

---

# Large-Scale Sparse Inverse Covariance Estimation via Thresholding and Max-Det Matrix Completion

---

Richard Y. Zhang<sup>1</sup> Salar Fattahi<sup>1</sup> Somayeh Sojoudi<sup>2</sup>

## Abstract

The sparse inverse covariance estimation problem is commonly solved using an  $\ell_1$ -regularized Gaussian maximum likelihood estimator known as “graphical lasso”, but its computational cost becomes prohibitive for large data sets. A recent line of results showed—under mild assumptions—that the graphical lasso estimator can be retrieved by soft-thresholding the sample covariance matrix and solving a maximum determinant matrix completion (MDMC) problem. This paper proves an extension of this result, and describes a Newton-CG algorithm to efficiently solve the MDMC problem. Assuming that the thresholded sample covariance matrix is sparse with a sparse Cholesky factorization, we prove that the algorithm converges to an  $\epsilon$ -accurate solution in  $O(n \log(1/\epsilon))$  time and  $O(n)$  memory. The algorithm is highly efficient in practice: we solve the associated MDMC problems with as many as 200,000 variables to 7-9 digits of accuracy in less than an hour on a standard laptop computer running MATLAB.

## 1. Introduction

Consider the problem of estimating an  $n \times n$  covariance matrix  $\Sigma$  (or its inverse  $\Sigma^{-1}$ ) of a  $n$ -variate probability distribution from  $N$  independently and identically distributed samples  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$  drawn from the same probability distribution. In applications spanning from computer vision, natural language processing, to economics (Li, 1994; Manning & Schütze, 1999; Durlauf, 1993), the matrix  $\Sigma^{-1}$  is often *sparse*, meaning that its matrix elements are mostly

zero. For Gaussian distributions, the statistical interpretation of sparsity in  $\Sigma^{-1}$  is that most of the variables are pairwise conditionally independent (Meinshausen & Bühlmann, 2006; Yuan & Lin, 2007; Friedman et al., 2008; Banerjee et al., 2008).

Imposing sparsity upon  $\Sigma^{-1}$  can regularize the associated estimation problem and greatly reduce the number of samples required. This is particularly important in high-dimensional settings where  $n$  is large, often significantly larger than the number of samples  $N \ll n$ . One popular approach regularizes the associated maximum likelihood estimation (MLE) problem by a sparsity-promoting  $\ell_1$  term, as in

$$\underset{X \succ 0}{\text{minimize}} \operatorname{tr} CX - \log \det X + \lambda \sum_{i=1}^n \sum_{j=1}^n |X_{i,j}|. \quad (1)$$

Here,  $C = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$  is the sample covariance matrix with sample mean  $\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$ , and  $X$  is the resulting estimator for  $\Sigma^{-1}$ . This approach, commonly known as the *graphical lasso* (Friedman et al., 2008), is known to enjoy a number of statistical guarantees (Rothman et al., 2008; Ravikumar et al., 2011), some of which are direct extensions of earlier work on the classical lasso (Obozinski et al., 2008; Negahban & Wainwright, 2008; Wainwright, 2009; Huang & Zhang, 2010). A variation on this theme is to only impose the  $\ell_1$  penalty on the off-diagonal elements of  $X$ , or to place different weights  $\lambda$  on the elements of the matrix  $X$ , as in the classical weighted lasso.

While the  $\ell_1$ -regularized problem (1) is technically convex, it is commonly considered intractable for large-scale datasets. The decision variable is an  $n \times n$  matrix, so simply fitting all  $O(n^2)$  variables into memory is already a significant issue. General-purpose algorithms have either prohibitively high complexity or slow convergence. In practice, (1) is solved using problem-specific algorithms. The state-of-the-art include GLASSO (Friedman et al., 2008), QUIC (Hsieh et al., 2014), and its “big-data” extension BIG-QUIC (Hsieh et al., 2013). These algorithms use between  $O(n)$  and  $O(n^3)$  time and between  $O(n^2)$  and  $O(n)$  memory per iteration, but the number of iterations needed to converge to an accurate solution can be very large.

---

MATLAB source code: <http://alum.mit.edu/www/ryz>. <sup>1</sup>Department of Industrial Engineering and Operations Research, University of California, Berkeley, USA. <sup>2</sup>Department of Electrical Engineering and Computer Science, University of California, Berkeley, USA.. Correspondence to: R.Y. Zhang <ryz@berkeley.edu>.

### 1.1. Graphical lasso, soft-thresholding, and MDMC

The high practical cost of graphical lasso has inspired a number of heuristics, which enjoy less guarantees but are significantly cheaper to use. Indeed, heuristics are often the only viable option once  $n$  exceeds the order of a few tens of thousands.

One simple idea is to threshold the sample covariance matrix  $C$ : to examine all of its elements and keep only the ones whose magnitudes exceed some threshold. In a recent line of work (Mazumder & Hastie, 2012; Sojoudi, 2016; Fattahi & Sojoudi, 2017; Fattahi et al., 2018), this simple heuristic was shown to enjoy some surprising guarantees. In particular, (Sojoudi, 2016; Fattahi & Sojoudi, 2017) proved that when the lasso weight is imposed over only the off-diagonal elements of  $X$  that—under some assumptions—the *sparsity pattern* of the associated graphical lasso estimator can be recovered by performing a soft-thresholding operation on  $C$ , as in

$$(C_\lambda)_{i,j} = \begin{cases} C_{i,j} & i = j, \\ C_{i,j} - \lambda & C_{i,j} > \lambda, i \neq j, \\ 0 & |C_{i,j}| \leq \lambda, i \neq j, \\ C_{i,j} + \lambda & -\lambda \leq C_{i,j}, i \neq j, \end{cases} \quad (2)$$

and recovering the sparsity pattern

$$G = \{(i, j) \in \{1, \dots, n\}^2 : (C_\lambda)_{i,j} \neq 0\}. \quad (3)$$

The associated graph (also denoted as  $G$  when there is no ambiguity) is obtained by viewing each nonzero element  $(i, j)$  in  $G$  as an edge between the  $i$ -th and  $j$ -th vertex in an undirected graph on  $n$  nodes. Moreover, they showed that the estimator  $X$  can be recovered by solving a version of (1) in which the sparsity pattern  $G$  is explicitly imposed, as in

$$\begin{aligned} & \underset{X \succ 0}{\text{minimize}} \quad \text{tr } C_\lambda X - \log \det X & (4) \\ & \text{subject to } X_{i,j} = 0 \quad \forall (i, j) \notin G. \end{aligned}$$

Recovering the exact value of  $X$  (and not just its sparsity pattern) is important because it provides a shrinkage MLE when the true MLE is ill-defined; for Gaussian fields, its nonzero values encode the partial correlations between variables. Problem (4) is named the *maximum determinant matrix completion* (MDMC) in the literature, for reasons explained below. The problem has a recursive closed-form solution whenever the graph of  $G$  is *acyclic* (i.e. a tree or forest) (Fattahi & Sojoudi, 2017), or more generally, if it is *chordal* (Fattahi et al., 2018). It is worth emphasizing that the closed-form solution is *extremely* fast to evaluate: a chordal example in (Fattahi et al., 2018) with 13,659 variables took just  $\approx 5$  seconds to solve on a laptop computer.

The assumptions needed for graphical lasso to be equivalent to thresholding are hard to check but relatively mild.

Indeed, (Fattahi & Sojoudi, 2017) proves that they are automatically satisfied whenever  $\lambda$  is sufficiently large relative to the sample covariance matrix. Their numerical study found “sufficiently large” to be a fairly loose criterion in practice, particularly in view of the fact that large values of  $\lambda$  are needed to induce a sufficiently sparse estimate of  $\Sigma^{-1}$ , e.g. with  $\approx 10n$  nonzero elements.

However, the requirement for  $G$  to be chordal is very strong. Aside from trivial chordal graphs like trees and cliques, thresholding will produce a chordal graph with probability zero. When  $G$  is nonchordal, no closed-form solution exists, and one must resort to an iterative algorithm. The state-of-the-art for nonchordal MDMC is to embed the nonchordal graph within a chordal graph, and to solve the resulting problem as a semidefinite program using an interior-point method.

### 1.2. Main results

The purpose of this paper is two-fold. First, we derive an extension of the guarantees derived by (Mazumder & Hastie, 2012; Sojoudi, 2016; Fattahi & Sojoudi, 2017; Fattahi et al., 2018) for a slightly more general version of the problem that we call *restricted* graphical lasso (RGL):

$$\begin{aligned} \hat{X} = \underset{X \succ 0}{\text{minimize}} \quad & \text{tr } CX - \log \det X & (5) \\ & + \sum_{i=1}^n \sum_{j=i+1}^n \lambda_{i,j} |X_{i,j}|. \\ \text{subject to } & X_{i,j} = 0 \quad \forall (i, j) \notin H. \end{aligned}$$

In other words, RGL is (1) penalized by a weighted lasso penalty  $\lambda_{i,j}$  on the off-diagonals, and with an *a priori* sparsity pattern  $H$  imposed as an additional constraint. We use the sparsity pattern  $H$  to incorporate prior information on the structure of the graphical model. For example, if the sample covariance  $C$  is collected over a graph, such as a communication system or a social network, then far-away variables can be assumed as pairwise conditionally independent (Park & Rilett, 1999; Honorio et al., 2009; Croft et al., 2010). Including these neighborhood relationships into  $H$  can regularize the statistical problem, as well as reduce the numerical cost for a solution.

In Section 2, we describe a procedure to transform RGL (5) into MDMC (4), in the same style as prior results by (Fattahi & Sojoudi, 2017; Fattahi et al., 2018) for graphical lasso. More specifically, we soft-threshold the sample covariance  $C$  and then project this matrix onto the sparsity pattern  $H$ . We give conditions for the resulting sparsity pattern to be equivalent to the one obtained by solving (5). Furthermore, we prove that the resulting estimator  $X$  can be recovered by solving the same MDMC problem (4) with  $C_\lambda$  appropriately modified.

The second purpose is to describe an efficient algorithm to solve MCDC when the graph  $G$  is *nonchordal*, based on the chordal embedding approach of (Dahl et al., 2008; Andersen et al., 2010; 2013b). We embed  $G$  within a chordal  $\tilde{G} \supset G$ , to result in a convex optimization problem over  $\mathbb{S}_{\tilde{G}}^n$ , the space of real symmetric matrices with sparsity pattern  $\tilde{G}$ . This way, the constraint  $X \in \mathbb{S}_{\tilde{G}}^n$  is *implicitly* imposed, meaning that we simply ignore the nonzero elements not in  $\tilde{G}$ . Next, we solve an optimization problem on  $\mathbb{S}_{\tilde{G}}^n$  using a custom Newton-CG method. The main idea is to use an inner conjugate gradients (CG) loop to solve the Newton subproblem of an outer Newton’s method. The actual algorithm has a number of features designed to exploit problem structure, including the sparse chordal property of  $\tilde{G}$ , duality, and the ability for CG and Newton to converge superlinearly; these are outlined in Section 3.

Assuming that the chordal embedding is sparse with  $|\tilde{G}| = O(n)$  nonzero elements, we prove in Section 3.4, that our algorithm converges to an  $\epsilon$ -accurate solution of MDMC (4) in

$$O(n \cdot \log \epsilon^{-1} \cdot \log \log \epsilon^{-1}) \text{ time and } O(n) \text{ memory.} \quad (6)$$

Most importantly, the algorithm is highly efficient in practice. In Section 4, we present computation results on a suite of test cases. Both synthetic and real-life graphs are considered. Using our approach, we solve sparse inverse covariance estimation problems containing as many as 200,000 variables, in less than an hour on a laptop computer.

### 1.3. Related Work

**Graphical lasso with prior information.** A number of approaches are available in the literature to introduce prior information to graphical lasso. The weighted version of graphical lasso mentioned before is an example, though RGL will generally be more efficient to solve due to a reduction in the number of variables. (Egilmez et al., 2017) introduced a class of graphical lasso in which the true graphical model is assumed to have Laplacian structure. This structure commonly appears in signal and image processing (Milanfar, 2013). For the *a priori* graph-based correlation structure described above, (Grechkin et al., 2015) introduced a *pathway* graphical lasso method similar to RGL.

**Algorithms for graphical lasso.** Algorithms for graphical lasso are usually based on some mixture of Newton (Oztoprak et al., 2012), proximal Newton (Hsieh et al., 2013; 2014), iterative thresholding (Rolfs et al., 2012), and (block) coordinate descent (Friedman et al., 2008; Treister & Turek, 2014). All of these suffer fundamentally from the need to keep track and act on all  $O(n^2)$  elements in the matrix  $X$  decision variable. Even if the final solution matrix were sparse with  $O(n)$  nonzeros, it is still possible for the

algorithm to traverse through a “dense region” in which the iterate  $X$  must be fully dense. Thresholding heuristics have been proposed to address issue, but these may adversely affect the outer algorithm and prevent convergence. It is generally impossible to guarantee a figure lower than  $O(n^2)$  time per-iteration, even if the solution contains only  $O(n)$  nonzeros. Most of the algorithms mentioned above actually have worst-case per-iteration costs of  $O(n^3)$ .

**Graphical lasso via thresholding.** The elementary estimator for graphical models (EE-GM) (Yang et al., 2014) is another thresholding-based low-complexity method that is able to recover the actual graphical lasso estimator. Both EE-GM and our algorithm have a similar level of performance in practice, because both algorithm are bottlenecked by the initial thresholding step, which is a quadratic  $O(n^2)$  time operation.

**Algorithms for MDMC.** Our algorithm is inspired by a line of results (Dahl et al., 2008; Andersen et al., 2010; 2013b; Li et al., 2017) for minimizing the log-det penalty on chordal sparsity patterns, culminating in the CVXOPT package (Andersen et al., 2013a). These algorithms all solve the Newton subproblem by explicitly forming and factoring the fully-dense Newton matrix in  $O(nm^2 + m^3)$  time, where  $m = |\tilde{G} \setminus G|$  is the number of edges added during chordal embedding. By comparison, our algorithm solves the Newton subproblem iteratively using CG, in  $O(n + m)$  time to machine precision (see Section 3.4).

### Notations

Let  $\mathbb{R}^n$  and  $\mathbb{S}^n$  be the set of  $n \times 1$  real vectors, and  $n \times n$  real symmetric matrices. We endow  $\mathbb{S}^n$  with the usual matrix inner product  $X \bullet Y = \text{tr } XY$  and Euclidean (i.e. Frobenius) norm  $\|X\|_F^2 = X \bullet X$ . Let  $\mathbb{S}_+^n \subset \mathbb{S}^n$  and  $\mathbb{S}_{++}^n \subset \mathbb{S}_+^n$  be the associated set of positive semidefinite and positive definite matrices. We will frequently write  $X \succeq 0$  to mean  $X \in \mathbb{S}_+^n$  and write  $X \succ 0$  to mean  $X \in \mathbb{S}_{++}^n$ . Given a sparsity pattern  $G$ , we define  $\mathbb{S}_G^n \subseteq \mathbb{S}^n$  as the set of  $n \times n$  real symmetric matrices with this sparsity pattern.

## 2. Restricted graphical lasso, soft-thresholding, and MDMC

Let  $P_H(X)$  denote the projection operator from  $\mathbb{S}^n$  onto  $\mathbb{S}_H^n$ , i.e. by setting all  $X_{i,j} = 0$  if  $(i, j) \notin H$ . Let  $C_\lambda$  be the sample covariance matrix  $C$  individually soft-thresholded by  $[\lambda_{i,j}]$ , as in

$$(C_\lambda)_{i,j} = \begin{cases} C_{i,j} & i = j, \\ C_{i,j} - \lambda_{i,j} & C_{i,j} > \lambda_{i,j}, i \neq j, \\ 0 & |C_{i,j}| \leq \lambda_{i,j} i \neq j, \\ C_{i,j} + \lambda_{i,j} & -\lambda_{i,j} \leq C_{i,j} i \neq j, \end{cases} \quad (7)$$

In this section, we state the conditions for  $P_H(C_\lambda)$ —the projection of the soft-thresholded matrix  $C_\lambda$  in (7) onto  $H$ —to have the same sparsity pattern as the RGL estimator  $\hat{X}$  in (5). Furthermore, the estimator  $\hat{X}$  can be explicitly recovered by solving the MDMC problem (4) while replacing  $C_\lambda \leftarrow P_H(C_\lambda)$  and  $G \leftarrow P_H(G)$ . For brevity, all proofs and remarks are omitted; these can be found in the supplementary materials.

Before we state the exact conditions, we begin by adopting the some definitions and notations from the literature.

**Definition 1.** (Fattahi & Sojoudi, 2017) Given a matrix  $M \in \mathbb{S}^n$ , define  $G_M = \{(i, j) : M_{i,j} \neq 0\}$  as its sparsity pattern. Then  $M$  is called **inverse-consistent** if there exists a matrix  $N \in \mathbb{S}^n$  such that

$$M + N \succ 0 \quad (8a)$$

$$N = 0 \quad \forall (i, j) \in G_M \quad (8b)$$

$$(M + N)^{-1} \in \mathbb{S}_{G_M}^n \quad (8c)$$

The matrix  $N$  is called an **inverse-consistent complement** of  $M$  and is denoted by  $M^{(c)}$ . Furthermore,  $M$  is called **sign-consistent** if for every  $(i, j) \in G_M$ , the  $(i, j)$ -th elements of  $M$  and  $(M + M^{(c)})^{-1}$  have opposite signs.

Moreover, we take the usual matrix max-norm to exclude the diagonal, as in  $\|M\|_{\max} = \max_{i \neq j} |M_{i,j}|$ , and adopt the  $\beta(G, \alpha)$  function defined with respect to the sparsity pattern  $G$  and scalar  $\alpha > 0$

$$\begin{aligned} \beta(G, \alpha) &= \max_{M \succ 0} \|M^{(c)}\|_{\max} \\ \text{s.t. } M &\in \mathbb{S}_G^n \text{ and } \|M\|_{\max} \leq \alpha \\ M_{i,i} &= 1 \quad \forall i \in \{1, \dots, n\} \\ M &\text{ is inverse-consistent.} \end{aligned}$$

We are now ready to state the conditions for soft-thresholding to be equivalent to RGL.

**Theorem 2.** Define  $C_\lambda$  as in (7), define  $C_H = P_H(C_\lambda)$  and let  $G_H = \{(i, j) : (C_H)_{i,j} \neq 0\}$  be its sparsity pattern. Then  $G_H$  coincides with sparsity pattern of the optimal solution  $\hat{X}$  of RGL (5) if the normalized matrix  $\tilde{C} = D^{-1/2} C_H D^{-1/2}$  where  $D = \text{diag}(C_H)$  satisfies the following conditions:

1.  $\tilde{C}$  is positive definite,
2.  $\tilde{C}$  is sign-consistent,
3. Let  $\beta_H = \beta(G_H, \|\tilde{C}\|_{\max})$ . Then

$$\beta_H \leq \min_{(k,l) \notin G_H} \frac{\lambda_{k,l} - |(C_H)_{k,l}|}{\sqrt{(C_H)_{k,k} \cdot (C_H)_{l,l}}} \quad (9)$$

*Proof.* See supplementary materials.  $\square$

Theorem 2 leads to the following corollary, which asserts that the optimal solution of RGL can be obtained by *maximum determinant matrix completion*: computing the matrix  $Z \succeq 0$  with the largest determinant that “fills-in” the zero elements of  $P_H(C_\lambda)$ .

**Corollary 3.** Suppose that the conditions in Theorem 2 are satisfied. Define  $\hat{Z}$  as the solution to the following

$$\begin{aligned} \hat{Z} &= \underset{Z \succeq 0}{\text{maximize}} \log \det Z \quad (10) \\ &\text{subject to } Z_{i,j} = P_H(C_\lambda) \text{ for all } (i, j) \\ &\quad \text{where } [P_H(C_\lambda)]_{i,j} \neq 0 \end{aligned}$$

Then  $\hat{Z} = \hat{X}^{-1}$ , where  $\hat{X}$  is the solution of (5).

*Proof.* See supplementary materials.  $\square$

Standard manipulations show that (10) is the Lagrangian dual of (4), thus explaining the etymology of (4) as MDMC.

### 3. Proposed Algorithm

This section describes an efficient algorithm to solve MDMC (4) in which the sparsity pattern  $G$  is *nonchordal*. If we assume that the input matrix  $C_\lambda$  is sparse, and that sparse Cholesky factorization is able to solve  $C_\lambda x = b$  in  $O(n)$  time, then our algorithm is guaranteed to compute an  $\epsilon$ -accurate solution in  $O(n \log \epsilon^{-1})$  time and  $O(n)$  memory.

The algorithm is fundamentally a Newton-CG method, i.e. Newton’s method in which the Newton search directions are computed using conjugate gradients (CG). It is developed from four key insights:

**1. Chordal embedding is easy via sparse matrix heuristics.** State-of-the-art algorithms for (4) begin by computing a chordal embedding  $\tilde{G}$  for  $G$ . The optimal chordal embedding with the fewest number of nonzeros  $|\tilde{G}|$  is NP-hard to compute, but a good-enough embedding with  $O(n)$  nonzeros is sufficient for our purposes. Computing a good  $\tilde{G}$  with  $|\tilde{G}| = O(n)$  is exactly the same problem as finding a sparse Cholesky factorization  $C_\lambda = LL^T$  with  $O(n)$  fill-in. Using heuristics developed for numerical linear algebra, we are able to find sparse chordal embeddings for graphs containing millions of edges and hundreds of thousands of nodes in seconds.

**2. Optimize directly on the sparse matrix cone.** Using log-det barriers for sparse matrix cones (Dahl et al., 2008; Andersen et al., 2010; 2013b; Vandenberghe et al., 2015), we can optimize directly in the space  $\mathbb{S}_{\tilde{G}}^n$ , while ignoring all matrix elements outside of  $\tilde{G}$ . If  $|\tilde{G}| = O(n)$ , then only  $O(n)$  decision variables must be explicitly optimized.

Moreover, each function evaluation, gradient evaluation, and matrix-vector product with the Hessian can be performed in  $O(n)$  time, using the numerical recipes in (Andersen et al., 2013b).

**3. The dual is easier to solve than the primal.** The primal problem starts with a feasible point  $X \in \mathbb{S}_G^n$  and seeks to achieve first-order optimality. The dual problem starts with an infeasible optimal point  $X \notin \mathbb{S}_G^n$  satisfying first-order optimality, and seeks to make it feasible. Feasibility is easier to achieve than optimality, so the dual problem is easier to solve than the primal.

**4. Conjugate gradients (CG) converges in  $O(1)$  iterations.** Our main result (Theorem 6) bounds the condition number of the Newton subproblem to be  $O(1)$ , independent of the problem dimension  $n$  and the current accuracy  $\epsilon$ . It is therefore cheaper to solve this subproblem using CG to machine precision  $\delta_{\text{mach}}$  in  $O(n \log \delta_{\text{mach}}^{-1})$  time than it is to solve for it directly in  $O(nm^2 + m^3)$  time using Cholesky factorization (Dahl et al., 2008; Andersen et al., 2010; 2013b). Moreover, CG is an optimal Krylov subspace method, and as such, it is often able to exploit clustering in the eigenvalues to converge superlinearly. Finally, computing the Newton direction to high accuracy further allows the outer Newton method to also converge quadratically.

The remainder of this section describes each consideration in further detail. We state the algorithm explicitly in Section 3.5.

### 3.1. Efficient chordal embedding

Following (Dahl et al., 2008), we begin by reformulating (4) into a sparse chordal matrix program

$$\begin{aligned} \hat{X} = \text{minimize } & \text{tr } CX - \log \det X & (11) \\ \text{subject to } & X_{i,j} = 0 \quad \forall (i,j) \in \tilde{G} \setminus G. \\ & X \in \mathbb{S}_{\tilde{G}}^n. \end{aligned}$$

in which  $\tilde{G}$  is a *chordal embedding* for  $G$ : a sparsity pattern  $\tilde{G} \supset G$  whose graph contains no induced cycles greater than three. This can be implemented using standard algorithms for large-and-sparse linear equations, due to the following result.

**Proposition 4.** *Let  $C \in \mathbb{S}_G^n$  be a positive definite matrix with sparsity pattern  $G$ . Compute its unique lower-triangular Cholesky factor  $L$  satisfying  $C = LL^T$ . Ignoring perfect numerical cancellation, the sparsity pattern of  $L + L^T$  is a chordal embedding  $\tilde{G} \supset G$ .*

*Proof.* The original proof is due to (Rose, 1970); see also (Vandenberghe et al., 2015).  $\square$

Note that  $\tilde{G}$  can be determined directly from  $G$  using a

```
p = amd(C); % fill-reducing ordering
[h,~,~,~,R] = symbfact(C(p,p));
Gt = R+R'; Gt(p,p) = Gt;
m = nnz(R) - nnz(tril(C));
```

Figure 1. MATLAB code for chordal embedding via its internal approximate minimum degree ordering. Given a sparse matrix  $(C)$ , compute a chordal embedding  $(Gt)$  and the number of added edges  $(m)$ .

*symbolic* Cholesky algorithm, which simulates the steps of Gaussian elimination using Boolean logic. Moreover, we can substantially reduce the number of edges added to  $G$  by reordering the columns and rows of  $C$  using a *fill-reducing ordering*.

**Corollary 5.** *Let  $\Pi$  be a permutation matrix. For the same  $C \in \mathbb{S}_G^n$  in Proposition 4, compute the unique Cholesky factor satisfying  $\Pi C \Pi^T = LL^T$ . Ignoring perfect numerical cancellation, the sparsity pattern of  $\Pi(L + L^T)\Pi^T$  is a chordal embedding  $\tilde{G} \supset G$ .*

The problem of finding the best choice of  $\Pi$  is known as the *fill-minimizing* problem, and is NP-complete (Yan-nakakis, 1981). However, good orderings are easily found using heuristics developed for numerical linear algebra, like minimum degree ordering (George & Liu, 1989) and nested dissection (Gilbert, 1988; Agrawal et al., 1993). If  $G$  admits sparse chordal embeddings, then a good-enough  $|\tilde{G}| = O(n)$  will usually be found using a simple minimum degree ordering; see the MATLAB code snippet in Figure 1.

### 3.2. Logarithmic barriers for sparse matrix cones

Define the cone of *sparse positive semidefinite matrices*  $\mathcal{K}$ , and the cone of *sparse matrices with positive semidefinite completions*  $\mathcal{K}_*$ , as the following

$$\begin{aligned} \mathcal{K} &= \mathbb{S}_+^n \cap \mathbb{S}_{\tilde{G}}^n, \\ \mathcal{K}_* &= \{S \bullet X \geq 0 : S \in \mathbb{S}_{\tilde{G}}^n\} = P_{\tilde{G}}(\mathbb{S}_+^n). \end{aligned}$$

Then (11) can be posed as the primal-dual pair:

$$\arg \min_{X \in \mathcal{K}} \{C \bullet X + f(X) : A^T(X) = 0\}, \quad (12)$$

$$\arg \max_{S \in \mathcal{K}_*, y \in \mathbb{R}^m} \{-f_*(S) : S = C - A(y)\}, \quad (13)$$

where the linear map  $A : \mathbb{R}^m \rightarrow \mathbb{S}_{\tilde{G} \setminus G}^n$  converts a list of  $m$  variables into the corresponding matrix in  $\tilde{G} \setminus G$ , and  $f$  and  $f_*$  are the “log-det” barrier functions on  $\mathcal{K}$  and  $\mathcal{K}_*$  as introduced by (Dahl et al., 2008; Andersen et al., 2010; 2013b)

$$f(X) = -\log \det X, \quad f_*(S) = -\min_{X \in \mathcal{K}} \{S \bullet X + f(X)\}.$$

Assuming that  $\tilde{G}$  is *sparse* and *chordal*, the functions  $f$  and  $f_*$ , their gradient evaluations, and Hessian matrix-vector products can all be efficiently evaluated in  $O(n)$  time and  $O(n)$  memory, using the numerical recipes described in (Andersen et al., 2013b).

### 3.3. Solving the dual problem

Our algorithm actually solves the dual problem (13), which can be rewritten as an unconstrained optimization problem

$$\hat{y} \equiv \arg \min_{y \in \mathbb{R}^m} g(y) \equiv f_*(C_\lambda - A(y)). \quad (14)$$

After the solution  $\hat{y}$  is found, we can recover the optimal estimator for the primal problem via  $\tilde{X} = -\nabla f_*(C_\lambda - A(\hat{y}))$ . The dual problem (13) is easier to solve than the primal (12) because the origin  $y = 0$  often lies very close to the solution  $\hat{y}$ . To see this, note that  $y = 0$  produces a candidate estimator  $\tilde{X} = -\nabla f_*(C_\lambda)$  that solves the *chordal* matrix completion problem

$$\tilde{X} = \arg \min \{ \text{tr } C_\lambda X - \log \det X : X \in \mathbb{S}_G^n \},$$

which is a relaxation of the nonchordal problem posed over  $\mathbb{S}_G^n$ . As observed by previous authors (Dahl et al., 2008), this relaxation is a high quality guess, and  $\tilde{X}$  is often “almost feasible” for the original nonchordal problem posed over  $\mathbb{S}_G^n$ , as in  $\tilde{X} \approx P_G(\tilde{X})$ . Some simple algebra shows that  $\|\nabla g(0)\| = \|\tilde{X} - P_G(\tilde{X})\|_F$ , so if  $\tilde{X} \approx P_G(\tilde{X})$  holds true, then the origin  $y = 0$  is close to optimal. Starting from this point, we can expect Newton’s method to rapidly converge at a quadratic rate.

### 3.4. CG converges in $O(1)$ iterations

The most computationally expensive part of Newton’s method is the solution of the Newton direction  $\Delta y$  via the  $m \times m$  system of equations

$$\nabla^2 g(y) \Delta y = -\nabla g(y). \quad (15)$$

The Hessian matrix  $\nabla^2 g(y)$  is fully dense, but matrix-vector products are linear  $O(n)$  time using the algorithms in Section 3.2. This insight motivates solving (15) using an iterative Krylov subspace method like conjugate gradients (CG), which is a *matrix-free* method that requires a single matrix-vector product with  $\nabla^2 g(y)$  at each iteration (Bartlett et al., 1994). Starting from the origin  $p = 0$ , the method converges to an  $\epsilon$ -accurate search direction  $p$  satisfying

$$(p - \Delta y)^T \nabla^2 g(y) (p - \Delta y) \leq \epsilon |\Delta y^T \nabla g(y)|$$

in at most

$$\lceil \sqrt{\kappa_g} \log(2/\epsilon) \rceil \text{ CG iterations,} \quad (16)$$

where  $\kappa_g = \|\nabla^2 g(y)\| \|\nabla^2 g(y)^{-1}\|$  is the condition number of the Hessian matrix (Greenbaum, 1997; Saad, 2003).

Below, we state our main result, which says that the condition number  $\kappa_g$  depends polynomially on the problem data and the quality of the initial point, but is *independent of the problem dimension  $n$  and the accuracy of the current iterate  $\epsilon$* .

**Theorem 6.** *At any  $y$  satisfying  $g(y) \leq g(y_0)$  and  $\nabla g(y)^T (y - y_0) \leq \phi_{\max}$ , the condition number  $\kappa_g$  of the Hessian matrix  $\nabla^2 g(y)$  is bound*

$$\kappa_g \leq 4 \left( 1 + \frac{\phi_{\max}^2 \lambda_{\max}(X_0)}{\lambda_{\min}(\hat{X})} \right)^2. \quad (17)$$

where  $\phi_{\max} = g(y_0) - g(\hat{y})$  is the initial infeasibility,  $\mathbf{A} = [\text{vec } A_1, \dots, \text{vec } A_m]$  is the vectorized data matrix,  $X_0 = -\nabla f_*(C - A(y_0))$ , and  $\hat{X} = -\nabla f_*(C - A(\hat{y}))$ .

*Proof.* See supplementary materials.  $\square$

*Remark 7.* Newton’s method is a descent method, so its  $k$ -th iterate  $y_k$  trivially satisfies  $g(y_k) \leq g(y_0)$ . Technically, the condition  $\nabla g(y_k)^T (y_k - y_0) \leq \phi_{\max}$  can be guaranteed by enclosing Newton’s method within an outer auxiliary path-following loop; see Section 4.3.5 of (Nesterov, 2013). In practice, naive Newton’s method will usually satisfy the condition on its own; see our numerical experiments in Section 4.

Applying Theorem 6 to (16) shows that CG solves each Newton subproblem to  $\epsilon$ -accuracy in  $O(\log \epsilon^{-1})$  iterations. Multiplying this figure by the  $O(\log \log \epsilon^{-1})$  Newton steps to converge yields a *global* iteration bound of  $O(\log \epsilon^{-1} \cdot \log \log \epsilon^{-1}) \approx O(1)$  CG iterations. Multiplying this figure by the  $O(n)$  cost of each CG iteration proves the claimed time complexity in (6). In practice, CG typically converges much faster than this worst-case bound, due to its ability to exploit the clustering of eigenvalues in  $\nabla^2 g(y)$ ; see (Greenbaum, 1997; Saad, 2003). Moreover, accurate Newton directions are only needed to guarantee quadratic convergence close to the solution. During the initial Newton steps, we may loosen the error tolerance for CG for a significant speed-up. Inexact Newton steps can be used to obtain a speed-up of a factor of 2-3.

### 3.5. The full algorithm

To summarize, we begin by computing a chordal embedding  $\tilde{G}$  for the sparsity pattern  $G$  of  $C_\lambda$ , using the code snippet in Figure 1. We use the embedding to reformulate (4) as (11), and solve the unconstrained problem  $\hat{y} = \min_y g(y)$  defined in (14), using Newton’s method

$$y_{k+1} = y_k + \alpha_k \Delta y_k, \quad \Delta y_k \equiv -\nabla^2 g(y_k)^{-1} \nabla g(y_k)$$

starting at the origin  $y_0 = 0$ . The function value  $g(y)$ , gradient  $\nabla g(y)$  and Hessian matrix-vector products are all

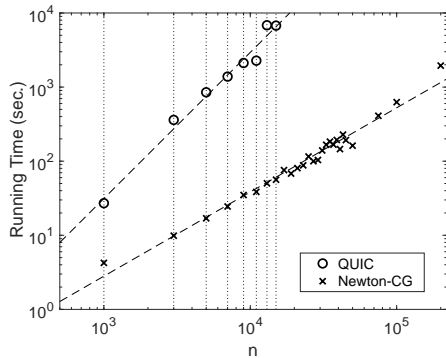


Figure 2. CPU time Newton-CG vs QUIC for case study 1.

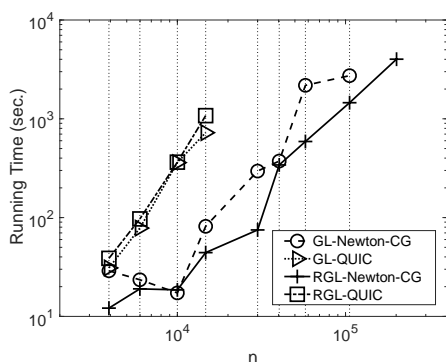


Figure 3. CPU time Newton-CG vs QUIC for case study 2.

evaluated using the numerical recipes described by (Andersen et al., 2013b).

At each  $k$ -th Newton step, we compute the Newton search direction  $\Delta y_k$  using conjugate gradients. A loose tolerance is used when the Newton decrement  $\delta_k = |\Delta y_k^T \nabla g(y_k)|$  is large, and a tight tolerance is used when the decrement is small, implying that the iterate is close to the true solution. Once a Newton direction  $\Delta y_k$  is computed with a sufficiently large Newton decrement  $\delta_k$ , we set the step-size  $\alpha_k$  to be the first instance of the sequence  $\{1, \rho, \rho^2, \rho^3, \dots\}$  that satisfies the Armijo–Goldstein condition

$$g(y + \alpha \Delta y) \leq g(y) + \gamma \alpha \Delta y^T \nabla g(y),$$

in which  $\gamma \in (0, 0.5)$  and  $\rho \in (0, 1)$  are line search parameters. Our implementation used  $\gamma = 0.01$  and  $\rho = 0.5$ . We complete the step and repeat the process, until convergence.

We terminate the outer Newton’s method if the Newton decrement  $\delta_k$  falls below a threshold. This implies either that the solution has been reached, or that CG is not converging to a good enough  $\Delta y_k$  to make significant progress. The associated estimator for  $\Sigma^{-1}$  is recovered by evaluating  $\hat{X} = -\nabla f_*(C_\lambda - A(\hat{y}))$ .

## 4. Numerical Results

Finally, we benchmark our algorithm against QUIC (Hsieh et al., 2014), commonly considered the fastest solver for graphical lasso or RGL<sup>1</sup>. We consider two case studies. The first case study numerically verifies the claimed  $O(n)$  complexity of our MDMC algorithm on problems with a nearly-banded structure. The second case study performs the full threshold-MDMC procedure for graphical lasso and RGL, on graphs collected from real-life applications. All experiments are performed on a laptop computer with an Intel Core i7 quad-core 2.50 GHz CPU and 16GB RAM. The reported results are based on a serial implementation in MATLAB-R2017b. Both our Newton decrement threshold and QUIC’s convergence threshold are  $10^{-7}$ .

### 4.1. Case Study 1: Banded Patterns

The first case study aims to verify the claimed  $O(n)$  complexity of our algorithm for MDMC. Here, we avoid the proposed thresholding step, and focus solely on the MDMC (4) problem. Each sparsity pattern  $G$  is a corrupted banded matrices with bandwidth 101. The off-diagonal nonzero elements of  $C$  are selected from the uniform distribution in  $[-2, 0)$  and then corrupted to zero with probability 0.3. The diagonal elements are fixed to 5. Our numerical experiments fix the bandwidth and vary the number of variables  $n$  from 1,000 to 200,000. A time limit of 2 hours is set for both algorithms.

Figure 2 compares the running time of both algorithms. A log-log regression results in an empirical time complexity of  $O(n^{1.1})$  for our algorithm, and  $O(n^2)$  for QUIC. The extra 0.1 in the exponent is most likely an artifact our MATLAB implementation. In either case, QUIC’s quadratic complexity limits it to  $n = 1.5 \times 10^4$ . By contrast, our algorithm solves an instance with  $n = 2 \times 10^5$  in less than 33 minutes. The resulting solutions are extremely accurate, with optimality and feasibility gaps of less than  $10^{-16}$  and  $10^{-7}$ , respectively.

### 4.2. Case Study 2: Real-Life Graphs

The second case study aims to benchmark the full thresholding-MDMC procedure for sparse inverse covariance estimation on real-life graphs. The actual graphs (i.e. the sparsity patterns) for  $\Sigma^{-1}$  are chosen from *SuiteSparse Matrix Collection* (Davis & Hu, 2011)—a publicly available dataset for large-and-sparse matrices collected from real-world applications. Our chosen graphs vary in size from  $n = 3918$  to  $n = 201062$ , and are taken from ap-

<sup>1</sup>Two other widely-used algorithms are GLASSO (Friedman et al., 2008) and BIGQUIC (Hsieh et al., 2013). On a serial machine and for the problem sizes that we consider, we found both to be slower than QUIC.

#	file name	type	$n$	$m$	$m/n$	Newton-CG			QUIC	diff. gap	speed-up
						sec	gap	feas	sec		
1	freeFlyingRobot-7	GL	3918	20196	5.15	28.9	5.7e-17	2.3e-7	31.0	3.9e-4	1.07
1	freeFlyingRobot-7	RGL	3918	20196	5.15	12.1	6.5e-17	2.9e-8	38.7	3.8e-5	3.20
2	freeFlyingRobot-14	GL	5985	27185	4.56	23.5	5.4e-17	1.1e-7	78.3	3.8e-4	3.33
2	freeFlyingRobot-14	RGL	5985	27185	4.56	19.0	6.0e-17	1.7e-8	97.0	3.8e-5	5.11
3	cryg10000	GL	10000	170113	17.0	17.3	5.9e-17	5.2e-9	360.3	1.5e-3	20.83
3	cryg10000	RGL	10000	170113	17.0	18.5	6.3e-17	1.0e-7	364.1	1.9e-5	19.68
4	epb1	GL	14734	264832	18.0	81.6	5.6e-17	4.3e-8	723.5	5.1e-4	8.86
4	epb1	RGL	14734	264832	18.0	44.2	6.2e-17	3.3e-8	1076.4	4.2e-4	24.35
5	bloweya	GL	30004	10001	0.33	295.8	5.6e-17	9.4e-9	*	*	*
5	bloweya	RGL	30004	10001	0.33	75.0	5.5e-17	3.6e-9	*	*	*
6	juba40k	GL	40337	18123	0.44	373.3	5.6e-17	2.6e-9	*	*	*
6	juba40k	RGL	40337	18123	0.44	341.1	5.9e-17	2.7e-7	*	*	*
7	bayer01	GL	57735	671293	11.6	2181.3	5.7e-17	5.2e-9	*	*	*
7	bayer01	RGL	57735	671293	11.6	589.1	6.4e-17	1.0e-7	*	*	*
8	hcircuit	GL	105676	58906	0.55	2732.6	5.8e-17	9.0e-9	*	*	*
8	hcircuit	RGL	105676	58906	0.55	1454.9	6.3e-17	7.3e-8	*	*	*
9	co2010	RGL	201062	1022633	5.08	4012.5	6.3e-17	4.6e-8	*	*	*

Table 1. Details of case study 2. Here, “ $n$ ” is the size of the covariance matrix, “ $m$ ” is the number of edges added to make its sparsity graph chordal, “sec” is the running time in seconds, “gap” is the optimality gap, “feas” is the feasibility the solution, “diff. gap” is the difference in duality gaps for the two different methods, and “speed-up” is the fact speed-up over QUIC achieved by our algorithm.

plications in chemical processes, material science, graph problems, optimal control and model reduction, thermal processes and circuit simulations.

For each sparsity pattern  $G$ , we design a corresponding  $\Sigma^{-1}$  as follows. For each  $(i, j) \in G$ , we select  $(\Sigma^{-1})_{i,j} = (\Sigma^{-1})_{j,i}$  from the uniform distribution in  $[-1, 1]$ , and then corrupt it to zero with probability 0.3. Then, we set each diagonal to  $(\Sigma^{-1})_{i,i} = 1 + \sum_j |(\Sigma^{-1})_{i,j}|$ . Using this  $\Sigma$ , we generate  $N = 5000$  samples i.i.d. as  $\mathbf{x}_1, \dots, \mathbf{x}_N \sim \mathcal{N}(0, \Sigma)$ . This results in a sample covariance matrix  $C = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T$ .

We solve graphical lasso and RGL with the  $C$  described above using our proposed soft-thresholding-MDMC algorithm and QUIC, in order to estimate  $\Sigma^{-1}$ . In the case of RGL, we assume that the graph  $G$  is known *a priori*, while noting that 30% of the elements of  $\Sigma^{-1}$  have been corrupted to zero. Our goal here is to discover the location of these corrupted elements. In all of our simulations, the threshold  $\lambda$  is set so that the number of nonzero elements in the estimator is roughly the same as the ground truth. We limit both algorithms to 3 hours of CPU time.

Figure 3 compares the CPU time of both two algorithms for this case study; the specific details are provided in Table 1. A log-log regression results in an empirical time complexity of  $O(n^{1.64})$  and  $O(n^{1.55})$  for graphical lasso and RGL using our algorithm, and  $O(n^{2.46})$  and  $O(n^{2.52})$  for the same using QUIC. The exponents of our algorithm are  $\geq 1$  due to the initial soft-thresholding step, which is quadratic-time on a serial computer, but  $\leq 2$  because procedure is dominated by the solution of the MDMC. Both algorithms solve graphs with  $n \leq 1.5 \times 10^4$  within the allotted time limit, though our algorithm is 11 times faster on average. Only our algorithm is able to solve the estimation problem with  $n \approx 2 \times 10^5$  in a little more than an hour.

To check whether thresholding-MDMC really does solve graphical lasso and RGL, we substitute the two sets of estimators back into their original problems (1) and (5). The corresponding objective values have a relative difference  $\leq 4 \times 10^{-4}$ , suggesting that both sets of estimators are about equally optimal. This observation verifies our claims in Theorem 2 and Corollary 3 that (1) and (5): thresholding-MDMC does indeed solve graphical lasso and RGL.

## 5. Conclusions

Graphical lasso is a widely-used approach for estimating a covariance matrix with a sparse inverse from limited samples. In this paper, we consider a slightly more general formulation called *restricted* graphical lasso (RGL), which additionally enforces a prior sparsity pattern to the estimation. We describe an efficient approach that substantially reduces the cost of solving RGL: 1) soft-thresholding the sample covariance matrix and projecting onto the prior pattern, to recover the estimator’s sparsity pattern; and 2) solving a maximum determinant matrix completion (MDMC) problem, to recover the estimator’s numerical values. The first step is quadratic  $O(n^2)$  time and memory but embarrassingly parallelizable. If the resulting sparsity pattern is *sparse* and *chordal*, then the second step can be performed using the Newton-CG algorithm described in this paper in linear  $O(n)$  time and memory. The algorithm is tested on both synthetic and real-life data, solving instances with as many as 200,000 variables to 7-9 digits of accuracy within an hour on a standard laptop computer.

**Acknowledgements.** This work was supported by the ONR grants N00014-17-1-2933 and N00014-15-1-2835, DARPA grant D16AP00002, and AFOSR grant FA9550-17-1-0163.



## References

- Agrawal, A., Klein, P., and Ravi, R. Cutting down on fill using nested dissection: Provably good elimination orderings. In *Graph Theory and Sparse Matrix Computation*, pp. 31–55. Springer, 1993. 3.1
- Andersen, M. S., Dahl, J., and Vandenberghe, L. Implementation of nonsymmetric interior-point methods for linear optimization over sparse matrix cones. *Mathematical Programming Computation*, 2(3):167–201, 2010. 1.2, 1.3, 3, 3.2
- Andersen, M. S., Dahl, J., and Vandenberghe, L. CVXOPT: A Python package for convex optimization. Available at [cvxopt.org](http://cvxopt.org), 54, 2013a. 1.3
- Andersen, M. S., Dahl, J., and Vandenberghe, L. Logarithmic barriers for sparse matrix cones. *Optimization Methods and Software*, 28(3):396–423, 2013b. 1.2, 1.3, 3, 3.2, 3.5
- Banerjee, O., Ghaoui, L. E., and d’Aspremont, A. Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data. *Journal of Machine Learning Research*, 9:485–516, 2008. 1
- Barrett, R., Berry, M. W., Chan, T. F., Demmel, J., Donato, J., Dongarra, J., Eijkhout, V., Pozo, R., Romine, C., and Van der Vorst, H. *Templates for the solution of linear systems: building blocks for iterative methods*, volume 43. Siam, 1994. 3.4
- Croft, D., O’Kelly, G., Wu, G., Haw, R., Gillespie, M., Matthews, L., Caudy, M., Garapati, P., Gopinath, G., Jassal, B., and Jupe, S. Reactome: a database of reactions, pathways and biological processes. *Nucleic acids research*, 39:691–697, 2010. 1.2
- Dahl, J., Vandenberghe, L., and Roychowdhury, V. Covariance selection for nonchordal graphs via chordal embedding. *Optimization Methods & Software*, 23(4):501–520, 2008. 1.2, 1.3, 3, 3.1, 3.2, 3.3
- Davis, T. A. and Hu, Y. The university of florida sparse matrix collection. *ACM Transactions on Mathematical Software (TOMS)*, 38(1):1, 2011. 4.2
- Durlauf, S. N. Nonergodic economic growth. *The Review of Economic Studies*, 60(2):349–366, 1993. 1
- Egilmez, H. E., Pavez, E., and Ortega, A. Graph learning from data under laplacian and structural constraints. *IEEE Journal of Selected Topics in Signal Processing*, 11(6):825–841, 2017. 1.3
- Fattahi, S. and Sojoudi, S. Graphical lasso and thresholding: Equivalence and closed-form solutions. <https://arxiv.org/abs/1708.09479>, 2017. 1.1, 1.1, 1.2, 1.2, 1
- Fattahi, S., Zhang, R. Y., and Sojoudi, S. Sparse inverse covariance estimation for chordal structures. <https://arxiv.org/abs/1711.09131>, 2018. 1.1, 1.1, 1.2, 1.2
- Friedman, J., Hastie, T., and Tibshirani, R. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441, 2008. 1, 1, 1.3, 1
- George, A. and Liu, J. W. The evolution of the minimum degree ordering algorithm. *Siam review*, 31(1):1–19, 1989. 3.1
- Gilbert, J. R. Some nested dissection order is nearly optimal. *Information Processing Letters*, 26(6):325–328, 1988. 3.1
- Grechkin, M., Fazel, M., Witten, D. M., and Lee, S.-I. Pathway graphical lasso. *AAAI*, pp. 2617–2623, 2015. 1.3
- Greenbaum, A. *Iterative methods for solving linear systems*, volume 17. Siam, 1997. 3.4, 3.4
- Honorio, J., Samaras, D., Paragios, N., Goldstein, R., and Ortiz, L. E. Sparse and locally constant gaussian graphical models. *Advances in Neural Information Processing Systems*, pp. 745–753, 2009. 1.2
- Hsieh, C.-J., Sustik, M. A., Dhillon, I. S., Ravikumar, P. K., and Poldrack, R. Big & quic: Sparse inverse covariance estimation for a million variables. In *Advances in neural information processing systems*, pp. 3165–3173, 2013. 1, 1.3, 1
- Hsieh, C. J., Sustik, M. A., Dhillon, I. S., and Ravikumar, P. Quic: quadratic approximation for sparse inverse covariance estimation. *Journal of Machine Learning Research*, 15(1):2911–2947, 2014. 1, 1.3, 4
- Huang, J. and Zhang, T. The benefit of group sparsity. *The Annals of Statistics*, 38(4):1978–2004, 2010. 1
- Li, J., Andersen, M. S., and Vandenberghe, L. Inexact proximal newton methods for self-concordant functions. *Mathematical Methods of Operations Research*, 85(1):19–41, 2017. 1.3
- Li, S. Z. Markov random field models in computer vision. *European conference on computer vision*, pp. 351–370, 1994. 1
- Manning, C. D. and Schütze, H. *Foundations of statistical natural language processing*. MIT press, 1999. 1
- Mazumder, R. and Hastie, T. Exact covariance thresholding into connected components for large-scale graphical lasso. *Journal of Machine Learning Research*, 13:781–794, 2012. 1.1, 1.2

- Meinshausen, N. and Bühlmann, P. High-dimensional graphs and variable selection with the lasso. *The annals of statistics*, 1436–1462, 2006. 1
- Milanfar, P. A tour of modern image filtering: New insights and methods, both practical and theoretical. *IEEE Signal Processing Magazine*, 30(1):106–128, 2013. 1.3
- Negahban, S. and Wainwright, M. J. Joint support recovery under high-dimensional scaling: Benefits and perils of  $l_{1,\infty}$ -regularization. *Proceedings of the 21st International Conference on Neural Information Processing Systems*, pp. 1161–1168, 2008. 1
- Nesterov, Y. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2013. 7
- Obozinski, G., Wainwright, M. J., and Jordan, M. I. Union support recovery in high-dimensional multivariate regression. *Communication, Control, and Computing, 2008 46th Annual Allerton Conference on*, pp. 21–26, 2008. 1
- Oztoprak, F., Nocedal, J., Rennie, S., and Olsen, P. A. Newton-like methods for sparse inverse covariance estimation. In *Advances in neural information processing systems*, pp. 755–763, 2012. 1.3
- Park, D. and Rilett, L. R. Forecasting freeway link travel times with a multilayer feedforward neural network. *Computer Aided Civil and Infrastructure Engineering*, 14(5):357–367, 1999. 1.2
- Ravikumar, P., Wainwright, M. J., Raskutti, G., and Yu, B. High-dimensional covariance estimation by minimizing  $l_1$ -penalized log-determinant divergence. *Electronic Journal of Statistics*, 5:935–980, 2011. 1
- Rolfs, B., Rajaratnam, B., Guillot, D., Wong, I., and Maleki, A. Iterative thresholding algorithm for sparse inverse covariance estimation. In *Advances in Neural Information Processing Systems*, pp. 1574–1582, 2012. 1.3
- Rose, D. J. Triangulated graphs and the elimination process. *Journal of Mathematical Analysis and Applications*, 32(3):597–609, 1970. 3.1
- Rothman, A. J., Bickel, P. J., Levina, E., and Zhu, J. Sparse permutation invariant covariance estimation. *Electronic Journal of Statistics*, 2:494–515, 2008. 1
- Saad, Y. *Iterative methods for sparse linear systems*, volume 82. Siam, 2003. 3.4, 3.4
- Sojoudi, S. Equivalence of graphical lasso and thresholding for sparse graphs. *Journal of Machine Learning Research*, 17(115):1–21, 2016. 1.1, 1.2
- Treister, E. and Turek, J. S. A block-coordinate descent approach for large-scale sparse inverse covariance estimation. In *Advances in neural information processing systems*, pp. 927–935, 2014. 1.3
- Vandenberghe, L., Andersen, M. S., et al. Chordal graphs and semidefinite optimization. *Foundations and Trends® in Optimization*, 1(4):241–433, 2015. 3, 3.1
- Wainwright, M. J. Sharp thresholds for high-dimensional and noisy sparsity recovery using  $l_1$ -constrained quadratic programming (lasso). *IEEE transactions on information theory*, 55(5):2183–2202, 2009. 1
- Yang, E., Lozano, A. C., and Ravikumar, P. K. Elementary estimators for graphical models. *Advances in neural information processing systems*, pp. 2159–2167, 2014. 1.3
- Yannakakis, M. Computing the minimum fill-in is NP-complete. *SIAM Journal on Algebraic Discrete Methods*, 2(1):77–79, 1981. 3.1
- Yuan, M. and Lin, Y. Model selection and estimation in the Gaussian graphical model. *Biometrika*, pp. 19–35, 2007. 1