# Learning Graphical Models With Hubs

Kean Ming Tan, Palma London, Karthik Mohan,
Su-In Lee, Maryam Fazel, Daniela Witten

**Abstract**

We consider the problem of learning a high-dimensional graphical model in which certain *hub* nodes are highly-connected to many other nodes. Many authors have studied the use of an $\ell_1$ penalty in order to learn a sparse graph in the high-dimensional setting. However, the $\ell_1$ penalty implicitly assumes that each edge is equally likely and independent of all other edges. We propose a general framework to accommodate more realistic networks with hub nodes, using a convex formulation that involves a row-column overlap norm penalty. We apply this general framework to three widely-used probabilistic graphical models: the Gaussian graphical model, the covariance graph model, and the binary Ising model. An alternating direction method of multipliers algorithm is used to solve the corresponding convex optimization problems. On synthetic data, we demonstrate that our proposed framework outperforms competitors that do not explicitly model hub nodes. We illustrate our proposal on a webpage data set and a gene expression data set.

Keywords: Gaussian graphical model, covariance graph, binary network, *lasso*, hub

## 1 Introduction

Graphical models are used to model a variety of biological and social processes, such as gene regulatory networks and social interaction networks. A graph consists of a set of $p$ nodes, each representing a variable, and a set of edges between pairs of nodes. The presence of an edge between two nodes indicates a relationship between the two variables. In this manuscript, we consider two types of graphs: conditional independence graphs and marginal independence graphs. In a conditional independence graph, an edge connects a pair of variables if and only if they are conditionally dependent — dependent conditional upon the other variables. In a marginal independence graph, two nodes are joined by an edge if and only if they are marginally dependent — dependent without conditioning on the other variables.

In recent years, many authors have studied the problem of learning a graphical model in the high-dimensional setting, in which the number of variables $p$ is larger than the number of observations $n$. Let $\mathbf{X}$ be a $n \times p$ matrix, with rows $\mathbf{x}_1, \ldots, \mathbf{x}_n$. Throughout the rest of the text, we will focus on three specific types of graphical models:

1. A *Gaussian graphical model*, where $\mathbf{x}_1, \ldots, \mathbf{x}_n \overset{\text{i.i.d.}}{\sim} N(\mathbf{0}, \mathbf{\Sigma})$. It is well-known that $(\mathbf{\Sigma}^{-1})_{jj'} = 0$ for some $j \neq j'$ if and only if the $j$th and $j'$th variables are conditionally

independent (Mardia et al. 1979); therefore, the sparsity pattern of $\mathbf{\Sigma}^{-1}$ determines the conditional independence graph.

2. A *Gaussian covariance graph model*, where $\mathbf{x}_1, \ldots, \mathbf{x}_n \overset{\text{i.i.d.}}{\sim} N(\mathbf{0}, \mathbf{\Sigma})$. Then $\Sigma_{jj'} = 0$ for some $j \neq j'$ if and only if the $j$th and $j'$th variables are marginally independent. Therefore, the sparsity pattern of $\mathbf{\Sigma}$ determines the marginal independence graph.

3. A *binary Ising graphical model*, where $\mathbf{x}_1, \ldots, \mathbf{x}_n \overset{\text{i.i.d.}}{\sim} p(\mathbf{x}, \mathbf{\Theta})$,

$$p(\mathbf{x}, \mathbf{\Theta}) = \frac{1}{Z(\mathbf{\Theta})} \exp \left[ \sum_{j=1}^{p} \theta_{jj} x_j + \sum_{1 \leq j < j' \leq p} \theta_{jj'} x_j x_{j'} \right],$$

$\mathbf{\Theta}$ is a $p \times p$ symmetric matrix, and $Z(\mathbf{\Theta})$ is the partition function, which ensures that the density sums to one. Here, $\mathbf{x}$ is a binary vector, and $\theta_{jj'} = 0$ if and only if the $j$th and $j'$th variables are conditionally independent. The sparsity pattern of $\mathbf{\Theta}$ determines the conditional independence graph.

To construct an interpretable graph when $p > n$, many authors have proposed applying an $\ell_1$ penalty to the parameter encoding each edge, in order to encourage sparsity. For instance, such an approach is taken by Yuan & Lin (2007*b*), Friedman et al. (2007), Rothman et al. (2008), and Yuan (2008) in the Gaussian graphical model; El Karoui (2008), Bickel & Levina (2008), Rothman et al. (2009), Bien & Tibshirani (2011), Cai & Liu (2011), and Xue et al. (2012) in the covariance graph model; and Lee et al. (2007), Höfling & Tibshirani (2009), and Ravikumar et al. (2010) in the binary model.

However, applying an $\ell_1$ penalty to each edge can be interpreted as an independent double-exponential prior on each edge. Consequently, such an approach implicitly assumes that each edge is equally likely and independent of all other edges; this corresponds to an Erdős-Rényi graph in which most nodes have approximately the same number of edges (Erdős & Rényi 1959). This is unrealistic in many real-world networks, in which we believe that certain nodes (which, unfortunately, are not known *a priori*) have many more edges than other nodes. An example is the network of webpages in the World Wide Web, where a relatively small number of webpages are highly connected to many other webpages (Barabási & Albert 1999). Many authors have shown that real-world networks are *scale-free*, in the sense that the number of edges for each node follows a power-law distribution; examples include gene-regulatory networks, social networks, and networks of collaborations among scientists (among others, Barabási & Albert 1999, Barabási 2009, Liljeros et al. 2001, Jeong et al. 2001, Newman 2000). More recently, Hao et al. (2012) have shown that certain genes, referred to as *super hubs*, regulate hundreds of downstream genes in a gene regulatory network, resulting in far denser connections than are typically seen in a scale-free network.

In this paper, we will be referring to nodes that are very densely-connected as hubs; this is a somewhat informal definition. The word *hub* has been used in a variety of contexts in the literature; for instance, some authors have used it to refer to relatively highly-connected nodes in a scale-free network. However, when we refer to hubs in this paper, we have in mind nodes that are *extremely* densely-connected to other nodes, such as the super hubs discussed by Hao et al. (2012), or webpages like Google in the World Wide Web. In many types of graphs, the hubs themselves are of interest.

Here we propose a hub penalty function for estimating graphs containing hubs. Our formulation simultaneously identifies the hubs and estimates the entire graph. The proposed penalty function is convex, and so yields a convex optimization problem when combined with a convex loss function. We consider the application of this hub penalty function in modeling Gaussian graphical models, covariance graph models, and binary Ising models. Our formulation does not require that we know *a priori* which nodes in the network are hubs.

In related work, several authors have proposed methods to estimate a scale-free Gaussian graphical model (Liu & Ihler 2011, Defazio & Caetano 2012). However, those methods do not model hub nodes — the most densely-connected nodes that arise in a scale-free network are far less densely-connected than the hubs that can be modeled using our formulation. Under a different framework, Hero & Rajaratnam (2012) proposed a screening-based procedure to identify hub nodes in the context of Gaussian graphical models. Our proposal outperforms such approaches when hub nodes are present.

In Figure 1, the performance of our proposed approach is shown in a toy example in the context of a Gaussian graphical model. We see that when the true network contains hub nodes (Figure 1(a)), our proposed approach (Figure 1(b)) is much better able to recover the signal than is the graphical lasso (Figure 1(c)), a well-studied approach that applies an $\ell_1$ penalty to each edge in the graph (Friedman et al. 2007).
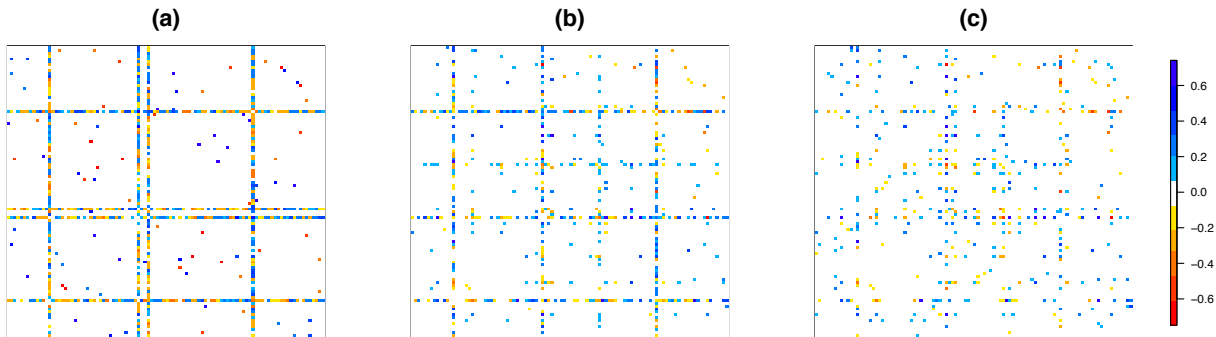


Figure 1: (a): Heatmap of the inverse covariance matrix in a toy example with four hub nodes. White elements are zero and colored elements are non-zero in the inverse covariance matrix. Thus, colored elements correspond to edges in the graph. (b): Hub graphical lasso estimate. (c): Graphical lasso estimate.

We present the hub penalty function in Section 2. We then apply it to the Gaussian graphical model, the covariance graph model, and the binary Ising model in Sections 3, 4, and 5, respectively. In Section 6, we apply our approach to a webpage data set and a gene expression data set. We close with a discussion in Section 7.

# 2 The General Formulation

## 2.1 The Hub Penalty Function

Let $\mathbf{X}$ be a $n \times p$ data matrix, $\mathbf{\Theta}$ a $p \times p$ symmetric matrix containing the parameters of interest, and $\ell(\mathbf{X}, \mathbf{\Theta})$ a loss function (assumed to be convex in $\mathbf{\Theta}$). In order to obtain a sparse and interpretable graph estimate, many authors have considered the problem

$$\underset{\mathbf{\Theta}}{\text{minimize}} \quad \{\ell(\mathbf{X}, \mathbf{\Theta}) + \lambda_1 \|\mathbf{\Theta} - \text{diag}(\mathbf{\Theta})\|_1\}, \tag{1}$$

where $\|\cdot\|_1$ is the sum of the absolute values of the matrix elements. For instance, in the case of a Gaussian graphical model, we could take $\ell(\mathbf{X}, \mathbf{\Theta}) = -\log \det \mathbf{\Theta} + \text{trace}(\mathbf{S}\mathbf{\Theta})$, the negative log-likelihood of the data, where $\mathbf{S}$ is the empirical covariance matrix. The solution to (1) can then be interpreted as an estimate of the inverse covariance matrix. The $\ell_1$ penalty in (1) encourages zeros in the solution. But it typically does not yield an estimate that contains hubs.
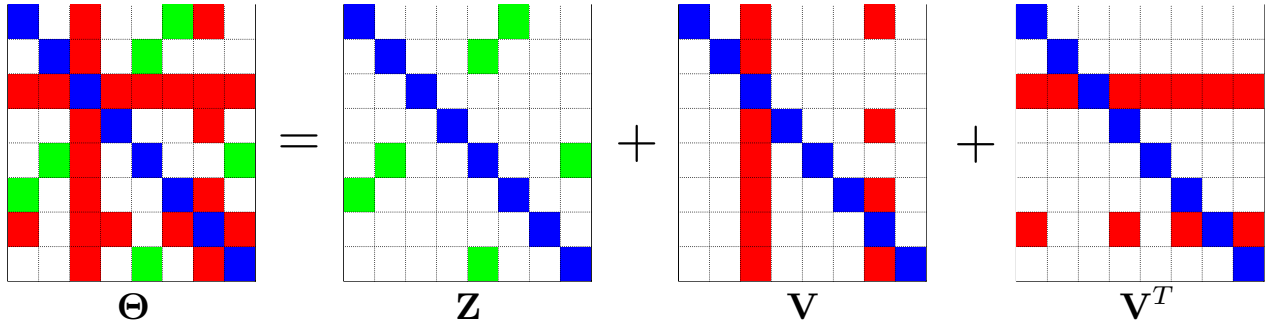


Figure 2: Decomposition of a symmetric matrix $\mathbf{\Theta}$ into $\mathbf{Z} + \mathbf{V} + \mathbf{V}^T$, where $\mathbf{Z}$ is sparse, and most columns of $\mathbf{V}$ are entirely zero. Blue, white, green, and red elements are diagonal, zero, non-zero in $\mathbf{Z}$, and non-zero due to two hubs in $\mathbf{V}$, respectively.

In order to explicitly model hub nodes in a graph, we wish to replace the $\ell_1$ penalty in (1) with a convex penalty that encourages a solution that can be decomposed as $\mathbf{Z} + \mathbf{V} + \mathbf{V}^T$, where $\mathbf{Z}$ is a sparse symmetric matrix, and $\mathbf{V}$ is a matrix whose columns are either entirely zero or almost entirely non-zero (see Figure 2). The sparse elements of $\mathbf{Z}$ represent edges between non-hub nodes, and the non-zero columns of $\mathbf{V}$ correspond to hub nodes that are highly connected to many other nodes. We achieve this goal via the *hub penalty function*, which takes the form

$$\text{P}(\mathbf{\Theta}) = \underset{\mathbf{V},\mathbf{Z}: \ \mathbf{\Theta}=\mathbf{V}+\mathbf{V}^T+\mathbf{Z}}{\min} \left\{ \lambda_1 \|\mathbf{Z} - \text{diag}(\mathbf{Z})\|_1 + \lambda_2 \|\mathbf{V} - \text{diag}(\mathbf{V})\|_1 + \lambda_3 \sum_{j=1}^{p} \|(\mathbf{V} - \text{diag}(\mathbf{V}))_j\|_q \right\}. \tag{2}$$

Here $\lambda_1, \lambda_2,$ and $\lambda_3$ are nonnegative tuning parameters. Sparsity in $\mathbf{Z}$ is encouraged via the $\ell_1$ penalty on its off-diagonal elements, and is controlled by the value of $\lambda_1$. The $\ell_1$ and $\ell_1/\ell_q$ norms on the columns of $\mathbf{V}$ induce group sparsity when $q = 2$ (Yuan & Lin 2007$a$, Simon

et al. 2012); $\lambda_3$ controls the selection of hub nodes, and $\lambda_2$ controls the sparsity of each hub node's connections to other nodes. The convex penalty (2) can be combined with $\ell(\mathbf{X}, \boldsymbol{\Theta})$ in order to yield the convex optimization problem

$$\underset{\boldsymbol{\Theta}}{\text{minimize}} \quad \{\ell(\mathbf{X}, \boldsymbol{\Theta}) + \mathrm{P}(\boldsymbol{\Theta})\}, \tag{3}$$

which can be rewritten as

$$\underset{\boldsymbol{\Theta}, \mathbf{V}, \mathbf{Z}}{\text{minimize}} \quad \left\{ \ell(\mathbf{X}, \boldsymbol{\Theta}) + \lambda_1 \|\mathbf{Z} - \mathrm{diag}(\mathbf{Z})\|_1 + \lambda_2 \|\mathbf{V} - \mathrm{diag}(\mathbf{V})\|_1 \right.$$
$$\left. + \lambda_3 \sum_{j=1}^{p} \|(\mathbf{V} - \mathrm{diag}(\mathbf{V}))_j\|_q \right\} \quad \text{subject to} \quad \boldsymbol{\Theta} = \mathbf{V} + \mathbf{V}^T + \mathbf{Z}. \tag{4}$$

Note that when $\lambda_2 \to \infty$ or $\lambda_3 \to \infty$, then (4) reduces to (1). In this paper, we take $q = 2$, which leads to estimation of a network containing dense hub nodes. Other values of $q$ such as $q = \infty$ are also possible (see e.g., Mohan et al. 2013). We note that the hub penalty function is closely related to recent work on overlapping group lasso penalties in the context of learning multiple sparse precision matrices (Mohan et al. 2013).

## 2.2   Algorithm

In order to solve (4) with $q = 2$, we use an *alternating direction method of multipliers* (ADMM) algorithm (see e.g., Eckstein & Bertsekas 1992, Boyd et al. 2010, Eckstein 2012). ADMM is an attractive algorithm for this problem, as it allows us to decouple some of the terms in (4) that are difficult to optimize jointly. In order to develop an ADMM algorithm for (4) with guaranteed convergence, we reformulate it as a consensus problem, as in Ma et al. (2013). Then the ADMM formulation involves only two blocks of primal variables, and convergence of the algorithm to the optimal solution follows from classical results (see e.g. the review papers, Boyd et al. 2010, Eckstein 2012).

In greater detail, we let $\mathbf{A} = (\boldsymbol{\Theta}, \mathbf{V}, \mathbf{Z})$, $\mathbf{B} = (\tilde{\boldsymbol{\Theta}}, \tilde{\mathbf{V}}, \tilde{\mathbf{Z}})$,

$$f(\mathbf{A}) = \ell(\mathbf{X}, \boldsymbol{\Theta}) + \lambda_1 \|\mathbf{Z} - \mathrm{diag}(\mathbf{Z})\|_1 + \lambda_2 \|\mathbf{V} - \mathrm{diag}(\mathbf{V})\|_1 + \lambda_3 \sum_{j=1}^{p} \|(\mathbf{V} - \mathrm{diag}(\mathbf{V}))\|_2,$$

and

$$g(\mathbf{B}) = \begin{cases} 0 & \text{if } \tilde{\boldsymbol{\Theta}} = \tilde{\mathbf{V}} + \tilde{\mathbf{V}}^T + \tilde{\mathbf{Z}} \\ \infty & \text{otherwise.} \end{cases}$$

Then, we can rewrite (4) as

$$\underset{\mathbf{A}, \mathbf{B}}{\text{minimize}} \ \{f(\mathbf{A}) + g(\mathbf{B})\} \qquad \text{subject to } \mathbf{A} = \mathbf{B}. \tag{5}$$

The scaled augmented Lagrangian for (5) takes the form

$$L(\mathbf{A}, \mathbf{B}, \mathbf{W}) = \ell(\mathbf{X}, \boldsymbol{\Theta}) + \lambda_1 \|\mathbf{Z} - \mathrm{diag}(\mathbf{Z})\|_1 + \lambda_2 \|\mathbf{V} - \mathrm{diag}(\mathbf{V})\|_1$$
$$+ \lambda_3 \sum_{j=1}^{p} \|(\mathbf{V} - \mathrm{diag}(\mathbf{V}))_j\|_2 + g(\mathbf{B}) + \frac{\rho}{2} \|\mathbf{A} - \mathbf{B} + \mathbf{W}\|_F^2, \tag{6}$$

5

where $\mathbf{W} = (\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3)$ is the dual variable. Note that the scaled augmented Lagrangian can be derived from the usual Lagrangian by adding a quadratic term and completing the square (Boyd et al. 2010). A general algorithm for solving (4) is provided in Algorithm 1. The derivation is in the Appendix. Note that only the update for $\mathbf{\Theta}$ (Step 2(a)i) depends on the form of the convex loss function $\ell(\mathbf{X}, \mathbf{\Theta})$.

---

**Algorithm 1** ADMM Algorithm for Solving (4).

---

1. **Initialize** the parameters:

   (a) primal variables $\mathbf{\Theta}, \mathbf{V}, \mathbf{Z}, \tilde{\mathbf{\Theta}}, \tilde{\mathbf{V}}$, and $\tilde{\mathbf{Z}}$ to the $p \times p$ identity matrix.

   (b) dual variables $\mathbf{W}_1, \mathbf{W}_2$, and $\mathbf{W}_3$ to the $p \times p$ zero matrix.

   (c) constants $\rho > 0$ and $\tau > 0$.

2. **Iterate** until the stopping criterion $\frac{\|\mathbf{\Theta}_t - \mathbf{\Theta}_{t-1}\|_F^2}{\|\mathbf{\Theta}_{t-1}\|_F^2} \leq \tau$ is met, where $\mathbf{\Theta}_t$ is the value of $\mathbf{\Theta}$ obtained at the $t$th iteration:

   (a) Update $\mathbf{\Theta}, \mathbf{V}, \mathbf{Z}$:

      i. $\mathbf{\Theta} = \arg\min_{\mathbf{\Theta}} \left\{ \ell(\mathbf{X}, \mathbf{\Theta}) + \frac{\rho}{2} \|\mathbf{\Theta} - \tilde{\mathbf{\Theta}} + \mathbf{W}_1\|_F^2 \right\}$.

      ii. $\mathbf{Z} = S(\tilde{\mathbf{Z}} - \mathbf{W}_3, \frac{\lambda_1}{\rho})$, $\mathrm{diag}(\mathbf{Z}) = \mathrm{diag}(\tilde{\mathbf{Z}} - \mathbf{W}_3)$. Here $S$ denotes the soft-thresholding operator, applied element-wise to a matrix: $S(A_{ij}, b) = \mathrm{sign}(A_{ij}) \max(|A_{ij}| - b, 0)$.

      iii. $\mathbf{C} = \tilde{\mathbf{V}} - \mathbf{W}_2 - \mathrm{diag}(\tilde{\mathbf{V}} - \mathbf{W}_2)$.

      iv. $\mathbf{V}_j = \max\left(1 - \frac{\lambda_3}{\rho\|S(\mathbf{C}_j, \lambda_2/\rho)\|_2}, 0\right) \cdot S(\mathbf{C}_j, \lambda_2/\rho)$ for $j = 1, \ldots, p$.

      v. $\mathrm{diag}(\mathbf{V}) = \mathrm{diag}(\tilde{\mathbf{V}} - \mathbf{W}_2)$.

   (b) Update $\tilde{\mathbf{\Theta}}, \tilde{\mathbf{V}}, \tilde{\mathbf{Z}}$:

      i. $\mathbf{\Gamma} = \frac{\rho}{6}\left[(\mathbf{\Theta} + \mathbf{W}_1) - (\mathbf{V} + \mathbf{W}_2) - (\mathbf{V} + \mathbf{W}_2)^T - (\mathbf{Z} + \mathbf{W}_3)\right]$.

      ii. $\tilde{\mathbf{\Theta}} = \mathbf{\Theta} + \mathbf{W}_1 - \frac{1}{\rho}\mathbf{\Gamma}$;     iii. $\tilde{\mathbf{V}} = \frac{1}{\rho}(\mathbf{\Gamma} + \mathbf{\Gamma}^T) + \mathbf{V} + \mathbf{W}_2$;     iv. $\tilde{\mathbf{Z}} = \frac{1}{\rho}\mathbf{\Gamma} + \mathbf{Z} + \mathbf{W}_3$.

   (c) Update $\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3$:

      i. $\mathbf{W}_1 = \mathbf{W}_1 + \mathbf{\Theta} - \tilde{\mathbf{\Theta}}$;     ii. $\mathbf{W}_2 = \mathbf{W}_2 + \mathbf{V} - \tilde{\mathbf{V}}$;     iii. $\mathbf{W}_3 = \mathbf{W}_3 + \mathbf{Z} - \tilde{\mathbf{Z}}$.

---

In the following sections, we consider special cases of (4) that lead to estimation of Gaussian graphical models, covariance graph models, and binary networks with hub nodes.

# 3   The Hub Graphical Lasso

Assume that $\mathbf{x}_1, \ldots, \mathbf{x}_n \overset{\text{i.i.d.}}{\sim} N(\mathbf{0}, \mathbf{\Sigma})$. The well-known graphical lasso problem (see e.g., Friedman et al. 2007) takes the form of (1) with $\ell(\mathbf{X}, \mathbf{\Theta}) = -\log\det\mathbf{\Theta} + \mathrm{trace}(\mathbf{S}\mathbf{\Theta})$, and $\mathbf{S}$

the empirical covariance matrix:

$$\underset{\boldsymbol{\Theta}}{\text{minimize}} \quad \left\{ -\log\det\boldsymbol{\Theta} + \text{trace}(\mathbf{S}\boldsymbol{\Theta}) + \lambda \sum_{j\neq j'} |\Theta_{jj'}| \right\}. \tag{7}$$

The solution to this optimization problem serves as an estimate for $\boldsymbol{\Sigma}^{-1}$, and is known as the *graphical lasso solution*. We now use the hub penalty function to extend the graphical lasso in order to accommodate hub nodes.

## 3.1 Formulation and Algorithm

We propose the *hub graphical lasso* (HGL) optimization problem, which takes the form

$$\underset{\boldsymbol{\Theta}}{\text{minimize}} \quad \left\{ -\log\det\boldsymbol{\Theta} + \text{trace}(\mathbf{S}\boldsymbol{\Theta}) + \text{P}(\boldsymbol{\Theta}) \right\}. \tag{8}$$

It encourages a solution that contains hub nodes, as well as edges that connect non-hubs (Figure 1). Problem (8) can be easily solved using Algorithm 1. The update for $\boldsymbol{\Theta}$ in Algorithm 1 (Step 2(a)i) can be derived by minimizing

$$-\log\det\boldsymbol{\Theta} + \text{trace}(\mathbf{S}\boldsymbol{\Theta}) + \frac{\rho}{2}\|\boldsymbol{\Theta} - \tilde{\boldsymbol{\Theta}} + \mathbf{W}_1\|_F^2 \tag{9}$$

with respect to $\boldsymbol{\Theta}$. This can be shown to have the solution

$$\boldsymbol{\Theta} = \frac{1}{2}\mathbf{U}\left(\mathbf{D} + \sqrt{\mathbf{D}^2 + \frac{4}{\rho}\mathbf{I}}\right)\mathbf{U}^T,$$

where $\mathbf{U}\mathbf{D}\mathbf{U}^T$ denotes the eigen-decomposition of $\tilde{\boldsymbol{\Theta}} - \mathbf{W}_1 - \frac{1}{\rho}\mathbf{S}$.

The complexity of the ADMM algorithm for HGL is $O(p^3)$ per iteration; this is the complexity of the eigen-decomposition for updating $\boldsymbol{\Theta}$. We now briefly compare the computational time for the ADMM algorithm to that of an interior point method (using the solver `Sedumi` called from `cvx`). On a 1.86 GHz Intel Core 2 Duo machine, the interior point method takes $\sim 3$ minutes, while ADMM takes only 1 second, on a data set with $p = 30$.

## 3.2 Conditions for HGL Solution to be Block Diagonal

We now present a necessary condition and a sufficient condition for the HGL solution to be block diagonal, subject to some permutation of the rows and columns. The conditions presented build upon similar results in the context of Gaussian graphical models from the recent literature (see e.g., Witten et al. 2011, Mazumder & Hastie 2012, Danaher et al. 2012, Mohan et al. 2013). Let $C_1, C_2, \ldots, C_K$ denote a partition of the $p$ features.

**Theorem 1.** *A sufficient condition for the HGL solution to be block diagonal with blocks given by $C_1, C_2, \ldots, C_K$ is that $\min\left\{\lambda_1, \frac{\lambda_2}{2}\right\} > |S_{jj'}|$ for all $j \in C_k, j' \in C_{k'}, k \neq k'$.*

**Theorem 2.** *A necessary condition for the HGL solution to be block diagonal with blocks given by $C_1, C_2, \ldots, C_K$ is that $\min\left\{\lambda_1, \frac{\lambda_2+\lambda_3}{2}\right\} > |S_{jj'}|$ for all $j \in C_k, j' \in C_{k'}, k \neq k'$.*

Theorem 1 implies that one can screen the empirical covariance matrix $\mathbf{S}$ to check if the HGL solution is block diagonal (using standard algorithms for identifying the connected components of an undirected graph; see e.g., Tarjan 1972). Suppose that the HGL solution is block diagonal with $K$ blocks, containing $p_1, \ldots, p_K$ features, and $\sum_{k=1}^{K} p_k = p$. Then, one can simply solve the HGL problem on the features within each block separately. Recall that the bottleneck of the HGL algorithm is the eigen-decomposition for updating $\mathbf{\Theta}$. The block diagonal condition leads to massive computational speed-ups for implementing the HGL algorithm: instead of computing an eigen-decomposition for a $p \times p$ matrix in each iteration of the HGL algorithm, we compute the eigen-decomposition of $K$ matrices of dimensions $p_1 \times p_1, \ldots, p_K \times p_K$. The computational complexity per-iteration is reduced from $O(p^3)$ to $\sum_{k=1}^{K} O(p_k^3)$.

We illustrate the reduction in computational time due to these results in an example with $p = 500$. Without exploiting Theorem 1, the ADMM algorithm for HGL (with a particular value of $\lambda$) takes 159 seconds; in contrast, it takes only 22 seconds when Theorem 1 is applied. The estimated precision matrix has 107 connected components, the largest of which contains 212 nodes.

## 3.3 Some Properties of HGL

We now present several properties of the HGL optimization problem (8), which can be used to provide guidance on the suitable range for tuning parameters $\lambda_1$, $\lambda_2$, and $\lambda_3$. In what follows, $\mathbf{Z}^*$ and $\mathbf{V}^*$ denote the optimal solutions for $\mathbf{Z}$ and $\mathbf{V}$ in (8). Let $\frac{1}{s} + \frac{1}{q} = 1$ (recall that $q$ appears in (2)).

**Lemma 1.** *A sufficient condition for $\mathbf{Z}^*$ to be a diagonal matrix is that $\lambda_1 > \frac{\lambda_2 + \lambda_3}{2}$.*

**Lemma 2.** *A sufficient condition for $\mathbf{V}^*$ to be a diagonal matrix is that $\lambda_1 < \frac{\lambda_2}{2} + \frac{\lambda_3}{2(p-1)^{1/s}}$.*

**Corollary 3.** *A necessary condition for both $\mathbf{V}^*$ and $\mathbf{Z}^*$ to be non-diagonal matrices is that $\frac{\lambda_2}{2} + \frac{\lambda_3}{2(p-1)^{1/s}} \leq \lambda_1 \leq \frac{\lambda_2 + \lambda_3}{2}$.*

Furthermore, (8) reduces to the graphical lasso problem (7) under a simple condition.

**Lemma 3.** *If $q = 1$, then (8) reduces to (7) with tuning parameter $\min\left\{\lambda_1, \frac{\lambda_2 + \lambda_3}{2}\right\}$.*

Note also that when $\lambda_2 \to \infty$ or $\lambda_3 \to \infty$, (8) reduces to (7) with tuning parameter $\lambda_1$. However, throughout this paper, we assume that $q = 2$, and $\lambda_2$ and $\lambda_3$ are finite.

## 3.4 Simulation Study

In this section, we compare HGL to two sets of proposals: proposals that learn an Erdős-Rényi Gaussian graphical model, and proposals that learn a Gaussian graphical model in which some nodes are highly-connected.

### 3.4.1 Notation and Measures of Performance

We start by defining some notation. Let $\hat{\boldsymbol{\Theta}}$ be the estimate of $\boldsymbol{\Theta} = \boldsymbol{\Sigma}^{-1}$ from a given proposal, and let $\hat{\boldsymbol{\Theta}}_j$ be its $j$th column. Let $\mathcal{H}$ denote the set of indices of the hub nodes in $\boldsymbol{\Theta}$ (that is, this is the set of true hub nodes in the graph), and let $|\mathcal{H}|$ denote the cardinality of the set. In addition, let $\hat{\mathcal{H}}_r$ be the set of *estimated hub nodes*: the set of nodes in $\hat{\boldsymbol{\Theta}}$ that are among the $|\mathcal{H}|$ most highly-connected nodes, and that have at least $r$ edges. The values chosen for $|\mathcal{H}|$ and $r$ depend on the simulation set-up, and will be specified in each simulation study.

We now define several measures of performance that will be used to evaluate the various methods.

- Number of correctly estimated edges: $\sum_{j<j'} \left( 1_{\{|\hat{\Theta}_{jj'}|>10^{-5} \text{ and } |\Theta_{jj'}|\neq 0\}} \right)$.

- Proportion of correctly estimated hub edges:

$$\frac{\sum_{j\in\mathcal{H},j'\neq j}\left(1_{\{|\hat{\Theta}_{jj'}|>10^{-5} \text{ and } |\Theta_{jj'}|\neq 0\}}\right)}{\sum_{j\in\mathcal{H},j'\neq j}\left(1_{\{|\Theta_{jj'}|\neq 0\}}\right)}.$$

- Proportion of correctly estimated hub nodes: $\frac{|\hat{\mathcal{H}}_r \cap \mathcal{H}|}{|\mathcal{H}|}$.

- Sum of squared errors: $\sum_{j<j'}\left(\hat{\Theta}_{jj'} - \Theta_{jj'}\right)^2$.

### 3.4.2 Data Generation

We consider three set-ups for generating a $p \times p$ adjacency matrix $\mathbf{A}$.

A - Network with hub nodes: for all $i < j$, we set $A_{ij} = 1$ with probability 0.02, and zero otherwise. We then set $A_{ji}$ equal to $A_{ij}$. Next, we randomly select $|\mathcal{H}| = 20$ hub nodes and set the elements of the corresponding rows and columns of $\mathbf{A}$ to equal one with probability 0.7 and zero otherwise.

B - Network with two connected components and hub nodes: the adjacency matrix is generated as $\mathbf{A} = \begin{pmatrix} \mathbf{A}_1 & 0 \\ 0 & \mathbf{A}_2 \end{pmatrix}$, with $\mathbf{A}_1$ and $\mathbf{A}_2$ as in Set-up A, each with 10 hub nodes.

C - Scale-free network: the probability that a given node has $k$ edges is proportional to $k^{-\alpha}$. Barabási & Albert (1999) observed that many real-world networks have $\alpha \in [2.1, 4]$; we took $\alpha = 2.5$. We consider any node with at least 50 edges to be a hub node.

We then use the adjacency matrix $\mathbf{A}$ to create a matrix $\mathbf{E}$, as

$$E_{ij} \overset{\text{i.i.d.}}{\sim} \begin{cases} 0 & \text{if } A_{ij} = 0 \\ \text{Unif}([-0.75, -0.25] \cup [0.25, 0.75]) & \text{otherwise} \end{cases},$$

and set $\bar{\mathbf{E}} = \frac{1}{2}(\mathbf{E} + \mathbf{E}^T)$. Given the matrix $\bar{\mathbf{E}}$, we set $\boldsymbol{\Sigma}^{-1}$ equal to $\bar{\mathbf{E}} + (0.1 - \Lambda_{\min}(\bar{\mathbf{E}}))\mathbf{I}$, where $\Lambda_{\min}(\bar{\mathbf{E}})$ is the smallest eigenvalue of $\bar{\mathbf{E}}$. We generate the data matrix $\mathbf{X}$ according to $\mathbf{x}_1, \ldots, \mathbf{x}_n \overset{\text{i.i.d.}}{\sim} N(\mathbf{0}, \boldsymbol{\Sigma})$. Then, variables are standardized to have standard deviation one.

9

### 3.4.3 Comparison to Proposals that Assume an Erdős-Rényi Graph

In this subsection, we compare the performance of HGL to three proposals that learn an Erdős-Rényi Gaussian graphical model, in which each edge is equally likely and independent of all other edges. The three proposals are as follows:

- The graphical lasso (7), implemented using the R package `glasso`.

- Neighborhood selection (Meinshausen & Bühlmann 2006), implemented using the R package `glasso`. This approach involves performing $p$ $\ell_1$-penalized regression problems, each of which involves regressing one feature onto the others.

- Sparse partial correlation estimation (Peng et al. 2009), implemented using the R package `space`. This extension of the neighborhood selection approach combines the $p$ $\ell_1$-penalized regression problems in order to obtain a symmetric estimator.

We consider the three simulation set-ups described in the previous section with $n = 500$ and $p = 1000$. Figure 3 displays the results, averaged over 50 simulated data sets. Note that the sum of squared errors is not computed for Peng et al. (2009) or Meinshausen & Bühlmann (2006), since those methods do not directly yield estimates of $\mathbf{\Theta} = \mathbf{\Sigma}^{-1}$.

HGL has three tuning parameters. To obtain the curves shown in Figure 3, we fixed $\lambda_1 = 0.4$, considered three values of $\lambda_3$ (each shown in a different color in Figure 3), and used a fine grid of values of $\lambda_2$. The proposals of Friedman et al. (2007), Peng et al. (2009), and Meinshausen & Bühlmann (2006) each involve one tuning parameter; we applied them using a fine grid of the tuning parameter to obtain the curves shown in Figure 3.

Results for Set-up A are displayed in Figures 3-A(i) through 3-A(iv), where we calculate the proportion of correctly estimated hub nodes as defined in Section 3.4.1 with $r = 200$. Since this simulation set-up exactly matches the assumptions of HGL, it is not surprising that HGL outperforms other methods. In particular, HGL is able to identify most of the hub nodes when the number of estimated edges is approximately equal to the true number of edges. We see similar results for Set-up B in Figures 3-B(i) through 3-B(iv), where the proportion of correctly estimated hub nodes is as defined in Section 3.4.1 with $r = 100$.

In Set-up C, recall that we define a node with at least 50 edges to be a hub. The proportion of correctly estimated hub nodes is then as defined in Section 3.4.1 with $r = 50$. The results are presented in Figures 3-C(i) through 3-C(iv). In this set-up, only one or two of the nodes (on average) have more than 50 edges, and the hub nodes are not as highly-connected as in Set-up A or Set-up B. Nonetheless, HGL outperforms Friedman et al. (2007) and Meinshausen & Bühlmann (2006), and is comparable to Peng et al. (2009), which (according to the authors of that paper) performs well on scale-free networks.

### 3.4.4 Comparison to Proposals that Do Not Assume an Erdős-Rényi Graph

In this subsection, we compare the performance of HGL to two proposals that model a Gaussian graphical model under the assumption that some nodes have many more edges than others. These proposals are as follows:

- The partial correlation screening procedure of Hero & Rajaratnam (2012). The elements of the partial correlation matrix are thresholded based on their absolute value,
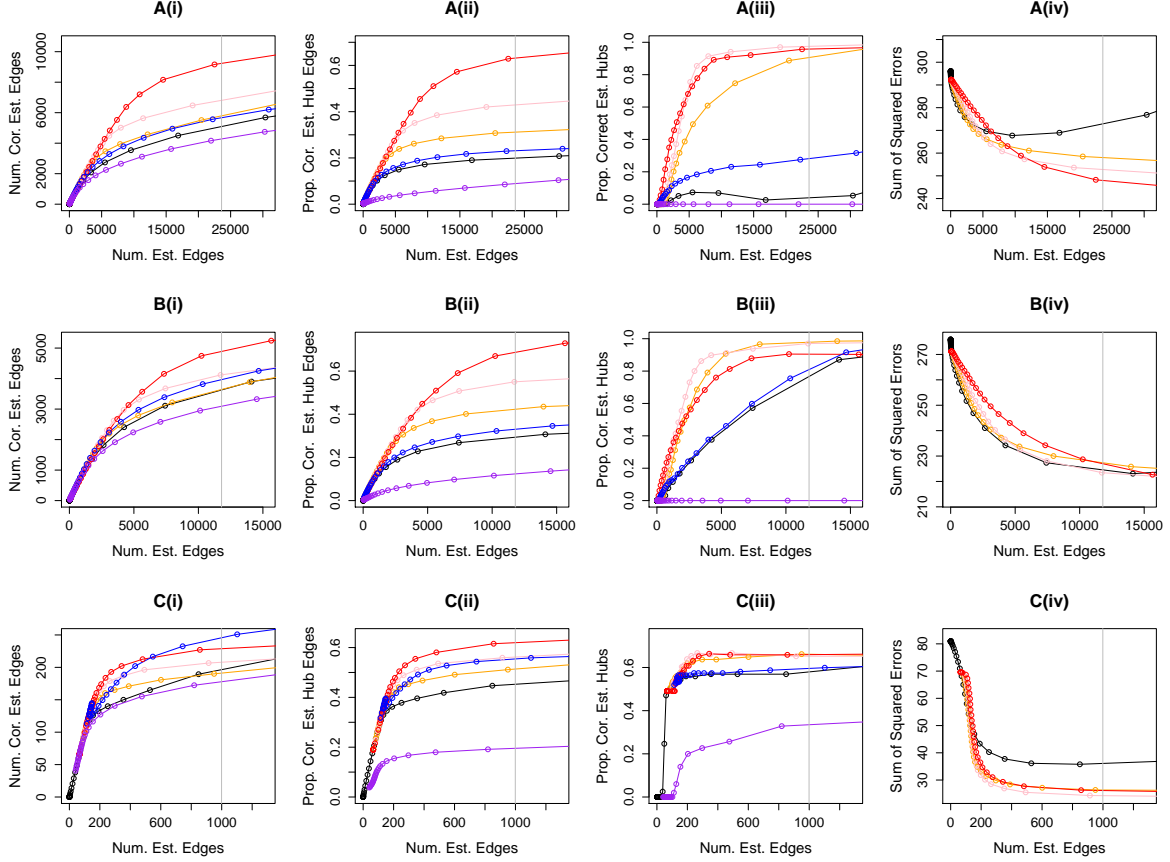
Figure 3: Row A: Results for Set-up A. Row B: Results for Set-up B. Row C: Results for Set-up C. In each panel, the $x$-axis displays the number of estimated edges, and the vertical gray line is the number of edges in the true network. The $y$-axes are as follows: Column (i): Number of correctly estimated edges; Column (ii): Proportion of correctly estimated hub edges; Column (iii): Proportion of correctly estimated hub nodes; Column (iv): Sum of squared errors. Colored lines correspond to the graphical lasso (Friedman et al. 2007) (——); HGL with $\lambda_3 = 0.5$ (——), $\lambda_3 = 1$ (——), and $\lambda_3 = 2$ (——); neighborhood selection (Meinshausen & Bühlmann 2006) (——); sparse partial correlation estimation (Peng et al. 2009) (——).

11

and a hub node is declared if the number of nonzero elements in the corresponding column of the thresholded partial correlation matrix is sufficiently large.

- The scale-free network estimation procedure of Liu & Ihler (2011). This is the solution to the non-convex optimization problem

$$\underset{\boldsymbol{\Theta}}{\text{minimize}} \quad \left\{ -\log\det\boldsymbol{\Theta} + \text{trace}(\mathbf{S}\boldsymbol{\Theta}) + \alpha\sum_{j=1}^{p}\log(\|\theta_{\backslash j}\|_1 + \epsilon_j) + \sum_{j=1}^{p}\beta_j|\theta_{jj}| \right\}, \qquad (10)$$

where $\theta_{\backslash j} = \{\theta_{jj'}|j' \neq j\}$, and $\epsilon_j$, $\beta_j$, and $\alpha$ are tuning parameters.

We generated data under Set-ups A and C (described in Section 3.4.2) with $n = 125$ and $p = 250$, with slight modifications: in Set-up A, we took $|\mathcal{H}| = 5$, and in Set-up C, we considered any node with at least 10 edges to be a hub. The results, averaged over 50 data sets, are displayed in Figures 4 and 5.

To obtain Figures 4 and 5, we applied (10) using a fine grid of $\alpha$ values, and using the choices for $\beta_j$ and $\epsilon_j$ specified by the authors: $\beta_j = 2\alpha/\epsilon_j$, where $\epsilon_j$ is a small constant specified in Liu & Ihler (2011). There are two tuning parameters in Hero & Rajaratnam (2012): (1) $\rho$, the value used to threshold the partial correlation matrix, and (2) $d$, the number of non-zero elements required for a column of the thresholded matrix to be declared a hub node. We used $d = \{10, 20\}$ in Figure 4 and $d = \{5, 10\}$ in Figure 5, and used a fine grid of values for $\rho$. Note that the value of $d$ has no effect on the results for Figures 4(i)-(ii) and Figures 5(i)-(ii), and that larger values of $d$ tend to yield worse results in Figures 4(iii) and 5(iii). The sum of squared errors was not calculated for Hero & Rajaratnam (2012) since it does not yield an estimate of the precision matrix. As a baseline reference, the graphical lasso is included in the comparison.

We see from Figure 4 that HGL outperforms the competitors when the underlying network contains hub nodes. It is not surprising that Liu & Ihler (2011) yields better results than the graphical lasso, since the former approach is implemented via an iterative procedure: in each iteration, the graphical lasso is performed with an updated tuning parameter based on the estimate obtained in the previous iteration. Hero & Rajaratnam (2012) has the worst results in Figures 4(i)-(ii); this is not surprising, since the purpose of Hero & Rajaratnam (2012) is to screen for hub nodes, rather than to estimate the individual edges in the network.

From Figure 5, we see that the performance of HGL is comparable to that of Liu & Ihler (2011) under the assumption of a scale-free network; note that this is the precise setting for which Liu & Ihler (2011)'s proposal is intended. Again, Liu & Ihler (2011) outperforms the graphical lasso, and Hero & Rajaratnam (2012) has the worst results in Figures 5(i)-(ii). We note that Liu & Ihler (2011) has the lowest MSE in both Figures 4 and 5. This is not surprising, since Liu & Ihler (2011)'s approach is based on a non-convex formulation, and thus avoids over-shrinkage of the non-zero estimates.

# 4   The Hub Covariance Graph

In this section, we consider estimation of a covariance matrix under the assumption that $\mathbf{x}_1, \ldots, \mathbf{x}_n \overset{\text{i.i.d.}}{\sim} N(\mathbf{0}, \boldsymbol{\Sigma})$; this is of interest because the sparsity pattern of $\boldsymbol{\Sigma}$ specifies the
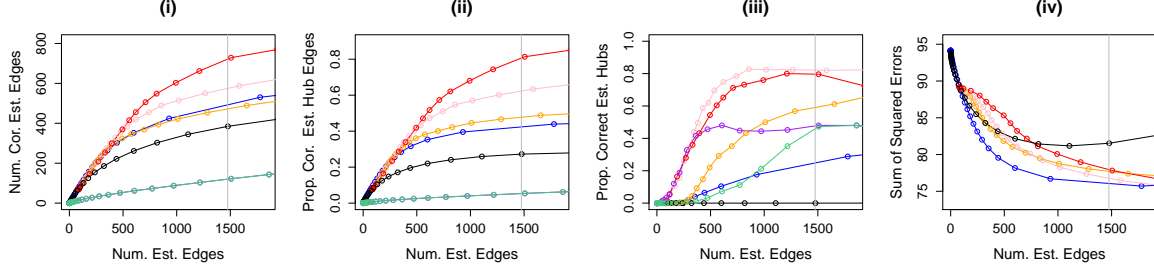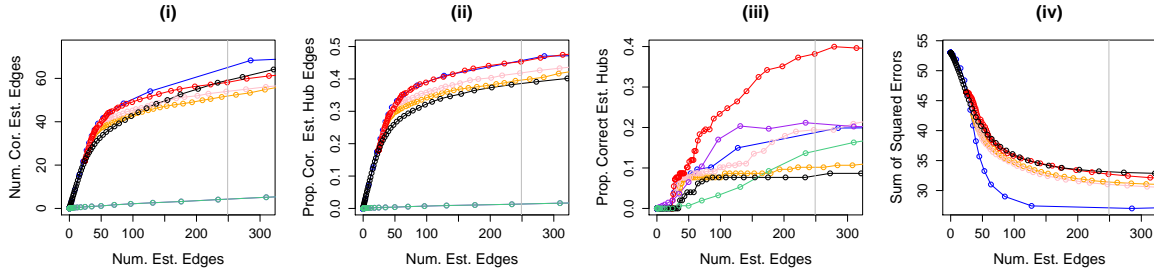
Figure 4: Simulation set-up A was applied with $n = 125$ and $p = 250$. Details of the axis labels are as in Figure 3. The colored lines correspond to the graphical lasso (Friedman et al. 2007) (——); HGL with $\lambda_3 = 1$ (——), $\lambda_3 = 2$ (——), and $\lambda_3 = 4$ (——); the hub screening procedure (Hero & Rajaratnam 2012) with $d = 10$ (——) and $d = 20$ (——); the scale-free network approach (Liu & Ihler 2011) (——).



Figure 5: Simulation set-up C was applied with $n = 125$ and $p = 250$. Details of the axis labels are as in Figure 3. The colored lines correspond to the graphical lasso (Friedman et al. 2007) (——); HGL with $\lambda_3 = 0.5$ (——), $\lambda_3 = 1$ (——), and $\lambda_3 = 2$ (——); the hub screening procedure (Hero & Rajaratnam 2012) with $d = 5$ (——) and $d = 10$ (——); the scale-free network approach (Liu & Ihler 2011) (——).

structure of the marginal independence graph. The idea of estimating a covariance graph is not new (see e.g., Drton & Richardson 2003, Chaudhuri et al. 2007, Drton & Richardson 2008). We extend the covariance estimator of Xue et al. (2012) to accommodate hub nodes.

## 4.1   Formulation and Algorithm

Xue et al. (2012) proposed to estimate $\boldsymbol{\Sigma}$ using

$$\hat{\boldsymbol{\Sigma}} = \underset{\boldsymbol{\Sigma} \succeq \epsilon \mathbf{I}}{\arg \min} \left\{ \frac{1}{2} \|\boldsymbol{\Sigma} - \mathbf{S}\|_F^2 + \lambda \|\boldsymbol{\Sigma}\|_1 \right\}, \tag{11}$$

where $\mathbf{S}$ is the empirical covariance matrix, and $\epsilon$ is a small positive constant; we take $\epsilon = 10^{-4}$. We extend (11) to accommodate hubs by imposing the hub penalty function (2)

on $\boldsymbol{\Sigma}$. This results in the *hub covariance graph* (HCG) optimization problem,

$$\underset{\boldsymbol{\Sigma} \succeq \epsilon \mathbf{I}}{\text{minimize}} \left\{ \frac{1}{2}\|\boldsymbol{\Sigma} - \mathbf{S}\|_F^2 + \mathrm{P}(\boldsymbol{\Sigma}) \right\}, \tag{12}$$

which can be solved via Algorithm 1. To update $\boldsymbol{\Theta} = \boldsymbol{\Sigma}$ in Step 2(a)i, we note that

$$\underset{\boldsymbol{\Sigma} \succeq \epsilon \mathbf{I}}{\arg\min} \left\{ \frac{1}{2}\|\boldsymbol{\Sigma} - \mathbf{S}\|_F^2 + \frac{\rho}{2}\|\boldsymbol{\Sigma} - \tilde{\boldsymbol{\Sigma}} + \mathbf{W}_1\|_F^2 \right\} = \frac{1}{1+\rho}(\mathbf{S} + \rho\tilde{\boldsymbol{\Sigma}} - \rho\mathbf{W}_1)^+,$$

where $(\mathbf{A})^+$ is the projection of a matrix $\mathbf{A}$ onto the convex cone $\{\boldsymbol{\Sigma} \succeq \epsilon \mathbf{I}\}$. That is, if $\sum_{j=1}^p d_j \mathbf{u}_j \mathbf{u}_j^T$ denotes the eigen-decomposition of the matrix $\mathbf{A}$, then $(\mathbf{A})^+$ is defined as $\sum_{j=1}^p \max(d_j, \epsilon)\mathbf{u}_j \mathbf{u}_j^T$. The complexity of the ADMM algorithm is $O(p^3)$ per iteration, due to the complexity of the eigen-decomposition for updating $\boldsymbol{\Sigma}$.

## 4.2 Simulation Study

We compare HCG to two competitors for obtaining a sparse estimate of $\boldsymbol{\Sigma}$:

1. The non-convex $\ell_1$-penalized log-likelihood approach of Bien & Tibshirani (2011), using the R package spcov. This approach solves

$$\underset{\boldsymbol{\Sigma} \succ 0}{\text{minimize}} \left\{ \log \det \boldsymbol{\Sigma} + \text{trace}(\boldsymbol{\Sigma}^{-1}\mathbf{S}) + \lambda\|\boldsymbol{\Sigma}\|_1 \right\}. \tag{13}$$

2. The convex $\ell_1$-penalized approach of Xue et al. (2012), given in (11).

We first generated an adjacency matrix $\mathbf{A}$ as in Set-up A in Section 3.4.2, modified to have 5 hub nodes. Then $\bar{\mathbf{E}}$ was generated as described in Section 3.4.2, and we set $\boldsymbol{\Sigma}$ equal to $\bar{\mathbf{E}} + (0.1 - \Lambda_{\min}(\bar{\mathbf{E}}))\mathbf{I}$. Next, we generated $\mathbf{x}_1, \ldots, \mathbf{x}_n \overset{\text{i.i.d.}}{\sim} N(\mathbf{0}, \boldsymbol{\Sigma})$. Finally, we standardized the variables to have mean zero and standard deviation one. In this simulation study, we set $n = 100$ and $p = 200$ due to computational constraints in the spcov R package; the HCG algorithm can be applied to much larger data sets.

Figure 6 displays the results, averaged over 50 simulated data sets. We calculated the proportion of correctly estimated hub nodes as defined in Section 3.3.1 with $r = 40$. We used a fine grid of tuning parameters for Bien & Tibshirani (2011) and Xue et al. (2012) in order to obtain the curves shown in each panel of Figure 6. HCG involves three tuning parameters, $\lambda_1$, $\lambda_2$, and $\lambda_3$. We fixed $\lambda_1 = 0.2$, considered three values of $\lambda_3$ (each shown in a different color), and varied $\lambda_2$ in order to obtain the curves shown in Figure 6.

We see that HCG outperforms the other methods. These results are not surprising, since the other methods do not explicitly model the hub nodes.

## 5 The Hub Binary Network

In this section, we focus on estimating a binary Ising Markov random field, which we refer to as a binary network. We refer the reader to Ravikumar et al. (2010) for an in-depth discussion of this type of graphical model and its applications.
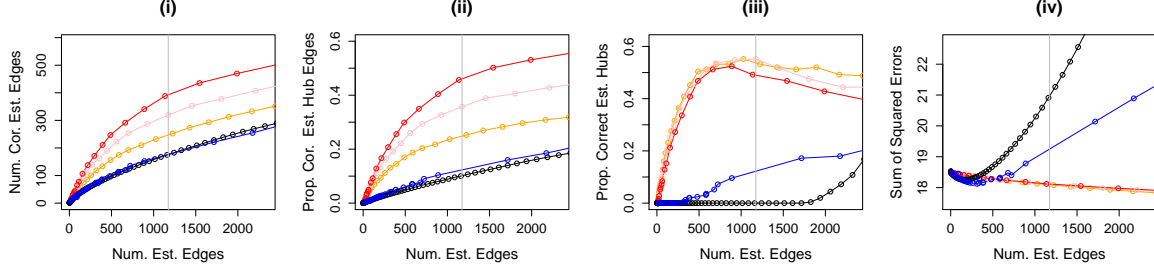
14

Figure 6: Details of the axis labels are as in Figure 3. The colored lines correspond to the proposal of Xue et al. (2012) (——); HCG with $\lambda_3 = 1$ (——), $\lambda_3 = 1.5$ (——), and $\lambda_3 = 2$ (——); and the proposal of Bien & Tibshirani (2011) (——).

In this set-up, each entry of the $n \times p$ data matrix $\mathbf{X}$ takes on a value of zero or one. We assume that each observation $\mathbf{x}_1, \ldots, \mathbf{x}_n \overset{\text{i.i.d.}}{\sim} p(\mathbf{x}, \mathbf{\Theta})$, where

$$p(\mathbf{x}, \mathbf{\Theta}) = \frac{1}{Z(\mathbf{\Theta})} \exp\left[\sum_{j=1}^{p} \theta_{jj} x_j + \sum_{1 \le j < j' \le p} \theta_{jj'} x_j x_{j'}\right], \tag{14}$$

and where $Z(\mathbf{\Theta})$ is the partition function, which ensures that the density sums to one. Here $\mathbf{\Theta}$ is a $p \times p$ symmetric matrix that specifies the network structure: $\theta_{jj'} = 0$ implies that the $j$th and $j'$th variables are conditionally independent.

In order to obtain a sparse graph, Lee et al. (2007) considered maximizing an $\ell_1$-penalized log-likelihood under this model. Due to the difficulty in computing the log-partition function, several authors have considered alternative approaches. For instance, Ravikumar et al. (2010) proposed a neighborhood selection approach, and Höfling & Tibshirani (2009) considered maximizing an $\ell_1$-penalized pseudo-likelihood.

## 5.1   Formulation and Algorithm

Under the model (14), the log-pseudo-likelihood for $n$ observations takes the form

$$\sum_{j=1}^{p} \sum_{j'=1}^{p} \theta_{jj'} (\mathbf{X}^T \mathbf{X})_{jj'} - \sum_{i=1}^{n} \sum_{j=1}^{p} \log(1 + \exp[\theta_{jj} + \sum_{j' \ne j} \theta_{jj'} x_{ij'}]), \tag{15}$$

where $\mathbf{x}_i$ is the $i$th row of the $n \times p$ matrix $\mathbf{X}$. The proposal of Höfling & Tibshirani (2009) involves maximizing (15) subject to an $\ell_1$ penalty on $\mathbf{\Theta}$. We propose to instead impose the hub penalty function (2) on $\mathbf{\Theta}$ in (15) in order to estimate a sparse binary network with hub nodes. This leads to the optimization problem

$$\underset{\mathbf{\Theta}}{\text{minimize}} \quad \left\{ -\sum_{j=1}^{p} \sum_{j'=1}^{p} \theta_{jj'} (\mathbf{X}^T \mathbf{X})_{jj'} + \sum_{i=1}^{n} \sum_{j=1}^{p} \log(1 + \exp[\theta_{jj} + \sum_{j' \ne j} \theta_{jj'} x_{ij'}]) + \mathrm{P}(\mathbf{\Theta}) \right\}. \tag{16}$$

15

We refer to the solution to (16) as the *hub binary network* (HBN). The ADMM algorithm for solving (16) is given in Algorithm 1. We solve the update for $\boldsymbol{\Theta}$ in Step 2(a)i using a local quadratic approximation with a diagonal Hessian matrix, as is detailed in Höfling & Tibshirani (2009). Other algorithms, such as L-BFGS, could instead be used.

## 5.2 Simulation Study

Here we compare the performance of HBN to the proposal of Höfling & Tibshirani (2009), implemented using the R package BMN.

We simulated a binary network with $p = 30$ and $|\mathcal{H}| = 5$ hub nodes. To generate the parameter matrix $\boldsymbol{\Theta}$, we created an adjacency matrix $\mathbf{A}$ as in Set-up A of Section 3.4.2 with five hub nodes. Then $\bar{\mathbf{E}}$ was generated as in Section 3.4.2, and we set $\boldsymbol{\Theta} = \bar{\mathbf{E}}$.

Each of $n = 100$ observations was generated using Gibbs sampling (Ravikumar et al. 2010, Guo et al. 2010). Suppose that $x_1^{(t)}, \ldots, x_p^{(t)}$ is obtained at the $t$th iteration of the Gibbs sampler. Then, the $(t+1)$th iteration is obtained according to

$$x_j^{(t+1)} \sim \text{Bernoulli} \left( \frac{\exp(\theta_{jj} + \sum_{j \neq j'} \theta_{jj'} x_{j'}^{(t)})}{1 + \exp(\theta_{jj} + \sum_{j \neq j'} \theta_{jj'} x_{j'}^{(t)})} \right) \qquad \text{for } j = 1, \ldots, p.$$

We took the first $10^5$ iterations as our burn-in period, and then collected an observation every $10^4$ iterations, such that the observations were nearly independent (Guo et al. 2010).

The results, averaged over 50 data sets, are shown in Figure 7. We used a fine grid of values for the $\ell_1$ tuning parameter for Höfling & Tibshirani (2009), resulting in curves shown in each panel of the figure. For HBN, we fixed $\lambda_1 = 5$, considered $\lambda_3 = \{10, 15, 20\}$, and used a fine grid of values of $\lambda_2$. The proportion of correctly estimated hub nodes was calculated using the definition in Section 3.4.1 with $r = 15$. Figure 7 indicates that HBN consistently outperforms the proposal of Höfling & Tibshirani (2009).
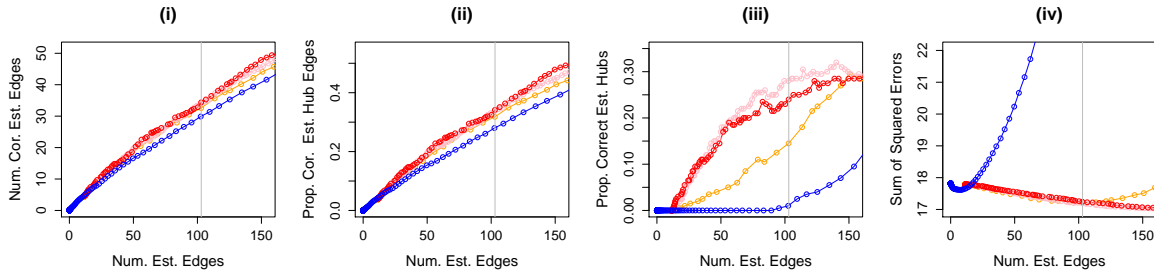


Figure 7: Details of the axis labels are as in Figure 3. The colored lines correspond to the $\ell_1$-penalized pseudo-likelihood proposal of Höfling & Tibshirani (2009) (———); and HBN with $\lambda_3 = 10$ (———), $\lambda_3 = 15$ (———), and $\lambda_3 = 20$ (———).

# 6   Real Data Application

We now apply HGL to a university webpage data set, and a brain cancer data set.

## 6.1 Application to University Webpage Data

We applied HGL to the university webpage data set from the "World Wide Knowledge Base" project at Carnegie Mellon University. This data set was pre-processed by Cardoso-Cachopo (2009). The data set consists of the occurrences of various terms (words) on webpages from four computer science departments at Cornell, Texas, Washington and Wisconsin. We consider only the 544 student webpages, and select 100 terms with the largest entropy for our analysis. The goal of the analysis is to understand the relationships among the terms that appear on the student webpages. To begin, we standardize each term to have standard deviation one. HGL is performed with $\lambda_1 = 0.55$, $\lambda_2 = 0.35$, and $\lambda_3 = 1.5$; these values were chosen in order to obtain a network estimate that is sparse enough that it can be easily interpreted. The estimated matrices are shown in Figure 8.

Figure 8(a) indicates that four hub nodes are detected: *compute*, *research*, *scienc*, and *system*. For instance, the fact that *compute* is a hub indicates that many terms' occurrences are fully explained by the occurrence of the word *compute*. From Figure 8(b), we see that several pairs of terms take on non-zero values in the matrix $(\mathbf{Z} - \mathrm{diag}(\mathbf{Z}))$. These include *(depart, universe)*; *(home, page)*; *(institut, technolog)*; and *(scienc, univers)*. These results provide an intuitive explanation of the relationships among the terms in the webpages.
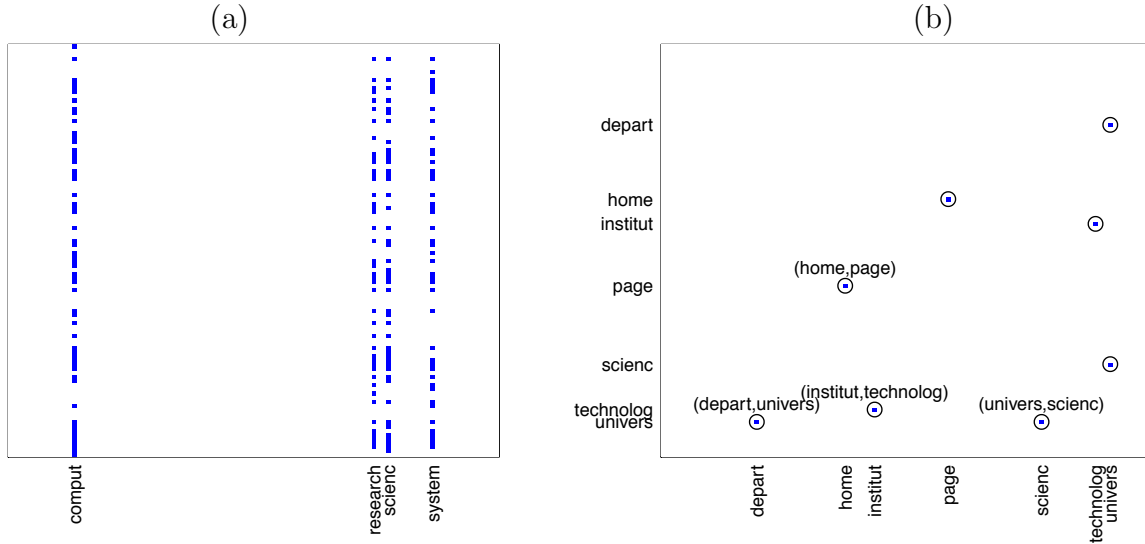


Figure 8: Results for HGL on the webpage data with $\lambda_1 = 0.55$, $\lambda_2 = 0.35$, $\lambda_3 = 1.5$. Non-zero estimated values are shown, for *(a):* $(\mathbf{V} - \mathrm{diag}(\mathbf{V}))$, and *(b):* $(\mathbf{Z} - \mathrm{diag}(\mathbf{Z}))$.

## 6.2 Application to Gene Expression Data

We applied HGL to a publicly available cancer gene expression data set (Verhaak et al. 2010). The data set consists of mRNA expression levels for 17,814 genes in 401 patients with glioblastoma multiforme (GBM), an extremely aggressive cancer with very poor patient prognosis. We aim to reconstruct the gene regulatory network that represents the interactions among the genes, as well as to identify hub genes that tend to have many interactions

with other genes. Such genes likely play an important role in regulating many other genes in the network. Identifying such regulatory genes will lead to a better understanding of brain cancer, and eventually may lead to new therapeutic targets.

Among 7,462 genes known to be associated with cancer (Rappaport et al. 2013), we selected 500 with the highest variance. We applied HGL to the empirical covariance matrix corresponding to the $401 \times 500$ data matrix, after standardizing each gene to have variance one. In Figure 9, we plotted the resulting network (for simplicity, only the 248 genes with at least two neighbors are displayed). We found that seven genes are identified as hubs and are connected to more than 50 other nodes. These genes are TBC1D2B, PTPN2, SLC45A1, PPP1R12C, ZNF763, AAAS, and PRH2, in decreasing order of estimated edges.

Interestingly, some of these genes have known regulatory roles. PTPN2 is known to be a signaling molecule that regulates a variety of cellular processes including cell growth, differentiation, mitotic cycle, and oncogenic transformation (Maglott et al. 2004). ZNF763 is a DNA-binding protein that regulates the transcription of other genes (Maglott et al. 2004). AAAS is a regulatory gene believed to be involved in normal development of the peripheral and central nervous system (Rappaport et al. 2013). These genes do not appear to be highly-connected to many other genes in the estimate that results from applying the graphical lasso (7) to this same data set (results not shown). These results indicate that HGL can be used to recover known regulators, as well as to suggest other potential regulators that may be targets for follow-up analysis.
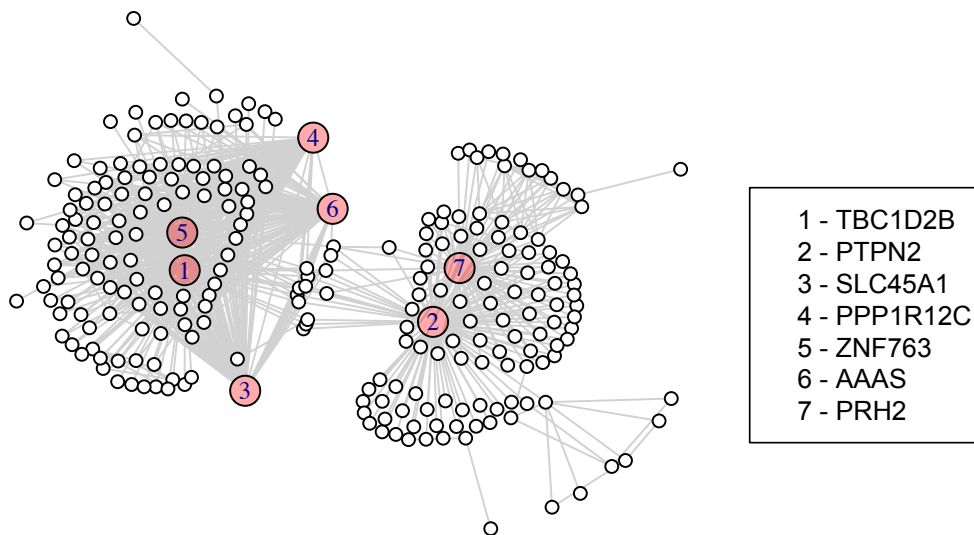


Figure 9: Results for HGL on the GBM data with $\lambda_1 = 0.7$, $\lambda_2 = 0.8$, $\lambda_3 = 2.8$. Only nodes with at least two edges in the estimated network are displayed. Nodes displayed in pink were found to be hubs by the HGL algorithm.

18

# 7 Discussion

We have proposed a general framework for estimating a network with hubs by way of a convex penalty function. The proposed framework has three tuning parameters, so that it can flexibly accommodate different numbers of hubs, sparsity levels within a hub, and connectivity levels among non-hubs. However, tuning parameter selection in unsupervised settings remains a challenging open problem (see e.g., Foygel & Drton 2010, Meinshausen & Bühlmann 2010). In practice, tuning parameters could be set based on domain knowledge, a desire for interpretability of the resulting estimates, or a BIC-type criterion.

The framework proposed in this paper assumes an underlying model involving a set of edges between non-hub nodes, as well as a set of hub nodes. For instance, it is believed that such hub nodes arise in biology, in which "super hubs" in transcriptional regulatory networks may play important roles (Hao et al. 2012). We note here that the underlying model of hub nodes assumed in this paper differs fundamentally from a scale-free network in which the degree of connectivity of the nodes follows a power law distribution — scale-free networks simply do not have such very highly-connected hub nodes. In fact, we have shown that existing techniques for estimating a scale-free network, such as Liu & Ihler (2011) and Defazio & Caetano (2012), cannot accommodate the very dense hubs for which our proposal is intended.

As discussed in Section 2, the hub penalty function involves decomposing a parameter matrix $\mathbf{\Theta}$ into $\mathbf{Z} + \mathbf{V} + \mathbf{V}^T$, where $\mathbf{Z}$ is a sparse matrix, and $\mathbf{V}$ is a matrix whose columns are entirely zero or (almost) entirely non-zero. In this paper, we used an $\ell_1$ penalty on $\mathbf{Z}$ in order to encourage it to be sparse. In effect, this amounts to assuming that the non-hub nodes obey an Erdős-Rényi network. But our formulation could be easily modified to accommodate a different network prior for the non-hub nodes. For instance, we could assume that the non-hub nodes obey a scale-free network, using the ideas developed in Liu & Ihler (2011) and Defazio & Caetano (2012). This would amount to modeling a scale-free network with some extra-hubby nodes.

In this paper, we applied the proposed framework to the tasks of estimating a Gaussian graphical model, a covariance graph model, and a binary network. The proposed framework can also be applied to other types of graphical models, such as the Poisson graphical model (Allen & Liu 2012) or the exponential family graphical model (Yang et al. 2012).

In future work, we will study the theoretical statistical properties of the HGL formulation. For instance, in the context of the graphical lasso, it is known that the rate of statistical convergence depends upon the maximal degree of any node in the network (Ravikumar et al. 2011). It would be interesting to see whether HGL theoretically outperforms the graphical lasso in the setting in which the true underlying network contains hubs.

# Appendix A: Derivation of Algorithm 1

Recall that the scaled augmented Lagrangian for (5) takes the form

$$L(\mathbf{A}, \mathbf{B}, \mathbf{W}) = \ell(\mathbf{X}, \mathbf{\Theta}) + \lambda_1 \|\mathbf{Z} - \text{diag}(\mathbf{Z})\|_1 + \lambda_2 \|\mathbf{V} - \text{diag}(\mathbf{V})\|_1$$
$$+ \lambda_3 \sum_{j=1}^{p} \|(\mathbf{V} - \text{diag}(\mathbf{V}))_j\|_2 + g(\mathbf{B}) + \frac{\rho}{2}\|\mathbf{A} - \mathbf{B} + \mathbf{W}\|_F^2. \tag{17}$$

The proposed ADMM algorithm requires the following updates:

1. $\mathbf{A}^{(t+1)} \leftarrow \underset{\mathbf{A}}{\text{argmin}}\ L(\mathbf{A}, \mathbf{B}^t, \mathbf{W}^t)$,

2. $\mathbf{B}^{(t+1)} \leftarrow \underset{\mathbf{B}}{\text{argmin}}\ L(\mathbf{A}^{(t+1)}, \mathbf{B}, \mathbf{W}^t)$,

3. $\mathbf{W}^{(t+1)} \leftarrow \mathbf{W}^t + \mathbf{A}^{(t+1)} - \mathbf{B}^{(t+1)}$.

We now proceed to derive the updates for $\mathbf{A}$ and $\mathbf{B}$.

## Updates for A

To obtain updates for $\mathbf{A} = (\mathbf{\Theta}, \mathbf{V}, \mathbf{Z})$, we exploit the fact that (17) is separable in $\mathbf{\Theta}, \mathbf{V}$, and $\mathbf{Z}$. Therefore, we can simply update with respect to $\mathbf{\Theta}, \mathbf{V}$, and $\mathbf{Z}$ one-at-a-time. Updates for $\mathbf{\Theta}$ depend on the form of the convex loss function, and are addressed in the main text. Updates for $\mathbf{V}$ and $\mathbf{Z}$ can be easily seen to take the form given in Algorithm 1.

## Update for B

Minimizing the function in (17) with respect to $\mathbf{B}$ is equivalent to

$$\underset{\tilde{\mathbf{\Theta}}, \tilde{\mathbf{V}}, \tilde{\mathbf{Z}}}{\text{minimize}} \quad \left\{ \frac{\rho}{2}\|\mathbf{\Theta} - \tilde{\mathbf{\Theta}} + \mathbf{W}_1\|_F^2 + \frac{\rho}{2}\|\mathbf{V} - \tilde{\mathbf{V}} + \mathbf{W}_2\|_F^2 + \frac{\rho}{2}\|\mathbf{Z} - \tilde{\mathbf{Z}} + \mathbf{W}_3\|_F^2 \right\}$$
$$\text{subject to} \quad \tilde{\mathbf{\Theta}} = \tilde{\mathbf{Z}} + \tilde{\mathbf{V}} + \tilde{\mathbf{V}}^T. \tag{18}$$

Let $\mathbf{\Gamma}$ be the $p \times p$ Lagrange multiplier matrix for the equality constraint. Then, the Lagrangian for (18) is

$$\frac{\rho}{2}\|\mathbf{\Theta} - \tilde{\mathbf{\Theta}} + \mathbf{W}_1\|_F^2 + \frac{\rho}{2}\|\mathbf{V} - \tilde{\mathbf{V}} + \mathbf{W}_2\|_F^2 + \frac{\rho}{2}\|\mathbf{Z} - \tilde{\mathbf{Z}} + \mathbf{W}_3\|_F^2 + \langle \mathbf{\Gamma}, \tilde{\mathbf{\Theta}} - \tilde{\mathbf{Z}} - \tilde{\mathbf{V}} - \tilde{\mathbf{V}}^T \rangle. \tag{19}$$

A little bit of algebra yields

$$\tilde{\mathbf{\Theta}} = \mathbf{\Theta} + \mathbf{W}_1 - \frac{1}{\rho}\mathbf{\Gamma},$$

$$\tilde{\mathbf{V}} = \frac{1}{\rho}(\mathbf{\Gamma} + \mathbf{\Gamma}^T) + \mathbf{V} + \mathbf{W}_2,$$

$$\tilde{\mathbf{Z}} = \frac{1}{\rho}\mathbf{\Gamma} + \mathbf{Z} + \mathbf{W}_3,$$

where $\mathbf{\Gamma} = \frac{\rho}{6}[(\mathbf{\Theta} + \mathbf{W}_1) - (\mathbf{V} + \mathbf{W}_2) - (\mathbf{V} + \mathbf{W}_2)^T - (\mathbf{Z} + \mathbf{W}_3)]$.

# Appendix B: Conditions for HGL Solution to be Block-Diagonal

We begin by introducing some notation.

Let $\|\mathbf{V}\|_{u,v}$ be the $\ell_u/\ell_v$ norm of a matrix $\mathbf{V}$. For instance, $\|\mathbf{V}\|_{1,q} = \sum_{j=1}^{p} \|(\mathbf{V} - \text{diag}(\mathbf{V}))_j\|_q$. We define the support of a matrix $\boldsymbol{\Theta}$ as follows: $\text{supp}(\boldsymbol{\Theta}) = \{(i,j) : \Theta_{ij} \neq 0\}$. We say that $\boldsymbol{\Theta}$ is supported on a set $\mathcal{G}$ if $\text{supp}(\boldsymbol{\Theta}) \subseteq \mathcal{G}$. Let $\{C_1, \ldots, C_K\}$ be a partition of the index set $\{1, \ldots, p\}$, and let $\mathcal{T} = \cup_{k=1}^{K}\{C_k \times C_k\}$. We let $\mathbf{A}_{\mathcal{T}}$ denote the restriction of the matrix $\mathbf{A}$ to the set $\mathcal{T}$: that is, $(\mathbf{A}_{\mathcal{T}})_{ij} = 0$ if $(i,j) \notin \mathcal{T}$ and $(\mathbf{A}_{\mathcal{T}})_{ij} = A_{ij}$ if $(i,j) \in \mathcal{T}$. Note that any matrix supported on $\mathcal{T}$ is block-diagonal with $K$ blocks, subject to some permutation of its rows and columns. Also, let $S_{\max} = \max\limits_{(i,j)\in\mathcal{T}^c} |S_{ij}|$. Define

$$\tilde{\mathbf{P}}(\boldsymbol{\Theta}) = \min_{\mathbf{V},\mathbf{Z}} \quad \|\mathbf{Z} - \text{diag}(\mathbf{Z})\|_1 + \hat{\lambda}_2\|\mathbf{V} - \text{diag}(\mathbf{V})\|_1 + \hat{\lambda}_3\|\mathbf{V} - \text{diag}(\mathbf{V})\|_{1,q} \tag{20}$$
$$\text{subject to} \quad \boldsymbol{\Theta} = \mathbf{Z} + \mathbf{V} + \mathbf{V}^T,$$

where $\hat{\lambda}_2 = \frac{\lambda_2}{\lambda_1}$ and $\hat{\lambda}_3 = \frac{\lambda_3}{\lambda_1}$. Then, optimization problem (8) is equivalent to

$$\underset{\boldsymbol{\Theta}}{\text{minimize}} \quad -\log\det(\boldsymbol{\Theta}) + \langle\boldsymbol{\Theta}, \mathbf{S}\rangle + \lambda_1\tilde{\mathbf{P}}(\boldsymbol{\Theta}). \tag{21}$$

## Proof of Theorem 1 (Sufficient Condition)

*Proof.* First, we note that if $(\boldsymbol{\Theta}, \mathbf{V}, \mathbf{Z})$ is a feasible solution to (8), then $(\boldsymbol{\Theta}_{\mathcal{T}}, \mathbf{V}_{\mathcal{T}}, \mathbf{Z}_{\mathcal{T}})$ is also a feasible solution to (8). Assume that $(\boldsymbol{\Theta}, \mathbf{V}, \mathbf{Z})$ is not supported on $\mathcal{T}$. We want to show that the objective value of (8) evaluated at $(\boldsymbol{\Theta}_{\mathcal{T}}, \mathbf{V}_{\mathcal{T}}, \mathbf{Z}_{\mathcal{T}})$ is smaller than the objective value of (8) evaluated at $(\boldsymbol{\Theta}, \mathbf{V}, \mathbf{Z})$. By Fischer's inequality (Horn & Johnson 1985),

$$-\log\det(\boldsymbol{\Theta}) \geq -\log\det(\boldsymbol{\Theta}_{\mathcal{T}}).$$

Therefore, it remains to show that

$$\langle\boldsymbol{\Theta}, \mathbf{S}\rangle + \lambda_1\|\mathbf{Z} - \text{diag}(\mathbf{Z})\|_1 + \lambda_2\|\mathbf{V} - \text{diag}(\mathbf{V})\|_1 + \lambda_3\|\mathbf{V} - \text{diag}(\mathbf{V})\|_{1,q} >$$
$$\langle\boldsymbol{\Theta}_{\mathcal{T}}, \mathbf{S}\rangle + \lambda_1\|\mathbf{Z}_{\mathcal{T}} - \text{diag}(\mathbf{Z}_{\mathcal{T}})\|_1 + \lambda_2\|\mathbf{V}_{\mathcal{T}} - \text{diag}(\mathbf{V}_{\mathcal{T}})\|_1 + \lambda_3\|\mathbf{V}_{\mathcal{T}} - \text{diag}(\mathbf{V}_{\mathcal{T}})\|_{1,q},$$

or equivalently, that

$$\langle\boldsymbol{\Theta}_{\mathcal{T}^c}, \mathbf{S}\rangle + \lambda_1\|\mathbf{Z}_{\mathcal{T}^c}\|_1 + \lambda_2\|\mathbf{V}_{\mathcal{T}^c}\|_1 + \lambda_3(\|\mathbf{V} - \text{diag}(\mathbf{V})\|_{1,q} - \|\mathbf{V}_{\mathcal{T}} - \text{diag}(\mathbf{V}_{\mathcal{T}})\|_{1,q}) > 0. \tag{22}$$

Since $\|\mathbf{V} - \text{diag}(\mathbf{V})\|_{1,q} \geq \|\mathbf{V}_{\mathcal{T}} - \text{diag}(\mathbf{V}_{\mathcal{T}})\|_{1,q}$, it suffices to show that

$$\langle\boldsymbol{\Theta}_{\mathcal{T}^c}, \mathbf{S}\rangle + \lambda_1\|\mathbf{Z}_{\mathcal{T}^c}\|_1 + \lambda_2\|\mathbf{V}_{\mathcal{T}^c}\|_1 > 0. \tag{23}$$

Note that $\langle\boldsymbol{\Theta}_{\mathcal{T}^c}, \mathbf{S}\rangle = \langle\boldsymbol{\Theta}_{\mathcal{T}^c}, \mathbf{S}_{\mathcal{T}^c}\rangle$. By the sufficient condition, $S_{\max} < \lambda_1$ and $2S_{\max} < \lambda_2$.

In addition, we have that

$$
\begin{aligned}
|\langle \mathbf{\Theta}_{\mathcal{T}^c}, \mathbf{S} \rangle| &= |\langle \mathbf{\Theta}_{\mathcal{T}^c}, \mathbf{S}_{\mathcal{T}^c} \rangle| \\
&= |\langle \mathbf{V}_{\mathcal{T}^c} + \mathbf{V}_{\mathcal{T}^c}^T + \mathbf{Z}_{\mathcal{T}^c}, \mathbf{S}_{\mathcal{T}^c} \rangle| \\
&= |\langle 2\mathbf{V}_{\mathcal{T}^c} + \mathbf{Z}_{\mathcal{T}^c}, \mathbf{S}_{\mathcal{T}^c} \rangle| \qquad , \\
&\leq (2\|\mathbf{V}_{\mathcal{T}^c}\|_1 + \|\mathbf{Z}_{\mathcal{T}^c}\|_1)S_{\max} \\
&< \lambda_2 \|\mathbf{V}_{\mathcal{T}^c}\|_1 + \lambda_1 \|\mathbf{Z}_{\mathcal{T}^c}\|_1
\end{aligned}
\tag{24}
$$

where the last inequality follows from the sufficient condition. We have shown (23) as desired.

$\square$

## Proof of Theorem 2 (Necessary Condition)

We first present a simple lemma for proving Theorem 2. Throughout the proof of Theorem 2, $\|\cdot\|_\infty$ indicates the maximal absolute element of a matrix and $\|\cdot\|_{\infty,s}$ indicates the dual norm of $\|\cdot\|_{1,q}$.

**Lemma 4.** *The dual representation of $\tilde{\mathbf{P}}(\mathbf{\Theta})$ in (20) is*

$$
\begin{aligned}
\tilde{\mathbf{P}}^*(\mathbf{\Theta}) \;=\; &\max_{\mathbf{X},\mathbf{Y},\mathbf{\Lambda}} \quad \langle \mathbf{\Lambda}, \mathbf{\Theta} \rangle \\
&\text{subject to} \quad \mathbf{\Lambda} + \mathbf{\Lambda}^T = \hat{\lambda}_2 \mathbf{X} + \hat{\lambda}_3 \mathbf{Y} \\
&\qquad\qquad\quad \|\mathbf{X}\|_\infty \leq 1, \|\mathbf{\Lambda}\|_\infty \leq 1, \|\mathbf{Y}\|_{\infty,s} \leq 1 \\
&\qquad\qquad\quad X_{ii} = 0, Y_{ii} = 0, \Lambda_{ii} = 0 \; for \; i = 1, \ldots, p,
\end{aligned}
\tag{25}
$$

*where $\frac{1}{s} + \frac{1}{q} = 1$.*

*Proof.* We first state the dual representations for the norms in (20):

$$
\begin{aligned}
\|\mathbf{Z} - \mathrm{diag}(\mathbf{Z})\|_1 \;=\; &\max_{\mathbf{\Lambda}} \quad \langle \mathbf{\Lambda}, \mathbf{Z} \rangle \\
&\text{subject to} \quad \|\mathbf{\Lambda}\|_\infty \leq 1, \Lambda_{ii} = 0 \text{ for } i = 1, \ldots, p,
\end{aligned}
$$

$$
\begin{aligned}
\|\mathbf{V} - \mathrm{diag}(\mathbf{V})\|_1 \;=\; &\max_{\mathbf{X}} \quad \langle \mathbf{X}, \mathbf{V} \rangle \\
&\text{subject to} \quad \|\mathbf{X}\|_\infty \leq 1, X_{ii} = 0 \text{ for } i = 1, \ldots, p,
\end{aligned}
$$

$$
\begin{aligned}
\|\mathbf{V} - \mathrm{diag}(\mathbf{V})\|_{1,q} \;=\; &\max_{\mathbf{Y}} \quad \langle \mathbf{Y}, \mathbf{V} \rangle \\
&\text{subject to} \quad \|\mathbf{Y}\|_{\infty,s} \leq 1, Y_{ii} = 0 \text{ for } i = 1, \ldots, p.
\end{aligned}
$$

Then,

$$
\begin{aligned}
\tilde{\mathbf{P}}(\boldsymbol{\Theta}) \;=\; & \min_{\mathbf{V},\mathbf{Z}} && \|\mathbf{Z} - \operatorname{diag}(\mathbf{Z})\|_1 + \hat{\lambda}_2 \|\mathbf{V} - \operatorname{diag}(\mathbf{V})\|_1 + \hat{\lambda}_3 \|\mathbf{V} - \operatorname{diag}(\mathbf{V})\|_{1,q} \\
& \text{subject to} && \boldsymbol{\Theta} = \mathbf{Z} + \mathbf{V} + \mathbf{V}^T \\
=\; & \min_{\mathbf{V},\mathbf{Z}} && \max_{\boldsymbol{\Lambda},\mathbf{X},\mathbf{Y}} \langle \boldsymbol{\Lambda}, \mathbf{Z} \rangle + \hat{\lambda}_2 \langle \mathbf{X}, \mathbf{V} \rangle + \hat{\lambda}_3 \langle \mathbf{Y}, \mathbf{V} \rangle \\
& \text{subject to} && \|\boldsymbol{\Lambda}\|_\infty \le 1, \|\mathbf{X}\|_\infty \le 1, \|\mathbf{Y}\|_{\infty,s} \le 1 \\
& && \Lambda_{ii} = 0, X_{ii} = 0, Y_{ii} = 0 \text{ for } i = 1, \dots, p \\
& && \boldsymbol{\Theta} = \mathbf{Z} + \mathbf{V} + \mathbf{V}^T \\
=\; & \max_{\boldsymbol{\Lambda},\mathbf{X},\mathbf{Y}} && \min_{\mathbf{V},\mathbf{Z}} \langle \boldsymbol{\Lambda}, \mathbf{Z} \rangle + \hat{\lambda}_2 \langle \mathbf{X}, \mathbf{V} \rangle + \hat{\lambda}_3 \langle \mathbf{Y}, \mathbf{V} \rangle && (26) \\
& \text{subject to} && \|\boldsymbol{\Lambda}\|_\infty \le 1, \|\mathbf{X}\|_\infty \le 1, \|\mathbf{Y}\|_{\infty,s} \le 1 \\
& && \Lambda_{ii} = 0, X_{ii} = 0, Y_{ii} = 0 \text{ for } i = 1, \dots, p \\
& && \boldsymbol{\Theta} = \mathbf{Z} + \mathbf{V} + \mathbf{V}^T \\
=\; & \max_{\boldsymbol{\Lambda},\mathbf{X},\mathbf{Y}} && \langle \boldsymbol{\Lambda}, \boldsymbol{\Theta} \rangle \\
& \text{subject to} && \boldsymbol{\Lambda} + \boldsymbol{\Lambda}^T = \hat{\lambda}_2 \mathbf{X} + \hat{\lambda}_3 \mathbf{Y} \\
& && \|\mathbf{X}\|_\infty \le 1, \|\boldsymbol{\Lambda}\|_\infty \le 1, \|\mathbf{Y}\|_{\infty,s} \le 1 \\
& && X_{ii} = 0, Y_{ii} = 0, \Lambda_{ii} = 0 \text{ for } i = 1, \dots, p.
\end{aligned}
$$

The third equality holds since the constraints on $(\mathbf{V}, \mathbf{Z})$ and on $(\boldsymbol{\Lambda}, \mathbf{X}, \mathbf{Y})$ are both compact convex sets and so by the minimax theorem, we can swap max and min. The last equality follows from the fact that

$$
\begin{aligned}
& \min_{\mathbf{V},\mathbf{Z}} && \langle \boldsymbol{\Lambda}, \mathbf{Z} \rangle + \hat{\lambda}_2 \langle \mathbf{X}, \mathbf{V} \rangle + \hat{\lambda}_3 \langle \mathbf{Y}, \mathbf{V} \rangle \\
& \text{subject to} && \boldsymbol{\Theta} = \mathbf{Z} + \mathbf{V} + \mathbf{V}^T && (27) \\
& = && \begin{cases} \langle \boldsymbol{\Lambda}, \boldsymbol{\Theta} \rangle & \text{if } \boldsymbol{\Lambda} + \boldsymbol{\Lambda}^T = \hat{\lambda}_2 \mathbf{X} + \hat{\lambda}_3 \mathbf{Y} \\ -\infty & \text{otherwise.} \end{cases}
\end{aligned}
$$

$\square$

We now present the proof of Theorem 2.

*Proof.* The optimality condition for (21) is given by

$$
\mathbf{0} = -\boldsymbol{\Theta}^{-1} + \mathbf{S} + \lambda_1 \boldsymbol{\Lambda}, \tag{28}
$$

where $\boldsymbol{\Lambda}$ is a subgradient of $\tilde{\mathbf{P}}(\boldsymbol{\Theta})$ in (20) and the left-hand side of the above equation is a zero matrix of size $p \times p$.

Now suppose that $\boldsymbol{\Theta}^*$ that solves (28) is supported on $\mathcal{T}$, i.e., $\boldsymbol{\Theta}^*_{\mathcal{T}^c} = 0$. Then for any $(i,j) \in \mathcal{T}^c$, we have that

$$
0 = S_{ij} + \lambda_1 \Lambda^*_{ij}, \tag{29}
$$

where $\boldsymbol{\Lambda}^*$ is a subgradient of $\tilde{\mathbf{P}}(\boldsymbol{\Theta}^*)$. Note that $\boldsymbol{\Lambda}^*$ must be an optimal solution to the optimization problem (25). Therefore, it is also a feasible solution to (25), implying that

$$
\begin{aligned}
|\Lambda^*_{ij} + \Lambda^*_{ji}| &\le \hat{\lambda}_2 + \hat{\lambda}_3, \\
|\Lambda^*_{ij}| &\le 1.
\end{aligned}
$$

From Equation 29, we have that $\Lambda_{ij}^* = -\frac{S_{ij}}{\lambda_1}$ and thus,

$$\lambda_1 \geq \lambda_1 \max_{(i,j) \in \mathcal{T}^c} |\Lambda_{ij}^*|$$

$$= \lambda_1 \max_{(i,j) \in \mathcal{T}^c} \frac{|S_{ij}|}{\lambda_1}$$

$$= S_{\max}.$$

Also, recall that $\hat{\lambda}_2 = \frac{\lambda_2}{\lambda_1}$ and $\hat{\lambda}_3 = \frac{\lambda_3}{\lambda_1}$. We have that

$$\lambda_2 + \lambda_3 \geq \lambda_1 \max_{(i,j) \in \mathcal{T}^c} |\Lambda_{ij}^* + \Lambda_{ji}^*|$$

$$= \lambda_1 \max_{(i,j) \in \mathcal{T}^c} \frac{2|S_{ij}|}{\lambda_1}$$

$$= 2S_{\max}.$$

Hence, we obtain the desired result.

$\square$

# Appendix C: Some Properties of HGL

## Proof of Lemma 1

*Proof.* Let $(\boldsymbol{\Theta}^*, \mathbf{Z}^*, \mathbf{V}^*)$ be the solution to (8) and suppose that $\mathbf{Z}^*$ is not a diagonal matrix. Note that $\mathbf{Z}^*$ is symmetric since $\boldsymbol{\Theta} \in \mathcal{S}^{++}$, where $\mathcal{S}^{++}$ denotes the set of positive definite matrices. Let $\hat{\mathbf{Z}} = \text{diag}(\mathbf{Z}^*)$, a matrix that contains the diagonal elements of the matrix $\mathbf{Z}^*$. Also, construct $\hat{\mathbf{V}}$ as follows,

$$\hat{\mathbf{V}}_{ij} = \begin{cases} \mathbf{V}_{ij}^* + \frac{\mathbf{z}_{ij}^*}{2} & \text{if } i \neq j \\ \mathbf{V}_{jj}^* & \text{otherwise.} \end{cases}$$

Then, we have that $\boldsymbol{\Theta}^* = \hat{\mathbf{Z}} + \hat{\mathbf{V}} + \hat{\mathbf{V}}^T$. Thus, $(\boldsymbol{\Theta}^*, \hat{\mathbf{Z}}, \hat{\mathbf{V}})$ is a feasible solution to (8). We now show that $(\boldsymbol{\Theta}^*, \hat{\mathbf{Z}}, \hat{\mathbf{V}})$ has a smaller objective than $(\boldsymbol{\Theta}^*, \mathbf{Z}^*, \mathbf{V}^*)$ in (8), giving us a contradiction. Note that

$$\begin{aligned} \lambda_1 \|\hat{\mathbf{Z}} - \text{diag}(\hat{\mathbf{Z}})\|_1 + \lambda_2 \|\hat{\mathbf{V}} - \text{diag}(\hat{\mathbf{V}})\|_1 &= \lambda_2 \|\hat{\mathbf{V}} - \text{diag}(\hat{\mathbf{V}})\|_1 \\ &= \lambda_2 \sum_{i \neq j} |\mathbf{V}_{ij}^* + \frac{\mathbf{z}_{ij}^*}{2}| \\ &\leq \lambda_2 \|\mathbf{V}^* - \text{diag}(\mathbf{V}^*)\|_1 + \frac{\lambda_2}{2} \|\mathbf{Z}^* - \text{diag}(\mathbf{Z}^*)\|_1, \end{aligned}$$

and

$$\begin{aligned} \lambda_3 \sum_{j=1}^p \|(\hat{\mathbf{V}} - \text{diag}(\hat{\mathbf{V}}))_j\|_q &\leq \lambda_3 \sum_{j=1}^p \|(\mathbf{V}^* - \text{diag}(\mathbf{V}^*))_j\|_q + \frac{\lambda_3}{2} \sum_{j=1}^p \|(\mathbf{Z}^* - \text{diag}(\mathbf{Z}^*))_j\|_q \\ &\leq \lambda_3 \sum_{j=1}^p \|(\mathbf{V}^* - \text{diag}(\mathbf{V}^*))_j\|_q + \frac{\lambda_3}{2} \|\mathbf{Z}^* - \text{diag}(\mathbf{Z}^*)\|_1, \end{aligned}$$

where the last inequality follows from the fact that for any vector $\mathbf{x} \in \mathbb{R}^p$ and $q \geq 1$, $\|\mathbf{x}\|_q$ is a nonincreasing function of $q$ (Gentle 2007).

Summing up the above inequalities, we get that

$$
\begin{aligned}
\lambda_1\|\hat{\mathbf{Z}} - \mathrm{diag}(\hat{\mathbf{Z}})\|_1 + \lambda_2\|\hat{\mathbf{V}} - \mathrm{diag}(\hat{\mathbf{V}})\|_1 + \lambda_3 \sum_{j=1}^p \|(\hat{\mathbf{V}} - \mathrm{diag}(\hat{\mathbf{V}}))_j\|_q &\leq \\
\tfrac{\lambda_2+\lambda_3}{2}\|\mathbf{Z}^* - \mathrm{diag}(\mathbf{Z}^*)\|_1 + \lambda_2\|\mathbf{V}^* - \mathrm{diag}(\mathbf{V}^*)\|_1 + \lambda_3 \sum_{j=1}^p \|(\mathbf{V}^* - \mathrm{diag}(\mathbf{V}^*))_j\|_q &< \qquad (30) \\
\lambda_1\|\mathbf{Z}^* - \mathrm{diag}(\mathbf{Z}^*)\|_1 + \lambda_2\|\mathbf{V}^* - \mathrm{diag}(\mathbf{V}^*)\|_1 + \lambda_3 \sum_{j=1}^p \|(\mathbf{V}^* - \mathrm{diag}(\mathbf{V}^*))_j\|_q,
\end{aligned}
$$

where the last inequality uses the assumption that $\lambda_1 > \frac{\lambda_2+\lambda_3}{2}$. We arrive at a contradiction and therefore the result holds. $\qquad\square$

## Proof of Lemma 2

*Proof.* Let $(\mathbf{\Theta}^*, \mathbf{Z}^*, \mathbf{V}^*)$ be the solution to (8) and suppose $\mathbf{V}^*$ is not a diagonal matrix. Let $\hat{\mathbf{V}} = \mathrm{diag}(\mathbf{V}^*)$, a diagonal matrix that contains the diagonal elements of $\mathbf{V}^*$. Also construct $\hat{\mathbf{Z}}$ as follows,

$$
\hat{\mathbf{Z}}_{ij} = \begin{cases} \mathbf{Z}_{ij}^* + \mathbf{V}_{ij}^* + \mathbf{V}_{ji}^* & \text{if } i \neq j \\ \mathbf{Z}_{ij}^* & \text{otherwise.} \end{cases}
$$

Then, we have that $\mathbf{\Theta}^* = \hat{\mathbf{V}} + \hat{\mathbf{V}}^T + \hat{\mathbf{Z}}$. We now show that $(\mathbf{\Theta}^*, \hat{\mathbf{Z}}, \hat{\mathbf{V}})$ has a smaller objective value than $(\mathbf{\Theta}^*, \mathbf{Z}^*, \mathbf{V}^*)$ in (8), giving us a contradiction. We start by noting that

$$
\begin{aligned}
\lambda_1\|\hat{\mathbf{Z}} - \mathrm{diag}(\hat{\mathbf{Z}})\|_1 + \lambda_2\|\hat{\mathbf{V}} - \mathrm{diag}(\hat{\mathbf{V}})\|_1 &= \lambda_1\|\hat{\mathbf{Z}} - \mathrm{diag}(\hat{\mathbf{Z}})\|_1 \\
&\leq \lambda_1\|\mathbf{Z}^* - \mathrm{diag}(\mathbf{Z}^*)\|_1 + 2\lambda_1\|\mathbf{V}^* - \mathrm{diag}(\mathbf{V}^*)\|_1.
\end{aligned}
$$

By Holder's Inequality, we know that $\mathbf{x}^T\mathbf{y} \leq \|\mathbf{x}\|_q\|\mathbf{y}\|_s$ where $\frac{1}{s} + \frac{1}{q} = 1$ and $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{p-1}$. Setting $\mathbf{y} = \mathrm{sign}(\mathbf{x})$, we have that $\|\mathbf{x}\|_1 \leq (p-1)^{\frac{1}{s}}\|\mathbf{x}\|_q$. Consequently,

$$
\frac{\lambda_3}{(p-1)^{\frac{1}{s}}}\|\mathbf{V}^* - \mathrm{diag}(\mathbf{V}^*)\|_1 \leq \lambda_3 \sum_{j=1}^p \|(\mathbf{V}^* - \mathrm{diag}(\mathbf{V}^*))_j\|_q.
$$

Combining these results, we have that

$$
\begin{aligned}
&\lambda_1\|\hat{\mathbf{Z}} - \mathrm{diag}(\hat{\mathbf{Z}})\|_1 + \lambda_2\|\hat{\mathbf{V}} - \mathrm{diag}(\hat{\mathbf{V}})\|_1 + \lambda_3 \sum_{j=1}^p \|(\hat{\mathbf{V}} - \mathrm{diag}(\hat{\mathbf{V}}))_j\|_q \\
&\leq \lambda_1\|\mathbf{Z}^* - \mathrm{diag}(\mathbf{Z}^*)\|_1 + 2\lambda_1\|\mathbf{V}^* - \mathrm{diag}(\mathbf{V}^*)\|_1 \\
&< \lambda_1\|\mathbf{Z}^* - \mathrm{diag}(\mathbf{Z}^*)\|_1 + \left(\lambda_2 + \frac{\lambda_3}{(p-1)^{\frac{1}{s}}}\right)\|\mathbf{V}^* - \mathrm{diag}(\mathbf{V}^*)\|_1 \\
&\leq \lambda_1\|\mathbf{Z}^* - \mathrm{diag}(\mathbf{Z}^*)\|_1 + \lambda_2\|\mathbf{V}^* - \mathrm{diag}(\mathbf{V}^*)\|_1 + \lambda_3 \sum_{j=1}^p \|(\mathbf{V}^* - \mathrm{diag}(\mathbf{V}^*))_j\|_q,
\end{aligned}
$$

where we use the assumption that $\lambda_1 < \frac{\lambda_2}{2} + \frac{\lambda_3}{2(p-1)^{\frac{1}{s}}}$. This leads to a contradiction. $\qquad\square$

## Proof of Lemma 3

In this proof, we consider the case when $\lambda_1 > \frac{\lambda_2 + \lambda_3}{2}$. A similar proof technique can be used to prove the case when $\lambda_1 < \frac{\lambda_2 + \lambda_3}{2}$.

*Proof.* Let $f(\boldsymbol{\Theta}, \mathbf{V}, \mathbf{Z})$ denote the objective of (8) with $q = 1$, and $(\boldsymbol{\Theta}^*, \mathbf{V}^*, \mathbf{Z}^*)$ the optimal solution. By Lemma 1, the assumption that $\lambda_1 > \frac{\lambda_2 + \lambda_3}{2}$ implies that $\mathbf{Z}^*$ is a diagonal matrix. Now let $\hat{\mathbf{V}} = \frac{1}{2}\left(\mathbf{V}^* + (\mathbf{V}^*)^T\right)$. Then

$$
\begin{aligned}
f(\boldsymbol{\Theta}^*, \hat{\mathbf{V}}, \mathbf{Z}^*) &= -\log\det\boldsymbol{\Theta}^* + \langle\boldsymbol{\Theta}^*, \mathbf{S}\rangle + \lambda_1\|\mathbf{Z}^* - \operatorname{diag}(\mathbf{Z}^*)\|_1 + (\lambda_2 + \lambda_3)\|\hat{\mathbf{V}} - \operatorname{diag}(\hat{\mathbf{V}})\|_1 \\
&= -\log\det\boldsymbol{\Theta}^* + \langle\boldsymbol{\Theta}^*, \mathbf{S}\rangle + \frac{\lambda_2 + \lambda_3}{2}\|\mathbf{V}^* + \mathbf{V}^{*T} - \operatorname{diag}(\mathbf{V}^* + \mathbf{V}^{*T})\|_1 \\
&\leq -\log\det\boldsymbol{\Theta}^* + \langle\boldsymbol{\Theta}^*, \mathbf{S}\rangle + (\lambda_2 + \lambda_3)\|\mathbf{V}^* - \operatorname{diag}(\mathbf{V}^*)\|_1 \\
&= f(\boldsymbol{\Theta}^*, \mathbf{V}^*, \mathbf{Z}^*) \\
&\leq f(\boldsymbol{\Theta}^*, \hat{\mathbf{V}}, \mathbf{Z}^*),
\end{aligned}
\tag{31}
$$

where the last inequality follows from the assumption that $(\boldsymbol{\Theta}^*, \mathbf{V}^*, \mathbf{Z}^*)$ solves (8). By strict convexity of $f$, this means that $\mathbf{V}^* = \hat{\mathbf{V}}$, i.e., $\mathbf{V}^*$ is symmetric. This implies that

$$
\begin{aligned}
f(\boldsymbol{\Theta}^*, \mathbf{V}^*, \mathbf{Z}^*) &= -\log\det\boldsymbol{\Theta}^* + \langle\boldsymbol{\Theta}^*, \mathbf{S}\rangle + \frac{\lambda_2 + \lambda_3}{2}\|\mathbf{V}^* + \mathbf{V}^{*T} - \operatorname{diag}(\mathbf{V}^* + \mathbf{V}^{*T})\|_1 \\
&= -\log\det\boldsymbol{\Theta}^* + \langle\boldsymbol{\Theta}^*, \mathbf{S}\rangle + \frac{\lambda_2 + \lambda_3}{2}\|\boldsymbol{\Theta}^* - \operatorname{diag}(\boldsymbol{\Theta}^*)\|_1 \\
&= g(\boldsymbol{\Theta}^*),
\end{aligned}
\tag{32}
$$

where $g(\boldsymbol{\Theta})$ is the objective of the graphical lasso optimization problem, evaluated at $\boldsymbol{\Theta}$, with tuning parameter $\frac{\lambda_2 + \lambda_3}{2}$. Suppose that $\tilde{\boldsymbol{\Theta}}$ minimizes $g(\boldsymbol{\Theta})$, and $\boldsymbol{\Theta}^* \neq \tilde{\boldsymbol{\Theta}}$. Then, by (32) and strict convexity of $g$, $g(\boldsymbol{\Theta}^*) = f(\boldsymbol{\Theta}^*, \mathbf{V}^*, \mathbf{Z}^*) \leq f(\tilde{\boldsymbol{\Theta}}, \tilde{\boldsymbol{\Theta}}/2, \mathbf{0}) = g(\tilde{\boldsymbol{\Theta}}) < g(\boldsymbol{\Theta}^*)$, giving us a contradiction. Thus it must be that $\tilde{\boldsymbol{\Theta}} = \boldsymbol{\Theta}^*$.

$\square$

# References

Allen, G. & Liu, Z. (2012), 'A log-linear graphical model for inferring genetic networks from high-throughput sequencing data', *IEEE International Conference on Bioinformatics and Biomedicine* .

Barabási, A. (2009), 'Scale-free networks: A decade and beyond', *Science* **325**, 412–413.

Barabási, A. & Albert, R. (1999), 'Emergence of scaling in random networks', *Science* **286**, 509–512.

Bickel, P. & Levina, E. (2008), 'Regularized estimation of large covariance matrices', *Annals of Statistics* **36(1)**, 199–227.

Bien, J. & Tibshirani, R. (2011), 'Sparse estimation of a covariance matrix', *Biometrika* **98**(4), 807–820.

Boyd, S., Parikh, N., Chu, E., Peleato, B. & Eckstein, J. (2010), 'Distributed optimization and statistical learning via the ADMM', *Foundations and Trends in Machine Learning* **3**(1), 1–122.

Cai, T. & Liu, W. (2011), 'Adaptive thresholding for sparse covariance matrix estimation', *Journal of the American Statistical Association* **106**(494), 672–684.

Cardoso-Cachopo, A. (2009). "http://web.ist.utl.pt/acardoso/datasets/".

Chaudhuri, S., Drton, M. & Richardson, T. (2007), 'Estimation of a covariance matrix with zeros', *Biometrika* **94**, 199–216.

Danaher, P., Wang, P. & Witten, D. (2012), 'The joint graphical lasso for inverse covariance estimation across multiple classes', *arXiv:1111.0324* .

Defazio, A. & Caetano, T. (2012), 'A convex formulation for learning scale-free network via submodular relaxation', *NIPS* .

Drton, M. & Richardson, T. (2003), 'A new algorithm for maximum likelihood estimation in Gaussian graphical models for marginal independence.', *Proceedings of the 19th Conference on Uncertainty in Artificial Intelligence* pp. 184–191.

Drton, M. & Richardson, T. (2008), 'Graphical methods for efficient likelihood inference in Gaussian covariance models', *Journal of Machine Learning Research* **9**, 893–914.

Eckstein, J. (2012), 'Augmented Lagrangian and alternating direction methods for convex optimization: A tutorial and some illustrative computational results', *RUTCOR Research Reports* **32**.

Eckstein, J. & Bertsekas, D. (1992), 'On the Douglas-Rachford splitting method and the proximal point algorithm for maximal monotone operators', *Mathematical Programming* **55**(3, Ser. A), 293–318.

El Karoui, N. (2008), 'Operator norm consistent estimation of large-dimensional sparse covariance matrices', *The Annals of Statistics* **36**, 2717–2756.

Erdős, P. & Rényi, A. (1959), 'On random graphs I', *Publ. Math. Debrecen* **6**, 290–297.

Foygel, R. & Drton, M. (2010), 'Extended Bayesian information criteria for Gaussian graphical models', *NIPS* .

Friedman, J., Hastie, T. & Tibshirani, R. (2007), 'Sparse inverse covariance estimation with the graphical lasso', *Biostatistics* **9**, 432–441.

Gentle, J. E. (2007), *Matrix Algebra: Theory, Computations, and Applications in Statistics*, Springer, New York.

Guo, J., Levina, E., Michailidis, G. & Zhu, J. (2010), Joint structure estimation for categorical Markov networks. Submitted, available at http://www.stat.lsa.umich.edu/~elevina.

Hao, D., Ren, C. & Li, C. (2012), 'Revisiting the variation of clustering coefficient of biological networks suggests new modular structure', *BMC System Biology* **6:34**, 1–10.

Hero, A. & Rajaratnam, B. (2012), 'Hub discovery in partial correlation graphs', *IEEE Transactions on Information Theory* **58**, 6064–6078.

Höfling, H. & Tibshirani, R. (2009), 'Estimation of sparse binary pairwise Markov networks using pseudo-likelihoods', *Journal of Machine Learning Research* **10**, 883–906.

Horn, R. A. & Johnson, C. R. (1985), *Matrix analysis*, Cambridge University Press, New York, NY.

Jeong, H., Mason, S., Barabási, A. & Oltvai, Z. (2001), 'Lethality and centrality in protein networks', *Nature* **411**, 41–42.

Lee, S.-I., Ganapathi, V. & Koller, D. (2007), 'Efficient structure learning of Markov networks using $\ell_1$-regularization', *NIPS* pp. 817–824.

Liljeros, F., Edling, C., Amaral, L., Stanley, H. & Y., A. (2001), 'The web of human sexual contacts', *Nature* **411**, 907–908.

Liu, Q. & Ihler, A. (2011), 'Learning scale free networks by reweighed $\ell_1$ regularization', *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics* **15**, 40–48.

Ma, S., Xue, L. & Zou, H. (2013), 'Alternating direction methods for latent variable Gaussian graphical model selection', *Neural Computation* .

Maglott et al. (2004), 'Entrez Gene: gene-centered information at NCBI', *Nucleic Acids Research* **33(D)**, 54–58.

Mardia, K., Kent, J. & Bibby, J. (1979), *Multivariate Analysis*, Academic Press.

Mazumder, R. & Hastie, T. (2012), 'Exact covariance thresholding into connected components for large-scale graphical lasso', *Journal of Machine Learning Research* **13**, 781–794.

Meinshausen, N. & Bühlmann, P. (2006), 'High dimensional graphs and variable selection with the lasso', *Annals of Statistics* **34**, 1436–1462.

Meinshausen, N. & Bühlmann, P. (2010), 'Stability selection (with discussion)', *Journal of the Royal Statistical Society, Series B* **72**, 417–473.

Mohan, K., London, P., Fazel, M., Witten, D. & Lee, S.-I. (2013), 'Node-based learning of Gaussian graphical models', *arXiv* .

Newman, M. (2000), 'The structure of scientific collaboration networks', *PNAS* **98**, 404–409.

Peng, J., Wang, P., Zhou, N. & Zhu, J. (2009), 'Partial correlation estimation by joint sparse regression model', *Journal of the American Statistical Association* **104(486)**, 735–746.

Rappaport et al. (2013), 'MalaCards: an integrated compendium for diseases and their annotation', *Database (Oxford)* .

Ravikumar, P., Wainwright, M. & Lafferty, J. (2010), 'High-dimensional Ising model selection using $\ell_1$-regularized logistic regression', *The Annals of Statistics* **38**(3), 1287–1319.

Ravikumar, P., Wainwright, M., Raskutti, G. & Yu, B. (2011), 'High-dimensional covariance estimation by minimizing $\ell_1$-penalized log-determinant divergence', *Electronic Journal of Statistics* **5**, 935–980.

Rothman, A., Bickel, P., Levina, E. & Zhu, J. (2008), 'Sparse permutation invariant covariance estimation', *Electronic Journal of Statistics* **2**, 494–515.

Rothman, A., Levina, E. & Zhu, J. (2009), 'Generalized thresholding of large covariance matrices', *Journal of the American Statistical Association* **104**, 177–186.

Simon, N., Friedman, J., Hastie, T. & Tibshirani, R. (2012), 'A sparse-group lasso', *Journal of Computational and Graphical Statistics* .

Tarjan, R. (1972), 'Depth-first search and linear graph algorithms', *SIAM journal on computing* **1**(2), 146–160.

Verhaak et al. (2010), 'Integrated genomic analysis identifies clinically relevant subtypes of glioblastoma characterized by abnormalities in PDGFRA, IDH1, EGFR, and NF1', *Cancer Cell* **17(1)**, 98–110.

Witten, D., Friedman, J. & Simon, N. (2011), 'New insights and faster computations for the graphical lasso', *Journal of Computational and Graphical Statistics* **20(4)**, 892–900.

Xue, L., Ma, S. & Zou, H. (2012), 'Positive definite $\ell_1$ penalized estimation of large covariance matrices', *Journal of the American Statistical Association* **107**(500), 1480–1491.

Yang, E., Ravikumar, P., Allen, G. & Liu, Z. (2012), 'Graphical models via generalized linear models', *NIPS* .

Yuan, M. (2008), 'Efficient computation of $\ell_1$ regularized estimates in Gaussian graphical models', *Journal of Computational and Graphical Statistics* **17(4)**, 809–826.

Yuan, M. & Lin, Y. (2007*a*), 'Model selection and estimation in regression with grouped variables', *Journal of the Royal Statistical Society, Series B* **68**, 49–67.

Yuan, M. & Lin, Y. (2007*b*), 'Model selection and estimation in the Gaussian graphical model', *Biometrika* **94(10)**, 19–35.