
Learning Bipartite Graphs: Heavy Tails and Multiple Components

José Vinícius de M. Cardoso¹, Jiayi Ying², Daniel P. Palomar^{1,3}

Department of Electronic and Computer Engineering¹

Department of Mathematics²

Department of Industrial Engineering and Decision Analytics³

The Hong Kong University of Science and Technology

Clear Water Bay, Hong Kong SAR China

{jvdmc, jx.ying}@connect.ust.hk, palomar@ust.hk

Abstract

We investigate the problem of learning an undirected, weighted bipartite graph under the Gaussian Markov random field model, for which we present an optimization formulation along with an efficient algorithm based on the projected gradient descent. Motivated by practical applications, where outliers or heavy-tailed events are present, we extend the proposed learning scheme to the case in which the data follow a multivariate Student- t distribution. As a result, the optimization program is no longer convex, but a verifiably convergent iterative algorithm is proposed based on the majorization-minimization framework. Finally, we propose an efficient and provably convergent algorithm for learning k -component bipartite graphs that leverages rank constraints of the underlying graph Laplacian matrix. The proposed estimators outperform state-of-the-art methods for bipartite graph learning, as evidenced by real-world experiments using financial time series data.

1 Introduction

Efficient optimization formulations alongside scalable optimization algorithms have been critical tools for performing inference in graphical models. In particular, learning sparse, *unconstrained* Gaussian Markov random fields (MRF) has enabled a wide range of practical applications across a number of fields, including brain network analysis [1], single-cell sequencing [2], time-varying network estimation [3], time-series clustering [4] and model selection [5], and financial networks [6, 7]. Such an impact may be attributed to the early development of the graphical lasso algorithm via iterative block coordinate descent methods [8–10].

More recently, there has been a growing interest in *constraining* Gaussian MRFs into a particular family such as bipartite and k -component graphs for data clustering [11–14], and considering robust statistics that account for heavy-tailed events [15, 16]. In addition, constraints on the sign of the elements of the precision matrix of a Gaussian MRF have been employed to learn total positivity models [17–20] as well as Laplacian models in sparse settings [21–23]. Estimating undirected graphs under smooth signal assumptions has been investigated in [24–27].

In this paper, we consider the problem of learning undirected, weighted bipartite graphs under the MRF assumption. The key contributions of this paper include:

1. We design an efficient optimization algorithm based on the projected gradient descent (PGD) framework for the problem of learning an undirected, weighted bipartite graph from the Gaussian MRF framework.

2. We extend our proposed formulation in order to accommodate the multivariate Student- t distribution, which is often needed for real datasets that contain outliers or are inherently heavy-tailed. In this case, the optimization formulation turns out to be nonconvex. We design an optimization algorithm based on the majorization-minimization (MM) framework to obtain a stationary point of such problem.
3. We propose an efficient and provably convergent algorithm based on the alternating direction method of multipliers (ADMM) for learning k -component bipartite graphs that leverages the rank properties of the Laplacian matrix.
4. We present empirical results, using real-world datasets from the US stock market, which reveal that our formulations outperform state-of-the-art ones in terms of graph modularity and node label accuracy.

Notation: Matrices (vectors) are denoted by bold, italic, capital (lowercase) roman letters like \mathbf{X} , \mathbf{x} . Vectors are assumed to be column vectors. We use $\mathbf{1}_n$ (resp. \mathbf{I}_n) to denote the n -dimensional all-one vector (resp. identity matrix). The (i, j) element of a matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ is denoted as X_{ij} . The i -th element of a vector \mathbf{x} is denoted as x_i . The i -th row of \mathbf{X} is denoted as $\mathbf{x}_i \in \mathbb{R}^{p \times 1}$. $\mathbf{X} \geq \mathbf{Y}$ implies that $X_{ij} \geq Y_{ij} \forall i, j$, whereas $\mathbf{X} \succeq \mathbf{Y}$ implies that $\mathbf{X} - \mathbf{Y}$ is positive semi-definite. For a square matrix \mathbf{X} , $\lambda_i(\mathbf{X})$ denotes the i -th smallest eigenvalue of \mathbf{X} . For a vector \mathbf{x} and indices $j > i$, $\mathbf{x}_{i:j}$ denotes $(x_i, x_{i+1}, \dots, x_j)^\top$. The operator $\text{Diag} : \mathbb{R}^p \rightarrow \mathbb{R}^{p \times p}$ creates a diagonal matrix with the elements of an input vector along its diagonal. The operator $\text{diag} : \mathbb{R}^{p \times p} \rightarrow \mathbb{R}^p$ extracts the diagonal of a square matrix. The optimal point of an optimization variable \mathbf{v} is denoted as \mathbf{v}^* .

2 Preliminaries & Related Works

An undirected, weighted, bipartite graph can be defined as a 4-tuple $\mathcal{G} \triangleq \{\mathcal{V}_r, \mathcal{V}_q, \mathcal{E}, \mathbf{W}\}$, where $\mathcal{V}_r \triangleq \{1, 2, \dots, r\}$ and $\mathcal{V}_q \triangleq \{r+1, r+2, \dots, r+q\}$ are the node sets associated to a group of objects and classes, respectively, where $p \triangleq r+q$ is the total number of nodes, and we assume that $r \gg q$, *i.e.*, there exist many more objects than classes; $\mathcal{E} \subseteq \{\{u, v\} : u \in \mathcal{V}_r, v \in \mathcal{V}_q\}$ is the edge set, that is, a subset of the set of all possible unordered pairs of nodes such that $\{u, v\} \in \mathcal{E}$ iff nodes u and v are connected; $\mathbf{W} \in \mathbb{R}_+^{p \times p}$ is the symmetric weighted adjacency matrix that satisfies $W_{ii} = 0, W_{ij} > 0$ iff $\{i, j\} \in \mathcal{E}$ and $W_{ij} = 0$, otherwise. The weighted Laplacian matrix of a graph is defined as $\mathbf{L} \triangleq \text{Diag}(\mathbf{W}\mathbf{1}) - \mathbf{W}$.

Our goal is to learn a bipartite graph from data under probabilistic assumptions. Thus, the data generating process is initially assumed to be a zero-mean improper Gaussian Markov random field (IGMRF) [28] $\mathbf{x} \in \mathbb{R}^p$, such that x_i is the random variable generating a signal measured at node i , whose rank-deficient precision matrix is modeled as a graph Laplacian matrix [21, 22, 29]. Assume we are given n observations of \mathbf{x} , *i.e.*, $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^\top$, $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{x}_i \in \mathbb{R}^{p \times 1}$. The goal of graph learning algorithms is to learn a Laplacian matrix given the data matrix \mathbf{X} .

To that end, the classical MLE of the Laplacian-constrained precision matrix of \mathbf{x} , on the basis of the observed data \mathbf{X} , may be formulated as the following optimization program [21, 22, 29]:

$$\underset{\mathbf{L} \succeq \mathbf{0}}{\text{minimize}} \quad \text{tr}(\mathbf{L}\mathbf{S}) - \log \det^*(\mathbf{L}), \text{ subject to } \mathbf{L}\mathbf{1} = \mathbf{0}, L_{ij} = L_{ji} \leq 0, i \neq j, \quad (1)$$

where \mathbf{S} is a similarity matrix, *e.g.*, the sample covariance (or correlation) matrix $\mathbf{S} \propto \mathbf{X}^\top \mathbf{X}$, and $\det^*(\mathbf{L})$ is the pseudo-determinant of \mathbf{L} , *i.e.*, the product of its positive eigenvalues [30].

The Laplacian matrix \mathbf{L} of a bipartite graph can be written as

$$\mathbf{L} = \begin{bmatrix} \text{Diag}(\mathbf{B}\mathbf{1}_q) & -\mathbf{B} \\ -\mathbf{B}^\top & \text{Diag}(\mathbf{B}^\top \mathbf{1}_r) \end{bmatrix}, \quad (2)$$

where $\mathbf{B} \in \mathbb{R}_+^{r \times q}$ contains the edge weights between the nodes of objects and the nodes of classes. Note that \mathbf{L} satisfies $\mathbf{L}\mathbf{1} = \mathbf{0}$ and $L_{ij} = L_{ji} \leq 0 \forall i \neq j$ by definition.

In what follows, we briefly discuss the two main approaches proposed in the literature to tackle the estimation of bipartite graphs.

Bipartite Structure. The authors in [12] proposed the following optimization problem to learn a k -component bipartite graph from a given bipartite graph weights $\mathbf{A} \in \mathbb{R}^{r \times q}$:

$$\underset{\mathbf{B}, \mathbf{V} \in \mathbb{R}^{p \times k}}{\text{minimize}} \quad \|\mathbf{B} - \mathbf{A}\|_{\text{F}}^2 + \eta \text{tr}(\mathbf{V}^{\top} \mathbf{L} \mathbf{V}), \text{ subject to } \mathbf{B} \geq \mathbf{0}, \mathbf{B} \mathbf{1}_q = \mathbf{1}_r, \mathbf{V}^{\top} \mathbf{V} = \mathbf{I}_k, \quad (3)$$

where \mathbf{L} depends on \mathbf{B} through (2), $\eta > 0$ is a hyperparameter that promotes the rank of \mathbf{L} to be $p - k$, and \mathbf{A} can be constructed from the correlation between nodes of objects and classes.

An alternating minimization approach was proposed in [12] to solve Problem (3). A shortcoming of this formulation is it lacks statistical support, which makes it difficult to draw statistical conclusions about the learned graph as well as to consider this model under different statistical assumptions.

Spectral Regularization. Properties associated with the spectral decomposition of graph matrices have demonstrated advantages that enable learning graphs with specific structures, such as bipartite and k -component graphs [11–14, 31]. The authors in [14] proposed a formulation to learn k -component bipartite graphs based on the GMRF framework. Their approach consists of including spectral decompositions of the Laplacian and adjacency matrices into the objective function as regularization terms. Although the formulation in [14] is derived from an statistical approach, we note that it does not always guarantee that the estimated graph is bipartite, further requiring postprocessing of the adjacency matrix. In addition, the lack of constraints on the node degrees allow trivial solutions with isolated nodes, which is also undesired in practice. The method in [14] does not require knowledge of the partition of the node set, however it does assume that the number of nodes in each set is known.

3 Proposed Formulations & Algorithms

In what follows, we present our proposed methods and algorithms, starting off with the simple case of a connected bipartite graph under Gaussian settings. It is important to look at this particular scenario first, as the proposed algorithm will serve as a building block for more elaborate formulations that include additional statistical assumptions, such as heavy tails, and its extension to k -component bipartite graphs.

3.1 Gaussian Bipartite Graphs

In this section, we propose a formulation to learn connected bipartite graphs from data along with an algorithm to find a global optimum. Noting that the adjacency matrix of a bipartite graph can be partitioned as $\mathbf{W} = [\mathbf{0} \ \mathbf{B}; \mathbf{B}^{\top} \ \mathbf{0}]$, where $\mathbf{B} \in \mathbb{R}_+^{r \times q}$, we can formulate its MLE as the following convex optimization program:

$$\underset{\mathbf{L}, \mathbf{B} \geq \mathbf{0}}{\text{minimize}} \quad \text{tr}(\mathbf{L} \mathbf{S}) - \log \det(\mathbf{L} + \mathbf{J}), \text{ subject to } \mathbf{L} = \begin{bmatrix} \text{Diag}(\mathbf{B} \mathbf{1}_q) & -\mathbf{B} \\ -\mathbf{B}^{\top} & \text{Diag}(\mathbf{B}^{\top} \mathbf{1}_r) \end{bmatrix}, \quad (4)$$

where we have used the fact that, for a connected graph, $\det^*(\mathbf{L}) = \det(\mathbf{L} + \mathbf{J})$, $\mathbf{J} \triangleq \frac{1}{p} \mathbf{1} \mathbf{1}^{\top}$ [29].

We can reformulate Problem (4) by taking advantage of the symmetry of the similarity matrix \mathbf{S} , *i.e.*, $\mathbf{S} = [\mathbf{S}_{rr} \ \mathbf{S}_{rq}; \mathbf{S}_{rq}^{\top} \ \mathbf{S}_{qq}]$ where $\mathbf{S}_{rq} \in \mathbb{R}^{r \times q}$ contains the pairwise similarities between the nodes in the set of objects \mathcal{V}_r and the nodes in the set of classes \mathcal{V}_q . In particular, we have

$$\text{tr}(\mathbf{L} \mathbf{S}) = \text{diag}(\mathbf{S})^{\top} \begin{bmatrix} \mathbf{B} \mathbf{1}_q \\ \mathbf{B}^{\top} \mathbf{1}_r \end{bmatrix} - 2 \text{tr}(\mathbf{B} \mathbf{S}_{rq}^{\top}) = \text{tr}(\mathbf{B} (\mathbf{1}_q \mathbf{s}_{1:r}^{\top} + \mathbf{s}_{r+1:p} \mathbf{1}_r^{\top} - 2 \mathbf{S}_{rq}^{\top})), \quad (5)$$

where $\mathbf{s} \triangleq \text{diag}(\mathbf{S})$.

Plugging (2) and (5) in (4), we arrive at the following formulation to learn a connected bipartite graph:

$$\underset{\mathbf{B} \geq \mathbf{0}}{\text{minimize}} \quad \text{tr}(\mathbf{B} (\mathbf{1}_q \mathbf{s}_{1:r}^{\top} + \mathbf{s}_{r+1:p} \mathbf{1}_r^{\top} - 2 \mathbf{S}_{rq}^{\top})) - \log \det \left(\begin{bmatrix} \text{Diag}(\mathbf{B} \mathbf{1}_q) + \mathbf{J}_{rr} & -\mathbf{B} + \mathbf{J}_{rq} \\ -\mathbf{B}^{\top} + \mathbf{J}_{qr} & \text{Diag}(\mathbf{B}^{\top} \mathbf{1}_r) + \mathbf{J}_{qq} \end{bmatrix} \right). \quad (6)$$

3.1.1 PGD Solution

We propose an optimization algorithm based on the PGD framework [32, 33] that amounts to two steps: (1) computation of a descent update and (2) projection onto the feasible set.

A common way to compute a descent direction is using the gradient of the objective function. The gradient of the objective function in Problem (6) would involve the computation of the inverse of the matrix inside the log-determinant term, which in general costs $O(p^3)$. By leveraging the block structure of the Laplacian matrix of bipartite graphs, and using the classical matrix determinant lemma for block matrices that relates the determinant of a block matrix with its Schur complement [34, pp. 4], we have the following equality:

$$\log \det \left(\begin{bmatrix} \text{Diag}(\mathbf{B}\mathbf{1}_q) + \mathbf{J}_{rr} & -\mathbf{B} + \mathbf{J}_{rq} \\ -\mathbf{B}^\top + \mathbf{J}_{qr} & \text{Diag}(\mathbf{B}^\top \mathbf{1}_r) + \mathbf{J}_{qq} \end{bmatrix} \right) = \log \det (\text{Diag}(\mathbf{B}\mathbf{1}_q) + \mathbf{J}_{rr}) + \log \det (\text{Diag}(\mathbf{B}^\top \mathbf{1}_r) + \mathbf{J}_{qq} - (-\mathbf{B}^\top + \mathbf{J}_{qr})(\text{Diag}(\mathbf{B}\mathbf{1}_q) + \mathbf{J}_{rr})^{-1}(-\mathbf{B} + \mathbf{J}_{rq})). \quad (7)$$

In practice, we also would like to normalize the degrees of the nodes in the objects group, so that the edge weights can be directly interpreted as the degree of membership of an object to a particular class. Mathematically, this can be achieved by a simple linear constraint $\mathbf{B}\mathbf{1}_q = \mathbf{1}_r$.

Plugging the above constraint and equality (7) into (6), we have the following formulation

$$\underset{\mathbf{B} \geq 0, \mathbf{B}\mathbf{1}_q = \mathbf{1}_r}{\text{minimize}} \quad \text{tr}(\mathbf{B}\mathbf{C}) - \log \det (\text{Diag}(\mathbf{B}^\top \mathbf{1}_r) + \mathbf{J}_{qq} - (\mathbf{B} - \mathbf{J}_{rq})^\top (\mathbf{I}_r + \mathbf{J}_{rr})^{-1} (\mathbf{B} - \mathbf{J}_{rq})), \quad (8)$$

where $\mathbf{C} \triangleq (\mathbf{s}_{r+1:p} \mathbf{1}_r^\top - 2\mathbf{S}_{rq}^\top)$.

Let $f(\mathbf{B}) \triangleq -\log \det(g(\mathbf{B})) + \text{tr}(\mathbf{B}\mathbf{C})$ be the objective function of Problem (8), where $g(\mathbf{B}) \triangleq \text{Diag}(\mathbf{B}^\top \mathbf{1}_r) + \mathbf{J}_{qq} - (\mathbf{B} - \mathbf{J}_{rq})^\top (\mathbf{I}_r + \mathbf{J}_{rr})^{-1} (\mathbf{B} - \mathbf{J}_{rq})$. Then the gradient of $f(\mathbf{B})$ can be computed as

$$\nabla f(\mathbf{B}) = \mathbf{1}_r [\text{diag}(-g(\mathbf{B})^{-1})]^\top - 2(\mathbf{I}_r + \mathbf{J}_{rr})^{-1} (\mathbf{B} - \mathbf{J}_{rq}) (-g(\mathbf{B})^{-1})^\top + \mathbf{C}^\top. \quad (9)$$

The PGD update is formulated as

$$\mathbf{B}^{l+1} = \underset{\mathbf{B} \geq 0, \mathbf{B}\mathbf{1}_q = \mathbf{1}_r}{\arg \min} \quad \|\mathbf{B} - (\mathbf{B}^l - \alpha_l \nabla f(\mathbf{B}^l))\|_{\mathbb{F}}^2 \triangleq P_{\Delta}(\mathbf{B}^l - \alpha_l \nabla f(\mathbf{B}^l)). \quad (10)$$

Problem (10) is an Euclidean projection of the rows of $\mathbf{B}^l - \alpha_l \nabla f(\mathbf{B}^l)$ onto the probability simplex. The unique solution to Problem (10) can be found efficiently via several algorithms [35–37] whose theoretical worst-case complexity is $O(rq^2)$ but their observed practical complexity is $O(rq)$ [38].

Finally, to validate the point \mathbf{B}^{l+1} and to adaptively update the learning rate α_l , we check the following backtracking condition [33]:

$$f(\mathbf{B}^{l+1}) \leq f(\mathbf{B}^l) + \langle \nabla f(\mathbf{B}^l), \mathbf{B}^{l+1} - \mathbf{B}^l \rangle + \frac{1}{2\alpha_l} \|\mathbf{B}^{l+1} - \mathbf{B}^l\|_{\mathbb{F}}^2. \quad (11)$$

In practice, we decrease the learning rate $\alpha_l \in (0, 1)$ until condition (11) is satisfied or convergence has been achieved. Algorithm 1 summarizes the proposed scheme for learning a bipartite graph from a Gaussian MRF assumption. Since the objective function of Problem (8) is convex and its feasible set is compact, Algorithm 1 converges to a global minimum [32].

3.2 Student- t Bipartite Graphs

Outliers and heavy-tailed events are pervasive in modern datasets [39]. To account for this phenomena, instead of assuming a Gaussian generative process, as done in [14] and in Problem (4), we assume the data generating process to be modeled by a multivariate zero-mean Student- t distribution, whose probability density function can be written as

$$p(\mathbf{x}) \propto \sqrt{\det^*(\Theta)} \left(1 + \frac{\mathbf{x}^\top \Theta \mathbf{x}}{\nu} \right)^{-(\nu+p)/2}, \quad (12)$$

Algorithm 1: Gaussian bipartite graph learning (GBG)

Data: Similarity matrix \mathbf{S} , initial feasible estimate of the graph weights \mathbf{B}^0 , initial learning rate $\alpha_0 > 0$, tolerance $\epsilon > 0$.

Result: Optimal bipartite graph: \mathbf{B}^*

```

1 while  $l \leq \text{maxiter}$  do
2    $\triangleright \mathbf{B}^{l+1} \leftarrow P_{\Delta}(\mathbf{B}^l - \alpha_l \nabla f(\mathbf{B}^l))$ , where  $\alpha_l$  is chosen such that (11) is satisfied
3   if  $\|\mathbf{B}^{l+1} - \mathbf{B}^l\|_{\text{F}} / \|\mathbf{B}^l\|_{\text{F}} \leq \epsilon$  then
4     | return  $\mathbf{B}^{l+1}$ 
5   end
6 end

```

where Θ is a positive-semidefinite inverse scatter matrix modeled as a combinatorial graph Laplacian matrix, and $\nu > 2$ is the number of degrees of freedom that controls the rate of decay of the tails. More precisely, as $\nu \rightarrow \infty$ the tails of (12) approach those of a Gaussian distribution.

This results in a robustified version of the MLE for connected graph learning, *i.e.*,

$$\begin{aligned}
 & \underset{\mathbf{L}, \mathbf{B}}{\text{minimize}} && -\log \det(\mathbf{L} + \mathbf{J}) + \frac{p + \nu}{n} \sum_{i=1}^n \log\left(1 + \frac{1}{\nu} \mathbf{x}_i^{\top} \mathbf{L} \mathbf{x}_i\right), \\
 & \text{subject to} && \mathbf{L} = \begin{bmatrix} \mathbf{I}_r & -\mathbf{B} \\ -\mathbf{B}^{\top} & \text{Diag}(\mathbf{B}^{\top} \mathbf{1}_r) \end{bmatrix}, \mathbf{B} \geq \mathbf{0}, \mathbf{B} \mathbf{1}_q = \mathbf{1}_r.
 \end{aligned} \tag{13}$$

We can further simplify the objective function of Problem (13) by noting that

$$\mathbf{x}_i^{\top} \mathbf{L} \mathbf{x}_i = \mathbf{x}_i^{\top} \begin{bmatrix} \mathbf{I}_r & -\mathbf{B} \\ -\mathbf{B}^{\top} & \text{Diag}(\mathbf{B}^{\top} \mathbf{1}_r) \end{bmatrix} \mathbf{x}_i = h_i + \text{tr}(\mathbf{B} \mathbf{G}_i), \tag{14}$$

where $\mathbf{G}_i \triangleq \text{diag}(\mathbf{x}_i \mathbf{x}_i^{\top})_{r+1:p} \mathbf{1}_r^{\top} - 2(\mathbf{x}_i \mathbf{x}_i^{\top})_{r,q}^{\top}$ and $h_i \triangleq \mathbf{1}_r^{\top} \text{diag}(\mathbf{x}_i \mathbf{x}_i^{\top})_{1:r}$.

Plugging (7) and (14) into (13), we have the following optimization program to learn a heavy-tailed bipartite graph:

$$\begin{aligned}
 & \underset{\mathbf{B} \geq \mathbf{0}, \mathbf{B} \mathbf{1}_q = \mathbf{1}_r}{\text{minimize}} && -\log \det(\text{Diag}(\mathbf{B}^{\top} \mathbf{1}_r) + \mathbf{J}_{qq} - (\mathbf{B} - \mathbf{J}_{rq})^{\top} (\mathbf{I}_r + \mathbf{J}_{rr})^{-1} (\mathbf{B} - \mathbf{J}_{rq})) \\
 & && + \frac{p + \nu}{n} \sum_{i=1}^n \log\left(1 + \frac{h_i + \text{tr}(\mathbf{B} \mathbf{G}_i)}{\nu}\right).
 \end{aligned} \tag{15}$$

Problem (15) is, in general, nonconvex due to the summation over concave terms and hence it is difficult to be dealt with directly. To tackle this issue, we leverage the MM framework that generates a sequence of updates, where in each iteration an upper-bounded subproblem is solved [40].

To find such upper-bound, we use the fact that the logarithm term in Problem (15) can be globally upper-bounded by its first-order Taylor expansion, *i.e.*, for all $a \geq 0, t \geq 0, b > 0$, we have:

$$\log\left(1 + \frac{t}{b}\right) \leq \log\left(1 + \frac{a}{b}\right) + \frac{t - a}{a + b}, \text{ hence at a point } \mathbf{B}^j \text{ we have that}$$

$$\log\left(1 + \frac{h_i + \text{tr}(\mathbf{B} \mathbf{G}_i)}{\nu}\right) \leq \frac{\text{tr}(\mathbf{B} \mathbf{G}_i)}{\nu + h_i + \text{tr}(\mathbf{B}^j \mathbf{G}_i)} + c,$$

where $c \triangleq \log\left(1 + \frac{h_i + \text{tr}(\mathbf{B}^j \mathbf{G}_i)}{\nu}\right) - \frac{\text{tr}(\mathbf{B}^j \mathbf{G}_i)}{\nu + h_i + \text{tr}(\mathbf{B}^j \mathbf{G}_i)}$ is a constant that does not depend on \mathbf{B} .

Hence, to find an update \mathbf{B}^{l+1} , we iterate the solution of the following upper-bounded convex subproblem:

$$\begin{aligned}
 \mathbf{B}^{j+1} = \underset{\mathbf{B}}{\arg \min} && -\log \det(\text{Diag}(\mathbf{B}^{\top} \mathbf{1}_r) + \mathbf{J}_{qq} - (\mathbf{B} - \mathbf{J}_{rq})^{\top} (\mathbf{I}_r + \mathbf{J}_{rr})^{-1} (\mathbf{B} - \mathbf{J}_{rq})) \\
 && + \text{tr}(\mathbf{B} \mathbf{M}^j), \text{ subject to } \mathbf{B} \geq \mathbf{0}, \mathbf{B} \mathbf{1}_q = \mathbf{1}_r,
 \end{aligned} \tag{16}$$

where $M^j \triangleq \frac{p + \nu}{n} \sum_{i=1}^n \frac{G_i}{\nu + h_i + \text{tr}(\mathbf{B}^j \mathbf{G}_i)}$. We observe in practice that a few (≈ 5) iterations are sufficient for convergence.

Problem (16) has the same format as Problem (8) and it can be readily solved by Algorithm 1. The proposed learning scheme for Student- t bipartite graphs is summarized in Algorithm 2. The sequence $\{\mathbf{B}^l\}$ generated by Algorithm 2 converges to the set of stationary points of Problem (15). This result directly follows from the convergence results of the MM framework in [40].

Algorithm 2: Student- t bipartite graph learning (SBG)

Data: Data matrix \mathbf{X} , initial feasible estimate of the graph weights \mathbf{B}^0 , initial learning rate $\alpha_0 > 0$, tolerance $\epsilon > 0$, degree of freedoms $\nu > 2$.

Result: Learned bipartite graph: \mathbf{B}^{l+1}

```

1 while  $l \leq \text{maxiter}$  do
2   ▷ update  $\mathbf{B}^{l+1}$  by iterating (16)
3   if  $\|\mathbf{B}^{l+1} - \mathbf{B}^l\|_{\text{F}} / \|\mathbf{B}^l\|_{\text{F}} \leq \epsilon$  then
4     | return  $\mathbf{B}^{l+1}$ 
5   end
6 end

```

3.3 k -component Student- t Bipartite Graphs

Graphs with multiple components can be learned by introducing rank constraints on the Laplacian matrix [11, 12, 14]. More precisely, the number of zero eigenvalues of the Laplacian matrix amounts to the number of graph components, *i.e.*, $\text{rank}(\mathbf{L}) = p - k$, where k is the number of components.

The approach taken by [12] approximates the rank constraint by applying Fan's theorem [41] as follows

$$\sum_{i=1}^k \lambda_i(\mathbf{L}) = \underset{\mathbf{V} \in \mathbb{R}^{p \times k}, \mathbf{V}^\top \mathbf{V} = \mathbf{I}_k}{\text{minimize}} \text{tr}(\mathbf{V}^\top \mathbf{L} \mathbf{V}), \quad (17)$$

where the right hand side of (17) can be used as a regularization term (*cf.* Problem (3)). On the other hand, the authors in [14] introduced the quadratic relaxation $\|\mathbf{L} - \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top\|_{\text{F}}^2$. Those formulations [12, 14] can be conveniently solved by coordinate descent methods. However, they introduce additional variables to be optimized, extra hyperparameters to be tuned, and do not guarantee that the rank constraint on \mathbf{L} is satisfied.

In our proposed method, we directly apply the rank constraint into the optimization formulation. More precisely, using the Student- t case, we formulate the following optimization problem to learn a k -component bipartite graph:

$$\begin{aligned} & \underset{\mathbf{L} \geq \mathbf{0}, \mathbf{B}}{\text{minimize}} && \frac{p + \nu}{n} \sum_{i=1}^n \log \left(1 + \frac{h_i + \text{tr}(\mathbf{B} \mathbf{G}_i)}{\nu} \right) - \log \det^*(\mathbf{L}), \\ & \text{subject to} && \mathbf{L} = \begin{bmatrix} \mathbf{I}_r & -\mathbf{B} \\ -\mathbf{B}^\top & \text{Diag}(\mathbf{B}^\top \mathbf{1}_r) \end{bmatrix}, \text{rank}(\mathbf{L}) = p - k, \mathbf{B} \geq \mathbf{0}, \mathbf{B} \mathbf{1}_q = \mathbf{1}_r. \end{aligned} \quad (18)$$

Unlike the case for connected bipartite graphs in Problem (6), the pseudo-determinant term in (18) cannot be simplified through the block matrix determinant lemma [34]. Therefore, we design an iterative algorithm based on the ADMM framework [42] for Problem (18).

3.3.1 ADMM Solution

The augmented Lagrangian of Problem (18) can be written as

$$\begin{aligned} L_\rho(\mathbf{L}, \mathbf{B}, \mathbf{Y}) &= \frac{p + \nu}{n} \sum_{i=1}^n \log \left(1 + \frac{h_i + \text{tr}(\mathbf{B} \mathbf{G}_i)}{\nu} \right) - \log \det^*(\mathbf{L}) \\ &+ \left\langle \mathbf{L} - \begin{bmatrix} \mathbf{I}_r & -\mathbf{B} \\ -\mathbf{B}^\top & \text{Diag}(\mathbf{B}^\top \mathbf{1}_r) \end{bmatrix}, \mathbf{Y} \right\rangle + \frac{\rho}{2} \left\| \mathbf{L} - \begin{bmatrix} \mathbf{I}_r & -\mathbf{B} \\ -\mathbf{B}^\top & \text{Diag}(\mathbf{B}^\top \mathbf{1}_r) \end{bmatrix} \right\|_{\text{F}}^2, \end{aligned} \quad (19)$$

where $\rho > 0$ is a penalty hyperparameter.

For fixed \mathbf{B} and \mathbf{Y} , the subproblem for \mathbf{L} can be written as

$$\mathbf{L}^* = \arg \min_{\text{rank}(\mathbf{L})=p-k} -\log \det^*(\mathbf{L}) + \langle \mathbf{L}, \mathbf{Y} \rangle + \frac{\rho}{2} \left\| \mathbf{L} - \begin{bmatrix} \mathbf{I}_r & -\mathbf{B} \\ -\mathbf{B}^\top & \text{Diag}(\mathbf{B}^\top \mathbf{1}_r) \end{bmatrix} \right\|_{\text{F}}^2, \quad (20)$$

whose closed-form solution is given by Lemma 1.

Lemma 1 *The global minimizer of problem (20) is [43, 44]*

$$\mathbf{L}^* = \frac{1}{2\rho} \mathbf{R} \left(\Gamma + \sqrt{\Gamma^2 + 4\rho \mathbf{I}} \right) \mathbf{R}^\top, \quad (21)$$

where $\mathbf{R}\Gamma\mathbf{R}^\top$ is the eigenvalue decomposition of $\rho \begin{bmatrix} \mathbf{I}_r & -\mathbf{B} \\ -\mathbf{B}^\top & \text{Diag}(\mathbf{B}^\top \mathbf{1}_r) \end{bmatrix} - \mathbf{Y}$, with Γ having the largest $p - k$ eigenvalues along its diagonal and $\mathbf{R} \in \mathbb{R}^{p \times (p-k)}$ containing the corresponding eigenvectors.

Fixing \mathbf{L} and \mathbf{Y} , the subproblem for \mathbf{B} is

$$\mathbf{B}^* = \arg \min_{\mathbf{B} \geq 0, \mathbf{B}\mathbf{1}=\mathbf{1}} \frac{p+\nu}{n} \sum_{i=1}^n \log \left(1 + \frac{h_i + \text{tr}(\mathbf{B}\mathbf{G}_i)}{\nu} \right) + \frac{\rho}{2} \left\| \mathbf{L} - \begin{bmatrix} \mathbf{I}_r & -\mathbf{B} \\ -\mathbf{B}^\top & \text{Diag}(\mathbf{B}^\top \mathbf{1}_r) \end{bmatrix} \right\|_{\text{F}}^2 - \left\langle \begin{bmatrix} \mathbf{I}_r & -\mathbf{B} \\ -\mathbf{B}^\top & \text{Diag}(\mathbf{B}^\top \mathbf{1}_r) \end{bmatrix}, \mathbf{Y} \right\rangle, \quad (22)$$

which can be simplified as

$$\mathbf{B}^* = \arg \min_{\mathbf{B} \geq 0, \mathbf{B}\mathbf{1}=\mathbf{1}} \frac{p+\nu}{n} \sum_{i=1}^n \log \left(1 + \frac{h_i + \text{tr}(\mathbf{B}\mathbf{G}_i)}{\nu} \right) + \text{tr}(\mathbf{B}\mathbf{H}) + \rho \|\mathbf{B}\|_{\text{F}}^2 + \frac{\rho}{2} \mathbf{1}_r^\top \mathbf{B} \mathbf{B}^\top \mathbf{1}_r \quad (23)$$

where $\mathbf{H} \triangleq 2\mathbf{Y}_{r,q}^\top - \mathbf{y}_{r+1:p} \mathbf{1}_r^\top - \rho (\mathbf{l}_{r+1:p} \mathbf{1}_r^\top - 2\mathbf{L}_{r,q}^\top)$, $\mathbf{y}_{r+1:p} \triangleq \text{diag}(\mathbf{Y})_{r+1:p}$, and $\mathbf{l}_{r+1:p} \triangleq \text{diag}(\mathbf{L})_{r+1:p}$.

Leveraging the MM framework to Problem (23), like it was applied to Problem (15), we can find an upperbounded subproblem at point \mathbf{B}^j as follows:

$$\mathbf{B}^{j+1} = \arg \min_{\mathbf{B} \geq 0, \mathbf{B}\mathbf{1}=\mathbf{1}} \text{tr}(\mathbf{B}(\mathbf{H} + \mathbf{M}^j)) + \rho \|\mathbf{B}\|_{\text{F}}^2 + \frac{\rho}{2} \mathbf{1}_r^\top \mathbf{B} \mathbf{B}^\top \mathbf{1}_r. \quad (24)$$

Subproblem (24) is strongly convex, thus it can be solved efficiently via convex solvers. Since the objective function of (24) is smooth and its feasible set is compact, we provide a PGD algorithm in the Supplementary Material. Finally, to find an update \mathbf{B}^* , we iterate the solution of Problem (24). The dual variable \mathbf{Y} is updated approximately via gradient ascent.

Algorithm 3 summarizes our proposed scheme to learn k -component bipartite graphs and its convergence is established through Theorem 2, whose proof is presented in the Supplementary Material.

Theorem 2 *Algorithm 3 converges subsequently for any sufficiently large ρ , that is, the sequence $\{(\mathbf{L}^l, \mathbf{B}^l, \mathbf{Y}^l)\}$ generated by Algorithm 3 has at least one limit point, and each limit point is a stationary point of (18).*

4 Experimental Results

We conduct experiments to illustrate the performance of the proposed algorithms in terms of quantitative measures such as node label accuracy and graph modularity. Accuracy is computed as the ratio between the number of correctly predicted node labels and the number of nodes in the objects set, whereas modularity measures the strength of division of a graph into groups [45]. A high modularity value means that the nodes from the same group are more likely to be connected.

Algorithm 3: Student- t k -component bipartite graph learning (kSBG)

Data: Data matrix \mathbf{X} , number of components k , initial feasible estimate of the graph weights \mathbf{B}^0 , penalty hyperparameter $\rho > 0$, tolerance $\epsilon > 0$, degrees of freedom $\nu > 2$.

Result: Learned bipartite graph: \mathbf{B}^{l+1}

```
1 while  $l \leq \text{maxiter}$  do
2   ▷ update  $\mathbf{L}^{l+1}$  via (21)
3   ▷ update  $\mathbf{B}^{l+1}$  by iterating (24)
4   ▷ update  $\mathbf{Y}^{l+1} = \mathbf{Y}^l - \rho \left( \mathbf{L}^{l+1} - \begin{bmatrix} \mathbf{I}_r & -\mathbf{B}^{l+1} \\ -\mathbf{B}^{l+1\top} & \text{Diag}(\mathbf{B}^{l+1\top} \mathbf{1}_r) \end{bmatrix} \right)$ 
5   if  $\|\mathbf{B}^{l+1} - \mathbf{B}^l\|_F / \|\mathbf{B}^l\|_F \leq \epsilon$  then
6     | return  $\mathbf{B}^{l+1}$ 
7   end
8 end
```

Benchmarks: We compare the proposed algorithms with state-of-the-art methods in two settings: (i) connected bipartite graphs, where SOB G ($k = 1$) [12] and SGA [14] are the benchmarks, and (ii) k -component bipartite graphs, where SOB G ($k > 1$) and SGLA [14] act as benchmarks. In our Algorithm 3, we set the hyperparameter $\rho = 1$ and the relative tolerance $\epsilon = 10^{-5}$. The hyperparameter of SOB G was tuned according to the heuristic procedure described in [12]. The proposed algorithms and benchmarks were implemented using the R programming language. For SGA and SGLA, we relied on their official implementation via the package spectralGraphTopology [46].

4.1 Returns of S&P500 Stocks

We perform experiments using publicly available, historical daily prices, queried from Yahoo! FinanceTM, of $r = 333$ stocks listed in the S&P500 Index and $q = 8$ sectors indexes, namely: Consumer Staples, Energy, Financials, Health Care, Industrials, Materials, Real Estate, and Utilities. Following our notation, the node set of classes \mathcal{V}_q represents the sectors indexes, whereas the stocks are represented by the node set of objects \mathcal{V}_r .

Stock sector labels are available through the Global Industry Classification Standard (GICS) [47, 48]. Arguably the most well-recognized, industry-standard stock classification system, GICS relies on a fixed classification structure that may not fully capture the actual impact of a particular stock at particular time. With the intent of solving this shortcoming, we employ the learned bipartite graphs to extend current industry stock classification system from a static classification to a dynamic soft clustering, in which a particular node (stock) may belong to different clusters (sectors) with different degrees of membership that can change over time. Our goal is to perform soft clustering of stocks according to the market sectors. To that end, we recall that each row of \mathbf{B} is a point belonging to the probability simplex in \mathbb{R}^q , which represents the level of membership of a stock to a particular set of sectors. The final sector assigned to the i -th stock corresponds to $\arg \max_{j \in \{1, \dots, q\}} B_{ij}$.

We start by constructing the log-returns data matrix, *i.e.*, a matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$, where n is the number of log-returns observations and p is the number of instruments, as $X_{i,j} = \log P_{i,j} - \log P_{i-1,j}$, where $P_{i,j}$ is the closing price of the j -th instrument at the i -th day. We collect data from Jan 5th 2016 to Jan 5th 2021, totalling $n = 1291$ daily observations. Then, we proceed to learn connected bipartite graphs via the benchmarks and proposed methods using the observations \mathbf{X} in a rolling-window fashion¹. More precisely, we use a window of length 504 days (2 years in terms of stock market days) and shift this window by 63 days (3 months in terms of stock market days).

Figure 1 shows quantitative measurements that reveal the higher performance of the proposed SBG algorithm both in terms of accuracy and modularity. That may be explained by the fact that SBG adopts a multivariate Student- t distribution, which is well-known to better model financial datasets in comparison to Gaussian assumptions.

¹An estimate for the degree of freedoms ν , required by Algorithms 2 and 3, is obtained by fitting a univariate Student- t distribution to the log-returns of the S&P500 index during the corresponding time period.

Furthermore, we evaluate the performance of the k -component bipartite estimators under the same settings as previously described, where we set $k = 8$, *i.e.*, k is the number of sectors. Figure 2 reveals that kSBG presents a higher performance across the whole time period both in terms of accuracy and modularity when compared to SGLA and SOBG. Interestingly, in both Figure 1 and 2, we observe a decline in accuracy and modularity of the methods GBG, SBG, and kSBG, starting in Jan 2020, which can be explained by the impact of the COVID-19 pandemic on the US stock market. This reveals that our methods may be able to identify periods of economic turmoil and suggests that, during such times, standards such as GICS may not be the best reference for stock market sector classification. This insight may be critical for investors who rely on diversification of their portfolios through information on the sectors [49].

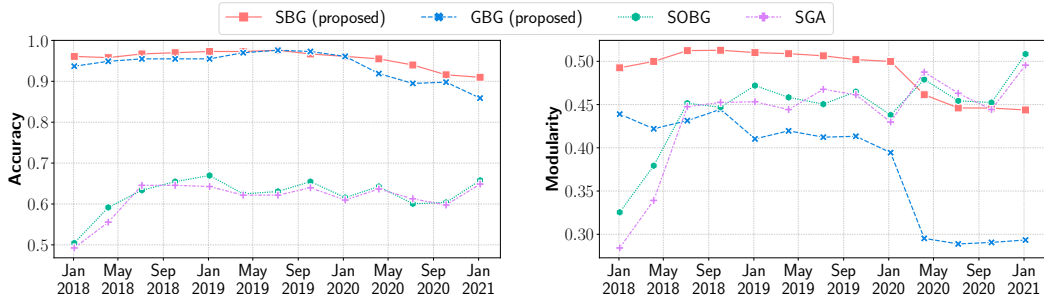


Figure 1: Performance of the estimators for connected bipartite graphs of S&P500 stocks.

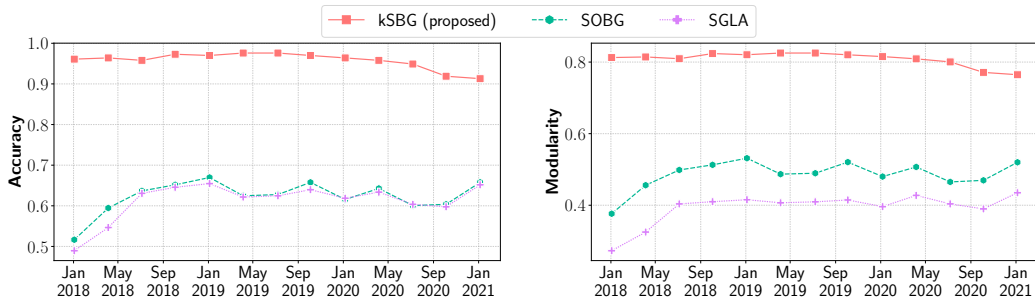


Figure 2: Performance of the estimators for 8-component bipartite graphs of S&P500 stocks.

Finally, Figure 3 shows the k -component bipartite models learned using the whole dataset. We observe that SGLA allows the existence of isolated nodes, which degrades performance. Our proposed method kSBG avoids that issue by imposing a linear constraint on the node degrees. In addition, kSBG yields a more intuitive representation of the financial network than SOBG as verified by the accuracy and modularity values.

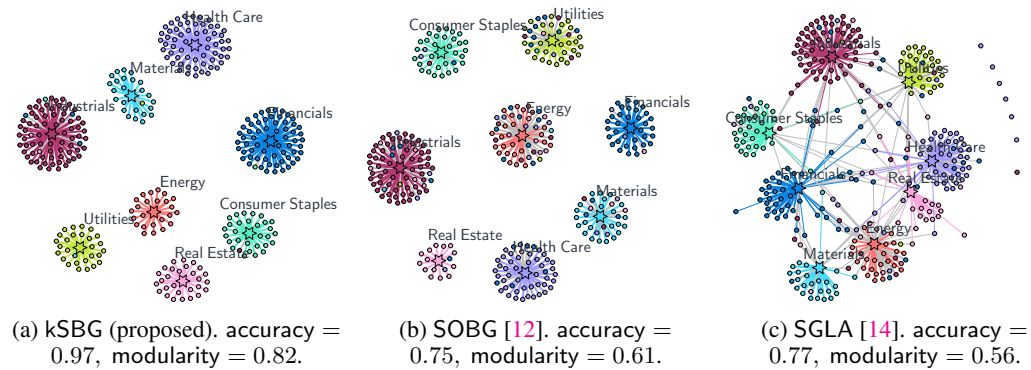


Figure 3: Graphs learned by different algorithms in the S&P500 stock returns dataset.

4.2 Robustness to Initialization

The performance of iterative algorithms may be affected by their initial point, especially when dealing with nonconvex problems as in Algorithms 2 and 3. Hence, we perform an experiment to compare two initialization schemes: (i) uniform, the graph weights are randomly sampled, *i.e.*, $B_{ij}^0 \sim \text{Uniform}(0, 1)$, and (ii) default, $B^0 \propto \max(\mathbf{0}, -\Xi_{rq})$, where $S^{-1} = [\Xi_{rr} \ \Xi_{rq}; \Xi_{rq}^\top \ \Xi_{qq}]$ is the inverse of the sample correlation matrix, as discussed in more details in [14]. In both schemes, we normalize the rows of B^0 such that $B^0 \mathbf{1} = \mathbf{1}$ holds. For this experiment, we consider log-returns of $r = 362$ stocks and $q = 9$ stock sectors from Oct. 5th 2005 to Dec. 30th 2015, totalling $n = 2577$ observations. Figure 4 shows the result of 100 realizations of Algorithms 2 and 3. It can be observed that Algorithms 2 and 3 are not sensitive to the choice of an initial point, which is an important feature for practical applications. Finally, Figures 5a and 5b show the learned graphs from kSBG with different initial points, where we observe that the performance of kSBG is largely unaffected.

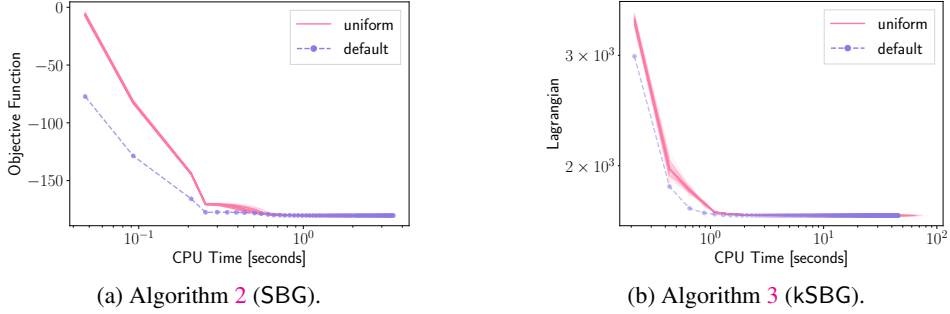


Figure 4: Convergence trend of the proposed algorithms for different initial points.

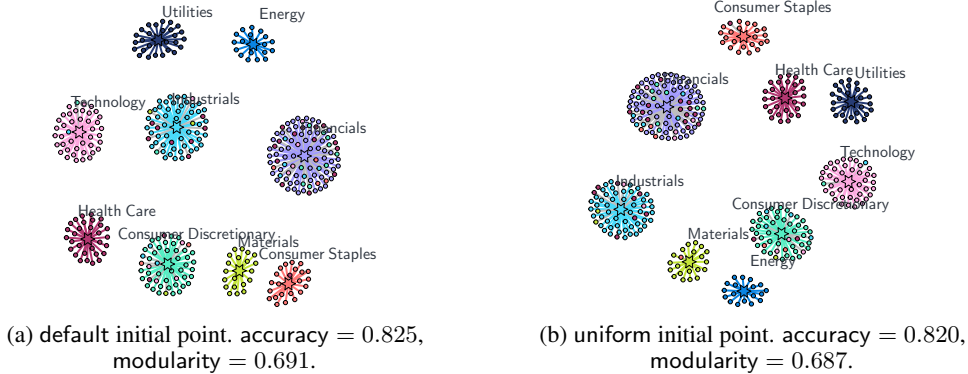


Figure 5: Visualizations of the graphs learned by kSBG with distinct initial points.

5 Conclusions

The design of numerical algorithms plays a critical role in disseminating the application of graphical models in practical, real-world instances. This paper has proposed an optimization formulation to the problem of learning a bipartite graph from data along with an algorithm based on the PGD framework. Envisioning applications in scenarios where outliers or heavy-tail events are present, we have extended the proposed formulation to the case where the data follow a multivariate Student- t distribution. In addition, we have proposed an optimization formulation, along with a provably convergent ADMM-based algorithm, for estimating k -component bipartite graphs. Experimental results with datasets of log-returns of S&P500 stocks have provided substantial evidence for the improved performance of the proposed algorithms in comparison to state-of-the-art methods.

Acknowledgments

This work was supported by the Hong Kong GRF 16207820 research grant.

References

- [1] S. M. Smith, K. L. Miller, G. Salimi-Khorshidi, M. Webster, C. F. Beckmann, T. E. Nichols, J. D. Ramsey, and M. W. Woolrich. Network modelling methods for FMRI. *NeuroImage*, 54(2):875–891, 2011.
- [2] O. Stegle, S. A. Teichmann, and J. C. Marioni. Computational and analytical challenges in single-cell transcriptomics. *Nature Reviews Genetics*, 16:133–145, 2015.
- [3] D. Hallac, Y. Park, S. Boyd, and J. Leskovec. Network inference via the time-varying graphical lasso. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '17, pages 205–213, 2017.
- [4] D. Hallac, S. Vare, S. Boyd, and J. Leskovec. Toeplitz inverse covariance-based clustering of multivariate time series data. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'17)*, page 215–223, 2017.
- [5] A. Jung, G. Hannak, and N. Goertz. Graphical lasso based model selection for time series. *IEEE Signal Processing Letters*, 22(10):1781–1785, 2015.
- [6] J. V. M. Cardoso and D. P. Palomar. Learning undirected graphs in financial markets. In *54th Annual Asilomar Conference on Signals, Systems, and Computers*, 2020.
- [7] J. V. M. Cardoso, J. Ying, and D. P. Palomar. Graphical models in heavy-tailed markets. In *Advances in Neural Information Processing Systems (NeurIPS'21)*, 2021.
- [8] J. Friedman, T. Hastie, and R. Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9:432–41, 2008.
- [9] O. Banerjee, L. E. Ghaoui, and A. d'Aspremont. Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data. *Journal of Machine Learning Research*, 9:485–516, 2008.
- [10] S. J. Wright. Coordinate descent algorithms. *Mathematical Programming*, 151:3–34, 2015.
- [11] F. Nie, X. Wang, M. I. Jordan, and H. Huang. The constrained Laplacian rank algorithm for graph-based clustering. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, AAAI'16, pages 1969–1976. AAAI Press, 2016.
- [12] F. Nie, X. Wang, C. Deng, and H. Huang. Learning a structured optimal bipartite graph for co-clustering. In *Advances on Neural Information Processing Systems (NeurIPS'17)*, page 4132–4141, 2017.
- [13] S. Kumar, J. Ying, J. V. M. Cardoso, and D. P. Palomar. Structured graph learning via Laplacian spectral constraints. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [14] S. Kumar, J. Ying, J. V. M. Cardoso, and D. P. Palomar. A unified framework for structured graph learning via spectral constraints. *Journal of Machine Learning Research*, 21:1–60, 2020.
- [15] M. Finegold and M. Drton. Robust graphical modeling with t -distributions. In *Conference on Uncertainty in Artificial Intelligence*, 2009.
- [16] Y. Wald, N. Noy, G. Elidan, and A. Wiesel. Globally optimal learning for structured elliptical losses. In *Advances in Neural Information Processing Systems (NeurIPS'19)*, 2019.
- [17] M. Slawski and M. Hein. Estimation of positive definite M -matrices and structure learning for attractive Gaussian Markov random fields. *Linear Algebra and its Applications*, pages 145–179, 2015.
- [18] R. Agrawal, U. Roy, and C. Uhler. Covariance matrix estimation under total positivity for portfolio selection. *Journal of Financial Econometrics*, 20:367–389, 2020.
- [19] Y. Wang, U. Roy, and C. Uhler. Learning high-dimensional Gaussian graphical models under total positivity without adjustment of tuning parameters. In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (AISTATS'20)*, volume 108, pages 2698–2708, 2020.
- [20] J.-F. Cai, J. V. M. Cardoso, D. P. Palomar, and J. Ying. Fast projected newton-like method for precision matrix estimation under total positivity. *arXiv e-prints: 2112.01939*, August 2022.
- [21] J. Ying, J. V. M. Cardoso, and D. P. Palomar. Does the ℓ_1 -norm learn a sparse graph under Laplacian constrained graphical models? *arXiv e-prints: 2006.14925*, June 2020.

- [22] J. Ying, J. V. M. Cardoso, and D. P. Palomar. Nonconvex sparse graph learning under Laplacian-structured graphical model. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, pages 7101–7113, 2020.
- [23] J. Ying, J. V. M. Cardoso, and D. Palomar. Minimax estimation of laplacian constrained precision matrices. In *24th International Conference on Artificial Intelligence and Statistics*, 2021.
- [24] C. Hu, L. Cheng, J. Sepulcre, K. A. Johnson, G. E. Fakhri, Y. M. Lu, and Q. Li. A spectral graph regression model for learning brain connectivity of alzheimer’s disease. *PLOS ONE*, 10(5):1–24, 05 2015.
- [25] X. Dong, D. Thanou, P. Frossard, and P. Vandergheynst. Learning Laplacian matrix in smooth graph signal representations. *IEEE Transactions on Signal Processing*, 64(23):6160–6173, 2016.
- [26] V. Kalofolias. How to learn a graph from smooth signals. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, volume 51, pages 920–929. PMLR, 2016.
- [27] G. Mateos, S. Segarra, A. G. Marques, and A. Ribeiro. Connecting the dots: Identifying network structure via graph signal processing. *IEEE Signal Processing Magazine*, 36(3):16–43, 2019.
- [28] H. Rue and L. Held. *Gaussian Markov Random Fields: Theory And Applications (Monographs on Statistics and Applied Probability)*. Chapman & Hall/CRC, 2005.
- [29] H. E. Egilmez, E. Pavez, and A. Ortega. Graph learning from data under Laplacian and structural constraints. *IEEE Journal of Selected Topics in Signal Processing*, 11(6):825–841, 2017.
- [30] O. Knill. Cauchy–Binet for pseudo-determinants. *Linear Algebra and its Applications*, 459:522–547, 2014.
- [31] F. Nie, X. Wang, and H. Huang. Clustering and projected clustering with adaptive neighbors. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD’14)*, 2014.
- [32] D. P. Bertsekas. *Nonlinear Programming*. Athena Scientific, 1999.
- [33] A. Beck. *First-Order Methods in Optimization*. MOS-SIAM Series on Optimization. Society for Industrial and Applied Mathematics (SIAM), 2017.
- [34] F. Zhang. *The Schur Complement and Its Applications*. Springer, 2005.
- [35] C. Michelot. A finite algorithm for finding the projection of a point onto the canonical simplex of \mathbb{R}^n . *Journal of Optimization Theory and Applications*, 50(1):195–200, July 1986.
- [36] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [37] D. P. Palomar and J. R. Fonollosa. Practical algorithms for a family of waterfilling solutions. *IEEE Transactions on Signal Processing*, 53(2):686–695, 2005.
- [38] L. Condat. Fast projection onto the simplex and the ℓ_1 ball. *Mathematical Programming*, 158, 2016.
- [39] S. I. Resnick. *Heavy-Tail Phenomena: Probabilistic and Statistical Modeling*. Springer-Verlag New York, 2007.
- [40] Y. Sun, P. Babu, and D. P. Palomar. Majorization-minimization algorithms in signal processing, communications, and machine learning. *IEEE Transactions on Signal Processing*, 65(3):794–816, 2017.
- [41] K. Fan. On a theorem of Weyl concerning eigenvalues of linear transformations I. *Proceedings of the National Academy of Sciences*, 35(11):652–655, 1949.
- [42] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning*, 3(1):1–122, 2011.
- [43] D. M. Witten and R. Tibshirani. Covariance-regularized regression and classification for high dimensional problems. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 71(3):615–636, 2009.

- [44] P. Danaher, P. Wang, and D. M. Witten. The joint graphical lasso for inverse covariance estimation across multiple classes. *Journal of the Royal Statistical Society Series B*, 76(2):373–397, 2014.
- [45] M. E. J. Newman. Modularity and community structure in networks. *Proceedings of the National Academy of Sciences of the United States of America*, 103, 2006.
- [46] J. V. M. Cardoso and D. Palomar. spectralGraphTopology: Learning graphs from data via spectral constraints. In *The Comprehensive R Archive Network*. <https://CRAN.R-project.org/package=spectralGraphTopology>, 2019.
- [47] Standard & Poor’s. Global Industry Classification Standard (GICS). *Tech Report*, 2006.
- [48] Morgan Stanley Capital International and S&P Dow Jones. Revisions to the global industry classification standard (GICS) structure, 2018.
- [49] Y. Yang, L. Zhao, L. Chen, C. Wang, and J. Han. Portfolio optimization with idiosyncratic and systemic risks for financial networks. *arXiv e-prints: 2111.11286*, 2021.
- [50] Y. Wang, W. Yin, and J. Zeng. Global convergence of admm in nonconvex nonsmooth optimization. *Journal of Scientific Computing*, 78:29–63, 2019.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes]
 - (c) Did you discuss any potential negative societal impacts of your work? [N/A]
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes]
 - (b) Did you include complete proofs of all theoretical results? [Yes] See Supplementary Material.
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] The code, data, and instructions to reproduce the experiments are available in <https://github.com/mirca/bipartite>.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes]
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes]
 - (b) Did you mention the license of the assets? [Yes]
 - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [N/A] The data is publicly available.
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A] The data does not contain personally identifiable information or offensive content.
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]

- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]