

# Graph Learning with Low-rank and Diagonal Structures: A Riemannian Geometric Approach

Xiang Zhang

*School of Information Science and Engineering  
Southeast University  
Nanjing, China  
xiangzhang369@seu.edu.cn*

Qiao Wang

*School of Information Science and Engineering  
Southeast University  
Nanjing, China  
qiaowang@seu.edu.cn*

**Abstract**—We study the problem of learning graphs in Gaussian graphic models by assuming that the underlying precision matrix has a “low-rank and diagonal” (LRaD) structure. This assumption enjoys a latent factor representation interpretation and can reduce the dimensionality of the problem, thus leading to superior performance with fewer samples. To address the optimization challenges caused by the LRaD constraint, we design a Riemannian quotient manifold to transform the proposed model into an unconstrained problem on the manifold. Then, we devise a Riemannian conjugate gradient algorithm to solve the proposed model. Experimental results on synthetic and real data illustrate the effectiveness of the proposed method.

**Index Terms**—Graph learning, Gaussian graphic model, Gaussian precision factor model, Riemannian manifold optimization

## I. INTRODUCTION

Graphs can flexibly represent irregular structures behind high dimensional data and have proven their power in many applications [1]. Inferring graph topology behind data, also known as graph learning (GL), has been increasingly studied and applied in numerous fields, such as social networks [2] and medical data [3]. Extensive studies are conducted on GL [4], [5], including models based on causality dependencies [6], [7] and graph signal processing [2], [8].

Here, we focus on GL in Gaussian graphical models (GGMs), where graphs encode conditional dependence between two variables in a multivariate Gaussian model. In GGMs, GL is equivalent to estimating precision (inverse covariance) matrices, which is challenging, especially for small sample sizes and high-dimensional data [9]. Thus, many works impose prior structures on the graphs to improve performance, such as sparsity [10], [11], separated clusters [12], non-negative correlations [13], [14], and tree structures [15].

In this study, we assume that the precision matrix enjoys a “low-rank and diagonal” (LRaD) structure. This structure is inspired by the Gaussian precision factor model [16], which provides a latent factor representation interpretation for data with an LRaD precision matrix. The LRaD assumption can reduce the dimensionality of the problem, thereby improving

the estimation performance under small sample sizes. In the literature, the LRaD structure is widely used to estimate covariance matrices, which is derived from the Gaussian factor analysis (GFA) model [17]–[19]. Our study differs from these models in that we aim to estimate precision matrices, on which the LRaD structure is imposed. The work [16] also uses the LRaD structure to construct a precision matrix estimator from a Bayesian perspective. Our model learns graphs by solving a penalized Gaussian maximum likelihood estimation problem.

However, the LRaD structure poses optimization challenges due to the non-convex low-rank constraint. Existing methods include relaxing the low-rank constraint by penalizing nuclear norm [20] so that the problem can be solved using convex methods [21]. Albeit feasible, the solutions of these methods are approximations of original problems and require additional tuning parameters. To this end, we propose an algorithm based on manifold optimization [22], [23]. Specifically, we design a Riemannian quotient manifold [24], [25] based on the LRaD structure. We then propose a Riemannian conjugate gradient (RCG) algorithm to solve our model on this manifold. The merits of our algorithm are two-fold. First, optimizing on the manifold transforms our model into an unconstrained problem, and the solution will exactly and automatically satisfy the LRaD constraint. Second, the solution is not an approximation of the original problem and is free of extra parameter tuning.

The main contributions of this study are summarized as follows. (i) We propose a GL model in GGMs by assuming the LRaD structure on precision matrices, which reduces the dimensionality of the problem and thus improves estimation performance under small sample sizes. (ii) We design a Riemannian manifold to describe the LRaD constraint, on which an RCG algorithm is proposed to solve our problem. (iii) Extensive experiments, including synthetic and real data, are conducted to illustrate the effectiveness of our approach.

## II. PROBLEM FORMULATION

In this section, we first introduce the background of GL in GGMs and then present the proposed model.

### A. Learning Graphs in GGMs

Let  $\mathcal{G} = \{\mathcal{A}, \mathcal{E}\}$  be an undirected graph with the set of  $p$  nodes  $\mathcal{A}$  and the set of edges  $\mathcal{E}$ . Assume that a random vector  $\mathbf{x} \in \mathbb{R}^p$  follows the Gaussian distribution  $\mathcal{N}(\mathbf{0}, \Sigma)$ , and let

This research is partially supported by Project Four of the National Natural Science Foundation of China (NSFC) Major Program, Prototype Construction of Carbon Neutrality Building Technology for High Density Urban Environment, Project Number: 52394224.

$\Theta = \Sigma^{-1}$  be the precision matrix. Associating  $\mathbf{x}$  with the graph  $\mathcal{G}$  forms a GGM satisfying

$$\Theta[ij] \neq 0 \Leftrightarrow (i, j) \in \mathcal{E}, \forall i \neq j \quad (1)$$

$$\Theta[ij] = 0 \Leftrightarrow \mathbf{x}[i] \perp \mathbf{x}[j] | \mathbf{x}[\mathcal{A} \setminus \{i, j\}], \quad (2)$$

where  $\Theta[ij]$  is the  $(i, j)$  entry of  $\Theta$ ,  $\mathbf{x}[i] \perp \mathbf{x}[j] | \mathbf{x}[\mathcal{A} \setminus \{i, j\}]$  means that  $\mathbf{x}[i]$  is conditionally independent of  $\mathbf{x}[j]$  given the other variables. Thus, the graph  $\mathcal{G}$  or  $\Theta$  characterizes conditional dependence between two variables in a multivariate vector  $\mathbf{x}$ . Within GGMs, GL aims to infer the precision matrix  $\Theta$  from  $N$  independently observed data  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{p \times N}$ , which boils down to solving the following penalized Gaussian maximum likelihood estimation problem:

$$\min_{\Theta \in \mathcal{S}_{++}^p} f(\Theta) = \min_{\Theta \in \mathcal{S}_{++}^p} \text{tr}(\Theta \mathbf{S}) - \log \det(\Theta) + \beta h(\Theta), \quad (3)$$

where  $\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top$  is the sample covariance matrix,  $\text{tr}(\cdot)$  is the trace operator,  $\det(\cdot)$  is the determinant of a matrix, and  $\mathcal{S}_{++}^p$  is the set of  $p \times p$  positive definite (PD) matrices. Moreover,  $h(\Theta) = \sum_{i \neq j} \psi(\Theta[ij])$ ,  $\psi(\cdot)$  is a sparsity penalty function applied to the off-diagonal elements of  $\Theta$ , and  $\beta \geq 0$  is a constant controlling sparsity.

### B. Learning Graphs with LRA D Structures

The model (3) imposes sparsity via the regularizer  $h(\Theta)$ , which is not sufficient and may lead to poor estimation results when  $N$  is small [12]. In this study, we assume that  $\Theta$  enjoys the LRA D structure. Specifically, for any  $\Theta \in \mathcal{S}_{++}^p$ , we have

$$\Theta = \mathbf{L} + \mathbf{D}, \quad \mathbf{D} \in \mathcal{D}_{++}^p, \quad \mathbf{L} \in \mathcal{S}_{+}^{p,r}, \quad (4)$$

where  $\mathcal{D}_{++}^p$  is the set of  $p \times p$  diagonal PD matrices, and  $\mathcal{S}_{+}^{p,r}$  is the set of all  $p \times p$  positive semi-definite (PSD) matrices with rank  $r$ ,  $0 \leq r \leq p$ . Note that the decomposition (4) always holds for a sufficiently large  $r$ . However, in practice, we empirically find that values of  $r \ll p$  suffice to approximate  $\Theta$  well even if  $\Theta$  does not exactly admit such an LRA D structure.

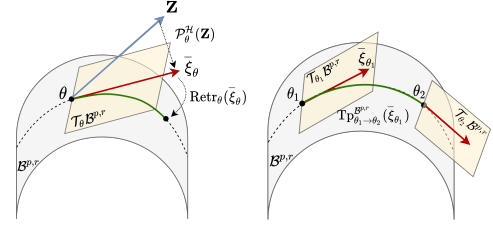
In the literature, the LRA D structure is used to estimate covariance matrices based on the GFA model [17], [18], while our model imposes the LRA D structure on precision matrices. Like the GFA model, our LRA D assumption also has a latent factor representation interpretation, known as the Gaussian precision factor model, see Proposition 1 in [16]. From a utilitarian perspective, we focus on leveraging the LRA D structures to facilitate the estimation of  $\Theta$ .

With the LRA D assumption, problem (3) becomes

$$\min_{\Theta \in \mathcal{M}^{p,r}} f(\Theta) \quad (5)$$

$$\mathcal{M}^{p,r} := \{\Theta : \Theta \in \mathcal{S}_{++}^p, \Theta = \mathbf{L} + \mathbf{D}, \mathbf{D} \in \mathcal{D}_{++}^p, \mathbf{L} \in \mathcal{S}_{+}^{p,r}\}.$$

Compared with (3), our model (5) has two merits. First, it reduces the dimensionality of the estimation problem from  $p(p+1)/2$  to  $p(r+1)$ , where we usually have  $r \ll p$ . Thus, our model requires fewer samples to recover the graph and can be applied to high-dimensional regimes. Second, the LRA D decomposition guarantees that the estimated  $\Theta$  lies in  $\mathcal{S}_{++}^p$  automatically, while traditional methods may require post-processing to ensure the positive definiteness of  $\Theta$ .



(a) Projection and retraction (b) Vector transport

Fig. 1: Illustration of operations in manifold  $\mathcal{B}^{p,r}$ .

### III. PROPOSED ALGORITHM

The proposed model (5) is non-convex and difficult to solve due to the rank constraint. Thus, we exploit manifold optimization to solve (5). The main idea is to view the LRA D constraint as a Riemann manifold, on which optimization algorithms are performed. To do this, we first construct a Riemannian manifold based on  $\mathcal{M}^{p,r}$ . Then, we propose an RCG algorithm over the constructed manifold to solve (5).

#### A. Riemannian Manifold of $\mathcal{M}^{p,r}$

We focus on using the particular structure of  $\mathbf{L}$  and  $\mathbf{D}$  to define a manifold. Specifically, for any  $\mathbf{L} \in \mathcal{S}_{+}^{p,r}$ , it allows a full rank decomposition  $\mathbf{L} = \mathbf{Y}\mathbf{Y}^\top$ , where  $\mathbf{Y} \in \mathcal{R}^{p,r}$ ,  $\mathcal{R}^{p,r} := \{\mathbf{Y} \in \mathbb{R}^{p \times r} : \det(\mathbf{Y}^\top \mathbf{Y}) \neq 0\}$  is the set of full-rank  $p \times r$  matrices. Let  $\mathcal{B}^{p,r} = \mathcal{R}^{p,r} \times \mathcal{D}_{++}^p$ , and we define

$$\phi : \mathcal{B}^{p,r} \rightarrow \mathcal{M}^{p,r} : (\mathbf{Y}, \mathbf{D}) \mapsto \phi(\mathbf{Y}, \mathbf{D}) = \mathbf{Y}\mathbf{Y}^\top + \mathbf{D}. \quad (6)$$

The optimization on  $\mathcal{M}^{p,r}$  can be transformed into the optimization on  $\mathcal{B}^{p,r}$  via  $\phi$ . However, for any  $\mathbf{O} \in \mathcal{O}^r$ , where  $\mathcal{O}^r$  is the orthogonal group in dimension  $r$ , we have an equivalence class w.r.t.  $\phi$ , i.e.,  $\phi(\mathbf{Y}\mathbf{O}, \mathbf{D}) = \phi(\mathbf{Y}, \mathbf{D})$ . This equivalence class of  $\theta := (\mathbf{Y}, \mathbf{D}) \in \mathcal{B}^{p,r}$  is denoted by  $[\theta] := \{\theta * \mathbf{O} : \mathbf{O} \in \mathcal{O}^r\}$ , where  $\theta * \mathbf{O} = (\mathbf{Y}\mathbf{O}, \mathbf{D})$ . The invariance property renders that the solution of (5) is not isolated. To address this issue, we define a quotient set  $\mathcal{B}^{p,r}/\mathcal{O}^r$  as

$$\mathcal{B}^{p,r}/\mathcal{O}^r := \{[\theta] : \theta \in \mathcal{B}^{p,r}\}. \quad (7)$$

The quotient map is then defined as

$$\pi : \mathcal{B}^{p,r} \rightarrow \mathcal{B}^{p,r}/\mathcal{O}^r : \theta \mapsto \pi(\theta) = [\theta]. \quad (8)$$

In practice, we can manipulate objects in  $\mathcal{B}^{p,r}$  and eliminate the effect of the equivalence class. As shown in Fig.1, we provide the operations required to design algorithms on  $\mathcal{B}^{p,r}$ .

(i) **Riemannian metric:** The premise of operations on manifolds is to define metrics. First, the tangent space at a point  $\theta \in \mathcal{B}^{p,r}$  is

$$\mathcal{T}_\theta \mathcal{B}^{p,r} = \{(\xi_Y, \xi_D) : \xi_Y \in \mathbb{R}^{p \times r}, \xi_D \in \mathcal{D}^p\}. \quad (9)$$

For  $\xi = (\xi_Y, \xi_D)$  and  $\zeta = (\zeta_Y, \zeta_D) \in \mathcal{T}_\theta \mathcal{B}^{p,r}$ , the metric is defined as sum of the metrics of  $\mathcal{R}^{p,r}$  and  $\mathcal{D}_{++}^p$ , i.e.,

$$\langle \xi, \zeta \rangle_\theta^{\mathcal{B}^{p,r}} = \text{tr}(\xi_Y^\top \zeta_Y) + \text{tr}(\mathbf{D}^{-1} \xi_D \mathbf{D}^{-1} \zeta_D). \quad (10)$$

The first term of the r.h.s. of (10) is from [26], and the second term is the affine-invariant one in [27].

(ii) **Projection:** The equivalence class  $[\theta]$  is a subspace of  $\mathcal{B}^{p,r}$ , and its tangent space  $\mathcal{T}_\theta [\theta]$  is also a subspace of  $\mathcal{T}_\theta \mathcal{B}^{p,r}$ , which is known as vertical space  $\mathcal{V}_\theta$

$$\mathcal{V}_\theta = \{(\mathbf{Y}\mathbf{O}, \mathbf{0}) : \mathbf{O} = -\mathbf{O}^\top, \mathbf{O} \in \mathbb{R}^{p \times p}\}, \quad (11)$$

where  $\mathbf{0}$  is the matrices of zeros. The space  $\mathcal{V}_\theta$  contains directions moving along the equivalence class  $[\theta]$ , which

should be eliminated. Thus, the vectors of interest in  $\mathcal{T}_\theta \mathcal{B}^{p,r}$  lie in the orthogonal complement of  $\mathcal{V}_\theta$ , termed horizontal space  $\mathcal{H}_\theta$

$$\mathcal{H}_\theta = \{(\mathbf{Z}_1, \mathbf{Z}_2) : \mathbf{Y}^\top \mathbf{Z}_1 = \mathbf{Z}_1^\top \mathbf{Y}, \mathbf{Z} \in \mathbb{R}^{p \times r}, \mathbf{Z}_2 \in \mathcal{D}^p\}. \quad (12)$$

Any tangent vector  $\bar{\xi}_{[\theta]} \in \mathcal{T}_{[\theta]} \mathcal{B}^{p,r} / \mathcal{O}^r$  can be represented by a unique vector  $\bar{\xi}_\theta = (\bar{\xi}_Y, \bar{\xi}_D) \in \mathcal{H}_\theta$  satisfying  $D\pi(\theta)[\bar{\xi}_\theta] = \bar{\xi}_{[\theta]}$ , where  $D\pi(\theta)[\bar{\xi}_\theta]$  is the directional derivative of  $\pi(\theta)$  in the direction  $\bar{\xi}_\theta$ . The  $\bar{\xi}_\theta$  is called horizontal lift of  $\bar{\xi}_{[\theta]}$  at  $\theta$  [23]. We can eliminate the effect of equivalence classes by projecting points onto  $\mathcal{H}_\theta$ . For any point  $\mathbf{Z} = (\mathbf{Z}_1, \mathbf{Z}_2)$  in the ambient space  $\mathbb{R}^{p \times r} \times \mathbb{R}^{p \times p}$ , the projection is defined as

$$\mathcal{P}_\theta^{\mathcal{H}}(\mathbf{Z}) = (\mathbf{Z}_1 - \mathbf{Y}\Omega, \text{ddiag}(\mathbf{Z}_2)), \quad (13)$$

where  $\text{ddiag}(\cdot)$  means only retaining diagonal elements, and  $\Omega$  is a skew-symmetric matrix satisfying Sylvester equation

$$\Omega \mathbf{Y}^\top \mathbf{Y} + \mathbf{Y}^\top \mathbf{Y} \Omega = \mathbf{Y}^\top \mathbf{Z}_1 - \mathbf{Z}_1^\top \mathbf{Y}. \quad (14)$$

(iii) **Retraction**: Given a tangent vector  $\bar{\xi}_\theta \in \mathcal{H}_\theta$ , retracting it back to the manifold  $\mathcal{B}^{p,r}$  is defined as

$$\text{Retr}_\theta(\bar{\xi}_\theta) = (\mathbf{Y} + \bar{\xi}_Y, \mathbf{D} + \bar{\xi}_D + \frac{1}{2} \bar{\xi}_D \mathbf{D}^{-1} \bar{\xi}_D). \quad (15)$$

(iv) **Vector transport**: Given  $\theta_1, \theta_2 \in \mathcal{B}^{p,r}$  and  $\bar{\xi}_{\theta_1} \in \mathcal{H}_{\theta_1}$ , the vector transport of  $\bar{\xi}_{\theta_1}$  to  $\mathcal{H}_{\theta_2}$  is

$$\text{Tp}_{\theta_1 \rightarrow \theta_2}^{\mathcal{B}^{p,r}}(\bar{\xi}_{\theta_1}) = \mathcal{P}_{\theta_2}^{\mathcal{H}}(\bar{\xi}_{\theta_1}). \quad (16)$$

We have defined all required operations on the manifold  $\mathcal{B}^{p,r}$ . Next, we will solve (5) on the  $\mathcal{B}^{p,r}$  using these operations.

### B. Riemannian Conjugate Gradient Algorithm

The objective function  $f(\Theta) : \mathcal{S}_{++}^p \rightarrow \mathbb{R}$  can be rewritten as  $\bar{f}(\theta) : \mathcal{B}^{p,r} \rightarrow \mathbb{R}$ , where  $\bar{f} = f \circ \phi$ . To solve (5) on  $\mathcal{B}^{p,r}$ , we need to derive the Riemannian gradient  $\nabla^{\mathcal{B}^{p,r}} \bar{f}(\theta)$ , which can be calculated from the Euclidean gradient  $\nabla \bar{f}(\theta) = (\mathbf{G}_Y, \mathbf{G}_D)$  via metric (10). Specifically,  $\nabla^{\mathcal{B}^{p,r}} \bar{f}(\theta)$  is equivalent to  $\nabla \bar{f}(\theta)$  projected onto  $\mathcal{T}_\theta \mathcal{B}^{p,r}$ , i.e.,

$$\nabla^{\mathcal{B}^{p,r}} \bar{f}(\theta) = (\mathbf{G}_Y, \text{Dddiag}(\mathbf{G}_D)\mathbf{D}), \quad (17)$$

where

$$\mathbf{G}_Y = 2\nabla f(\phi(\theta))\mathbf{Y}, \mathbf{G}_D = \text{ddiag}(\nabla f(\phi(\theta))), \quad (18)$$

and

$$\nabla f(\phi(\theta)) = \nabla f(\Theta) = \mathbf{S} - \Theta^{-1} + \beta \nabla h(\Theta). \quad (19)$$

Here, we use  $\psi(z) = \nu \log(\cosh(z/\nu))$ ,  $\nu > 0$ , as an alternative to widely-used  $\ell_1$  norm since the algorithm requires  $\psi(z)$  to be smooth [18]. When  $\nu \rightarrow 0$ ,  $\psi(z)$  yields  $\ell_1$  norm. Empirically, we find that the results are insensitive to  $\nu$ , and we let  $\nu = 10^{-12}$ . Then,  $\nabla h(\Theta)$  is calculated as

$$\nabla h(\Theta) = \begin{cases} 0, & i = j \\ \tanh(\Theta[ij]/\nu), & i \neq j. \end{cases} \quad (20)$$

Since the objective function of interest is invariant to the equivalence class,  $\nabla^{\mathcal{B}^{p,r}} \bar{f}(\theta)$  naturally lies in  $\mathcal{H}_\theta$ , which can be used directly. After obtaining  $\nabla^{\mathcal{B}^{p,r}} \bar{f}(\theta)$ , the update of our RCG algorithm [28] at the  $t$ -th iteration is

$$\boldsymbol{\eta}^{(k)} = -\nabla^{\mathcal{B}^{p,r}} \bar{f}(\theta^{(k)}) + \mu^{(k)} \text{Tp}_{\theta^{(k-1)} \rightarrow \theta^{(k)}}^{\mathcal{B}^{p,r}}(\boldsymbol{\eta}^{(k-1)}) \quad (21)$$

$$\theta^{(k+1)} = \text{Retr}_{\theta^{(k)}}(\gamma^{(k)} \boldsymbol{\eta}^{(k)}), \quad (22)$$

where  $\gamma^{(k)}$  is the stepsize determined by linesearch [22], and  $\mu^{(k)}$  can be computed via the method in [29]. Eq.(21) calculates the descent direction  $\boldsymbol{\eta}^{(k)}$ , which is composed of  $-\nabla^{\mathcal{B}^{p,r}} \bar{f}(\theta)$  plus the previous descent direction  $\boldsymbol{\eta}^{(k-1)}$ . Moreover,  $\boldsymbol{\eta}^{(k-1)}$  is transported to  $\mathcal{H}_{\theta^{(k)}}$  so that it can be added to  $-\nabla^{\mathcal{B}^{p,r}} \bar{f}(\theta^{(k)})$ . Eq.(22) performs updates using the

### Algorithm 1 RCG algorithm to solve (5)

**Input:** Data  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ ,  $\beta$ , rank  $r$

**Output:** The learned graph (precision matrix)  $\Theta$

- 1: Initialize  $\theta^{(0)} \in \mathcal{B}^{p,r}$  randomly, and let  $\boldsymbol{\eta}^{(0)} = -\nabla^{\mathcal{B}^{p,r}} \bar{f}(\theta^{(0)})$
- 2: **for**  $k = 0, 1, 2, \dots$ , until convergence **do**
- 3:   Compute the stepsize  $\gamma^{(k)}$  using the linesearch in [22]
- 4:   Update  $\theta^{(k+1)}$  using (22)
- 5:   Compute the stepsize  $\mu^{(k+1)}$  using the rule in [29]
- 6:   Let  $\Theta^{(k+1)} = \phi(\theta^{(k+1)})$
- 7:   Compute  $\nabla^{\mathcal{B}^{p,r}} \bar{f}(\theta^{(k+1)})$  using (17) - (20)
- 8:   Update  $\boldsymbol{\eta}^{(k+1)}$  using (21)
- 9: **end for**

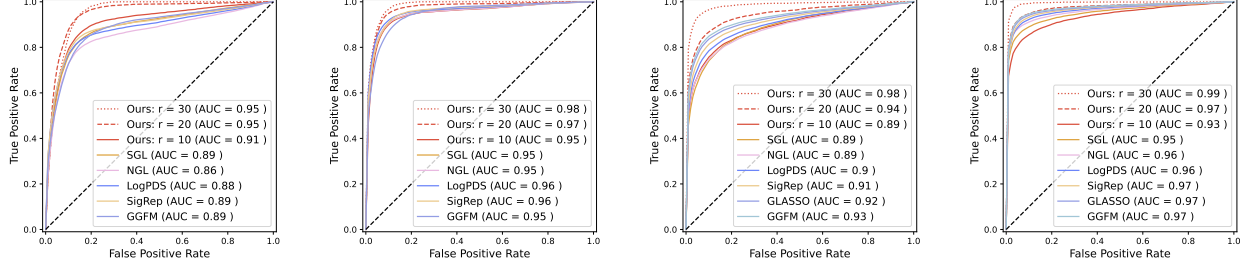
calculated  $\boldsymbol{\eta}^{(k)}$  and retracts the updates onto the manifold so that  $\theta^{(k+1)}$  still lies in  $\mathcal{B}^{p,r}$ . The complete procedure of the proposed algorithm is displayed in Algorithm 1.

Our algorithm calculates the Euclidean gradient  $\nabla f(\phi(\theta))$  to obtain the Riemannian gradient, which costs  $O(p^3)$  per iteration due to the inversion of  $\Theta$ . All other operations are quite efficient thanks to the LRaD structure. Specifically, the vector transport on the manifold needs  $O(pr^2 + pr + r^3)$ , and the retraction requires  $O(pr)$ . Usually, we have  $r \ll p$ . Thus, the computational burden lies in the operations of Euclidean space, and the manifold-related operations are efficient.

## IV. EXPERIMENTS

### A. Synthetic Data

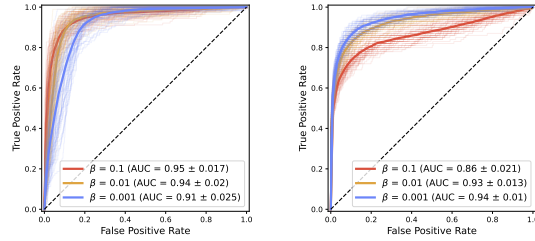
In this study, we generate two synthetic graphs with  $p = 50$ , i.e., the graph with LRaD structures and the Watts-Strogatz (WS) random graph. For the LRaD graph, we first generate a sparse matrix  $\mathbf{Y}^* \in \mathbb{R}^{50 \times 10}$  and let  $\Theta^* = \mathbf{Y}^* \mathbf{Y}^{*\top} + \mathbf{D}^*$ , where  $\mathbf{D}^* \in \mathcal{D}_{++}^p$  is a random diagonal matrix such that  $\Theta^* \in \mathcal{S}_{++}^{50 \times 50}$ . We generate the adjacency matrix  $\mathbf{W}^*$  of WS random graphs, where each node has five neighbours, and the edges are randomly reconnected with probability 0.1. The edge weights are sampled from the uniform distribution  $\mathcal{U}(2, 5)$  by following [13]. Then, we let  $\Theta^* = \Xi^* + 0.1 \cdot \mathbf{I}$ , where  $\Xi^*$  is the Laplacian matrix corresponding to  $\mathbf{W}^*$ , and  $\mathbf{I}$  is the identity matrix. Based on the generated  $\Theta^*$ ,  $N$  samples are drawn from  $\mathcal{N}(\mathbf{0}, (\Theta^*)^{-1})$ . We compare our method with six baselines, i.e., GGFM [18], GLASSO [9], SGL [12], NGL [13], SigRep [2], and LogPDS [8]. The first four methods are based on GGMs, and the last two are based on graph signal processing. In GL, detecting the existence of edges in a graph can be regarded as a binary classification problem. Thus, we evaluate results using the receiver operating characteristic (ROC) curves obtained from the estimated precision matrices. The ROC curves depict the relationship between the true positive rate (TPR) and the false positive rate (FPR), where TPR denotes the capacity to recover real edges, and FPR denotes the false recovery of non-existent edges. We calculate the area under curve (AUC) of each ROC curve, whose value is in  $[0, 1]$ , where 1 represents a perfect recovery of the true edges. We do not compare edge weights learned by different baselines since they use different normalization methods, which may not be comparable. The parameters of all baselines are selected



(a) LRA graph,  $N/p = 2, \beta = 0.05$  (b) LRA graph,  $N/p = 4, \beta = 0.05$  (c) WS graph,  $N/p = 2, \beta = 0.05$  (d) WS graph,  $N/p = 4, \beta = 0.05$   
 Fig. 2: ROC curves of the learned graphs. The GLASSO is ignored in the LRA graph due to numerical divergence.

TABLE I: Modularity of the learned graphs.

|     | GLASSO | SGL  | NGL  | LogPDS | SigRep | Ours        |
|-----|--------|------|------|--------|--------|-------------|
| Mod | 0.57   | 0.74 | 0.42 | 0.70   | 0.74   | <b>0.85</b> |



(a) LRA graph,  $N/p = 2, r = 20$  (b) WS graph,  $N/p = 2, r = 20$   
 Fig. 3: Effect of sparsity parameter  $\beta$  of our method.

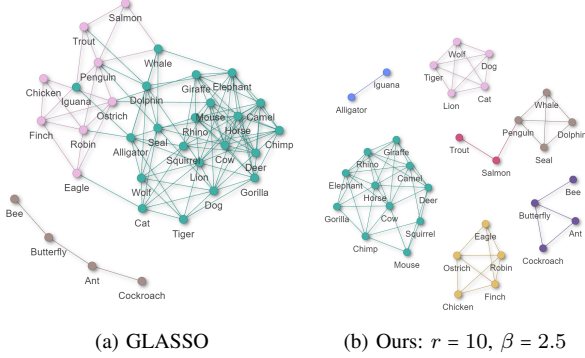


Fig. 4: The graphs of the *animal* dataset. Colors of the nodes represent communities detected from the learned graphs.

as those achieving the best AUC values. All results are the average of 50 independent experiments.

Figure 2 depicts ROC curves of the graphs estimated by different methods. Overall, our method outperforms all baselines, especially when  $N$  is small. As  $N$  increases, the performance improvement of our method becomes less significant. In the LRA graph, our method achieves competitive AUC results when  $r = 10$  due to the low-rank structure of the underlying precision matrix. The performance of our method improves as  $r$  increases. In the WS graph, our method is inferior to the baselines when  $r = 10$  since it may not suffice to describe the real WS graph. However, our method obtains better AUC results when  $r$  increases to 15 and 20 despite WS graphs do not exactly have LRA structures. The results show that we can reasonably reduce the dimensionality by assuming the LRA structure without significantly sacrificing performance.

Figure 3 depicts the effect of the parameter  $\beta$ . It is observed that our method can achieve satisfactory performance in a wide range of  $\beta$ . In the WS graph, our method is more sensitive to  $\beta$ , meaning that we should carefully choose this parameter.

### B. Real Data

We test the proposed method using the *animal* dataset [11], which contains the descriptions of different animals. For example, is the animal poisonous or can the animal fly? In total, there are  $N = 102$  questions for 33 animals, and hence  $\mathbf{X} \in \mathbb{R}^{33 \times 102}$ . Since the ground-truth graph is unavailable in the real dataset, we cannot use ROC curves to evaluate the learned graphs. Note that animals of different species (such as mammals and birds) are distinct, while animals of the same species share similar descriptions. The learned graphs should have several communities. Thus, we employ the algorithm in [30] to detect communities in the learned graphs and evaluate the detection results using *modularity* (Mod) [31], a metric to measure the separability of a graph into sub-components. High Mod values mean better separability.

Table I lists the modularity of the graphs learned by different methods. Our method obtains the highest Mod values, implying that the corresponding graph better preserves the underlying communities. We visualize the learned graphs in Fig.4. Due to space limitations, only the results of our method and GLASSO are provided. In our graph, the nodes of the animals in the same species are closely connected, while those corresponding to different species (such as mammals, birds, and insects) are disconnected. Compared with GLASSO, our method can better recover the underlying communities.

## V. CONCLUSION

We proposed a graph estimator in GGMs by assuming an LRA structure, which can reduce the dimensionality of the problem. A Riemannian manifold describing the LRA constraint was constructed, and an RCG algorithm was designed to solve the proposed model. Finally, experiments on synthetic and real data showed the superiority of our method.

## REFERENCES

- [1] D. I. Shuman, S. K. Narang, P. Frossard, A. Ortega, and P. Vandergheynst, "The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains," *IEEE Signal Process. Mag.*, vol. 30, no. 3, pp. 83–98, 2013.
- [2] X. Dong, D. Thanou, P. Frossard, and P. Vandergheynst, "Learning laplacian matrix in smooth graph signal representations," *IEEE Trans. Signal Process.*, vol. 64, no. 23, pp. 6160–6173, 2016.
- [3] X. Pu, T. Cao, X. Zhang, X. Dong, and S. Chen, "Learning to learn graph topologies," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 34, pp. 4249–4262, 2021.
- [4] G. Mateos, S. Segarra, A. G. Marques, and A. Ribeiro, "Connecting the dots: Identifying network structure via graph signal processing," *IEEE Signal Process. Mag.*, vol. 36, no. 3, pp. 16–43, 2019.
- [5] X. Dong, D. Thanou, M. Rabbat, and P. Frossard, "Learning graphs from data: A signal representation perspective," *IEEE Signal Process. Mag.*, vol. 36, no. 3, pp. 44–63, 2019.
- [6] Y. Shen, G. B. Giannakis, and B. Baingana, "Nonlinear structural vector autoregressive models with application to directed brain networks," *IEEE Trans. Signal Process.*, vol. 67, no. 20, pp. 5325–5339, 2019.
- [7] A. Bolstad, B. D. Van Veen, and R. Nowak, "Causal network inference via group sparse regularization," *IEEE Trans. Signal Process.*, vol. 59, no. 6, pp. 2628–2641, 2011.
- [8] V. Kalofolias, "How to learn a graph from smooth signals," in *Proc. Int. Conf. Artif. Intell. Stat.* PMLR, 2016, pp. 920–929.
- [9] J. Friedman, T. Hastie, and R. Tibshirani, "Sparse inverse covariance estimation with the graphical lasso," *Biostatistics*, vol. 9, no. 3, pp. 432–441, 2008.
- [10] M. Yuan and Y. Lin, "Model selection and estimation in the Gaussian graphical model," *Biometrika*, vol. 94, no. 1, pp. 19–35, 2007.
- [11] B. Lake and J. Tenenbaum, "Discovering structure by learning sparse graphs," in *Proc. 33rd Annu. Cogn. Sci. Conf.*, vol. 32, no. 32, 2010.
- [12] S. Kumar, J. Ying, J. V. d. M. Cardoso, and D. P. Palomar, "A unified framework for structured graph learning via spectral constraints," *J. Mach. Learn. Res.*, vol. 21, no. 22, pp. 1–60, 2020.
- [13] J. Ying, J. V. de Miranda Cardoso, and D. Palomar, "Noncon-vex sparse graph learning under laplacian constrained graphical model," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 33, pp. 7101–7113, 2020.
- [14] J. Ying, J. V. De Miranda Cardoso, and D. P. Palomar, "Adaptive estimation of graphical models under total positivity," in *Proc. Int. Conf. Mach. Learn.*, vol. 202. PMLR, 23–29 Jul 2023, pp. 40 054–40 074.
- [15] A. Anandkumar, V. Y. Tan, F. Huang, and A. S. Willsky, "High-dimensional gaussian graphical model selection: Walk summability and local separation criterion," *J. Mach. Learn. Res.*, vol. 13, no. 76, pp. 2293–2337, 2012.
- [16] N. K. Chandra, P. Mueller, and A. Sarkar, "Bayesian scalable precision factor analysis for massive sparse gaussian graphical models," *arXiv preprint arXiv:2107.11316*, 2021.
- [17] R. Zhou, J. Ying, and D. P. Palomar, "Covariance matrix estimation under low-rank factor model with nonnegative correlations," *IEEE Trans. Signal Process.*, vol. 70, pp. 4020–4030, 2022.
- [18] A. Hippert-Ferrer, F. Bouchard, A. Mian, T. Vayer, and A. Breloy, "Learning graphical factor models with Riemannian optimization," in *Proc. Eur. Conf. Mach. Learn. (ECML)*. Springer, 2023, pp. 349–366.
- [19] P. Stoica and P. Babu, "Low-rank covariance matrix estimation for factor analysis in anisotropic noise: Application to array processing and portfolio selection," *IEEE Trans. Signal Process.*, vol. 71, pp. 1699–1711, 2023.
- [20] Y. Wu, Y. Qin, and M. Zhu, "High-dimensional covariance matrix estimation using a low-rank and diagonal decomposition," *Can. J. Stat.*, vol. 48, no. 2, pp. 308–337, 2020.
- [21] S. Boyd, S. P. Boyd, and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [22] P.-A. Absil, R. Mahony, and R. Sepulchre, *Optimization algorithms on matrix manifolds*. Princeton University Press, 2008.
- [23] N. Boumal, *An introduction to optimization on smooth manifolds*. Cambridge University Press, 2023.
- [24] E. Massart and P.-A. Absil, "Quotient geometry with simple geodesics for the manifold of fixed-rank positive-semidefinite matrices," *Siam J. Matrix Anal. A.*, vol. 41, no. 1, pp. 171–198, 2020.
- [25] M. Journée, F. Bach, P.-A. Absil, and R. Sepulchre, "Low-rank optimization on the cone of positive semidefinite matrices," *SIAM J. Optim.*, vol. 20, no. 5, pp. 2327–2351, 2010.
- [26] S. Zheng, W. Huang, B. Vandereycken, and X. Zhang, "Riemannian optimization using three different metrics for hermitian psd fixed-rank constraints: an extended version," *arXiv preprint arXiv:2204.07830*, 2023.
- [27] R. Bhatia, *Positive definite matrices*. Princeton University Press, 2009.
- [28] H. Sato, "Riemannian conjugate gradient methods: General framework and specific algorithms with convergence analyses," *SIAM J. Optimiz.*, vol. 32, no. 4, pp. 2690–2717, 2022.
- [29] M. R. Hestenes, E. Stiefel et al., *Methods of conjugate gradients for solving linear systems*. NBS Washington, DC, 1952, vol. 49, no. 1.
- [30] G. Cordasco and L. Gargano, "Community detection via semi-synchronous label propagation algorithms," *arXiv preprint arXiv:1103.4550*, 2011.
- [31] M. E. Newman, "Modularity and community structure in networks," in *Proc. Natl. Acad. Sci.*, vol. 103, no. 23, pp. 8577–8582, 2006.