distributionally distinct.

The framework used in [21] to test this hypothesis, which we shall repeat here, is as follows: we begin by drawing $\mathbf{X}_1, \dots, \mathbf{X}_n, \mathbf{Z}_1, \dots, \mathbf{Z}_n \in \mathbb{R}^3$ identically according to a Dirichlet distribution

with parameter $\alpha = (1, 1, 1)^\top$, select a subset I of some fixed size uniformly at random among all such subsets of $\{1, \dots, n\}$, and define

$$\mathbf{Y}_i = \begin{cases} \mathbf{Z}_i & \text{if } i \in I \\ \mathbf{X}_i & \text{otherwise} \end{cases} \quad (23)$$

We generate two graphs \mathcal{G}_1 and \mathcal{G}_2 with adjacency matrices $\mathbf{A}^{(1)}$ and $\mathbf{A}^{(2)}$ satisfying $\mathbf{A}_{ij}^{(1)} \sim \text{Bern}(\mathbf{X}_i^\top \mathbf{X}_j)$ and $\mathbf{A}_{ij}^{(2)} \sim \text{Bern}(\mathbf{Y}_i^\top \mathbf{Y}_j)$, and estimate the latent positions $\hat{\mathbf{X}}$ and $\hat{\mathbf{Y}}$ of the two graphs by using the right embeddings as described above, and (in the case of the omnibus embedding) the first and last n rows of the spectral embedding of the matrix

$$\mathbf{M} = \begin{pmatrix} \mathbf{A}^{(1)} & \frac{1}{2}(\mathbf{A}^{(1)} + \mathbf{A}^{(2)}) \\ \frac{1}{2}(\mathbf{A}^{(1)} + \mathbf{A}^{(2)}) & \mathbf{A}^{(2)} \end{pmatrix}. \quad (24)$$

We note that this is only possible due to our prior knowledge of the matrix \mathbf{P} , which allows us to construct the required transformations.

In both cases we use the test statistic $T = \sum_{i=1}^n \|\hat{\mathbf{X}}_i - \hat{\mathbf{Y}}_i\|^2$; and accept or reject based on an estimate of the critical value of T under the null hypothesis obtained by using 2000 Monte Carlo iterates to estimate the distribution of T .

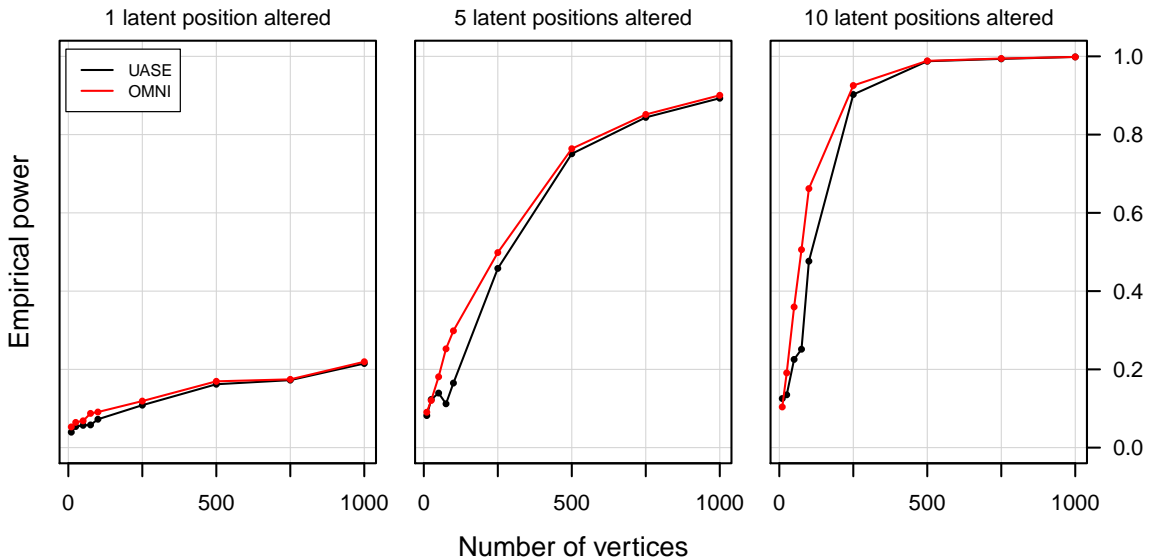


Figure 10: Empirical power of the UASE (black) and omnibus (red) tests to detect when the two graphs being tested differ in the specified number of their latent positions. Each point is the proportion of 2000 trials for which the given technique correctly rejected the null hypothesis.

Figure 10 shows the power of each method for testing the null hypothesis for different sized graphs and for different numbers of altered latent positions, by calculating the proportion (out of 2000 trials conducted for each sized graph) of trials for which we correctly reject the null hypothesis. For smaller graphs, the omnibus embedding provides the most effective method (although there is not much difference between the two where only one latent position is altered - however in this case the empirical power of both methods does not exceed 0.25). For larger graphs, particularly those with more than 500 vertices, the two methods are almost indistinguishable, and in such cases the UASE might be preferred based on size considerations, due to only requiring an $n \times kn$ rather than an $kn \times kn$ matrix.

5 Real data: Link prediction on a computer network

5.1 Dynamic link prediction

The Los Alamos National Laboratory computer network [37] was studied in [29], in which it was demonstrated empirically that the disassortative connectivity behaviour inherent in the network leads to the GRDPG offering a marked modelling improvement over the RDPG in the task of out-of-sample link prediction between computers in the network. For a large-scale dynamic network such as this, the MRDPG offers the possibility of further refinement by allowing us to consider multiple “snapshots” of communication behaviour at different points in time simultaneously.

As an example, we extract a ten minute sample at random from the “Network Event Data” dataset, which we divide into two separate five minute samples. From the first sample we generate five graphs, each one describing the communication behaviour of the computers in the network over a period of one minute, by assigning each IP address to a node (with this assignation being kept consistent across all graphs), and recording an edge between two nodes if the corresponding edges are observed to communicate at least once within this period, and then construct the corresponding adjacency matrices $\mathbf{A}^{(r)}$. Setting our embedding dimension $d = 10$ (an admittedly arbitrary choice) we then generate estimates $\hat{\mathbf{P}}^{(r)}$ for the probability matrices as described in Section 4.3. We then use the average of these matrices to give us estimates of the probabilities of a link being generated between any given pair of computers.

In a similar manner, we generate estimates of the link probabilities for the mean adjacency matrix $\bar{\mathbf{A}}$. We also construct adjacency matrices $\mathbf{A}_{[1]}$ and $\mathbf{A}_{[5]}$ from the connectivity graphs for the first minute and the full five minute period respectively (both of which follow a standard GRDPG) and generate link probability estimates accordingly.

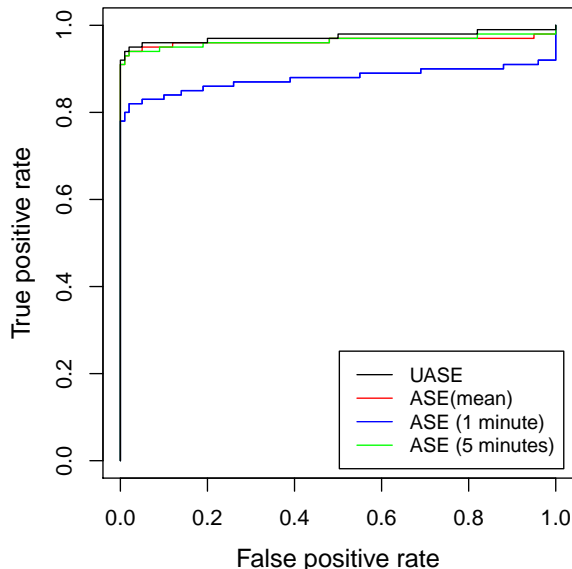


Figure 11: Receiver Operating Characteristic curves for the UASE (black), mean embedding (red) and ASE (blue and green) methods for out-of-sample link prediction on the Los Alamos National Laboratory computer network. See main text for details.

Using these estimates, we attempt to predict which *new* edges will occur within the second five minute window, disregarding those involving new nodes. Figure 11 shows the receiver operating characteristic (ROC) curves for each model and for each port, where we treat the prediction task as a binary classification problem whose outcomes are either the presence or absence of an edge between nodes, which we predict by thresholding the estimated link probabilities. As one would expect, using the ASE of the graph corresponding to only a single minute of communication (the blue curve) produces the least accurate predictions, but the UASE (black), mean embedding (red) and the ASE for a five minute sample (green) all produce similar results, with the UASE slightly outperforming the other two methods. This can be confirmed numerically by calculating the area

under each ROC curve (AUC) which is equivalent to the probability that a given classifier will rank a randomly chosen positive instance higher than a randomly chosen negative instance [11]. We calculate that the UASE has an AUC of 0.9734, which is a small improvement over the mean embedding and the 5-minute ASE, which have AUC values of 0.9627 and 0.9635 respectively (the 1-minute ASE, by contrast, yields an AUC of 0.8767).

5.2 Port-specific link prediction

The LANL network data presents the opportunity to demonstrate another significant improvement offered by the MRDPG over the GRDPG, namely its ability to integrate data from sources which do not necessarily behave similarly. Each communication within the LANL network passes through a source and destination port, the latter of which indicates the type of service being used, and it is natural to expect that different services may exhibit different communication behaviours.

We consider the first five minute sample from the previous section. During the first minute alone, there are a total of 121,737 (not necessarily unique) communications between 10,762 computers, with 4,379 different destination ports being used. Of these, the 8 most commonly-used ports account for over 75% of communications, and Table 1 lists these, together with the purpose to which each port is assigned.

Table 1: Purpose and proportion of traffic utilizing the 8 most frequently used ports during 1 minute of activity on the Los Alamos National Laboratory computer network.

Port	Purpose	Proportion of traffic
53	DNS	27.7%
443	HTTPS	14.8%
80	HTTP	11.9%
514	Syslog	7.2%
389	Lightweight Directory Access Protocol	4.3%
427	Service Location Protocol	4.1%
88	Kerberos authentication system	3.4%
445	Microsoft-DS Active Directory	1.8%

For each of these 8 ports, we generate a graph of the communications made between computers within the network through this specific port over the first minute of our sample as in the previous example. Figure 12 visualizes the adjacency matrix of each of these graphs as a 2-dimensional plot, together with the adjacency matrix of the full network graph.

As before, we calculate estimates of the link probabilities for each port using the UASE, but now rather than averaging them we consider each port individually. For comparison, we also estimate link probabilities for each individual port using the corresponding GRDPG, and then use both estimates to attempt to predict which *new* edges will occur within the remainder of our five-minute window. Figure 13 shows the ROC curves for each model and for each port, while Table 2 gives the AUC values for each curve. We note that for Port 53 (the busiest port) the standard ASE actually outperforms the UASE, but for every other port the UASE is the superior method, offering a significant improvement over the ASE for the less active ports.

Table 2: AUC values for the ROC curves in Figure 13.

Embedding	Port 53	Port 443	Port 80	Port 514	Port 389	Port 427	Port 88	Port 445
UASE	0.8391	0.7850	0.9010	0.8931	0.8534	0.8580	0.8949	0.9836
ASE	0.8568	0.6370	0.7185	0.7806	0.5532	0.6566	0.5668	0.5720

5.3 Link prediction using mixed data sources

Our final example demonstrates how combining data from graphs with *different* nodes can inform our knowledge of a common subset. We begin by extracting a five minute sample at random from the “Network Event Data” dataset similarly to Section 5.1, and construct the ASE of the adjacency matrix $\mathbf{A}^{(1)}$ obtained by restricting our attention to the first minute of computer-to-computer communication. For the UASE, we augment this data by constructing the submatrix $\mathbf{A}^{(2)}$ of computer-to-port communications, in which $\mathbf{A}_{ij}^{(2)} = 1$ if there is a connection involving

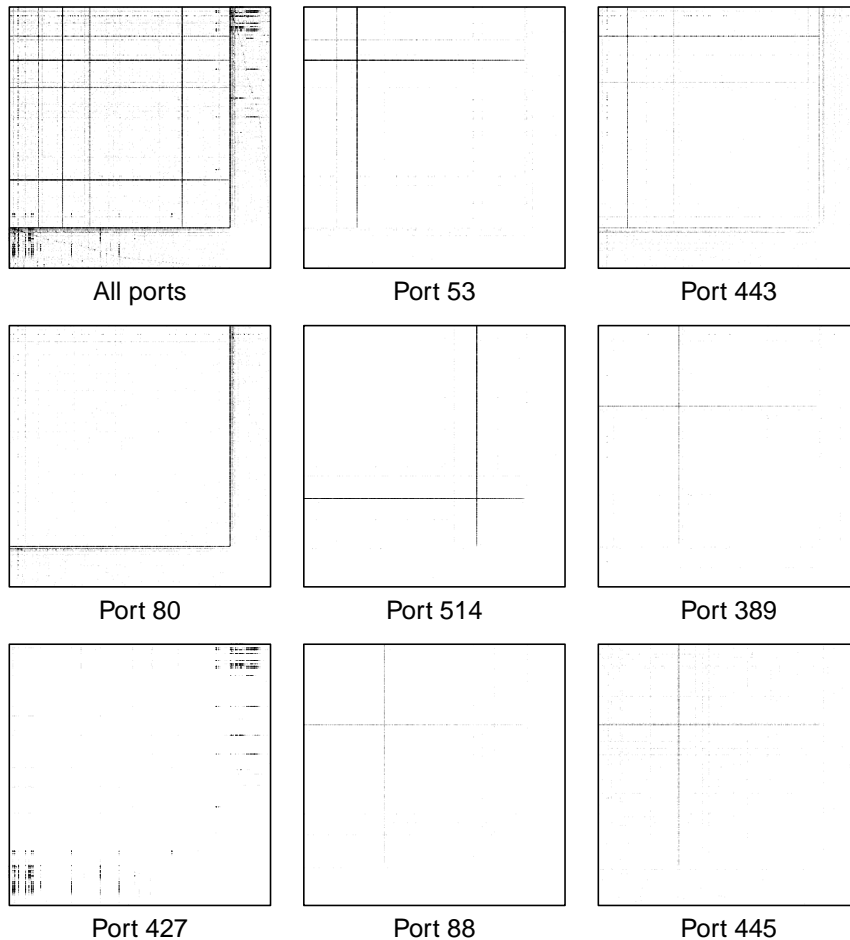


Figure 12: Visualization of adjacency matrices showing connections between computers on the Los Alamos National Laboratory computer network during 1 minute of activity. The top-left image shows all connections during this time, while the remaining images show only connections via the specified port.

the i th computer during this first minute for which the j th port is the destination port. We then generate estimates of the link probabilities for the ASE of $\mathbf{A}^{(1)}$ and UASE of $\mathbf{A} = [\mathbf{A}^{(1)}|\mathbf{A}^{(2)}]$, and generate link probability estimates accordingly. Figure 14 plots the resulting ROC curves, in which we observe that augmenting the computer-to-computer connectivity data with computer-to-port connectivity data yields an improvement in prediction power (similarly, we note the AUC values of 0.949 for the UASE compared to 0.905 for the ASE).

6 Chernoff information and the GMSBM

Given a collection of matrices $\mathbf{A}^{(r)}$ that are distributed according to a GMSBM, it is reasonable to ask whether there is any tangible benefit to studying the UASE as opposed to the ASEs of the individual matrices, and how one might quantify this. Tang and Priebe [33], in the context of comparing the performance of the spectral embeddings of the Laplacian and adjacency matrix in recovering block assignments from a stochastic block model graph, proposed using the *Chernoff information* [14] of the limiting Gaussian distributions obtained from the Central Limit Theorem associated to each embedding as a means of doing so. In a two-cluster problem, the Chernoff information is the exponential rate at which the Bayes error (from the decision rule which assigns each data point to its most likely cluster *a posteriori*) decreases asymptotically. The Chernoff information is an example of a f -divergence [2], [9] and therefore possesses the desirable attribute

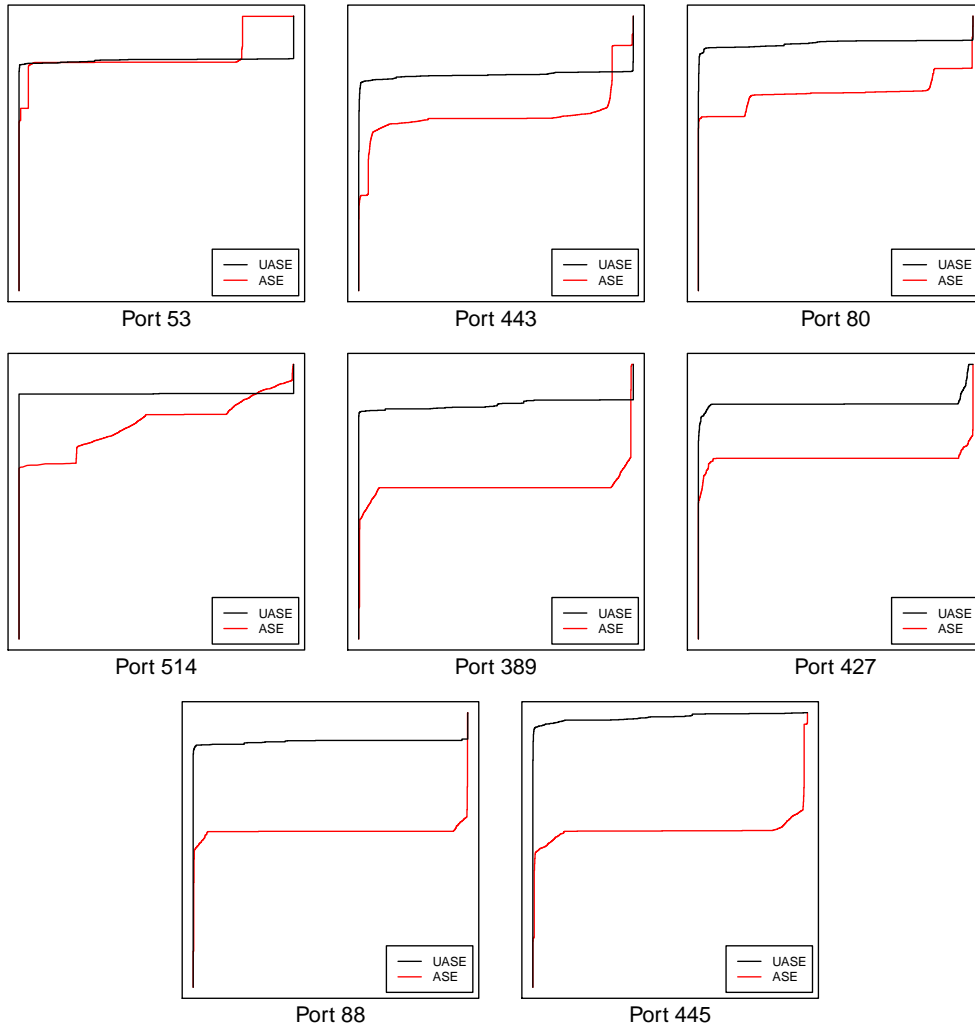


Figure 13: Receiver Operating Characteristic curves for the UASE (black) and ASE (red) for out-of-sample link prediction via the specified port on the Los Alamos National Laboratory computer network. See main text for details.

of being invariant under invertible linear transformations [22].

If F_1 and F_2 are two absolutely continuous multivariate distributions supported on $\Omega \subset \mathbb{R}^d$, with density functions f_1 and f_2 respectively, the Chernoff information between F_1 and F_2 is defined by $C(F_1, F_2) = \sup_{t \in (0,1)} C_t(F_1, F_2)$ [14] where the *Chernoff divergence*

$$C_t(F_1, F_2) = -\log\left(\int_{\Omega} f_1(\mathbf{x})^t f_2(\mathbf{x})^{1-t} d\mathbf{x}\right). \quad (25)$$

For a K -cluster problem, in which we have distributions F_1, \dots, F_K with corresponding density functions f_1, \dots, f_K , we consider the Chernoff information of the critical pair $\min_{i \neq j} C(F_i, F_j)$.

If the F_i are multivariate normal distributions, it is known (see [27]) that the Chernoff information can be expressed as $C(F_i, F_j) = \sup_{t \in (0,1)} C_t(F_i, F_j)$, where

$$C_t(F_i, F_j) = \left(\frac{t(1-t)}{2}(\mathbf{x}_i - \mathbf{x}_j)^\top \boldsymbol{\Sigma}_t^{-1}(\mathbf{x}_i - \mathbf{x}_j) + \frac{1}{2} \log\left(\frac{|\boldsymbol{\Sigma}_t|}{|\boldsymbol{\Sigma}_i|^t |\boldsymbol{\Sigma}_j|^{1-t}}\right)\right), \quad (26)$$

where $F_i \sim \mathcal{N}(\mathbf{x}_i, \boldsymbol{\Sigma}_i)$ and $\boldsymbol{\Sigma}_t = t\boldsymbol{\Sigma}_1 + (1-t)\boldsymbol{\Sigma}_2$.

We now apply this to obtain an expression for the Chernoff information of the left UASE of a

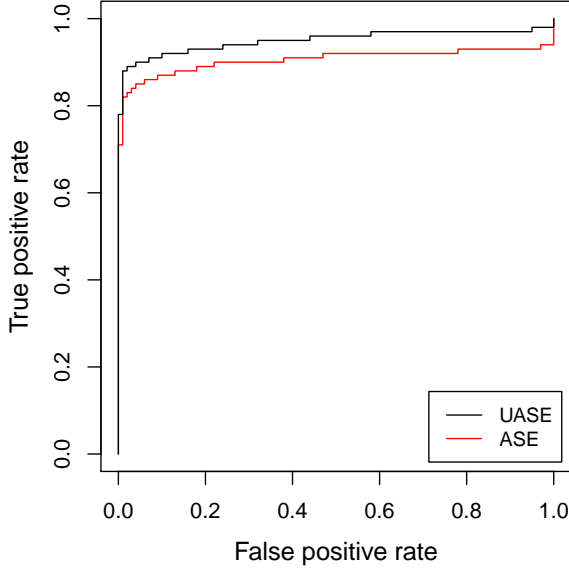


Figure 14: Receiver Operating Characteristic curves for the UASE (black) and ASE (red) methods for out-of-sample link prediction on the Los Alamos National Laboratory computer network, in which the UASE is augmented with computer-to-port connectivity data. See main text for details.

GMSBM. Let $(\mathbf{A}, \mathbf{X}, \mathbf{Y}) \sim \text{GMSBM}(\mathcal{F}, \mathbf{B})$, and define $C_{\mathbf{A}} = \min_{i \neq j} \sup_{t \in (0,1)} C_t(F_i, F_j)$, where

$$F_i \sim \mathcal{N}(\mathbf{e}_i, \frac{1}{n} \Delta_{\mathbf{B}, \mathbf{Y}}^{-1} \mathbf{B} \Sigma_{\mathbf{Y}}(\mathbf{e}_i) \mathbf{B}^{\top} \Delta_{\mathbf{B}, \mathbf{Y}}^{-1}) \quad (27)$$

and $\Delta_{\mathbf{B}, \mathbf{Y}}$ and $\Sigma_{\mathbf{Y}}(\mathbf{e}_i)$ are as defined in Theorem 3 (for simplicity, we will assume that our graphs have the same sparsity factor).

Some standard algebraic manipulation shows that this quantity is invariant under the transformations of the latent positions listed in Proposition 1 (and thus the Chernoff information of the underlying MRDPG is well-defined), and similarly that it is invariant under invertible linear transformations of the UASE $\mathbf{X}_{\mathbf{A}}$. Thus we may study $\mathbf{X}_{\mathbf{A}}$ rather than the estimate $\mathbf{X}_{\mathbf{A}} \mathbf{L}$ of the latent positions, which requires knowledge - that we will not typically possess - of the underlying matrices of probabilities $\mathbf{P}^{(r)}$. We note that we can similarly define the Chernoff information for the right UASEs (where they are defined) and that this too would be invariant under invertible transformations of the latent positions.

Given a collection $\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(k)}$ of matrices and a subset \mathcal{K} of $\{1, \dots, k\}$, let $\mathbf{A}_{\mathcal{K}}$ denote the matrix obtained by concatenating the matrices $\mathbf{A}^{(r)}$ for $r \in \mathcal{K}$. If we have a collection of matrices $\mathbf{A}^{(r)}$ which are *identically* distributed as a GMSBM, the following result shows that it is *always* preferable to embed as many matrices as possible:

Proposition 6. *Let $(\mathbf{A}, \mathbf{X}, \mathbf{Y}) \stackrel{\text{id}}{\sim} \text{GMSBM}(\mathcal{F}, \mathbf{B})$ be identically distributed as a GMSBM. Then $C_{\mathbf{A}} \geq C_{\mathbf{A}_{\mathcal{K}}}$ for any subset \mathcal{K} of $\{1, \dots, k\}$.*

Proof. Let $\Sigma_i = \mathbf{B}^{-\top} \Delta_{\mathbf{Y}}^{-1} \Sigma_{\mathbf{Y}}(\mathbf{e}_i) \Delta_{\mathbf{Y}}^{-1} \mathbf{B}^{-1}$ denote the covariance matrix present in the Central Limit Theorem for the single graph ASE given in Corollary 5 for our particular case. Then

$$C_{\mathbf{A}} = \min_{i \neq j} \sup_{t \in (0,1)} \left(\frac{kt(1-t)}{2c} (\mathbf{e}_i - \mathbf{e}_j)^{\top} \Sigma_{i,j,t}^{-1} (\mathbf{e}_i - \mathbf{e}_j) + \frac{1}{2} \log \left(\frac{|\Sigma_{i,j,t}|}{|\Sigma_i|^t |\Sigma_j|^{1-t}} \right) \right) \quad (28)$$

where $\Sigma_{i,j,t} = t\Sigma_i + (1-t)\Sigma_j$, while for a subset \mathcal{K} of size r we have

$$C_{\mathbf{A}_{\mathcal{K}}} = \min_{i \neq j} \sup_{t \in (0,1)} \left(\frac{rt(1-t)}{2c} (\mathbf{e}_i - \mathbf{e}_j)^{\top} \Sigma_{i,j,t}^{-1} (\mathbf{e}_i - \mathbf{e}_j) + \frac{1}{2} \log \left(\frac{|\Sigma_{i,j,t}|}{|\Sigma_i|^t |\Sigma_j|^{1-t}} \right) \right). \quad (29)$$

Since the matrix $\Sigma_{i,j,t}$ is positive semi-definite for any i and j , the first term in each pair of brackets is positive, and thus the term $C_{\mathbf{A}}$ clearly dominates. \square

If the adjacency matrices $\mathbf{A}^{(r)}$ are *not* identically distributed, however, the situation is not so clear-cut, as it is entirely possible to encounter situations for which $C_{\mathbf{A}_{\mathcal{K}}} > C_{\mathbf{A}}$ for some subset \mathcal{K} . For example, if we consider the two-graph multilayer stochastic block model with matrices

$$\mathbf{B}^{(1)} = \begin{pmatrix} 0.67 & 0.46 \\ 0.46 & 0.36 \end{pmatrix}, \quad \mathbf{B}^{(2)} = \begin{pmatrix} 0.98 & 0.49 \\ 0.49 & 0.10 \end{pmatrix}, \quad (30)$$

then while the ratio $C_{\mathbf{A}}/C_{\mathbf{A}^{(1)}}$ tends to 11.98, the ratio $C_{\mathbf{A}}/C_{\mathbf{A}^{(2)}}$ tends to 0.96.

Before proceeding further, we note that any analysis of the Chernoff information for large-scale GMSBMs can be simplified by observing that the logarithmic term in the definition is independent of the number of vertices n , and so becomes insignificant as $n \rightarrow \infty$ if we impose the simplifying assumption that the covariance matrices be non-singular. To this end, we consider instead the truncated terms $\rho_{\mathbf{A}}$, in which we simply omit the logarithmic term from the definition of $C_{\mathbf{A}}$. Considering the function $\rho_{\mathbf{A}}$ gives an accurate means of comparison between the large-scale behaviour of two multilayer stochastic block models, as the ratio $C_{\mathbf{A}}/C_{\mathbf{A}_{\mathcal{K}}}$ tends to $\rho_{\mathbf{A}}/\rho_{\mathbf{A}_{\mathcal{K}}}$ as n increases for any subset \mathcal{K} .

To observe the effect of embedding additional graphs on the Chernoff information of a GMSBM, we conducted the following experiment: for each $k \in \{2, \dots, 10\}$, we performed 1000 trials in which a $(2, 2, \dots, 2)$ -community GMSBM was generated by choosing k matrices $\mathbf{B}^{(r)} \in [0, 1]^{2 \times 2}$ and a probability vector $\boldsymbol{\pi}$ at random (with the entries $\mathbf{B}_{ij}^{(r)} \sim \text{Uniform}[0, 1]$, while $\boldsymbol{\pi} \sim \text{Dirichlet}(1, 1)$). We then calculated the ratio $\rho_{\mathbf{A}}/\rho_{\mathbf{A}_{\mathcal{K}}}$ for each subset \mathcal{K} of $\{1, \dots, k\}$ under the assumption that the vectors $\mathbf{Y}^{(r)}$ were identically distributed. We then repeated the experiment for a $(2, 2, \dots, 2)$ -community GMSBM in which we allowed different probability vectors $\boldsymbol{\pi}_r$ for each embedding, and finally repeated both experiments for a $(2, 3, \dots, 3)$ -community GMSBM.

The results of these simulations are presented in Figure 15. Entries above the diagonal (in blue) correspond to trials in which the parameters $\boldsymbol{\pi}_r$ are equal, while entries below the diagonal (in red) correspond to trials in which the parameters $\boldsymbol{\pi}_r$ are allowed to differ; the (r, k) th and (k, r) th coordinates in the respective cases indicate the proportion of embeddings of k matrices for which $\rho_{\mathbf{A}}/\rho_{\mathbf{A}_{\mathcal{K}}} > 1$ for *every* r -element subset \mathcal{K} .

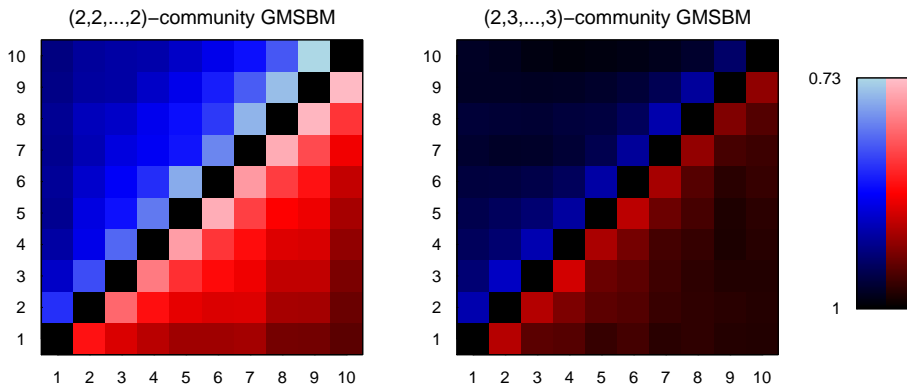


Figure 15: Proportion of GMSBMs for which $\rho_{\mathbf{A}}/\rho_{\mathbf{A}_{\mathcal{K}}} > 1$ for all subsets \mathcal{K} of a given size. See main text for details.

We make the following observations:

- On average, embedding more matrices seems to improve performance, with $\rho_{\mathbf{A}}/\rho_{\mathbf{A}_{\mathcal{K}}} > 1$ for at least 73% of the trials we conducted, and this improvement in performance seems to increase as we allow the latent positions $\mathbf{Y}^{(r)}$ to be taken from more communities.
- In general, $\rho_{\mathbf{A}}/\rho_{\mathbf{A}_{\mathcal{K}}} > 1$ more often the greater the difference in size between the subset \mathcal{K} and the total number of embeddings k .
- We detected little distinguishable difference between allowing the parameters $\boldsymbol{\pi}_r$ to differ as opposed to keeping them all the same.

- In each trial, we noted that there is always at least *one* subset \mathcal{K} of any given size for which $\rho_{\mathbf{A}}/\rho_{\mathbf{A}_{\mathcal{K}}} > 1$, and we conjecture that this should always hold. If true, this would mean in particular that the UASE is always better than the worst of our embeddings (as opposed to, say, the mean embedding, which we have seen can lead to degeneracy when the matrices $\mathbf{P}^{(r)}$ have differing signatures).

7 Conclusion

The multilayer random dot product graph is a vast yet natural extension of the random dot product graph, granting us insight into the behaviour of a common subset of nodes across a series of graphs—both undirected and directed—in which we allow a mixture of assortativity behaviours. Its simplicity and flexibility make it an ideal model for a variety of situations, and it can be seen to perform as well as (and in many cases better than) existing models at multiple graph inference tasks such as community detection and graph-to-graph comparison, while allowing inference in a much wider range of situations than current methods permit. These experimental results are supported by theoretical results showing that the node representations obtained by the left- and right-sided spectral embeddings converge uniformly in the Euclidean norm to the latent positions with Gaussian error, in particular providing us with the first known examples of such results for bipartite graphs. Finally, we demonstrate the practical effectiveness of our model by applying it to the task of link prediction within a computer network, indicating its usefulness to the field of cyber-security.

Bibliography

- [1] L. Akoglu and C. Faloutsos. Anomaly, event, and fraud detection in large network datasets. In *Proceedings of the sixth ACM international conference on Web search and data mining*, pages 773–774, 2013.
- [2] S. M. Ali and S. D. Shelvey. A general class of coefficients of divergence of one distribution from another. *Journal of the Royal Statistical Society: Series B*, 28:121–132, 1966.
- [3] J. Arroyo, A. Athreya, J. Cape, G. Chen, C. E. Priebe, and J. T. Vogelstein. Inference for multiple heterogeneous networks with a common invariant subspace. *arXiv preprints arXiv:1906.10026*, 2019.
- [4] A. Athreya, C. E. Priebe, M. Tang, V. Lyzinski, D. J. Marchette, and D. L. Sussman. A limit theorem for scaled eigenvectors of random dot product graphs. *Sankhya A*, 78(1):1–18, 2016.
- [5] A. Athreya, D. E. Fishkind, M. Tang, C. E. Priebe, Y. Park, J. T. Vogelstein, K. Levin, V. Lyzinski, and Y. Qin. Statistical inference on random dot product graphs: a survey. *The Journal of Machine Learning Research*, 18(1):8393–8484, 2017.
- [6] J. Cape, M. Tang, and C. E. Priebe. The two-to-infinity norm and singular subspace geometry with applications to high-dimensional statistics. *The Annals of Statistics*, 2017. To appear; preprint available at <http://arxiv.org/abs/1705.10735>.
- [7] J. Cape, M. Tang, and C. E. Priebe. On spectral embedding performance and elucidating network structure in stochastic blockmodel graphs. *Network Science*, 7(3):269–291, 2019.
- [8] J. Cape, M. Tang, and C. E. Priebe. Signal-plus-noise matrix models: eigenvector deviations and fluctuations. *Biometrika*, 106(1):243–250, 2019.
- [9] I. Csizár. Information-type measures of difference of probability distributions and indirect observations. *Studia Scientiarum Mathematicarum Hungarica*, 2:229–318, 1967.
- [10] P. Erdős and A. Rényi. On the evolution of random graphs. *Proceedings of the Hungarian Academy of Sciences*, pages 17–61, 1960.
- [11] T. Fawcett. An introduction to ROC analysis. *Pattern Recogn. Lett.*, 27(8):861–874, 2006.

- [12] A. Fornito, A. Zalesky, and E. Bullmore. *Fundamentals of brain network analysis*. Academic Press, 2016.
- [13] E. N. Gilbert. Random graphs. *Annals of Mathematical Statistics*, 30(4):1141–1144, 1959.
- [14] Chernoff. H. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *The Annals of Mathematical Statistics*, 23(4):493–507, 1952.
- [15] P. D. Hoff, A. E. Raftery, and M. S. Handcock. Latent space approaches to social network analysis. *Journal of the American Statistical Association*, 97:1090–1098, 2002.
- [16] P. W. Holland, K. Laskey, and S. Leinhardt. Stochastic blockmodels: First steps. *Social Networks*, 5:109–137, 1983.
- [17] R. Horn and C. Johnson. *Matrix Analysis (Second Edition)*. Cambridge University Press, New York, NY, 2012.
- [18] S. Khor. Concurrency and network disassortativity. *Artificial life*, 16(3):225–232, 2010.
- [19] M. Kivelä, A. Arenas, M. Barthelemy, J. P. Gleeson, Y. Moreno, and M. A. Porter. Multilayer networks. *Journal of complex networks*, 2(3):203–271, 2014.
- [20] L. Lathauwer, B. Moor, and J. Vandewalle. A multilinear singular value decomposition. *SIAM J. Matrix Anal. Appl.*, 21(4):1253–1278, 2000.
- [21] K. Levin, A. Athreya, M. Tang, V. Lyzinski, and C. E. Priebe. A central limit theorem for an omnibus embedding of random dot product graphs. *arXiv preprint arXiv:1705.09355*, 2017.
- [22] F. Liese and I. Vajda. On divergences and informations in statistics and information theory. *IEEE Transactions on Information Theory*, 52:4394–4412, 2006.
- [23] V. Lyzinski, M. Tang, A. Athreya, Y. Park, and C. E. Priebe. Community detection and classification in hierarchical stochastic blockmodels. *IEEE Transactions in Network Science and Engineering*, 4(1):13–26, 2017.
- [24] M. E. J. Newman. Assortative mixing in networks. *Physical Review Letters*, 89:208701, 2002.
- [25] C. Nickel. *Random dot product graphs: a model for social networks*. PhD thesis, Johns Hopkins University, 2006.
- [26] A. M. Nielsen and D. Witten. The multiple random dot product graph model. *arXiv preprint arXiv:1811.12172*, 2018.
- [27] L. Pardo. *Statistical inference based on divergence measures*. Chapman and Hall/CRC, 2005.
- [28] F. S. Passino, A. S. Bertiger, J. C. Neil, and N. A. Heard. Link prediction in dynamic networks using random dot product graphs. *arXiv preprint arXiv:1912.10419*, 2019.
- [29] P. Rubin-Delanchy, J. Cape, C. E. Priebe, and M. Tang. A statistical interpretation of spectral embedding: the generalised random dot product graph. *arXiv preprint arXiv:1709.05506v3*, 2020.
- [30] P. Sarkar and P. J. Bickel. Role of normalization in spectral clustering for stochastic blockmodels. *The Annals of Statistics*, 43(3):962–990, 2015.
- [31] L. Scrucca, M. Fop, T. B. Murphy, and A. E. Raftery. mclust 5: clustering, classification and density estimation using Gaussian finite mixture models. *The R Journal*, 8(1):205–233, 2016.
- [32] D. L. Sussman, M. Tang, D. E. Fishkind, and C. E. Priebe. A consistent adjacency spectral embedding for stochastic blockmodel graphs. *Journal of the American Statistical Association*, 107(499):1119–1128, 2012.
- [33] M. Tang and C. Priebe. Limit theorems for eigenvectors of the normalized Laplacian for random graphs. *The Annals of Statistics*, 2019. To appear.

- [34] M. Tang, A. Athreya, D. L. Sussman, V. Lyzinski, Y. Park, and C. E. Priebe. A semiparametric two-sample hypothesis testing problem for random graphs. *Journal of Computational and Graphical Statistics*, 26(2):344–354, 2017.
- [35] R. Tang, M. Ketcha, A. Badea, E. Calabrese, D. Margulies, J. Vogelstein, C. Priebe, and D. Sussman. Connectome smoothing via low-rank approximations. *IEEE Transactions on Medical Imaging*, 38(6):1446–1456, 2019.
- [36] J. A. Tropp. An introduction to matrix concentration inequalities. *Foundations and Trends in Machine Learning*, 8(1–2):1–230, 2015.
- [37] M. Turcotte, A. Kent, and C. Hash. Unified host and network data set. *Data Science for Cyber-Security*, pages 1–22, 2018.
- [38] U. Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007.
- [39] S. Wang, J. D. Arroyo-Reli3n, J. T. Vogelstein, and C. E. Priebe. Joint embedding of graphs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019. To appear.
- [40] Stephen J. Young and Edward R. Scheinerman. Random dot product graph models for social networks. In *International Workshop on Algorithms and Models for the Web-Graph.*, pages 138–149. Springer, 2007.
- [41] Y. Yu, T. Wang, and R. J. Samworth. A useful variant of the Davis-Kahan theorem for statisticians. *Biometrika*, 102:315–323, 2015.
- [42] M. Zhu and A. Ghodsi. Automatic dimensionality selection from the scree plot via the use of profile likelihood. *Computational Statistics & Data Analysis*, 51(2):918–930, 2006.

Appendix

Before proving Theorems 2 and 3, we first require some control over the asymptotic behaviour of the singular values of the matrices \mathbf{P} , \mathbf{A} and $\mathbf{A} - \mathbf{P}$, which we establish using a series of results. Throughout this section, we will assume that $(\mathbf{A}, \mathbf{X}, \mathbf{Y}) \sim \text{MRDPG}(\mathcal{F}_\rho, \mathbf{\Lambda}_\epsilon)$ for a distribution \mathcal{F} and sparsity factors ρ and ϵ_r satisfying the criteria stated in Section 2.1.

Proposition 7. *The non-zero singular values $\sigma_i(\mathbf{P})$ for $i \in \{1, \dots, d\}$ satisfy*

$$\frac{\sigma_i(\mathbf{P})}{\rho n} \longrightarrow \sqrt{\lambda_i(\Delta_X \mathbf{\Lambda}_* \Delta_Y \mathbf{\Lambda}_*^\top)} \quad (31)$$

almost surely, where $\Delta_Y = c_1 \Delta_{Y,1} \oplus \dots \oplus c_\kappa \Delta_{Y,\kappa}$, and consequently $\sigma_i(\mathbf{P}) = O(\epsilon \rho n)$ and $\sigma_i(\mathbf{P}) = \Omega(\rho n)$ almost surely.

Proof. Since $\mathbf{P} = \mathbf{X} \mathbf{\Lambda}_\epsilon \mathbf{Y}^\top$, we see that

$$\sigma_i(\mathbf{P}) = \sqrt{\lambda_i(\mathbf{X} \mathbf{\Lambda}_\epsilon \mathbf{Y}^\top \mathbf{Y} \mathbf{\Lambda}_\epsilon^\top \mathbf{X}^\top)} = \sqrt{\lambda_i(\mathbf{X}^\top \mathbf{X} \mathbf{\Lambda}_\epsilon \mathbf{Y}^\top \mathbf{Y} \mathbf{\Lambda}_\epsilon^\top)}, \quad (32)$$

as the non-zero eigenvalues of a product of matrices are invariant under cyclic permutations of its factors.

Our assumptions regarding the distribution \mathcal{F} ensure that the spectral norms $\|\mathbf{X}^\top \mathbf{X} - \rho n \Delta_X\|$ and $\|\mathbf{Y}^\top \mathbf{Y} - \rho n \Delta_Y\|$ are of order $O(\rho n^{1/2} \log^{1/2}(n))$ mutually almost surely, and we note that consequently

$$\|\mathbf{Y}^\top \mathbf{Y}\| \leq \rho \|\Delta_Y\| + \|\mathbf{Y}^\top \mathbf{Y} - \rho \Delta_Y\| = O(\rho n) \quad (33)$$

almost surely by a standard application of the triangle inequality.

Next, note that

$$\begin{aligned} \mathbf{X}^\top \mathbf{X} \mathbf{\Lambda}_\epsilon \mathbf{Y}^\top \mathbf{Y} \mathbf{\Lambda}_\epsilon^\top - \rho^2 n \Delta_X \mathbf{\Lambda}_\epsilon \Delta_Y \mathbf{\Lambda}_\epsilon^\top &= (\mathbf{X}^\top \mathbf{X} - \rho n \Delta_X) \mathbf{\Lambda}_\epsilon \mathbf{Y}^\top \mathbf{Y} \mathbf{\Lambda}_\epsilon^\top \\ &\quad + \rho n \Delta_X \mathbf{\Lambda}_\epsilon (\mathbf{Y}^\top \mathbf{Y} - \rho \Delta_Y) \mathbf{\Lambda}_\epsilon^\top \end{aligned} \quad (34)$$

and so

$$\begin{aligned} \|\mathbf{X}^\top \mathbf{X} \mathbf{\Lambda}_\epsilon \mathbf{Y}^\top \mathbf{Y} \mathbf{\Lambda}_\epsilon^\top - \rho^2 n \Delta_X \mathbf{\Lambda}_\epsilon \Delta_Y \mathbf{\Lambda}_\epsilon^\top\| &\leq \|\mathbf{\Lambda}_\epsilon\|^2 (\|\mathbf{X}^\top \mathbf{X} - \rho n \Delta_X\| \|\mathbf{Y}^\top \mathbf{Y}\| \\ &\quad + \rho n \|\Delta_X\| \|\mathbf{Y}^\top \mathbf{Y} - \rho \Delta_Y\|) \end{aligned} \quad (35)$$

meaning that

$$\|\mathbf{X}^\top \mathbf{X} \mathbf{\Lambda}_\epsilon \mathbf{Y}^\top \mathbf{Y} \mathbf{\Lambda}_\epsilon^\top - \rho^2 n \Delta_X \mathbf{\Lambda}_\epsilon \Delta_Y \mathbf{\Lambda}_\epsilon^\top\| = O(\epsilon^2 \rho^2 n^{3/2} \log^{1/2}(n)) \quad (36)$$

almost surely.

As a result, we see that $\frac{1}{\rho^2 n^2} \mathbf{X}^\top \mathbf{X} \mathbf{\Lambda}_\epsilon \mathbf{Y}^\top \mathbf{Y} \mathbf{\Lambda}_\epsilon^\top$ converges to $\Delta_X \mathbf{\Lambda}_\epsilon \Delta_Y \mathbf{\Lambda}_\epsilon^\top$ in the spectral norm, and thus also in the Frobenius norm, implying in particular that the entries of the two matrices converge in absolute value. Now, the eigenvalues of any matrix are the roots of its characteristic polynomial, the coefficients of which are polynomial functions of the entries of the matrix. In particular, by continuity of the roots of polynomials, the eigenvalues of $\frac{1}{\rho^2 n} \mathbf{X}^\top \mathbf{X} \mathbf{\Lambda}_\epsilon \mathbf{Y}^\top \mathbf{Y} \mathbf{\Lambda}_\epsilon^\top$ and $\Delta_X \mathbf{\Lambda}_\epsilon \Delta_Y \mathbf{\Lambda}_\epsilon^\top$ must converge, giving the first result.

For the second, note that $\lambda_i(\Delta_X \mathbf{\Lambda}_\epsilon \Delta_Y \mathbf{\Lambda}_\epsilon^\top) = \lambda_i(\Delta_X^{1/2} \mathbf{\Lambda}_\epsilon \Delta_Y \mathbf{\Lambda}_\epsilon^\top \Delta_X^{1/2})$ (where $\Delta_X^{1/2}$ is the unique positive definite matrix square root of Δ_X) and that the latter matrix is a sum of symmetric matrices, allowing us to apply Weyl's inequalities to obtain the upper and lower bounds. \square

Proposition 8. $\|\mathbf{A} - \mathbf{P}\| = O(\epsilon^{1/2} \rho^{1/2} k^{1/4} n^{1/2} \log^{1/2}(n))$ almost surely.

Proof. Condition on some choice of latent positions. We will make use of a matrix analogue of the Bernstein inequality (see [36], Theorem 1.6.2):

Theorem 9 (Matrix Bernstein). *Let $\mathbf{M}_1, \dots, \mathbf{M}_n$ be independent random matrices with common dimensions $m_1 \times m_2$, satisfying $\mathbb{E}[\mathbf{M}_k] = 0$ and $\|\mathbf{M}_k\| \leq L$ for each $1 \leq k \leq n$, for some fixed value L .*

Let $\mathbf{M} = \sum \mathbf{M}_k$ and let $v(\mathbf{M}) = \max\{\|\mathbb{E}[\mathbf{M}\mathbf{M}^\top]\|, \|\mathbb{E}[\mathbf{M}^\top \mathbf{M}]\|\}$ denote the matrix variance statistic of \mathbf{M} . Then for all $t \geq 0$, we have

$$\mathbb{P}(\|\mathbf{M}\| \geq t) \leq (m_1 + m_2) \exp\left(\frac{-t^2/2}{v(\mathbf{M}) + Lt/3}\right). \quad (37)$$

We apply this as follows: given $r \in \{1, \dots, k\}$, define a matrix $\mathbf{T}_{ij}^{(r)} \in \mathbb{R}^{n \times (n_1 + \dots + n_k)}$ for each $i \in \{1, \dots, n\}$ and $j \in \{1, \dots, n_r\}$ whose $(i, n_1 + \dots + n_{r-1} + j)$ th entry is equal to $\mathbf{A}_{ij}^{(r)} - \mathbf{P}_{ij}^{(r)}$, with all other entries equal to 0 (in other words, if we divide $\mathbf{T}_{ij}^{(r)}$ into distinct $n \times n_t$ blocks, then the r th block is the only non-zero one, and within this block only the (i, j) th entry is non-zero).

We then define matrices $\mathbf{M}_{ij}^{(r)}$ for each $r \in \{1, \dots, k\}$ as follows:

- If $\mathbf{A}^{(r)}$ is bipartite, we define $\mathbf{M}_{ij}^{(r)} = \mathbf{T}_{ij}^{(r)}$ for all $i \in \{1, \dots, n\}$ and $j \in \{1, \dots, n_r\}$;
- If $\mathbf{A}^{(r)}$ is directed, we define $\mathbf{M}_{ij}^{(r)} = \mathbf{T}_{ij}^{(r)}$ for all $i, j \in \{1, \dots, n_r\}$ with $i \neq j$;
- If $\mathbf{A}^{(r)}$ is undirected, we define $\mathbf{M}_{ij}^{(r)} = \mathbf{T}_{ij}^{(r)} + \mathbf{T}_{ji}^{(r)}$ for all $i, j \in \{1, \dots, n_r\}$ with $i < j$.

Observe that $\|\mathbf{M}_{ij}^{(r)}\| = |\mathbf{A}_{ij}^{(r)} - \mathbf{P}_{ij}^{(r)}| < 1$, and by definition $\mathbb{E}[\mathbf{M}_{ij}^{(r)}] = 0$, so the matrix sum $\mathbf{M} = \sum_{i,j,r} \mathbf{M}_{ij}^{(r)}$ satisfies the criteria for Bernstein's Theorem (where we sum over the variables i, j and r according to the cases above). To bound the matrix variance statistic $v(\mathbf{M})$, let $\mathbf{M}_r = \sum_{i,j} \mathbf{M}_{ij}^{(r)}$, and note that

$$\mathbf{M}_r \mathbf{M}_r^\top = \begin{cases} \sum_{l=1}^{n_r} (\mathbf{A}_{il}^{(r)} - \mathbf{P}_{il}^{(r)}) (\mathbf{A}_{jl}^{(r)} - \mathbf{P}_{jl}^{(r)}) & \text{if } \mathbf{A}^{(r)} \text{ is bipartite} \\ \sum_{l \neq i,j} (\mathbf{A}_{il}^{(r)} - \mathbf{P}_{il}^{(r)}) (\mathbf{A}_{jl}^{(r)} - \mathbf{P}_{jl}^{(r)}) & \text{otherwise,} \end{cases} \quad (38)$$

and therefore that

$$\mathbb{E}[\mathbf{M}_r \mathbf{M}_r^\top]_{ij} = \begin{cases} \sum_{l=1}^{n_r} \mathbf{P}_{il}^{(r)} (1 - \mathbf{P}_{il}^{(r)}) & \text{if } \mathbf{A}^{(r)} \text{ is bipartite and } i = j \\ \sum_{l \neq i} \mathbf{P}_{il}^{(r)} (1 - \mathbf{P}_{il}^{(r)}) & \text{if } \mathbf{A}^{(r)} \text{ is not bipartite and } i = j \\ 0 & \text{if } i \neq j. \end{cases} \quad (39)$$

By definition, $\mathbf{P}_{il}^{(r)}(1 - \mathbf{P}_{il}^{(r)}) \leq \epsilon_r \rho$ for all i and l , and so since $\mathbb{E}[\mathbf{M}_r \mathbf{M}_r^\top]$ is diagonal, we see that $\|\mathbb{E}[\mathbf{M}_r \mathbf{M}_r^\top]\| \leq \epsilon_r \rho n_r$. Since $\mathbf{M} \mathbf{M}^\top = \sum \mathbf{M}_r \mathbf{M}_r^\top$ and the \mathbf{M}_r are independent, it follows that $\|\mathbb{E}[\mathbf{M} \mathbf{M}^\top]\| \leq (\epsilon_1 n_1 + \dots + \epsilon_k n_k) \rho \leq \epsilon \rho k^{1/2} n_{\max}$ by the Cauchy–Schwarz inequality.

Similarly,

$$[\mathbf{M}_r^\top \mathbf{M}_r]_{ij} = \begin{cases} \sum_{l=1}^n (\mathbf{A}_{li}^{(r)} - \mathbf{P}_{li}^{(r)}) (\mathbf{A}_{lj}^{(r)} - \mathbf{P}_{lj}^{(r)}) & \text{if } \mathbf{A}^{(r)} \text{ is bipartite} \\ \sum_{l \neq i, j} (\mathbf{A}_{li}^{(r)} - \mathbf{P}_{li}^{(r)}) (\mathbf{A}_{lj}^{(r)} - \mathbf{P}_{lj}^{(r)}) & \text{otherwise,} \end{cases} \quad (40)$$

and so

$$\mathbb{E}[\mathbf{M}_r^\top \mathbf{M}_r]_{ij} = \begin{cases} \sum_{l=1}^n \mathbf{P}_{li}^{(r)} (1 - \mathbf{P}_{li}^{(r)}) & \text{if } \mathbf{A}^{(r)} \text{ is bipartite and } i = j \\ \sum_{l \neq i} \mathbf{P}_{li}^{(r)} (1 - \mathbf{P}_{li}^{(r)}) & \text{if } \mathbf{A}^{(r)} \text{ is not bipartite and } i = j \\ 0 & \text{if } i \neq j, \end{cases} \quad (41)$$

and as before we see that $\|\mathbb{E}[\mathbf{M}_r^\top \mathbf{M}_r]\| \leq \epsilon_r \rho n$, while $\|\mathbb{E}[\mathbf{M}_s^\top \mathbf{M}_t]\| = 0$ if $s \neq t$ since the matrices \mathbf{M}_r are independent. Thus the matrix $\mathbb{E}[\mathbf{M}^\top \mathbf{M}]$ is block diagonal, and so it follows that $\|\mathbb{E}[\mathbf{M}^\top \mathbf{M}]\| \leq \rho n$, and so certainly $v(\mathbf{M}) = O(\epsilon \rho k^{1/2} n)$.

Substituting these into Bernstein’s Theorem and rearranging, we find that, for any $t \geq 0$,

$$\mathbb{P}(\|\mathbf{M}\| \geq t) \leq (n + n_1 + \dots + n_k) \exp\left(\frac{-3t^2}{6\rho k^{1/2} n + 2t}\right). \quad (42)$$

The numerator of the exponential term dominates for n sufficiently large if $t = O(\epsilon^{1/2} \rho^{1/2} k^{1/4} n^{1/2} \log^{1/2}(n))$, and so $\|\mathbf{M}\| = O(\epsilon^{1/2} \rho^{1/2} k^{1/4} n^{1/2} \log^{1/2}(n))$ almost surely.

Finally, we note that $\mathbf{M} = \mathbf{A} - \mathbf{P} + \mathbf{P}_0$, where \mathbf{P}_0 comprises k distinct blocks of sizes $n \times n_r$, with the r th such block either identically zero or diagonal, with entries equal to the diagonal entries of $\mathbf{P}^{(r)}$, depending on whether or not $\mathbf{A}^{(r)}$ is bipartite, and satisfies $\|\mathbf{P}_0\| \leq \epsilon \rho$. The result then follows from subadditivity of the spectral norm and integrating over all possible choices of latent positions. \square

Corollary 10. *The leading d singular values of \mathbf{A} satisfy*

$$\frac{\sigma_i(\mathbf{A})}{\rho n} \rightarrow \sqrt{\lambda_i(\Delta_X \Lambda_* \Delta_Y \Lambda_*^\top)} \quad (43)$$

almost surely, where $\Delta_Y = c_1 \Delta_{Y,1} \oplus \dots \oplus c_\kappa \Delta_{Y,\kappa}$, and consequently $\sigma_i(\mathbf{A}) = O(\epsilon \rho n)$ and $\sigma_i(\mathbf{A}) = \Omega(\rho n)$ almost surely.

Proof. A corollary of Weyl’s inequalities (see, for example, [17], Corollary 7.3.5) states that $|\sigma_i(\mathbf{A}) - \sigma_i(\mathbf{P})| \leq \|\mathbf{A} - \mathbf{P}\|$, and so in particular (applying the reverse triangle inequality where necessary)

$$\sigma_i(\mathbf{P}) - \|\mathbf{A} - \mathbf{P}\| \leq \sigma_i(\mathbf{A}) \leq \sigma_i(\mathbf{P}) + \|\mathbf{A} - \mathbf{P}\|. \quad (44)$$

The result then follows from Propositions 7 and 8. \square

The next few results provide bounds on the asymptotic growth of a number of residual terms in the proofs of our main theorems. While the proofs are similar in nature to a number of results in [23], there are some minor differences to account for the fact that the matrices \mathbf{A} and \mathbf{P} are not symmetric, and so we reproduce them in full. We begin an analogue of a bound appearing in Lemma 17 of [23]:

Proposition 11. $\|\mathbf{U}_{\mathbf{P}}^{\top}(\mathbf{A} - \mathbf{P})\mathbf{V}_{\mathbf{P}}\|_F = O(\log^{1/2}(n))$ almost surely.

Proof. Condition on some choice of latent positions. For any $i, j \in \{1, \dots, d\}$ and $r \in \{1, \dots, k\}$, let

$$\mathbf{E}_{ij}^{(r)} = \begin{cases} \sum_{p=1}^n \sum_{q=1}^{n_r} u_p v_q^{(r)} (\mathbf{A}_{pq}^{(r)} - \mathbf{P}_{pq}^{(r)}) & \text{if } \mathbf{A}^{(r)} \text{ is bipartite} \\ \sum_{p \neq q} u_p v_q^{(r)} (\mathbf{A}_{pq}^{(r)} - \mathbf{P}_{pq}^{(r)}) & \text{if } \mathbf{A}^{(r)} \text{ is directed} \\ \sum_{p < q} (u_p v_q^{(r)} + u_q v_p^{(r)}) (\mathbf{A}_{pq}^{(r)} - \mathbf{P}_{pq}^{(r)}) & \text{if } \mathbf{A}^{(r)} \text{ is undirected} \end{cases} \quad (45)$$

and

$$\mathbf{F}_{ij}^{(r)} = \begin{cases} \mathbf{0} & \text{if } \mathbf{A}^{(r)} \text{ is bipartite} \\ \sum_{p=1}^n u_p v_p^{(r)} \mathbf{P}_{pp}^{(r)} & \text{otherwise,} \end{cases} \quad (46)$$

where u and $v^{(r)}$ denote the i th and j th columns of $\mathbf{U}_{\mathbf{P}}$ and $\mathbf{V}_{\mathbf{P}}^{(r)}$ respectively, so that

$$(\mathbf{U}_{\mathbf{P}}^{\top}(\mathbf{A} - \mathbf{P})\mathbf{V}_{\mathbf{P}})_{ij} = \sum_{r=1}^k \mathbf{E}_{ij}^{(r)} - \sum_{r=1}^k \mathbf{F}_{ij}^{(r)}. \quad (47)$$

We can bound the latter term by applying the Cauchy–Schwarz inequality to find that

$$\left| \sum_{r=1}^k \mathbf{F}_{ij}^{(r)} \right| \leq \left(\sum_{r=1}^k \sum_{p=1}^n |u_p \mathbf{P}_{pp}^{(r)}|^2 \right)^{1/2} \left(\sum_{r=1}^k \sum_{p=1}^n |v_p^{(r)}|^2 \right)^{1/2} = O(\epsilon \rho), \quad (48)$$

and thus can be discounted in our asymptotic analysis.

Each of the $\mathbf{E}_{ij}^{(r)}$ is a sum of independent zero-mean random variables, with each of the individual terms bounded in absolute value by $|u_p v_q^{(r)}|$ in the bipartite and directed cases and $|u_p v_q^{(r)} + u_q v_p^{(r)}|$ in the undirected case. Applying Hoeffding’s inequality, we thus observe that

$$\mathbb{P} \left(\left| \sum_{r=1}^k \mathbf{E}_{ij}^{(r)} \right| > t \right) \leq 2 \exp \left(\frac{-2t^2}{4 \left(\sum_{r_1} \sum_{p,q} |u_p v_q^{(r_1)}|^2 + \sum_{r_2} \sum_{p < q} |u_p v_q^{(r_2)} + u_q v_p^{(r_2)}|^2 \right)} \right), \quad (49)$$

where r_1 sums over the bipartite and directed cases and r_2 sums over the undirected cases.

Note that $|u_p v_q^{(r)} + u_q v_p^{(r)}|^2 \leq |u_p v_q^{(r)}|^2 + |u_q v_p^{(r)}|^2 + 2|u_p u_q v_q^{(r)} v_p^{(r)}|$, and so

$$\mathbb{P} \left(\left| \sum_{r=1}^k \mathbf{E}_{ij}^{(r)} \right| > t \right) \leq 2 \exp \left(\frac{-2t^2}{4 \left(\sum_{r=1}^k \sum_{p,q} |u_p v_q^{(r)}|^2 + 2 \sum_{r_2} \sum_{p < q} |u_p u_q v_q^{(r_2)} v_p^{(r_2)}| \right)} \right). \quad (50)$$

Both summands are at most 1; the first is clear, while for the second we apply the Cauchy–Schwarz inequality and note that

$$\sum_{r_2} \sum_{p < q} |u_p u_q v_q^{(r_2)} v_p^{(r_2)}| \leq \left(\sum_{r_3} \sum_{p < q} |u_p v_q^{(r_3)}|^2 \right)^{1/2} \left(\sum_{r_2} \sum_{p < q} |u_q v_p^{(r_2)}|^2 \right)^{1/2} \quad (51)$$

$$\leq \left(\sum_{r_2} \sum_{p, q} |u_p|^2 |v_q^{(r_2)}|^2 \right) = 1. \quad (52)$$

Thus $\sum_{r=1}^k \mathbf{E}_{ij}^{(r)} = O(\log^{1/2}(n))$ almost surely, and the result follows after integrating over all possible choices of latent positions. \square

Before establishing the next set of bounds (which relate to the left and right singular vectors of the matrices \mathbf{A} and \mathbf{P}), we state the following variation of the Davis–Kahan theorem (see [41], Theorem 4):

Theorem 12 (Variant of Davis–Kahan). *Let $\mathbf{M}_1, \mathbf{M}_2 \in \mathbb{R}^{m \times n}$ have singular value decompositions*

$$\mathbf{M}_i = \mathbf{U}_i \boldsymbol{\Sigma}_i \mathbf{V}_i^\top + \mathbf{U}_{i,\perp} \boldsymbol{\Sigma}_{i,\perp} \mathbf{V}_{i,\perp}^\top, \quad (53)$$

where $\mathbf{U}_i \in \mathbb{O}(m \times d)$ has orthonormal columns corresponding to the d greatest singular values of \mathbf{M}_i , for some $1 \leq d \leq n$. Then, if $|\sigma_d(\mathbf{M}_1)^2 - \sigma_{d+1}(\mathbf{M}_1)^2| > 0$, we have

$$\|\sin \Theta(\mathbf{U}_2, \mathbf{U}_1)\| \leq \frac{2\sqrt{d}(2\sigma_1(\mathbf{M}_1) + \|\mathbf{M}_2 - \mathbf{M}_1\|) \|\mathbf{M}_2 - \mathbf{M}_1\|}{\sigma_d(\mathbf{M}_1)^2 - \sigma_{d+1}(\mathbf{M}_1)^2}. \quad (54)$$

where we take $\sigma_{n+1}(\mathbf{M}_1) = -\infty$.

We note that the same inequality holds for $\|\sin \Theta(\mathbf{V}_2, \mathbf{V}_1)\|$ since the right-hand side of (54) is invariant under matrix transposition. Using this result, we can prove the following:

Proposition 13. *The following bounds hold almost surely:*

- (i) $\|\mathbf{U}_\mathbf{A} \mathbf{U}_\mathbf{A}^\top - \mathbf{U}_\mathbf{P} \mathbf{U}_\mathbf{P}^\top\|, \|\mathbf{V}_\mathbf{A} \mathbf{V}_\mathbf{A}^\top - \mathbf{V}_\mathbf{P} \mathbf{V}_\mathbf{P}^\top\| = O\left(\frac{\epsilon^{3/2} k^{1/4} \log^{1/2}(n)}{\rho^{1/2} n^{1/2}}\right);$
- (ii) $\|\mathbf{U}_\mathbf{A} - \mathbf{U}_\mathbf{P} \mathbf{U}_\mathbf{P}^\top \mathbf{U}_\mathbf{A}\|_F, \|\mathbf{V}_\mathbf{A} - \mathbf{V}_\mathbf{P} \mathbf{V}_\mathbf{P}^\top \mathbf{V}_\mathbf{A}\|_F = O\left(\frac{\epsilon^{3/2} k^{1/4} \log^{1/2}(n)}{\rho^{1/2} n^{1/2}}\right);$
- (iii) $\|\mathbf{U}_\mathbf{P}^\top \mathbf{U}_\mathbf{A} \boldsymbol{\Sigma}_\mathbf{A} - \boldsymbol{\Sigma}_\mathbf{P} \mathbf{V}_\mathbf{P}^\top \mathbf{V}_\mathbf{A}\|_F, \|\boldsymbol{\Sigma}_\mathbf{P} \mathbf{U}_\mathbf{P}^\top \mathbf{U}_\mathbf{A} - \mathbf{V}_\mathbf{P}^\top \mathbf{V}_\mathbf{A} \boldsymbol{\Sigma}_\mathbf{A}\|_F = O(\epsilon^2 k^{1/2} \log(n));$
- (iv) $\|\mathbf{U}_\mathbf{P}^\top \mathbf{U}_\mathbf{A} - \mathbf{V}_\mathbf{P}^\top \mathbf{V}_\mathbf{A}\|_F = O\left(\frac{\epsilon^2 k^{1/2} \log(n)}{\rho n}\right)$

Proof.

- (i) Let $\sigma_1, \dots, \sigma_d$ denote the singular values of $\mathbf{U}_\mathbf{P}^\top \mathbf{U}_\mathbf{A}$, and let $\theta_i = \cos^{-1}(\sigma_i)$ be the principal angles. It is a standard result that the non-zero eigenvalues of the matrix $\mathbf{U}_\mathbf{A} \mathbf{U}_\mathbf{A}^\top - \mathbf{U}_\mathbf{P} \mathbf{U}_\mathbf{P}^\top$ are precisely the $\sin(\theta_i)$ (each occurring twice) and so by Davis–Kahan we have

$$\|\mathbf{U}_\mathbf{A} \mathbf{U}_\mathbf{A}^\top - \mathbf{U}_\mathbf{P} \mathbf{U}_\mathbf{P}^\top\| = \max_{i \in \{1, \dots, d\}} |\sin(\theta_i)| \leq \frac{2\sqrt{d}(2\sigma_1(\mathbf{P}) + \|\mathbf{A} - \mathbf{P}\|) \|\mathbf{A} - \mathbf{P}\|}{\sigma_d(\mathbf{P})^2} \quad (55)$$

for n sufficiently large.

Applying the bounds from Propositions 7 and 8 then shows that

$$\|\mathbf{U}_\mathbf{A} \mathbf{U}_\mathbf{A}^\top - \mathbf{U}_\mathbf{P} \mathbf{U}_\mathbf{P}^\top\| = O\left(\frac{\epsilon^{3/2} k^{1/4} \log^{1/2}(n)}{\rho^{1/2} n^{1/2}}\right). \quad (56)$$

An identical argument gives the result for $\|\mathbf{V}_\mathbf{A} \mathbf{V}_\mathbf{A}^\top - \mathbf{V}_\mathbf{P} \mathbf{V}_\mathbf{P}^\top\|$.

- (ii) Using the bound from part (i), we find that

$$\|\mathbf{U}_\mathbf{A} - \mathbf{U}_\mathbf{P} \mathbf{U}_\mathbf{P}^\top \mathbf{U}_\mathbf{A}\|_F = \|(\mathbf{U}_\mathbf{A} \mathbf{U}_\mathbf{A}^\top - \mathbf{U}_\mathbf{P} \mathbf{U}_\mathbf{P}^\top) \mathbf{U}_\mathbf{A}\|_F = O\left(\frac{\epsilon^{3/2} k^{1/4} \log^{1/2}(n)}{\rho^{1/2} n^{1/2}}\right). \quad (57)$$

An identical argument bounds the term $\|\mathbf{V}_\mathbf{A} - \mathbf{V}_\mathbf{P} \mathbf{V}_\mathbf{P}^\top \mathbf{V}_\mathbf{A}\|_F$.

(iii) Observe that

$$\mathbf{U}_P^\top \mathbf{U}_A \Sigma_A - \Sigma_P \mathbf{V}_P^\top \mathbf{V}_A = \mathbf{U}_P^\top (\mathbf{A} - \mathbf{P}) \mathbf{V}_A \quad (58)$$

and that we may rewrite the right-hand term to find that

$$\begin{aligned} \mathbf{U}_P^\top \mathbf{U}_A \Sigma_A - \Sigma_P \mathbf{V}_P^\top \mathbf{V}_A &= \mathbf{U}_P^\top (\mathbf{A} - \mathbf{P}) (\mathbf{V}_A - \mathbf{V}_P \mathbf{V}_P^\top \mathbf{V}_A) \\ &\quad + \mathbf{U}_P^\top (\mathbf{A} - \mathbf{P}) \mathbf{V}_P \mathbf{V}_P^\top \mathbf{V}_A. \end{aligned} \quad (59)$$

These terms satisfy

$$\|\mathbf{U}_P^\top (\mathbf{A} - \mathbf{P}) (\mathbf{V}_A - \mathbf{V}_P \mathbf{V}_P^\top \mathbf{V}_A)\|_F = O(\epsilon^2 k^{1/2} \log(n)) \quad (60)$$

and

$$\mathbf{U}_P^\top (\mathbf{A} - \mathbf{P}) \mathbf{V}_P \mathbf{V}_P^\top \mathbf{V}_A = O(\log^{1/2}(n)) \quad (61)$$

by Propositions 8, 11 and the result from part (ii), and thus

$$\|\mathbf{U}_P^\top \mathbf{U}_A \Sigma_A - \Sigma_P \mathbf{V}_P^\top \mathbf{V}_A\|_F = O(\epsilon^2 k^{1/2} \log(n)). \quad (62)$$

An identical argument bounds the term $\|\Sigma_P \mathbf{U}_P^\top \mathbf{U}_A - \mathbf{V}_P^\top \mathbf{V}_A \Sigma_A\|_F$.

(iv) Note that

$$\begin{aligned} \mathbf{U}_P^\top \mathbf{U}_A - \mathbf{V}_P^\top \mathbf{V}_A &= ((\mathbf{U}_P^\top \mathbf{U}_A \Sigma_A - \Sigma_P \mathbf{V}_P^\top \mathbf{V}_A) + (\Sigma_P \mathbf{U}_P^\top \mathbf{U}_A \\ &\quad - \mathbf{V}_P^\top \mathbf{V}_A \Sigma_A)) \Sigma_A^{-1} - \Sigma_P (\mathbf{U}_P^\top \mathbf{U}_A - \mathbf{V}_P^\top \mathbf{V}_A) \Sigma_A^{-1}. \end{aligned} \quad (63)$$

For any i, j we find (after rearranging and bounding the absolute value of the right-hand terms by the Frobenius norm):

$$\begin{aligned} |(\mathbf{U}_P^\top \mathbf{U}_A - \mathbf{V}_P^\top \mathbf{V}_A)_{ij}| \left(1 + \frac{\sigma_i(\mathbf{P})}{\sigma_j(\mathbf{A})}\right) &\leq (\|\mathbf{U}_P^\top \mathbf{U}_A \Sigma_A - \Sigma_P \mathbf{V}_P^\top \mathbf{V}_A\|_F \\ &\quad + \|\Sigma_P \mathbf{U}_P^\top \mathbf{U}_A - \mathbf{V}_P^\top \mathbf{V}_A \Sigma_A\|_F) \|\Sigma_A^{-1}\|_F \end{aligned} \quad (64)$$

where we have used the result from part (iii) and Corollary 10.

Consequently, we find that

$$|(\mathbf{U}_P^\top \mathbf{U}_A - \mathbf{V}_P^\top \mathbf{V}_A)_{ij}| = O\left(\frac{\epsilon^2 k^{1/2} \log(n)}{\rho n}\right) \quad (65)$$

by noting that $\left(1 + \frac{\sigma_i(\mathbf{P})}{\sigma_j(\mathbf{A})}\right) \geq 1$.

□

The following result (an analogue of [23], Proposition 16) relates to orthogonal matrix \mathbf{W} used to perform a simultaneous Procrustes alignment of \mathbf{X}_A with \mathbf{X}_P and \mathbf{Y}_A with \mathbf{Y}_P .

Proposition 14. *Let $\mathbf{U}_P^\top \mathbf{U}_A + \mathbf{V}_P^\top \mathbf{V}_A$ admit the singular value decomposition*

$$\mathbf{U}_P^\top \mathbf{U}_A + \mathbf{V}_P^\top \mathbf{V}_A = \mathbf{W}_1 \Sigma \mathbf{W}_2^\top, \quad (66)$$

and let $\mathbf{W} = \mathbf{W}_1 \mathbf{W}_2^\top$. Then

$$\max\{\|\mathbf{U}_P^\top \mathbf{U}_A - \mathbf{W}\|_F, \|\mathbf{V}_P^\top \mathbf{V}_A - \mathbf{W}\|_F\} = O\left(\frac{\epsilon^3 k^{1/2} \log(n)}{\rho n}\right) \quad (67)$$

almost surely.

Proof. A standard argument shows that \mathbf{W} minimises the term $\|\mathbf{U}_P^\top \mathbf{U}_A - \mathbf{Q}\|_F^2 + \|\mathbf{V}_P^\top \mathbf{V}_A - \mathbf{Q}\|_F^2$ among all $\mathbf{Q} \in \mathbb{O}(d)$. Let $\mathbf{U}_P^\top \mathbf{U}_A = \mathbf{W}_{U,1} \boldsymbol{\Sigma}_U \mathbf{W}_{U,2}^\top$ be the singular value decomposition of $\mathbf{U}_P^\top \mathbf{U}_A$, and define $\mathbf{W}_U \in \mathbb{O}(d)$ by $\mathbf{W}_U = \mathbf{W}_{U,1} \mathbf{W}_{U,2}^\top$. Then

$$\begin{aligned} \|\mathbf{U}_P^\top \mathbf{U}_A - \mathbf{W}_U\|_F &= \|\boldsymbol{\Sigma} - \mathbf{I}\|_F = \left(\sum_{i=1}^d (1 - \sigma_i)^2 \right)^{1/2} \leq \sum_{i=1}^d (1 - \sigma_i) \\ &\leq \sum_{i=1}^d (1 - \sigma_i^2) = \sum_{i=1}^d \sin^2(\theta_i) \leq d \|\mathbf{U}_A \mathbf{U}_A^\top - \mathbf{U}_P \mathbf{U}_P^\top\|^2 \end{aligned} \quad (68)$$

and so

$$\|\mathbf{U}_P^\top \mathbf{U}_A - \mathbf{W}_U\|_F = O\left(\frac{\epsilon^3 k^{1/2} \log(n)}{\rho n}\right). \quad (69)$$

Also,

$$\|\mathbf{V}_P^\top \mathbf{V}_A - \mathbf{W}_U\|_F \leq \|\mathbf{V}_P^\top \mathbf{V}_A - \mathbf{U}_P^\top \mathbf{U}_A\|_F + \|\mathbf{U}_P^\top \mathbf{U}_A - \mathbf{W}_U\|_F \quad (70)$$

and so

$$\|\mathbf{V}_P^\top \mathbf{V}_A - \mathbf{W}_U\|_F = O\left(\frac{\epsilon^3 k^{1/2} \log(n)}{\rho n}\right) \quad (71)$$

by Proposition 13.

Combining these shows that

$$\|\mathbf{U}_P^\top \mathbf{U}_A - \mathbf{W}\|_F^2 + \|\mathbf{V}_P^\top \mathbf{V}_A - \mathbf{W}\|_F^2 \leq \|\mathbf{U}_P^\top \mathbf{U}_A - \mathbf{W}_U\|_F^2 + \|\mathbf{V}_P^\top \mathbf{V}_A - \mathbf{W}_U\|_F^2 \quad (72)$$

and thus

$$\max\{\|\mathbf{U}_P^\top \mathbf{U}_A - \mathbf{W}\|_F, \|\mathbf{V}_P^\top \mathbf{V}_A - \mathbf{W}\|_F\} = O\left(\frac{\epsilon^3 k^{1/2} \log(n)}{\rho n}\right) \quad (73)$$

as required. \square

The following bounds are a straightforward adaptation of [23], Lemma 17:

Proposition 15. *The following bounds hold almost surely:*

- (i) $\|\mathbf{W} \boldsymbol{\Sigma}_A - \boldsymbol{\Sigma}_P \mathbf{W}\|_F = O(\epsilon^4 k^{1/2} \log(n));$
- (ii) $\|\mathbf{W} \boldsymbol{\Sigma}_A^{1/2} - \boldsymbol{\Sigma}_P^{1/2} \mathbf{W}\|_F = O\left(\frac{\epsilon^4 k^{1/2} \log(n)}{\rho^{1/2} n^{1/2}}\right);$
- (iii) $\|\mathbf{W} \boldsymbol{\Sigma}_A^{-1/2} - \boldsymbol{\Sigma}_P^{-1/2} \mathbf{W}\|_F = O\left(\frac{\epsilon^4 k^{1/2} \log(n)}{\rho^{3/2} n^{3/2}}\right).$

Proof.

- (i) Observe that

$$\mathbf{W} \boldsymbol{\Sigma}_A - \boldsymbol{\Sigma}_P \mathbf{W} = (\mathbf{W} - \mathbf{U}_P^\top \mathbf{U}_A) \boldsymbol{\Sigma}_A + \mathbf{U}_P^\top \mathbf{U}_A \boldsymbol{\Sigma}_A - \boldsymbol{\Sigma}_P \mathbf{W} \quad (74)$$

and that the right-hand expression may be rewritten as

$$(\mathbf{W} - \mathbf{U}_P^\top \mathbf{U}_A) \boldsymbol{\Sigma}_A + (\mathbf{U}_P^\top \mathbf{U}_A \boldsymbol{\Sigma}_A - \boldsymbol{\Sigma}_P \mathbf{V}_P^\top \mathbf{V}_A) + \boldsymbol{\Sigma}_P (\mathbf{V}_P^\top \mathbf{V}_A - \mathbf{W}). \quad (75)$$

The terms $\|(\mathbf{W} - \mathbf{U}_P^\top \mathbf{U}_A) \boldsymbol{\Sigma}_A\|_F$ and $\|\boldsymbol{\Sigma}_P (\mathbf{V}_P^\top \mathbf{V}_A - \mathbf{W})\|_F$ are both $O(\epsilon^4 k^{1/2} \log(n))$ (as shown by Propositions 7, 14 and Corollary 10), while the term $\|\mathbf{U}_P^\top \mathbf{U}_A \boldsymbol{\Sigma}_A - \boldsymbol{\Sigma}_P \mathbf{V}_P^\top \mathbf{V}_A\|_F$ is $O(\epsilon^2 k^{1/2} \log(n))$ by Proposition 13, and so $\|\mathbf{W} \boldsymbol{\Sigma}_A - \boldsymbol{\Sigma}_P \mathbf{W}\|_F = O(\epsilon^4 k^{1/2} \log(n))$ as required.

(ii) Note that

$$(\mathbf{W}\Sigma_{\mathbf{A}}^{1/2} - \Sigma_{\mathbf{P}}^{1/2}\mathbf{W})_{ij} = \mathbf{W}_{ij}(\sigma_j(\mathbf{A})^{1/2} - \sigma_i(\mathbf{P})^{1/2}) \quad (76)$$

$$= \frac{\mathbf{W}_{ij}(\sigma_j(\mathbf{A}) - \sigma_i(\mathbf{P}))}{\sigma_j(\mathbf{A})^{1/2} + \sigma_i(\mathbf{P})^{1/2}} \quad (77)$$

$$= \frac{(\mathbf{W}\Sigma_{\mathbf{A}} - \Sigma_{\mathbf{P}}\mathbf{W})_{ij}}{\sigma_j(\mathbf{A})^{1/2} + \sigma_i(\mathbf{P})^{1/2}}, \quad (78)$$

and so we find that $\|\mathbf{W}\Sigma_{\mathbf{A}}^{1/2} - \Sigma_{\mathbf{P}}^{1/2}\mathbf{W}\|_F = O\left(\frac{\epsilon^4 k^{1/2} \log(n)}{\rho^{1/2} n^{1/2}}\right)$ by applying part (i) and summing over all $i, j \in \{1, \dots, d\}$.

(iii) Note that

$$(\mathbf{W}\Sigma_{\mathbf{A}}^{-1/2} - \Sigma_{\mathbf{P}}^{-1/2}\mathbf{W})_{ij} = \frac{\mathbf{W}_{ij}(\sigma_i(\mathbf{P})^{1/2} - \sigma_j(\mathbf{A})^{1/2})}{\sigma_i(\mathbf{P})^{1/2}\sigma_j(\mathbf{A})^{1/2}} \quad (79)$$

$$= \frac{(\mathbf{W}\Sigma_{\mathbf{A}}^{1/2} - \Sigma_{\mathbf{P}}^{1/2}\mathbf{W})_{ij}}{\sigma_i(\mathbf{P})^{1/2}\sigma_j(\mathbf{A})^{1/2}} \quad (80)$$

and so we find that $\|\mathbf{W}\Sigma_{\mathbf{A}}^{-1/2} - \Sigma_{\mathbf{P}}^{-1/2}\mathbf{W}\|_F = O\left(\frac{\epsilon^4 k^{1/2} \log(n)}{\rho^{3/2} n^{3/2}}\right)$ by applying part (ii) and summing over all $i, j \in \{1, \dots, d\}$. \square

The next results establish the existence of, and some properties relating to, the matrices $\tilde{\mathbf{L}}$ and $\tilde{\mathbf{R}}_r$ which map the latent position matrices \mathbf{X} and $\mathbf{Y}^{(r)}$ to the embeddings $\mathbf{X}_{\mathbf{P}}$ and $\mathbf{Y}_{\mathbf{P}}^{(r)}$ respectively. In particular, where they exist, we bound the growth of the spectral norm of the inverses (or pseudo-inverses where appropriate) of these matrices, which is necessary in order to be able to recover the latent positions from the UASE.

Proposition 16. *If \mathbf{X} is of rank d then there exists a matrix $\tilde{\mathbf{L}} \in \text{GL}(d)$ such that $\mathbf{X}_{\mathbf{P}} = \mathbf{X}\tilde{\mathbf{L}}$. In addition, if $\mathbf{Y}^{(r)}$ is of rank d_r then $\mathbf{Y}_{\mathbf{P}}^{(r)} = \mathbf{Y}^{(r)}\tilde{\mathbf{R}}_r$, where the matrix $\tilde{\mathbf{R}}_r \in \mathbb{R}^{d_r \times d}$ satisfies $\tilde{\mathbf{L}}\tilde{\mathbf{R}}_r^\top = \Lambda_{\epsilon, r}$. In particular, $\text{rank}(\tilde{\mathbf{R}}_r) = \text{rank}(\Lambda_r)$.*

Proof. Let $\Pi_{\mathbf{X}}, \Pi_{\mathbf{Y}} \in \text{GL}(d)$ satisfy $\Pi_{\mathbf{X}} = (\mathbf{X}^\top \mathbf{X})^{1/2}$ and $\Pi_{\mathbf{Y}} = (\Lambda_\epsilon \mathbf{Y}^\top \mathbf{Y} \Lambda_\epsilon^\top)^{1/2}$, where we take the unique positive-definite matrix square root in each case.

Observe that

$$(\mathbf{X}_{\mathbf{P}}\Sigma_{\mathbf{P}}^{1/2})(\mathbf{X}_{\mathbf{P}}\Sigma_{\mathbf{P}}^{1/2})^\top = \mathbf{U}_{\mathbf{P}}\Sigma_{\mathbf{P}}^2\mathbf{U}_{\mathbf{P}}^\top = \mathbf{P}\mathbf{P}^\top = (\mathbf{X}\Pi_{\mathbf{Y}})(\mathbf{X}\Pi_{\mathbf{Y}})^\top \quad (81)$$

and similarly (noting that $\mathbf{U}_{\mathbf{P}}\Sigma_{\mathbf{P}}\mathbf{V}_{\mathbf{P}}^{(r)\top} = \mathbf{P}^{(r)}$) that

$$(\mathbf{Y}_{\mathbf{P}}^{(r)}\Sigma_{\mathbf{P}}^{1/2})(\mathbf{Y}_{\mathbf{P}}^{(r)}\Sigma_{\mathbf{P}}^{1/2})^\top = \mathbf{P}^{(r)\top}\mathbf{P}^{(r)} = (\mathbf{Y}^{(r)}\Lambda_{\epsilon, r}^\top\Pi_{\mathbf{X}})(\mathbf{Y}^{(r)}\Lambda_{\epsilon, r}\Pi_{\mathbf{X}})^\top. \quad (82)$$

Thus there exist orthogonal matrices $\mathbf{Q}_{\mathbf{P}}, \mathbf{Q}_{\mathbf{P}}^{(r)} \in \mathbb{O}(d)$ such that

$$\mathbf{X}_{\mathbf{P}}\Sigma_{\mathbf{P}}^{1/2} = \mathbf{X}\Pi_{\mathbf{Y}}\mathbf{Q}_{\mathbf{P}}, \quad \mathbf{Y}_{\mathbf{P}}^{(r)}\Sigma_{\mathbf{P}}^{1/2} = \mathbf{Y}^{(r)}\Lambda_{\epsilon, r}^\top\Pi_{\mathbf{X}}\mathbf{Q}_{\mathbf{P}}^{(r)} \quad (83)$$

and so

$$\tilde{\mathbf{L}} = \Pi_{\mathbf{Y}}\mathbf{Q}_{\mathbf{P}}\Sigma_{\mathbf{P}}^{-1/2} \in \text{GL}(d), \quad \tilde{\mathbf{R}}_r = \Lambda_{\epsilon, r}^\top\Pi_{\mathbf{X}}\mathbf{Q}_{\mathbf{P}}^{(r)}\Sigma_{\mathbf{P}}^{-1/2} \in \mathbb{R}^{d_r \times d} \quad (84)$$

are our desired matrices.

For the final statement, observe that

$$\mathbf{X}\tilde{\mathbf{L}}\tilde{\mathbf{R}}_r^\top\mathbf{Y}^{(r)\top} = \mathbf{X}_{\mathbf{P}}\mathbf{Y}_{\mathbf{P}}^{(r)\top} = \mathbf{P}^{(r)} = \mathbf{X}\Lambda_{\epsilon, r}\mathbf{Y}^{(r)\top}, \quad (85)$$

and so the result follows after multiplying by $(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{X}^\top$ and $\mathbf{Y}^{(r)}(\mathbf{Y}^{(r)\top}\mathbf{Y}^{(r)})^{-1}$ on the left and right respectively. \square

Corollary 17. *The matrices $\tilde{\mathbf{L}}$ and $\tilde{\mathbf{R}}_r$ satisfy $\|\tilde{\mathbf{L}}\| = O(\epsilon)$, $\|\tilde{\mathbf{L}}^{-1}\| = O(\epsilon^{1/2})$ and $\|\tilde{\mathbf{R}}_r\| = O(\epsilon_r)$ almost surely.*

Moreover, if the matrix $\mathbf{\Lambda}_r$ is of rank d_r , then $\|\tilde{\mathbf{R}}_r^+\| = O(\frac{1}{\epsilon_r})$ almost surely, where $\tilde{\mathbf{R}}_r^+ = \tilde{\mathbf{R}}_r^\top (\tilde{\mathbf{R}}_r \tilde{\mathbf{R}}_r^\top)^{-1}$ is the Moore-Penrose inverse of $\tilde{\mathbf{R}}_r$.

Proof. Proposition 7 shows us that $\|\Sigma_{\mathbf{P}}\| = O(\epsilon \rho n)$ and $\|\Sigma_{\mathbf{P}}^{-1}\| = O(\frac{1}{\rho n})$, and an identical line of reasoning shows that $\|\Pi_{\mathbf{\Lambda}_{\epsilon}, \mathbf{Y}}\| = O(\epsilon \rho^{1/2} n^{1/2})$ and $\|\Pi_{\mathbf{\Lambda}_{\epsilon}, \mathbf{Y}}^{-1}\| = O(\frac{1}{\rho^{1/2} n^{1/2}})$, and similarly that $\|\Pi_{\mathbf{X}}\| = O(\rho^{1/2} n^{1/2})$ and $\|\Pi_{\mathbf{X}}^{-1}\| = O(\frac{1}{\rho^{1/2} n^{1/2}})$.

The first three bounds then follow from submultiplicativity and unitary invariance of the spectral norm. For the final result, note that $\tilde{\mathbf{R}}_r \tilde{\mathbf{R}}_r^\top = \mathbf{\Lambda}_{\epsilon, r}^\top \Pi_{\mathbf{X}} \mathbf{Q}_{\mathbf{P}}^{(r)} \Sigma_{\mathbf{P}}^{-1} \mathbf{Q}_{\mathbf{P}}^{(r)\top} \Pi_{\mathbf{X}}^\top \mathbf{\Lambda}_{\epsilon, r}$, and so we see that

$$\sigma_{d_r}(\tilde{\mathbf{R}}_r \tilde{\mathbf{R}}_r^\top) = \lambda_{d_r}(\mathbf{Q}_{\mathbf{P}}^{(r)\top} \Pi_{\mathbf{X}}^\top \mathbf{\Lambda}_{\epsilon, r} \mathbf{\Lambda}_{\epsilon, r}^\top \Pi_{\mathbf{X}} \mathbf{Q}_{\mathbf{P}}^{(r)} \Sigma_{\mathbf{P}}^{-1}) \quad (86)$$

$$\geq \sigma_1(\mathbf{P})^{-1} \lambda_{d_r}(\mathbf{Q}_{\mathbf{P}}^{(r)\top} \Pi_{\mathbf{X}}^\top \mathbf{\Lambda}_{\epsilon, r} \mathbf{\Lambda}_{\epsilon, r}^\top \Pi_{\mathbf{X}} \mathbf{Q}_{\mathbf{P}}^{(r)}) \quad (87)$$

by a standard application of the min-max theorem for eigenvalues of Hermitian matrices.

Now,

$$\lambda_{d_r}(\mathbf{Q}_{\mathbf{P}}^{(r)\top} \Pi_{\mathbf{X}}^\top \mathbf{\Lambda}_{\epsilon, r} \mathbf{\Lambda}_{\epsilon, r}^\top \Pi_{\mathbf{X}} \mathbf{Q}_{\mathbf{P}}^{(r)}) = \lambda_{d_r}(\mathbf{\Lambda}_{\epsilon, r}^\top \Pi_{\mathbf{X}} \Pi_{\mathbf{X}}^\top \mathbf{\Lambda}_{\epsilon, r}) = \lambda_{d_r}(\mathbf{\Lambda}_{\epsilon, r}^\top \mathbf{X}^\top \mathbf{X} \mathbf{\Lambda}_{\epsilon, r}), \quad (88)$$

and a similar line of reasoning to that in Proposition 7 shows that $\sigma_{d_r}(\tilde{\mathbf{R}}_r \tilde{\mathbf{R}}_r^\top) = \Omega(\epsilon_r^2)$ almost surely. The result then follows from submultiplicativity and unitary invariance of the spectral norm. \square

Proposition 18. *If \mathbf{X} is of rank d and each $\mathbf{Y}^{(r)}$ is of rank d_r then*

$$(\tilde{\mathbf{R}}_r \Sigma_{\mathbf{P}}^{-1} \tilde{\mathbf{L}}^{-1})^\top = (\mathbf{\Lambda}_{\epsilon} \mathbf{Y}^\top \mathbf{Y} \mathbf{\Lambda}_{\epsilon}^\top)^{-1} \mathbf{\Lambda}_{\epsilon, r}. \quad (89)$$

Moreover, if $d_r = d$ and the matrix $\mathbf{\Lambda}_r$ is invertible, then

$$(\tilde{\mathbf{L}} \Sigma_{\mathbf{P}}^{-1} \tilde{\mathbf{R}}_r^{-1})^\top = \mathbf{\Lambda}_{\epsilon, r}^{-1} (\mathbf{X}^\top \mathbf{X})^{-1}. \quad (90)$$

Proof. Recall from Proposition 16 that $\tilde{\mathbf{L}} \tilde{\mathbf{R}}_r^\top = \mathbf{\Lambda}_{\epsilon, r}$. Similarly, since

$$\mathbf{X} \tilde{\mathbf{L}} \Sigma_{\mathbf{P}} \tilde{\mathbf{L}}^\top \mathbf{X}^\top = \mathbf{X}_{\mathbf{P}} \Sigma_{\mathbf{P}} \mathbf{X}_{\mathbf{P}}^\top = \mathbf{P} \mathbf{P}^\top = \mathbf{X} \mathbf{\Lambda}_{\epsilon} \mathbf{Y}^\top \mathbf{Y} \mathbf{\Lambda}_{\epsilon}^\top \mathbf{X}^\top, \quad (91)$$

we find that $\tilde{\mathbf{L}} \Sigma_{\mathbf{P}} \tilde{\mathbf{L}}^\top = \mathbf{\Lambda}_{\epsilon} \mathbf{Y}^\top \mathbf{Y} \mathbf{\Lambda}_{\epsilon}^\top$. Thus

$$(\tilde{\mathbf{R}}_r \Sigma_{\mathbf{P}}^{-1} \tilde{\mathbf{L}}^{-1})^\top = \tilde{\mathbf{L}}^{-\top} \Sigma_{\mathbf{P}}^{-1} \tilde{\mathbf{L}}^{-1} \tilde{\mathbf{L}} \tilde{\mathbf{R}}_r^\top \quad (92)$$

$$= (\mathbf{\Lambda}_{\epsilon} \mathbf{Y}^\top \mathbf{Y} \mathbf{\Lambda}_{\epsilon}^\top)^{-1} \mathbf{\Lambda}_{\epsilon, r} \quad (93)$$

and

$$(\tilde{\mathbf{L}} \Sigma_{\mathbf{P}}^{-1} \tilde{\mathbf{R}}_r^{-1})^\top = \tilde{\mathbf{R}}_r^{-\top} \Sigma_{\mathbf{P}}^{-1} \mathbf{X}_{\mathbf{P}}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \quad (94)$$

$$= \tilde{\mathbf{R}}_r^{-\top} \Sigma_{\mathbf{P}}^{-1} \mathbf{X}_{\mathbf{P}}^\top \mathbf{X}_{\mathbf{P}} \tilde{\mathbf{L}}^{-1} (\mathbf{X}^\top \mathbf{X})^{-1} \quad (95)$$

$$= \tilde{\mathbf{R}}_r^{-\top} \Sigma_{\mathbf{P}}^{-1} \Sigma_{\mathbf{P}} \tilde{\mathbf{L}}^{-1} (\mathbf{X}^\top \mathbf{X})^{-1} \quad (96)$$

$$= \tilde{\mathbf{R}}_r^{-\top} \tilde{\mathbf{L}}^{-1} (\mathbf{X}^\top \mathbf{X})^{-1} \quad (97)$$

$$= \mathbf{\Lambda}_{\epsilon, r}^{-1} (\mathbf{X}^\top \mathbf{X})^{-1}, \quad (98)$$

where we have used the identities $\tilde{\mathbf{L}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}_{\mathbf{P}}$ and $\mathbf{X}_{\mathbf{P}}^\top \mathbf{X}_{\mathbf{P}} = \Sigma_{\mathbf{P}}$. \square

The final result before proving our main theorems is an adaptation of Lemma 4 in [5] to the MRDPG, and utilises the previous results to establish upper bounds for the two-to-infinity norms of a number of residual terms that will appear in the proofs of Theorems 2 and 3:

Proposition 19. *Let*

$$\mathbf{R}_{1,1} = \mathbf{U}_P(\mathbf{U}_P^\top \mathbf{U}_A \Sigma_A^{1/2} - \Sigma_P^{1/2} \mathbf{W}) \quad (99)$$

$$\mathbf{R}_{1,2} = (\mathbf{I} - \mathbf{U}_P \mathbf{U}_P^\top)(\mathbf{A} - \mathbf{P})(\mathbf{V}_A - \mathbf{V}_P \mathbf{W}) \Sigma_A^{-1/2} \quad (100)$$

$$\mathbf{R}_{1,3} = -\mathbf{U}_P \mathbf{U}_P^\top (\mathbf{A} - \mathbf{P}) \mathbf{V}_P \mathbf{W} \Sigma_A^{-1/2} \quad (101)$$

$$\mathbf{R}_{1,4} = (\mathbf{A} - \mathbf{P}) \mathbf{V}_P (\mathbf{W} \Sigma_A^{-1/2} - \Sigma_P^{-1/2} \mathbf{W}) \quad (102)$$

and

$$\mathbf{R}_{2,1} = \mathbf{V}_P(\mathbf{V}_P^\top \mathbf{V}_A \Sigma_A^{1/2} - \Sigma_P^{1/2} \mathbf{W}) \quad (103)$$

$$\mathbf{R}_{2,2} = (\mathbf{I} - \mathbf{V}_P \mathbf{V}_P^\top)(\mathbf{A} - \mathbf{P})^\top (\mathbf{U}_A - \mathbf{U}_P \mathbf{W}) \Sigma_A^{-1/2} \quad (104)$$

$$\mathbf{R}_{2,3} = -\mathbf{V}_P \mathbf{V}_P^\top (\mathbf{A} - \mathbf{P})^\top \mathbf{U}_P \mathbf{W} \Sigma_A^{-1/2} \quad (105)$$

$$\mathbf{R}_{2,4} = (\mathbf{A} - \mathbf{P})^\top \mathbf{U}_P (\mathbf{W} \Sigma_A^{-1/2} - \Sigma_P^{-1/2} \mathbf{W}) \quad (106)$$

Then the following bounds hold almost surely:

- (i) $\|\mathbf{R}_{1,1}\|_{2 \rightarrow \infty} = O\left(\frac{\epsilon^5 k^{1/2} \log(n)}{\rho^{1/2} n}\right)$ and $\|\mathbf{R}_{2,1}\|_{2 \rightarrow \infty} = O\left(\frac{\epsilon^4 k^{1/2} \log(n)}{\rho^{1/2} n}\right)$;
- (ii) $\|\mathbf{R}_{1,2}\|_{2 \rightarrow \infty}, \|\mathbf{R}_{2,2}\|_{2 \rightarrow \infty} = O\left(\frac{\epsilon^2 k^{1/2} \log(n)}{\rho^{1/2} n^{3/4}}\right)$;
- (iii) $\|\mathbf{R}_{1,3}\|_{2 \rightarrow \infty} = O\left(\frac{\epsilon \log^{1/2}(n)}{\rho^{1/2} n}\right)$ and $\|\mathbf{R}_{2,3}\|_{2 \rightarrow \infty} = O\left(\frac{\log^{1/2}(n)}{\rho^{1/2} n}\right)$;
- (iv) $\|\mathbf{R}_{1,4}\|_{2 \rightarrow \infty}, \|\mathbf{R}_{2,4}\|_{2 \rightarrow \infty} = O\left(\frac{\epsilon^{9/2} k^{3/4} \log^{3/2}(n)}{n}\right)$.

Proof. We give full proofs of the bounds only for the terms $\mathbf{R}_{1,i}$, noting any differences for the proofs for the terms $\mathbf{R}_{2,i}$.

- (i) Recall that $\mathbf{U}_P \Sigma_P^{1/2} = \mathbf{X} \tilde{\mathbf{L}}$, where the matrix $\tilde{\mathbf{L}} \in \text{GL}(d)$ satisfies $\|\tilde{\mathbf{L}}\| = O(\epsilon)$ by Corollary 17. Using the relation $\|\mathbf{A}\mathbf{B}\|_{2 \rightarrow \infty} \leq \|\mathbf{A}\|_{2 \rightarrow \infty} \|\mathbf{B}\|$ (see, for example, [6], Proposition 6.5) we find that $\|\mathbf{U}_P\|_{2 \rightarrow \infty} \leq \|\mathbf{X}\|_{2 \rightarrow \infty} \|\tilde{\mathbf{L}}\| \|\Sigma_P^{-1/2}\|$, and thus $\|\mathbf{U}_P\|_{2 \rightarrow \infty} = O\left(\frac{\epsilon}{n^{1/2}}\right)$ as the rows of \mathbf{X} are by definition of order $O(\rho^{1/2})$ (similarly, we find that $\|\mathbf{V}_P\|_{2 \rightarrow \infty} = O\left(\frac{1}{n^{1/2}}\right)$ by splitting $\mathbf{V}_P \Sigma_P^{1/2}$ into the separate terms $\mathbf{V}_P^{(r)} \Sigma_P^{1/2}$ and evaluating each separately, and noting that $\epsilon_r \leq 1$ for all r).

Thus

$$\|\mathbf{R}_{1,1}\|_{2 \rightarrow \infty} \leq \|\mathbf{U}_P\|_{2 \rightarrow \infty} \|\mathbf{U}_P^\top \mathbf{U}_A \Sigma_A^{1/2} - \Sigma_P^{1/2} \mathbf{W}\| \quad (107)$$

$$\leq \|\mathbf{U}_P\|_{2 \rightarrow \infty} (\|(\mathbf{U}_P^\top \mathbf{U}_A - \mathbf{W}) \Sigma_A^{1/2}\|_F + \|\mathbf{W} \Sigma_A^{1/2} - \Sigma_P^{1/2} \mathbf{W}\|_F) \quad (108)$$

The first summand is $O\left(\frac{\epsilon^{7/2} k^{1/2} \log(n)}{\rho^{1/2} n^{1/2}}\right)$ by Proposition 14 and Corollary 10, while Proposition 15 shows that the second is $O\left(\frac{\epsilon^4 k^{1/2} \log(n)}{\rho^{1/2} n^{1/2}}\right)$, and so

$$\|\mathbf{R}_{1,1}\|_{2 \rightarrow \infty} = O\left(\frac{\epsilon^5 k^{1/2} \log(n)}{\rho^{1/2} n}\right). \quad (109)$$

- (ii) We begin by splitting the term $\mathbf{R}_{1,2} = \mathbf{M}_1 + \mathbf{M}_2$, where

$$\mathbf{M}_1 = \mathbf{U}_P \mathbf{U}_P^\top (\mathbf{A} - \mathbf{P})(\mathbf{V}_A - \mathbf{V}_P \mathbf{W}) \Sigma_A^{-1/2} \quad (110)$$

$$\mathbf{M}_2 = (\mathbf{A} - \mathbf{P})(\mathbf{V}_A - \mathbf{V}_P \mathbf{W}) \Sigma_A^{-1/2} \quad (111)$$

Now,

$$\|\mathbf{M}_1\|_{2 \rightarrow \infty} \leq \|\mathbf{U}_P\|_{2 \rightarrow \infty} \|\mathbf{A} - \mathbf{P}\| \|\mathbf{V}_A - \mathbf{V}_P \mathbf{W}\| \|\Sigma_A^{-1/2}\|, \quad (112)$$

where the term

$$\|\mathbf{U}_P\|_{2 \rightarrow \infty} \|\mathbf{A} - \mathbf{P}\| \|\Sigma_{\mathbf{A}}^{-1/2}\| = O\left(\frac{\epsilon^{3/2} k^{1/4} \log^{1/2}(n)}{n^{1/2}}\right) \quad (113)$$

by Proposition 8 and Corollary 10, while

$$\|\mathbf{V}_{\mathbf{A}} - \mathbf{V}_P \mathbf{W}\| \leq \|\mathbf{V}_{\mathbf{A}} - \mathbf{V}_P \mathbf{V}_P^\top \mathbf{V}_{\mathbf{A}}\| + \|\mathbf{V}_P (\mathbf{V}_P^\top \mathbf{V}_{\mathbf{A}} - \mathbf{W})\| \quad (114)$$

$$= O\left(\frac{\epsilon^{3/2} k^{1/4} \log^{1/2}(n)}{\rho^{1/2} n^{1/2}}\right) \quad (115)$$

by Propositions 13 and 14 and the asymptotic growth conditions imposed on ρ . Thus

$$\|\mathbf{M}_1\|_{2 \rightarrow \infty} = O\left(\frac{\epsilon^3 k^{1/2} \log(n)}{\rho^{1/2} n}\right). \quad (116)$$

Next, note that

$$\mathbf{M}_2 = (\mathbf{A} - \mathbf{P})(\mathbf{I} - \mathbf{V}_P \mathbf{V}_P^\top) \mathbf{V}_{\mathbf{A}} \Sigma_{\mathbf{A}}^{-1/2} + (\mathbf{A} - \mathbf{P}) \mathbf{V}_P (\mathbf{V}_P^\top \mathbf{V}_{\mathbf{A}} - \mathbf{W}) \Sigma_{\mathbf{A}}^{-1/2} \quad (117)$$

where

$$\|(\mathbf{A} - \mathbf{P}) \mathbf{V}_P (\mathbf{V}_P^\top \mathbf{V}_{\mathbf{A}} - \mathbf{W}) \Sigma_{\mathbf{A}}^{-1/2}\|_{2 \rightarrow \infty} \leq \|\mathbf{A} - \mathbf{P}\| \|\mathbf{V}_P^\top \mathbf{V}_{\mathbf{A}} - \mathbf{W}\| \|\Sigma_{\mathbf{A}}^{-1/2}\| \quad (118)$$

$$= O\left(\frac{\epsilon^{7/2} k^{3/4} \log^{3/2}(n)}{\rho^{3/2} n}\right) \quad (119)$$

by Propositions 8, 14 and Corollary 10.

To bound the remaining term, let $\mathbf{M} = (\mathbf{A} - \mathbf{P})(\mathbf{I} - \mathbf{V}_P \mathbf{V}_P^\top) \mathbf{V}_{\mathbf{A}} \mathbf{V}_{\mathbf{A}}^\top$, so that

$$(\mathbf{A} - \mathbf{P})(\mathbf{I} - \mathbf{V}_P \mathbf{V}_P^\top) \mathbf{V}_{\mathbf{A}} \Sigma_{\mathbf{A}}^{-1/2} = \mathbf{M} \mathbf{V}_{\mathbf{A}} \Sigma_{\mathbf{A}}^{-1/2} \quad (120)$$

and so

$$\|(\mathbf{A} - \mathbf{P})(\mathbf{I} - \mathbf{V}_P \mathbf{V}_P^\top) \mathbf{V}_{\mathbf{A}} \Sigma_{\mathbf{A}}^{-1/2}\|_{2 \rightarrow \infty} \leq \|\mathbf{M}\|_{2 \rightarrow \infty} \|\mathbf{V}_{\mathbf{A}} \Sigma_{\mathbf{A}}^{-1/2}\|. \quad (121)$$

The term $\|\mathbf{V}_{\mathbf{A}} \Sigma_{\mathbf{A}}^{-1/2}\|$ is $O\left(\frac{1}{\rho^{1/2} n^{1/2}}\right)$ by Corollary 10, so it suffices to bound $\|\mathbf{M}\|_{2 \rightarrow \infty}$. To do this, we claim that the Frobenius norms of the rows of the matrix \mathbf{M} are exchangeable, and thus have the same expectation, which implies that $\mathbb{E}(\|\mathbf{M}\|_F^2) = n \mathbb{E}(\|\mathbf{M}_i\|^2)$ for any $i \in \{1, \dots, n\}$. Applying Markov's inequality, we therefore see that

$$\mathbb{P}(\|\mathbf{M}_i\| > t) \leq \frac{\mathbb{E}(\|\mathbf{M}_i\|^2)}{t^2} = \frac{\mathbb{E}(\|\mathbf{M}\|_F^2)}{nt^2}. \quad (122)$$

Now,

$$\|\mathbf{M}\|_F \leq \|\mathbf{A} - \mathbf{P}\| \|\mathbf{V}_{\mathbf{A}} - \mathbf{V}_P \mathbf{V}_P^\top \mathbf{V}_{\mathbf{A}}\|_F \|\mathbf{V}_{\mathbf{A}}^\top\|_F \quad (123)$$

$$= O(\epsilon^2 k^{1/2} \log(n)) \quad (124)$$

by Propositions 8 and 13. It follows that

$$\mathbb{P}\left(\|\mathbf{M}_i\| > \frac{\epsilon^2 k^{1/2} \log(n)}{n^{1/4}}\right) = O\left(\frac{1}{n^{1/2}}\right) \quad (125)$$

and thus

$$\|\mathbf{M}\|_{2 \rightarrow \infty} = O\left(\frac{\epsilon^2 k^{1/2} \log(n)}{n^{1/4}}\right) \quad (126)$$

and

$$\|(\mathbf{A} - \mathbf{P})(\mathbf{I} - \mathbf{V}_P \mathbf{V}_P^\top) \mathbf{V}_{\mathbf{A}} \Sigma_{\mathbf{A}}^{-1/2}\|_{2 \rightarrow \infty} = O\left(\frac{\epsilon^2 k^{1/2} \log(n)}{\rho^{1/2} n^{3/4}}\right) \quad (127)$$

almost surely.

We must therefore show that the Frobenius norms of the rows of \mathbf{M} are exchangeable. Let $\mathbf{Q}_L \in \mathbb{O}(n)$ and, for each $r \in \{1, \dots, k\}$, $\mathbf{Q}_{R,r} \in \mathbb{O}(n_r)$ be permutation matrices (where we require that $\mathbf{Q}_{R,r} = \mathbf{Q}_L$ if $\mathbf{Y}^{(r)} = \mathbf{X}\mathbf{G}_r$ with probability one for some matrix $\mathbf{G}_r \in \mathbb{R}^{d \times d_r}$, and similarly that $\mathbf{Q}_{R,r} = \mathbf{Q}_{R,s}$ if $\mathbf{Y}^{(s)} = \mathbf{Y}^{(r)}\mathbf{G}_{r,s}$ with probability one for some matrix $\mathbf{G}_{r,s} \in \mathbb{R}^{d_r \times d_s}$), and let $\mathbf{Q}_R = \mathbf{Q}_{R,1} \oplus \dots \oplus \mathbf{Q}_{R,k}$. For any matrix \mathbf{G} , let $\mathcal{R}_d(\mathbf{G})$ denote the projection onto the subspace spanned by the right singular vectors corresponding to the leading d singular values of \mathbf{G} , and let $\mathcal{R}_d^\perp(\mathbf{G})$ denote the projections onto the orthogonal complements of this subspace.

Note that

$$\mathcal{R}_d(\mathbf{P}) = \mathbf{V}_P \mathbf{V}_P^\top \quad \text{and} \quad \mathcal{R}_d(\mathbf{A}) = \mathbf{V}_A \mathbf{V}_A^\top, \quad (128)$$

while for any permutation matrices $\mathbf{Q}_L \in \mathbb{O}(n)$ and $\mathbf{Q}_{R,r} \in \mathbb{O}(n_r)$ we have

$$\mathcal{R}_d(\mathbf{Q}_L \mathbf{P} \mathbf{Q}_R^\top) = \mathbf{Q}_R \mathbf{V}_P \mathbf{V}_P^\top \mathbf{Q}_R^\top \quad \text{and} \quad \mathcal{R}_d(\mathbf{Q}_L \mathbf{A} \mathbf{Q}_R^\top) = \mathbf{Q}_R \mathbf{V}_A \mathbf{V}_A^\top \mathbf{Q}_R^\top. \quad (129)$$

For any commensurate pair of matrices \mathbf{G} and \mathbf{H} , define an operator

$$\mathcal{P}_{\mathcal{R},d}(\mathbf{G}, \mathbf{H}) = (\mathbf{G} - \mathbf{H}) \mathcal{R}_d^\perp(\mathbf{H}) \mathcal{R}_d(\mathbf{G}) \quad (130)$$

and note that $\mathcal{P}_{\mathcal{R},d}(\mathbf{A}, \mathbf{P}) = \mathbf{M}$, while

$$\mathcal{P}_{\mathcal{R},d}(\mathbf{Q}_L \mathbf{A} \mathbf{Q}_R^\top, \mathbf{Q}_L \mathbf{P} \mathbf{Q}_R^\top) = \mathbf{Q}_L (\mathbf{A} - \mathbf{P}) \mathbf{Q}_R^\top \mathbf{Q}_R (\mathbf{I} - \mathbf{V}_P \mathbf{V}_P^\top) \mathbf{Q}_R^\top \mathbf{Q}_R \mathbf{V}_A \mathbf{V}_A^\top \mathbf{Q}_R^\top \quad (131)$$

$$= \mathbf{Q}_L \mathbf{M} \mathbf{Q}_R^\top. \quad (132)$$

Our assumptions regarding the distribution of the latent positions ensure that the entries of the pair (\mathbf{A}, \mathbf{P}) have the same joint distribution as those of the pair $(\mathbf{Q}_L \mathbf{A} \mathbf{Q}_R^\top, \mathbf{Q}_L \mathbf{P} \mathbf{Q}_R^\top)$, since \mathbf{Q}_R^\top permutes the columns of each $\mathbf{A}^{(r)}$ and $\mathbf{P}^{(r)}$ separately. Therefore, the entries of the matrix $\mathcal{P}_{\mathcal{R},d}(\mathbf{A}, \mathbf{P})$ have the same joint distribution as those of the matrix $\mathcal{P}_{\mathcal{R},d}(\mathbf{Q}_L \mathbf{A} \mathbf{Q}_R^\top, \mathbf{Q}_L \mathbf{P} \mathbf{Q}_R^\top)$, which implies that \mathbf{M} has the same distribution as $\mathbf{Q}_L \mathbf{M} \mathbf{Q}_R^\top$, and consequently the Frobenius norms of the rows of \mathbf{M} have the same distribution as those of $\mathbf{Q}_L \mathbf{M}$, which proves our claim.

Combining these results, we see that

$$\|\mathbf{R}_{1,2}\|_{2 \rightarrow \infty} = O\left(\frac{\epsilon^2 k^{1/2} \log(n)}{\rho^{1/2} n^{3/4}}\right) \quad (133)$$

almost surely, as required.

The proof of the bound for the term $\mathbf{R}_{2,2}$ follows similarly, and culminates in showing that the term

$$\mathbf{N} = (\mathbf{A} - \mathbf{P})^\top (\mathbf{I} - \mathbf{U}_P \mathbf{U}_P^\top) \mathbf{U}_A \mathbf{U}_A^\top \quad (134)$$

satisfies

$$\|\mathbf{N}\|_{2 \rightarrow \infty} = O\left(\frac{\epsilon^2 k^{1/2} \log(n)}{n^{1/4}}\right) \quad (135)$$

almost surely. The matrix $\mathbf{N} \in \mathbb{R}^{(n_1 + \dots + n_k) \times n}$ splits into k distinct matrices $\mathbf{N}^{(r)} \in \mathbb{R}^{n_r \times n}$. This time, we must show that the Frobenius norms of the rows of each $\mathbf{N}^{(r)}$ are interchangeable, from which a similar argument to our previous one will allow us to derive our desired bound for $\|\mathbf{N}\|_{2 \rightarrow \infty}$.

Note that for any $r \in \{1, \dots, k\}$

$$\mathcal{R}_d(\mathbf{P}^{(r)\top}) = \mathbf{U}_P \mathbf{U}_P^\top \quad \text{and} \quad \mathcal{R}_d(\mathbf{A}^{(r)\top}) = \mathbf{U}_A \mathbf{U}_A^\top, \quad (136)$$

while for any permutation matrices $\mathbf{Q}_L \in \mathbb{O}(n_r)$ and $\mathbf{Q}_R \in \mathbb{O}(n)$ we have

$$\mathcal{R}_d(\mathbf{Q}_L \mathbf{P}^{(r)\top} \mathbf{Q}_R^\top) = \mathbf{Q}_R \mathbf{U}_P \mathbf{U}_P^\top \mathbf{Q}_R^\top \quad \text{and} \quad \mathcal{R}_d(\mathbf{Q}_L \mathbf{A}^{(r)\top} \mathbf{Q}_R^\top) = \mathbf{Q}_R \mathbf{V}_A \mathbf{V}_A^\top \mathbf{Q}_R^\top. \quad (137)$$

Also, $\mathcal{P}_{\mathcal{R},d}(\mathbf{A}^{(r)\top}, \mathbf{P}^{(r)\top}) = \mathbf{N}^{(r)}$, while we find that

$$\mathcal{P}_{\mathcal{R},d}(\mathbf{Q}_L \mathbf{A}^{(r)\top} \mathbf{Q}_R^\top, \mathbf{Q}_L \mathbf{P}^{(r)\top} \mathbf{Q}_R^\top) = \mathbf{Q}_L \mathbf{N}^{(r)} \mathbf{Q}_R^\top, \quad (138)$$

and so it follows from a similar argument to before that the rows of $\mathbf{N}^{(r)}$ are interchangeable, and thus we obtain our desired bound.

(iii) Similarly to part (i), we see that

$$\|\mathbf{R}_{1,3}\|_{2 \rightarrow \infty} \leq \|\mathbf{U}_P\|_{2 \rightarrow \infty} \|\mathbf{U}_P^\top (\mathbf{A} - \mathbf{P}) \mathbf{V}_P \mathbf{W} \Sigma_A^{-1/2}\| \quad (139)$$

$$\leq \|\mathbf{U}_P\|_{2 \rightarrow \infty} \|\mathbf{U}_P^\top (\mathbf{A} - \mathbf{P}) \mathbf{V}_P\|_F \|\mathbf{W} \Sigma_A^{-1/2}\|_F \quad (140)$$

$$= O\left(\frac{\epsilon \log^{1/2}(n)}{\rho^{1/2} n}\right) \quad (141)$$

by Proposition 11 and Corollary 10.

(iv) Observe that

$$\|\mathbf{R}_{1,4}\|_{2 \rightarrow \infty} \leq \|\mathbf{R}_{1,4}\|_F \quad (142)$$

$$\leq \|\mathbf{A} - \mathbf{P}\| \|\mathbf{V}_P\|_F \|\mathbf{W} \Sigma_A^{-1/2} - \Sigma_P^{-1/2} \mathbf{W}\|_F \quad (143)$$

$$= O\left(\frac{\epsilon^{9/2} k^{3/4} \log^{3/2}(n)}{n}\right) \quad (144)$$

by Propositions 8 and 15.

□

Proof of Theorem 2

Proof. We first consider the left embedding \mathbf{X}_A . Observe that

$$\mathbf{X}_A - \mathbf{X}_P \mathbf{W} = \mathbf{U}_A \Sigma_A^{1/2} - \mathbf{U}_P \Sigma_P^{1/2} \mathbf{W} \quad (145)$$

$$= \mathbf{U}_A \Sigma_A^{1/2} - \mathbf{U}_P \mathbf{U}_P^\top \mathbf{U}_A \Sigma_A^{1/2} + \mathbf{U}_P (\mathbf{U}_P^\top \mathbf{U}_A \Sigma_A^{1/2} - \Sigma_P^{1/2} \mathbf{W}) \quad (146)$$

$$= \mathbf{U}_A \Sigma_A^{1/2} - \mathbf{U}_P \mathbf{U}_P^\top \mathbf{U}_A \Sigma_A^{1/2} + \mathbf{R}_1. \quad (147)$$

Noting that $\mathbf{U}_A \Sigma_A^{1/2} = \mathbf{A} \mathbf{V}_A \Sigma_A^{-1/2}$ and $\mathbf{U}_P \mathbf{U}_P^\top \mathbf{P} = \mathbf{P}$, we see that

$$\mathbf{X}_A - \mathbf{X}_P \mathbf{W} = \mathbf{A} \mathbf{V}_A \Sigma_A^{-1/2} - \mathbf{U}_P \mathbf{U}_P^\top \mathbf{A} \mathbf{V}_A \Sigma_A^{-1/2} + \mathbf{R}_{1,1} \quad (148)$$

$$= (\mathbf{A} - \mathbf{P}) \mathbf{V}_A \Sigma_A^{-1/2} - (\mathbf{U}_P \mathbf{U}_P^\top \mathbf{A} - \mathbf{P}) \mathbf{V}_A \Sigma_A^{-1/2} + \mathbf{R}_{1,1} \quad (149)$$

$$= (\mathbf{A} - \mathbf{P}) \mathbf{V}_A \Sigma_A^{-1/2} - \mathbf{U}_P \mathbf{U}_P^\top (\mathbf{A} - \mathbf{P}) \mathbf{V}_A \Sigma_A^{-1/2} + \mathbf{R}_{1,1} \quad (150)$$

$$= (\mathbf{I} - \mathbf{U}_P \mathbf{U}_P^\top) (\mathbf{A} - \mathbf{P}) \mathbf{V}_A \Sigma_A^{-1/2} + \mathbf{R}_{1,1} \quad (151)$$

$$= (\mathbf{I} - \mathbf{U}_P \mathbf{U}_P^\top) (\mathbf{A} - \mathbf{P}) (\mathbf{V}_P \mathbf{W} + (\mathbf{V}_A - \mathbf{V}_P \mathbf{W})) \Sigma_A^{-1/2} + \mathbf{R}_{1,1} \quad (152)$$

$$= (\mathbf{A} - \mathbf{P}) \mathbf{V}_P \mathbf{W} \Sigma_A^{-1/2} + \mathbf{R}_{1,3} + \mathbf{R}_{1,2} + \mathbf{R}_{1,1} \quad (153)$$

$$= (\mathbf{A} - \mathbf{P}) \mathbf{V}_P \Sigma_P^{-1/2} \mathbf{W} + \mathbf{R}_{1,4} + \mathbf{R}_{1,3} + \mathbf{R}_{1,2} + \mathbf{R}_{1,1}. \quad (154)$$

Applying Proposition 19, we find that

$$\|\mathbf{X}_A - \mathbf{X}_P \mathbf{W}\|_{2 \rightarrow \infty} = \|(\mathbf{A} - \mathbf{P}) \mathbf{V}_P \Sigma_P^{-1/2}\|_{2 \rightarrow \infty} + O\left(\frac{\epsilon^2 k^{1/2} \log(n)}{\rho^{1/2} n^{3/4}}\right) \quad (155)$$

$$\leq \sigma_d(\mathbf{P})^{-1/2} \|(\mathbf{A} - \mathbf{P}) \mathbf{V}_P\|_{2 \rightarrow \infty} + O\left(\frac{\epsilon^2 k^{1/2} \log(n)}{\rho^{1/2} n^{3/4}}\right). \quad (156)$$

almost surely.

We condition on some set of latent positions. For any $i \in \{1, \dots, n\}$, $j \in \{1, \dots, d\}$ and $r \in \{1, \dots, k\}$, let

$$\mathbf{E}_{ij}^{(r)} = \begin{cases} \sum_{l=1}^{n_r} (\mathbf{A}_{il}^{(r)} - \mathbf{P}_{il}^{(r)}) v_l^{(r)} & \text{if } \mathbf{A}^{(r)} \text{ is bipartite} \\ \sum_{l \neq i} (\mathbf{A}_{il}^{(r)} - \mathbf{P}_{il}^{(r)}) v_l^{(r)} & \text{otherwise} \end{cases} \quad (157)$$

and

$$\mathbf{F}_{ij}^{(r)} = \begin{cases} \mathbf{0} & \text{if } \mathbf{A}^{(r)} \text{ is bipartite} \\ \mathbf{P}_{ii}^{(r)} v_i^{(r)} & \text{otherwise,} \end{cases} \quad (158)$$

where $v^{(r)}$ denotes the j th column of $\mathbf{V}_{\mathbf{P}}^{(r)}$, so that

$$((\mathbf{A} - \mathbf{P})\mathbf{V}_{\mathbf{P}})_{ij} = \sum_{r=1}^k \mathbf{E}_{ij}^{(r)} - \sum_{r=1}^k \mathbf{F}_{ij}^{(r)}. \quad (159)$$

The latter term is of order $O(\epsilon \rho k^{1/2})$ (as can be seen by applying the Cauchy–Schwarz inequality) and thus can be discounted from our asymptotic analysis. The former is a sum of independent, zero-mean random variables, with each individual term bounded in absolute value by $|v_l^{(r)}|$, and thus we can apply Hoeffding’s inequality to see that

$$\mathbb{P}\left(\left|\sum_{r=1}^k \mathbf{E}_{ij}^{(r)}\right| > t\right) \leq 2 \exp\left(\frac{-2t^2}{4 \sum_{r=1}^k \sum_{l=1}^{n_r} |v_l^{(r)}|^2}\right) = 2 \exp\left(\frac{-t^2}{2}\right). \quad (160)$$

Thus $((\mathbf{A} - \mathbf{P})\mathbf{V}_{\mathbf{P}})_{ij} = O(\log^{1/2}(n))$ almost surely, and hence $|((\mathbf{A} - \mathbf{P})\mathbf{V}_{\mathbf{P}})_i| = O(\log^{1/2}(n))$ almost surely by summing over all $j \in \{1, \dots, d\}$. Taking the union bound over all n rows then shows that

$$\sigma_d(\mathbf{P})^{-1/2} \|(\mathbf{A} - \mathbf{P})\mathbf{V}_{\mathbf{P}}\|_{2 \rightarrow \infty} = O\left(\frac{\log^{1/2}(n)}{\rho^{1/2} n^{1/2}}\right), \quad (161)$$

almost surely and consequently that

$$\|\mathbf{X}_{\mathbf{A}} - \mathbf{X}_{\mathbf{L}}\|_{2 \rightarrow \infty} = O\left(\frac{\log^{1/2}(n)}{\rho^{1/2} n^{1/2}}\right) \quad (162)$$

almost surely by setting $\mathbf{L} = \tilde{\mathbf{L}}\mathbf{W}$. The second bound follows from Corollary 17 and the fact that $\|\mathbf{A}\mathbf{B}\|_{2 \rightarrow \infty} \leq \|\mathbf{A}\|_{2 \rightarrow \infty} \|\mathbf{B}\|$. Integrating over all possible sets of latent positions gives the result.

A similar argument is used for the right embedding. Observe that

$$\mathbf{Y}_{\mathbf{A}} - \mathbf{Y}_{\mathbf{P}}\mathbf{W} = \mathbf{V}_{\mathbf{A}}\Sigma_{\mathbf{A}}^{1/2} - \mathbf{V}_{\mathbf{P}}\Sigma_{\mathbf{P}}^{1/2}\mathbf{W} \quad (163)$$

$$= \mathbf{V}_{\mathbf{A}}\Sigma_{\mathbf{A}}^{1/2} - \mathbf{V}_{\mathbf{P}}\mathbf{V}_{\mathbf{P}}^{\top}\mathbf{V}_{\mathbf{A}}\Sigma_{\mathbf{A}}^{1/2} + \mathbf{V}_{\mathbf{P}}(\mathbf{V}_{\mathbf{P}}^{\top}\mathbf{V}_{\mathbf{A}}\Sigma_{\mathbf{A}}^{1/2} - \Sigma_{\mathbf{P}}^{1/2}\mathbf{W}) \quad (164)$$

$$= \mathbf{V}_{\mathbf{A}}\Sigma_{\mathbf{A}}^{1/2} - \mathbf{V}_{\mathbf{P}}\mathbf{V}_{\mathbf{P}}^{\top}\mathbf{V}_{\mathbf{A}}\Sigma_{\mathbf{A}}^{1/2} + \mathbf{R}_{2,1}. \quad (165)$$

Noting that $\mathbf{A}^\top \mathbf{U}_\mathbf{A} \Sigma_\mathbf{A}^{-1/2} = \mathbf{V}_\mathbf{A} \Sigma_\mathbf{A}^{1/2}$ and $\mathbf{V}_\mathbf{P} \mathbf{V}_\mathbf{P}^\top \mathbf{P}^\top = \mathbf{P}^\top$, we see that

$$\mathbf{Y}_\mathbf{A} - \mathbf{Y}_\mathbf{P} \mathbf{W} = \mathbf{A}^\top \mathbf{U}_\mathbf{A} \Sigma_\mathbf{A}^{-1/2} - \mathbf{V}_\mathbf{P} \mathbf{V}_\mathbf{P}^\top \mathbf{A}^\top \mathbf{U}_\mathbf{A} \Sigma_\mathbf{A}^{-1/2} + \mathbf{R}_{2,1} \quad (166)$$

$$= (\mathbf{A} - \mathbf{P})^\top \mathbf{U}_\mathbf{A} \Sigma_\mathbf{A}^{-1/2} - (\mathbf{V}_\mathbf{P} \mathbf{V}_\mathbf{P}^\top \mathbf{A}^\top - \mathbf{P}^\top) \mathbf{U}_\mathbf{A} \Sigma_\mathbf{A}^{-1/2} + \mathbf{R}_{2,1} \quad (167)$$

$$= (\mathbf{A} - \mathbf{P})^\top \mathbf{U}_\mathbf{A} \Sigma_\mathbf{A}^{-1/2} - \mathbf{V}_\mathbf{P} \mathbf{V}_\mathbf{P}^\top (\mathbf{A} - \mathbf{P})^\top \mathbf{U}_\mathbf{A} \Sigma_\mathbf{A}^{-1/2} + \mathbf{R}_{2,1} \quad (168)$$

$$= (\mathbf{I} - \mathbf{V}_\mathbf{P} \mathbf{V}_\mathbf{P}^\top) (\mathbf{A} - \mathbf{P})^\top \mathbf{U}_\mathbf{A} \Sigma_\mathbf{A}^{-1/2} + \mathbf{R}_{2,1} \quad (169)$$

$$= (\mathbf{I} - \mathbf{V}_\mathbf{P} \mathbf{V}_\mathbf{P}^\top) (\mathbf{A} - \mathbf{P})^\top (\mathbf{U}_\mathbf{P} \mathbf{W} + (\mathbf{U}_\mathbf{A} - \mathbf{U}_\mathbf{P} \mathbf{W})) \Sigma_\mathbf{A}^{-1/2} + \mathbf{R}_{2,1} \quad (170)$$

$$= (\mathbf{A} - \mathbf{P})^\top \mathbf{U}_\mathbf{P} \mathbf{W} \Sigma_\mathbf{A}^{-1/2} + \mathbf{R}_{2,3} + \mathbf{R}_{2,2} + \mathbf{R}_{2,1} \quad (171)$$

$$= (\mathbf{A} - \mathbf{P})^\top \mathbf{U}_\mathbf{P} \Sigma_\mathbf{P}^{-1/2} \mathbf{W} + \mathbf{R}_{2,4} + \mathbf{R}_{2,3} + \mathbf{R}_{2,2} + \mathbf{R}_{2,1}. \quad (172)$$

Applying Proposition 19 once more, we find that

$$\|\mathbf{Y}_\mathbf{A} - \mathbf{Y}_\mathbf{P} \mathbf{W}\|_{2 \rightarrow \infty} = \|(\mathbf{A} - \mathbf{P})^\top \mathbf{U}_\mathbf{P} \Sigma_\mathbf{P}^{-1/2}\|_{2 \rightarrow \infty} + O\left(\frac{\epsilon^2 k^{1/2} \log(n)}{\rho^{1/2} n^{3/4}}\right) \quad (173)$$

almost surely and consequently that

$$\|\mathbf{Y}_\mathbf{A}^{(r)} - \mathbf{Y}_\mathbf{P}^{(r)} \mathbf{W}\|_{2 \rightarrow \infty} \leq \sigma_d(\mathbf{P})^{-1/2} \|(\mathbf{A}^{(r)} - \mathbf{P}^{(r)})^\top \mathbf{U}_\mathbf{P}\|_{2 \rightarrow \infty} + O\left(\frac{\epsilon^2 k^{1/2} \log(n)}{\rho^{1/2} n^{3/4}}\right). \quad (174)$$

almost surely.

If we condition on a set of latent positions, an identical argument to before shows that

$$\sigma_d(\mathbf{P})^{-1/2} \|(\mathbf{A}^{(r)} - \mathbf{P}^{(r)})^\top \mathbf{U}_\mathbf{P}\|_{2 \rightarrow \infty} = O\left(\frac{\log^{1/2}(n)}{\rho^{1/2} n^{1/2}}\right), \quad (175)$$

and consequently that

$$\|\mathbf{Y}_\mathbf{A}^{(r)} - \mathbf{Y} \mathbf{R}_r\|_{2 \rightarrow \infty} = O\left(\frac{\log^{1/2}(n)}{\rho^{1/2} n^{1/2}}\right) \quad (176)$$

by setting $\mathbf{R}_r = \tilde{\mathbf{R}}_r \mathbf{W}$. Proposition 16 shows that $\tilde{\mathbf{L}} \tilde{\mathbf{R}}_r^\top = \mathbf{\Lambda}_r$, and the remaining bounds follow from Corollary 17 as for the left embedding. As before, integrating over all possible sets of latent positions gives the final result. \square

Proof of Theorem 3

Proof. We first consider the left embedding $\mathbf{X}_\mathbf{A}$. Recall from the proof of Theorem 2 that

$$n^{1/2} (\mathbf{X}_\mathbf{A} \mathbf{L}^{-1} - \mathbf{X}) = n^{1/2} (\mathbf{A} - \mathbf{P}) \mathbf{V}_\mathbf{P} \Sigma_\mathbf{P}^{-1/2} \tilde{\mathbf{L}}^{-1} + n^{1/2} \mathbf{R}, \quad (177)$$

where the residual term \mathbf{R} satisfies $\|n^{1/2} \mathbf{R}\|_{2 \rightarrow \infty} \rightarrow 0$ by Proposition 19 and our assumptions in Section 2.1.

The first of the right-hand terms may be rewritten as

$$n^{1/2} (\mathbf{A} - \mathbf{P}) \mathbf{V}_\mathbf{P} \Sigma_\mathbf{P}^{-1/2} \tilde{\mathbf{L}}^{-1} = n^{1/2} \sum_{r=1}^k (\mathbf{A}^{(r)} - \mathbf{P}^{(r)}) \mathbf{Y}^{(r)} \tilde{\mathbf{R}}_r \Sigma_\mathbf{P}^{-1} \tilde{\mathbf{L}}^{-1} \quad (178)$$

by splitting $(\mathbf{A} - \mathbf{P}) \mathbf{V}_\mathbf{P}$ into the individual terms $(\mathbf{A}^{(r)} - \mathbf{P}^{(r)}) \mathbf{V}_\mathbf{P}^{(r)}$ and noting that

$$\mathbf{V}_\mathbf{P}^{(r)} \Sigma_\mathbf{P}^{-1/2} = \mathbf{Y}_\mathbf{P}^{(r)} \Sigma_\mathbf{P}^{-1} = \mathbf{Y}^{(r)} \tilde{\mathbf{R}}_r \Sigma_\mathbf{P}^{-1}. \quad (179)$$

Consequently,

$$n^{1/2} (\mathbf{X}_\mathbf{A} \mathbf{L}^{-1} - \mathbf{X})_i^\top \rightarrow n^{1/2} \sum_{r=1}^k (\tilde{\mathbf{R}}_r \Sigma_\mathbf{P}^{-1} \tilde{\mathbf{L}}^{-1})^\top [(\mathbf{A}^{(r)} - \mathbf{P}^{(r)}) \mathbf{Y}^{(r)}]_i^\top \quad (180)$$

$$= n (\mathbf{\Lambda}_\epsilon \mathbf{Y}^\top \mathbf{Y} \mathbf{\Lambda}_\epsilon^\top)^{-1} \sum_{r=1}^k \left[\frac{1}{n^{1/2}} \sum_{j=1}^{n_r} (\mathbf{A}_{ij}^{(r)} - \mathbf{P}_{ij}^{(r)}) \mathbf{\Lambda}_{\epsilon,r} \mathbf{Y}_j^{(r)} \right] \quad (181)$$

almost surely by Proposition 18.

Noting that $\mathbf{Y}_j^{(r)} = \rho^{1/2} v_j^{(r)}$ and that $n(\boldsymbol{\Lambda}_\epsilon \mathbf{Y}^\top \mathbf{Y} \boldsymbol{\Lambda}_\epsilon^\top)^{-1} \rightarrow \frac{1}{\rho} \Delta_{\boldsymbol{\Lambda}, Y}^{-1}$ almost surely by the law of large numbers, we see that

$$n^{1/2}(\mathbf{X}_A \mathbf{L}^{-1} - \mathbf{X})_i^\top \rightarrow \Delta_{\boldsymbol{\Lambda}, Y}^{-1} \sum_{r=1}^k \left[\frac{1}{\rho^{1/2} n^{1/2}} \sum_{j=1}^{n_r} (\mathbf{A}_{ij}^{(r)} - \mathbf{P}_{ij}^{(r)}) \boldsymbol{\Lambda}_{\epsilon, r} v_j^{(r)} \right] \quad (182)$$

almost surely.

If $\mathbf{A}^{(r)}$ is not bipartite, we may disregard the diagonal terms in our asymptotic analysis, as each of the k such terms satisfies

$$\left\| \frac{1}{\rho^{1/2} n^{1/2}} \mathbf{P}_{ii}^{(r)} \boldsymbol{\Lambda}_{\epsilon, r} v_i^{(r)} \right\| \rightarrow 0 \quad (183)$$

almost surely.

Assume first that each of the matrices $\mathbf{Y}^{(r)}$ is independent of the others. Conditional on $\xi_i = \mathbf{x}$, we have $\mathbf{P}_{ij}^{(r)} = \rho \mathbf{x}^\top \boldsymbol{\Lambda}_{\epsilon, r} v_j^{(r)}$, and so

$$\frac{1}{\rho^{1/2} n^{1/2}} \sum_j (\mathbf{A}_{ij}^{(r)} - \mathbf{P}_{ij}^{(r)}) \boldsymbol{\Lambda}_{\epsilon, r} v_j^{(r)} \quad (184)$$

is a scaled sum of independent, identically distributed, zero-mean random variables, each with covariance matrix given by

$$\mathbb{E}[\mathbf{x}^\top \boldsymbol{\Lambda}_{\epsilon, r} v_r (1 - \rho \mathbf{x}^\top \boldsymbol{\Lambda}_{\epsilon, r} v_r) \cdot \boldsymbol{\Lambda}_{\epsilon, r} v_r v_r^\top \boldsymbol{\Lambda}_{\epsilon, r}^\top] \quad (185)$$

which implies (by the multivariate central limit theorem) that

$$\frac{1}{\rho^{1/2} n^{1/2}} \sum_j (\mathbf{A}_{ij}^{(r)} - \mathbf{P}_{ij}^{(r)}) \boldsymbol{\Lambda}_{\epsilon, r} v_j^{(r)} \rightarrow \mathcal{N}(\mathbf{0}, c_r \boldsymbol{\Lambda}_r \boldsymbol{\Sigma}_Y^{(r)}(\mathbf{x}) \boldsymbol{\Lambda}_r^\top) \quad (186)$$

if $\epsilon_r = 1$, and vanishes otherwise, and thus we find that

$$n^{1/2}(\mathbf{X}_A \mathbf{L}^{-1} - \mathbf{X})_i^\top \rightarrow \mathcal{N}(\mathbf{0}, \Delta_{\boldsymbol{\Lambda}, Y}^{-1} \boldsymbol{\Lambda}_* \boldsymbol{\Sigma}_Y(\mathbf{x}) \boldsymbol{\Lambda}_*^\top \Delta_{\boldsymbol{\Lambda}, Y}^{-1}) \quad (187)$$

almost surely, where $\boldsymbol{\Sigma}_Y(\mathbf{x}) = c_1 \boldsymbol{\Sigma}_Y^{(1)}(\mathbf{x}) \oplus \dots \oplus c_k \boldsymbol{\Sigma}_Y^{(k)}(\mathbf{x})$. We deduce the Central Limit Theorem by integrating over all possible values of $\mathbf{x} \in \mathcal{X}$.

The same statement holds even if there is dependence between any of the matrices $\mathbf{Y}^{(r)}$. Indeed, suppose that the matrices $\mathbf{Y}^{(r)}$ are dependent for all $r \in \mathcal{K}$, for some subset \mathcal{K} of $\{1, \dots, k\}$. Then for any i and j ,

$$\frac{1}{\rho^{1/2} n^{1/2}} \sum_j \sum_{r \in \mathcal{K}} (\mathbf{A}_{ij}^{(r)} - \mathbf{P}_{ij}^{(r)}) \boldsymbol{\Lambda}_{\epsilon, r} v_j^{(r)} \quad (188)$$

is a sum of independent, identically distributed, zero-mean random variables, for which the covariance matrices

$$\mathbb{E} \left[\sum_{r, s \in \mathcal{K}} (\mathbf{A}_{ij}^{(r)} - \mathbf{P}_{ij}^{(r)}) (\mathbf{A}_{ij}^{(s)} - \mathbf{P}_{ij}^{(s)}) \cdot \boldsymbol{\Lambda}_{\epsilon, r} v_j^{(r)} v_j^{(s)\top} \boldsymbol{\Lambda}_{\epsilon, s}^\top \right] \quad (189)$$

reduce to

$$\mathbb{E} \left[\sum_{r \in \mathcal{K}} \mathbf{x}^\top \boldsymbol{\Lambda}_{\epsilon, r} v_r (1 - \rho \mathbf{x}^\top \boldsymbol{\Lambda}_{\epsilon, r} v_r) \cdot \boldsymbol{\Lambda}_{\epsilon, r} v_r v_r^\top \boldsymbol{\Lambda}_{\epsilon, r}^\top \right] \quad (190)$$

since the terms

$$\mathbb{E} \left[(\mathbf{A}_{ij}^{(r)} - \mathbf{P}_{ij}^{(r)}) (\mathbf{A}_{ij}^{(s)} - \mathbf{P}_{ij}^{(s)}) \cdot \boldsymbol{\Lambda}_{\epsilon, r} v_j^{(r)} v_j^{(s)\top} \boldsymbol{\Lambda}_{\epsilon, s}^\top \right] \quad (191)$$

vanish if r and s are not equal.

Similarly, for the right embedding we observe that

$$n^{1/2}(\mathbf{Y}_{\mathbf{A}}^{(r)} \mathbf{R}_r^{-1} - \mathbf{Y}^{(r)}) = n^{1/2}(\mathbf{A}^{(r)} - \mathbf{P}^{(r)})^\top \mathbf{U}_{\mathbf{P}} \boldsymbol{\Sigma}_{\mathbf{P}}^{-1/2} \tilde{\mathbf{R}}_r^{-1} + n^{1/2} \mathbf{R}, \quad (192)$$

where the residual term \mathbf{R} satisfies $\|n^{1/2} \mathbf{R}\|_{2 \rightarrow \infty} \rightarrow 0$.

Again, we may rewrite the first of the right-hand terms, obtaining the expression

$$n^{1/2}(\mathbf{A}^{(r)} - \mathbf{P}^{(r)})^\top \mathbf{U}_{\mathbf{P}} \boldsymbol{\Sigma}_{\mathbf{P}}^{-1/2} \tilde{\mathbf{R}}_r^{-1} = n^{1/2}(\mathbf{A}^{(r)} - \mathbf{P}^{(r)})^\top \mathbf{X} \tilde{\mathbf{L}} \boldsymbol{\Sigma}_{\mathbf{P}}^{-1} \tilde{\mathbf{R}}_r^{-1} \quad (193)$$

by noting that $\mathbf{U}_{\mathbf{P}} \boldsymbol{\Sigma}_{\mathbf{P}}^{-1/2} = \mathbf{X} \tilde{\mathbf{L}} \boldsymbol{\Sigma}_{\mathbf{P}}^{-1}$, and consequently find that

$$n^{1/2}(\mathbf{Y}_{\mathbf{A}}^{(r)} \mathbf{R}_r^{-1} - \mathbf{Y}^{(r)}) \rightarrow n^{1/2}(\tilde{\mathbf{L}} \boldsymbol{\Sigma}_{\mathbf{P}}^{-1} \tilde{\mathbf{R}}_r^{-1})^\top [(\mathbf{A}^{(r)} - \mathbf{P}^{(r)})^\top \mathbf{X}]_i^\top \quad (194)$$

$$= n \boldsymbol{\Lambda}_{\epsilon, r}^{-1} (\mathbf{X}^\top \mathbf{X})^{-1} \left[\frac{1}{n^{1/2}} \sum_{j=1}^n (\mathbf{A}_{ji}^{(r)} - \mathbf{P}_{ji}^{(r)}) \mathbf{X}_j \right] \quad (195)$$

almost surely by Proposition 18.

Noting that $\mathbf{X}_j = \rho^{1/2} \xi_j$ and that $n(\mathbf{X}^\top \mathbf{X})^{-1} \rightarrow \frac{1}{\rho} \Delta_X^{-1}$ almost surely by the law of large numbers, we see that

$$n^{1/2}(\mathbf{Y}_{\mathbf{A}}^{(r)} \mathbf{R}_r^{-1} - \mathbf{Y}^{(r)}) \rightarrow \boldsymbol{\Lambda}_r^{-1} \Delta_X^{-1} \left[\frac{1}{\rho^{1/2} n^{1/2}} \sum_{j=1}^n (\mathbf{A}_{ji}^{(r)} - \mathbf{P}_{ji}^{(r)}) \xi_j \right] \quad (196)$$

almost surely.

As before, if $\mathbf{A}^{(r)}$ is not bipartite we may disregard the diagonal term in our asymptotic analysis, as it satisfies

$$\left\| \frac{1}{\rho^{1/2} n^{1/2}} \mathbf{P}_{ii}^{(r)} \xi_i \right\| \rightarrow 0 \quad (197)$$

almost surely.

Conditional on $v_i^{(r)} = \mathbf{y}$, we have $\mathbf{P}_{ji}^{(r)} = \rho \xi_j^\top \boldsymbol{\Lambda}_{\epsilon, r} \mathbf{y}$, and so

$$\frac{1}{\rho^{1/2} n^{1/2}} \sum_j (\mathbf{A}_{ji}^{(r)} - \mathbf{P}_{ji}^{(r)}) \xi_j \quad (198)$$

is a scaled sum of independent, zero-mean random variables, each with covariance matrix given by

$$\mathbb{E}[\xi^\top \boldsymbol{\Lambda}_{\epsilon, r} \mathbf{y} (1 - \rho \xi^\top \boldsymbol{\Lambda}_{\epsilon, r} \mathbf{y}) \cdot \xi \xi^\top] \quad (199)$$

which implies (by the multivariate central limit theorem) that, provided $\epsilon_r = 1$,

$$\frac{1}{\rho^{1/2} n^{1/2}} \sum_j (\mathbf{A}_{ji}^{(r)} - \mathbf{P}_{ji}^{(r)}) \xi_j \rightarrow \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_X^{(r)}(\mathbf{y})), \quad (200)$$

and thus we find that, provided $\epsilon_r = 1$,

$$n^{1/2}(\mathbf{Y}_{\mathbf{A}}^{(r)} \mathbf{R}_r^{-1} - \mathbf{Y}^{(r)})_i^\top \rightarrow \mathcal{N}(\mathbf{0}, \boldsymbol{\Lambda}_r^{-1} \Delta_X^{-1} \boldsymbol{\Sigma}_X^{(r)}(\mathbf{y}) \Delta_X^{-1} \boldsymbol{\Lambda}_r^{-\top}) \quad (201)$$

almost surely. We deduce the Central Limit Theorem by integrating over all possible values of $\mathbf{y} \in \mathcal{Y}_r$. \square