

Gradient-Based Spectral Embeddings of Random Dot Product Graphs

Marcelo Fiori, Bernardo Marengo, Federico Larroca, Paola Bermolen, and Gonzalo Mateos

Abstract—The Random Dot Product Graph (RDPG) is a generative model for relational data, where nodes are represented via latent vectors in low-dimensional Euclidean space. RDPGs crucially postulate that edge formation probabilities are given by the dot product of the corresponding latent positions. Accordingly, the *embedding* task of estimating these vectors from an observed graph is typically posed as a low-rank matrix factorization problem. The workhorse Adjacency Spectral Embedding (ASE) enjoys solid statistical properties, but it is formally solving a surrogate problem and can be computationally intensive. In this paper, we bring to bear recent advances in non-convex optimization and demonstrate their impact to RDPG inference. We advocate first-order gradient descent methods to better solve the embedding problem, and to organically accommodate broader network embedding applications of practical relevance. Notably, we argue that RDPG embeddings of directed graphs lose interpretability unless the factor matrices are constrained to have orthogonal columns. We thus develop a novel feasible optimization method in the resulting manifold. The effectiveness of the graph representation learning framework is demonstrated on reproducible experiments with both synthetic and real network data. Our open-source algorithm implementations are scalable, and unlike the ASE they are robust to missing edge data and can track slowly-varying latent positions from streaming graphs.

Index Terms—Graph representation learning, gradient descent, manifold optimization, random dot product graphs.

I. INTRODUCTION

DURING the last few years the Random Dot Product Graph (RDPG) model has emerged as a popular generative model for random graphs [3]. This latent position model associates each node $i \in \{1, \dots, N\}$ with a vector $\mathbf{x}_i \in \mathcal{X} \subset \mathbb{R}^d$, a so-termed node embedding where typically $d \ll N$. In its simplest rendition for undirected and unweighted graphs without self loops, RDPGs specify an edge exists between nodes i and j with probability given by the inner-product of the corresponding embeddings; independently of all other edges. That is to say, entry A_{ij} of the random adjacency matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ has the Bernoulli($\mathbf{x}_i^\top \mathbf{x}_j$) distribution.

This apparent simplicity does not compromise expressive power. Indeed, one can verify that Erdős-Rényi (ER) graphs or Stochastic Block Models (SBMs) with a positive semi-definite

(PSD) probability matrix are particular cases of an RDPG. Several other more sophisticated models may be included as particular cases of RDPG [3], being this expressiveness one of the main reasons behind its popularity. A second reason is the model’s interpretability. Since the connection probability is given by the dot product of the embeddings, the affinity between the corresponding nodes is directly captured by their alignment. We may for instance rely on visual inspection of the nodes’ vector representations (possibly after further dimensionality reduction if $d > 3$) to screen for community structure, or carry out angle-based clustering of nodes [4], [5].

The restriction to undirected graphs is overcome by considering the directed version of RDPG, where each node is assigned *two* vectors $\mathbf{x}_i^l, \mathbf{x}_i^r \in \mathcal{X} \subset \mathbb{R}^d$. A directed edge from node i to j exists with probability $(\mathbf{x}_i^l)^\top \mathbf{x}_j^r$ [3]. The interpretation is analogous to the undirected case, with \mathbf{x}_i^l and \mathbf{x}_i^r now representing the propensity of node i towards establishing outgoing and accepting incoming directed edges, respectively. For extensions to weighted graphs, see e.g., [6].

Rather than generating random graphs from vectors, focus in this paper is on the inverse *embedding problem* at the heart of graph representation learning (GRL) [7]: given the adjacency matrix \mathbf{A} of a graph adhering to an RDPG model, estimate the latent position vectors for each node. Because of the RDPG definition in terms of dot products, the latent position vectors are only identifiable up to a common rotation. For both undirected and directed graphs (digraphs), the workhorse approach is based on a spectral decomposition of \mathbf{A} – an estimator known as Adjacency Spectral Embedding (ASE) [3].

A. Challenges facing the ASE

Although the ASE is widely adopted and its statistical properties (such as consistency and its limiting distribution as $N \rightarrow \infty$) are well documented [3], it does present drawbacks which we seek to overcome.

Large data. The first challenge pertains to scalability. Computing the spectrum of a large adjacency matrix \mathbf{A} , even only the d dominant components, is computationally intensive and constitutes a bottleneck of state-of-the-art ASE implementations [8], especially when multiple graphs are to be embedded. Recent work explicitly comments on the difficulty of scaling spectral-based inference of RDPGs to large graph settings [9].

Missing data. A second drawback of ASE is its inability to properly account for missing data, meaning unobserved entries in \mathbf{A} . On a related note, the ASE neglects the all-zeros diagonal in the adjacency matrix. These limitations were recognized more than a decade ago [4], yet to the best of our knowledge they have not been satisfactorily addressed in the RDPG literature. Indeed, repeated ASE computation to

Work in this paper is supported in part by CSIC (I+D project 22520220100076UD) and the NSF awards CCF-1750428, CCF-1934962. Part of the results in this paper appeared at the 2022 *EUSIPCO* and *Asilomar Conferences* [1], [2]. (Corresponding author: Gonzalo Mateos).

Marcelo Fiori, Bernardo Marengo, Federico Larroca, and Paola Bermolen are with the Facultad de Ingeniería, Universidad de la República, Montevideo 11000, Uruguay (e-mail: mfiori@fing.edu.uy; bmarengo@fing.edu.uy; flarroca@fing.edu.uy; paola@fing.edu.uy). Fiori, Marengo and Bermolen are also with the Centro Interdisciplinario en Ciencia de Datos y Aprendizaje Automático (CICADA), Universidad de la República, Uruguay.

Gonzalo Mateos is with the Department of Electrical and Computer Engineering, University of Rochester, Rochester, NY 14627 USA (e-mail: gmateosb@ece.rochester.edu).

iteratively impute the unknown entries of \mathbf{A} using the inner-product of the embeddings estimated in the previous step lacks convergence guarantees, and multiplies the ASE complexity by the number of iterations [4].

Streaming data. A third scenario that ASE cannot address satisfactorily arises with streaming data from a dynamic network; i.e., when we observe a sequence of graphs over time and would like to track the evolution of the corresponding embeddings, ideally without having to store past observations. Network dynamics may include changes in the edges between a fixed set of nodes (e.g., monitoring a wireless network), the addition of new information (e.g., a user that ranks an item in a recommender system), or the deletion/addition of nodes (e.g., a new user in a social network). Especially for large graphs, re-computing the ASE from scratch each time step is computationally demanding. Given the rotational ambiguity inherent to RDPGs, independently obtaining the ASE after each modification to the graph will likely result in misaligned embeddings that can hinder the assessment of changes.

B. Contributions and paper outline

We seek to address these limitations by (i) re-considering the underlying optimization problem of which ASE is a solution (Section II); and (ii) developing iterative embedding algorithms for the refined formulations (Sections III and IV).

Unlike the ASE recipe of performing a spectral decomposition of \mathbf{A} , borrowing recent low-rank matrix-factorization advances we advocate solving the non-convex embedding problem using gradient descent (GD) [10], [11]. Explicitly solving the optimization problem facilitates more precise and flexible GRL; e.g., it is straightforward to accommodate unobserved edges by including a mask matrix. Given the iterative nature of GD, warm-restarting the estimates of either new or existing nodes allows to embed multiple (possibly streaming) graphs, while preserving the alignment of consecutive embeddings as a byproduct. Discarding the residuals corresponding to the diagonal of \mathbf{A} offers better nodal representations as well as favorable problem structure, which we leverage in Section III-B to derive block-coordinate descent (BCD) iterations for efficient undirected RDPG inference.

Applying GD to embed digraph nodes requires special care. As we argue in Section IV-A, inherent ambiguities in the directed RDPG model extend beyond a global rotation, and they may compromise representation quality and the interpretability benefits discussed earlier. We show that an effective way of retaining these desirable features is to impose orthogonality constraints on the matrix factors in the decomposition of \mathbf{A} – a novel GRL formulation for digraphs. This constraint in turn defines a smooth manifold, over which we optimize using a custom-made feasible method. We stress this is not the well-known Stiefel manifold, where matrices are constrained to be *orthonormal* (and not just orthogonal as here¹). This is no minor point. Algorithm construction thus requires careful definition of the tangent space, the Riemannian gradient and the retraction [12], [13], all of which we derive in Section

¹We will henceforth use the term *orthonormal matrix* to refer to any matrix \mathbf{T} such that $\mathbf{T}^\top \mathbf{T} = \mathbf{I}$ (i.e., the columns of \mathbf{T} are orthonormal vectors). The term *orthogonal matrix* will be reserved for those matrices whose columns are mutually orthogonal, but not necessarily of unit norm.

IV-B. Comprehensive synthetic and real-world (wireless network and United Nations voting) data experiments in Section V demonstrate the interpretability, robustness, and versatility of the novel GRL framework. In the interest of reproducible research, the code and datasets used to generate all figures in this paper is publicly available at <https://github.com/marfiori/efficient-ASE>. Concluding remarks are in Section VI. Non-essential mathematical details and supplementary experimental results are deferred to the Appendix.

All in all, relative to prior our RDPG embedding framework offers a *better* representation at a *competitive* computational cost, and it is applicable to *more general* settings. This full-blown journal paper extends our preliminary results [1], [2] in multiple significant ways. In addition to a more thorough treatment with extended discussions, unpublished proofs, accompanying software, and expanded numerical studies with new datasets; the BCD algorithm for undirected graphs as well as the treatment of digraphs in Section IV are completely new.

II. PROBLEM STATEMENT AND RELATED WORK

Let us formulate the embedding problem, beginning with undirected graphs. Consider stacking all the nodes’ latent position vectors in the matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times d}$. Given an observed graph \mathbf{A} and a prescribed embedding dimension d (typically obtained using an elbow rule on \mathbf{A} ’s eigenvalue scree plot [8]), the goal is to estimate \mathbf{X} . Recalling the definition of the RDPG model, the edge-wise formation probabilities are the entries $P_{ij} = \mathbf{x}_i^\top \mathbf{x}_j$ of the rank- d , PSD matrix $\mathbf{P} = \mathbf{X}\mathbf{X}^\top$. Since we do not allow for self-loops, the diagonal entries in \mathbf{A} should be zero and we thus have $\mathbb{E}[\mathbf{A} | \mathbf{X}] = \mathbf{M} \circ \mathbf{P}$, where \circ is the entry-wise or Hadamard product and $\mathbf{M} = \mathbf{1}_N \mathbf{1}_N^\top - \mathbf{I}_N$ is a mask matrix with ones everywhere except in the diagonal where it is zero.

In lieu of a maximum likelihood estimator that is computationally challenging and may not be unique [?], here we advocate a least-squares (LS) approach [4] between

$$\hat{\mathbf{X}} \in \underset{\mathbf{X} \in \mathbb{R}^{N \times d}}{\operatorname{argmin}} \left\| \mathbf{M} \circ (\mathbf{A} - \mathbf{X}\mathbf{X}^\top) \right\|_F^2. \quad (1)$$

In words, $\hat{\mathbf{P}} = \hat{\mathbf{X}}\hat{\mathbf{X}}^\top$ is the best rank- d PSD approximation to the off-diagonal entries of the adjacency matrix \mathbf{A} , in the Frobenius-norm sense. The RDPG model is only identifiable up to rotations, and accordingly the solution of (1) is not unique. Indeed, for any orthonormal matrix $\mathbf{T} \in \mathbb{R}^{d \times d}$ we have that $\mathbf{X}\mathbf{T}(\mathbf{X}\mathbf{T})^\top = \mathbf{X}\mathbf{X}^\top = \mathbf{P}$.

Entrywise multiplication with $\mathbf{M} = \mathbf{1}_N \mathbf{1}_N^\top - \mathbf{I}_N$ effectively discards the residuals corresponding to the diagonal entries of \mathbf{A} . If suitably redefined, the binary mask \mathbf{M} may be used for other purposes, such as modeling unknown edges if data are missing. For instance, in a recommender system we typically have the rating of each user over a limited number of items. This (dis)information may be captured in (1) by zeroing out the entries of \mathbf{M} corresponding to the unknown edges.

Remark 1 (Adjacency Spectral Embedding) Typically the mask \mathbf{M} is ignored (and sometimes non-zero values are iteratively imputed to the diagonal of \mathbf{A} [4]), which results in a closed-form solution for $\hat{\mathbf{X}}$. Indeed, if we let $\mathbf{M} = \mathbf{1}_N \mathbf{1}_N^\top$ in (1), we have that $\hat{\mathbf{X}} = \hat{\mathbf{V}} \hat{\mathbf{\Lambda}}^{1/2}$, where $\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^\top$ is the eigendecomposition of \mathbf{A} , $\hat{\mathbf{\Lambda}} \in \mathbb{R}^{d \times d}$ is a diagonal matrix with

the d largest-magnitude eigenvalues of \mathbf{A} , and $\hat{\mathbf{V}} \in \mathbb{R}^{N \times d}$ are the associated eigenvectors. This workhorse estimator is known as the Adjacency Spectral Embedding (ASE); see also [3] for consistency and asymptotic normality results, as well as applications of statistical inference with RDPGs.

Given the ASE limitations outlined in Section I-A, we develop efficient gradient-based iterative solvers for the embedding problem (1). Beyond scalability, the algorithmic framework offers a natural pathway to facilitate embedding graph sequences. In the applications we study in Section V, said dynamic network data may be only partially observed, they could be acquired in a streaming fashion, while both the number of nodes and edges are allowed to vary across time.

Embedding digraphs. Moving on to digraphs [14], recall that we now embed each node with two vectors, $\mathbf{x}_i^l, \mathbf{x}_i^r \in \mathcal{X} \subset \mathbb{R}^d$. Existence of a directed edge from node i to j is modeled as the outcome of a Bernoulli trial with success probability $(\mathbf{x}_i^l)^\top \mathbf{x}_j^r$ [3]. Again, vertically stacking the embeddings into two matrices $\mathbf{X}^l, \mathbf{X}^r \in \mathbb{R}^{N \times d}$, we introduce the probability matrix $\mathbf{P} = \mathbf{X}^l (\mathbf{X}^r)^\top$ such that the expected value of the random adjacency matrix is $\mathbb{E}[\mathbf{A} \mid \mathbf{X}^l, \mathbf{X}^r] = \mathbf{M} \circ \mathbf{P}$.

If we ignore the mask \mathbf{M} , the embedding problem boils down to finding the best rank- d approximation to the (possibly asymmetric) adjacency matrix. One such solution may be obtained via the singular value decomposition (SVD) of \mathbf{A} ; i.e., $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$. We thus have that $\hat{\mathbf{X}}^l = \hat{\mathbf{U}} \hat{\mathbf{\Sigma}}^{1/2}$ and $\hat{\mathbf{X}}^r = \hat{\mathbf{V}} \hat{\mathbf{\Sigma}}^{1/2}$, with $\hat{\mathbf{\Sigma}}$ containing only the d largest singular values, and $\hat{\mathbf{U}}$ and $\hat{\mathbf{V}}$ the associated singular vectors. This procedure is also known as ASE.

As noted in [6], ambiguities with directed RDPGs can be more problematic than in the undirected case. Now given *any* invertible matrix \mathbf{T} (not necessarily orthonormal), the embeddings $\mathbf{Y}^l = \mathbf{X}^l \mathbf{T}$ and $\mathbf{Y}^r = \mathbf{X}^r \mathbf{T}^{-\top}$ generate the same probability matrix as before; i.e., $\mathbf{Y}^l (\mathbf{Y}^r)^\top = \mathbf{X}^l \mathbf{T} (\mathbf{X}^r \mathbf{T}^{-\top})^\top = \mathbf{X}^l (\mathbf{X}^r)^\top = \mathbf{P}$. As we show in Section IV-A, constraining matrices \mathbf{X}^l and \mathbf{X}^r to being orthogonal and having the same column-wise norms² effectively limits this ambiguity without compromising the model's expressivity (now an admissible \mathbf{T} may only be orthonormal; see Proposition 2), all while preserving its interpretability. Given these considerations, our approach to embedding digraphs is to solve the following manifold-constrained optimization problem

$$\begin{aligned} \{\hat{\mathbf{X}}^l, \hat{\mathbf{X}}^r\} \in \underset{\{\mathbf{X}^l, \mathbf{X}^r\} \in \mathbb{R}^{N \times d}}{\operatorname{argmin}} \quad & \|\mathbf{M} \circ (\mathbf{A} - \mathbf{X}^l (\mathbf{X}^r)^\top)\|_F^2 \\ \text{s. to } & (\mathbf{X}^l)^\top \mathbf{X}^l = (\mathbf{X}^r)^\top \mathbf{X}^r \text{ diagonal.} \end{aligned} \quad (2)$$

In the absence of a mask, a solution of (2) is the legacy ASE. Indeed, $\hat{\mathbf{X}}^l$ and $\hat{\mathbf{X}}^r$ are obtained from orthonormal singular vectors and have the same column-wise norms because both $\hat{\mathbf{U}}$ and $\hat{\mathbf{V}}$ are right-multiplied by $\hat{\mathbf{\Sigma}}^{1/2}$. To tackle the general case, a novel Riemannian GD method over the manifold of matrices with orthogonal columns is developed in Section IV-B.

A. Related work

The low-rank matrix factorization problem (1) has a long history, with applications to recommender systems (where the

objective is to complete a matrix of user-item ratings which is assumed to have low rank) [15]; or, in sensor localization from pairwise distances (the so-called Euclidean distance matrix) [16], just to name a couple of examples. Solution methods typically rely on spectral decomposition of the full data matrix (as in ASE), or by considering a convex relaxation via nuclear-norm minimization [17]. The latter is not best suited for our problem, where we are interested in the actual factors (not in \mathbf{P}), and their dimensionality could change with time due to e.g., node additions. Alternatively, over the last few years we have witnessed increased interest in non-convex optimization approaches for matrix factorization problems [10]. Our work may be seen as an effort in this direction. In particular, we bring to bear recent advances in first-order methods for matrix factorization problems and demonstrate impact to GRL (specifically, RDPG inference). The formulation (2) is novel to the best of our knowledge. To solve it we derive GD iterations over the manifold of orthogonal matrices, which is different from the Stiefel manifold and thus requires careful treatment given the unique geometric properties of our problem.

The scalability of ASE, or any other spectral embedding method for that matter, has long been considered an issue [18]. This challenge is compounded when multiple graphs are to be embedded, especially in *batch* settings where all graphs in the sequence are stored in memory [9]. Existing approaches seeking aligned embeddings rely on the spectral decomposition of some matrix whose dimension(s) grow linearly with the number of graphs [9], [19], [20]. In addition to increasing the computation cost of ASE, these methods are not applicable in streaming scenarios, where a possibly infinite sequence of graphs $\{\mathbf{A}_t\}$ is observed and we want to recursively update the embeddings ‘on-the-fly’ as new graphs are acquired.

There is an extensive collection of numerical linear algebra approaches to recursively update an eigendecomposition or SVD when the (adjacency) matrix is perturbed; e.g., [18]. However, these do not offer major computational savings except for specific types of changes (e.g., rank-1 updates), and they may be prone to error accumulation as t increases [21]. Moreover, they can yield misaligned embeddings due to the rotational ambiguity of RDPGs. The sketching literature offers highly-scalable randomized algorithms [22]. Other than to initialize our iterative methods we do not consider those here, because we are interested in exact solutions to (1) and (2).

In dynamic environments, not only does \mathbf{A}_t change over time, but new nodes may be added to the graph and current ones removed. Embedding previously unseen nodes corresponds to an inductive learning setting, generally regarded as beyond the capabilities of shallow embeddings as the one we are discussing here [7, Ch. 3.4], [23]. Previous efforts in this direction (that avoid re-computing eigendecompositions from scratch) either assume that the connections between the existing nodes, or, their current embeddings, do not change [24], [25]. In the latter case, a projection scheme onto the space of current embeddings produces an asymptotically ($N \rightarrow \infty$) consistent ASE for the new node [25]. However, even if latent positions were time invariant, the estimation error of current nodes’ embeddings propagates to the new estimates. We will use the projection-based estimate in [25] as initialization for new nodes in our GD algorithms, demonstrating benefits in accuracy and stability especially as several nodes are added,

²Let $\bar{\mathbf{x}}_i^l, \bar{\mathbf{x}}_i^r \in \mathbb{R}^N$ be the i -th columns of \mathbf{X}^l and \mathbf{X}^r , respectively. When we say \mathbf{X}^l and \mathbf{X}^r have the same column-wise norms we mean that $\|\bar{\mathbf{x}}_i^l\|_2 = \|\bar{\mathbf{x}}_i^r\|_2$ holds for all $i = 1, \dots, d$.

while at the same time refining previous nodes' embeddings.

III. EMBEDDING ALGORITHMS FOR UNDIRECTED GRAPHS

We start with the embedding problem for undirected graphs. Recognizing limitations of state-of-the-art ASE implementations, here we first review a GD algorithm with well-documented merits for symmetric matrix completion, yet so far unexplored in RDPG inference. GD offers flexible computation of embeddings and a pathway towards tracking nodal representations in a streaming graph setting. We then show that the particular structure of the problem lends itself naturally to more efficient BCD iterations, and discuss the relative merits of the different approaches in terms of convergence properties, complexity, and empirical execution time.

A. Back to basics: Estimation via gradient descent

Recall the embedding problem for undirected graphs in (1), and denote by $f : \mathbb{R}^{N \times d} \rightarrow \mathbb{R}$ its smooth objective function $f(\mathbf{X}) = \|\mathbf{M} \circ (\mathbf{A} - \mathbf{X}\mathbf{X}^\top)\|_F^2$. Although the problem is not convex with respect to \mathbf{X} , it is convex with respect to $\mathbf{P} = \mathbf{X}\mathbf{X}^\top$. In the broad context of matrix factorization problems where the objective function depends on the product $\mathbf{X}\mathbf{X}^\top$, the GD approach is often referred to as *factored GD* [26].

The workhorse GD algorithm generates embedding updates

$$\mathbf{X}_{k+1} = \mathbf{X}_k - \alpha \nabla f(\mathbf{X}_k), \quad k = 0, 1, 2, \dots \quad (3)$$

where $\alpha > 0$ is the stepsize, and the gradient of f is $\nabla f(\mathbf{X}) = 4 [\mathbf{M} \circ (\mathbf{X}\mathbf{X}^\top - \mathbf{A})] \mathbf{X}$. Recall that \mathbf{A} and \mathbf{M} are known symmetric matrices that specify the problem instance.

There have been several noteworthy advances in the study of GD's convergence (including rates) for this non-convex setting, as well as accelerated variants [10], [11], [26]–[29]. For the RDPG embedding problem dealt with here, it follows that if the initial condition \mathbf{X}_0 is close to the solution of (1), the GD iteration (3) converges linearly to $\hat{\mathbf{X}}$; see [26, Corollary 7], [11, Theorem 1], as well as [10, Lemma 4] and references therein for a similar version of this proposition.

Proposition 1 *Let $\hat{\mathbf{X}}$ be a solution of (1). Then there exist $\delta > 0$ and $0 < \kappa < 1$ such that, if $d(\mathbf{X}_0, \hat{\mathbf{X}}) \leq \delta$, we have*

$$d(\mathbf{X}_k, \hat{\mathbf{X}}) \leq \delta \kappa^k, \quad \forall k > 0 \quad (4)$$

where $\{\mathbf{X}_k\}$ are GD iterates (3) with appropriate constant stepsize, and $d(\mathbf{X}, \hat{\mathbf{X}}) := \min_{\mathbf{W}} \|\mathbf{X}\mathbf{W} - \hat{\mathbf{X}}\|_F^2$ s. to $\mathbf{W}^\top \mathbf{W} = \mathbf{I}_d$ is a matrix distance accounting for the rotational invariance.

Although there are specific \mathbf{X}_0 which correspond to sub-optimal stationary points, in our experience GD converges to the global optimum when initialized randomly. For further details on the initialization of factored GD, strong convexity assumptions, and the choice of the stepsize α , see e.g., [26, Section 5].

Remark 2 (Warm restarts to embed multiple graphs)

On top of being flexible to handle missing data encoded in \mathbf{M} , this approach also allows to track the latent positions of a graph sequence $\{\mathbf{A}_t\}$ using warm restarts. That is, after computing the embeddings \mathbf{X}_t of the graph with adjacency

matrix \mathbf{A}_t , we initialize the GD iterations (3) with \mathbf{X}_t to compute \mathbf{X}_{t+1} corresponding to \mathbf{A}_{t+1} . This way, unless the graphs at times t and $t+1$ are uncorrelated, the embedding of each node will transition smoothly to its new position (up to noise). Moreover, if the embeddings of the graph at time $t+1$ are sufficiently close to the embeddings at time t , say for slowly time-varying graphs where $d(\mathbf{X}_t, \mathbf{X}_{t+1}) \leq \delta$, then the GD iterates for computing \mathbf{X}_{t+1} also have the same convergence guarantees and rates ($\delta \kappa^k$), by virtue of Proposition 1. This was also observed empirically, indeed experiments demonstrating this longitudinal stability property [9] are presented in Sections V-B and V-C.

B. Block coordinate descent

Here we develop a BCD method for solving (1), which turns out to be quite efficient. The algorithm updates one row of \mathbf{X} at a time in a cyclic fashion, by minimizing the objective function f with respect to the corresponding row (treating all other entries of \mathbf{X} as constants, evaluated at their most recent updates). In general, this row-wise sub-problem may not admit a simple solution; however, we show that due to the structure of the mask matrix \mathbf{M} the updates are given in closed form.

Let $f(\mathbf{X}) = \|\mathbf{M} \circ (\mathbf{A} - \mathbf{X}\mathbf{X}^\top)\|_F^2$ and recall \mathbf{x}_i^\top is the i -th row of \mathbf{X} . The gradient $\nabla_i f$ of f with respect to \mathbf{x}_i is

$$\begin{aligned} \nabla_i f(\mathbf{X}) &= (-(\mathbf{M} \circ \mathbf{A})_i \mathbf{X} + ((\mathbf{M} \circ \mathbf{X}\mathbf{X}^\top)_i \mathbf{X})_i)^\top \\ &= -\mathbf{X}^\top (\mathbf{M} \circ \mathbf{A})_i^\top + \mathbf{X}^\top (\mathbf{M} \circ \mathbf{X}\mathbf{X}^\top)_i^\top, \end{aligned} \quad (5)$$

where $(\cdot)_i$ stands for the i -th row of the matrix argument.

Note that since the graph has no self-loops (i.e., the diagonal entries of \mathbf{A} are zero), then the entry-wise product of \mathbf{A} with \mathbf{M} is vacuous over the diagonal. Also because $A_{ii} = 0$, the term $\mathbf{X}^\top (\mathbf{M} \circ \mathbf{A})_i^\top = \mathbf{X}^\top (\mathbf{A})_i^\top$ in (5) does not depend on \mathbf{x}_i . More importantly, $\mathbf{X}\mathbf{X}^\top$ clearly depends on \mathbf{x}_i , and this would challenge solving $\nabla_i f(\mathbf{X}) = \mathbf{0}_d$ to obtain a minimizer [due to the trilinear form of the second term in (5)]. However, a close re-examination of $\mathbf{X}^\top (\mathbf{M} \circ \mathbf{X}\mathbf{X}^\top)_i^\top$ suggests this purported challenge can be overcome. First, observe that

$$(\mathbf{M} \circ \mathbf{X}\mathbf{X}^\top)_i = \mathbf{x}_i^\top \mathbf{X}^\top - \mathbf{r}_i^\top,$$

where $\mathbf{r}_i \in \mathbb{R}^N$ is a column vector with zeros everywhere except in entry i , where it takes the value $\mathbf{x}_i^\top \mathbf{x}_i$. Hence,

$$\mathbf{X}^\top (\mathbf{M} \circ \mathbf{X}\mathbf{X}^\top)_i^\top = \mathbf{X}^\top (\mathbf{X}\mathbf{x}_i - \mathbf{r}_i) = (\mathbf{X}^\top \mathbf{X} - \mathbf{x}_i \mathbf{x}_i^\top) \mathbf{x}_i.$$

All in all, we have a simplified expression for the gradient

$$\nabla_i f(\mathbf{X}) = -\mathbf{X}^\top (\mathbf{A}_i)^\top + (\mathbf{X}^\top \mathbf{X} - \mathbf{x}_i \mathbf{x}_i^\top) \mathbf{x}_i. \quad (6)$$

Now, define $\mathbf{R} = \mathbf{X}^\top \mathbf{X} - \mathbf{x}_i \mathbf{x}_i^\top$ and notice this matrix does not depend on \mathbf{x}_i . Therefore, from (6) it follows that the equation $\nabla_i f(\mathbf{x}_i) = \mathbf{0}_d$ is *linear* in \mathbf{x}_i , namely $\mathbf{R}\mathbf{x}_i = \mathbf{X}^\top (\mathbf{A}_i)^\top$. The pseudo-code of the algorithm is tabulated under Algorithm 1.

The $d \times d$ matrix \mathbf{R} is invertible provided that \mathbf{X} has rank d . This also implies that the row-wise sub-problem has a unique minimizer, which is key to establish convergence of BCD to a stationary point [30, Prop. 2.7.1]. It is worth reiterating that this favorable linear structure is lost in the absence of a mask matrix \mathbf{M} (cf. ASE in Remark 1). Since in RDPG embeddings we typically have $d \ll N$, solving multiple $d \times d$ linear systems is affordable; especially when compared to matrix-vector operations of order $\Theta(N^\gamma)$, $\gamma > 1$, like in GD.

Algorithm 1 Block coordinate descent (BCD)

Require: Initial $\mathbf{X} \leftarrow \mathbf{X}_0$

- 1: Compute $\mathbf{R} = \mathbf{X}^\top \mathbf{X}$
- 2: **repeat**
- 3: **for** $i = 1 : N$ **do**
- 4: $\mathbf{R} \leftarrow \mathbf{R} - \mathbf{x}_i \mathbf{x}_i^\top$
- 5: $\mathbf{b} \leftarrow \mathbf{X}^\top (\mathbf{A}_i)^\top$
- 6: $\mathbf{x}_i \leftarrow$ solution of $\mathbf{R}\mathbf{x} = \mathbf{b}$
- 7: $\mathbf{R} \leftarrow \mathbf{R} + \mathbf{x}_i \mathbf{x}_i^\top$
- 8: **end for**
- 9: **until** convergence
- 10: **return** \mathbf{X} .

C. Complexity and execution time analyses

We compare four computational methods to obtain RDPG embeddings of undirected graphs: the ASE based on (i) full eigendecomposition, and (ii) truncated SVD as implemented in *Graspologic* [8]; (iii) GD initialized with the randomized-SVD (RSVD) [22] (we account for the RSVD in the execution time); and (iv) randomly initialized BCD as in Algorithm 1.

The full eigendecomposition of \mathbf{A} has worst-case $\Theta(N^3)$ complexity, while for sparse graphs the d -dominant components can be obtained with $\Theta(Nd)$ per-iteration cost. For GD, the per-iteration computational cost incurred to evaluate $\nabla f(\mathbf{X})$ is dominated by the matrix multiplications, which is $\Theta(N^2d)$ for a naïve implementation. The number of iterations depends on \mathbf{X}_0 , but even with a favorable initialization the runtime is still higher than the $\Theta(Nd)$ of truncated SVD-based ASE. A refined convergence-rate analysis of GD for the symmetric matrix completion problem is presented in [11].

Although in general it is tricky to compare the complexity of GD against BCD approaches, we can evaluate the per-iteration computational cost of both methods (for BCD this means a whole cycle over the rows of \mathbf{X} is completed). In both cases, each entry of the matrix \mathbf{X} is updated exactly once. Each cycle consist of N instances of $d \times d$ linear systems, so this is $\Theta(Nd^3)$ in the worst case. In addition, in our experience Algorithm 1 converges in fewer iterations than the GD method.

In Fig. 1 we compare the execution times of methods (i)-(iv) as a function of N , all the way to $N = 24000$. For ASE, we use the *SciPy* optimized implementation of the eigendecomposition in Python, as in state-of-the-art RDPG inference packages such as *Graspologic* [8]. Our GD and BCD implementations are in pure Python, not optimized for performance. For each N , we sampled several 2-block SBM graphs, with connection probabilities of $p = 0.5$ (within block) and $q = 0.2$ (across blocks). Community sizes are $N/3$ and $2N/3$. We let $d = 2$ in all cases. Results are averaged over 10 Monte Carlo replicates, and corresponding standard deviations are depicted in Fig. 1. In all cases, the methods converge to a solution of (1). The obtained cost function is very similar for each run, with slightly lower values for the GD and BCD methods because they are solving the problem with the zero-diagonal restriction. As expected, BCD exhibits competitive scaling with the truncated SVD-based ASE, and can embed graphs with $N = 20000$ nodes in just over a minute.

A moderately large graph, such as the one with $N = 24000$, is ideal to assess the effect of d on the computation time. Graphs of this scale are expected to have several communities,

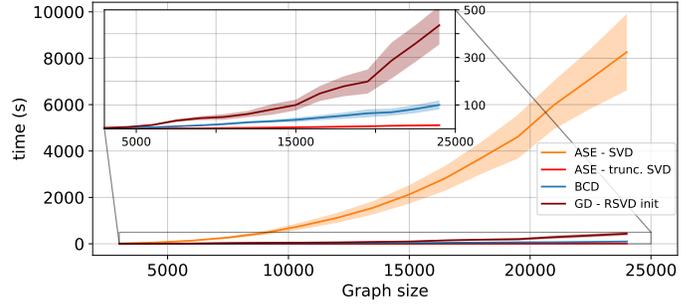


Fig. 1. Execution time for embedding SBM graphs with up to $N = 24000$ nodes. As N grows, BCD exhibits competitive scaling to the state-of-the-art ASE algorithm implemented in the *Graspologic* package.

TABLE I
EXECUTION TIME (IN SECONDS) AS A FUNCTION OF THE EMBEDDING DIMENSION d , FOR d -BLOCK SBM GRAPHS WITH $N = 24000$ NODES.

	$d = 10$	$d = 50$	$d = 100$
BCD - Algorithm 1	67 ± 4	98 ± 9	154 ± 32
ASE - Truncated SVD	14 ± 1	247 ± 21	466 ± 47

and thus values larger than $d = 2$ (as before) will likely be required. We thus explore this scenario in more detail, embedding a d -block SBM graph using d dimensions, and measure the execution time of BCD (Algorithm 1) and the truncated SVD methods as d increases. Results are reported in Table III-C. Interestingly, BCD yields faster results than truncated SVD when $d \geq 50$, gracefully scaling with the embedding dimension. As we mentioned before, in our experience BCD converges in very few iterations and offers competitive computation performance.

IV. EMBEDDING ALGORITHMS FOR DIGRAPHS

Shifting gears to embedding nodes in digraphs, we start with a close examination of the ambiguities inherent to the directed RDPG model and justify the need for orthogonality constraints on the factors' columns. To compute the desired nodal representations, we then develop a first-order feasible optimization method in the manifold of matrices with orthogonal columns.

A. On the interpretability of the directed RDPG

Recall that if $\{\hat{\mathbf{X}}^l, \hat{\mathbf{X}}^r\}$ is a minimizer of $f(\mathbf{X}^l, \mathbf{X}^r) = \|\mathbf{M} \circ (\mathbf{A} - \mathbf{X}^l (\mathbf{X}^r)^\top)\|_F^2$, then so is $\{\hat{\mathbf{X}}^l \mathbf{T}, \hat{\mathbf{X}}^r \mathbf{T}^{-\top}\}$ for any invertible matrix \mathbf{T} . Let us now discuss why \mathbf{X}^l and \mathbf{X}^r should be constrained to be orthogonal and have the same column-wise norms. In other words, why do we need the constraints in (2) to obtain useful embeddings when the graph is directed.

To gain some insights, suppose we ignore these constraints altogether and use GD to minimize $f(\mathbf{X}^l, \mathbf{X}^r)$. Similar to (3), at iteration $k + 1$ we update $\{\mathbf{X}_{k+1}^l, \mathbf{X}_{k+1}^r\}$ as follows

$$\mathbf{X}_{k+1}^l = \mathbf{X}_k^l - \alpha \nabla_{\mathbf{X}^l} f(\mathbf{X}_k^l, \mathbf{X}_k^r), \quad (7)$$

$$\mathbf{X}_{k+1}^r = \mathbf{X}_k^r - \alpha \nabla_{\mathbf{X}^r} f(\mathbf{X}_k^l, \mathbf{X}_k^r), \quad (8)$$

where $\nabla_{\mathbf{X}^l} f(\mathbf{X}^l, \mathbf{X}^r) = 4 [\mathbf{M} \circ (\mathbf{X}^l (\mathbf{X}^r)^\top - \mathbf{A})] \mathbf{X}^l$ and a similar expression holds for $\nabla_{\mathbf{X}^r} f(\mathbf{X}^l, \mathbf{X}^r)$.

The ASE offers an alternative baseline, which requires to discard the mask \mathbf{M} . ASE estimates $\{\hat{\mathbf{X}}^l, \hat{\mathbf{X}}^r\}$ have orthogonal columns because they are derived from the SVD of \mathbf{A} . Same

index columns in $\hat{\mathbf{X}}^l$ and $\hat{\mathbf{X}}^r$ have the same norm as well, since the orthonormal matrices $\hat{\mathbf{U}}$ and $\hat{\mathbf{V}}$ are *both* right-multiplied by $\hat{\Sigma}^{1/2}$. However, if we minimize $f(\mathbf{X}^l, \mathbf{X}^r)$ iteratively as in (7)-(8) to accommodate missing and streaming data, we may lose column-wise orthogonality with detrimental consequences we illustrate in the following example.

Example 1 (Bipartisan senate) Consider a synthetic political dataset of votes cast by senators in support of laws, over a certain period of time. This may be regarded as a bipartite digraph where nodes correspond to senators and laws, and the fact that senator i has voted affirmatively for law j is indicated by the existence of edge (i, j) . We examine a polarized scenario, where two political parties put forth laws for voting. Affirmative votes are more likely for senators from the party that introduced the law, and less likely for senators from the opposing party. There are also a few bipartisan laws, for which affirmative votes tend to be more balanced across parties. We simulated such a graph with 50 senators of each party, where Party 1 proposed 50 laws and Party 2 proposed 200 laws, and there were 40 additional bipartisan laws under consideration (i.e., $N = 390$ in total). Furthermore, the inter-community probability matrix is

$$\mathbf{\Pi} = \begin{pmatrix} 0 & 0 & 0.9 & 0.01 & 0.2 \\ 0 & 0 & 0.1 & 0.8 & 0.3 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix},$$

where communities are ordered as Party 1 senator, Party 2 senator, Party 1 law, Party 2 law, and finally bipartisan law. In other words, senators of Party 1 are very supportive of their own laws and unlikely to vote for those introduced in the other side of aisle, whereas Party 2 senators are more moderate.

We compare the embeddings estimated through ASE and GD [i.e., iterating (7) and (8) until convergence]. Recall that interpretation of these results should rely on the geometry induced by the RDPG model. Similarity among nodes is captured by their colinearity in latent space, not by their Euclidean distance being small (as it is the case with e.g., Laplacian eigenmaps [?]). Accordingly, in this particular example we expect that the embeddings of Party 1 senators and laws will be almost perfectly aligned, while slightly less so for Party 2. ASE results using $d = 2$ are shown in Figure 2 (left). As expected, the outward embeddings for laws and inward embeddings for senators are zero (since the former do not vote and the latter are not voted). Furthermore, the outward embeddings corresponding to senators of each party are close and roughly orthogonal to senators of the other party, reflecting the polarized landscape. Finally, the inward embeddings of laws submitted by each party are aligned with the corresponding cluster of senators, whereas bipartisan laws lie somewhere between both groups. The difference in magnitude between embeddings of senators and laws is due to the different number of such nodes in the graph, and the column-wise norm constraint imposed to the embeddings.

On the other hand, inspection of Fig. 2 (right) reveals that GD converges to a solution where laws are not aligned with the corresponding party senators. Accordingly, the affinity of parties to their laws is less evident than before. In fact, it appears as if Party 1 is not as supportive of its laws as in the ASE-based visualization which, as we discussed before, is the opposite of what we expected. While the input graph is the

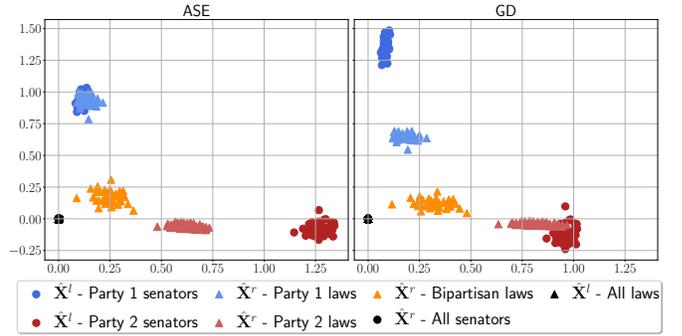


Fig. 2. Bipartisan senate example. ASE (left) and GD (right). Since ASE implicitly imposes equally normed orthogonal columns (as it is derived from an SVD), it produces interpretable embeddings. On the other hand, GD may result in laws and parties that are not aligned, and thus loses interpretability if no further constraints are imposed in the formulation.

same for both methods and the total cost $f(\hat{\mathbf{X}}^l, \hat{\mathbf{X}}^r)$ is smaller for the GD method, interpretability is hindered because of the larger ambiguity set in the absence of additional constraints.

Example 2 (Digraph with symmetric expectation) To further justify why orthogonality constraints are essential, consider a digraph sampled from a symmetric \mathbf{P} (i.e., the probability of both edge directions is the same, but each arc is independently sampled). It would be desirable that in this case the model enforced $\mathbf{X}^l = \mathbf{X}^r$, since the outgoing and incoming behaviour of the nodes is the same. The general directed model should recover the subsumed undirected one and, naturally, the same should hold for the embedding method.

However, these desiderata are not necessarily met. As stated earlier, given an invertible matrix \mathbf{T} , the embeddings $\mathbf{Y}^l = \mathbf{X}^l \mathbf{T}$ and $\mathbf{Y}^r = \mathbf{X}^r \mathbf{T}^{-\top}$ yield the same probability matrix $\mathbf{P} = \mathbf{X}^l (\mathbf{X}^r)^\top$. This implies that unless $\mathbf{T} = \mathbf{T}^{-\top}$ (meaning \mathbf{T} is an orthonormal matrix corresponding to a rotation), we could apply *different* transformations to the inward and outward embeddings and still obtain the same RDPG.

Given these observations, consider a directed RDPG model where the embedding matrices \mathbf{X}^l and \mathbf{X}^r are constrained to be orthogonal and of the same column-wise norm. The following result asserts that this suffices to ensure an admissible \mathbf{T} is orthonormal, hence reducing the model's ambiguity to a global rotation – just like in the undirected case.

Proposition 2 *Let $\mathbf{P} = \mathbf{X}^l (\mathbf{X}^r)^\top$ be the probability matrix of a directed RDPG model, where $\{\mathbf{X}^l, \mathbf{X}^r\}$ are $N \times d$ matrices with rank d such that $(\mathbf{X}^l)^\top \mathbf{X}^l = (\mathbf{X}^r)^\top \mathbf{X}^r = \mathbf{D}_X$ is diagonal. Let $\mathbf{T} \in \mathbb{R}^{d \times d}$ be an invertible matrix such that $\mathbf{Y}^l = \mathbf{X}^l \mathbf{T}$ and $\mathbf{Y}^r = \mathbf{X}^r \mathbf{T}^{-\top}$ are also orthogonal with the same column-wise norms; i.e. $(\mathbf{Y}^l)^\top \mathbf{Y}^l = (\mathbf{Y}^r)^\top \mathbf{Y}^r = \mathbf{D}_Y$ is diagonal. Then, \mathbf{T} is an orthonormal matrix.*

Proof: Combining $\mathbf{Y}^r = \mathbf{X}^r \mathbf{T}^{-\top}$ with $(\mathbf{X}^r)^\top \mathbf{X}^r = \mathbf{D}_X$ and $(\mathbf{Y}^r)^\top \mathbf{Y}^r = \mathbf{D}_Y$ we find that $\mathbf{D}_X = \mathbf{T} \mathbf{D}_Y \mathbf{T}^\top$. Proceeding analogously with \mathbf{Y}^l we further obtain that $\mathbf{D}_X = \mathbf{T}^{-\top} \mathbf{D}_Y \mathbf{T}^{-1}$. Multiplying both identities results in $\mathbf{D}_X^2 = \mathbf{T} \mathbf{D}_Y^2 \mathbf{T}^{-1}$. Thus, the columns of \mathbf{T} are linearly independent eigenvectors of a diagonal matrix. Furthermore, given

$\mathbf{D}_X = \mathbf{T}\mathbf{D}_Y\mathbf{T}^\top$ it follows that the above eigendecomposition is necessarily one with orthonormal eigenvectors. ■

The constraints in (2) do not limit the expressiveness of the model, since they are compatible with those ASE implicitly imposes. Next, we develop a feasible first-order method that enforces the orthogonality constraint at all iterations. After convergence it is straightforward to equalize the resulting column-wise norms so that they are the same for both $\hat{\mathbf{X}}^l$ and $\hat{\mathbf{X}}^r$, without affecting the generated \mathbf{P} ; see Remark 3.

B. Optimizing on a manifold

We have concluded that for the sake of interpretability and quality of the representation, it is prudent to impose the matrices \mathbf{X}^l and \mathbf{X}^r have orthogonal columns. One classical way to tackle this is by adding these constraints to the optimization problem as in (2), and solve it via Lagrangian-based methods. For some constraints with geometric properties, a more suitable and timely approach is to pose the optimization problem on a smooth manifold. One can then resort to feasible methods that search exclusively over the manifold, i.e., the constraints are satisfied from the start and throughout the entire iterative minimization process [12], [13]. This way, we can think of the optimization as being *unconstrained* because the manifold is all there is. In the sequel we explore this last idea.

Interestingly, the space of matrices having orthogonal columns does not form any known and well-studied manifold. Yet, we show the required geometric structure is present in our problem and thus we have to define several objects as well as compute various operators to facilitate optimization [12], [13]. The conceptual roadmap is as follows. Recall that a smooth manifold \mathcal{M} can be locally approximated by a linear space, the so-called *tangent space*. If we consider the objective function $f : \mathcal{M} \mapsto \mathbb{R}$ defined from the (Riemannian) manifold to \mathbb{R} , then the Riemannian gradient of the function is an element of the tangent space. This *Riemannian gradient*, which will be denoted as $\text{grad } f$, can be computed as the projection of the Euclidean gradient ∇f to the tangent space. Having computed the gradient, a classical descent method consists of taking a certain step in the opposite direction. However, this step likely results in a point outside of the manifold, so we have to project it back to \mathcal{M} . This projection might be computationally intensive, so the *retraction* alternative is used instead.

Next, we define more precisely our manifold, and derive the tangent space, the projection and finally the retraction. The manifold that resembles the most to ours is the so-called Stiefel manifold, which consist of matrices with orthogonal and unit-norm (i.e., orthonormal) columns

$$St(d, N) := \{\mathbf{X} \in \mathbb{R}^{N \times d} : \mathbf{X}^\top \mathbf{X} = \mathbf{I}_d\}. \quad (9)$$

But here we do not require unit-norm columns. Thus, let $\mathbb{R}_*^N = \mathbb{R}^N \setminus \{\mathbf{0}_N\}$ be the set of N dimensional vectors without the null vector, and let $\mathbb{R}_*^{N \times d}$ be the product of d copies of \mathbb{R}_*^N . This open set is the set of $N \times d$ matrices without null columns. We are interested in matrices with orthogonal columns, namely

$$\begin{aligned} \mathcal{M}(d, N) &:= \{\mathbf{X} \in \mathbb{R}_*^{N \times d} : \mathbf{X}^\top \mathbf{X} \text{ is diagonal}\} \\ &= \{\mathbf{X} \in \mathbb{R}_*^{N \times d} : \mathbf{M} \circ (\mathbf{X}^\top \mathbf{X}) = \mathbf{0}_{d \times d}\}, \end{aligned} \quad (10)$$

where once more $\mathbf{M} = \mathbf{1}_d \mathbf{1}_d^\top - \mathbf{I}_d$ is a particular mask matrix, with zeros in the diagonal and ones everywhere else.

The following proposition establishes that \mathcal{M} is actually a manifold (for notational convenience, we henceforth use \mathcal{M} instead of $\mathcal{M}(d, N)$ since both d and N are fixed throughout). Moreover, \mathcal{M} is a Riemannian manifold since $\mathcal{M} \subset \mathbb{R}^{N \times d}$ is a vector space equipped with the usual trace inner product. The proofs of subsequent results can be found in the Appendix.

Proposition 3 *The set \mathcal{M} in (10) is a differential manifold and its dimension is $Nd - d(d-1)/2$. Furthermore, the tangent space at $\mathbf{X} \in \mathcal{M}$ is given by*

$$T_{\mathbf{X}}\mathcal{M} = \{\zeta \in \mathbb{R}^{N \times d} : \mathbf{M} \circ (\zeta^\top \mathbf{X} + \mathbf{X}^\top \zeta) = \mathbf{0}_{N \times d}\}.$$

To perform a manifold GD step, one needs to compute the Riemannian gradient of the function defined in \mathcal{M} . We obtain $\text{grad } f$ as the projection of the Euclidean gradient [cf. (7)-(8)] onto the tangent space. A natural way to compute said projection is to first characterize and compute the projection to the normal space. Given $\mathbf{X} \in \mathcal{M}$, the normal space at \mathbf{X} is $T_{\mathbf{X}}\mathcal{M}^\perp = \{\mathbf{N} \in \mathbb{R}^{N \times d} : \langle \mathbf{N}, \zeta \rangle = \text{tr}(\mathbf{N}^\top \zeta) = 0, \forall \zeta \in T_{\mathbf{X}}\mathcal{M}\}$. A useful alternative characterization is given next.

Lemma 1 *The normal space at \mathbf{X} is*

$$T_{\mathbf{X}}\mathcal{M}^\perp = \{\mathbf{X}\mathbf{\Lambda} \in \mathbb{R}^{N \times d} : \mathbf{\Lambda} \in \mathcal{S}_d\},$$

where $\mathcal{S}_d = \{\mathbf{X} \in \mathbb{R}^{d \times d} : \mathbf{X} = \mathbf{X}^\top, \text{diag}(\mathbf{X}) = \mathbf{0}_{d \times d}\}$ is the set of $d \times d$ symmetric matrices with null diagonal.

Computing the projection to the normal space requires some work due to the null diagonal constraint in \mathcal{S}_d , which is not present in the normal space to $St(d, N)$ [13, p. 161]. The result is given in the next lemma.

Lemma 2 *Let $\mathbf{X} \in \mathcal{M}$ and let $\pi_{\mathbf{X}}^\perp : \mathbb{R}^{N \times d} \mapsto T_{\mathbf{X}}\mathcal{M}^\perp$ be the projection to the normal space. Then*

$$\pi_{\mathbf{X}}^\perp(\mathbf{Z}) = \mathbf{X}s(2\mathbf{D}\mathbf{L}), \quad (11)$$

where $s : \mathbb{R}^{d \times d} \mapsto \mathcal{S}_d$ is a symmetrizing function $s(\mathbf{Z}) = \frac{\mathbf{Z} + \mathbf{Z}^\top}{2} - \text{diag}(\mathbf{Z})$, $\mathbf{D} = (\mathbf{X}^\top \mathbf{X})^{1/2}$ and $\mathbf{L} = (\mathbf{D}^{-1} \mathbf{X}^\top \mathbf{Z}) \circ \mathbf{F}$, where $\mathbf{E} = \mathbf{1}_d \mathbf{1}_d^\top \mathbf{D}^2 + \mathbf{D}^2 \mathbf{1}_d \mathbf{1}_d^\top$ and \mathbf{F} has entries $F_{ij} = E_{ij}^{-1}$.

Note that (11) is of the form $\mathbf{X}\mathbf{\Lambda}$, with $\mathbf{\Lambda} \in \mathcal{S}_d$. It thus belongs to the normal space by virtue of the characterization in Lemma 1. The calculations to show it is indeed the projection are detailed in the Appendix, and boil down to proving that $\mathbf{Z} - \pi_{\mathbf{X}}^\perp(\mathbf{Z})$ lives in the tangent space. Specifically, to establish (11) we take $\mathbf{X}\mathbf{\Lambda}$ and derive conditions that $\mathbf{\Lambda}$ had to verify when imposing that $\mathbf{Z} - \mathbf{X}\mathbf{\Lambda} \in T_{\mathbf{X}}\mathcal{M}$. After some derivations, we find $\mathbf{\Lambda} = s(2\mathbf{D}\mathbf{L})$, with the auxiliary matrices \mathbf{D} , \mathbf{E} , \mathbf{F} and \mathbf{L} defined in Lemma 2. Finally, the desired projection to the tangent space is given as follows.

Proposition 4 *Let $\mathbf{X} \in \mathcal{M}$. The projection to the tangent space $\pi_{\mathbf{X}} : \mathbb{R}^{N \times d} \mapsto T_{\mathbf{X}}\mathcal{M}$ can be computed as:*

$$\pi_{\mathbf{X}}(\mathbf{Z}) = \mathbf{Z} - \pi_{\mathbf{X}}^\perp(\mathbf{Z}) = \mathbf{Z} - \mathbf{X}s(2\mathbf{D}\mathbf{L}).$$

When we take a small step in the opposite direction of $\text{grad } f$, in general we fall outside \mathcal{M} and we have to project back to it. We need a projection from the tangent bundle to the manifold, or a retraction, which is more efficient in general.

Algorithm 2 Riemannian Gradient Descent (GD) on \mathcal{M}

Require: Initial \mathbf{X}_0^l and \mathbf{X}_0^r

- 1: **repeat**
 - 2: Compute Euclidean gradients
 $\nabla f_{\mathbf{X}^l}(\mathbf{X}_k^l, \mathbf{X}_k^r)$ and $\nabla f_{\mathbf{X}^r}(\mathbf{X}_k^l, \mathbf{X}_k^r)$
 - 3: Compute Riemannian gradients as
 $\text{grad } f_{\mathbf{X}^l}(\mathbf{X}_k^l, \mathbf{X}_k^r) = \pi_{\mathbf{X}_k^l}(\nabla f_{\mathbf{X}^l}(\mathbf{X}_k^l, \mathbf{X}_k^r))$
 $\text{grad } f_{\mathbf{X}^r}(\mathbf{X}_k^l, \mathbf{X}_k^r) = \pi_{\mathbf{X}_k^r}(\nabla f_{\mathbf{X}^r}(\mathbf{X}_k^l, \mathbf{X}_k^r))$
 - 4: Compute retraction with α chosen via the Armijo rule
 $\mathbf{X}_{k+1}^l = R_{\mathbf{X}_k^l}(-\alpha \text{grad } f_{\mathbf{X}^l}(\mathbf{X}_k^l, \mathbf{X}_k^r))$,
 $\mathbf{X}_{k+1}^r = R_{\mathbf{X}_k^r}(-\alpha \text{grad } f_{\mathbf{X}^r}(\mathbf{X}_k^l, \mathbf{X}_k^r))$
 - 5: **until** convergence
 - 6: **return** $\{\mathbf{X}_k^l, \mathbf{X}_k^r\}$.
-

Given a full rank matrix $\mathbf{Z} \in \mathbb{R}^{N \times d}$, consider its decomposition $\mathbf{Z} = \tilde{\mathbf{Q}}\tilde{\mathbf{R}}$, where $\tilde{\mathbf{Q}}$ is a matrix with orthogonal columns and $\tilde{\mathbf{R}}$ is upper triangular with ones in the diagonal. This decomposition is unique. Indeed, one may obtain $\tilde{\mathbf{Q}}$ by a Gram-Schmidt process, but skipping the normalization steps. A more efficient approach is to consider the classical QR decomposition ($\mathbf{Z} = \mathbf{Q}\mathbf{R}$, with \mathbf{Q} orthonormal and \mathbf{R} upper triangular), and compute $\tilde{\mathbf{Q}} = \mathbf{Q}\mathbf{D}_R$, where $\mathbf{D}_R = \text{diag}(\mathbf{R})$ is the diagonal matrix with the diagonal entries of \mathbf{R} . In a way, this modification of the QR decomposition shifts the “normalization” of the columns from the upper triangular factor towards the orthogonal factor.

Note that $\tilde{\mathbf{Q}} \in \mathcal{M}$ and this decomposition will serve to define a retraction to the manifold in the next proposition. Again, this procedure differs from the popular Q -factor retraction to the Stiefel manifold [13, p. 160].

Proposition 5 Let $\mathbf{X} \in \mathcal{M}$ and $\zeta \in T_{\mathbf{X}}\mathcal{M}$ a tangent vector. Then, the mapping

$$R_{\mathbf{X}}(\zeta) = \tilde{q}f(\mathbf{X} + \zeta)$$

is a retraction, where $\tilde{q}f(\mathbf{A})$ denotes the $\tilde{\mathbf{Q}}$ factor of the modified QR decomposition described above, and the sum $\mathbf{X} + \zeta$ stands for the usual abuse of notation for embedded manifolds on vector spaces.

We now have all the ingredients for the GD method to minimize $f(\mathbf{X}^l, \mathbf{X}^r) = \|\mathbf{M} \circ (\mathbf{A} - \mathbf{X}^l(\mathbf{X}^r)^\top)\|_F^2$ over \mathcal{M} , which is tabulated under Algorithm 2. The convergence rate of Riemannian GD is the same as the unconstrained counterpart (i.e., producing points with $\text{grad } f$ smaller than ε in $\mathcal{O}(1/\varepsilon^2)$ iterations) [31]. The computational complexity of each iteration is dominated by the QR decomposition in the retraction.

We extended the Pymanopt package [32] with a class for the manifold \mathcal{M} defined in Proposition 3, which forms part of the code available for this paper.

Remark 3 (Rescaling the factors’ columns) Algorithm 2 does not quite solve (2). While both $\{\mathbf{X}_k^l, \mathbf{X}_k^r\}$ belong to \mathcal{M} , the constraint $(\mathbf{X}_k^l)^\top \mathbf{X}_k^l = (\mathbf{X}_k^r)^\top \mathbf{X}_k^r$ will in general not be satisfied upon convergence. Dropping the iteration index for simplicity, let $\bar{\mathbf{x}}_i^l, \bar{\mathbf{x}}_i^r \in \mathbb{R}^N$ be the i -th columns of \mathbf{X}^l and \mathbf{X}^r , respectively. To obtain a feasible solution from the output of Algorithm 2, for each dimension $i = 1, \dots, d$ we

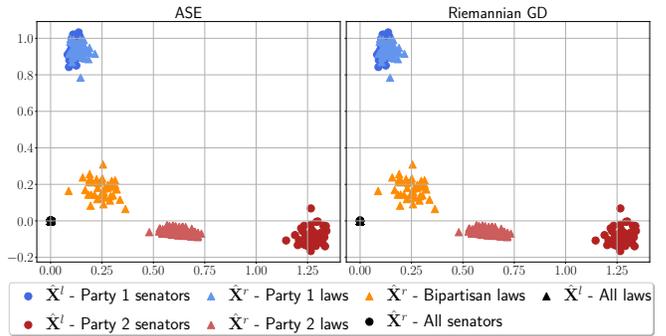


Fig. 3. Solution to the embedding problem (2) for the bipartisan senate example. ASE (left) and Riemannian GD (right). Notice how both solutions are nearly identical [cf. unconstrained GD in Fig. 2 (right)], underscoring the importance of the orthogonality constraints.

define scaling factors $s_i = \|\bar{\mathbf{x}}_i^l\|_2 / \|\bar{\mathbf{x}}_i^r\|_2$ and collect them in the diagonal matrix $\mathbf{S} = \text{diag}(s_1, \dots, s_d)$. We then rescale the columns of the embedding matrices via $\mathbf{X}_k^l \leftarrow \mathbf{X}_k^l \mathbf{S}^{-1/2}$ and $\mathbf{X}_k^r \leftarrow \mathbf{X}_k^r \mathbf{S}^{1/2}$, without affecting the value of the objective function but now satisfying the constraint in (2).

Example 3 (Bipartisan senate revisited) Going back to the bipartisan senate from Example 1, Fig. 3 depicts the solution of (2) for the same simulated bipartite senator-law digraph (imposing the orthogonality constraints and rescaling in Remark 3). Unlike in Example 1, the Riemannian GD algorithm on the manifold \mathcal{M} is able to recover the same structure as the ASE. Laws are now correctly aligned with their corresponding party, thus faithfully revealing the structure in the data.

V. NUMERICAL EXPERIMENTS AND APPLICATIONS

In this section we illustrate our embedding algorithms’ ability to produce accurate and informative estimates of nodal latent position vectors. We explore a variety of GRL applications and consider synthetic and real network data. Our test cases are designed to target ASE challenges outlined in Section I-A, namely: i) missing data; ii) embedding multiple networks; and iii) graph streams (with fixed and varying number of nodes). For each case we assess the results with respect to estimation accuracy, interpretability, and stability/alignment in dynamic environments. The suitability of spectral embeddings (rooted in the RDPG model) for downstream tasks has already been well-documented [3], [25]. For this reason, the goal here is to demonstrate the effectiveness of our algorithms in generating node embeddings that faithfully represent network structure in novel settings i)-iii). Supplementary results exploring algorithm sensitivity to random initialization are in Appendix E. The code for all these experiments is publicly available at <https://github.com/marfiori/efficient-ASE>.

A. Inference with missing data

First we illustrate how GD-based inference can be useful for GRL with missing data. The setup is similar to that of Example 1, but here we rely on real United Nations (UN) General Assembly voting data [33]. For each roll call and country, the dataset includes if the country was present and

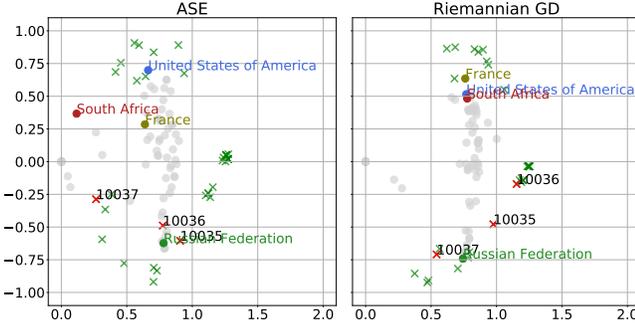


Fig. 4. UN General Assembly voting data for 1955. ASE (left) and Riemannian GD (i.e., Algorithm 2) with mask matrix encoding present and absent (or abstained) voters (right). Our approach is able to assign the absent voters to the correct group (e.g., South Africa) and offers a more clear clustering of roll calls.

if so the corresponding vote (either ‘Yes’, ‘No’, or ‘Abstain’) for each proposal. We analyze the associated bipartite digraph pertaining to a particular year, where nodes correspond to countries and roll calls, and an edge from a country to a roll call exists if it voted affirmatively. If the country was absent or abstained, we will tag that edge as unknown ($M_{ij} = 0$).

Fig. 4 depicts the node embeddings ($d = 2$) of the graph from 1955, estimated by ASE (naively assuming unknown edges do not exist, $A_{ij} = 0$) and Riemannian GD (i.e. Algorithm 2). Consider the countries, which are displayed as circles. We highlight four interesting cases: Russia, USA, France, and South Africa. At the time, the first two represented two poles of the world, and are naturally almost orthogonal to each other for both methods. Note furthermore how the ASE seems to indicate that South Africa is less likely to vote in agreement with Russia than (even) the USA, whereas the opposite is true for France.

The problem comes from equating an absence or abstention to a negative vote. For instance, South Africa was only present in roughly one third of the roll calls, and it voted almost identically to the USA. The Riemannian GD method, which acknowledges unknown edges via the mask \mathbf{M} , provides an embedding that reflects this agreement. Something similar happens with France, which differed from the USA only in six roll calls. Four correspond to USA abstentions and France voting ‘Yes’, another one where the opposite happened (and thus both cases should not be accounted for in the optimization problem), and finally the roll call 10036 where France was one of only two countries to vote ‘No’ (the USA voted ‘Yes’).

Regarding the embeddings of roll calls marked with a cross in Fig. 4, note how 10036 is aligned with the Russian block of countries by ASE, but it is better placed as an intermediate proposal in Fig. 4 (right) – equally likely to be voted by all countries. Something similar occurs with roll call 10035, which dealt with the same subject of 10036, but met resistance from more countries (roughly a dozen, including the USA and France). In both cases several countries were not present or abstained during the voting. Incorrectly assuming these votes as negative by ASE leads to biased results. Much more can be said about the roll calls and their associated UN resolutions, but let us conclude the discussion by noting that roll call embeddings generated by Algorithm 2 form three

clusters reflecting the geopolitical landscape at the time. There is a cluster for each pole (American and Russian), plus an intermediate one where both poles tend to vote similarly. On the other hand, ASE generates roll call embeddings that are incorrectly aligned (e.g., 10036), and a loose grouping of intermediate roll calls with shared voting from both poles.

B. Embedding multiple graphs: the batch case

Suppose now that we observe $m > 1$ graphs $\{\mathbf{A}_t\}_{t=1}^m$ and we are interested, for instance, in testing whether they are drawn from the same RDPG model, or, in tracking the embeddings over time. Assume that we can identify nodes across different observations; e.g., they correspond to labeled users in a social network and so a matching algorithm is not needed. Independently obtaining the ASE for each graph is undesirable because it yields arbitrarily rotated embeddings, a challenge that has motivated several recent research efforts.

Indeed, a hypothesis test which involves solving a Procrustes problem to align the embeddings was put forth in [34]. An alignment alternative is to jointly embed all m graphs via a single ‘super-matrix’ decomposition. The so-called *Omnibus* embedding first forms an $mN \times mN$ matrix derived from all $\{\mathbf{A}_t\}_{t=1}^m$, and then computes its ASE which enjoys asymptotic normality [19]. The Unfolded ASE (UASE) also constructs an auxiliary matrix, but by horizontally stacking all $\{\mathbf{A}_t\}_{t=1}^m$ [9], [20]. Nodal representations are then extracted from the SVD of this $N \times mN$ matrix. Under some technical assumptions, the UASE provably offers desirable longitudinal and cross-sectional stability [9]. However, the complexity and memory footprint of these *batch* approaches grow linearly with m , and they are only applicable to undirected graphs.

In the context of the algorithms proposed in this paper, we may leverage their iterative nature and initialize them using the estimated embeddings of another related (e.g., contiguous in time) graph. Unless radical changes take place from one graph to the other, this so-called warm restart is expected to produce embeddings that are closely aligned, with the added benefit of converging in few iterations.

Stability of BCD estimates. Let us illustrate this (desirable) behaviour through a numerical example. We borrow the setting and code from [9]. Consider two graph samples drawn from a dynamic SBM with inter-community probability matrices

$$\mathbf{\Pi}_1 = \begin{pmatrix} 0.08 & 0.02 & 0.18 & 0.10 \\ 0.02 & 0.20 & 0.04 & 0.10 \\ 0.18 & 0.04 & 0.02 & 0.02 \\ 0.10 & 0.10 & 0.02 & 0.06 \end{pmatrix}, \quad \mathbf{\Pi}_2 = \begin{pmatrix} 0.16 & 0.16 & 0.04 & 0.10 \\ 0.16 & 0.16 & 0.04 & 0.10 \\ 0.04 & 0.04 & 0.09 & 0.02 \\ 0.10 & 0.10 & 0.02 & 0.06 \end{pmatrix}.$$

Initially there are four communities. At time 2, the first two communities merge, community 3 moves, and community 4 has its connection probabilities unchanged.

Ideally, when embedding both graphs: i) the representations of nodes in community 4 should not change (longitudinal stability); and ii) the time 2 embeddings of members of communities 1 and 2 should be similar, up to noise (cross-sectional stability). Fig. 5 displays the results for UASE [9], Omnibus embedding [19], independent ASE for each graph, and BCD (i.e. Algorithm 1 warm-restarted at time 2 with the result of time 1). As expected, independent ASE lacks longitudinal stability, and the Omnibus embedding fails to exhibit cross-sectional stability. Note how the time 2 Omnibus estimates of communities 1 and 2 remain separate, due to time 1 ‘interference’ affecting this joint embedding.

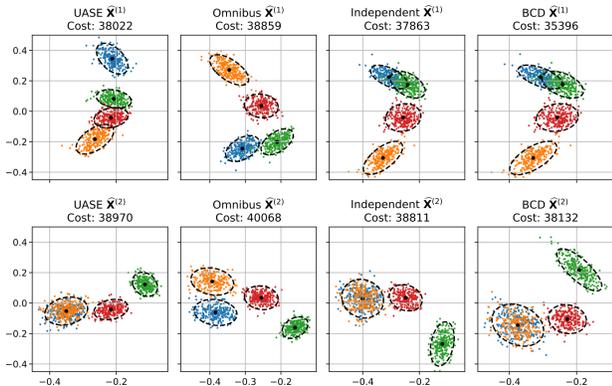


Fig. 5. Embeddings of two SBM graph realizations, where communities 1 and 2 merge, while community 4 keeps the connection probabilities with other groups. Observe how the BCD approach (far right) manages to capture this behaviour, while providing the best representation for each graph individually (quantified by the smallest cost function values). Example adapted from [9].

UASE and BCD produce embeddings that fulfill both stability requirements i) and ii). However, BCD yields a better overall representation for both graphs. This is quantified via the cost function (1) evaluated at each solution; see above each plot for the numerical values. Unlike the batch UASE, gradient descent methods as the ones we present here offer a pathway towards tracking nodal representations in a streaming graph setting – the subject dealt with next.

C. Model tracking for graph streams

Consider now a monitoring scenario, where we observe a stream of time-indexed graphs $\{\mathbf{A}_t\}$ and the goal is to track the underlying model. Different from the batch setting of the previous section, we are now unable to jointly process the entire graph sequence. This may be due to memory constraints or stringent delay requirements. We will still assume that nodes are identifiable across time, but the algorithm’s computational cost and memory footprint may not increase with t .

1) *Fixed vertex set*: We first consider the setting where N is fixed and we would like to track the latent vectors $\mathbf{X}_t \in \mathbb{R}^{N \times d}$.³ Previous efforts in this direction have been mainly motivated by the change-point detection problem; i.e., detecting if and when the generative model of the observed graph sequence changes [6], [35]–[37]. Our focus is on the related problem of estimating the embeddings’ evolution. A couple noteworthy applications include recommender systems (where rankings are revealed, or even change, over time) [38] or, as we discuss below, monitoring wireless networks [39].

Independent ASE computation for each \mathbf{A}_t suffers from the alignment issue already discussed. Instead, and supposing for now that \mathbf{M} can be ignored, it may well be the case that recursive methods to update the SVD of a perturbed matrix $\mathbf{A}_t = \mathbf{A}_{t-1} + \Delta_t$ suffice [18]. However, as we show in the following synthetic example, these approaches may also produce arbitrarily rotated estimates from one time-step to the next, and suffer from catastrophic error accumulation [21].

Tracking of a dynamic SBM. Our idea is instead to proceed as in Remark 2, and warm-restart the GD iterations with

³We stick to undirected graphs for ease of exposition, but extensions to digraphs are straightforward and presented in the numerical experiments.

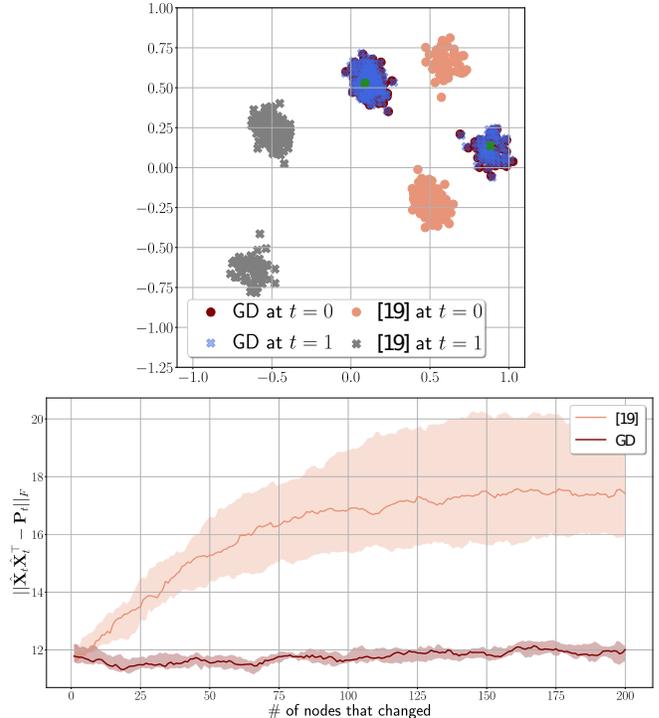


Fig. 6. Two-block dynamic SBM in which a single node changes affiliation at each t . Comparison between online GD and recursive SVD [18]. (top) Embeddings for the first two time-steps ($d = 2$); the node that changed communities is highlighted in green. Best viewed in a color display. Note how the change of a single node produces markedly different results for [18], whereas online GD offers stable estimates. (bottom) Evolution of $\|\hat{\mathbf{X}}_t \hat{\mathbf{X}}_t^T - \mathbf{P}_t\|_F$. Solid line indicates median across ten realizations, with the range between first and third quartiles shown in a lighter color. Online GD exhibits uniformly bounded error, whereas [18] accumulates error as t grows.

the previous time-step’s estimate $\hat{\mathbf{X}}_{t-1}$ (analogously to the example in Fig. 5). Consider a dynamic SBM graph with $N = 200$ nodes and two communities. At each time-step $t = 0, 1, 2, \dots$ a single randomly chosen node changes its community affiliation. We compare the tracking performance of warm-restarted GD [i.e., several iterations of GD in (3) initialized with the previous time-step’s estimate] and the fast, recursive SVD algorithm in [18]. The nodal embeddings for $t = 0$ and 1 (i.e., a single node changed affiliation) are depicted in Fig. 6 (top). Notice how online GD produces stable results, with a single vector moving from one cluster to the other. The rest of the nodes’ embeddings remain virtually unchanged. On the other hand, the recursive SVD in [18] fails to preserve a common angular reference for $\hat{\mathbf{X}}_0$ and $\hat{\mathbf{X}}_1$. Another well-documented drawback of these incremental SVD methods is that, since they update only the d most significant components, the error $\|\hat{\mathbf{X}}_t \hat{\mathbf{X}}_t^T - \mathbf{P}_t\|_F$ increases with t [21]. Fig. 6 (bottom) illustrates this error-accumulation behavior, to be contrasted with online GD that keeps the error in check for all $t \geq 0$.

Wireless network monitoring. We may further leverage the fact that \mathbf{X}_t is typically correlated with \mathbf{X}_{t-1} in order to improve the embeddings’ accuracy over time. For example, suppose $\mathbf{X}_{t-m} = \dots = \mathbf{X}_t = \mathbf{X}$ over some interval of length m . It is then prudent to estimate \mathbf{X} by solving (1), but using the average $\bar{\mathbf{A}}_t = 1/m \sum_{k=t-m}^t \mathbf{A}_k$ instead [40]. Note that $\bar{\mathbf{A}}_t$ may be interpreted as the adjacency matrix of a weighted

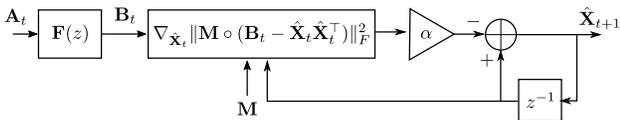


Fig. 7. A diagram of the proposed tracking system. The entry-wise filter $\mathbf{F}(z)$ implements an averaging operator, e.g., a fixed-length moving average.

graph. Edge weights can also be modeled by an RDPG, where now the embeddings are such that $\mathbf{X}\mathbf{X}^T = \mathbb{E}[\mathbf{A}]$. Unlike the unweighted case, $\mathbb{E}[\mathbf{A}]$ are not probabilities. Still, under mild assumptions on the weights' distribution, the solution of (1) for weighted \mathbf{A} is a consistent estimator of \mathbf{X} as $N \rightarrow \infty$ [6].

These observations motivate well the two-stage tracking system depicted in Fig. 7. The stream of adjacency matrices $\{\mathbf{A}_t\}$ is fed to the entry-wise filter $\mathbf{F}(z)$, which outputs $\{\mathbf{B}_t\}$. For instance, $\mathbf{F}(z)$ may be a moving average of fixed length m as before. If memory is at a premium, we may use a single-pole IIR filter instead so that $\{\mathbf{B}_t\}$ is an exponentially-weighted moving average of the input adjacency matrices. We may even drop the filtering stage altogether (setting $m = 1$) to yield a least mean squares (LMS)-type online GD algorithm.

We now empirically demonstrate this simple tracking system yields accurate and stable embeddings of dynamic RDPG graphs. Consider a Wi-Fi network from which a monitoring system periodically acquires the Received Signal Strength Indicator (RSSI) between Access Points (APs) – a feature typically available in enterprise-level deployments. We will use our GRL framework to flag network changes and eventually diagnose them. We analyze graphs \mathbf{A}_t whose nodes are the APs and the edge weights are the measured RSSI values (plus a constant offset so that all values are positive). Since these measurements are typically not symmetric, we have a digraph sequence. We rely on the dataset described in [41], which consists of hourly measurements between $N = 6$ APs at a Uruguayan school, over almost four weeks ($m = 655$ graphs). During the monitoring period, the network administrator moved an AP ($i = 4$) at $t \approx 310$.

To track the AP embeddings, we run an online version of Algorithm 2 as schematically shown in the diagram of Fig. 7, but adapted to digraphs. This entails a retraction after the Riemannian GD step, not shown in the diagram. We use an IIR filter $\mathbf{F}(z)$ with a pole at 0.9. Furthermore, we adopt a fix stepsize $\alpha = 0.01$ instead of choosing it via the Armijo rule.

The evolution of the online Riemannian GD estimates $\hat{\mathbf{X}}_t^l$ and $\hat{\mathbf{X}}_t^r$ for $d = 2$ is shown in Fig. 8. Different color palettes are used to distinguish the nodes, and as t increases the colors become lighter. Note how at all times there are two (almost) orthogonal cluster of nodes: c1) APs 1 and 2 (in the lower part of the plots); and c2) APs 3, 5 and (to a lesser extent) 6. AP 4 is embedded between both communities for all t . Moreover, note how the *trajectory of each AP* can be split into a couple clear states, discernable as the colors transition from darker to lighter. This is indicative of the change in AP 4's position at roughly the middle of the monitoring period. Finally, movement within both AP clusters is mostly radial, hence dot products between cluster members are preserved. On the other hand, AP 4 moves transversally closer to c2, consistent with the information provided by the network administrator. It appears as if it was moved closer to AP 5, and AP 1 remains

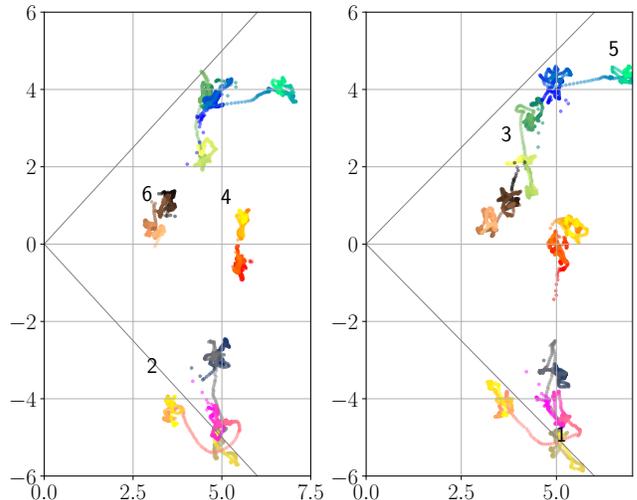


Fig. 8. Embeddings $\hat{\mathbf{X}}_t^l$ (left) and $\hat{\mathbf{X}}_t^r$ (right) for the RSSI digraph ($d = 2$). Color palettes distinguish the APs and a lighter tone indicates larger values of t . Best viewed in a color display. The network's change at $t \approx 310$ is apparent. AP 4 was moved ($i = 4$) closer to the upper cluster of APs.

its closest member from c1.

2) *Time-varying node set*: In dynamic environments it is not uncommon for nodes to join or leave the network. Going back to the wireless network test case, the question remains on how to proceed should an AP fail, or if the administrator decides to add a new one to improve coverage. Dealing with the former case is straightforward; if a node leaves the network at time t , we simply drop the corresponding row in $\hat{\mathbf{X}}_{t-1}$ and re-run the GD algorithm (warm-restarted from there).

Node additions require more thought. Suppose that a single node $i = N + 1$ joins the network at time t . Let $\mathbf{a}_{N+1} = [A_{1,N+1}, \dots, A_{N,N+1}]^T \in \{0, 1\}^N$ be the $(N + 1)$ -th column of $\mathbf{A}_t \in \{0, 1\}^{N+1 \times N+1}$, excluding $A_{N+1,N+1} = 0$ and dropping the subindex t for notational convenience. Then given $\hat{\mathbf{X}}_{t-1} \in \mathbb{R}^{N \times d}$, we can embed node i by solving

$$\hat{\mathbf{x}}_{N+1} = \operatorname{argmin}_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{a}_{N+1} - \hat{\mathbf{X}}_{t-1} \boldsymbol{\theta}\|_2^2. \quad (12)$$

This simple but intuitive out-of-sample embedding procedure was studied in [25], and shown to recover the true latent positions as $N \rightarrow \infty$. If several nodes are added at a given time-step, they can all be embedded by solving multiple LS problems like (12). However, this procedure disregards the information from the connections between new nodes. Furthermore, if the embeddings of existing nodes are not updated, their growing inaccuracies as \mathbf{A}_t evolves will negatively impact future nodes' representations.

As we show in the following numerical experiments, these drawbacks can be overcome by running our online GD-based algorithms to update *all embeddings* $\hat{\mathbf{X}}_t$, initializing existing nodes with $\hat{\mathbf{X}}_{t-1}$ and new one(s) with $\hat{\mathbf{x}}_{N+1}$ as in (12).

Dynamic random graph with growing vertex set. Consider an Erdős-Rényi graph with a fixed connection probability $p = 0.1$, and initial number of nodes $N_0 = 100$. At each time-step t we add a single node so that $N_t = N_{t-1} + 1$. The evolution of the error $\|\hat{\mathbf{X}}_t \hat{\mathbf{X}}_t^T - \mathbf{P}_t\|_F / \sqrt{N_t}$ is shown in Fig. 9. Note how

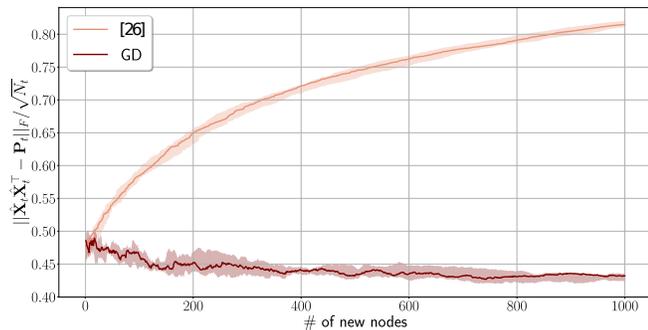


Fig. 9. Dynamic Erdős-Rényi graph in which a single node is added at each t . Comparison between online GD and out-of-sample LS embedding [25]. Evolution of $\|\hat{\mathbf{X}}_t \hat{\mathbf{X}}_t^T - \mathbf{P}\|_F / \sqrt{N_t}$. Solid line indicates median across ten realizations, with range between first and third quartiles shown in a lighter color. Once more, online GD exhibits uniformly bounded error, whereas the baseline method [25] accumulates error as t grows.

(carefully warm-restarted) online GD exhibits bounded error behavior, in stark contrast with repeated LS-based embeddings as in [25]. Admittedly, this gain in accuracy comes with a modest increase in computation (few GD steps), and identical memory footprint (i.e., storing the current embeddings and the new adjacency matrix) as the baseline in [25].

Tracking international relations from UN voting data. Here we revisit the UN General Assembly voting data from Section V-A. Following the same bipartite digraph construction procedure, we study all yearly graphs from 1955 to 2015. In this dynamic network we have a time-varying node set. Roll calls change from one year to the next, and also several countries joined the UN later (while others have ceased to exist). We embed the first graph from 1955 using Riemannian GD initialized at random (as before, using $d = 2$). For each successive year, we warm-restart Algorithm 2 with the embeddings from the previous year, while new nodes are initialized using the LS solution (12).

Fig. 10 depicts the embeddings of four countries: USA, Israel, Cuba, and the USSR (later, the Russian Federation). We use a similar visualization style as in Fig. 8, with different color palettes used to distinguish among countries, and lighter tones indicating more recent years. Observe how the representations for the USA and Israel remain strongly aligned over the entire time horizon, which is consistent with their longstanding agreement on UN resolution matters. The embedding for the USSR is initially (nearly) orthogonal to the USA and Israel, with Cuba initially showing a greater affinity to the USA/Israel block. This is consistent with Cold War geopolitics of the time. Then, after 1959, Cuba’s position shifts to the lower half-plane, becoming more aligned with the USSR. This is expected given Cuba’s sharp shift in foreign policy as a result of the Cuban revolution, with its ideology being in agreement with that of the USSR. This polarized scenario remained unchanged until 1991. That year the embedding for the USSR (now the Russian Federation) moves closer to the USA/Israel block, which reflects the politics of the Russian Federation in the aftermath of USSR’s dissolution. Cuba remains at an (almost) orthogonal position from the USA/Israel block, with Russia eventually shifting to a middle ground after the mid-2000’s.

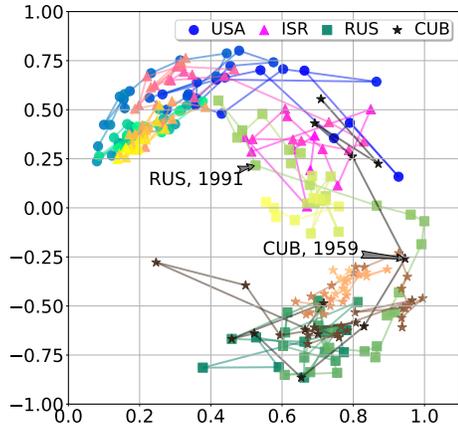


Fig. 10. UN General Assembly voting data from 1955 to 2015. Evolution of nodal positions for the USA, Israel, Cuba, and the USSR (or, after 1991, the Russian Federation) estimated via online Riemannian GD. Color palettes distinguish the countries and a lighter tone indicates later years. Best viewed in a color display. Note how the USA and Israel remain strongly aligned over the entire span, with Cuba and the USSR shifting alignments depending of their political views.

VI. CONCLUDING REMARKS

We developed a gradient-based spectral-embedding framework to estimate latent positions of RDPGs. Relative to prior art our algorithmic approaches offer better representation at a competitive computational cost, and they are more broadly applicable to settings with incomplete, dynamic, and directed network data. We motivated and proposed a novel manifold-constrained formulation to embed directed RDPGs, and developed novel Riemannian GD iterations to estimate interpretable latent nodal positions. The effectiveness of the GRL framework is demonstrated via reproducible experiments with both synthetic and real (wireless network and United Nations voting) data. We made all our codes publicly available.

This work and its current limitations open up several exciting and challenging directions for future research. For even better scalability, in the near future we plan to migrate our Python implementations to a faster language such as C. Exploring the viability of parallelizing the BCD iterations is another direction in our future research agenda. With regards to the streaming scenario, it would be of interest to develop lightweight online rules to adaptively determine the embedding dimension. Performing dynamic regret analyses of the online GD methods would be a valuable contribution, since such guarantees in non-convex settings are so far quite rare.

REFERENCES

- [1] M. Fiori, B. Marenco, F. Larroca, P. Bermolen, and G. Mateos, “Algorithmic advances for the adjacency spectral embedding,” in *Proc. of European Signal Process. Conf.*, August 2022.
- [2] F. Larroca, P. Bermolen, M. Fiori, B. Marenco, and G. Mateos, “Tracking the Adjacency Spectral Embedding for Streaming Graphs,” in *Proc. Asilomar Conf. on Signals, Systems, Computers*, 2022.
- [3] 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,” *J. Mach. Learn. Res.*, vol. 18, no. 1, p. 8393–8484, January 2017.
- [4] E. Scheinerman and K. Tucker, “Modeling graphs using dot product representations,” *Comput. Stat.*, vol. 25, pp. 1–16, 2010.

- [5] V. Lyzinski, M. Tang, A. Athreya, Y. Park, and C. E. Priebe, "Community detection and classification in hierarchical stochastic blockmodels," *IEEE Trans. Netw. Sci. Eng.*, vol. 4, no. 1, pp. 13–26, 2017.
- [6] B. Marenco, P. Bermolen, M. Fiori, F. Larroca, and G. Mateos, "Online change point detection for weighted and directed random dot product graphs," *IEEE Trans. Signal Inf. Process. Netw.*, vol. 8, pp. 144–159, 2022.
- [7] W. L. Hamilton, "Graph representation learning," *Synthesis Lectures on Artificial Intelligence and Machine Learning*, vol. 14, pp. 1–159, 2020.
- [8] J. Chung, B. D. Pedigo, E. W. Bridgeford, B. K. Varjavand, H. S. Helm, and J. T. Vogelstein, "GraSPy: Graph statistics in Python." *J. Mach. Learn. Res.*, vol. 20, no. 158, pp. 1–7, 2019.
- [9] I. Gallagher, A. Jones, and P. Rubin-Delanchy, "Spectral embedding for dynamic networks with stability guarantees," *Proc. Adv. Neural. Inf. Process. Syst.*, 2021.
- [10] Y. Chi, Y. Lu, and Y. Chen, "Nonconvex optimization meets low-rank matrix factorization: An overview," *IEEE Trans. Signal Process.*, vol. 67, no. 20, pp. 5239–5269, 2019.
- [11] T. Vu and R. Raich, "Exact linear convergence rate analysis for low-rank symmetric matrix completion via gradient descent," in *Proc. Int. Conf. Acoustics, Speech, Signal Process.*, 2021, pp. 3240–3244.
- [12] P.-A. Absil, R. Mahony, and R. Sepulchre, *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, 2009.
- [13] N. Boumal, *An Introduction to Optimization on Smooth Manifolds*. Cambridge University Press, 2023.
- [14] A. G. Marques, S. Segarra, and G. Mateos, "Signal processing on directed graphs," *IEEE Signal Process. Mag.*, vol. 37, no. 6, pp. 99–116, 2020.
- [15] Y. Koren, R. Bell, and C. Volinsky, "Matrix factorization techniques for recommender systems," *Computer*, vol. 42, no. 8, pp. 30–37, 2009.
- [16] I. Dokmanic, R. Parhizkar, J. Ranieri, and M. Vetterli, "Euclidean distance matrices: Essential theory, algorithms, and applications," *IEEE Signal Process. Mag.*, vol. 32, no. 6, pp. 12–30, 2015.
- [17] M. A. Davenport and J. Romberg, "An overview of low-rank matrix recovery from incomplete observations," *IEEE J. Sel. Topics Signal Process.*, vol. 10, no. 4, pp. 608–622, 2016.
- [18] M. Brand, "Fast low-rank modifications of the thin singular value decomposition," *Linear Algebra Appl.*, vol. 415, no. 1, pp. 20–30, 2006, special Issue on Large Scale Linear and Nonlinear Eigenvalue Problems.
- [19] K. Levin, A. Athreya, M. Tang, V. Lyzinski, and C. E. Priebe, "A central limit theorem for an omnibus embedding of multiple random dot product graphs," in *Int. Conf. on Data Mining Workshops*, 2017, pp. 964–967.
- [20] A. Jones and P. Rubin-Delanchy, "The multilayer random dot product graph," *arXiv:2007.10455 [stat.ML]*, 2020.
- [21] Z. Zhang, P. Cui, J. Pei, X. Wang, and W. Zhu, "TIMERS: Error-bounded SVD restart on dynamic networks," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, Apr. 2018.
- [22] N. Halko, P. G. Martinsson, and J. A. Tropp, "Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions," *SIAM Rev.*, vol. 53, no. 2, pp. 217–288, 2011.
- [23] I. Chami, S. Abu-El-Haija, B. Perozzi, C. Ré, and K. Murphy, "Machine learning on graphs: A model and comprehensive taxonomy," *J. Mach. Learn. Res.*, vol. 23, no. 89, pp. 1–64, 2022.
- [24] V. Kalantzis and P. Traganitis, "Matrix resolvent eigenembeddings for dynamic graphs," in *Proc. Int. Conf. Acoustics, Speech, Signal Process.*, 2023.
- [25] K. Levin, F. Roosta, M. Mahoney, and C. Priebe, "Out-of-sample extension of graph adjacency spectral embedding," in *Proc. Int. Conf. Mach. Learn.*, vol. 80, 2018, pp. 2975–2984.
- [26] S. Bhojanapalli, A. Kyriillidis, and S. Sanghavi, "Dropping convexity for faster semi-definite optimization," in *Proc. Conf. Learn. Theory*, 2016, pp. 530–582.
- [27] Y. Chen and M. Wainwright, "Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees," *arXiv:1509.03025 [math.ST]*, 2015.
- [28] R. Sun and Z. Luo, "Guaranteed matrix completion via non-convex factorization," *IEEE Trans. Inf. Theory*, vol. 62, pp. 6535–6579, 2016.
- [29] D. Zhou, Y. Cao, and Q. Gu, "Accelerated factored gradient descent for low-rank matrix factorization," in *Proc. Int. Conf. Artif. Intell. Statist.*, 2020, pp. 4430–4440.
- [30] D. Bertsekas, *Nonlinear Programming*. Athena Scientific, 1999.
- [31] N. Boumal, P.-A. Absil, and C. Cartis, "Global rates of convergence for nonconvex optimization on manifolds," *IMA Journal of Numerical Analysis*, vol. 39, no. 1, pp. 1–33, 02 2018.
- [32] J. Townsend, N. Koep, and S. Weichwald, "Pymanopt: A Python toolbox for optimization on manifolds using automatic differentiation," *J. Mach. Learn. Res.*, vol. 17, no. 137, pp. 1–5, 2016.
- [33] E. Voeten, A. Strehznev, and M. Bailey, "United Nations General Assembly Voting Data," 2009. [Online]. Available: <https://doi.org/10.7910/DVN/LEJUOZ>
- [34] M. Tang, A. Athreya, D. Sussman, V. Lyzinski, Y. Park, and C. Priebe, "A semiparametric two-sample hypothesis testing problem for random graphs," *J. Comput. Graph. Stat.*, vol. 26, no. 2, 2017.
- [35] Y. Yu, O. H. M. Padilla, D. Wang, and A. Rinaldo, "Optimal network online change point localisation," *arXiv:2101.05477 [math.ST]*, 2021.
- [36] H. Chen, "Sequential change-point detection based on nearest neighbors," *Ann. Stat.*, vol. 47, no. 3, pp. 1381 – 1407, 2019.
- [37] M. Zhang, L. Xie, and Y. Xie, "Online community detection by spectral CUSUM," in *Proc. Int. Conf. Acoustics, Speech, Signal Process.*, 2020.
- [38] P. Campos, F. Díez, and I. Cantador, "Time-aware recommender systems: A comprehensive survey and analysis of existing evaluation protocols," *User Model. User-adapt. Interact.*, vol. 24, pp. 67–119, 2014.
- [39] G. Mateos and K. Rajawat, "Dynamic network cartography: Advances in network health monitoring," *IEEE Signal Process. Mag.*, vol. 30, no. 3, pp. 129–143, 2013.
- [40] R. Tang, M. Tang, J. T. Vogelstein, and C. E. Priebe, "Robust estimation from multiple graphs under gross error contamination," *arXiv:1707.03487 [stat.ME]*, 2017.
- [41] G. Capdehourat, F. Larroca, and G. Morales, "A nation-wide Wi-Fi RSSI dataset: Statistical analysis and resulting insights," in *Proc. IFIP Networking Conf.*, 2020, pp. 370–378.

APPENDIX

A. Proof of Proposition 3

In order to show that \mathcal{M} is a manifold and to further understand its differential structure, consider the function $F : \mathbb{R}^{N \times d} \mapsto \mathcal{S}_d$ defined as $F(\mathbf{X}) = \mathbf{M} \circ (\mathbf{X}^\top \mathbf{X} - \mathbf{I}_d)$. Observe that \mathcal{M} is defined as the preimage of zero through F , so we will prove that this is a regular value.

For $\zeta \in \mathbb{R}^{N \times d}$, the derivative of F in $\mathbf{X} \in F^{-1}(\mathbf{0}_{d \times d})$ along ζ is $DF(\mathbf{X})\zeta = \mathbf{M} \circ (\zeta^\top \mathbf{X} + \mathbf{X}^\top \zeta)$. Next we establish that $DF(\mathbf{X})$ is onto. Indeed, let η be a matrix in the orthogonal complement of the image, i.e., $\eta \in \text{Im}DF(\mathbf{X})^\perp \subset \mathcal{S}_d$. Then

$$\langle \eta, \mathbf{M} \circ (\zeta^\top \mathbf{X} + \mathbf{X}^\top \zeta) \rangle = 0, \forall \zeta \in \mathbb{R}^{N \times d}.$$

Now, since the diagonal of η is null, we may drop the Hadamard product with \mathbf{M} and obtain

$$\langle \eta, \mathbf{M} \circ (\zeta^\top \mathbf{X} + \mathbf{X}^\top \zeta) \rangle = \langle \eta, \zeta^\top \mathbf{X} + \mathbf{X}^\top \zeta \rangle = 0, \forall \zeta \in \mathbb{R}^{N \times d}.$$

So we have $\text{tr}(\eta \zeta^\top \mathbf{X}) + \text{tr}(\eta \mathbf{X}^\top \zeta) = 0$ and these two summands are equal to each other by virtue of the circular property of the trace operator. Hence, we obtain $2\text{tr}(\eta \mathbf{X}^\top \zeta) = 0, \forall \zeta \in \mathbb{R}^{N \times d}$, and since this trace vanishes for all ζ , we have that $\eta \mathbf{X}^\top = \mathbf{0}_{d \times N}$. Multiplying by \mathbf{X} we obtain $\eta \mathbf{X}^\top \mathbf{X} = \mathbf{0}_{d \times d}$. Because $\mathbf{X}^\top \mathbf{X}$ is diagonal, necessarily $\eta = \mathbf{0}_{d \times d}$ and therefore $DF(\mathbf{X})$ is onto. The conclusion is that \mathcal{M} is a differential manifold, of dimension $Nd - \frac{d(d-1)}{2}$.

The tangent space at \mathbf{X} can be obtained as the kernel of $DF(\mathbf{X})$, so we have

$$T_{\mathbf{X}}\mathcal{M} = \{\zeta \in \mathbb{R}^{N \times d} : \mathbf{M} \circ (\zeta^\top \mathbf{X} + \mathbf{X}^\top \zeta) = \mathbf{0}_{d \times d}\}, \quad (13)$$

completing the proof. \blacksquare

B. Proof of Lemma 1

Consider a matrix of the form $\mathbf{X}\Lambda$ with $\Lambda \in \mathcal{S}_d$, and let us show that it is orthogonal to a matrix of the tangent space. Now, observe that $\text{tr}((\mathbf{X}\Lambda)^\top \zeta) = \text{tr}(\Lambda \zeta^\top \mathbf{X})$. Therefore,

$$\text{tr}((\mathbf{X}\Lambda)^\top \zeta) = \frac{1}{2} \text{tr}(\Lambda (\mathbf{X}^\top \zeta + \zeta^\top \mathbf{X})) = 0.$$

The last trace is zero since $\Lambda \in \mathcal{S}_d$ and $\mathbf{X}^\top \zeta + \zeta^\top \mathbf{X}$ is diagonal, because $\zeta \in T_{\mathbf{X}}\mathcal{M}$. \blacksquare

As expected, the dimension of the normal space is $\frac{N(N-1)}{2}$, which is the dimension of \mathcal{S}_d .

C. Proof of Lemma 2 and Proposition 4

To compute the projection to the normal space, recall some auxiliary matrices defined in Lemma 2. Let $\mathbf{X} \in \mathcal{M}$. Then $\mathbf{X}^\top \mathbf{X}$ is diagonal, with positive entries. Define $\mathbf{D} = (\mathbf{X}^\top \mathbf{X})^{1/2}$ and let $\mathbf{E} = \mathbf{1}_d \mathbf{1}_d^\top \mathbf{D}^2 + \mathbf{D}^2 \mathbf{1}_d \mathbf{1}_d^\top$. These matrices allow us to re-write the operation $\varphi(\mathbf{A}) = \mathbf{A} \mathbf{D}^2 + \mathbf{D}^2 \mathbf{A}$ as $\varphi(\mathbf{A}) = \mathbf{A} \circ \mathbf{E}$. In particular, this allows us to obtain an expression for the inverse operation, which will be needed. Indeed, if \mathbf{F} is the matrix with entries $F_{ij} = E_{ij}^{-1}$, then $(\mathbf{A} \mathbf{D}^2 + \mathbf{D}^2 \mathbf{A}) \circ \mathbf{F} = \mathbf{A}$, for all $\mathbf{A} \in \mathbb{R}^{d \times d}$. We can now prove the expression of the projection as follows.

From the characterization of Lemma 1, it is clear that $\mathbf{X}_s(2\mathbf{D}\mathbf{L}) \in T_{\mathbf{X}}\mathcal{M}^\perp$, since $s(2\mathbf{D}\mathbf{L}) \in \mathcal{S}_d$. Let us see that $\mathbf{Z} - \mathbf{X}_s(2\mathbf{D}\mathbf{L}) \in T_{\mathbf{X}}\mathcal{M}$. From (13), we have to show that

$$(\mathbf{Z} - \mathbf{X}_s(2\mathbf{D}\mathbf{L}))^\top \mathbf{X} + \mathbf{X}^\top (\mathbf{Z} - \mathbf{X}_s(2\mathbf{D}\mathbf{L})) \text{ is diagonal.}$$

Indeed,

$$\begin{aligned} & (\mathbf{Z} - \mathbf{X}_s(2\mathbf{D}\mathbf{L}))^\top \mathbf{X} + \mathbf{X}^\top (\mathbf{Z} - \mathbf{X}_s(2\mathbf{D}\mathbf{L})) = \\ & \mathbf{Z}^\top \mathbf{X} - (s(2\mathbf{D}\mathbf{L}))^\top \mathbf{X}^\top \mathbf{X} + \mathbf{X}^\top \mathbf{Z} - \mathbf{X}^\top \mathbf{X}_s(2\mathbf{D}\mathbf{L}) = \\ & \mathbf{Z}^\top \mathbf{X} + \mathbf{X}^\top \mathbf{Z} - \frac{1}{2} [2\mathbf{D}\mathbf{L} + 2(\mathbf{D}\mathbf{L})^\top] \mathbf{X}^\top \mathbf{X} + \text{diag}(2\mathbf{D}\mathbf{L}) \mathbf{X}^\top \mathbf{X} - \\ & \frac{1}{2} \mathbf{X}^\top \mathbf{X} [2\mathbf{D}\mathbf{L} + 2(\mathbf{D}\mathbf{L})^\top] + \mathbf{X}^\top \mathbf{X} \text{diag}(2\mathbf{D}\mathbf{L}). \end{aligned}$$

Now, since $\text{diag}(2\mathbf{D}\mathbf{L})$ and $\mathbf{X}^\top \mathbf{X}$ are diagonal matrices, they commute, and their product is diagonal. So we can forget those two terms in the expression, and continue with the rest. We will use the expression of \mathbf{L} and the fact that $\mathbf{X}^\top \mathbf{X} = \mathbf{D}^2$. Hence,

$$\begin{aligned} & \mathbf{Z}^\top \mathbf{X} + \mathbf{X}^\top \mathbf{Z} - \frac{1}{2} [2\mathbf{D}\mathbf{L} + 2(\mathbf{D}\mathbf{L})^\top] \mathbf{X}^\top \mathbf{X} - \\ & \frac{1}{2} \mathbf{X}^\top \mathbf{X} [2\mathbf{D}\mathbf{L} + 2(\mathbf{D}\mathbf{L})^\top] = \\ & \mathbf{Z}^\top \mathbf{X} + \mathbf{X}^\top \mathbf{Z} - [\mathbf{D}\mathbf{L}\mathbf{D}^2 + \mathbf{L}^\top \mathbf{D}\mathbf{D}^2 + \mathbf{D}^2 \mathbf{D}\mathbf{L} + \mathbf{D}^2 \mathbf{L}^\top \mathbf{D}] = \\ & \mathbf{Z}^\top \mathbf{X} + \mathbf{X}^\top \mathbf{Z} - [\mathbf{D}(\mathbf{L}\mathbf{D}^2 + \mathbf{D}^2 \mathbf{L}) + (\mathbf{L}^\top \mathbf{D}^2 + \mathbf{D}^2 \mathbf{L}^\top) \mathbf{D}] = \\ & \mathbf{Z}^\top \mathbf{X} + \mathbf{X}^\top \mathbf{Z} - [\mathbf{D}(\mathbf{L} \circ \mathbf{E}) + (\mathbf{L}^\top \circ \mathbf{E}) \mathbf{D}]. \end{aligned}$$

Now, observe that $\mathbf{L} \circ \mathbf{E} = ((\mathbf{D}^{-1} \mathbf{X}^\top \mathbf{Z}) \circ \mathbf{F}) \circ \mathbf{E} = \mathbf{D}^{-1} \mathbf{X}^\top \mathbf{Z}$, and the same happens with \mathbf{L}^\top . We end up with

$$\mathbf{Z}^\top \mathbf{X} + \mathbf{X}^\top \mathbf{Z} - [\mathbf{D}(\mathbf{D}^{-1} \mathbf{X}^\top \mathbf{Z}) + (\mathbf{Z}^\top \mathbf{X} \mathbf{D}^{-1}) \mathbf{D}] = \mathbf{0}_{d \times d},$$

which in particular is diagonal. Therefore, $\mathbf{Z} - \mathbf{X}_s(2\mathbf{D}\mathbf{L}) \in T_{\mathbf{X}}\mathcal{M}$ and this completes the proof. \blacksquare

The proof of Proposition 4 is straightforward now. We have all we need to compute the projection to the tangent space $\pi_{\mathbf{X}} : \mathbb{R}^{N \times d} \mapsto T_{\mathbf{X}}\mathcal{M}$, since $\pi_{\mathbf{X}}(\mathbf{Z}) + \pi_{\mathbf{X}}^\perp(\mathbf{Z}) = \mathbf{Z}$. \blacksquare

D. Proof of Proposition 5

Denoting by $\mathbb{R}_{fr}^{N \times d}$ the set of $N \times d$ full-rank matrices, and by $\text{Supp}_1(d)$ the set of upper triangular matrices with ones in the diagonal, let us consider the mapping

$$\phi : \mathcal{M} \times \text{Supp}_1(d) \mapsto \mathbb{R}_{fr}^{N \times d}, \text{ with } \phi(\tilde{\mathbf{Q}}, \tilde{\mathbf{R}}) = \tilde{\mathbf{Q}}\tilde{\mathbf{R}}.$$

From the discussion immediately preceding the statement of Proposition 5, we have that ϕ is bijective. Furthermore, ϕ is smooth since its the restriction of the matrix multiplication to a submanifold. Now, given a full rank matrix \mathbf{M} , the first component of ϕ^{-1} can be obtained as the result of a modified

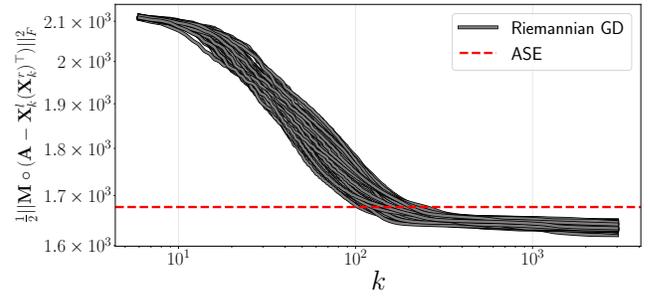


Fig. 11. The evolution of $f(\mathbf{X}_k^l, \mathbf{X}_k^r) = \frac{1}{2} \|\mathbf{M} \circ (\mathbf{A} - \mathbf{X}_k^l(\mathbf{X}_k^r)^\top)\|_F^2$ using Algorithm 2 to embed an LFR graph, starting from 75 different random initializations. The first 5 iterations are omitted for clarity. Note how the algorithm systematically produces estimates of the embeddings with a lower cost than ASE, and marginal variability regardless of the initialization.

Gram-Schmidt process, which is differentiable. The second component can then be obtained as $\tilde{\mathbf{R}} = \tilde{\mathbf{Q}}^{-1} \mathbf{M}$, and therefore it is also differentiable. It follows that ϕ is a diffeomorphism.

We also have that $\phi(\tilde{\mathbf{Q}}, \mathbf{I}_d) = \tilde{\mathbf{Q}}$. Following [12, Prop. 4.1.2] we have that the projection onto the first component of ϕ^{-1} is a retraction, which is exactly the \tilde{qf} mapping defined in Proposition 5. \blacksquare

E. Robustness to initialization

Since the objective functions we optimize are non-convex and convergence guarantees provided are to stationary solutions, it is prudent to study the algorithms' sensitivity to the initialization. As discussed in Section III-A, except for specific problematic initializations corresponding to sub-optimal stationary points, in our experience all algorithms converge to the optimum when they are initialized at random. The following experiment illustrates this desirable property, in particular for the Riemannian GD (i.e., Algorithm 2) method developed to embed digraphs. Similar results are obtained for the other algorithms, but not included here to avoid repetition.

We consider a Lancichinetti–Fortunato–Radicchi (LFR) [?] benchmark graph with $N = 1000$ nodes, randomly initialize Algorithm 2 and plot the evolution of the cost function $f(\mathbf{X}_k^l, \mathbf{X}_k^r)$ in (2). The LFR model is a widely adopted benchmark that produces graphs with properties observed in real-world networks, particularly in terms of the resulting degree distribution and community sizes. Here, the LFR graph was generated using parameters $\tau_1 = 3$ and $\tau_2 = 2$ as exponents of the power law distributions for the degree and community size, respectively, and mixing parameter $\mu = 0.1$. The resulting graph has 16 communities, the larger one with 142 nodes, and the smaller one with 30 nodes. The largest hub has 157 neighbors, and there are several nodes with degree 2. Fig. 11 shows the results over 75 randomly initialized runs where $d = 16$, and also the cost function value 1676.49 obtained by ASE. Regardless of the initialization, the limiting objective values obtained via Riemannian GD exhibit marginal variability (mean = 1635.66, std = 6.52) and always outperform ASE. In terms of timing, embedding each of these LFR graphs with $N = 1000$ nodes takes 40 seconds on average.